

An Introduction to Bioinformatics Infrastructures:

Text Mining and Information Extraction Applications for Bioinformatics and Systems Biology

Plant Bioinformatics, Systems and Synthetic Biology Summer School 27-31 July 2009 - University of Nottingham, UK

> Martin Krallinger, Spanish National Cancer Research Centre - CNIO *mkrallinger@cnio.es*



Talk Outline / Topcis (I)

- Bioinformatics infrastructures
- Integration of heterogeneous data types
- Bioinformatics resources
- Importance and use of scientific literature data
- Manual literature curation process for building systems biology resources
- Annotation types
- Building literature curation workflows
- Relevance of text mining strategies in the context of SB 2



Talk Outline / Topics (II)

- Short intro to text mining and NLP
- Short overview of existing BioNLP application types
- Implementing a text mining system: basic steps
- The PLAN2L literature mining tool

Bioinformatics & biological projects

BIOINFORMATICS

Vol. 18 no. 12 2002 Pages 1551–1552

Editorial

BIOINFORMATICS: BIOLOGY BY OTHER MEANS

The success of bioinformatics in its application to genomics and proteomics has complicated the relationship of computation with experimental biology. There is a need to attend to our pressing needs of bioinformatics applications without forgetting other, perhaps less evident but equally important, aspects of computation in biology.

BIOINFORMATICS IN THE STUDY OF GENERAL BIOLOGICAL PROBLEMS

A much deeper aspect of bioinformatics extends towards the study of fundamental biological questions, such as gene assembly, protein folding and the nature of functional specificity. Such issues extend beyond the current perception of bioinformatics as a support discipline and address aspects of biological complexity, including the simulation of cellular systems and molecular interaction networks. The contribution of bioinformatics to these areas is related to the development

ALL biological projects need or will need Bioinformatics (.. as soon as they enter into genomics):

- as **resource** (databases and software)
- as support for design, organization & interpretation of the data
- in the research team for the specific scientific project

Bioinformaticians are scientists working in:

- developing methods (Bioinformatics as a research area)
- developing resources e.g. databases (Bioinformatics as technology)
- Embedded in biology/Biotech/Biomed (the single bioinformatician syndrome)



- To construct and operate a <u>sustainable infrastructure</u> for biological information in Europe,
- To <u>support</u> life science research and its translation to medicine and the environment, the bio-industries and society.
- Partners: 32 partners, 13 member states
- **Funding:** 4.5 M€ from EU FP7
- Deliverable: Consortium agreement to <u>define the</u> <u>scope</u> of the infrastructure and how it will be constructed







- Optimal Data Management
 - Coordinated Data Resources with improved <u>access</u>
 - Integration and interoperability of diverse heterogeneous data
 - Good Value for Money
- Forge Links to data in other related domains
- A single European voice in international collaborations to influence global decisions and <u>maintain open access to data</u>
- Enhance European competitiveness in bioscience industries
- Address need for Increased Funding & its Coordination



The Preparatory Phase project



Elixir is organised into 14 work packages which have committees of (mainly) European experts associated with them.

- 1. Project management
- 2. Data providers
- 3. User communities
- 4. Organisation and Legal
- 5. Funding
- 6. Physical infrastructure
- 7. Data interoperability

8. Literature

- 9. Healthcare
- 10.Chemistry, Plants, Agriculture & Environment
- 11.Training
- 12.Tools integration
- 13. Feasibility studies
- 14.Reporting and negotiation

Why do we need ELIXIR?

(Why do we need bioinformatics infrastructures)

USA

Europ

- Data Growth
- Global context
- Very large user community:
 - 3.3 m web hits/day
 - 20,000 unique users per day
- Need to preserve data and make <u>accessible</u> to all

Japan

- Impact on Medicine, Agriculture & Biotechnology
- Impact on society & bioindustries
- Need for increased funding for biodata resources







EBI Hits in 2008



WP3: User Communities



- User Survey: 1000 responses
 - Long term support essential
 - Top 3 challenges:
 - Data integration; Format compatibility; Website usability
 - Concerns
 - Data <u>quality</u> and measures; Quality of tools; Training
- Need to consider different needs in different countries
- Need for a plan for long-term <u>maintenance</u> of computational tools
 - Create mechanisms for long-term maintenance of bioinformatics tools
 - user-friendly & machine-friendly interfaces
- Need for <u>standards</u> for formats and integration
 - Increased integration of databases, tools and between infrastructure domains
- Need to provide mechanisms for prioritisation of need for resources





Total European effort



- 200 Databases
- 700 People
- 100 Institutions
- 60 million web hits per month
- Total investment to date €308 million
- Annual cost €35 million

RECOMMENDATION

Coordination and prioritisation, as well as stable funding, is needed for many of these resources

ESFRI



Biology Research Infrastructure proposals.



INSTRUCT	Infrafrontier	EATRIS	BBMRI	ECRIN			
(Structural biology)	(Mouse)	(Translational Research)	(Biobanking) ↓	(Clinical Trials) ↓			
ELIXIR							

(Biological Information)

WP 10: Chemistry, Plants, Agriculture & Environment



- Support / extend current core resources for
 - Nucleotide/protein sequence, genomes, structures, interactions etc.
- Selected specialist resources migrated to Elixir infrastructure
 - Reduce complexity of informatics landscape, maintain functionality
 - Integration allows mining of combined data
- Adopt key data standards and work for common infrastructure
 - Link to other ESFRI, non ESFRI European projects
 - Link to non European initiatives (NSF/iPlant, DOE/Camera)
- Free access to Elixir data and core analysis tools
 - Web based queries, programmatic access, download

WP11: Training



Identified training issues in Europe:

- Little or no coordination
- Rapid evolution of bioinformatics resources
- Lack of a **centralised** body for guidance;
- Lack of recognition of the importance of bioinformatics user training, even within the bioinformatics community.

Elixir recommendations:

Link the **development** of data resources to the provision of **training** materials;

Create a training support unit that will:

a) provide a centralised training registry;

b) provide support for trainers throughout Europe

c) develop benchmarking and evaluation systems;

d) provide mechanisms for developing new training programmes

e) act as a single point of contact for national and pan-European training

Elixir WP8: Scientific Literature Interdisciplinary Interactions



Chair: Alfonso Valencia (CNIO)

Co-Chairs: Dietrich Rebholz-Schuhmann & Peter Stoehr (EMBL-EBI)

Initial committee

- Robert Kiley, Wellcome Trust
- Carole Goble, U. Manchester
- Larry Hunter, UCHSColorado
- Manuel Peitsch, SIB
- Matthew Cockerill, BMC
- Jun'ichi Tsujii, NaCTeM and U. Tokyo
- Timo Hannay, Nature PG

Addtional Contributions

Ian Dix, Astrazeneca Ian Harrow, Pfizer Udo Hahn, U. Jena Sophia Ananiadou , NacTeM Patrick Ruch, Geneva University Christopher Bake, New Brunswick U. Juliane Fluck, Fraunhofer Anita Burgun, Rennes University and Kostas Repanas (CNIO) WP Coordinator

European Life-science Infrastructure for Biological Information (Elixir) WP 8: Scientific Literature Interdisciplinary Interactions

D8.1 A report summarising the current

(1) <u>status</u> of literature repositories throughout Europe and recommendations for the future

(2) infrastructure <u>needs</u> in Europe to establish an informationsharing platform to integrate databases and literature for (*) experts and non-experts, with

(3a) specific reference to the <u>provision</u> of literature from repositories commonly used in biological information extraction and

(3bi) tools for access to the literature, for

(3bii) data representation and for

(3biii) interaction with end users.





Emergence IT layout



MIT Repository of Parts



article **Protein Coding Regions**

The coding sequence has been modified from wild type c2 to contain an LVA tag.

	-?-	Name	Protein	Description		Tag -?-	Lengt
	AW	BBa_J23012		SpecR open reading frame basic part	Forward		88
	AW	BBa_J31000	Hin	Generates Hin from Salmonella typhimurium			57
_	AW	BBa_J31001		Hin invertase tagged with LVA (HinLVA)			61
	AW	BBa_J31002		Kanamycin Resistance backwards (KanB)			81
	ΑW	BBa_J31006		Tet Resistance Backwards (TetB)			119
-	AW	BBa_J31007		Tet Resistance Forwards (TetF)			119
	AW	BBa_J45001		SAMT enzyme			115
	AW	BBa_J45002		BAMT (SAM:Benzoic Acid Carboxyl Methyltransferase)			109
	AW	BBa_J45004		BSMT: converts salicylic acid to methyl salicylate (wintergreen)			107
	AW	BBa 145008		Branched-chain amino acid transaminase (BAT2); used in biosynthesis of banana scent			112
			•	2 Create an account or log	n		- 113
scussio	n	edit history		DNA Available			158
3a_0	C005	53		Experience: Non	9		173
Maia N	lahoney			Entered: Antiq	uity		42
	D 22 0	o					124
501,1	-22 0	2			rward	LVA	112
epresso	or proteir	n coding sequence is a	720 base-pair seque	nce with the standard RBS-compatible BioBrick prefix and the standard BioBrick suffix sections on its ends. It binds to the P22 c2	rward	None	105
quence	e, BBa_F	10053. The sequence co	ontains a LVA tag for	faster degradation.	nward	1.1/4	66

iump to part BBa_

jump to part BBa_ navigation Main Page Browse Part Type

Usage and Biology P22 c2 is a member of the lamboid cl protein family.

and two TAA stop codons.

Part Main Page Part Design = Experience Hard Information

Physical DNA navigation

BBa_C0053

Main Page Browse Part Types

iGEM Wiki Community portal

= Recent changes

Recent part changes

resources

User Accounts

What links here

Related changes

Upload file

Special pages





article discussion

Part:BBa_C0053:Physical DNA

Designed by Maia Mahoney

Repressor, P22 c2

Physical DNA statistics

	Plasmic	Plasmid Length	Part and Plasmid	VF2 - VR
	pSB1A2	2079	2766	925
	BBa_J6	1031 10662	11349	None
	BBa_152	2001 1090	1777	None
	BBa_J6	1003 3016	3703	1862
	BBa_J2	3018 2977	3664	1823
Part Length: 687bp	BBa_J5	2017 4101	4788	1915
,	BBa_J6	1002 3002	3689	1848
ACTG Ratios	BBa_J6	1009 5130	5817	None
a: 0.323, c: 0.199, t: 0.231, g: 0.246	BBa_J6	1007 2888	3575	None
	BBa_I51	001 2438	3125	2057
260:280 ratio: 2.378	pSB1A3	2157	2844	1003
	pSB1A7	2431	3118	1277
	pSB1AC	3 3055	3742	1003
	pSB1AK	3 3189	3876	1003
	pSB1AT	3 3446	4133	1003
	pSB2K3	4425	5112	1003
	pSB4A3	3339	4026	1003
	BBa_150	040 2226	2913	None

Existing versions

Validated collection and assembly version exists Validated collection exists curated

Start a new sequence analysis

arch	This part may		
	Library		
Co. Search	IGEM 200		
uu search			

part may be found in these wells/tubes			Show all locations				
_ibrary	Well	Plate	Plasmid	Cell			
GEM 2006	5K	iGEM2006 DNA-1	pSB1A2	V1004	Available		

MaDAS principal features

- MaDAS allows users to <u>add</u>, <u>edit</u>, or <u>remove</u> self generated sequence annotations
- Allows to <u>upload</u> multiple annotations from different sources.
- Provides a <u>security</u> system based on projects.
 The annotations could be public or only available for the project members.
- Provides an <u>interface</u> to manage projects, users and collections of annotations.

Collaborative features

- **Project based** system. Users can create their own projects or participate in projects hosted in MaDAS.
- Projects can be public or private, in private projects the project leader decide who can view or edit the project annotations.
- The notification system inform about: new projects, new annotations, new users or new plugins.
- Searches by: category, project leader, institutions, etc

MaDas Manual Sequence Annotation System



MaDAS modules

MaDAS is composed by:

- •"The core" which provide different APIs in order to facilitated the development of plug-ins and the communication between them.
- Data Source plug-ins
- DAS server plug-ins
- Visualization plug-ins



Data source plug-ins

Manage Reference plug-in: We use the DAS reference sequence concept (http://www.biodas.org/wiki/DAS/1/Overview#.5BReference.5D_Sequence) to describe a biological sequence that will be annotated.

Setup Ensembl genome, a collection of proteins , a new sequenced genome or just a DNA/protein fragment.

Load GFF plug-in: This plug-in allows users to upload GFF files to the system.

Manage DAS Tracks plug-in: Through this plug-in users can add annotations provided by any DAS server

Load chip plug-in: This plug-in allows experimentalist to map Affymetrix or Illumina microarray probes to a human reference sequence stored in MaDAS. Probe associated genes and proteins are also mapped.

Load Gene expression plug-in: Allows users to upload data from a gene expression experiments.

Map Annotations plug-in: Using this plug-in is possible to add new annotations just mapping existing annotations to other online resource. For example if we have a gene track is possible to setup a disease track mapping these genes to OMIM diseases. This plug-in use several mapping services to map the annotations (Biomart, Uniprot Database mapping, PICR, ID converter)

Treefam plug-in: This is an example of a very specific plug-in, which allows to information form Treefam).

Bionemo plug-in: import information stored in the Bionemo database (Bopdegradation and gene control reactions)

Manage annotations plug-in: to remove or inactivate an entire set of annotations.



Introducing expert annotations and consolidating them in databases/visualization



How to exchange annotations



- ✓ Distributed annotation system (DAS) protocol. (MR)
- ✓ Web services. (MR)
- ✓ Database dump. (MR)
- ✓ Biological Web Elements and Registry Embed Code. (HR)
- MR = Machine readable HR = Human readable

Integration of heterogeneous data types



Text mining covers multiple topics



Importance of literature data for Biology

- Life sciences -> generates heterogeneous data types (sequence, structure,..)
- Natural language used for **communicating** scientific discoveries.
- Natural language texts amenable for direct human interpretation
- Natural language not only in scientific articles, but also patents, reports, newswire, database records, controlled vocabularies (GO terms),...
- Functional information & annotations directly or indirectly derived from the literature (curation and electronic annotation).
- Databases are generally only capable of covering a small fraction of the biological context information that can be encountered in the literature.
- **Contextual information** of experimental results (cell line, tissue, conditions).
- User demands of better information access (beyond keyword searches)
- Rapid growth of information, manual information extraction not efficient.

Literature and the scientific discovery process

- Define the biological question
 Biology
- Select the actual target being studied
- Extract information relevant for experimental set up
- Locate relevant resources
- Essential to understand and interpret the resulting data
- Draw conclusions about new discoveries
- Communicated to the scientific community using publications in peer-reviewed journals
- Resource for clinical decision support in evidencebased clinical practice
 Leaful information for diagnostic side
- Useful information for diagnostic aids





Drug discovery and target selection

- Identifying adverse drug effect
- Competitive intelligence and knowledge management
- Global view of the current research state & monitor trends to ensure optimal resource allocation Funding
- Find domain experts for specific topics for the peer-review process & detecting potential cases of plagiarism Publ.

34

Pharma

Literature Gold Standard datasets / DBs



Biocuration: manual literature annotations & databases


Curation challenge I: growing number of CV terms





Curation challenge II: growing <u>number of on</u>tologies

number of ontologies						
Domain	- <u>File</u>					
Biological imaging methods	image.obo	Mosquito insecticide resistance				
Biological process	<u>gene ontology.obo</u> 🌋	Mouse adult gross anatomy				
BRENDA tissue / enzyme source		Mouse gross anatomy and develop	mosquito insecticide resistance.obo			
C. elegans development	worm development.obo	Multiple alignment	adult mouse anatomy.obo			
C. elegans gross anatomy	worm anatomy.obo	NCBL organismal classification	EMAP.obo			
C. elegans phenotype	worm phenotype.obo	NCI Thesaurus	mouse pathology.obo			
Cell type	cell.obo	NMR-instrument specific component	mao.obo			
Cellular component	gene ontology.obo 🎳	investigations	taxonomy.dat			
Cereal plant development	cereals development.obo	OBO relationship types	EVS ftp site			
Cereal plant trait	plant trait.obo	Ontology for biomedical investigation	n <u>mr.owl</u>			
Chemical entities of biological interest	chebi.obo	Pathway ontology	ro.obo 🌋			
Common Anatomy Reference Ontology	caro.obo	Phenotypic quality	OBLow A			
Dictyostelium discoideum anatomy	dictyostelium anatomy.obo	Physico-chemical methods and prop	nathway obo			
Drosophila development	fly development.obo	Physico-chemical process	guality.obo			
<u>Drosophila gross anatomy</u>	fly anatomy.obo	Plant environmental conditions	fix.obo			
Environment Ontology	envo.obo	Plant growth and developmental sta	rex.obo			
Event (INOH pathway ontology)	event.obo	Plant structure	environment ontology.obo			
Evidence codes	evidence code.obo	<u>Plasmodium life cycle</u>	po temporal.obo			
eVOC (Expressed Sequence Annotation for Humans)	evoc.obo.tar (v2./)	Protein domain	po anatomy.obo			
<u>Fly taxonomy</u>	flybase controlled vecabulary obe	Protein modification	PLO ontology			
FlyBase Controlled Vocabulary	fma obo obo	protein ontology	InterPro ETP directory			
Foundational Model of Anatomy (subset)	fungal anatomy obo	Protein-protein interaction	psi-mod.obo			
Fungal gross anatomy	protege source	Proteomics data and process prove	pro.obo			
Habronattus courtship	human dev anat abstract.obo	Sample processing and separation	psi-mi.obo			
Human developmental anatomy, abstract version	human dev anat staged.obo	Sequence types and features	ProPreO.owl			
Human developmental anatomy, timed version	human disease.obo	Suddested Ontology for Pharmacod	sep.obo			
<u>Human disease</u>	<u>B84670B0</u>	- I I I I I I I I I I I I I I I I I I I	so.obo 🍅			
Loggerhead nesting	zea mays anatomy.obo		<u>sopharm.owl</u>			
<u>Maize gross anatomy</u>	mammalian phenotype.obo	Tick gross anatomy	SBO OWL.owl			
<u>Mammalian phenotype</u>	medaka ontology.obo	Unit	teleost anatomy.obo 🌋			
Medaka fish anatomy and development	MGEDOntology.owl	Xenopus anatomy and developmen	tick_apatomy.obo			
Microarray experimental conditions	gene ontology.obo 🍣		unit abo	,		
Molecular function	molecule role.obo	Correcto				
Molecule role (INOH Protein name/family name		Formats	(OBO, OVVL, XIVIL, RDF)			
	mosquito anatomy.obo	(http://w	www.obofoundry.org)	28		
Mosquito gross anatomy	mosquito insecticide resistance.obc	• (IIIIp.//w	ww.obolouliury.org/	50		
Mosquito insecticide resistance	adult mouse anatomy.obo			1		



Computational prediction of cancer-gene function Pingzhao Hu, Gary Bader, Dennis A. Wigle and Andrew Emili Nature Reviews Cancer 7, 23-34 (January 2007)

.ე9

Creating reference datasets for Systems Biology applications using text mining

- Manually annotated data repositories: incomplete, fraction of knowledge in literature
- Text mining: to extract, organize and present information for topic of interest
- Enable topic-centric literature navigation
- Assist in construction of manually revised data repositories
- Prioritization of biological entities for experimental characterization
- Facilitate human interpretation of large scale experiments by providing direct literature pointers
- Automatic retrieval of information relevant to human kinases.
- Linking kinase protein mentions to database records (i.e. sequences): protein mention normalization
- Extraction of Kinase mutations described in the literature
- Integration of information from full text articles, databases and genomic studies

Krallinger, M et al. Creating reference datasets for Systems Biology applications using text mining. 40 Ann N Y Acad Sci., (2009) 1158:14-28.

3rd International Biocuration Conference

April 16-19 Berlin, Germany

Meeting Schedule Registration

Abstract Submission Venue & Lodging

Sponsors

Contact

Workshop

Text Mining for the BioCuration Workflow

20 09

Organizers:

Lynette Hirschman, MITRE: lynette@mitre.org Gully APC Burns, ISI/USC: GullyBurns@gmail.com K. Bretonnel Cohen, University of Colorado: kevin.cohen@gmail.com Martin Krallinger, CNIO: mkrallinger@cnio.es Cathy Wu, Georgetown: wuc@georgetown.edu biocurator.org

The goals of this workshop are to update the BioCurator community on the state of the art in text mining and to elicit the requirements from the BioCurator community for enhanced tools to support the curation workflow.

The workshop will be divided into two parts. The first part will be tutorial in nature and will cover what tools are available, how to integrate components into a curation workflow, and what kind of performance to expect based on available resources. We will also discuss models for curation, including structured digital abstracts.

The second part of the workshop will be interactive, with a focus on understanding the diversity of curation workflows and requirements. For this part, we will invite participants to submit short presentations (5-10 min) on their requirements and their experiences or needs integrating text mining into their curation workflow. We will also discuss how to create partnerships between the bio-text mining tool developers and the BioCurator community.

BIOCURATION WORKFLOW TASKS



WORKFLOW TASKS AND TEXT MINING

- DEFINE & FORMALIZE INDIVIDUAL STEPS IN THE WORKFLOW
- DETECT WHICH STEPS CAN BE HANDLED THROUGH TEXT MINING ASSISTANCE
- PRIORITIZE MOST TIME CONSUMING STEPS
- FIND SUITABLE TEXT MINING APPROACH FOR EACH PARTICULAR TASK
- EVALUATE ANNOTATION EFFICICIENCY USING TEXT MINING ASSISTANCE
- USER FEEDBACK AND POTENTIAL ITERATIVE IMPROVEMENTS

ARTICLE IDENTIFICATION: TRIAGE TASK (1)



1



ARTICLE IDENTIFICATION: TRIAGE TASK (3)

- Traditionally addressed using keyword searches (e.g. Species names, interaction keywords, gene names, etc,..).
- Importance of triage task depends strongly on the annotation type and criteria used, organism source and literature volume.
- Potential text mining approaches for this task:
- More sophisticated keyword searches and Information retrieval (term weightings, Boolean queries, MeSH terms).
- Use of rules, regular expressions and pattern mining
- Document similarity (eTBLAST, vector space model)
- Machine learning and text categorization approaches (usually requires some sort of labeled text, e.g. PPI relevant articles) to learn which words

are useful to classify articles as relevant to the topic.

- For full text articles often retrieval is done at the level of text passages
- Sometime the triage task is combined with the bio-entity identification task
- Examples: BCMS, Genomics TREC, PreBIND,...

ANNOTATION EVENT IDENTIFICATION TASK

- Often consist in extraction of some kind of biological relation, e.g. Between two proteins (PPI), proteins and genes (TF and regulated genes),
- Between gene products and functional terms (GO, phenotypes) or between proteins and compounds.
- Often require the identification of some evidential text passages for the annotation event
- Is a very complex process, often domain export knowledge inference.
- Based on interpretation of author provided articles by curator
- Often requires mapping to controlled vocabulary terms and ontologies
- Text Mining approaches for this task:
- Automatic extraction of annotations, often based on sentence co-occurence assumption
- Article, passage, sentence classifiers
- Provide ranked collection of evidence passages
- Some approaches use patterns (trigger words), regular expressions or syntactic relations.

EVIDENTIAL QUALIFIER IDENTIFICATION TASK

- Evidential support for a given annotation important for interpretation.
- Indicative of the reliability of a given annotation and useful also for bioinformatics analysis
- Examples: GO evidence codes, PSI-MI interaction detection methods,

Oreganno evidence codes, ...

- Text mining approaches
- Either addressed as additional information for a given annotation event or through labeling the articles with evidence qualifiers
- Some NLP approaches more concerned with linguistic cues expressing uncertainty or negation
- Example: BioCreative II IMS task

PPI ANNOTATION OF BIOGRID



Pre-processing scientific articles

- Document Standardization: variety of formats (ASCII, HTML, XML, PDF, scanned PDF, SGML), convert them into a common format and encoding.
- XML /Extensible Markup language, standard way to insert tags onto a text to identify its parts
- OCR (Optical Character Recognition), used to digitalize older literature (PMC Back Issue Digitization initiative).
- Recover article Structure and content
- pdftotext, PDFLib, PDF Concerter
- Tokenization: break a stream of characters into words (tokens), e.g. white space, special chars.
- Each token is an instance of a type
- Stemming and lemmatization: standardize word tokens (e.g. Morphological analysis and
- Inflectional stemming, convert words to their corresponding root form)
- Lexical analysis of the text with the objective of treating digits, hyphens, punctuation marks, and the case of letters
- Elimination of stop-words
- Selection of index terms

Basic characteristics: exploring textual data



IMB

Considerations of Journal-specific characteristics:

- Journal/article Format (for pre-processing)
- Paper structure (section types)
- Article type (review, clinical study, etc.)
- Target audience of journal/article.



Processing levels of natural language texts



Krallinger M, et al. Analysis of biological processes and diseases using text mining approaches. *Methods Mol Biol.* (2009), to appear

Basic characteristics: biomedical literature

□ Heavy use of domain specific terminology (12%) biochemistry related technical terms^{*}), examples: chemoattractant, fibroblasts, angiogenesis Polysemic words (word sense disambiguation), examples: APC: (1) Argon Plasma Coagulation (2) Activated Protein C; or **teashirt:** (1) a type of cloth (2) a gene name (tsh). □ Heavy use of acronyms, examples: Activated protein C (APC), or vascular endothelial growth factor (VEGF) □ Most words with low frequency (data sparseness)

Word morphology and gene symbols



Krallinger M, et al. **Analysis of biological processes and diseases using text mining approaches**. *Methods Mol Biol*. (2009), to appear

Basic characteristics: biomedical literature

□ New names and terms created (novelty), example:

'This disorder maps to chromosome 7q11-21, and this locus was named *CLAM*. '[PMID:12771259]

□ Typographical variants (e.g. in writing gene names), example: TNF-alpha and TNF alpha (without hyphen)

□ Different writing styles (native languages): syntactic and semantic and word usage implications.

□ Heavy use of referring expressions (anaphora, cataphora and ellipsis) and inference, example:

Glycogenin is a glycosyltransferase.
 It functions as the autocatalytic initiator for the synthesis of glycogen in eukaryotic organisms.

Variability in Biomedical language

Country	Adjectives	Nouns	Verbs	Adverbs	Example sentence	PMID ref
Spain	Infrequent, bibliographic	Repercussion, evolution, existence, sunflower, olive, wine	-	Basically	Prevalence of CYP2D6 gene duplication and its <u>repercussion</u> on the oxidative phenotype in a white population.	7697944
Japan	Useful	Bullfrog, shadow (in radiography)	Clarify	Faintly, next, suddenly, scarcely	MDR-1 protein was <u>faintly</u> expressed in one of four chemoresistant patients, but Bcl-2 were [sic] clearly detected in four patients.	12538495
UK	Unsuitable, unlinked, unfamiliar	Marmoset, consultant, questionnaire	Lie, mirror, arise, tackle	Wholly, principally, particularly	The morphology of these projection neurons was revealed in great detail and confirmed that the projection <u>arises wholly</u> from pyramidal cells.	11602231
Russia	Gravitational	(Space) mission, quantum, hibernate, peculiarity, regularity, realization	-	Thermo- dynamically	The article is devoted to the question of peculiarity of bronchopulmonary system's pathology in the workers of the animal fodder production [sic].	10341521
India	Malarial, -wise (as in stepwise), ascorbic	Malaria, buffalo, peanut, garlic, catfish,	Impart (convey)	Appreciable the agglomera	Hydroxypropylmethylcellulose (HPMC) was used to <u>impart</u> strength and sphericity to tes.	12476867
France	Exceptional, digestive	Trouble	Envisage (imagine)	Successively (sequentially), essentially, sometimes	These 2 cells [sic] lines being able to clone, it is hard to <u>envisage</u> clonogenic assays.	3051563
China	Medicinal, radiant (heat), noxious (heat)	Acupuncture, coal, tea	Burn, replenish, alleviate	Obviously, meanwhile	Because only a catalytic amount of ERK2/pTpY is required, this method <u>alleviates</u> the need for large quantities of phospho-ERK2.	12056917
Germany	Satisfying practicable, unremarkable	Hint, precondition multitude	-	Additionally, exactly,	In clinically presumed spontaneous spinal cord infarction and <u>unremarkable</u> signaling of the spinal cord during sequential MRI investigations vertebral body infarction may serve as the only confirmatory sign of spinal cord ischemic stroke.	11987007
US	Federal, investigational, supplemental Residency, cocaine, payment, veteran, reimbursement, physician, care, plan, noncompliance, effort, profit Sponsor, mandate		Sponsor, mandate	-	Loss of revenue, mainly from noncompliance with charge capture resulted in the hospital billing only U\$\$386,794.32 with a total reimbursement of U\$\$165,779.86.	12488156

Literature repositories for life sciences

□ NLP: need electronically accessible texts.

Main scientific textual data types: e-books and earticles and the Web (online reports, etc).

□ e-Books: NCBI bookshelf.

□ Biomedical article citations (abstracts): PubMed

- □ Full text articles: PubMed Central (PMC)
- □ Repositories such as HighWire Press, BioMed Central

□ AGRICOLA, BIOSIS, Conference proceedings,...

PubMed database



- □ Scientific articles: new scientific discoveries.
- Citation entries of scientific articles of all biomedical sciences, nursing, biochemistry, engineering, chemistry, environmental sciences, psychology, etc,...
- Developed at the NCBI (NIH).
- Digital library contains more than 16 million citations
- □ From over 4,800 biomedical journals
- □ Most articles (over 12 million) articles in English.
- Each entry is characterized by a unique identifier, the PubMed identifier: PMID.
- □ More than half of them (over 7,000,000) have abstracts
- □ Often links to the full text articles are displayed.

PubMed database



- Approx. one million entries (with abstracts) refer to gene descriptions.
- □ Author, journal and title information of the publication.
- Some records with gene symbols and molecular sequence databank numbers
- □ Indexed with Medical Subject Headings (MeSH)
- Accessed online through a text-based search query system called Entrez
- Offers additional programming utilities, the Entrez Programming Utilities (eUtils)
- NLM also leases the content of the PubMed/ Medline database on a yearly basis

PubMed growth



Krallinger M, et al. Analysis of biological processes and diseases using text mining approaches. Methods Mol Biol. (2009), to appear

PubMed is accumulating over 600,000 new entries every year

60

Arabidopsis articles in PubMed



PubMed XML record

```
<PubmedArticle>
<MedlineCitation Status="Publisher" Owner="NLM">
   <PMID>18642075</PMID>
   <DateCreated>
       <Year>2008</Year>
       <Month>7</Month>
       <Dav>21</Dav>
   </DateCreated>
   <Article PubModel="Print-Electronic">
       <Journal>
          <ISSN IssnType="Print">0167-6806</ISSN>
          <JournalIssue CitedMedium="Print">
             <PubDate>
                 <Year>2008</Year>
                 <Month>Jul</Month>
                 <Dav>19</Dav>
              </PubDate>
          </JournalIssue>
          <Title>Breast cancer research and treatment</Title>
          <ISOAbbreviation>Breast Cancer Res. Treat.</ISOAbbreviation>
       </Journal>
       <ArticleTitle>Promoter methylation patterns of ATM, ATR, BRCA1, BRCA2 and P53 as putative cancer
       risk modifiers in Jewish BRCA1/BRCA2 mutation carriers.</ArticleTitle>
       <Pagination>
          <MedlinePgn/>
       </Pagination>
       <Abstract>
          <<u>AbstractText>BRCA1/BRCA2</u> germline mutations substantially increase breast and ovarian cancer
          risk, yet penetrance is incomplete. We hypothesized that germline epigenetic gene silencing may
          affect mutant BRCA1/2 penetrance. To test this notion, we determined the methylation status,
          using methylation-specific quantitative PCR of the promoter in putative modifier genes: BRCA1,
          BRCA2, ATM, ATR and P53 in Jewish BRCA1/BRCA2 mutation carriers with (n = 41) or without (n =
          48) breast cancer, in sporadic breast cancer (n = 52), and healthy controls (n = 89). Promoter
          hypermethylation was detected only in the BRCA1 promotor in 5.6-7.3% in each of the four
          subsets of participants, regardless of health and BRCA1/2 status.Germline promoter
          hypermethylation in the BRCA1 gene can be detected in about 5% of the female Israeli Jewish
          population, regardless of the BRCA1/2 status. The significance of this observation is yet to be
          determined. </AbstractText>
       </Abstract>
       <Affiliation>The Susanne Levy Gertner Oncogenetics Unit, The Danek Gertner Institute of Human
       Genetics, The Chaim Sheba medical Center, Tel-Hashomer, 52621, Israel. </ Affiliation>
       <AuthorList>
          <Author>
              <LastName>Kontorovich</LastName>
              <FirstName>Tair</FirstName>
             <Initials>T</Initials>
          </Author>
          <Author>
             <LastName>Cohen</LastName>
              <FirstName>Yoram</FirstName>
```

Biomedical corpora and text collections

- Medtag corpus, includes the Abgene, MedPost and GENETAG corpora
- Trec Genomics Track collections
- BioCreative corpus
- GENIA corpus
- Yapex corpus
- •Others, e.g. LL05 dataset, BioText Data, PennBioIE, OHSUMED text collection, Medstract corpus,...

Features for Natural Language Processing

- Techniques that analyze, understand and generate language (free text, speech).
- Multidisciplinary field: information technology, computational linguistics, AI, statistics, psychology, language studies, etc,.
- Strongly language dependent.
- Create computational models of language.
- Learn statistical properties of language.
- Methods: statistical analysis, machine learning, rule-based, pattern-matching, AI, etc...
- Explore the grammatical, morphological, syntactical and semantic features of well-structured language
- The statistical analysis of these features in large text collections is generally the basic approach used by NLP techniques.

Word	Base Form	Part-Of-Speech	Chunk	Named Entity
HAX-1	HAX-1	NN	B-NP	B-protein
associates	associate	VBZ	B-VP	0
with	with	IN	B-PP	0
cortactin	cortactin	NN	B-NP	B-protein
in	in	IN	B-PP	0
the	the	DT	B-NP	0
apical	apical	ງງ	I-NP	0
membrane	membrane	NN	I-NP	0
of	of	IN	B-PP	0
hepatocytes	hepatocyte	NNS	B-NP	B-cell_type
			0	0
Word	Morphology	Grammar	Syntax	Semantics

Krallinger M, et al Linking genes to literature: text mining, information extraction, and retrieval applications for biology. Genome Biol. 2008;9 Suppl 2:S8

Grammatical features

- Grammar: rules governing a particular language.
- Rules for correct formulation of a specific language
- Grammatical features in NLP, e.g. part of speech (POS)
- POS of a word depends on sentence context
- Examples: noun, verb, adjective, adverb or preposition.
- Programs label words with POS: POS taggers.
- Example:

Caspase-3 Proper noun, sing. *was* Verb, past tense *partially* Adverb *activated* Verb, past part. *by* Prep. or subord. Conjunction *IFN-gamma* Proper noun, sing. [PMID 12700631].

- POS taggers are usually based on machine learning
- Trained with a set of manually POS-tagged sentences.
- POS useful for gene name identification and protein interactions
- detection from text,
- MedPost {Smith, 2004} a POS for biomedical domain
- MedPost: 97% accuracy in PubMed abstracts (86.8% gen. POS tagger)

GENIA Tagger

GENIA GENIA GENIA GENIA GENIA Tagger Demo
(using <u>GENIA tagger</u> version 3.0)
Please enter a text that you want to analyze.
We show that Skbl interacts with a region of the N-terminal regulatory domain of Shkl distinct from that to which Cdc42 binds, and that Shkl, Cdc42, and Skbl are able to form a ternary complex in vivo.
Submit Text Restablecer

GENIA POS Tagger output

Chunking (shallow parsing)

(chunk types: ADJP, ADVP, CONJP, INTJ, LST, NP, PP, PRT, SBAR, VP)

[NP We] [VP show] [SBAR that] [NP Skb1] [VP interacts] [PP with] [NP a region] [PP of] [NP the N-terminal regulatory domain] [PP of] [NP Shk1] [ADJP distinct] [PP from] [NP that] [PP to] [NP which] [NP Cdc42] [VP binds] , and [SBAR that] [NP Shk1] , [NP Cdc42] , and [NP Skb1] [VP are] [ADJP able] [VP to form] [NP a ternary complex] [ADVP in vivo] .

Named entity recognition

(entity types: protein, DNA, RNA, cell_line, cell_type)

We show that Skb1 interacts with a region of the N-terminal regulatory domain of Shk1 distinct from that to which Cdc42 binds , and that Shk1 , Cdc42 , and Skb1 are able to form a ternary complex in vivo .

http://text0.mib.man.ac.uk/software/geniatagger/index.html

Morphological features

- Word structure analysis
- Rules of how words relate to each other.
- Example 1: plural formation rules, e.g.:

gene and genes or caspase and caspases

- Example 2: verb inflection rules, e.g. phosphorylate, phosphorylates and phosphorylating all have the same verb stem, word root.
- Stemmer algorithms to standardize word forms to a common stem
- Linking different words to the same entity.
- Different algorithms, e.g. Porter stemmer {Porter, 1980}
- Problem: collapse two semantically different words, e.g. gallery and gall.

Stemmer example results



💌 🏢 Porter's Stemming Algorit) 💶 🛋 🗙 File Edit View Go Bookmarks Tools Window 9 2 Search Home Bookmarks & Yahoo & Google **Porter's Stemming** Results Original Word Stemmed Word glycogenin glycogenin selfqlycosyl selfglycosylating protein protein primer primer initiates initi glycogen glycogen granule granul formation format examine examin role role protein protein during dure glycogen glycogen resynthesis resynthesi eight eight

http://maya.cs.depaul.edu/~classes/ds575/porter.htm

Syntactic features

- Relationships between words in a sentence: syntactic structure
- Shallow parsers analyze such relations at a coarse level, identification of phrases (groups of words which function as a syntactic unit).

• Example: Connexor shallow parser output: *Caspase-3* <: nominal head, noun, single-word noun phrase,> *was,* <auxiliary verb, indicative past> *partially* <adverbial head, adverb> *activated*<main verb, past participle, perfect> *by* <preposed marker, preposition> *IFN*- <premodifier, noun, noun phrase begins,> *gamma* <nominal head, noun, noun phrase ends>.

- Word labeled to corresponding phrase.
- Noun phrases (head is a noun, NP) e.g. 'Caspase-3' and 'INF-gamma' and verbal phrases (head is a verb, VP).

Protein interaction & Syntactic features



Krallinger M, et al. Analysis of biological processes and diseases using text mining approaches. *Methods Mol***7**1 *Biol*. (2009), to appear

Semantic features

- Associations of words with their corresponding meaning in a given context.
- Semantics (meanings) of a word -> understand meaning sentence.
- Dictionaries and thesauri provide such associations
- Gene Ontology (GO) provides concepts for biological aspects of genes
- Gene names and symbols contained in SwissProt (symbol dict.)
- Example: Caspase-3 /GENE PRODUCT was partially activated /INTERACTION VERB by *IFN-gamma* /GENE PRODUCT.
- Caspase-3 and INF-gamma are identified as gene products
- The verb 'activated' refers in this context to a certain type of interaction

PMID	Regulator	Regulated	Туре	Regulation association evidence sentence	0	b	X		Q WikiG		Source
<u>16055635</u>	E2FB <u>AT5G22220</u>	CDKB1;1 <u>AT3G54180</u>	Activation	Indeed, our recent data show that E2FB can directly induce the promoter of the Arabidopsis CDKB1;1 gene (Z. Magyar, unpublished results).	0	None		Q ihop	Q WikiG	TAIR	
<u>16055635</u>	E2FA <u>At2g36010</u>	CDKB1;1 <u>AT3G54180</u>	Activation	It is not clear how E2FA could promote the expression of CDKB1;1 , which has a separate expression window in S- and M-phases (Magyar et al., 1997, 2000).	0	None	X	्र ihop	Q WikiG		
NLP Tasks

- Information Retrieval (IR)
- Text clustering
- Text classification
- Information extraction (IE)
- Question Answering (QA)
- Automatic summarization

Main task types which have been addressed by Bio-NLP systems

Additional task types

- Natural Language Generation
- Anaphora resolution
- Text zoning
- Machine translation
- Text proofing
- Speech recognition

Information Retrieval (IR) and Search Engines

- IR: process of recovery of those documents from a collection of documents which satisfy a given information demand.
- Information demand often posed in form of a search query.
- Example: retrieval of web-pages using search engines, e.g.
 Google.
- Important steps for indexing document collection:
 - Tokenization
 - Case folding
 - Stemming
 - Stop word removal
- Efficient indexing to reduce vocabulary of terms and query formulations.
- Example: 'Glycogenin AND binding' and 'glycogenin AND bind'.
- Query types: Boolean query and Vector Space Model based query.

VECTOR SPACE MODEL

- Measure similarity between query and documents.
- (1) Document indexing,
- (2) Term weighting,
- (3) Similarity coefficient
- Query: a list of terms or even whole documents.
- Query as vectors of terms.
- Term weighting (w) according to their frequency:
 within the document (i) & within the document collection (d)

w: term weight

tf: term frequency

idf: inverted document frequency

- Widespread term weighting: tf x idf.
- Calculate similarity between those vectors
- Cosine similarity often used.
- Return a ranked list.
- Example: related article search in PubMed

ctors.

$$sim(Q, D) = \frac{\sum_{j=1}^{V} w_{Q,j} \times w_{i,j}}{\sqrt{\sum_{i=1}^{V} w_{Q,j} \times \sum_{i=1}^{V} w_{i,i}^2}}$$

sim(Q,D): similarity between query 75 and document

 $idf_{i,j} = \log\left(\frac{N}{df_j}\right)$

 $w_{i,j} = tf_{i,j} \times idf_j$

eTBLAST



eTBLAST results: high scoring words

🕹 PMID: 8529663 - Mozilla Firefox	_ 7 🗙
<u>A</u> rchivo <u>E</u> ditar <u>V</u> er <u>I</u> r <u>M</u> arcadores Herramien <u>t</u> as Ay <u>u</u> da	0
փ 🔹 🚽 🗧 🚱 🚱 👔 🗋 http://invention.swmed.edu/cgi-bin/etblast/abstract_local?pmid=8529663&use 💟 🔇 Ir 💽	
BLAST	
	61111111
Eur J Biochem 1995 Nov;234(1);343-9.	
Glycogen metabolism in quail embryo muscle. The role of the glycogenin primer and the intermediate proglycogen.	
J Lomako W M Lamako	
W J Whelan Terms with hig	jh weight
Department of Biochemistry and Molecular Biology, University of Miami School of Medicine, FL 33101, USA.	
Cultured quail embryo muscle has proven to be an excellent model system for studying the synthesis of macromolecular glycogen fro	om.
and its degradation to, glycogenin, the autocatalytic, self-glucosylating primer for glycogen synthesis. We recently demonstrated that	1
proglycogen, a low-M(r) form of glycogen, is an intermediate in the synthesis. Here we show that proglycogen also functions as an	
intermediate in macroglycogen degradation and, in one set of circumstances, represents an arrest point in glycogen breakdown, whi	ich
between prodycogen and macromolecular dycogen and are not normally depredent turnover of giveogen in tissues, the molecules cycle	open.
the released glycogenin is active, capable of re-initiating glycogen synthesis. Under culture conditions where the conversion of	
proglycogen into glycogenin does take place, the intermediates lying between form a discrete rather than a continuous series, sugge	stive
of a cluster structure for proglycogen and indicating that breakdown is stepwise. Evidence of post-translational modification of glyco	igenin
was obtained by the finding that, in glycogen from cultured muscle, glycogenin is phosphorylated.	
MedlinelD: 0	
PMID: 8529663	

Text clustering

•Find which documents have many words in common, and place the documents with the most words in common into the same groups.

•Similarity of documents instead of similarity of sequences, expression profiles or structures

•Cluster documents into topics, for instance: clinical, biochemical and microbiology articles

•A clustering program tries to find the groups in the data.

•Clustering programs often choose first the documents that seem representative of the middle of each of the clusters (candidate centers of the clusters).

•Then it compares all the documents to these initial representatives.

•Each documents is assigned to the cluster it is most similar to.

•Similarity is based on how many words the documents have in common, and how strongly they are weighted.

•The topical terms of the clusters are chosen from words that represent the center of the cluster.

•The best clustering is one in which the average difference of the documents to their cluster centers smallest.

•Agglomerative clustering: first comparing every pair of documents, and finding the pair of documents which are most similar to each other.

Clustering documents, genes, terms



Krallinger M, et al. Analysis of biological processes and diseases using text mining approaches. *Methods Mol***7**9 *Biol*. (2009), to appear

Text classification

•Common problem in information science.

•Assignment of an electronic document to one or more categories, based on its contents (words).

•Can be divided into two sorts: <u>supervised</u> document classification where some external mechanism (such as human feedback) provides information on the correct classification for documents, and unsupervised document classification.

• Document classification techniques include:

- * naive Bayes classifier
- * tf-idf
- * latent semantic indexing
- * support vector machines
- * artificial neural network
- * kNN
- * decision trees, such as ID3
- * Concept Mining
- Classification techniques have been applied to spam filtering
- Cane use the bow toolkit, SVMlight, LibSVM etc,...

Text classification & supervised learning





Cell cycle protein ranking

AT5G11300 92.501942618 0.872659836019 108.181946101 106 Image: Construction of the second se	
TAIR db gene identifier AT3G48750 42.268552948 0.862623529551 50.830108648 49 Image: Constraint of the second s	
identifier AT5G51330 11.97953993 0.855681423571 17.12329801 14 Image: Constrained and the constrained and	ore
AT5G20850 11.948959364 0.702879962588 15.18793846 17 Image: Constraint of the second	d lots
AT5G05490 11.57761703 1.92960283833 12.3889299 6 Image: Comparison of the second seco	
	ction ces
ADSTRACT A14G37490 8.65024547 0.786385951818 11.52331461 11	
AT4G21270 6.44660105 1.6116502625 6.99035612 4 C & C & C & C & C & C & C & C & C & C	tion
AT2G31970 6.20142349 1.03357058167 6.20142349 6 S	
3G54180 5.59239631 0.932066051667 6.43610856 6 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2	ord Co- ence
ATIG08560 5.18091117 0.370065083571 7.49510317 14	
Diamonds E0 150 5.1173723 2.55868615 5.1173723 2 2 Experi	iment
AT3G24810 5.0762994 1.6920998 5.0762994 3 S Keywo	

Protein abstract associations

The following results were retrieved for your query: AT3G48750

[PubMed-ID] 7523194

[TAIRID] AT3G48750

[NAMES] cdc2 # cdk2

[TITLE] Olomoucine, an inhibitor of the cdc2/cdk2 kinases activity, blocks plant cells at the G1 to S and G2 to M cell cycle transitions.

[ABSTRACT] The cdc2/cdk2 protein kinases play key roles in the cell cycle at two control points: the G1/S transition and the entry into mitosis. Olomoucine, a specific inhibitor of these kinases, was tested in two plant cell systems: Petunia mesophyll protoplasts induced to divide and Arabidopsis thaliana cell suspension cultures. The cell cycle status was analysed from DNA histograms or through continuous labelling of cells with 5-bromodeoxyuridine (BrdUrd) followed by double staining with bis-benzimide (Hoechst 33258) and propidium iodide (PI). Such analyses resolve cells from several generations according to the extent of their DNA replication. Olomoucine was shown to reversibly arrest differentiated Petunia cells induced to divide at G1 phase and cycling Arabidopsis cells in late G1 and G2. A comparison of the effects of aphidicolin, oryzalin and olomoucine suggests that in the Arabidopsis cell suspension culture, a cdc2/cdk2-like kinase is activated at a restriction point in late G1.

Cell cycle terms: G1 phase; mitosis; cell cycle

Species ambiguity scores: 0.666667

Cell cycle scores: 3.76766

[PubMed-ID] 9428718

[TAIRID] AT3G48750

[NAMES] CDC2 # p34cdc2

[TITLE] Plant CDC2 is not only targeted to the pre-prophase band, but also co-localizes with the spindle, phragmoplast, and chromosomes.

[ABSTRACT] A polyclonal antiserum against the p34cdc2 homologue of Arabidopsis thaliana, CDC2aAt, was used in parallel with a polyclonal antiserum against the PSTAIRE motif to study the subcellular localization of CDC2 during the cell cycle of isolated root tip cells of Medicago sativa. During interphase, CDC2 was located in the nucleus and in the cytoplasm. The cytoplasmic localization persisted during the complete cell cycle, whereas the nuclear signal disappeared at nuclear envelope breakdown. At the beginning of anaphase, the anti-CDC2aAt antibody transiently co-localized with condensed chromosomes. The chromosomal co-localization disappeared as anaphase continued and remained excluded from the separated chromosomes until cytokinesis, when CDC2 re-located to the newly forming nuclei. We also observed a co-localization of CDC2 with three microtubular structures, the pre-prophase band, the spindle, and the phragmoplast.

Cell cycle terms: <u>cell cycle;</u> interphase; prophase; anaphase

Species ambiguity scores: 0.5

Cell cycle scores: 3.4199

Query flower development

Search Arabidopsis bibliome Clear

Searching the Arabidopsis

literature: abstracts (1)

[PubMed-ID] 15923332

[TITLE] The AtRAD51C gene is required for normal meiotic chromosome synapsis and double-stranded break repair in Arabidopsis.

[ABSTRACT] Meiotic prophase I is a complex process involving homologous chromosome (homolog) pairing, synapsis, and recombination. The budding yeast (Saccharomyces cerevisiae) RAD51 gene is known to be important for recombination and DNA repair in the mitotic cell cycle. In addition, RAD51 is required for meiosis and its Arabidopsis (Arabidopsis thaliana) ortholog is important for normal neiotic DNA meiotic homolog pairing, synapsis, and repair of double-stranded breaks. In vertebrate cell cultures, the RAD51 paralog RAD51C is also important for mitotic homologous recombination and ucial for the maintenance of genome integrity. However, the function of RAD51C in meiosis is not well understood. vith a meiotic Here we describe the identification and analysis of a mutation in the Arabidopsis RAD51C ortholog, AtRAD51C. Although the atrad51c-1 mutant has normal vegetative and flower development and to plants and has no detectable abnormality in mitosis, it is completely male and female sterile. During early vegetative he. meiosis, homologous chromosomes in atrad51c-1 fail to undergo synapsis and become severely imented fragmented. In addition, analysis of the atrad51c-1 atspo11-1 double mutant showed that eiosis are that basic us. fragmentation was nearly completely suppressed by the atspoll-1 mutation, indicating that the bidopsis These results ripts fragmentation largely represents a defect in processing double-stranded breaks generated by roughout repair during AtSPO11-1. Fluorescence in situ hybridization experiments suggest that homolog juxtaposition might G1-phase tMnd1 also be abnormal in atrad51c-1 meiocytes. These results demonstrate that AtRAD51C is essential for ility to ants also late normal meiosis and is probably required for homologous synapsis. SB repair in B1:1. Cell cycle terms: synapsis; mitosis; meiosis I; cell cycle; prophase; meiosis; meiotic r that cells prophase I; mitotic cell cycle ologs of the male cdc2 ant Species ambiguity scores: 0.6 ral Cell cycle scores: 3.04013 istem Cell cycle scores: 1.16984 Cell cycle scores: 1.06576 Cell cycle scores: 1.05989





Experimentally characterized in more detail.

For each protein links to their associated articles as well as keywords, e.g. Spindle and Cell Cycle GO terms and experimental keywords are provided



Information Extraction

- Identification of semantic structures within free text.
- Use of syntactic and Part of Speech (POS) information.
- Integration of domain specific knowledge (e.g. ontologies).
- Identification of textual patterns.
- Extraction of predefined entities (NER), relations, facts.
- Entities like: companies, places or proteins, drugs.
- Relations like: protein interactions
- Methods: heuristics, rule-based systems, machine learning and statistical techniques, regular expressions,.



TAGGING BIO-ENTITIES IN TEXT

- Aim: <u>Identify</u> biological entities in articles and to <u>link</u> them to entries in biological databases.
- Generic NER: corporate names and places (0.9 f-score),
 Message Understanding Conferences (MUC).
- Biology NER: more complex (synonyms, disambiguation, typographical variants, official symbols not used,..).
- Bioinformatics vs. NLP approach.
- Performance organism dependent.
- Methods: POS tagging, rule-based, flexible matching, statistics, ML (naïve Bayes, ME, SVM, CRF, HMM).
- Important for down-stream text mining.

SOME TRICKY CASES OF GENE TAGGING

(1) The nightcap mutation caused severe defects in these cells [PMID:12399306].

(2) In the present investigation, we have discovered that Piccolo, a CAZ (cytoskeletal matrix associated with the active zone) protein in neurons that is structurally related to Rim2, [PMID:12401793]

(3) The Drosophila takeout gene is regulated by the somatic sex-determination pathway and affects male courtship behavior. [PMID:12435630]

(4) This function is independent of Chico, the Drosophila insulin receptor substrate (IRS) homolog [PMID:12702880].

(5) A new longevity gene, Indy (for I'm not dead yet), which doubles the average [PMID:12391301]

(6) The Drosophila peanut gene is required for cytokinesis and encodes a protein similar to yeast putative bud neck filament proteins [PMID 8181057].

(7) Ambiguity of PKC: Protein kinase C and Pollution kerato-conjunctivitis

Choose Classes Virus **Tissue** RNA Protein Polynucleotide Peptide OtherOrganicCompound ✓ OtherName ✓ OtherArtificialSource 🔽 Ørganism ✓ Nucleotide MultiCell MonoCell - Lipid **Inorganic** DNA CellType CellLine CellComponent Carbohydrate BodyPart - Atom AminoAcidMonomer Select All Deselect All

Biomedical Named Entity Recognizer

Analysis of murine Brca2 reveals conservation of protein-protein interactions but differences in nuclear localization signals. In this report, we have analyzed the protein encoded by the murine Brca2 locus. We find that murine Brca2 shares multiple properties with human BRCA2 including its regulation during the cell cycle, localization to nuclear foci, and interaction with Brca1 and Rad51. Murine Brca2 stably interacts with human BRCA1, and the amino terminus of Brca2 is sufficient for this interaction. Exon 11 of murine Brca2 is required for its stable association with RAD51, whereas the carboxyl terminus of Brca2 is dispensable for this interaction. Finally, in contrast to human BRCA2, we demonstrate that carboxyl-terminal truncations of murine Brca2 localize to the nucleus. This finding may explain the apparent inconsistency between the cytoplasmic localization of carboxyl-terminal truncations of human BRCA2 and the hypomorphic phenotype of mice homozygous for similar carboxyl-terminal truncating mutatio

- Based on Machine learning
- Good results in the COLING Bio-NER contest (Geneva)

• Many classes (entity types), including Virus, Tissue, RNA, Protein, Polynucleotide, Peptide, Organism, Nucleotide, Lipid, DNA, Cell Type, Cell Line, Cell Component, Carbohydrate, Body Part Atom and Amino Acid Monomer

Analysis of murine **Brca2** reveals conservation of **protein protein interactions** but differences in **nuclear localization** signals. In this report , we have analyzed the protein encoded by the murine **Brca2** locus. We find that murine **Brca2** shares multiple properties with human **BRCA2** including its regulation during the **cell cycle**, localization to nuclear foci, and interaction with **Brca1** and Rad51. Murine **Brca2** stably interacts with human **BRCA1**, and the **amino terminus** of **Brca2** is sufficient for this interaction. Exon 11 of murine **Brca2** is required for its stable association with **RAD51**, whereas the **carboxyl terminus** of **Brca2** is dispensable for this interaction. Finally , in contrast to human **BRCA2**, we demonstrate that **carboxyl terminal truncations** of murine **Brca2** localization of carboxyl terminal truncations of human **BRCA2** and the hypomorphic phenotype of **mice** homozygous for similar carboxyl terminal truncating mutatio

PLAN2L: a web tool for integrated text mining & literature-based bioentity relation extraction



http://zope.bioinfo.cnio.es/plan2l

Krallinger, M. et al . **PLAN2L: a web tool for integrated text mining and literature-based bioentity relation extraction**. To appear *in Nucl. Acids Res.*, Web Server Issue, 2009.

PLAN2L

PMID	Regulator	Regulated	Туре	Regulation association evidence sentence	0	b	X		Q WikiG	Source
<u>16055635</u>	E2FB <u>AT5G22220</u>	CDKB1;1 <u>AT3G54180</u>	Activation	Indeed, our recent data show that E2FB can directly induce the promoter of the Arabidopsis CDKB1;1 gene (Z. Magyar, unpublished results).	0	None	X		Q WikiG	
<u>16055635</u>	E2FA <u>At2g36010</u>	CDKB1;1 <u>AT3G54180</u>	Activation	It is not clear how E2FA could promote the expression of CDKB1;1 , which has a separate expression window in S- and M-phases (Magyar et al., 1997, 2000).	0	None	X	Q ihop	Q WikiG	

PMID	Proteins	Location description evidence sentence	Location terms	Location words	X	
<u>11058164</u>	LEUNIG	The nuclear localization of LEUNIG -GFP is consistent with a role of LEUNIG as a transcriptional regulator.	nuclear	localization		Q ihop

PMID	PLAN2L Protein association evidence	٩		Ŷ	.•	*	\gg	1		X
<u>11041883</u>	Additional analyses indicate that the absence of marginal tissues in leunig aintegumenta double mutants is not mediated by ectopic AGAMOUS .	-0.384139	-1.26458	-0.121035	-0.430636	2.57862	-0.0047048	-1.25884	-0.937231	X
<u>16625397</u>	SEUSS and LEUNIG encode components of a putative transcriptional regulatory complex that controls organ identity specification through the repression of the floral organ identity gene AGAMOUS .	-1.5747	1.60754	-0.190263	-0.104091	3.33405	-0.294006	-1.00271	-0.170982	X
<u>11782418</u>	The effects of seuss mutations are most striking when combined with mutations in LEUNIG , a previously identified repressor of AGAMOUS .	-0.885257	-0.754272	-0.239062	-1.4516	1.33617	-0.847629	-1.31295	-0.751337	X
<u>11058164</u>	LEUNIG , a putative transcriptional corepressor that regulates AGAMOUS expression during flower development.	-1.3353	3.11992	-0.338721	-0.641414	5.46814	-0.307351	-0.943231	-0.865053	
<u>16854969</u>	Previously, we identified and isolated two Arabidopsis transcription co-repressors LEUNIG (LUG) and SEUSS (SEU) that function together in a putative co-repressor complex to prevent ectopic AGAMOUS (AG) transcription in flowers.	-0.726197	2.60201	-0.39487	0.303502	1.99957	-0.326139	-1.42551	-1.02715	X

http://zope.bioinfo.cnio.es/plan2l

PLAN2L flowchart



PLAN2L protein mention normalization



97

PLAN2L mutation extraction



iHOP system



99

iHOP system: query to DB record



iHOP system: Defining information

	Find in this P	age
Sentences in this view contain definitions for BRCA2 - Definitions are available whenever you see this symbol []] - Read more. For a summary overview of the information in this page click here. new	Show all Order by releva	_ ance _
PALB2, which encodes a BRCA2 Interacting protein, is a <u>breast cancer</u> susceptibility gene	e.[2007]	📓 🛨
Inheritance of one defective BRCA2 🎓 allele predisposes humans to <u>breast cancer</u> . [2001]		📓 土
A common variant in BRCA2 🎲 is associated with both <u>breast cancer</u> risk and prenatal viab	ility. [2000]	📑 🛨
Inherited mutations in the gene BRCA2 🎓 predispose carriers to early onset <u>breast cancer</u> , account for fewer than 2% of all cases in East Anglia. [2000]	but such mutations	≝ <u></u>
Mutations in BRCA2 😭 are thought to account for as much as 35% of all inherited <u>breast ca</u> proportion of inherited <u>ovarian cancer</u> . [1996]	<mark>ncer</mark> as wall as a	≝
Two of the five BRCA2 😭 mutation carriers reported a family history of <u>breast_cancer</u> , and n history of <u>ovarian_cancer</u> . [2002]	one reported a family	₫ ±
Our results indicate that BRCA2 🎲 confers a very high risk of <u>breast_cancer</u> and is responsil fraction of breast and <u>ovarian_cancer</u> in Iceland, but only a small proportion of other cancers	ble for a substantial s. [1996]	≝ ±
Recent studies have identified mutations in the breast and (<u>ovarian cancer</u> susceptibility ge which has been found in the germline of several males and one female affected with <u>breast (</u>	ne 2 (BRCA2 🏠), one <mark>cancer</mark> . [1996]	<u>∎</u> ±
The <u>breast_cancer</u> susceptibility gene BRCA2 🎓 on <u>chromosome</u> 13q12-13 has recently b	een identified. [1997]	📑 🛨
The <u>breast cancer</u> susceptibility gene, BRCA2 ừ on <u>chromosome</u> 13q12-13, was recently	isolated. [1996]	📑 🛨
The BRCA2 🎡 gene on <u>chromosome</u> 13 has been shown to be associated with familial male <u>cancer</u> . [1996]	and female <u>breast</u>	

Main gene Colour Associated genes legend Relevant Biomedical terms Compounds

Defining Information for this Gene



iHOP system: interaction information

Sentences in this view contain interactions of BRCA2 - Interaction Information is available whenever you see this symbol 2 - Read more. Show all Show all For a summary overview of the information in this page click here. Order by relevance Order by relevance
RESULTS: Definite BRCA2 and mutations were found in 2 of the 73 women with early-onset <u>breast cancer</u> (2.7 percent; 95 percent <u>confidence interval</u> , 0.4 to 9.6 percent), suggesting that BRCA2 is associated with fewer cases than <u>BRCA1</u> (P=0.03). [1997]
Age penetrance is greater for BRCA1 🏾 -linked than for BRCA2 🎓 -linked cancers in this population. [2000] 👘 🗮 👗
Tumors lacking BRCA1 & mRNA were more likely to lack BRCA2 & mRNA than tumors expressing BRCA1 & 📔 🛃 MRNA (P<.001). [2002]
We evaluate current knowledge of BRCA1 😭 and BRCA2 🎓 functions to explain why mutations in BRCA1 😭 and 🧱 去 👗 BRCA2 🎓 lead specifically to breast and ovarian cancer. [2001]
PURPOSE: Morphologic and immunohistochemical studies of familial breast cancers have identified specific characteristics associated with BRCA1 & mutation-associated tumors when compared with BRCA2 & and non-BRCA1/2 tumors, but have not identified differences between BRCA2 & and non-BRCA1/2 tumors. [2005]
What you don't know can hurt you: adverse psychologic effects in members of BRCA1 🔄 -linked and BRCA2 🎓 🛛 🧱 📥 -linked families who decline genetic testing. [1998]
Here we report the chromosomal gains and losses as measured by CGH in 25 BRCA2 &-associated <u>breast</u> [] 🛃 👗 👗 <u>tumors</u> and compared them with our existing 36 <u>BRCA1</u> & and 30 control profiles. [2005]
Germline mutations of BRCA1 & are also associated with ovarian cancer and mutations of BRCA2 & are associated with an increased risk of male breast cancer, ovarian cancer, prostate cancer and pancreatic cancer. [1997]
As these studies concerned sporadic cancer cases, we investigated whether N372H and another common variant 📗 ± 👗 located in the 5'-untranslated region (203G > A) of the BRCA2 🎓 gene modify breast or ovarian cancer risk in BRCA1 🎓 mutation carriers. [2005]
The identification of molecules that interact with Brcal 🖈 and Brca2 🖈 has greatly enhanced our knowledge of 🛛 🧱 📥 how BRCA1 🛊 and BRCA2 🛊 may function as tumor suppressors. [1998]
BRCA1 🎓 mutations are more commonly associated with ovarian cancer than BRCA2 😭 mutations. [2001] 🛛 🧱 📥



iHOP system: recent information

	"
recent information is available whenever you see this symbol 7 - Read more.	-
For a summary overview of the information in this page click here. new Order by relevance	a 🔽
Mutations in the BRCA2 🏫 interacting DSS1 🏫 are not a <u>risk factor</u> for <u>male breast cancer</u> . [2007]	1 🛨
Constitutive activation of MAPK [?] @/ERK [?] @ inhibits prostate cancer cell proliferation through upregulation of BRCA2 @. [2007]	1 📥
BRCA2 🎓 is central to an utterly diverse biological behavior elicited after <u>integrin</u> -mediated normal and <u>prostate</u> cancer cell adhesion to basement membrane (BM) and extracellular matrix (ECM) proteins. [2007]	1 📥
We investigated ERK [?] 🎓 and AKT phosphorylation in normal (PNT1A) and cancer (PC-3) prostate cells after adhesion to ECM and the effects upon BRCA2 🏫 and cell proliferation. [2007]	1 🛨
PNT1A <u>cell_adhesion</u> to <u>ECM</u> triggered <u>MAPK [?] @/ERK [?]</u> @ signaling resulting in <u>upregulation</u> of BRCA2 @ [[mRNA and protein, with negligible effects upon <u>cell proliferation</u> . [2007]	1 📥
The BRCA2 🎓 mutation c.3531-3534delCAGC (3758del4) is novel and the BRCA1 🎓 mutation c.1840A>T (K614X) is 📑 reported for the first time in Cypriot patients. [2007]	1 🛨
METHODS: 277 families with pathogenic BRCA1 @/BRCA2 @ mutations were reviewed and 28 breast cancer phenocopies identified. [2007]	1 🛨
FINDINGS: Questionnaires were completed by 799 women with a history of invasive <u>ovarian cancer</u> (670 with BRCA1 mutations, 128 with BRCA2 mutations, and one with a mutation in both genes), and controls were 2424 women without <u>ovarian cancer</u> (2043 with <u>BRCA1</u> are mutations, 380 with <u>BRCA2</u> mutations, and one with a mutation in both genes). [2007]	<u>+</u>
Contribution of BRCA1 @ and BRCA2 @ germline mutations to the incidence of early-onset breast cancer in Cyprus. [2007]	1 🛨
The <u>Fanconi anemia</u> and <u>BRCA</u> anetworks are considered interconnected, as BRCA2 and gene defects have been a discovered in individuals with <u>Fanconi anemia</u> subtype D1. [2007]	1 📥
In particular, the genetic testing is limited in its ability to determine which of the many missense mutations identified in BRCA1 & and BRCA2 & actually predispose to cancer and which are simply neutral alterations. [2007]	1 🛨
METHODS: We did a <u>matched case-control study</u> in women who were found to carry a pathogenetic mutation in BRCA1 & or BRCA2 . [2007]	1 🛨



iHOP system: gene model/ graph

Symbol Name		Synonyms	Organism
BRCA2 breast cancer 2, e	early onset	BRCC2, Breast cancer type 2 susceptibility protein, FACD, FAD, FAD1, FANCB, FANCD, FANCD1, Fanconi anemia group D1 protein	Homo sapiens
UniProt P51587, Q5TB Q8IU82 IntAct P51587 PDB Structure 1N0W OMIM 114480, 15522 NCBI Gene 675	J7, 55 more than 1,500 organi	sms. 80,000 genes. 12 million se	ntences.
NCBI RefSeq NP_000050 NCBI RefSeq NM_000059	г	aiways up-to-	uate.
NCBI UniGene 675 NCBI Accession CAA98995,		Gene Model - the logbook	
AAQ97181 Homologues of BRCA2		In the course of your navigation through added to the <i>Gene Model</i> by clicking o	i iHOP, interesting sentences can be n the 🏨 icon beside the sentence.
Interaction information for	BRCA2 🛐	The Gene Model stores these sentence graph. More about the Gene Model	es and represents their relation in a
Most recent information for	or BRCA2 🔯 new	graphi Hore about the Cene Houeini	
Enhanced PubMed/Google WARNING: Please keep in mind that ger confidence value 😪 😭	e query ne detection is done automatically and can exhit	e.g.	
Gene model is a	interactive	(PTENB)	
graph where you	can add		
interesting senter	nces and	(PTPMT)	CDBA
interactions.			



iHOP system: confidence



The synomnym ambiguity limitation

Many gene or protein synonyms are ambiguous, thus one and the same synonym is often used for different genes. Even human experts can have difficulties to resolve such ambiguities and automatic systems, like iHOP, will therefore always exhibit certain errors.

The iHOP confidence value

Although no definite solution for the problem of synonym ambiguity is in sight, it is possible to put an automatically derived confidence value to specific gene references.

This iHOP confidence value is illustrated through the colour intensity of a star



The absence of a star does not mean that a certain term could not be a gene, but simply that supporting evidence is not available.

EBIMed



You can explore a total of 1846 permutations for this HitPair table arrangement. Click on the secondary columns' headers to rearrange the table.
 Rows 1 to 5 (out of 199).

first << 1/40 >> last								
Protein/Gene	Protein/Gene	<u>Cellular</u> <u>component</u>	<u>Biological</u> <u>process</u>	<u>Molecular</u> <u>function</u>	Drug	<u>Species</u>		
BRCA2 or	BRCA1 (<u>244/966</u>)	chromosome (<u>40/61</u>)	DNA repair (<u>45/52</u>)	binding (<u>17/24</u>)	gel (<u>15/16</u>)	cancer (<u>391/1088</u>)		
FANCD1	RAD51 (<u>26/59</u>)	chromatin (<u>8/14</u>)	development (<u>29/35</u>)	DNA-binding (<u>6/8</u>)	spectrum (<u>13/18</u>)	human <i>or</i>		
(score: 1603)	PCR (<u>21/21</u>)	nucleus (<u>8/9</u>)	localization (<u>15/20</u>)	E2 (<u>2/2</u>)	via (<u>9/12</u>)	man (<u>71/111</u>)		
	brca2 (<u>19/22</u>)	replication forks (<u>4/5</u>)	cell cycle (<u>14/20</u>)	CDK (<u>1/1</u>)	trigger <i>or</i>	mouse (<u>15/19</u>)		
	Rad51 (<u>18/41</u>)	endoplasmic	transcription (<u>14/17</u>)		labels (<u>6/6</u>)	anemia (<u>13/22</u>)		
	a protein (<u>13/13</u>)	reticulum or	pathogenesis (<u>12/12</u>)		mitomycin (<u>6/6</u>)	codons (<u>10/11</u>)		
	recombinase <i>or</i> recombinases (12/15)	midbody (<u>2/2</u>)	double-strand break repair (9/10)		lines (<u>5/9</u>) For women (5/5)	mice (<u>8/10</u>) veast (6/7)		
	FANCD2 (8/20)	intracellular (<u>2/2</u>)	cell proliferation (8/20)		Adriamycin <i>or</i>	chicken (5/12)		
	p53 (<u>8/11</u>)	Golgi vesicles (<u>2/2</u>)	S-phase (7/9)		doxorubicin (<u>3/7</u>)	MCF (5/7)		
	estrogen receptor or	extracellular matrix (<u>2/2</u>)	RNA interference or		estrogen (<u>3/6</u>) Jumen or	murine (<u>5/6</u>)		
	green fluorescent	centrosome (<u>1/3</u>)	recombinational		luminal (<u>3/4</u>)	Caenorhabditis		
	protein <i>or</i>	collagen type I (<u>1/2</u>)	repair (<u>7/7</u>)		del (<u>2/6</u>)	mammals (3/3)		
	GFP (<u>5/5</u>)	buds (<u>1/1</u>)	phosphorylation (<u>6/10</u>)		tamoxifen (<u>2/3</u>)	beta (2/3)		
	DSS1 (<u>4/14</u>)	spindle (<u>1/1</u>)	DNA		maps (<u>2/2</u>)	Castilla (2/3)		
	RB1 (<u>4/13</u>)	plasma	recombination (<u>6/9</u>)		eleven (<u>2/2</u>)	aa (2/3)		
	FANCG or	membrane (<u>1/1</u>)	DNA replication (<u>5/6</u>)		etoposide (<u>2/2</u>)	thymus (2/3)		
	XRCC9 (<u>4/12</u>)	microtubules (<u>1/1</u>)	death (<u>4/6</u>)		vincristine (<u>2/2</u>)	dogs or		
	MBC (<u>4/8</u>)	cytoplasm (<u>1/1</u>)	behavior or		docetaxel (<u>1/2</u>)	Canis canis (<u>2/2</u>)		
	Brcal (<u>4/4</u>)	nuclear matrix (<u>1/1</u>)	benaviour (<u>4/4</u>)		prenatal (<u>1/1</u>)	helix (<u>2/2</u>)		
		micronucleus (<u>1/1</u>)	cytokinesis (<u>3/6</u>)		cisplatin (<u>1/1</u>)	Arabidopsis		
	PALB2 (3/10)	basement	meiosis (<u>3/6</u>)		Ets (<u>1/1</u>)	thaliana (<u>2/2</u>)		
	PARP (3/9)	membrane (<u>1/1</u>)	M phases or M phase (3/5)		compounds (<u>1/1</u>)	Chinese		
FANCC (<u>3/7</u>)		nucleoplasm (<u>1/1</u>)	Cell cycle control (3/4)		Inc (<u>1/1</u>)	Hamster (<u>Z/Z</u>)		
			cell division (3/3)		mutagen <i>or</i>	ostilago mavdis <i>or</i>		
	receptor (<u>3/3</u>)		pregnancy or		nitrogen mustard (<u>1/1</u>)	U. maydis (<u>1/3</u>)		
	ЕСМ (<u>2/8</u>)		gestation (<u>3/3</u>)		retinoic acid (1/1)	rat or		

GOPubMed


BioCreative

∳ ▼ <u>⇒</u> ▼	୧ 😣		😼 http	://ww	w.bio	creativ	e.org/													N v	• G • Go	ogle	Q. 3
PubMed Home	LiMTox	https	://webr	nail.cni	i o	CNIO	BioC	CreAtIvE	Yaho	o Sp	pindled	ome	planetom	ix W	S_	Rad	UniProt	Com	pendium	of Text	SNP2L: SNP	to Literat	
🔀 Genome Bi	iology Ful	text	Lin (3 0)	Те	xt Minir	ng lectu	re	\otimes		Bio	Creative -	News	- La	test	0						
Bio were fu	hose <u>lo</u>	<u>gin re</u>	giste	<u>c</u>																		Searc	:h 🔝
shoving active (cells,	r i	. t j	L C i	a 1	I	s	s e s	s s n	ı e	n t	C	o f	I n	f	o 1	•						
to deterrine who	ether m	a t	ic	n	E	x t	ra	a c t	: i c	n	i	n	B i	o 1	o	g y	,						
particulari asso with ive res	scent	N	ews		1	Abo	ut	1	Eve	ents	;	1	Task	s	1	Re	sourc	es					
Year		Org	aniz	ers																			
2003																							
2004		Web	site ı	upgra	ades	5 [200	8-12-03	3]															
2005		The	e new	BioC	reati	vew	ebsite	just h	as rec	eive	d son	ne ex	tra funct	ionali	ity:								
2006			1. Log	gged i	n us	ers n	ow ca	n not o	only ch	nang	e the	pass	word, b	ut als	o th	eir e	mail ac	Idress					
2007		2	2. RS	S 2.0	Fee	ds fo	r all m	ayor s	ection	s (se	e hei	re) ar	d a favi	con fo	or th	e we	ebpage						
2008		;	3. Tea	am reg	gistra	ation/	mana	gemei	nt has	beer	n add	ed. T	his is no	ot the	offic	cial s	tart of	the reg	gistratio	on proce	ss, but you	can	
			alre	eady r	egis	ter a	team t	for Bio	Creati	ve II.	5 her	e if y	ou feel l	ike it	- or	use	the nev	w "tear	m page	e" link on	the top of	the page	
Content			in t	he us	er m	enu.																	
BioCreative I		4	4. Tor	ns of r	nino	r upd	ates a	and a b	ougfix	have	beer	n add	ed.										
BioCreative II		We	hope	, this	mak	es vis	siting	this sit	e an e	venl	better	expe	erience!										
BioCreative II.5																							
Organizers		Bio	Crea	ative	e II.	5																	
Publications		BioC	reati	ve II.	.5 A	nno	unce	ment	(Eve	nts)	[2008	-11-18	1										

Why community assessments?

- Compare different methods and strategies
- □ Reproduce performance of systems on common data
- Provide useful data collections: Gold Standard data
- □ Explore meaningful evaluation strategies and tools
- Determine the state of the art
- □ Monitor improvements in the field
- □ Point out needs of the user community
- Promote collaborative efforts





Community assessments



CASP: Critical assessment of Protein Structure Prediction CAMDA: Critical Assessment of Microarray Data Analysis CAPRI: Critical Assessment of Prediction of Interactions GASP: Genome Annotation Assessment Project GAW: Genome Access Workshop PTC: Predictive Toxicology Challenge

INLPBA: Joint workshop on Natural Language Processing in Biomedicine TREC: Text Retrieval conference MUC: Message Understanding conference LLL05: Genic interaction extraction challenge RTE: Textual Entailment challenge







Participants - Annotation Servers

- Alias I, New York, Bob Carpenter
- Georgetown University, Hongfang Liu
- Humboldt Univ., Berlin, Jörg Hakenberg
- Inst. of Biomed. Inf., Taiwan, Cheng-Ju Kuo
- Inst. of Inform. Sci., Taiwan, Richard Tsai
- Jena Univ., Germany, Kathrin Tomanek
- Milwaukee Marquette Univ., Craig Struble
- National Inst. of Health, William Lau
- Norweg. Univ. of Sci. and Tech., Janny Chen
- Seoul National University, Sun Kim
- Univ. of Colorado, William Baumgartner
- University of Edinburgh, Barry Haddow
- University of Geneva, Patrick Ruch
- University of Michigan, Arzucan Ozgur
- Univ. of Pennsylvania, Kuzman Ganchev
- Yale University, ThaiBinh Luong

Main advantages of BCMS

- Data Integration: multi-site annotations
- Simplicity of usage: single API with many annotations
- ✤ User-oriented: TM & biologist
- Novel/ unique: first system in biomedical text mining
- Scalability: additional systems
- Extensibility: additional annotation types
- Flexibility: additional input text types, e.g. full-text articles

GM Predictions								
Mention	#	Conf.						
Muc4	2	0.998						
ErbB2	2	0.994						
ASGP-2	2	0.963						
neu Ab1	2	0.924						
anti-ErbB2	2	0.863						
ErbB2	2	0.808						
sialomucin	1	0.704						
SMC	2	0.555						
Muc4/SMC	1	0.173						
sialomucin	1	-						
anti-phospho-Er	1	-						
Neomarkers	1	-						

GN Predictions

Normalization	#	Conf.
Mucin-4 precurs	1	1.000
Transmembrane p	1	1.000
Receptor tyrosi	1	1.000
S-layer protein	1	0.886
Chromosome part	1	0.884
Matrix protein	1	0.500

PPI Predictions

Home | XML-RPC

Differential localization of ErbB2 in different tissues of the rat female reproductive tract: implications for the use of specific antibodies for ErbB2 analysis.

ErbB2 has been implicated in numerous functions, including normal and aberrant development of a variety of tissues. Although no soluble ligand has been identified for ErbB2, we have recently shown that ASGP-2, the transmembrane subunit of the cell surface glycoprotein Muc4 (also called sialomucin complex, SMC), can act as an intramembrane ligand for ErbB2 and modulate its activity. Muc4/SMC is abundantly expressed at the apical surface of most epithelia of the rat female reproductive tract. Since Muc4/SMC can interact with ErbB2 when they are expressed in the same cell and membrane, we investigated whether these two proteins are co-expressed and co-localized in tissues of the female reproductive tract. Using an anti-ErbB2 antibody from Dako, we found moderate staining at the basolateral surface of the oviduct and also around the cell membrane of the most superficial and medial layers of the stratified epithelia of the vagina. In contrast, Neomarkers neu Ab1 antibody intensely stained the apical surface of the epithelium of the oviduct and the medial and basal layers of the stratified epithelia of the vagina, substantially overlapping the distribution of Muc4/SMC. Furthermore, Muc4/SMC and ErbB2 association in different tissues of the female reproductive tract was demonstrated by co-immunoprecipitation analysis. Interestingly, phosphorylated ErbB2 detected by anti-phospho-ErbB2 is primarily present at the apical surface of the oviduct. Thus, our results show that differentially localized forms of ErbB2 are recognized by different antibodies and raise interesting questions about the nature of the different forms of ErbB2, the mechanism for differential localization, and possible functions of ErbB2 in the female reproductive tract. They also raise a cautionary note about the use of different ErbB2 antibodies for expression and localization studies.

PubMed ID: 11598901

MEDLINE creation date: 2001-10-12

Acknowledgements



Prof. Alfonso Valencia & Structural Computational Biology group at CNIO. 116