

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 / Topics (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 SB₂

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

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*)



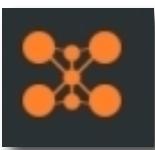
ELIXIR

EUROPEAN LIFE SCIENCES INFRASTRUCTURE FOR BIOLOGICAL INFORMATION

To construct and operate a sustainable infrastructure for biological information in Europe,

To support 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 define the scope of the infrastructure and how it will be constructed





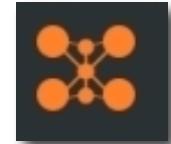
ELIXIR

EUROPEAN LIFE SCIENCES INFRASTRUCTURE FOR BIOLOGICAL INFORMATION

- **Optimal Data Management**
 - Coordinated Data Resources with improved access
 - 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 maintain open access to data
- Enhance European competitiveness in bioscience industries
- Address need for Increased Funding & its Coordination

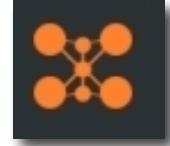


The Preparatory Phase project



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

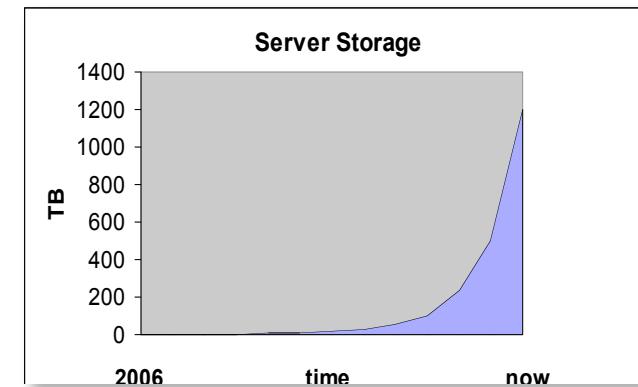
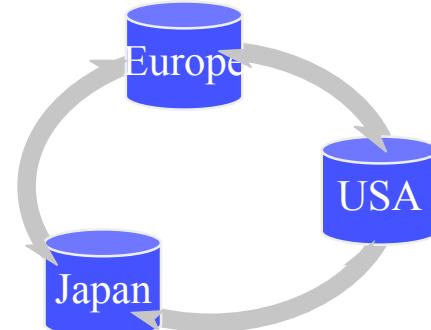
- 1. Project management
- 2. Data providers
- 3. User communities**
- 4. Organisation and Legal
- 5. Funding
- 6. Physical infrastructure
- 7. Data interoperability
- 8. Literature**
- 9. Healthcare
- 10. Chemistry, Plants, Agriculture & Environment**
- 11. Training**
- 12. Tools integration
- 13. Feasibility studies
- 14. Reporting and negotiation



Why do we need ELIXIR?

(Why do we need bioinformatics infrastructures)

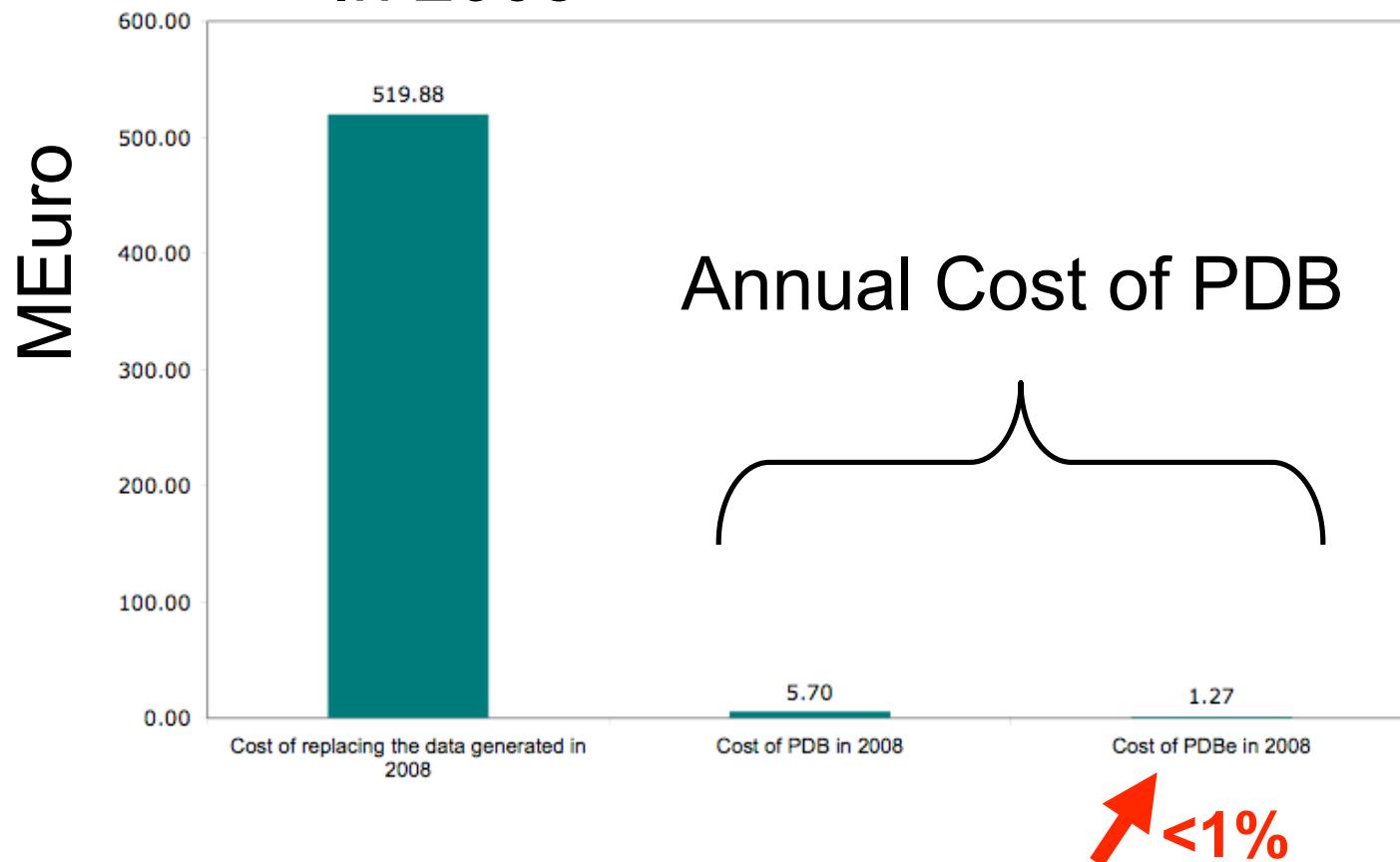
- Data Growth
- Global context
- Very large user community:
 - 3.3 m web hits/day
 - 20,000 unique users per day
- Need to preserve data and make accessible to all
- Impact on Medicine, Agriculture & Biotechnology
- Impact on society & bioindustries
- Need for increased funding for biodata resources

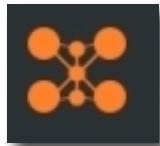




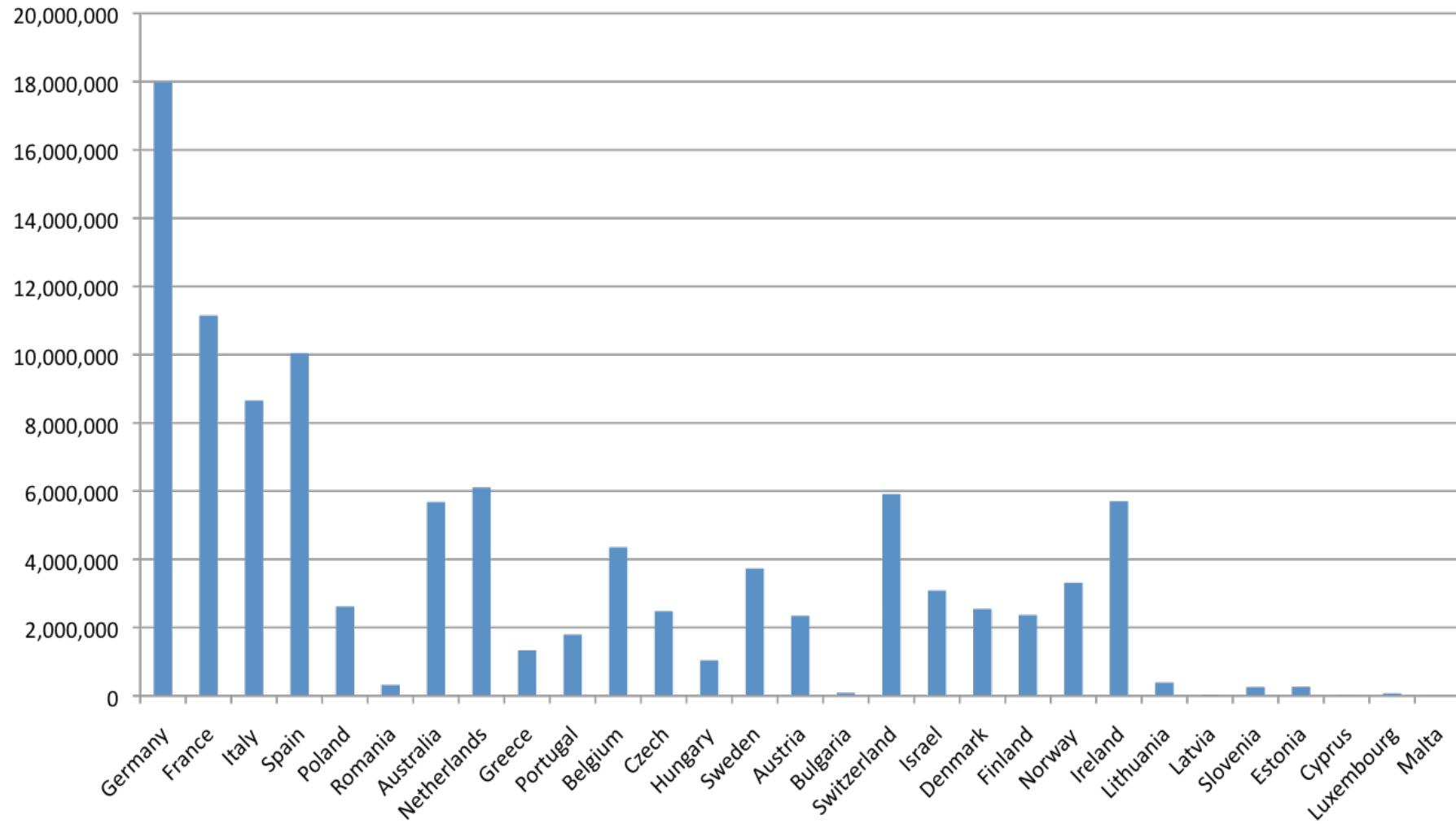
Good Value for Money e.g. PDB

Data collection
In 2008

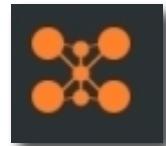




EBI Hits in 2008



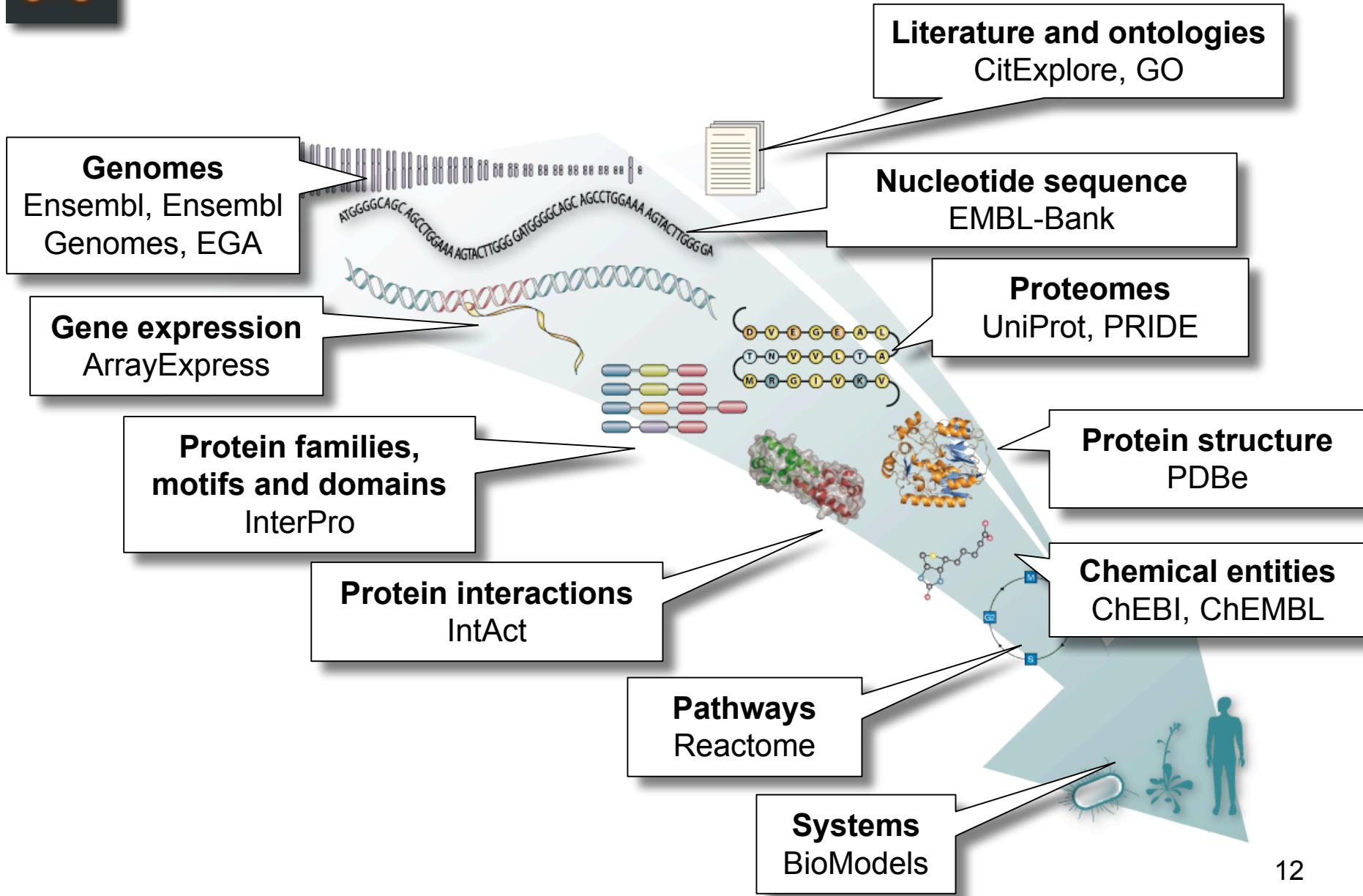
WP3: User Communities



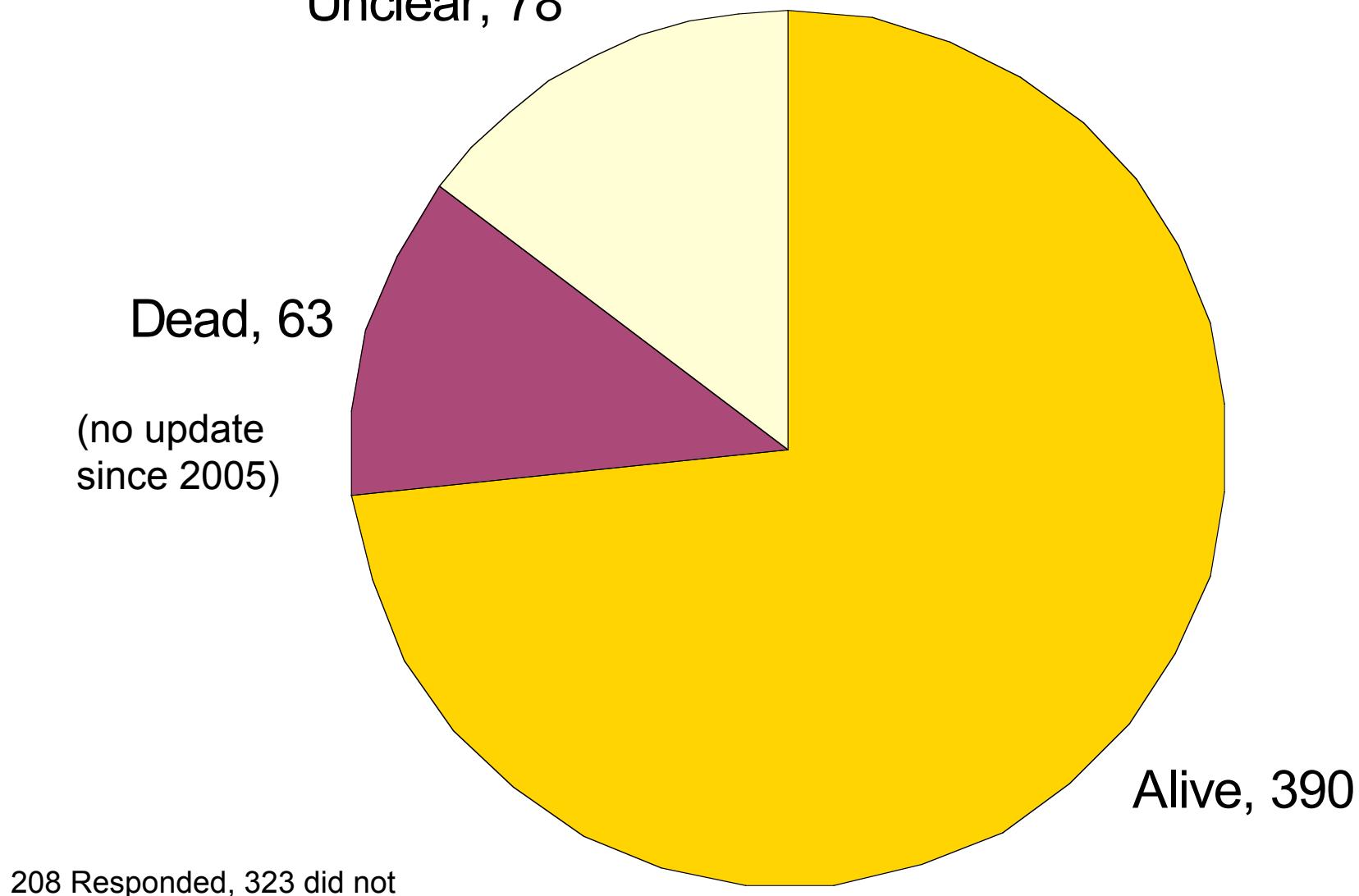
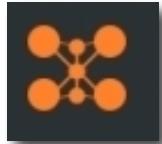
- User Survey: 1000 responses
 - Long term support essential
 - Top 3 challenges:
 - **Data integration; Format compatibility; Website usability**
 - Concerns
 - Data quality and measures; Quality of tools; Training
- Need to consider different needs in different countries
- Need for a plan for long-term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



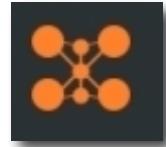
Databases: molecules to systems



531 Databases surveyed



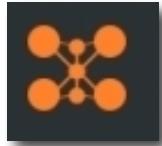
Total European effort



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

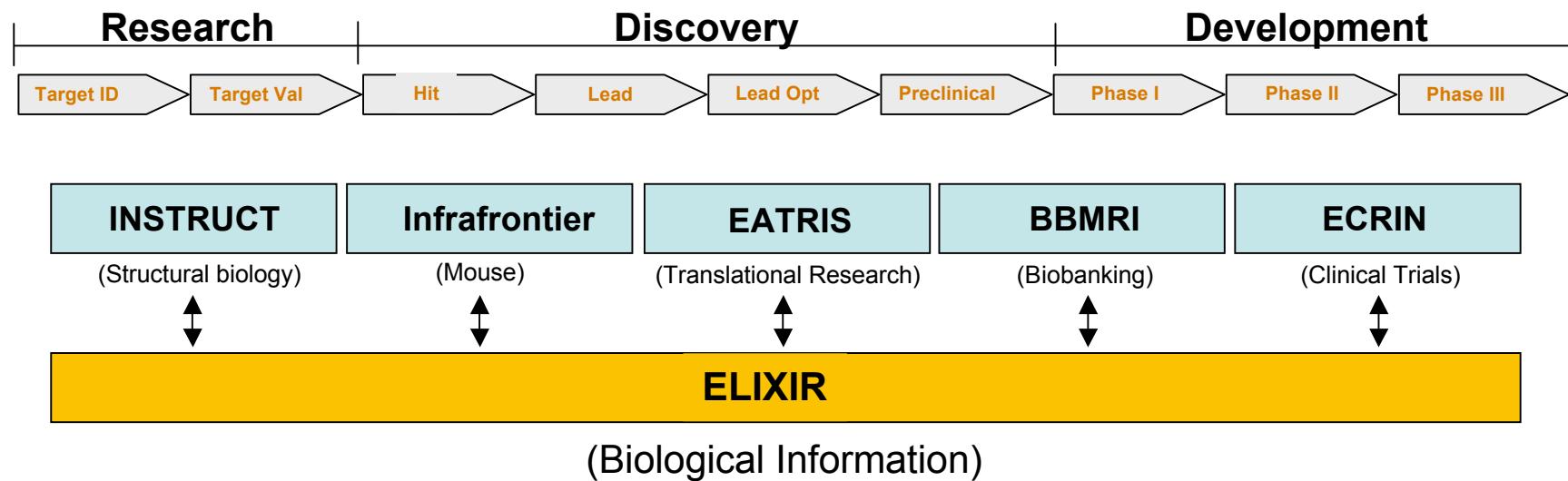
RECOMMENDATION

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



ESFRI

Biology Research Infrastructure proposals.

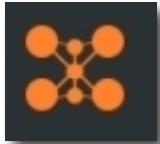




WP 10: Chemistry, Plants, Agriculture & Environment

- 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.

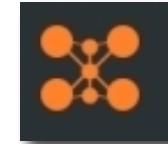
Elixir recommendations:

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, UCHSCColorado
- 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) status of literature repositories throughout Europe and

recommendations for the future

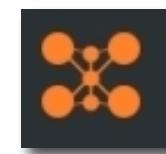
(2) infrastructure needs in Europe to establish an information-sharing platform to integrate databases and literature for (*) experts and non-experts, with

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

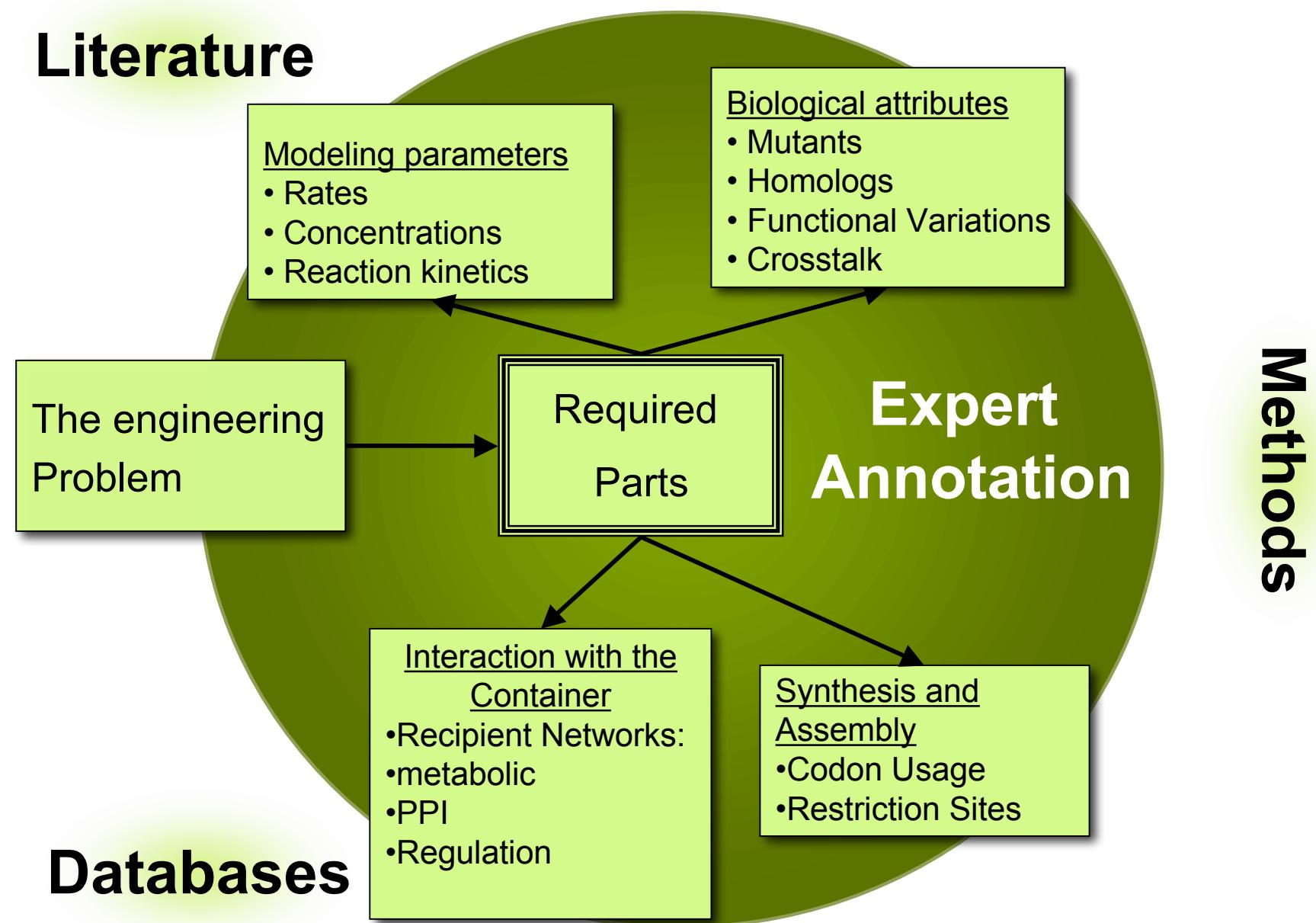
(3bi) tools for access to the literature, for

(3bii) data representation and for

(3biii) interaction with end users.



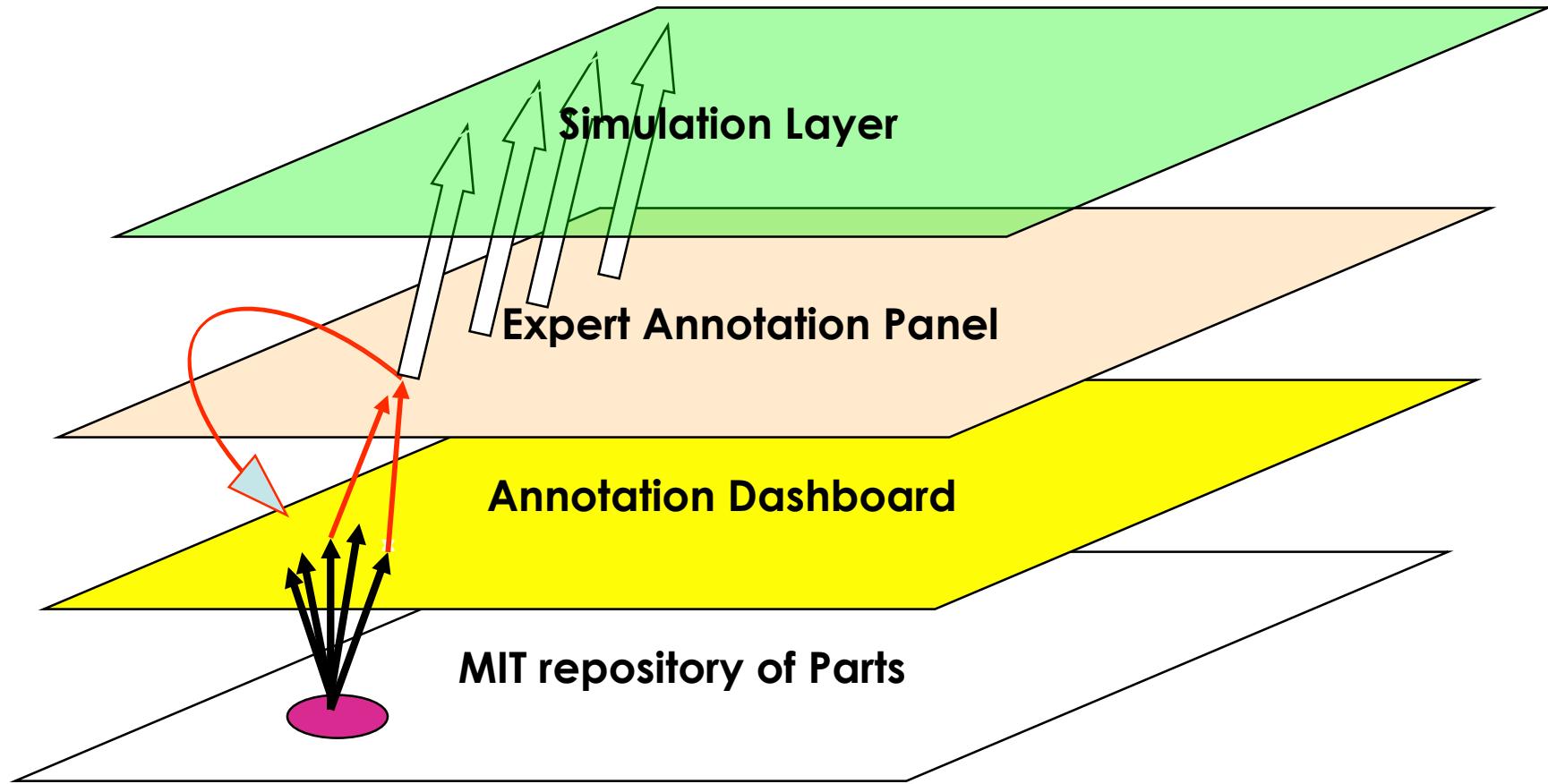
Literature



Databases

Adapted from I. Cases

Emergence IT layout



MIT Repository of Parts



jump to part
BBa_

navigation
Main Page
Browse Part Types

Protein Coding Regions

Available Protein Coding Regions

-?	Name	Protein	Description	Direction	Tag	Length
A W	BBa_J23012		SpecR open reading frame basic part	Forward		888
A W	BBa_J31000	Hin	Generates Hin from Salmonella typhimurium			573
A W	BBa_J31001		Hin invertase tagged with LVA (HinLVA)			612
A W	BBa_J31002		Kanamycin Resistance backwards (KanB)			816
A W	BBa_J31006		Tet Resistance Backwards (TetB)			1191
A W	BBa_J31007		Tet Resistance Forwards (TetF)			1191
A W	BBa_J45001		SAMT enzyme			1155
A W	BBa_J45002		BAMT (SAM:Benzoic Acid Carboxyl Methyltransferase)			1098
A W	BBa_J45004		BSMT: converts salicylic acid to methyl salicylate (wintergreen)			1074
A W	BBa_J45008		Branched-chain amino acid transaminase (BAT2); used in biosynthesis of banana scent			1134

Show 314 more parts [Edit](#)

Create an account or log in



article | discussion | edit | history

Designed by Maia Mahoney

Part:BBa_C0053

Repressor, P22 c2

The P22 c2 repressor protein coding sequence is a 720 base-pair sequence with the standard RBS-compatible BioBrick prefix and the standard BioBrick suffix sections on its ends. It binds to the P22 c2 regulatory sequence, BBa_R0053. The sequence contains a LVA tag for faster degradation.

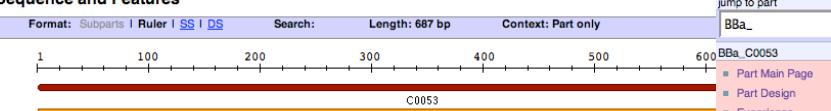
Usage and Biology

P22 c2 is a member of the lambdoid cl protein family.

The coding sequence has been modified from wild type c2 to contain an LVA tag, and two TAA stop codons.

Sequence and Features

Format: Subparts | Ruler | SS | DS Search: Length: 687 bp Context: Part only



jump to part
BBa_C0053

Part Main Page
Part Design
Experience
Hard Information
Physical DNA

navigation
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
Printable version



article | discussion

Designed by Maia Mahoney

Part:BBa_C0053:Physical DNA

Repressor, P22 c2

Physical DNA statistics

Plasmid Plasmid Length Part and Plasmid VF2 - VR

pSB1A2	2079	2766	925
BBa_J61031	10662	11349	None
BBa_I52001	1090	1777	None
BBa_J61003	3016	3703	1862
BBa_J23018	2977	3664	1823
BBa_J52017	4101	4788	1915
BBa_J61002	3002	3689	1848
BBa_J61009	5130	5817	None
BBa_J61007	2888	3575	None
BBa_I51001	2438	3125	2057
pSB1A3	2157	2844	1003
pSB1A7	2431	3118	1277
pSB1AC3	3055	3742	1003
pSB1AK3	3169	3876	1003
pSB1AT3	3446	4133	1003
pSB2K3	4425	5112	1003
pSB4A3	3339	4026	1003
BBa_I50040	2226	2913	None

Part Length: 687bp

ACTG Ratios
 $a: 0.323, c: 0.199, t: 0.231, g: 0.246$
 $260:280 \text{ ratio: } 2.378$

Existing versions
 Only the BioBrick Standard Assembly version exists

Existing sequence analyses
 No sequence analyses are available

curated / validated collection of artificial parts

[Start a new sequence analysis](#)

This part may be found in these wells/tubes [Show all locations](#)

Library	Well	Plate	Plasmid	Cell
iGEM 2006	5K	iGEM2006 DNA-1	pSB1A2	V1004
				Available

MaDAS principal features

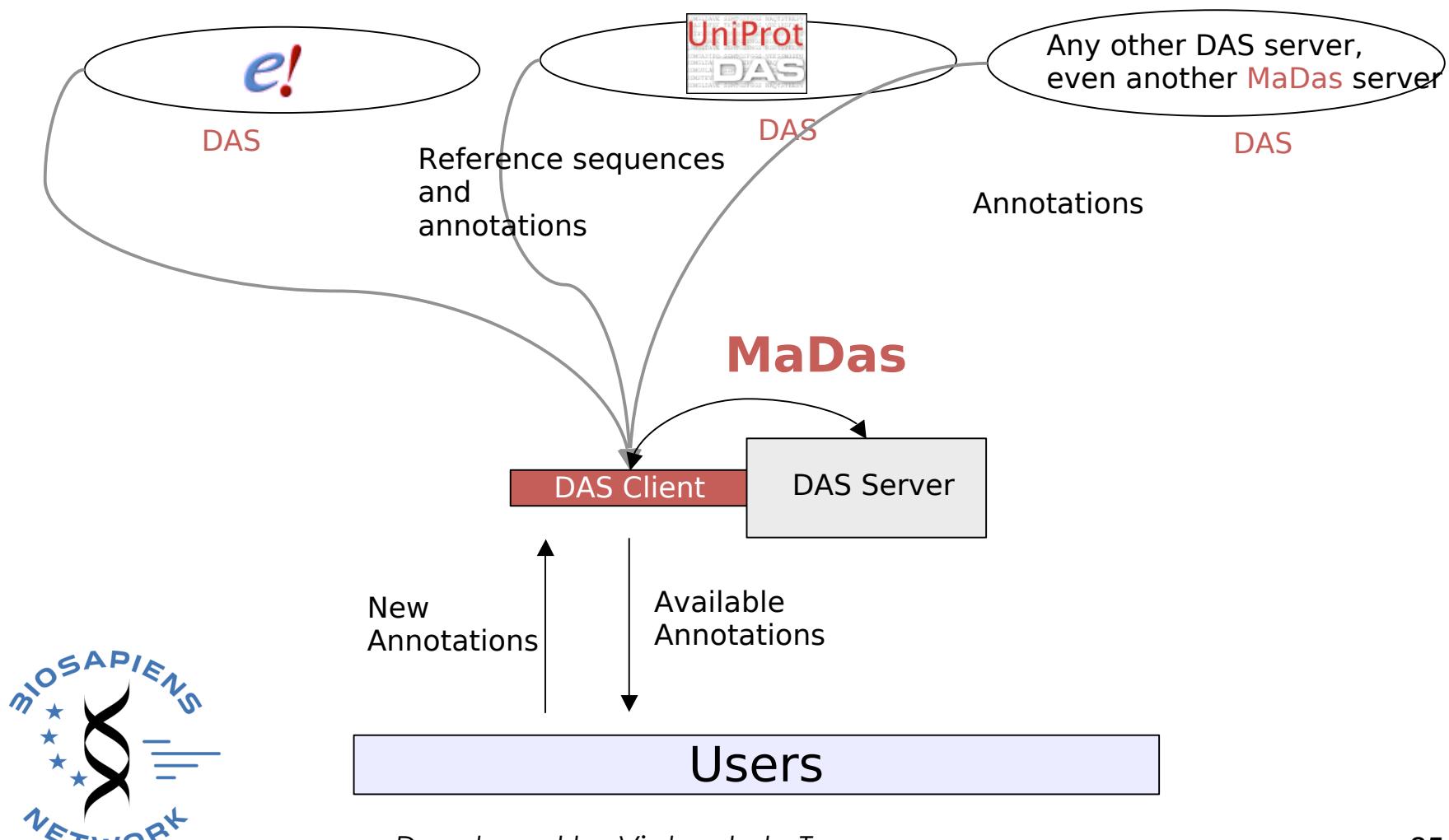
- ✓ MaDAS allows users to add, edit, or remove self generated sequence annotations
- ✓ Allows to upload multiple annotations from different sources.
- ✓ Provides a security system based on projects. The annotations could be public or only available for the project members.
- ✓ Provides an interface to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



Developed by Victor de la Torre

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 deactivate an entire set of annotations.

[HOME](#) [SIGN IN](#) [UPLOAD DATA](#) [DOWN](#)

MaDas

Version 0.5 (For testing purposes only)

Welcome Guest !!!

MaDas is a manual sequence annotation system.

The general purpose is to provide users with a simple or public repository. The system is based on the DAS modules: a DAS annotations server and a DAS and proteins. In addition, the FunCUT generator is displayed.

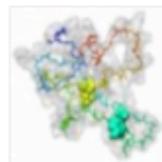
You need to register to use it. For more info

This is a testing version. Therefore comments
[here](#).

Sequences to annotate



Genomes



Proteins

The screenshot shows the MaDAS 2.0 Released interface. On the left, there's a sidebar with a diamond icon containing a computer monitor and the text "VISUALIZATION PLUGINS" and a link "[<< Back](#)". The main area features a large green header bar. Below it, a central box displays a genomic feature visualization with various colored tracks (red, blue, green) representing different genomic elements. A tooltip box is overlaid on the visualization, listing project details:

TYPE:	CDS
START:	86517
END:	89135
METHOD:	Genbank
SCORE:	
ORIENTATION:	+
PHASE:	

Below the visualization, the text "MaDAS 2.0 - DAS Feature" is visible. To the right, a sidebar titled "Now Working In" shows a "MaDAS Test" project with details: By: Victor De La Torre, At: cnio, Created: 2008-02-22, Category: test projects, Security: public. It also includes a "Description" section stating "Test project to show MaDAS 2.0 capabilities" and a "Delete Project" button. Further down, a "Project Members (1)" section lists "Osvaldo GraÃ±a". At the bottom right, there's a "Join Similar Projects" button.



MaDAS

2.0 Released

Welcome guest 

[Home](#) [Projects](#) [Plugins](#) [Help](#)

Now Working In

[your projects](#)

Protein Test

By: Victor De La Torre
At: cni
Created: 2008-02-23
Category: test projects
Security: public

Description:
test project to show proteins in MAJAX

Project Members (1)
guest

Join Similar Projects

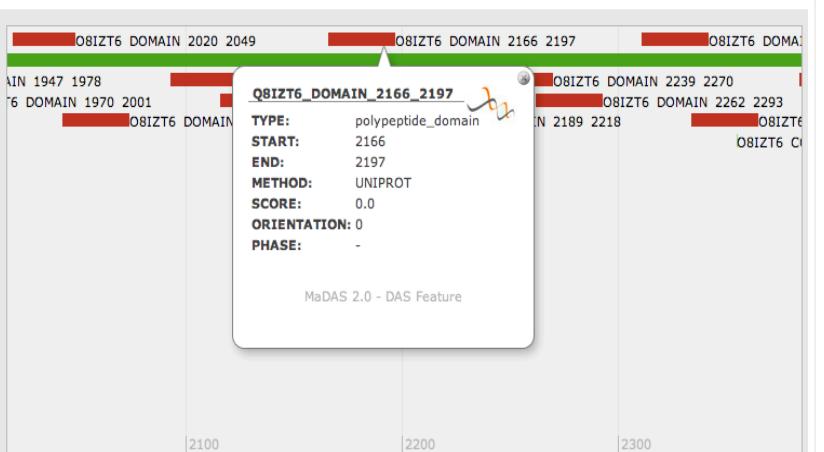
In category: test projects
By: Victor De La Torre
At: cni

[!\[\]\(662ea80a7c7ea5080da432a79ec6ddd4_img.jpg\) Back to your projects](#)

MAJAX
PLUGINS

<< Back

MaDAS 2.0 - DAS Feature



Introducing expert annotations and consolidating them in databases/visualization

The screenshot shows the MaDasM web application interface. On the left, there is a genome browser visualization of a chromosome with markers at 62 MB and 93 MB. A circular navigation button with arrows and a plus sign is positioned above the browser. Below the browser, the text "Created by Victor De La Torre" is visible. On the right, a modal window titled "Add/Edit Feature" is open. The form fields include:

- ID/Label *: SNP1
- Type *: MUTATED_VARIANT_SITE
- Start *: 24700567
- End *: 24700568
- Score: (empty)
- Orientation *: + (selected)
- Phase *: 0

The "Note" section contains the text: "The effect of the mutation have been described in (Sample et al. 2008)".

The "Link text" field contains "PubMi".

The "Link url" field contains "http://".

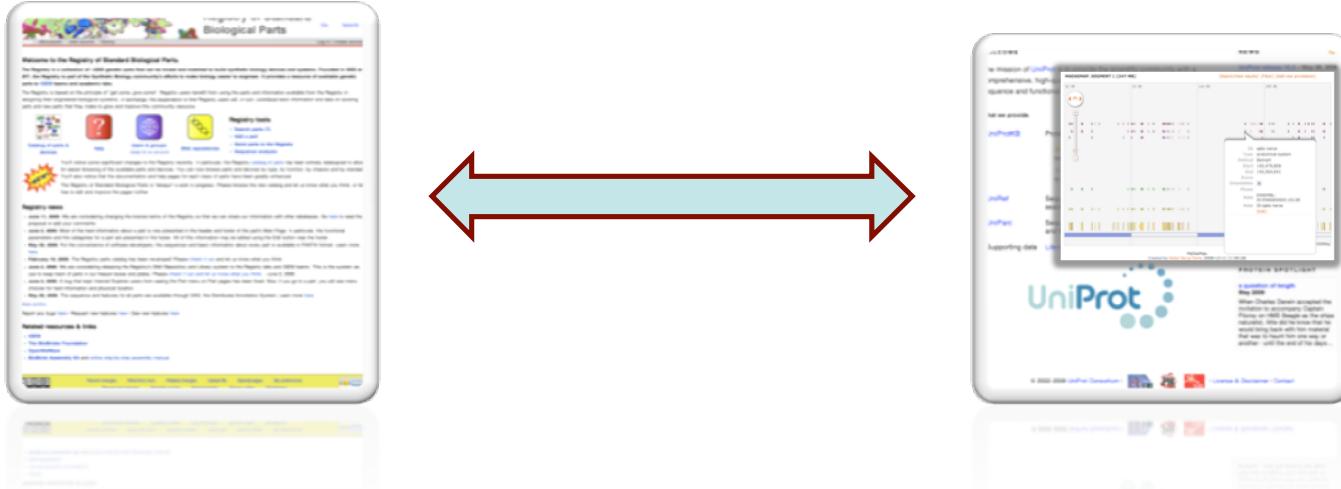
The "Save" button is located at the bottom of the modal.

At the bottom of the modal, there is a large block of XML code representing the DAS (Data Atomization Specification) annotations:

```
<DASDSN>
- <DSN>
  - <SOURCE id="Rhodococcus_RHA1" version="6" project="25">Rhodococcus_RHA1</SOURCE>
  - <MAPMASTER>
    http://madas2.bioinfo.cnio.es/plugins_dir/dasServers/das_common-server/das.php/Rhodococcus_RHA1
  </MAPMASTER>
  <DESCRIPTION/>
</DSN>
- <DSN>
  - <SOURCE id="Nocardoides_KP7" version="6" project="25">Nocardoides_KP7</SOURCE>
  - <MAPMASTER>
    http://madas2.bioinfo.cnio.es/plugins_dir/dasServers/das_common-server/das.php/Nocardoides_KP7
  </MAPMASTER>
  <DESCRIPTION/>
</DSN>
- <DSN>
  - <SOURCE id="Pseudomonas_SY5" version="6" project="25">Pseudomonas_SY5</SOURCE>
  - <MAPMASTER>
    http://madas2.bioinfo.cnio.es/plugins_dir/dasServers/das_common-server/das.php/Pseudomonas_SY5
  </MAPMASTER>
  <DESCRIPTION/>
</DSN>
- <DSN>
  - <SOURCE id="Rhodobacter_sphaeroides_IL106" version="6" project="25">Rhodobacter_sphaeroides_IL106</SOURCE>
  - <MAPMASTER>
    http://madas2.bioinfo.cnio.es/plugins_dir/dasServers/das_common-server/das.php/Rhodobacter_sphaeroides_IL106
  </MAPMASTER>
  <DESCRIPTION/>
</DSN>
```

Added annotations are also available through DAS

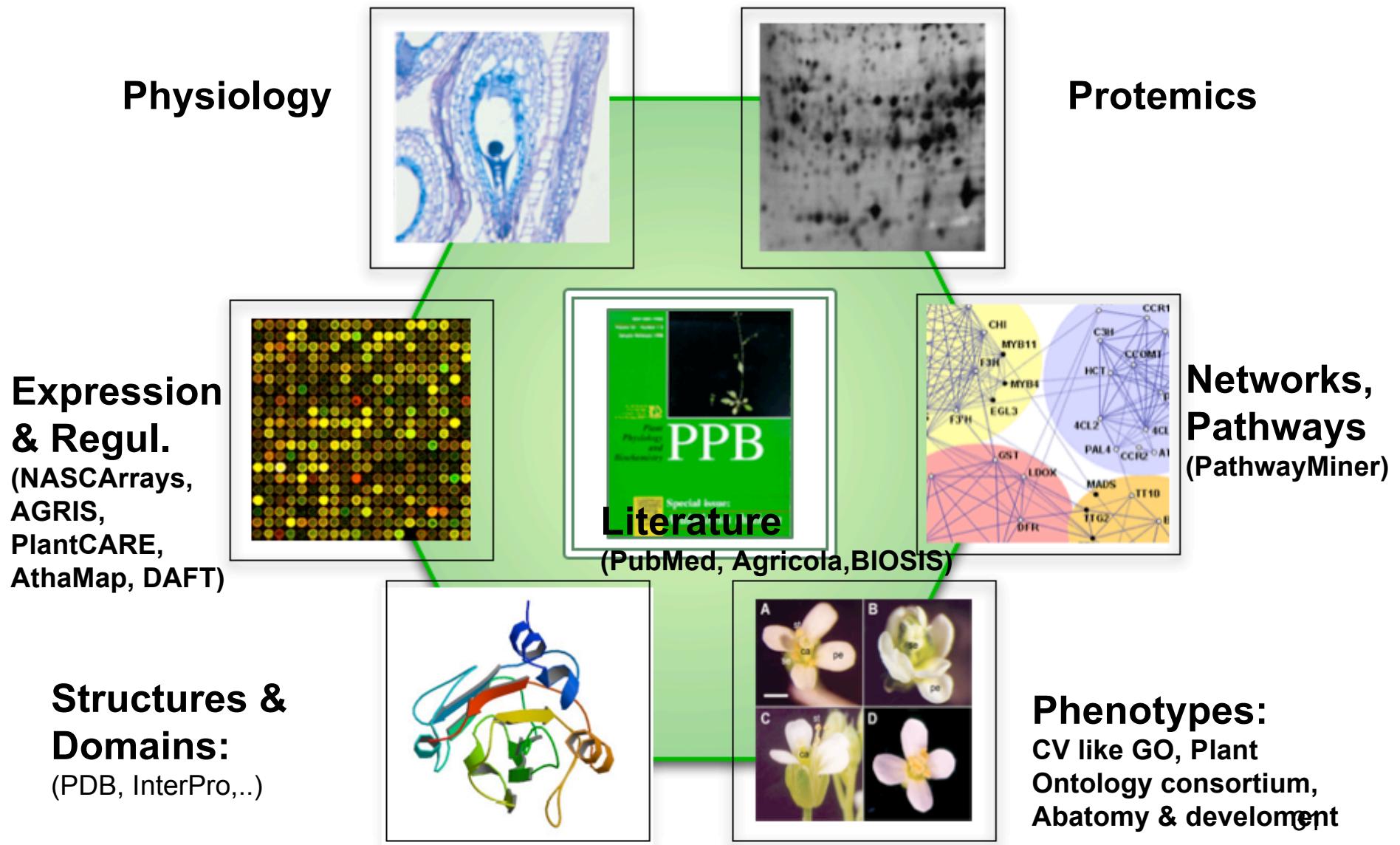
How to exchange annotations



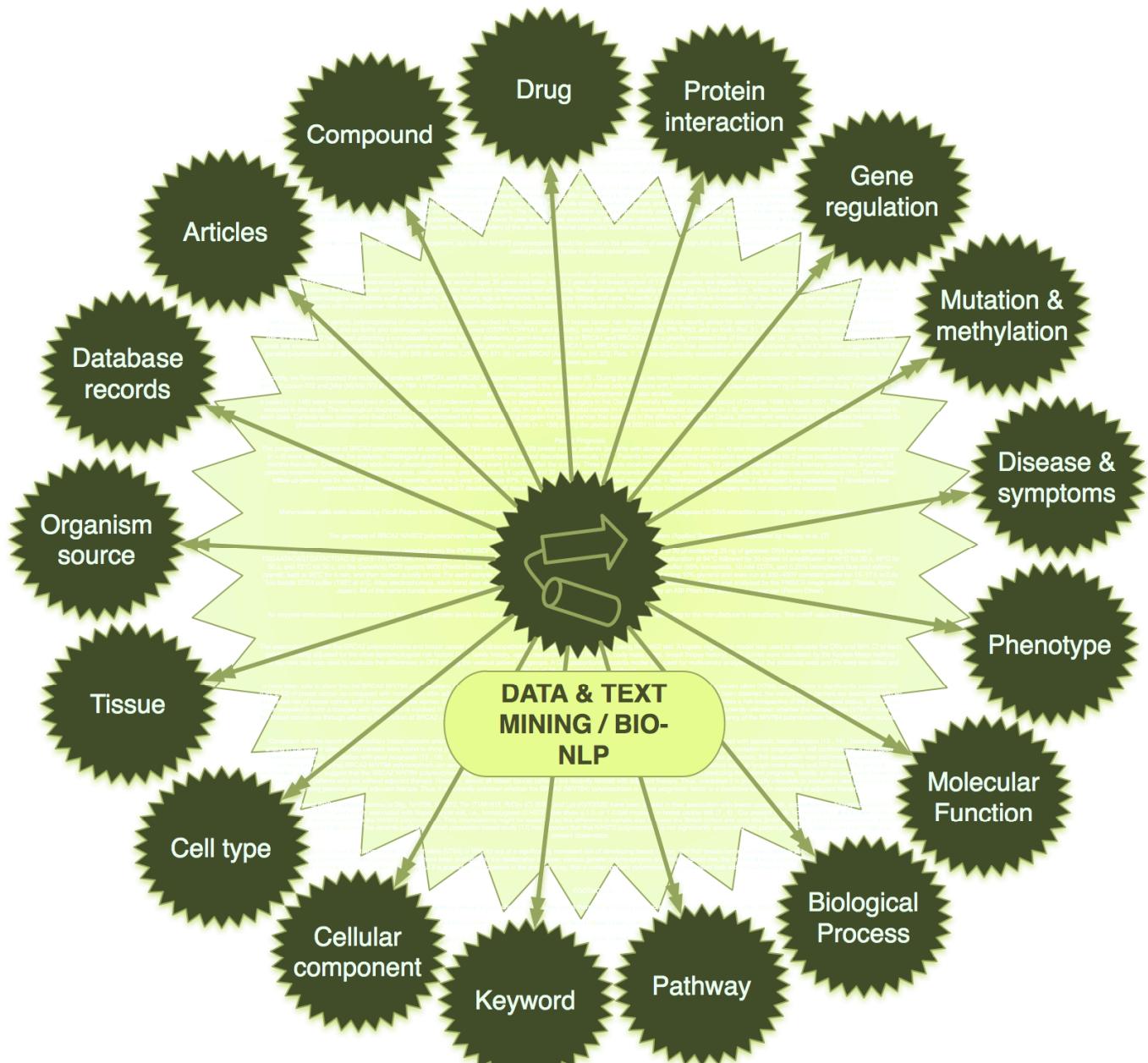
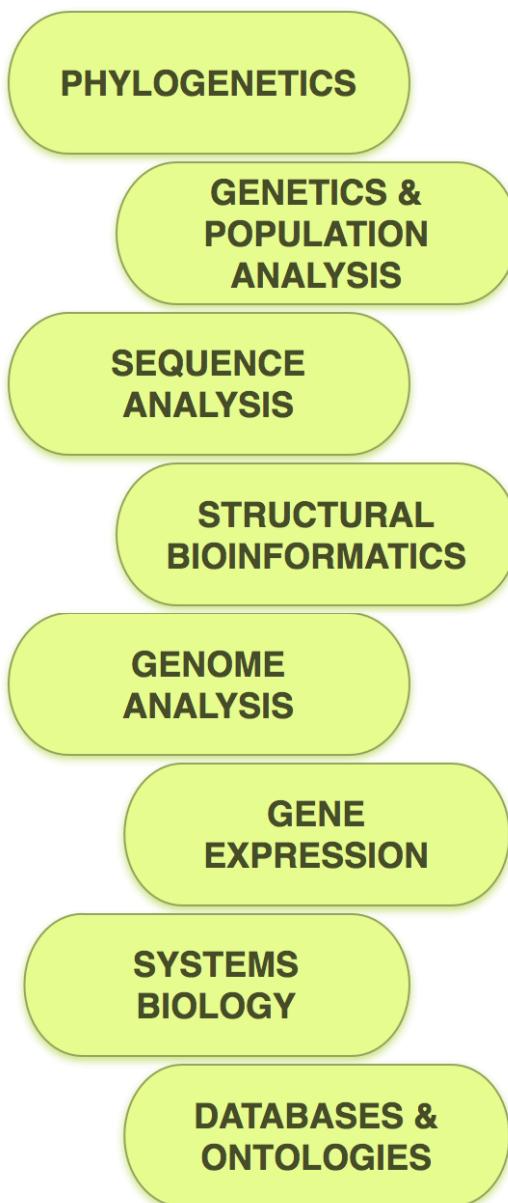
- ✓ Distributed annotation system (DAS) protocol. (MR)
- ✓ Web services. (MR)
- ✓ Database dump. (MR)
- ✓ Biological Web Elements and Registry Embed Code. (HR)

MR = Machine readable
HR = Human readable

Integration of heterogeneous data types



Text mining covers multiple topics



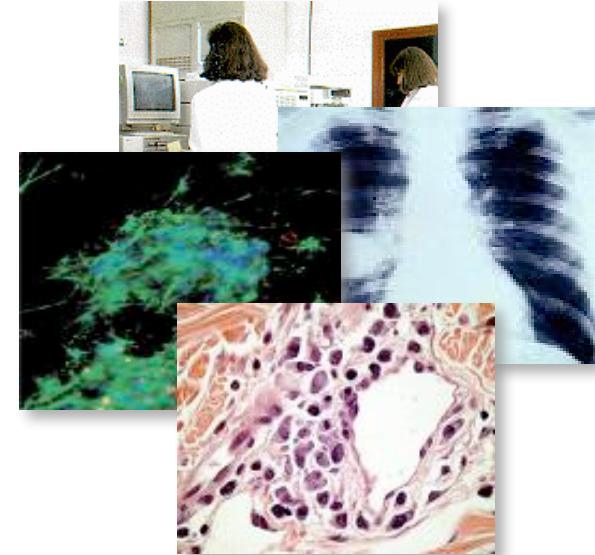
Importance of literature data for Biology

- 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
- 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

Biology



- Resource for clinical decision support in evidence-based clinical practice
- Useful information for diagnostic aids

Clinics



- Drug discovery and target selection
- Identifying adverse drug effect
- Competitive intelligence and knowledge management

Pharma

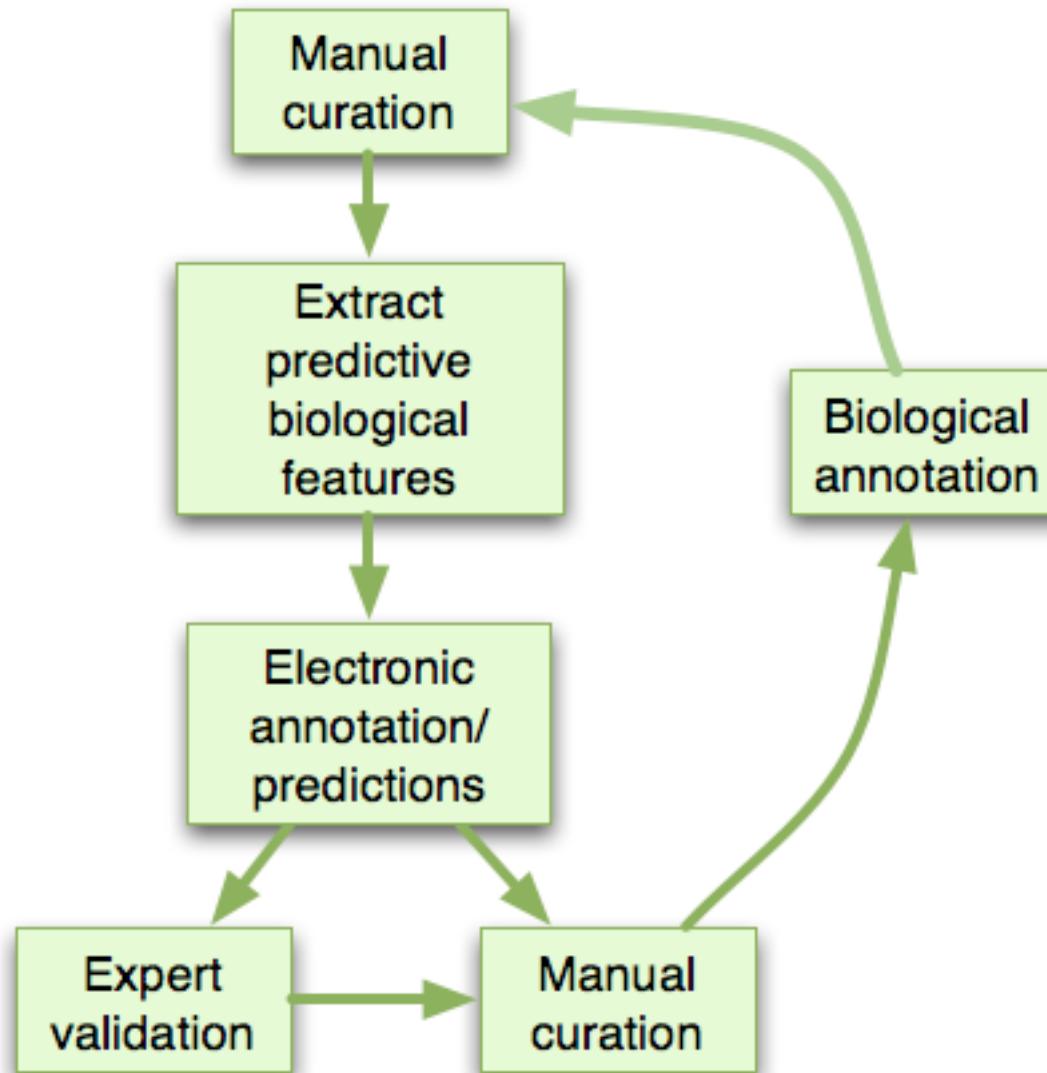
- Global view of the current research state & monitor trends to ensure optimal resource allocation

Funding

- Find domain experts for specific topics for the peer-review process & detecting potential cases of plagiarism

Publ.

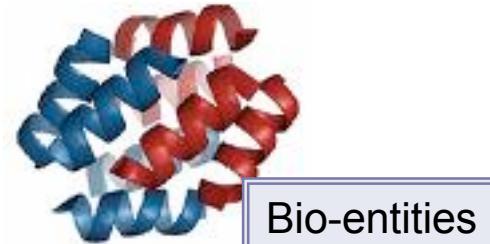
Literature Gold Standard datasets / DBs



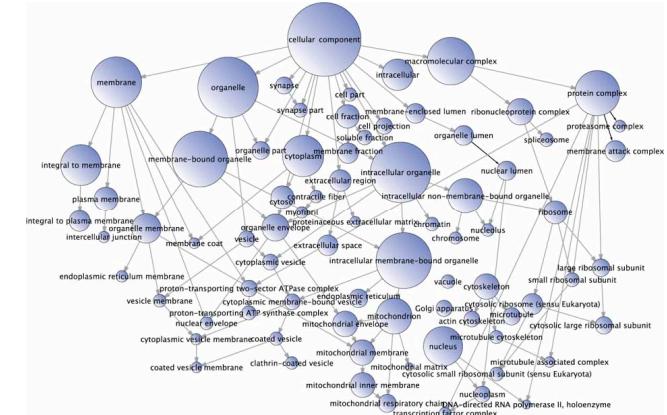
Biocuration: manual literature annotations & databases



Scientific Literature



Database curator



Controlled vocabularies

(*Lycopersicon esculentum*). Here, we demonstrate that two *Arabidopsis thaliana* MAF1 homologs, **WPP1** and WPP2, are associated with the NE specifically in undifferentiated cells of the root tip. Reentry into cell cycle after

Locus: AT5G43070

Date last modified 2003-05-02
TAIR Accession Locus:2167831
Representative Gene AT5G43070.1
Model Other names: MMG4.9, MMG4_9, WPP DOMAIN PROTEIN 1,



FlyBase

MaizeGDB
Maize Genetics and Genomics Database



EcoCyc

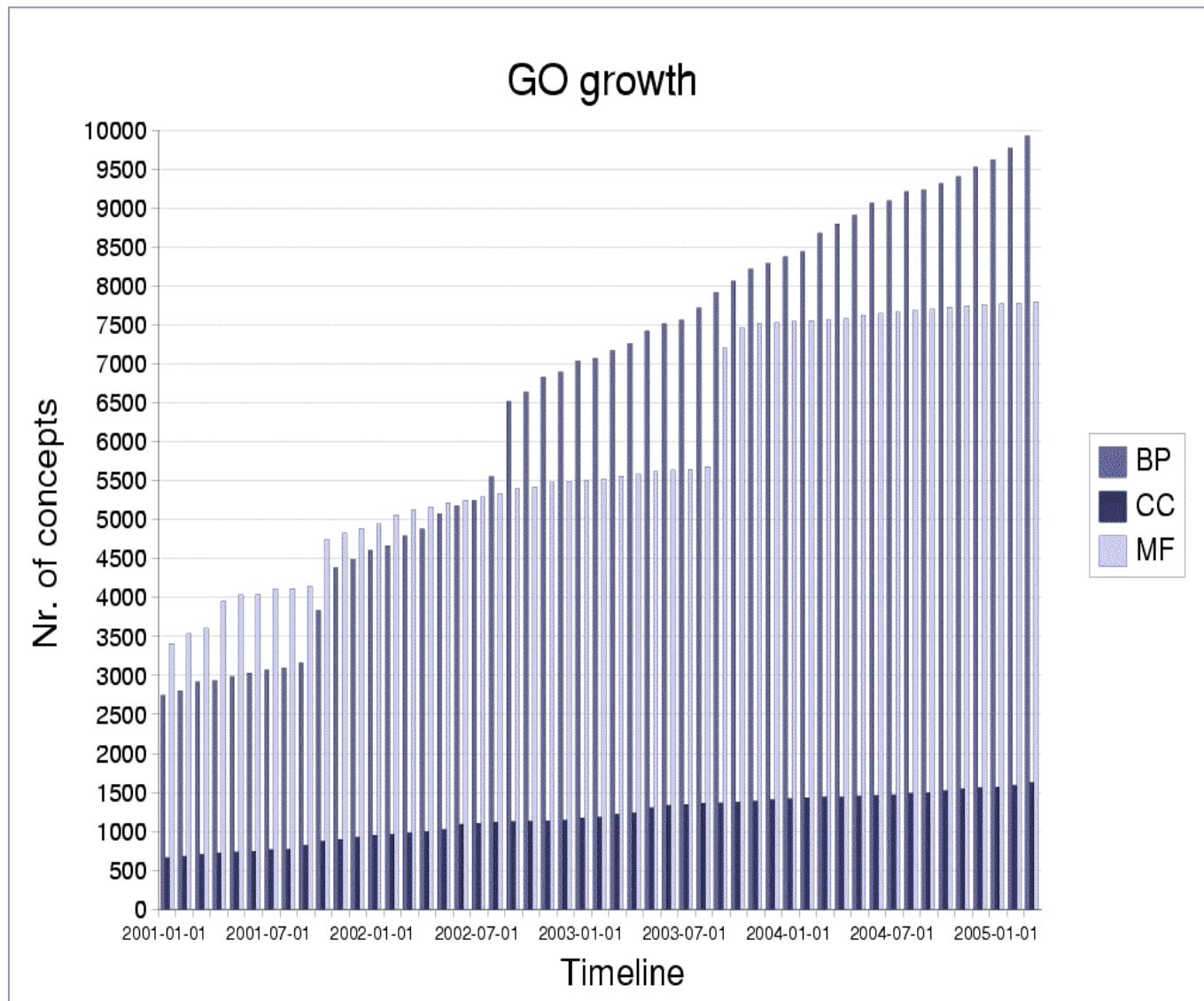


MINT



36

Curation challenge I: growing number of CV terms



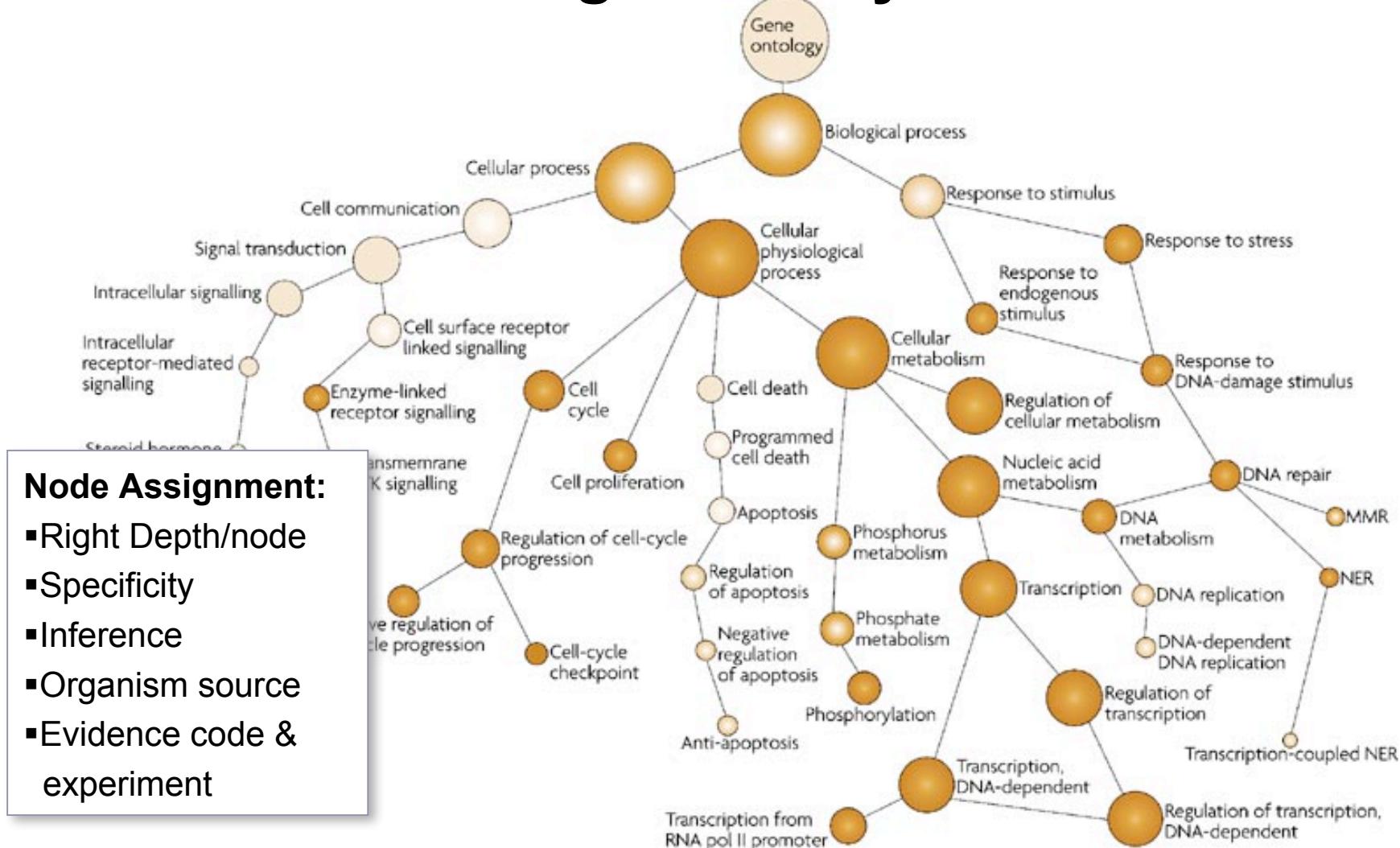
Curation challenge II: growing number of ontologies

> 130



Domain	File
Biological imaging methods	image.obo
Biological process	gene ontology.obo
BRENDA tissue / enzyme source	worm development.obo
C. elegans development	worm anatomy.obo
C. elegans gross anatomy	worm phenotype.obo
C. elegans phenotype	cell.obo
Cell type	gene ontology.obo
Cellular component	cereals development.obo
Cereal plant development	plant trait.obo
Cereal plant trait	chebi.obo
Chemical entities of biological interest	caro.obo
Common Anatomy Reference Ontology	dictyostelium anatomy.obo
Dictyostelium discoideum anatomy	fly development.obo
Drosophila development	fly anatomy.obo
Drosophila gross anatomy	envo.obo
Environment Ontology	event.obo
Event (INOH pathway ontology)	evidence code.obo
Evidence codes	evoc.obo.tar (v2.7)
eVOC (Expressed Sequence Annotation for Humans)	fly taxonomy.obo
Fly taxonomy	flybase controlled vocabulary.obo
FlyBase Controlled Vocabulary	fma.obo.obo
Foundational Model of Anatomy (subset)	fungal anatomy.obo
Fungal gross anatomy	protege source
Habronattus courtship	human dev anat abstract.obo
Human developmental anatomy, abstract version	human dev anat staged.obo
Human developmental anatomy, timed version	human disease.obo
Human disease	B8467OBO
Loggerhead nesting	zea mays anatomy.obo
Maize gross anatomy	mammalian phenotype.obo
Mammalian phenotype	medaka ontology.obo
Medaka fish anatomy and development	MGEDOntology.owl
Microarray experimental conditions	gene ontology.obo
Molecular function	molecule role.obo
Molecule role (INOH Protein name/family name ontology)	mosquito anatomy.obo
Mosquito gross anatomy	mosquito insecticide resistance.obo
Mosquito insecticide resistance	adult mouse anatomy.obo
Formats (OBO, OWL, XML, RDF) (http://www.obofoundry.org)	

Curation challenge III: annotation granularity



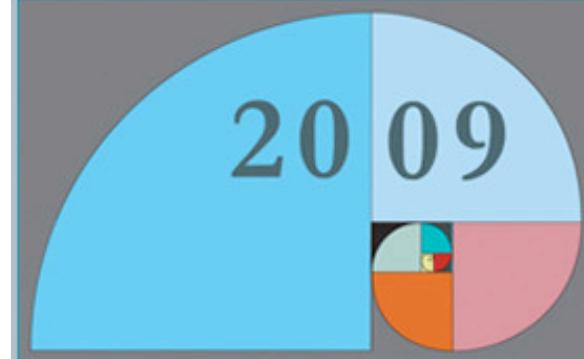
Node Assignment:

- Right Depth/node
 - Specificity
 - Inference
 - Organism source
 - Evidence code & experiment

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

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**



2009

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

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

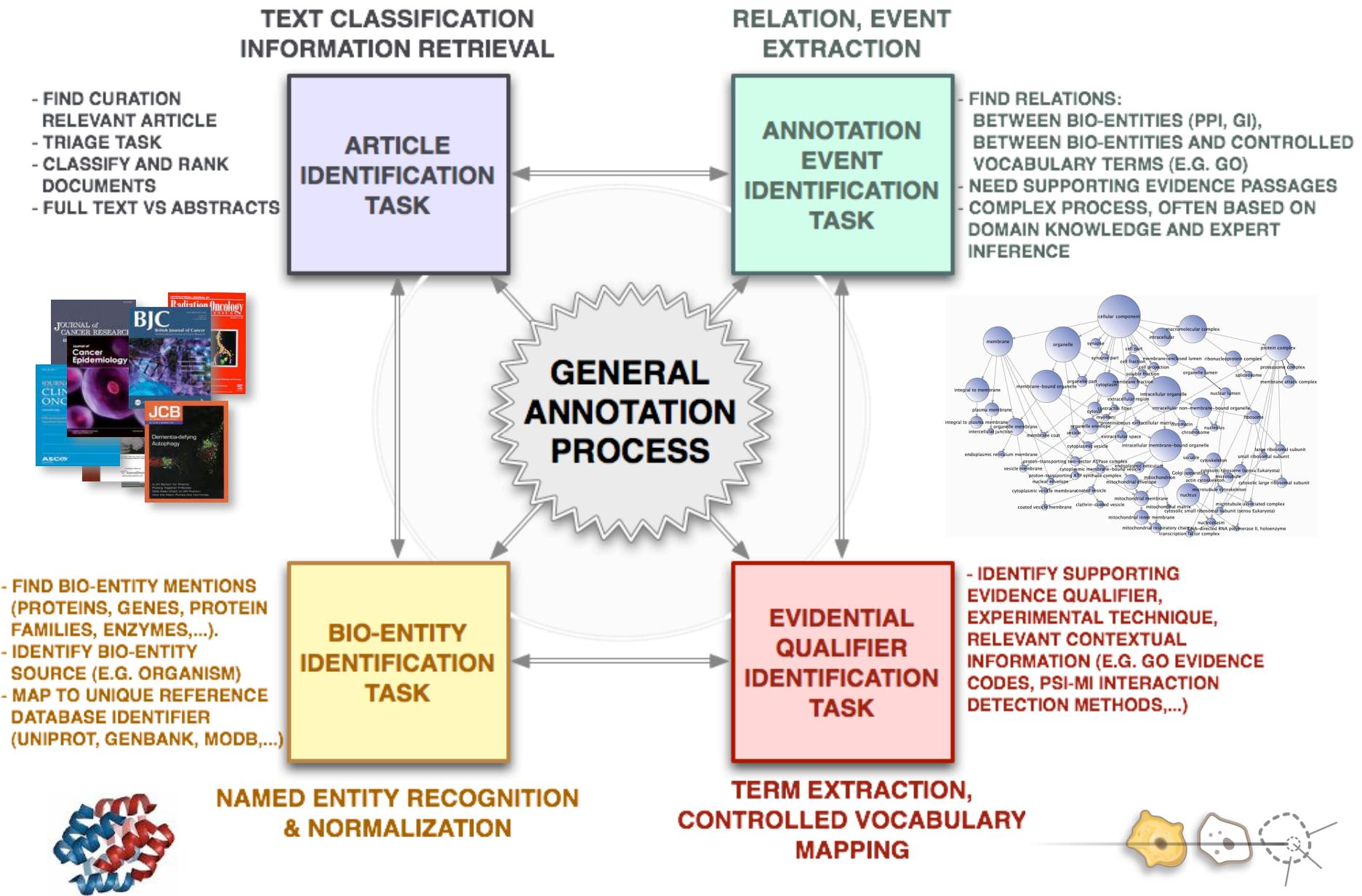
biocurator.org

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.

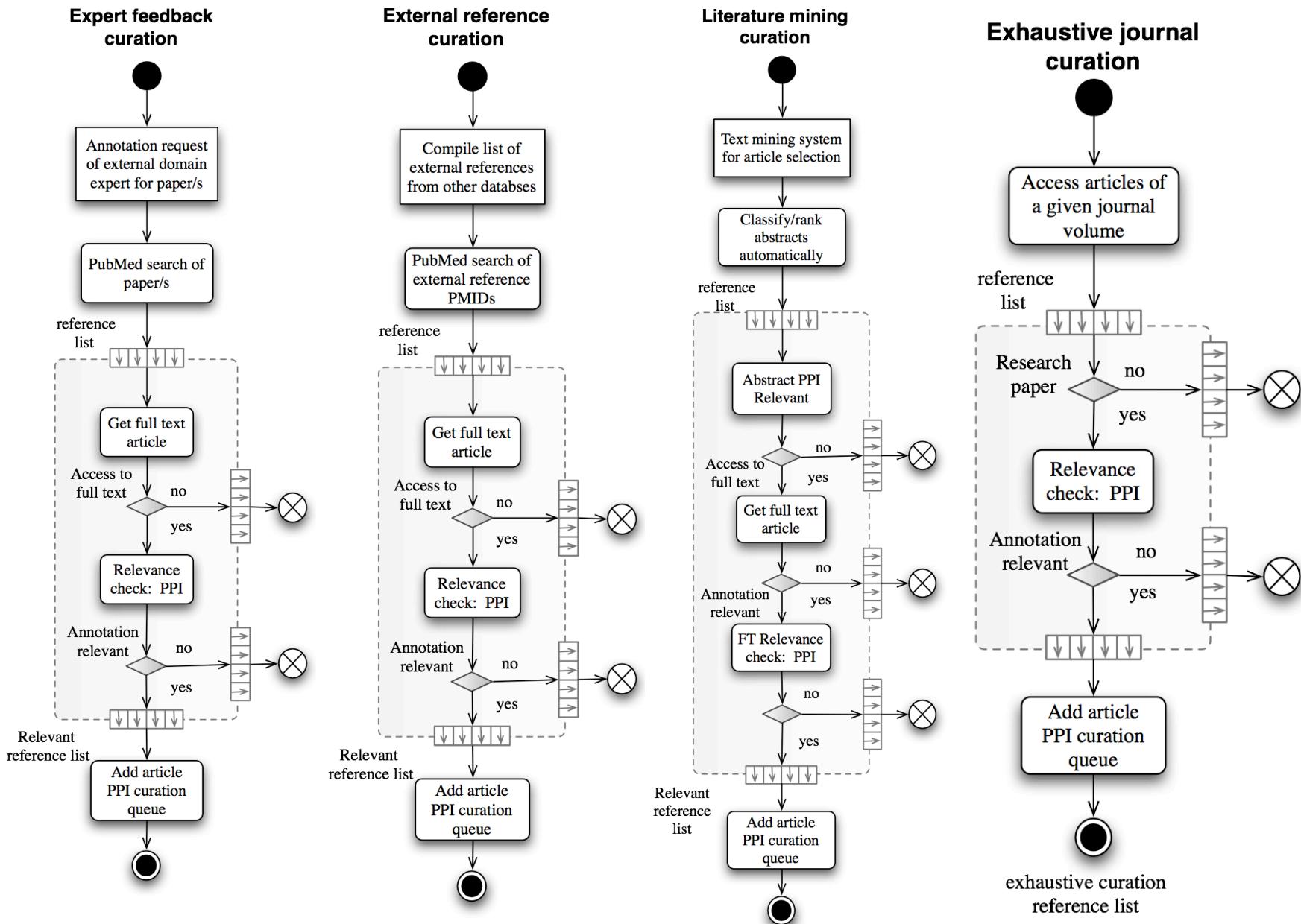
BIOCURATION WORKFLOW TASKS



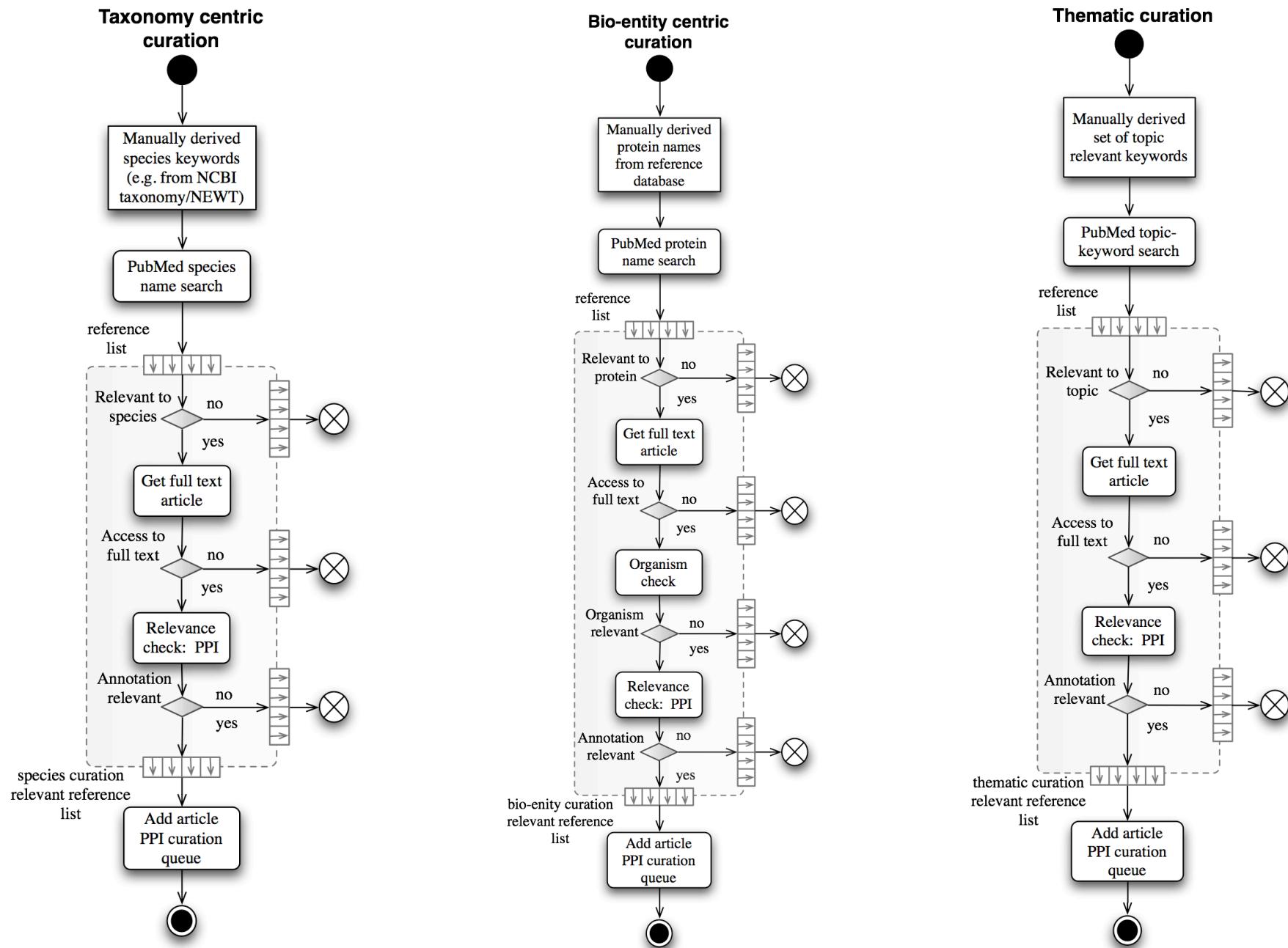
WORKFLOW TASKS AND TEXT MINING

- 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 EFFICIENCY USING TEXT MINING ASSISTANCE
- USER FEEDBACK AND POTENTIAL ITERATIVE IMPROVEMENTS

ARTICLE IDENTIFICATION:TRIAGE TASK (1)



ARTICLE IDENTIFICATION:TRIAGE TASK (2)



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,...

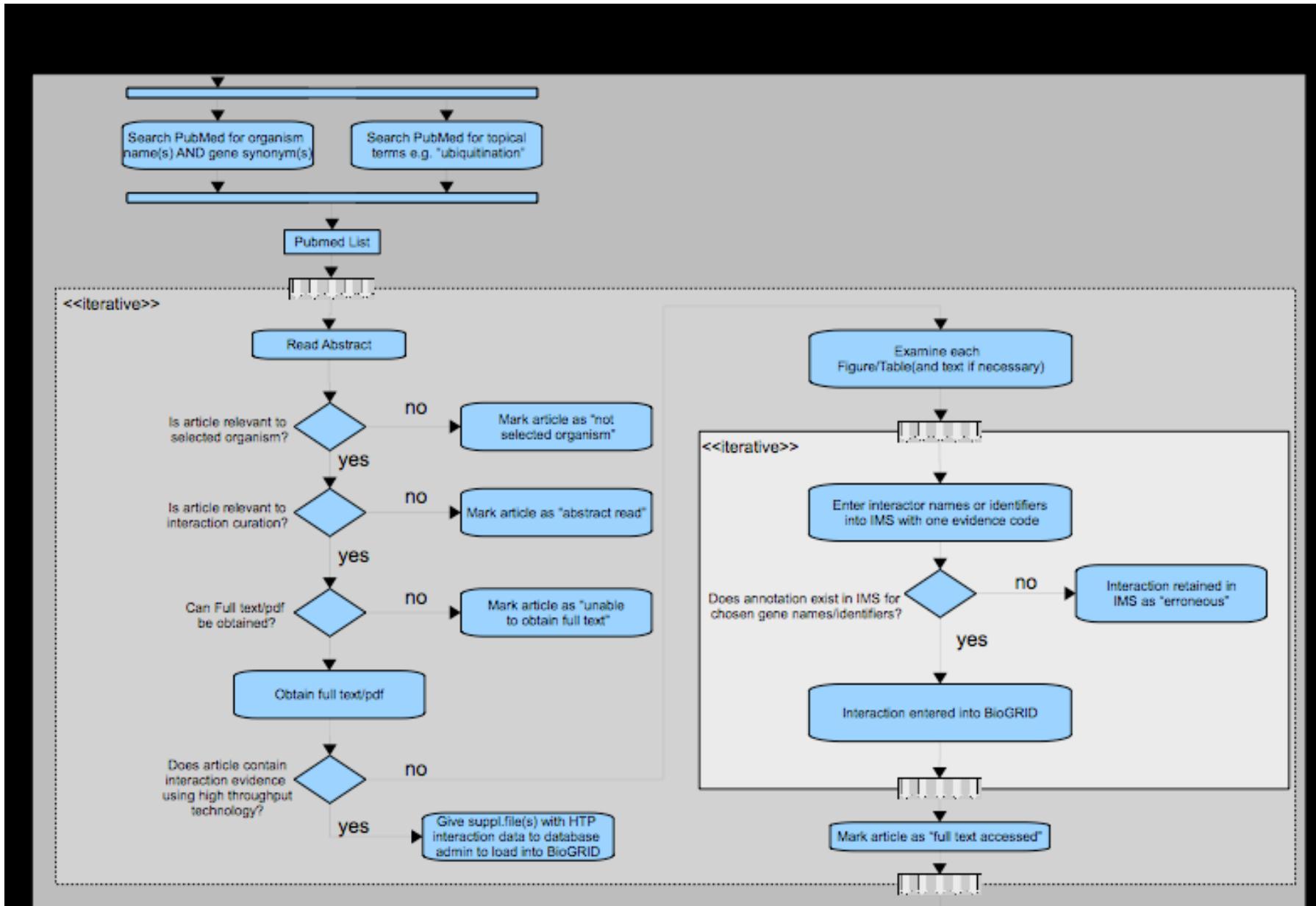
ANNOTATION EVENT IDENTIFICATION TASK

- 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-occurrence 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



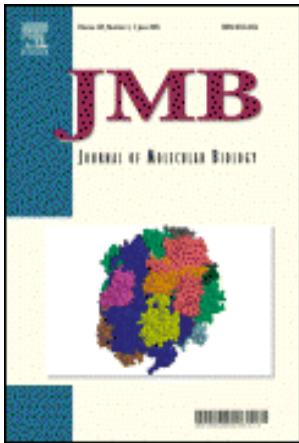
Many thanks to Andrew Winter



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



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.

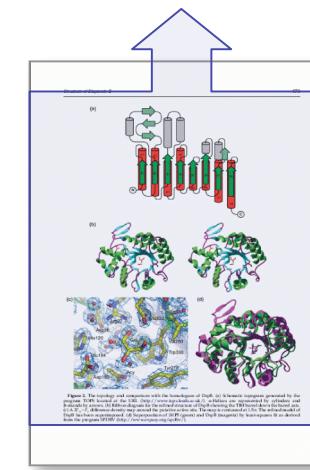
Full text:

- Title
 - Authors
 - Abstract
 - Text Body
 - References

Tables & table legends



Figures & figure legends



Processing levels of natural language texts

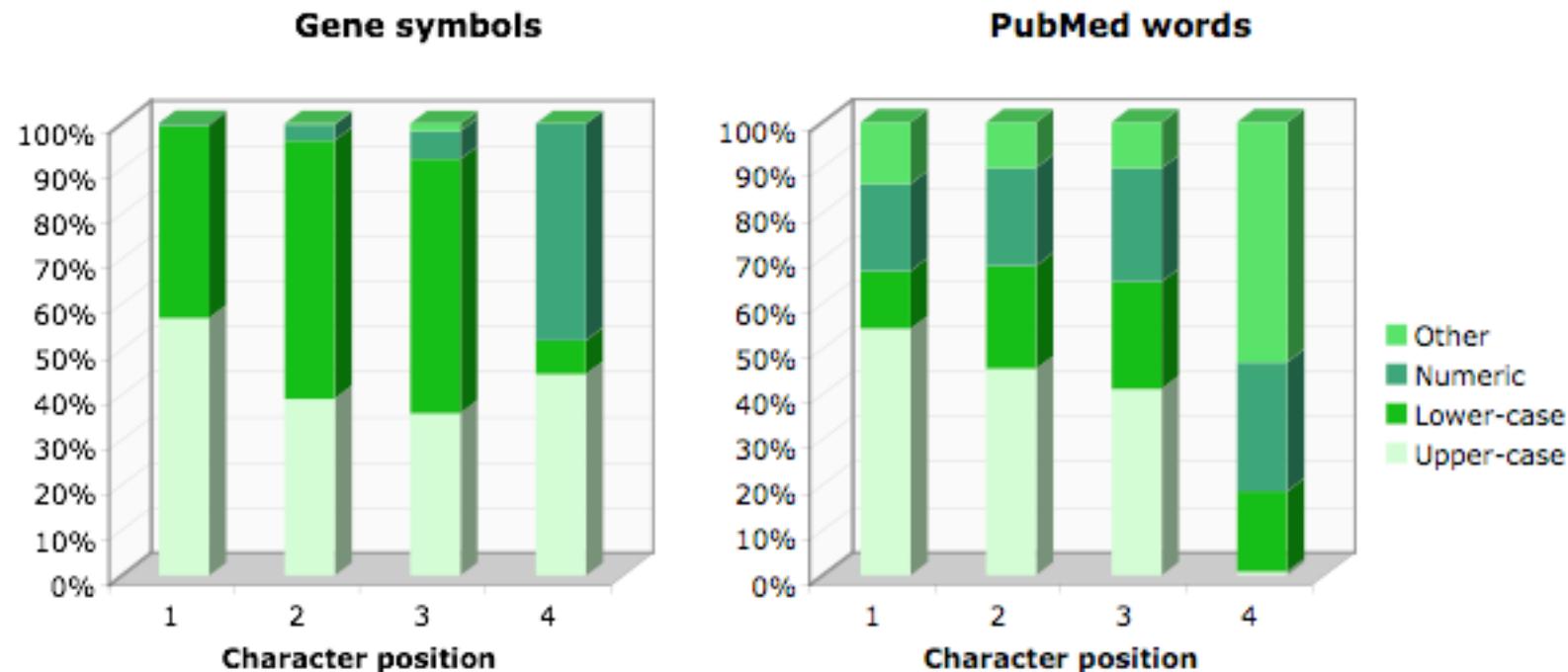
PROCESSING LEVEL	TASKS AND APPLICATIONS
Character & strings level	Word tokenization, sentence boundary detection, gene symbol recognition, text pattern extraction
Word token level	POS-tagging, parsing, chunking, term extraction, gene mention recognition
Sentence level	Sentence classification and retrieval and ranking, question answering, automatic summarization
Sentence window level	Anaphora resolution
Paragraph & passages level	Detection of rhetorical zones
Whole document level	Document similarity calculation
Multi-document collection level	Document clustering, multi-document summarization

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



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

Table 2 | Most frequently used words in various countries

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 repercussion on the oxidative phenotype in a white population.	7697944
Japan	Useful	Bullfrog, shadow (in radiography)	Clarify	Faintly, next, suddenly, scarcely	MDR-1 protein was faintly 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 arises wholly from pyramidal cells.	11602231
Russia	Gravitational	(Space) mission, quantum, hibernate, peculiarity, regularity, realization	–	Thermodynamically	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	Hydroxypropylmethylcellulose (HPMC) was used to impart strength and sphericity to the agglomerates.	12476867
France	Exceptional, digestive	Trouble	Envisage (imagine)	Successively (sequentially), essentially, sometimes	These 2 cells [sic] lines being able to clone, it is hard to envisage 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 alleviates 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 unremarkable 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	–	Loss of revenue, mainly from noncompliance with charge capture resulted in the hospital billing only US\$386,794.32 with a total reimbursement of US\$165,779.86.	12488156

Words in bold typeface have specific meanings and are probably related to local research rather than to local language usage. The bold and underlined words in the example sentences indicate the most abundant country-specific terms. The words shown were found to be more common in the abstracts of the corresponding country than in the abstracts of any other of the 19 representative countries (as in Fig. 2). Note that most of the sentences are grammatically correct, but the usage of the marked (bold and underlined) words is unusual. PMID ref, PubMed reference number.

Literature repositories for life sciences

- NLP: need electronically accessible texts.
- Main scientific textual data types: e-books and e-articles 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



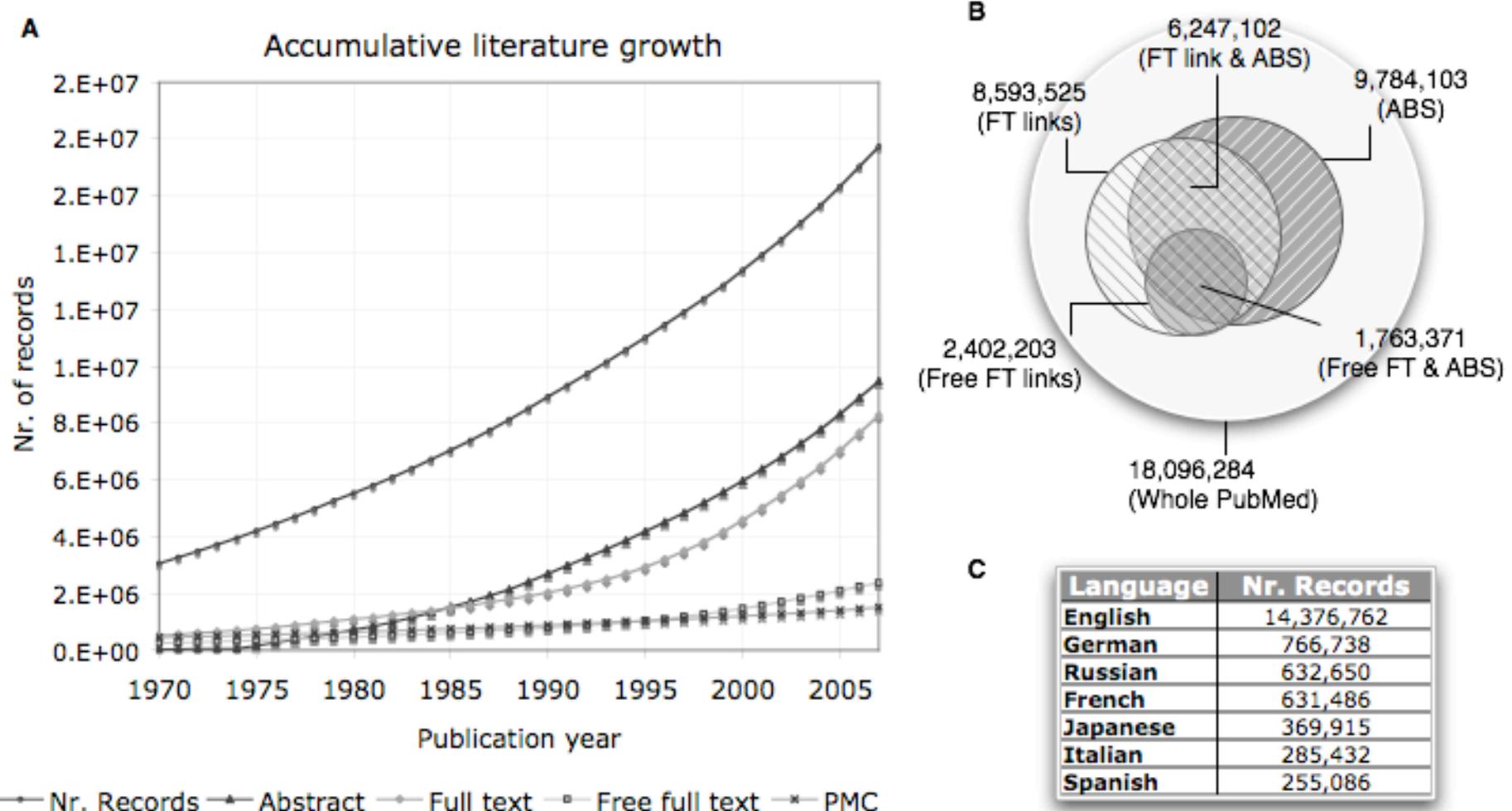
- Scientific articles: new scientific discoveries.
- Citation entries of scientific articles of all biomedical sciences, nursing, biochemistry, engineering, chemistry, environmental sciences, psychology, etc,...
- Developed at the NCBI (NIH).
- Digital library contains more than 16 million citations
- 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

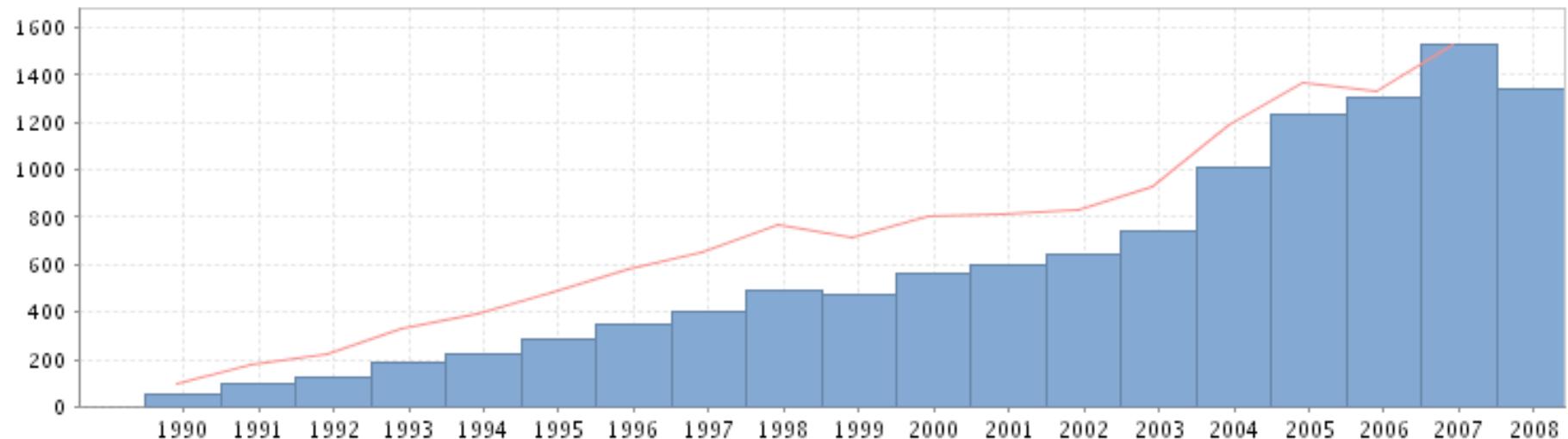
PubMed growth



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

Arabidopsis articles in PubMed



PubMed XML record

```
<PubmedArticle>
  <MedlineCitation Status="Publisher" Owner="NLM">
    <PMID>18642075</PMID>
    <DateCreated>
      <Year>2008</Year>
      <Month>7</Month>
      <Day>21</Day>
    </DateCreated>
    <Article PubModel="Print-Electronic">
      <Journal>
        <ISSN IssnType="Print">0167-6806</ISSN>
        <JournalIssue CitedMedium="Print">
          <PubDate>
            <Year>2008</Year>
            <Month>Jul</Month>
            <Day>19</Day>
          </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>
        <AbstractText>BRCA1/BRCA2 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>
        </Author>
      </AuthorList>
    </Article>
  </MedlineCitation>
</PubmedArticle>
```

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	O
with	with	IN	B-PP	O
cortactin	cortactin	NN	B-NP	B-protein
in	in	IN	B-PP	O
the	the	DT	B-NP	O
apical	apical	JJ	I-NP	O
membrane	membrane	NN	I-NP	O
of	of	IN	B-PP	O
hepatocytes	hepatocyte	NNS	B-NP	B-cell_type
.	.	.	O	O
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

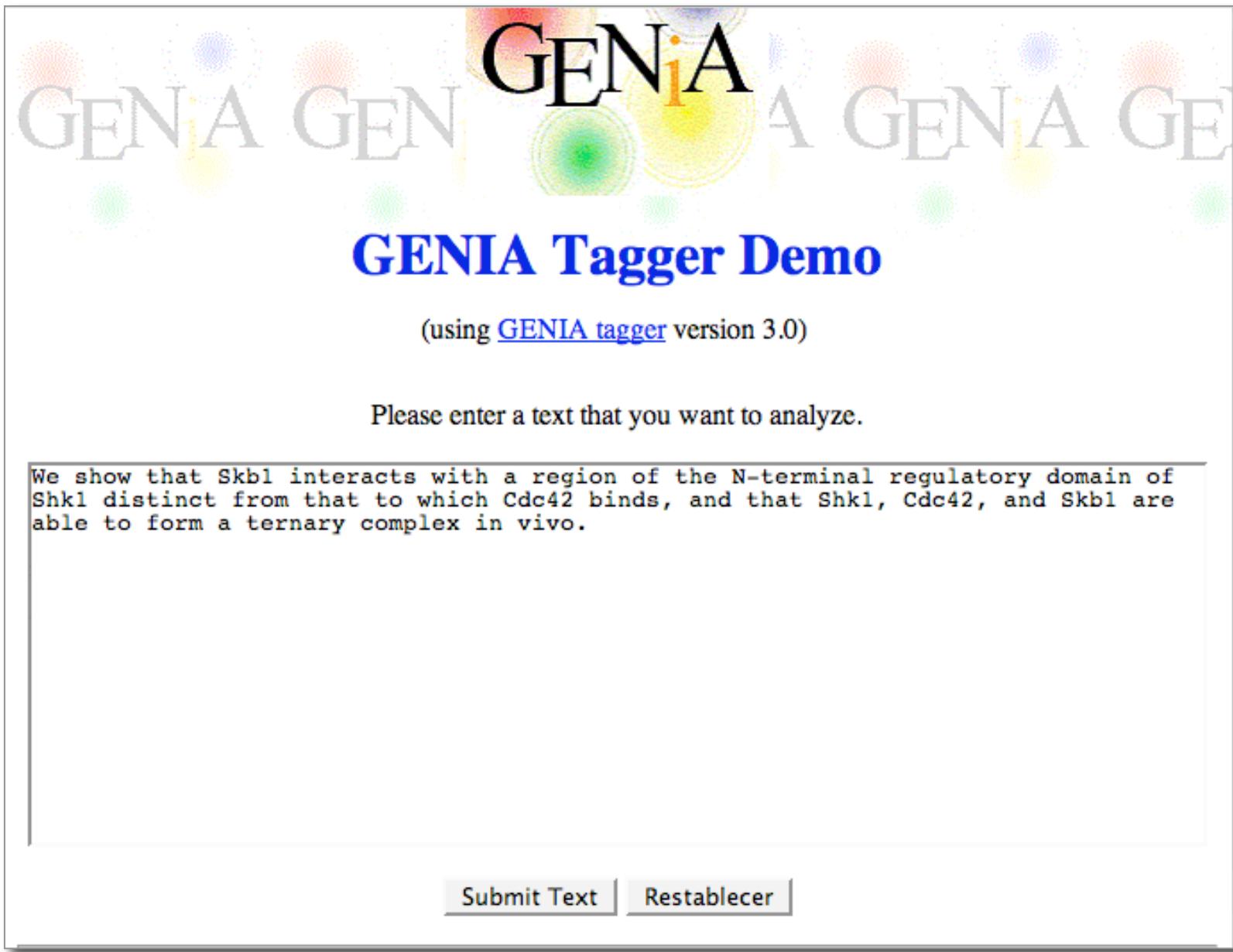
Grammatical features

- 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



The image shows the GENIA Tagger Demo interface. At the top, there is a decorative banner with the word "GENIA" repeated in a grayscale dotted pattern. Below the banner, the title "GENIA Tagger Demo" is displayed in large blue capital letters. Underneath the title, the text "(using [GENIA tagger](#) version 3.0)" is shown in smaller blue font. A prompt "Please enter a text that you want to analyze." is followed by a text input area containing the following text: "We show that Skbl interacts with a region of the N-terminal regulatory domain of Shk1 distinct from that to which Cdc42 binds, and that Shk1, Cdc42, and Skbl are able to form a ternary complex *in vivo*." At the bottom of the interface, there are two buttons: "Submit Text" and "Restablecer".

GENIA Tagger Demo

(using [GENIA tagger](#) 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 Shk1 distinct from that to which Cdc42 binds, and that Shk1, 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 Algorithm Online - Mozilla Firefox

File Edit View Go Bookmarks Tools Window Help

Home Bookmarks Yahoo Google MK Homepage

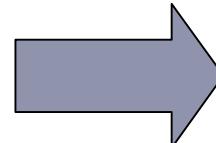
Porter's Stemming Algorithm Online

Enter a sequence of words in the box below to stem
(Note: "stop" words and punctuation are automatically removed)

Glycogenin is the self-glycosylating protein primer that initiates glycogen granule formation. To examine the role of this protein during glycogen resynthesis, eight, male subjects exercised to exhaustion on a cycle ergometer at 75% VO₂ max followed by 5 x 30s sprints at maximal capacity to further deplete glycogen stores. During recovery, carbohydrate (75g/h) was supplied to promote rapid glycogen repletion and muscle biopsies were obtained from the vastus lateralis at 0, 30, 120 and 300min post-exercise. At time 0,

Stem!

M Done



Porter's Stemming Algorithm

File Edit View Go Bookmarks Tools Window

Home Bookmarks Yahoo Google

Porter's Stemming Results

Original Word	Stemmed Word
glycogenin	glycogenin
selfglycosylating	selfglycosyl
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>

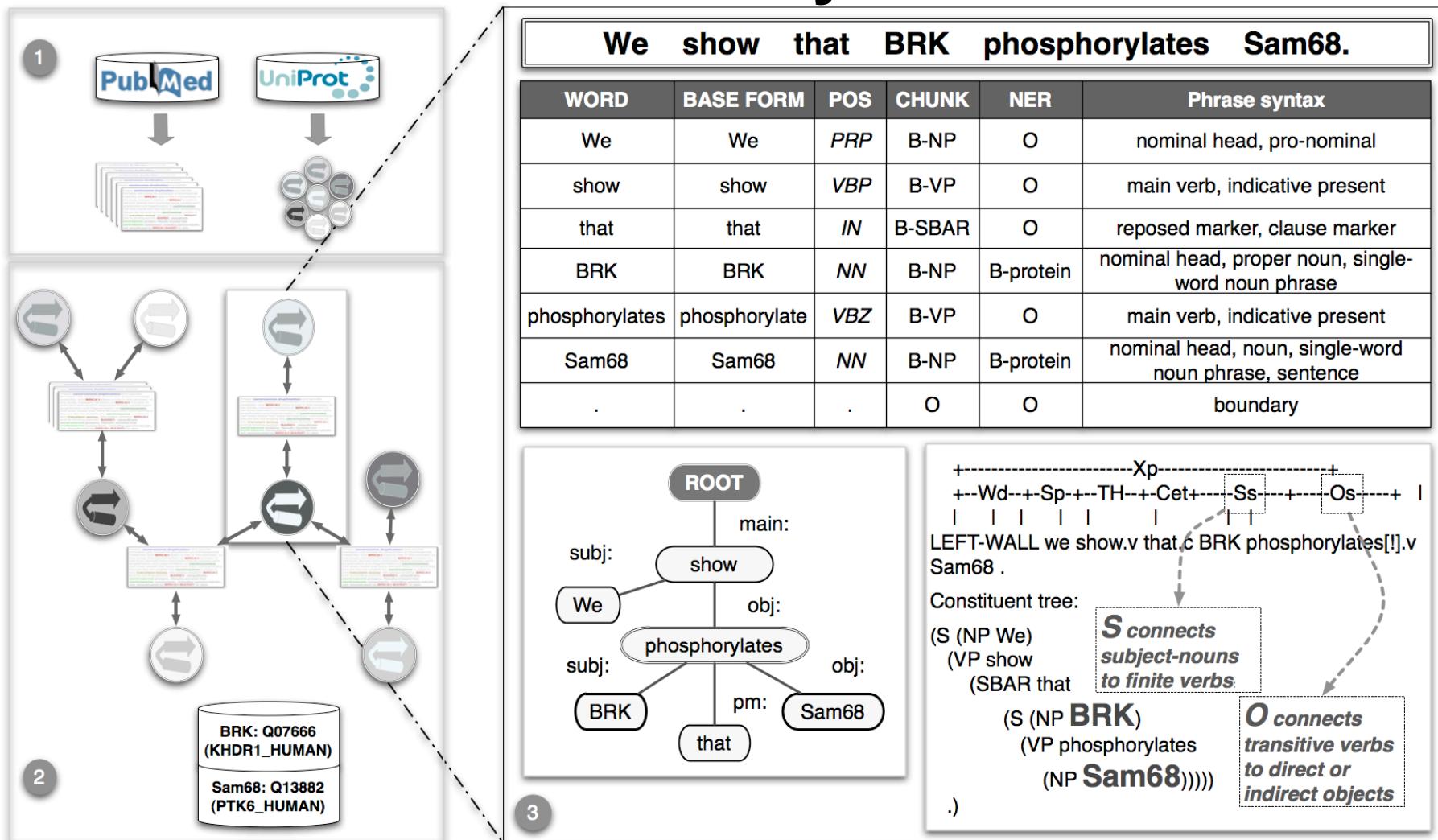
Syntactic features

- 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



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	Type	Regulation association evidence sentence	iHOP	WikiG	TAIR	Source	
16055635	E2FB AT5G22220	CDKB1;1 AT3G54180	Activation	Indeed, our recent data show that E2FB can directly induce the promoter of the <i>Arabidopsis CDKB1;1</i> gene (Z. Magyar, unpublished results).	0	None	 iHOP	 WikiG	 TAIR
16055635	E2FA At2g36010	CDKB1;1 AT3G54180	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	 iHOP	 WikiG	 TAIR

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 , w: term weight
- (2) Term weighting, tf: term frequency $w_{i,j} = tf_{i,j} \times idf_j$
- (3) Similarity coefficient idf: inverted document frequency

- Query: a list of terms or even whole documents.

- Query as vectors of terms.

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

- Term weighting (w) according to their frequency:
within the document (i) & within the document collection (d)

- Widespread term weighting: tf x idf.

- Calculate similarity between those vectors.

- Cosine similarity often used.

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

- Return a ranked list.

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

- Example: related article search in PubMed

eTBLAST

The image shows the eTBLAST search interface. At the top left is the eTBLAST logo and a grid of smaller logos for Bibus, TRITE, eTSNAP, FRISC, VER.2, RIC, ARGH, and De Ja Vu. To the right is the text: "eTBLAST: A text similarity-based engine for searching literature collections".

Input your text

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 mutations.

--OR--

Upload a [text only](#) file [Browse...](#)

Optional Email:

If you would like your results emailed to you, please enter an email address. Your address will be kept strictly confidential, and will not be used for any other purpose.

Search Database

- MEDLINE
- NASA
- IOP
- CRISP
- USPTO (coming)
- PMC "Methods" (coming)
- OMIM (coming)
- DrugBank (coming)

**•Ranked list of abstracts
•Visualize Pairwise Comparisons
•Find an Expert in this Field
•Find a Journal for your Manuscript
•Publication History of this Topic**

eTBLAST results: high scoring words

The screenshot shows a Mozilla Firefox window with the title bar "PMID: 8529663 - Mozilla Firefox". The menu bar includes "Archivo", "Editar", "Ver", "Ir", "Marcadores", "Herramientas", and "Ayuda". The toolbar includes standard icons for back, forward, search, and refresh. The address bar shows the URL "http://invention.swmed.edu/cgi-bin/etblast/abstract_local?pmid=8529663&use". The main content area displays the following text:

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 Lomako
W J Whelan

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 from, and its degradation to, glycogenin, the autocatalytic, self-glucosylating primer for glycogen synthesis. We recently demonstrated that 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, which does not continue to glycogenin. We suggest that in the nutritionally dependent turnover of glycogen in tissues, the molecules cycle between proglycogen and macromolecular glycogen and are not normally degraded to glycogenin. Nevertheless, when this does happen, 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, suggestive of a cluster structure for proglycogen and indicating that breakdown is stepwise. Evidence of post-translational modification of glycogenin was obtained by the finding that, in glycogen from cultured muscle, glycogenin is phosphorylated.

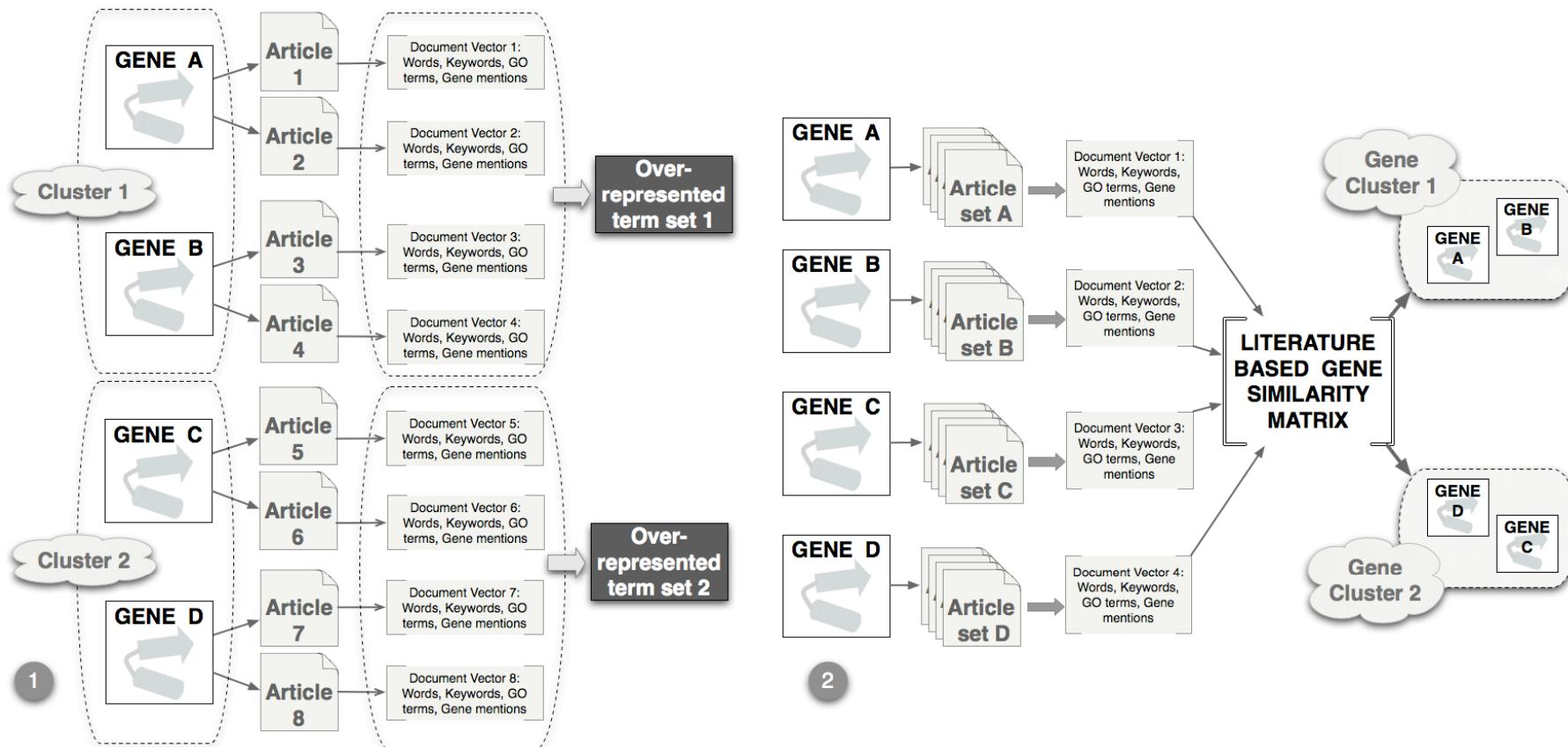
MedlineID: 0
PMID: 8529663

Terms with high weight

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 document 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.

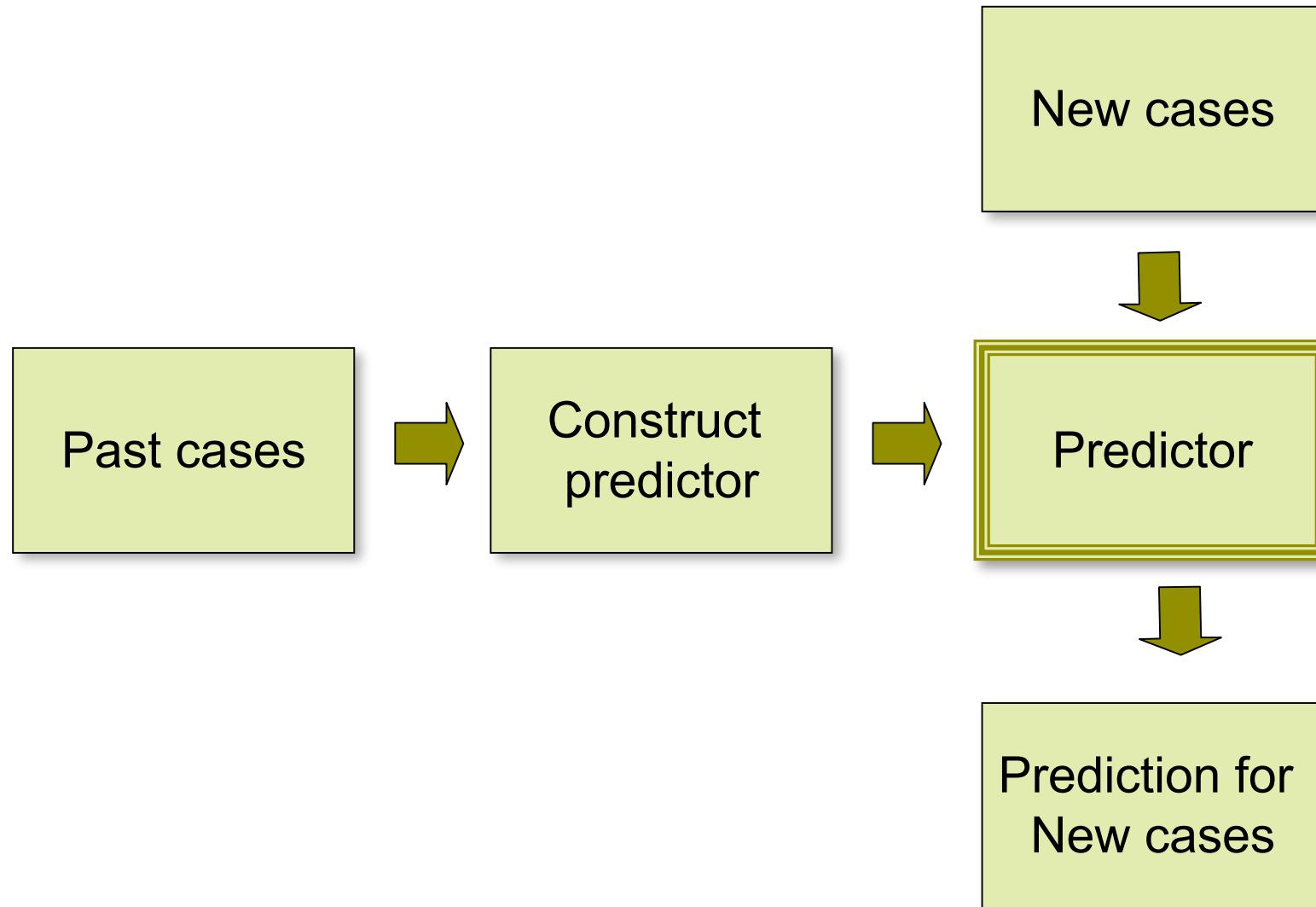
Clustering documents, genes, terms



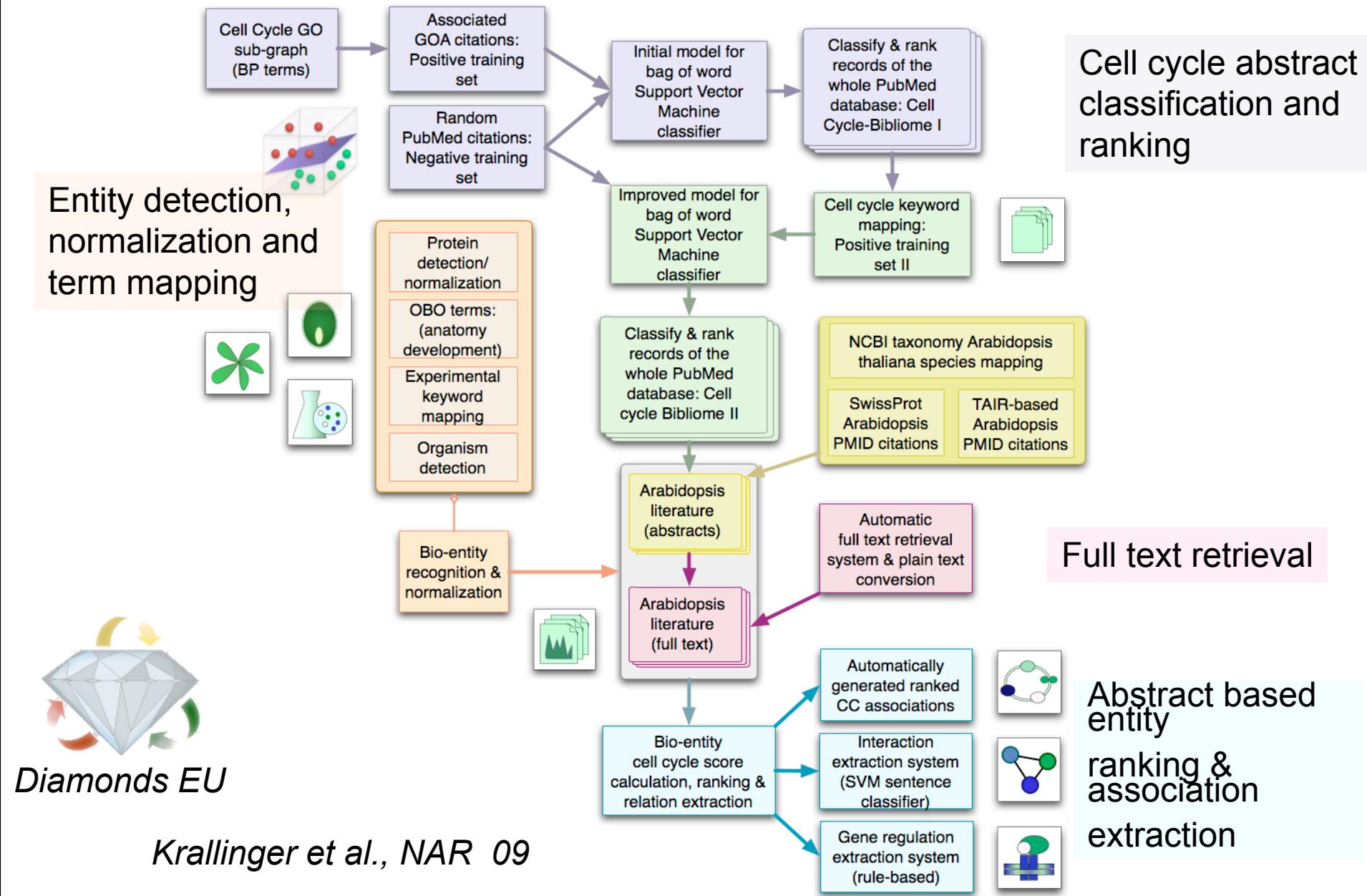
Text classification

- 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: supervised 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
- Can use the bow toolkit, SVMlight, LibSVM etc,..

Text classification & supervised learning



System overview



Cell cycle protein ranking



	TAIR-Id	CC-SVM sum	CC-SVM mean	CC-SVM pos. sum	Nr. abstracts	Mentions	PA	GR	KW	EV
TAIR db gene identifier	AT5G11300	92.501942618	0.872659836019	108.181946101	106					
	AT3G48750	42.268552948	0.862623529551	50.830108648	49					
	AT5G51330	11.97953993	0.855681423571	17.12329801	14					
	AT5G20850	11.948959364	0.702879962588	15.18793846	17					
	AT5G05490	11.57761703	1.92960283833	12.3889299	6					
	AT4G37490	8.65024547	0.786385951818	11.52331461	11					
	AT4G21270	6.44660105	1.6116502625	6.99035612	4					
	AT2G31970	6.20142349	1.03357058167	6.20142349	6					
	AT3G54180	5.59239631	0.932066051667	6.43610856	6					
	AT1G08560	5.18091117	0.370065083571	7.49510317	14					
<i>Diamonds EU</i>	AT5G06150	5.1173723	2.55868615	5.1173723	2					
	AT3G24810	5.0762994	1.6920998	5.0762994	3					

Annotations from left to right:

- TAIR db gene identifier
- Sum of CC abstract scores
- Diamonds EU
- CC score ranked abstracts
- Interaction sentences
- Gene regulation
- Keyword Co-occurrence
- Experiment keywords

Krallinger et al., NAR 09

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 <i>Arabidopsis thaliana</i> 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 <i>Arabidopsis</i> cells in late G1 and G2. A comparison of the effects of aphidicolin, oryzalin and olomoucine suggests that in the <i>Arabidopsis</i> 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 <i>Arabidopsis thaliana</i> , 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 <i>Medicago sativa</i> . 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: cell cycle ; interphase ; prophase ; anaphase
Species ambiguity scores: 0.5
Cell cycle scores: 3.4199

Searching the Arabidopsis literature: abstracts (1)

Query flower development

Search Arabidopsis biome Clear

[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 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 maintenance of genome integrity. However, the function of RAD51C in meiosis is not well understood. 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 has no detectable abnormality in mitosis, it is completely male and female sterile. During early meiosis, homologous chromosomes in *atrad51c-1* fail to undergo synapsis and become severely fragmented. In addition, analysis of the *atrad51c-1 atspo11-1* double mutant showed that fragmentation was nearly completely suppressed by the *atspo11-1* mutation, indicating that the fragmentation largely represents a defect in processing double-stranded breaks generated by AtSPO11-1. Fluorescence *in situ* hybridization experiments suggest that homolog juxtaposition might also be abnormal in *atrad51c-1* meiocytes. These results demonstrate that AtRAD51C is essential for normal meiosis and is probably required for homologous synapsis.

Cell cycle terms: [synapsis](#); [mitosis](#); [meiosis I](#); [cell cycle](#); [prophase](#); [meiosis](#); [meiotic prophase I](#); [mitotic cell cycle](#)

Species ambiguity scores: 0.6

Cell cycle scores: 3.04013

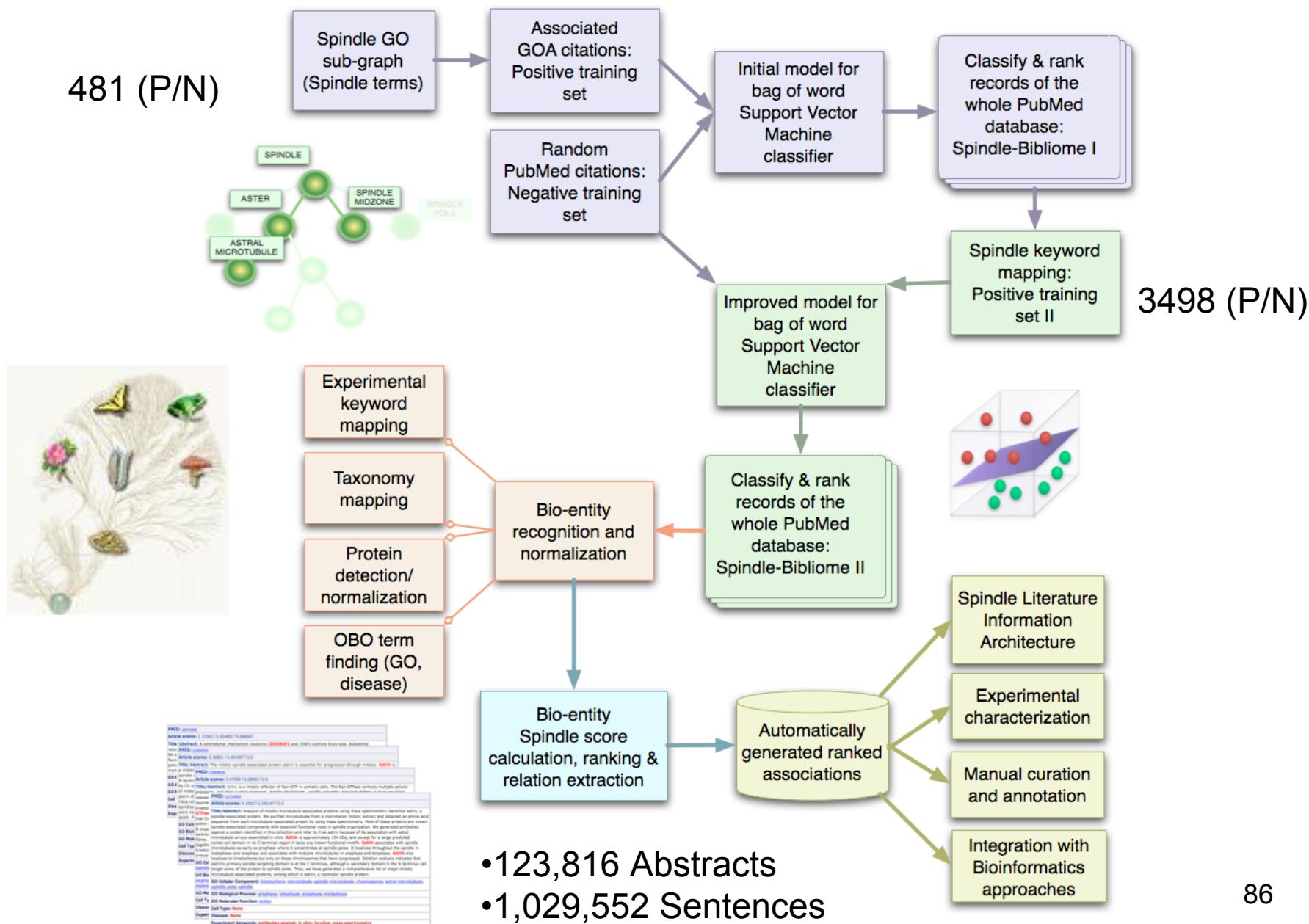
Cell cycle scores: 1.16984

Cell cycle scores: 1.06576

Cell cycle scores: 1.05989

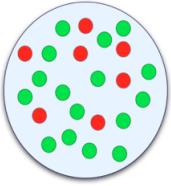
Mitotic spindle relevance protein ranking

481 (P/N)



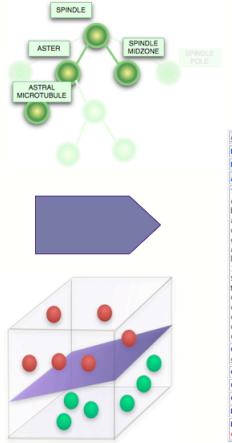
- 123,816 Abstracts
- 1,029,552 Sentences

3. OVERVIEW



- Small collection of known spindle-relevant proteins
- Large set of potentially spindle-relevant proteins (from large scale experiments)
- Experimentally verified new spindle-relevant proteins

Large scale experiments often result in a list of potentially relevant genes or proteins which require further more detailed characterization.



Q96B08_1 SPAG5_HUMAN
Names: **Astin**
PMID: 12256810

Q96CWS_1 GCP3_HUMAN
Names: **KOCPA_HUMAN**
PMID: 10542288

Q96CWS_1 GCP3_HUMAN
Names: **KOCPA_HUMAN**
PMID: 10542288

Article scores: 3.21078 | 0.547917 | 0.5

[Title/Abstract:] The mammalian gamma-tubulin complex contains homologues of the yeast spindle pole body components spc37p and spc80p. gamma-Tubulin is a universal component of microtubule organizing centers where it is believed to play an important role in the nucleation of microtubule polymerization. gamma-Tubulin also exists as part of a cytoplasmic complex whose size and complexity varies in different organisms. To investigate the composition of the gamma-tubulin complex in humans, we have used immunoprecipitation and sequencing to identify the components of human gamma-tubulin were made. The epitope-tagged gamma-tubulin expressed in these cells localize to the centrosome and are incorporated into the cytoplasmic gamma-tubulin complex. Immunoprecipitation of this complex identifies at least 10 proteins, including the gamma-tubulin itself. We have identified two additional proteins which colocalize with the 100- and 101-kD components of the gamma-tubulin complex as homologues of the yeast spindle pole body proteins Spc37p and Spc80p, and named the corresponding human proteins hGCP2 and hGCP3. Sequence analysis revealed that hGCP2 and hGCP3 are homologous to the yeast proteins Spc37p and Spc80p, respectively. Both hGCP2 and hGCP3 colocalize with gamma-tubulin at the centrosome, cosediment with gamma-tubulin in sucrose gradients, and coimmunoprecipitate with gamma-tubulin. Thus, they are part of the gamma-tubulin complex. The conservation of the gamma-tubulin complex between yeast and the animal centrosome share a common molecular mechanism for microtubule nucleation.

GO Cellular Component: cytoplasm; microtubule; gamma-tubulin complex; spindle pole; spindle; spindle pole body; microtubule organizing center; tubulin complex

GO Biological Process: microtubule nucleation; microtubule polymerization

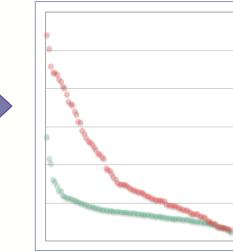
GO Molecular Function: None

Cell Type: None

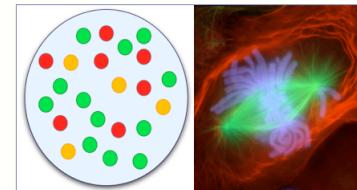
Disease: None

Experiment keywords: Immunoprecipitation; precipitate; coimmunoprecipitate; localize; colocalize; cosediment; immunoprecipitate; sucrose gradient; epitope-tag

Using supervised learning techniques, such as Support Vector Machines it is possible to automatically categorize large collections of textual data. Combined with protein normalization and keyword mapping, it is possible to both, rank proteins in terms of their contextual association to spindle poles as well as provide ranked pointers to the corresponding textual evidences.



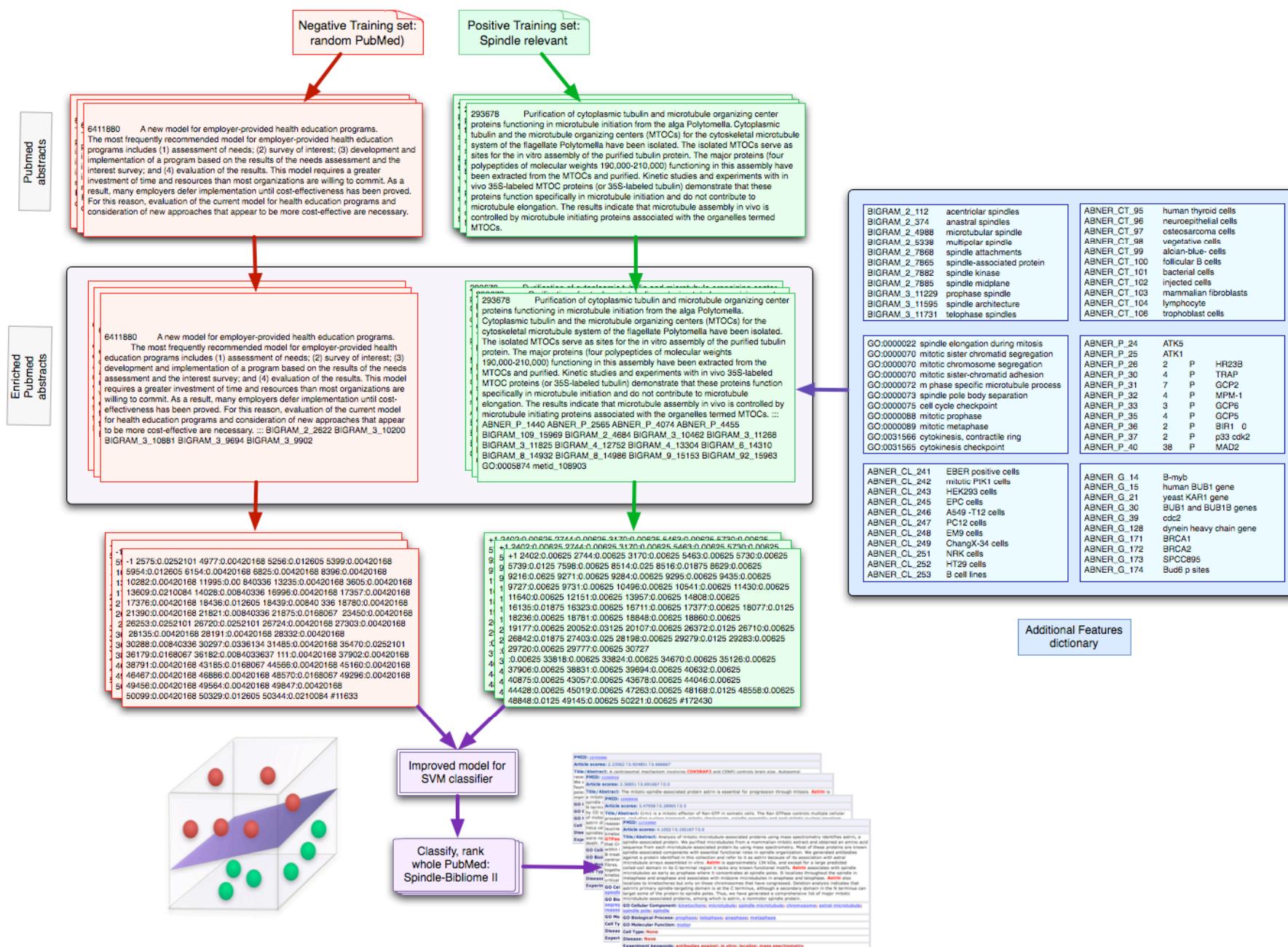
Based on the average document spindle scores
Proteins can be ranked according to their spindle associations



The prioritized proteins can then be
Experimentally characterized in more detail.

Type	Accession	UniProt-Id	Score
K	Q9BVA0	KTNB1_HUMAN	4.58482
K	Q75330	HMMR_HUMAN	3.31712
K	Q96R06	SPAG5_HUMAN	3.23685
K	Q96CWS	GCP3_HUMAN	2.77059
K	Q14980	NUMA1_HUMAN	2.63128
K	Q9UGI1	GCP4_HUMAN	2.55049
K	Q9BSI2	GCP2_HUMAN	2.54923
K	Q7Z460	CLAP1_HUMAN	2.52104
P	Q9BZE4	NOG1_HUMAN	2.3877
K	Q96SN8	CK5P2_HUMAN	2.23562
K	Q95L13	PCNT_HUMAN	2.02624
K	Q13257	MD2L1_HUMAN	2.02318
K	Q9NQ57	INCE_HUMAN	1.92961
K	Q43264	ZW10_HUMAN	1.90836
K	Q02224	CENPE_HUMAN	1.89177
K	P46060	RGP1_HUMAN	1.86112
P	Q95347	SMC2_HUMAN	1.85301
P	Q9NT13	SMC4_HUMAN	1.8333
K	Q14965	STK6_HUMAN	1.8002
K	Q43683	BUB1_HUMAN	1.79463
K	Q95229	ZWINT_HUMAN	1.75663
K	Q9H410	CT172_HUMAN	1.74528
K	P53350	PLK1_HUMAN	1.66669
K	P07900	HS90A_HUMAN	1.64744
K	Q43684	BUB3_HUMAN	1.64454
K	Q15021	CND1_HUMAN	1.61889
P	P39023	RL3_HUMAN	1.60029
P	Q14137	BOP1_HUMAN	1.60029
K	Q60566	BUB1B_HUMAN	1.59503

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,,



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

TAGGING BIO-ENTITIES IN TEXT

- Aim: Identify biological entities in articles and to link 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**

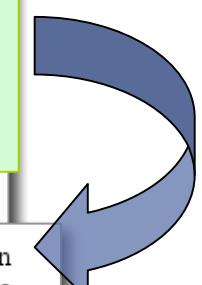
Choose Classes

Virus
 Tissue
 RNA
 Protein
 Polynucleotide
 Peptide
 OtherOrganicCompound
 OtherName
 OtherArtificialSource
 Organism
 Nucleotide
 MultiCell
 MonoCell
 Lipid
 Inorganic
 DNA
 CellType
 CellLine
 CellComponent
 Carbohydrate
 BodyPart
 Atom
 AminoAcidMonomer

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 mutations.

- 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** 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 mutations.

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

PMID	PLAN2L sentence results															
15377755	The plant-specific cyclin-dependent kinase CDKB1;1 and transcription factor E2Fa-DPa control the balance of mitotically dividing and endoreduplicating cells in Arabidopsis.	3.38874	1.15593	-0.916281	-0.354371	-0.117159	-0.536632	-1.06705	-1.17651							
11457971	CdkB1 mRNA accumulates through S until M phase and its associated kinase activity peaks at the G2/M boundary, confirming that transcription of PPTALRE CDKs is cell cycle regulated.	3.12938	1.21085	0.306721	-0.0958637	-0.798032	-0.713507	-0.987208	-0.858957							
15863515	KRP2 phosphorylation by the mitotic cell cycle-specific CDKB1;1 kinase suggests a mechanism in which CDKB1;1 controls the level of CDKA;1 activity through regulating KRP2 protein abundance.	2.64175	-0.729148	0.0441319	-0.354698	-0.990122	-0.686953	-1.14144	-1.22541							
15377755	Plan-specific cyclin-dependent kinase CDKB1;1 and transcription factor E2Fa-DPa control the balance of mitotically dividing and endoreduplicating cells in Arabidopsis.	16055635	E2FB AT5G22220	CDKB1;1 AT3G54180	Activation	Regulation association evidence sentence									Source	
11457971	CdkB1 mRNA accumulates through S until M phase and its associated kinase activity peaks at the G2/M boundary, confirming that transcription of PPTALRE CDKs is cell cycle regulated.	16055635	E2FA At2g36010	CDKB1;1 AT3G54180		Indeed, our recent data show that E2FB can directly induce the promoter of the Arabidopsis CDKB1;1 gene (Z. Magyar, unpublished results).	0	None						Source		
		16055635	E2FA At2g36010	CDKB1;1 AT3G54180	Activation	It is not clear how E2FA could promote the expression of CDKB1;1 , which has a separate	0	None						Source		
		15863515	CDKB1;1 AT3G54180	CDKA;1 AT3G48750	Protein A	Protein B		Protein interaction evidence sentence		Experimental evidence						Source
		15863515	KRP2 AT3G50630	CDKB1;1 AT3G54180	Protein A	Protein B		KRP2 phosphorylation by the mitotic cell cycle-specific CDKB1;1 kinase suggests a mechanism in which CDKB1;1 controls the level of CDKA;1 activity through regulating KRP2 protein abundance.							Source	
		15863515	KRP2 AT3G50630	CDKB1;1 AT3G54180	Protein A	Protein B		KRP2 phosphorylation by the mitotic cell cycle-specific CDKB1;1 kinase suggests a mechanism in which CDKB1;1 controls the level of CDKA;1 activity through regulating KRP2 protein abundance.							Source	

<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.

CDKB1;1:Arabidopsis homolog of yeast cdc2, a protein kinase (cyclin-dependent kinase) that plays a central role in control of the mitotic cell cycle.

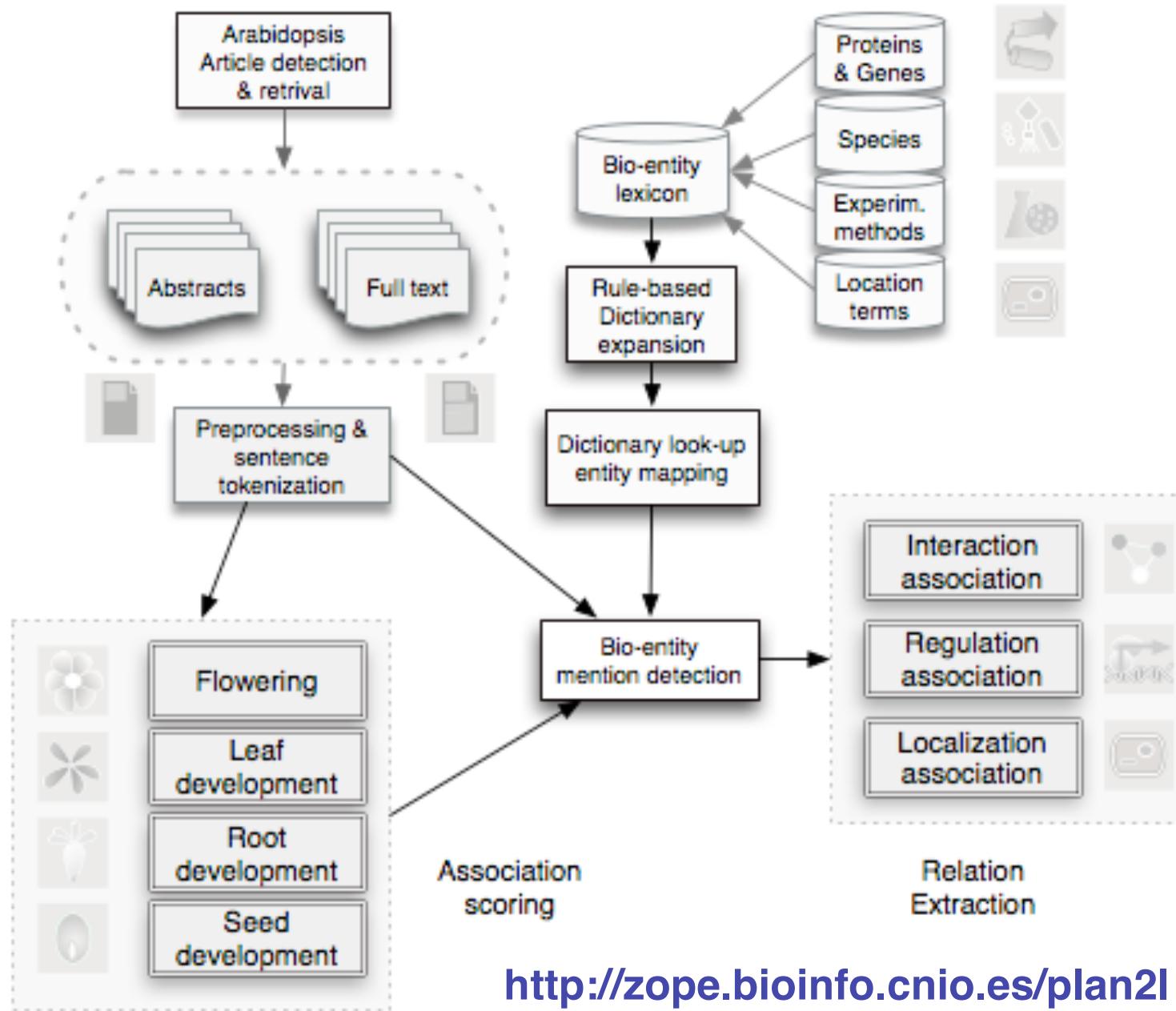
PLAN2L

PMID	Regulator	Regulated	Type	Regulation association evidence sentence							Source
16055635	E2FB AT5G22220	CDKB1;1 AT3G54180	Activation	Indeed, our recent data show that E2FB can directly induce the promoter of the Arabidopsis CDKB1;1 gene (Z. Magyar, unpublished results).		None					
16055635	E2FA At2g36010	CDKB1;1 AT3G54180	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).		None					

PMID	Proteins	Location description evidence sentence	Location terms	Location words
11058164	LEUNIG	The nuclear localization of LEUNIG -GFP is consistent with a role of LEUNIG as a transcriptional regulator.	nuclear	localization

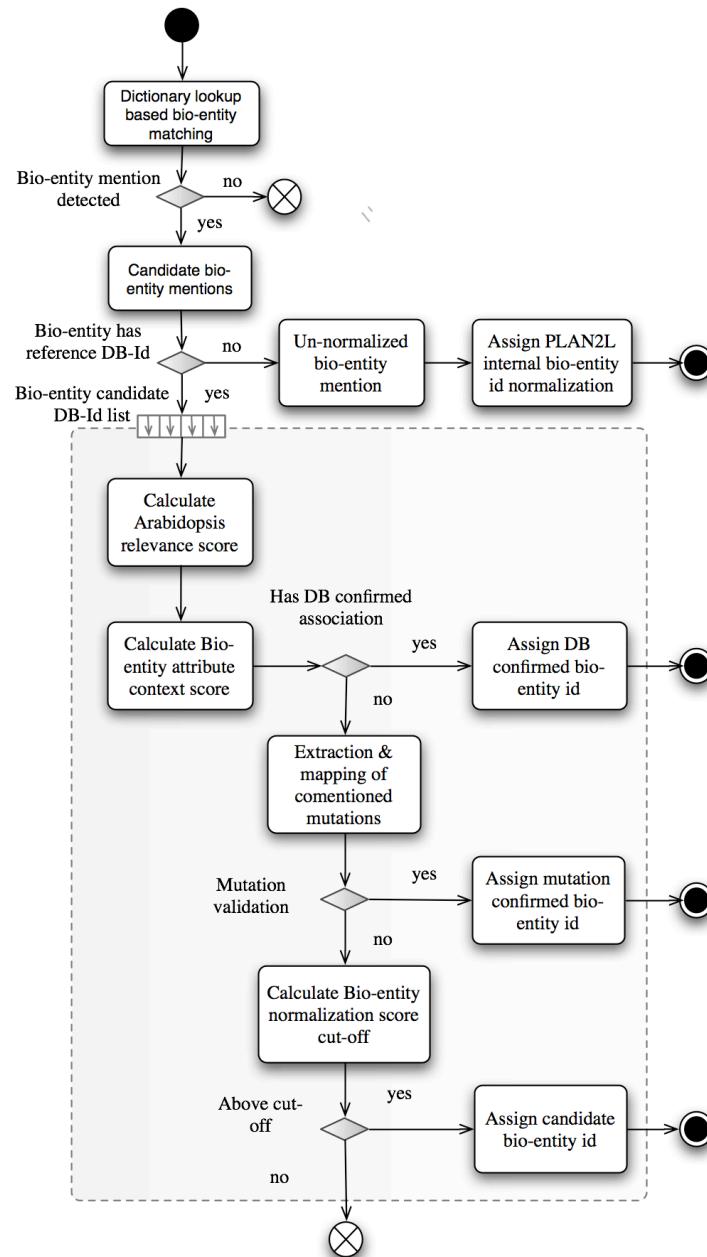
PMID	PLAN2L Protein association evidence											
11041883	Additional analyses indicate that the absence of marginal tissues in leunigaintegumenta double mutants is not mediated by ectopic AGAMOUS .											
16625397	SEUSS and LEUNIG encode components of a putative transcriptional regulatory complex that controls organ identity specification through the repression of the floral organ identity gene AGAMOUS .											
11782418	The effects of seuss mutations are most striking when combined with mutations in LEUNIG , a previously identified repressor of AGAMOUS .											
11058164	LEUNIG , a putative transcriptional corepressor that regulates AGAMOUS expression during flower development.											
16854969	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.											

PLAN2L flowchart

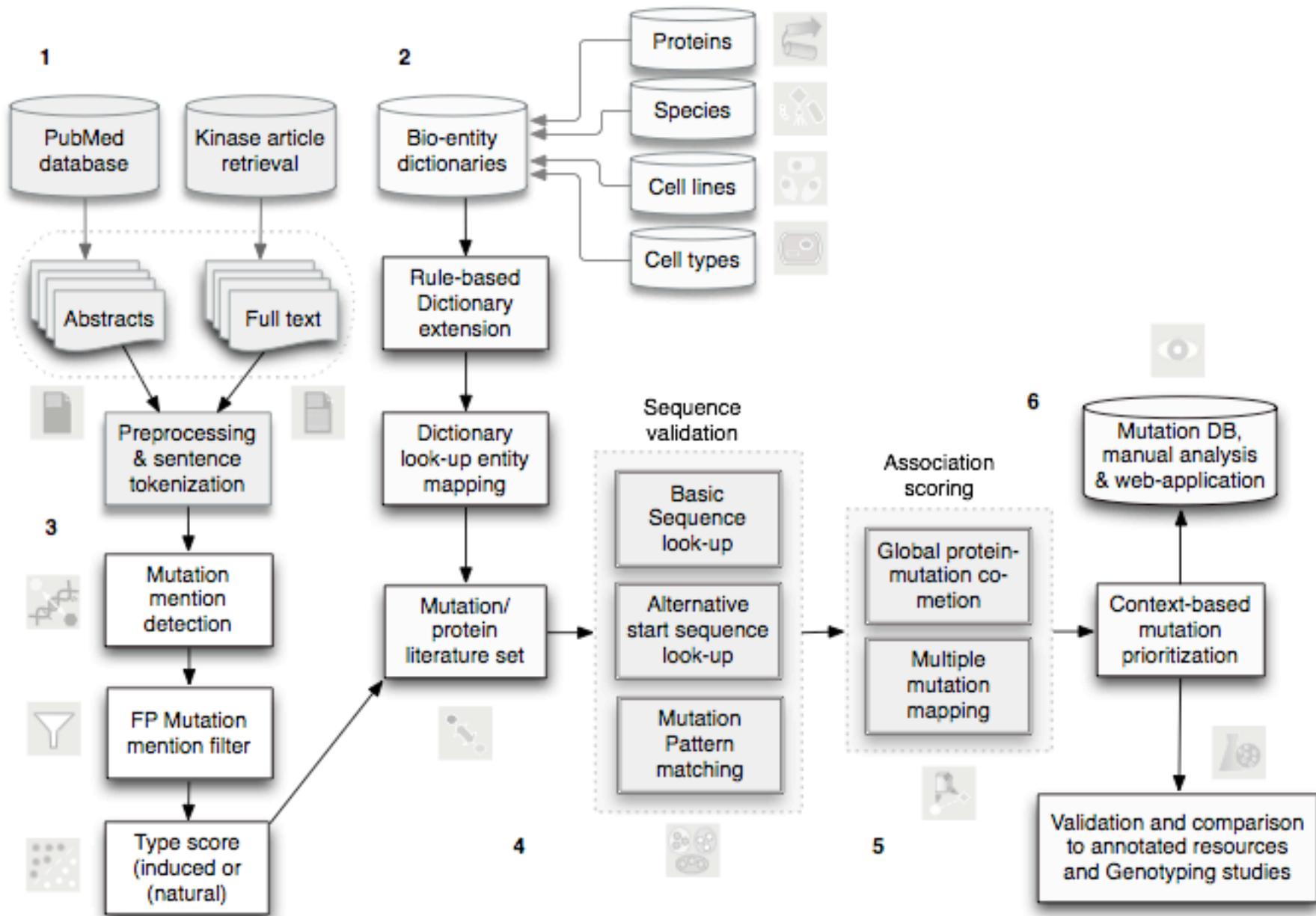


<http://zope.bioinfo.cnio.es/plan2l>

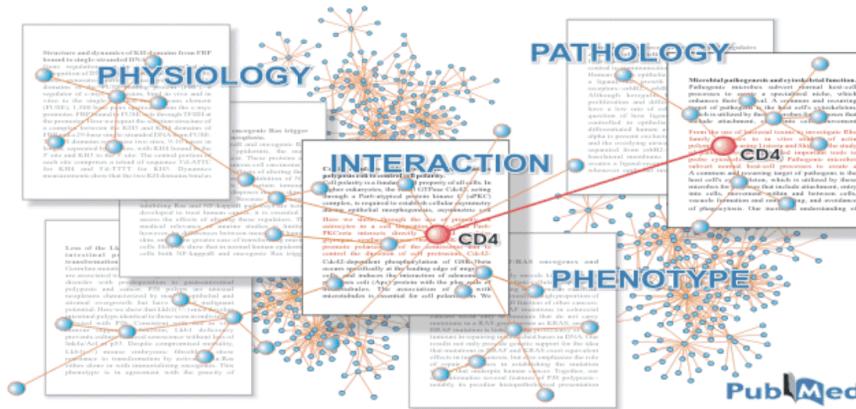
PLAN2L protein mention normalization



PLAN2L mutation extraction



iHOP system



iHOP
Information hyperlinked over proteins

Search Gene

Show overview Find in this Page

Filter and options

Gene Model

Developer's Zone

Help

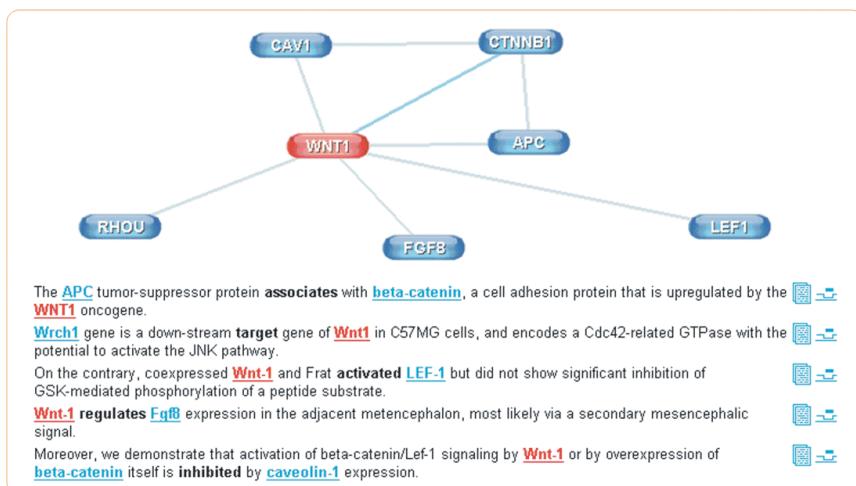
Concept & Implementation by Robert Hoffmann

Symbol	Name	Synonyms	Organism
WNT1	Wnt-1 proto-oncogene protein precursor	INT1	Homo sapiens
UniProt	P39428		
OMIM	164020		
NCBI Gene	7471		
NCBI RefSeq	NP_005421		
NCBI Accession	CAA28874, X03072		
Homologues of WNT1 ...	new		
Definitions for WNT1 ...	new		
Enhanced PubMed/Google query ...	new		

WARNING: Please keep in mind that gene deletion is done automatically and can exhibit a certain error. [Read more](#).

Find in this Page

However, mAkt could act synergistically with **Wnt-1** or Frat to activate **LEF-1**.
Beta-catenin: a common target for the regulation of cell adhesion by **Wnt-1** and Src signaling pathways.
Wnt-1 regulates **Fgf** expression in the adjacent metencephalon, most likely via a secondary mesencephalic signal.
Cultured cells transfected with a membrane-tethered form of **Wnt-1** bind epitope-tagged **Frab-1** in the 10(-10) M range.
In mammalian cells, **Axin** inhibits **Wnt-1** stimulation of beta-catenin/lymphoid enhancer factor-1-dependent transcription.
Furthermore, **beta-catenin** is the target of two signal transduction pathways mediated by the proto-oncogenes **src** and **wnt-1**.
Ectopic expression of **Wnt-1** in ST3-L1 preadipocytes stabilizes **beta-catenin**, activates TCF-dependent gene transcription, and blocks adipogenesis.
Wrfch1 gene is a down-stream target gene of **Wnt-1** in C57MG cells, and encodes a Cdc42-related GTPase with the potential to activate the JNK pathway.
Wnt-1 induces morphological transformation of C57MG mammary epithelial cells and accumulation of cytosolic **beta-catenin** whereas **Wnt-5a** has no effect.
On the contrary, coexpressed **Wnt-1** and **Frat** activated **LEF-1** but did not show significant inhibition of GSK-mediated phosphorylation of a peptide substrate.
The specificity of the approach enabled us to identify an **Max-1** consensus **DNA** site within the transcriptional control region of the developmental regulatory gene **Wnt-1**.
Fra2 efficiently inhibited the **Wnt-1** mediated increase in cytoplasmic (beta)-catenin levels as well as the **Wnt-1** induction of transcription from a **Lef/tcf** reporter gene.
Furthermore, a similar phenotype is not observed in **Wnt1/RCas**-infected **brain**, demonstrating that ectopic **Wnt1** is insufficient to mediate the effect of **Lmx1b**.



iHOP
Information hyperlinked over proteins

Search Gene

Show overview Find in this Page

Filter and options

Gene Model

Developer's Zone

Help

Concept & Implementation by Robert Hoffmann

Symbol	Name	Synonyms	Organism
LEF1	Lymphoid enhancer binding factor 1	LEF-1, lymphoid enhancer-binding factor 1-alpha, TCF1ALPHA, TCF1-alpha	Homo sapiens
UniProt	Q9HAZ0, Q9UJU2		
OMIM	153245		
NCBI Gene	5179		
NCBI RefSeq	NP_057363		
NCBI Accession	AAF13268, AAG01022, AAG26866		
Homologues of LEF1 ...	new		
Definitions for LEF1 ...	new		
Enhanced PubMed/Google query ...	new		

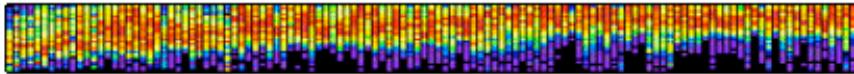
WARNING: Please keep in mind that gene deletion is done automatically and can exhibit a certain error. [Read more](#).

Find in this Page

However, mAkt could act synergistically with **Wnt-1** or Frat to activate **LEF-1**.
On the contrary, coexpressed **Wnt-1** and **Frat** activated **LEF-1** but did not show significant inhibition of GSK-mediated phosphorylation of a peptide substrate.
Addition of **Wnt-1** to normal **epithelial** cell lines stabilizes cytosolic **beta-catenin** that **LEF-1** then transports to nuclei, causing a small amount of EMT.
Here we study the mechanism of transcriptional regulation by **LEF-1** in response to a **Wnt-1** signal under conditions of endogenous **beta-catenin** in NIH 3T3 cells, and we examine whether association with **beta-catenin** is obligatory for the function of **LEF-1**.

In **Wnt-1**-transfected C57MG cells, free **beta-catenin** accumulated and was able to associate with **LEF-1**.
Beta-catenin forms complexes with **Tcf** and **Lef-1** and functions as a transcriptional activator in the **Wnt** signalling pathway.
Thus, the apoptotic effects of overexpressed exogenous **beta-catenin** do not rely on its transactivating function with nuclear **LEF-1**.
Beta-catenin forms complexes with **Tcf** and **Lef-1** and functions as a transcriptional activator downstream of the **Wnt** signalling pathway.
NICD stimulation of **LEF-1** activity was context dependent and occurred on a subset of promoters distinct from those activated by **beta-catenin**.
The **Wnt**-responsive **transcription factor** **LEF1** can activate transcription in association with **beta-catenin** and repress transcription in association with **Groucho**.
Among others, **LEF-1** regulates expression of cytokeratin genes involved in formation of hair follicles and the gene encoding the cell-adhesion molecule **Ecadherin**.

iHOP system: query to DB record

Symbol	Name	Synonym/ DB-reference	Organism	Results
				Life cycles of successful genes
BRCA2	breast cancer 2, early onset		Homo sapiens	
HMG20B	high-mobility group 20B	BRCA2-associated factor 35	Homo sapiens	
BCCIP	BRCA2 and CDKN1A interacting protein		Homo sapiens	

more than 1,500 organisms. 80,000 genes. 12 million sentences.
...always up-to-date.

Search for a gene **synonym** or **accession number**...

[SEARCH]

Results
options



iHOP system: Defining information

Find in this Page

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 relevance

PALB2, which encodes a **BRCA2**★-interacting protein, is a [breast cancer](#) susceptibility gene. [2007]

Inheritance of one defective **BRCA2**★ allele predisposes humans to [breast cancer](#). [2001]

A common variant in **BRCA2**★ is associated with both [breast cancer](#) risk and prenatal viability. [2000]

Inherited mutations in the gene **BRCA2**★ predispose carriers to early onset [breast cancer](#), but such mutations account for fewer than 2% of all cases in East Anglia. [2000]

Mutations in **BRCA2**★ are thought to account for as much as 35% of all inherited [breast cancer](#) as well as a proportion of inherited [ovarian cancer](#). [1996]

Two of the five **BRCA2**★ mutation carriers reported a family history of [breast cancer](#), and none reported a family history of [ovarian cancer](#). [2002]

Our results indicate that **BRCA2**★ confers a very high risk of [breast cancer](#) and is responsible for a substantial fraction of breast and [ovarian cancer](#) in Iceland, but only a small proportion of other cancers. [1996]

Recent studies have identified mutations in the breast and ([ovarian cancer](#) susceptibility gene 2 (**BRCA2**★)), one which has been found in the germline of several males and one female affected with [breast cancer](#). [1996]

The [breast cancer](#) susceptibility gene **BRCA2**★ on [chromosome](#) 13q12-13 has recently been identified. [1997]

The [breast cancer](#) susceptibility gene, **BRCA2**★ on [chromosome](#) 13q12-13, was recently isolated. [1996]

The **BRCA2**★ gene on [chromosome](#) 13 has been shown to be associated with familial male and female [breast cancer](#). [1996]

Colour legend	Main gene
	Associated genes
	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 - [Read more](#).

Show all ▾

Order by relevance ▾

For a summary overview of the information in this page [click here](#). new

RESULTS: Definite **BRCA2**★ mutations were found in 2 of the 73 women with early-onset **breast cancer** (2.7 percent; 95 percent confidence interval, 0.4 to 9.6 percent), suggesting that **BRCA2**★ is **associated** with fewer cases than **BRCA1**★ ($P=0.03$). [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 **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 **breast tumors** and compared them with our existing 36 **BRCA1**★ 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

Sentences in this view contain the most recent information on BRCA2 - Most recent information is 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 relevance

Mutations in the **BRCA2** interacting **DSS1** are not a [risk factor](#) for [male breast cancer](#). [2007]



Constitutive activation of [MAPK \[?\]](#)/[ERK \[?\]](#) inhibits [prostate cancer cell proliferation](#) through [upregulation](#) of **BRCA2**. [2007]



BRCA2 is central to an utterly diverse biological behavior elicited after [integrin-mediated](#) normal and [prostate cancer cell adhesion](#) to [basement membrane](#) (BM) and [extracellular matrix](#) (ECM) proteins. [2007]



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]



PNT1A [cell adhesion](#) to [ECM](#) triggered [MAPK \[?\]](#)/[ERK \[?\]](#) signaling resulting in [upregulation](#) of **BRCA2** mRNA and protein, with negligible effects upon [cell proliferation](#). [2007]



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]



METHODS: 277 families with pathogenic **BRCA1**/**BRCA2** mutations were reviewed and 28 [breast cancer](#) phenocopies identified. [2007]



FINDINGS: Questionnaires were completed by 799 women with a history of invasive [ovarian cancer](#) (670 with **BRCA1** mutations, 128 with **BRCA2** mutations, and one with a mutation in both genes), and controls were 2424 women without [ovarian cancer](#) (2043 with **BRCA1** mutations, 380 with **BRCA2** mutations, and one with a mutation in both genes). [2007]



Contribution of **BRCA1** and **BRCA2** [germline mutations](#) to the incidence of early-onset [breast cancer](#) in Cyprus. [2007]



The [Fanconi anemia](#) and **BRCA** networks are considered interconnected, as **BRCA2** gene defects have been discovered in individuals with [Fanconi anemia](#) subtype D1. [2007]



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]



METHODS: We did a [matched case-control study](#) in women who were found to carry a pathogenetic mutation in **BRCA1** or **BRCA2**. [2007]





iHOP system: gene model/ graph

Symbol	Name	Synonyms	Organism
BRCA2	breast cancer 2, early onset	BRCC2, Breast cancer type 2 susceptibility protein, FACD, FAD, FAD1, FANCB, FANCD, FANCD1, Fanconi anemia group D1 protein	Homo sapiens
UniProt	P51587, Q5TBJ7, Q8IU82		
IntAct	P51587		
PDB Structure	1N0W		
OMIM	114480, 155255	more than 1,500 organisms. 80,000 genes. 12 million sentences. ...always up-to-date.	
NCBI Gene	675		
NCBI RefSeq	NP_000050		
NCBI RefSeq	NM_000059		
NCBI UniGene	675		
NCBI Accession	CAA98995, AAQ97181		
Homologues of BRCA2 ...			
Interaction information for BRCA2 ...			
Most recent information for BRCA2 ... new			
Enhanced PubMed/Google query ...			
WARNING: Please keep in mind that gene detection is done automatically and can exhibit confidence value .			

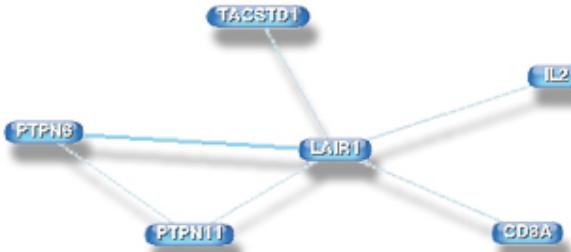
Gene model is a interactive graph where you can add interesting sentences and interactions.

Gene Model - the logbook

In the course of your navigation through iHOP, interesting sentences can be added to the *Gene Model* by clicking on the icon beside the sentence.

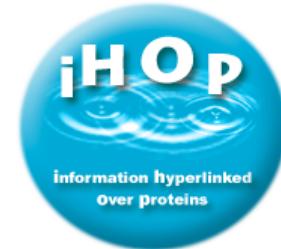
The Gene Model stores these sentences and represents their relation in a graph. [More about the Gene Model...](#)

e.g.





iHOP system: confidence



The synonym 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

The screenshot shows the EBIMed search interface. At the top, the EBIMed logo is displayed, featuring the word "EBIMed" in large, bold letters where the "E" is blue and the "M" is green, with a cartoon character of a person reading books integrated into the letter "e". Below the logo is a search bar with the term "query" followed by a text input containing "BRCA2", a "Search" button, and links for "[Advanced Search]" and "[Query Syntax]". To the right of the search bar is a small icon with a question mark. The main content area contains a descriptive text about EBIMed's functionality, mentioning its integration of Medline, UniProt, GO, Drugs, and Species data. To the right of this text is a graphic titled "Semantic Mining in Biomedicine" showing a network of nodes and connections, and the logo for the "SIXTH FRAMEWORK PROGRAMME".

Term query **Search**

[\[Advanced Search\]](#) [\[Query Syntax\]](#)

EBIMed is a web application that combines Information Retrieval and Extraction from [Medline](#). EBIMed finds Medline abstracts in the same way PubMed does. Then it goes a step beyond and analyses them to offer a complete overview on associations between [UniProt](#) protein/gene names, [GO](#) annotations, [Drugs](#) and [Species](#). The results are shown in a table that displays all the associations and links to the sentences that support them and to the original abstracts.



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

GoPubMed

what

- Top 5 categories
 - Mutation [766]
 - Breast Neoplasms [700]
 - Humans [588]
 - Genes [751]
 - Genes, BRCA2 [487]
- Top categories of GO
 - Anatomy [366]
 - Biological Sciences [9]
 - biological_process [4]
 - cellular_component [1]
 - Chemicals and Drugs
 - Diseases [941]
 - Health Care [796]
 - molecular_function [1]
 - Named Groups [627]
 - Natural Sciences [791]
 - Organisms [641]
 - Psychiatry and Psychotherapy
 - Techniques and Equipment
 - Technology, Industry,
- Found categories of G
- Find related categories
- My last 5 queries
 - BRCA2

who

where

?

BRCA2

find it!

goPubMed

! 1,000 articles

?

Expand your query with synonyms for BRCA2

PubMed has found 3,428 citations for the query BRCA2. The 1,000 latest documents were used by GoPubMed.

Show statistics for these 1,000 articles.

1: Genetic variants and haplotype analyses of the ZBRK1/ZNF350 gene in high-risk non BRCA1/2 French Canadian breast and ovarian cancer families.

PMID: 17764113 Related Articles
Desjardins S, Belleau P, Labrie Y, Ouellette G, Bessette P, Chiquette J, Laframboise R, Lépine J, Lespérance B, Pichette R, Plante M, Durocher F.
Int J Cancer: , 2007

Our current understanding of breast cancer susceptibility involves mutations in the 2 major genes BRCA1 and BRCA2, found in about 25% of high-risk families, as well as few other low penetrance genes such as ATM and CHEK2. Approximately two-thirds of the multiple cases families remain to be explained by mutations in still unknown genes. In a candidate gene approach to identify new genes potentially involved in breast cancer susceptibility, we analyzed genomic variants in the ZBRK1 gene, a co-repressor implicated in BRCA1-mediated repression of GADD45. Direct sequencing of ZBRK1 entire coding region in affected breast cancer individuals from 97 high-risk French Canadian breast/ovarian cancer families and 94 healthy controls led to the identification of 18 genomic variants. Haplotype analyses, using PHASE, COCAPHASE and HaploStats programs, put in evidence 3 specific haplotypes which could potentially modulate breast cancer

BioCreative

The screenshot shows the BioCreative website homepage. At the top, there is a navigation bar with links to various bioinformatics resources like PubMed Home, LiMTox, and different CNIO/BioCreative projects. Below the navigation bar, the main content area features a banner for "Critical Assessment of Information Extraction in Biology". The banner includes a small diagram showing nodes labeled "Bio", "cre", "at", and "ive" connected by lines, and the text "were further divided into those showing active (cells, fibroblasts) or inactive (skeletal muscle) to determine whether IL-2R + cells have a particular association with the IL-2R receptor". Below the banner is a menu bar with five tabs: News, About, Events, Tasks, and Resources. On the left side, there is a sidebar with a "Year" section containing links for 2003, 2004, 2005, 2006, 2007, and 2008, and a "Content" section containing links for BioCreative I, BioCreative II, BioCreative II.5, Organizers, and Publications.

Organizers

Website upgrades [2008-12-03]

The new BioCreative website just has received some extra functionality:

1. Logged in users now can not only change the password, but also their email address.
2. RSS 2.0 Feeds for all major sections (see [here](#)) and a favicon for the webpage.
3. Team registration/management has been added. This is not the official start of the registration process, but you can already register a team for BioCreative II.5 [here](#) if you feel like it - or use the new "team page" link on the top of the page in the user menu.
4. Tons of minor updates and a bugfix have been added.

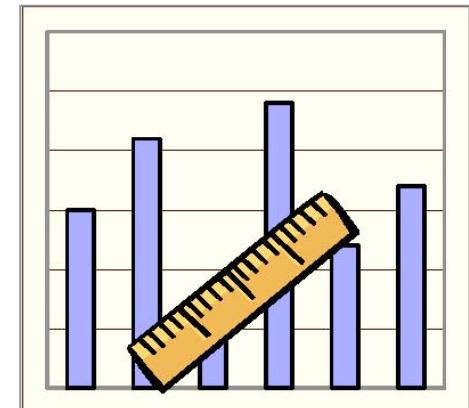
We hope, this makes visiting this site an even better experience!

BioCreative II.5

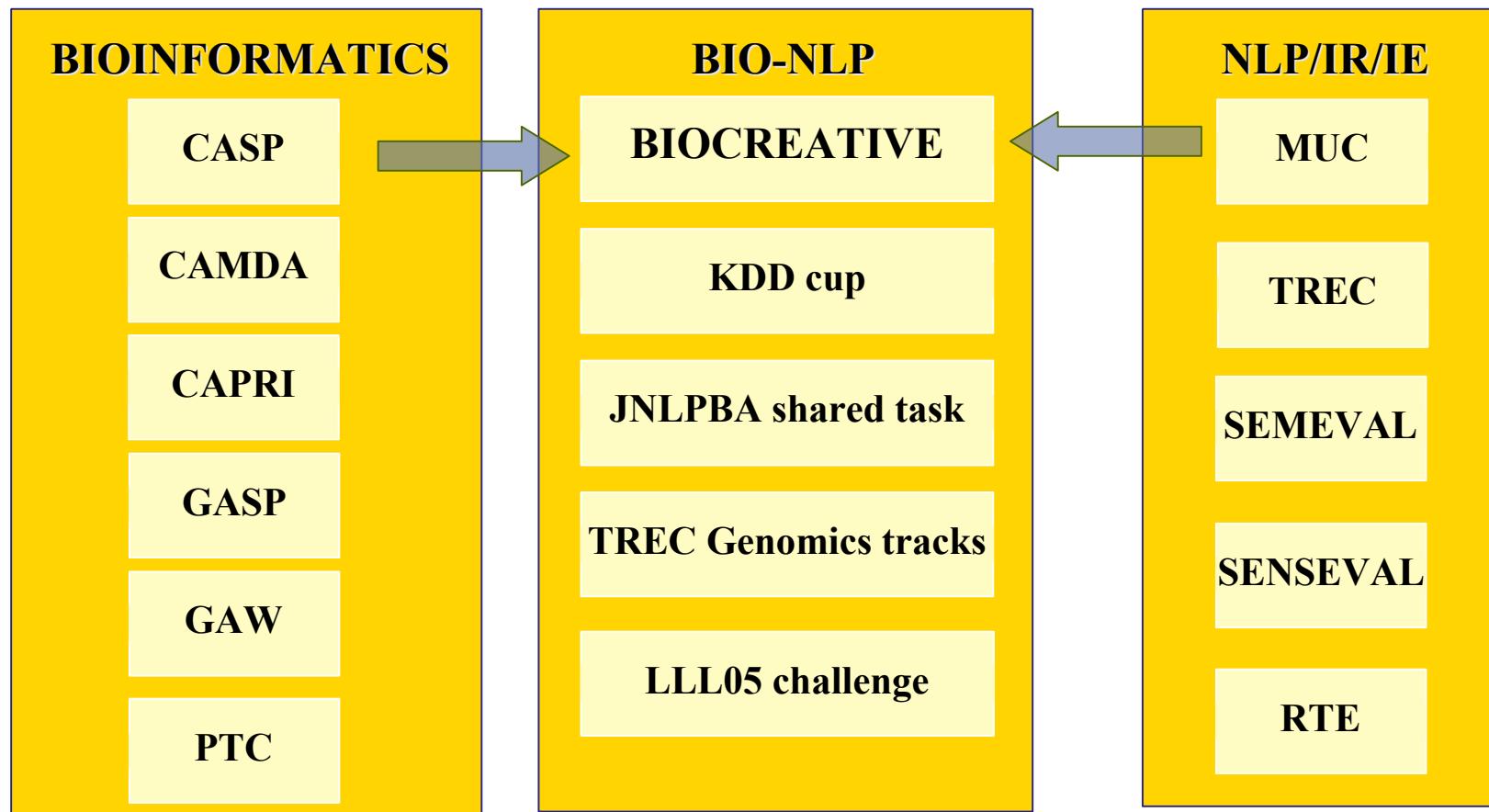
BioCreative II.5 Announcement (Events) [2008-11-18]

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

KDD: Knowledge Discovery and Data mining

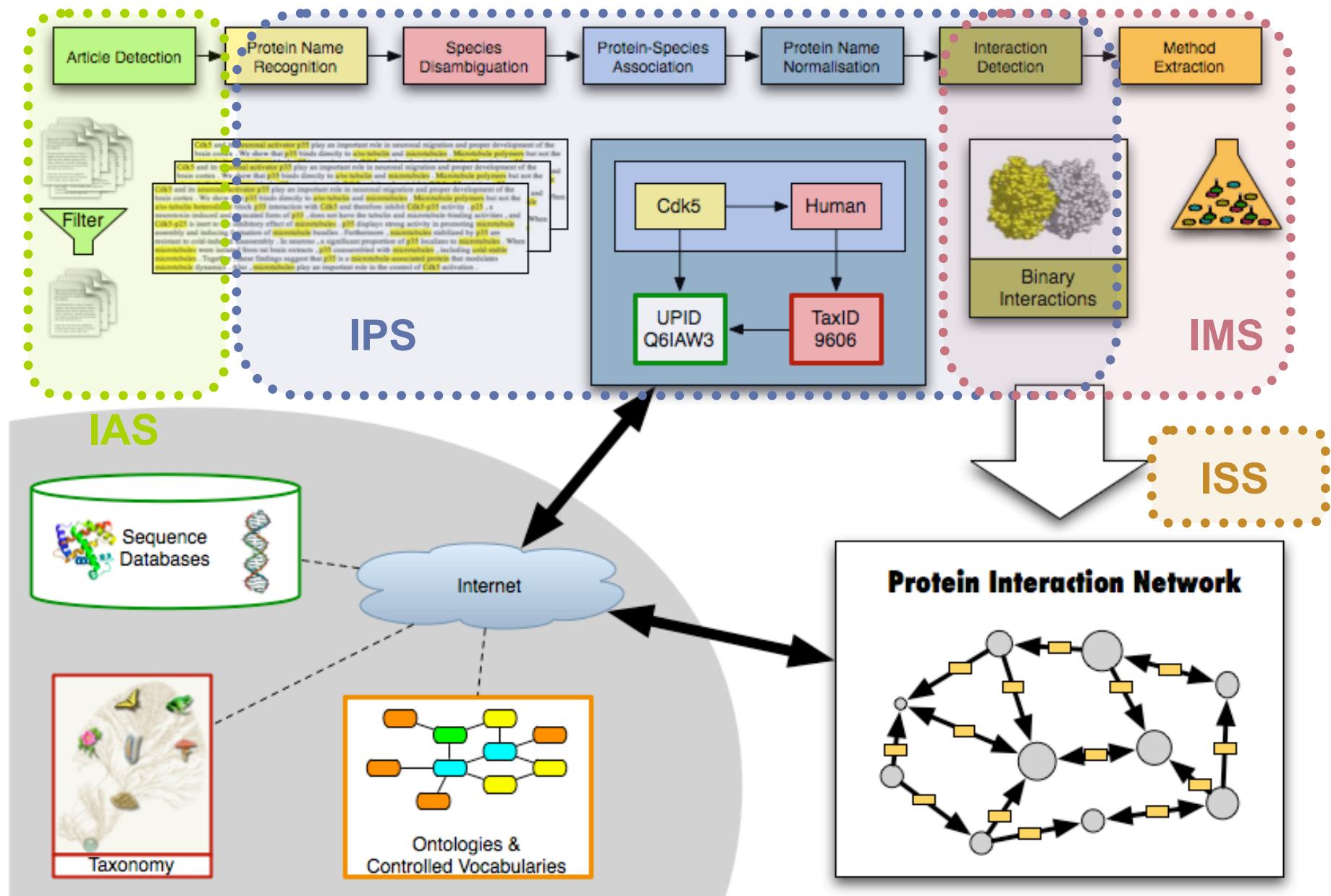
JNLPBA: Joint workshop on Natural Language Processing in Biomedicine

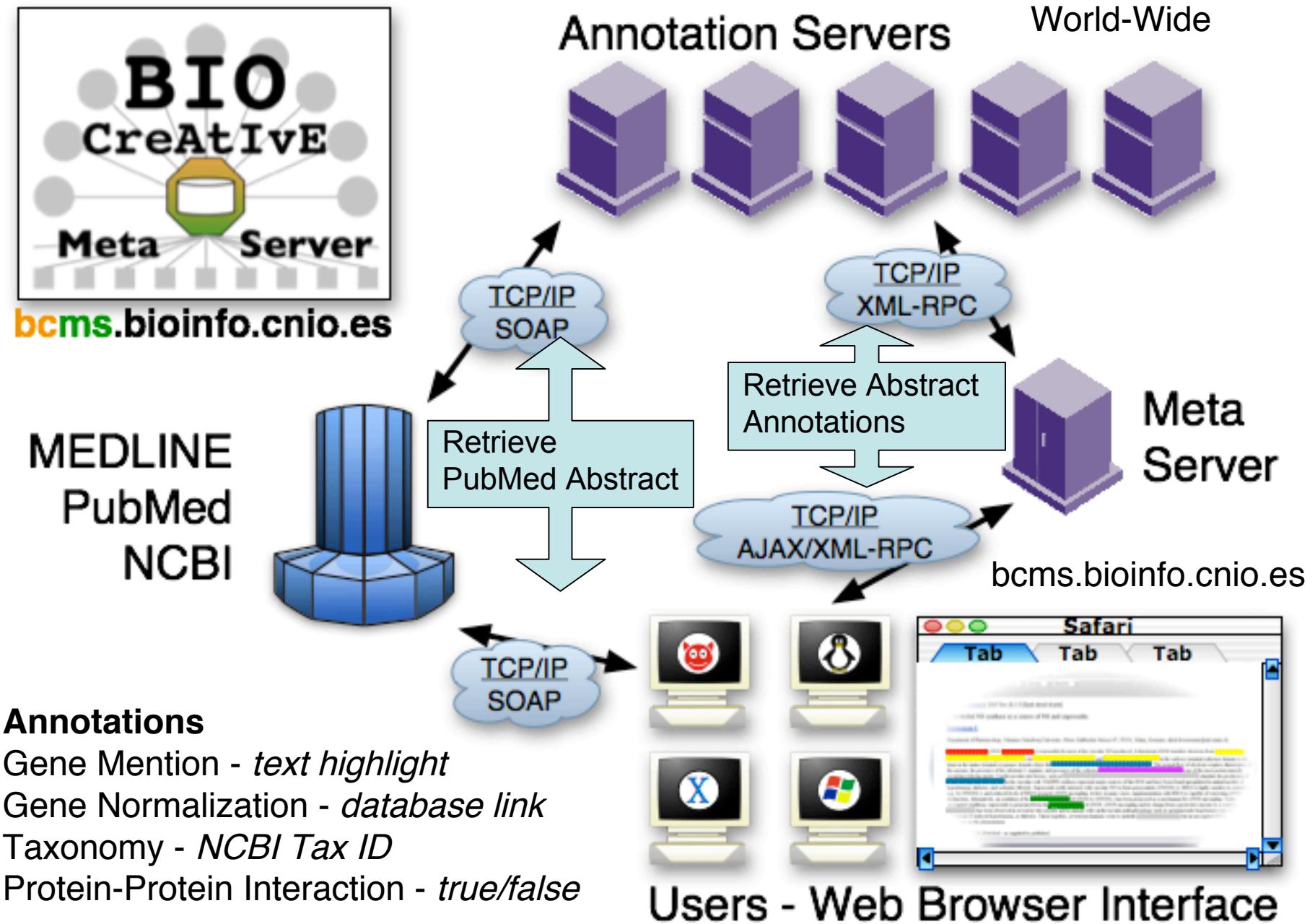
TREC: Text Retrieval conference

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**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.**