



FACULDADE · DE · CIÊNCIAS UNIVERSIDADE · DE · LISBOA

ProFAL: Protein Functional Annotation through Literature

Francisco M. Couto, Mário J. Silva
{fjmc,mjs}@di.fc.ul.pt
Universidade de Lisboa

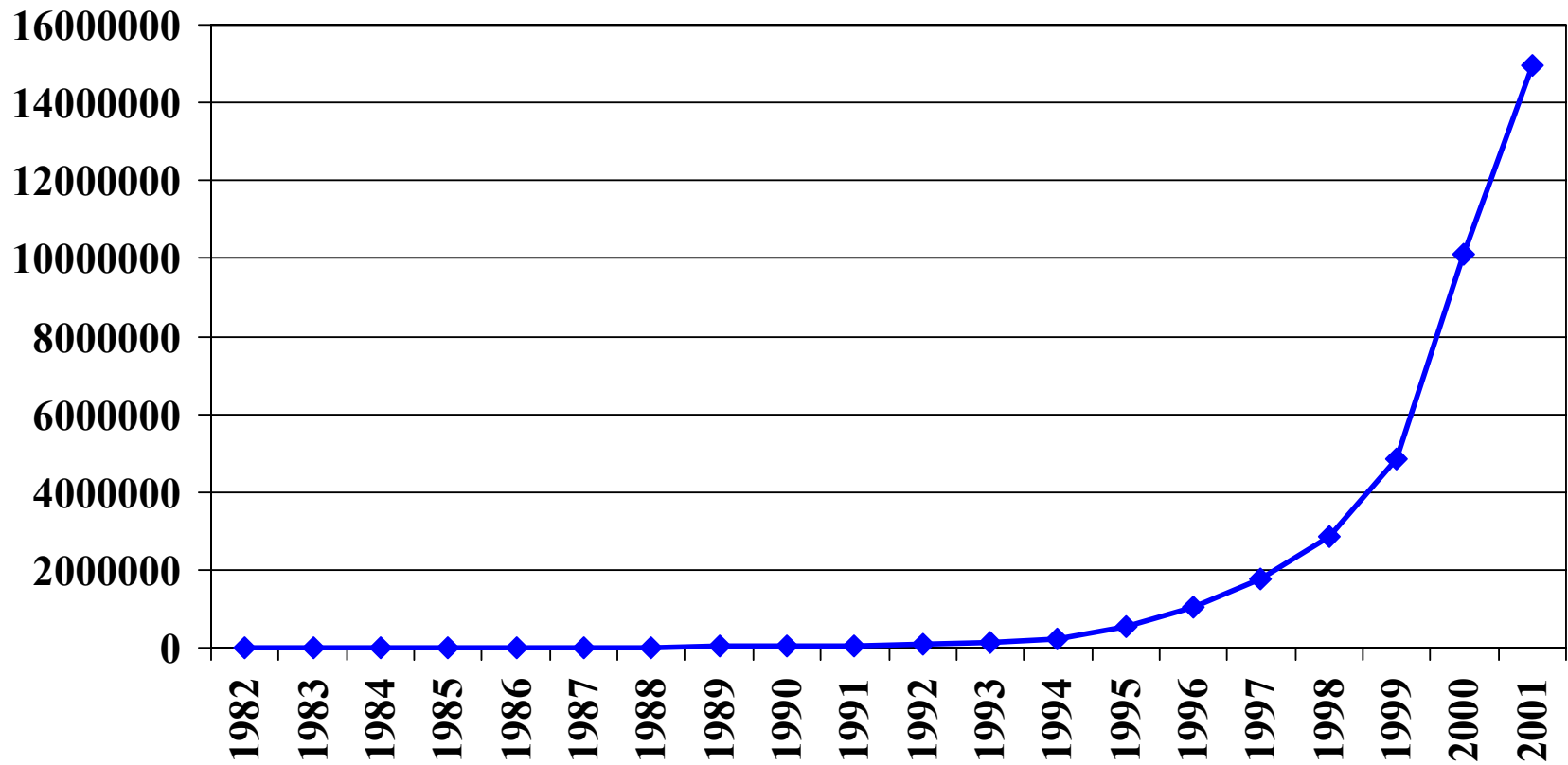
Pedro.Coutinho@afmb.cnrs-mrs.fr
(Université de Provence/CNRS, France)

Outline

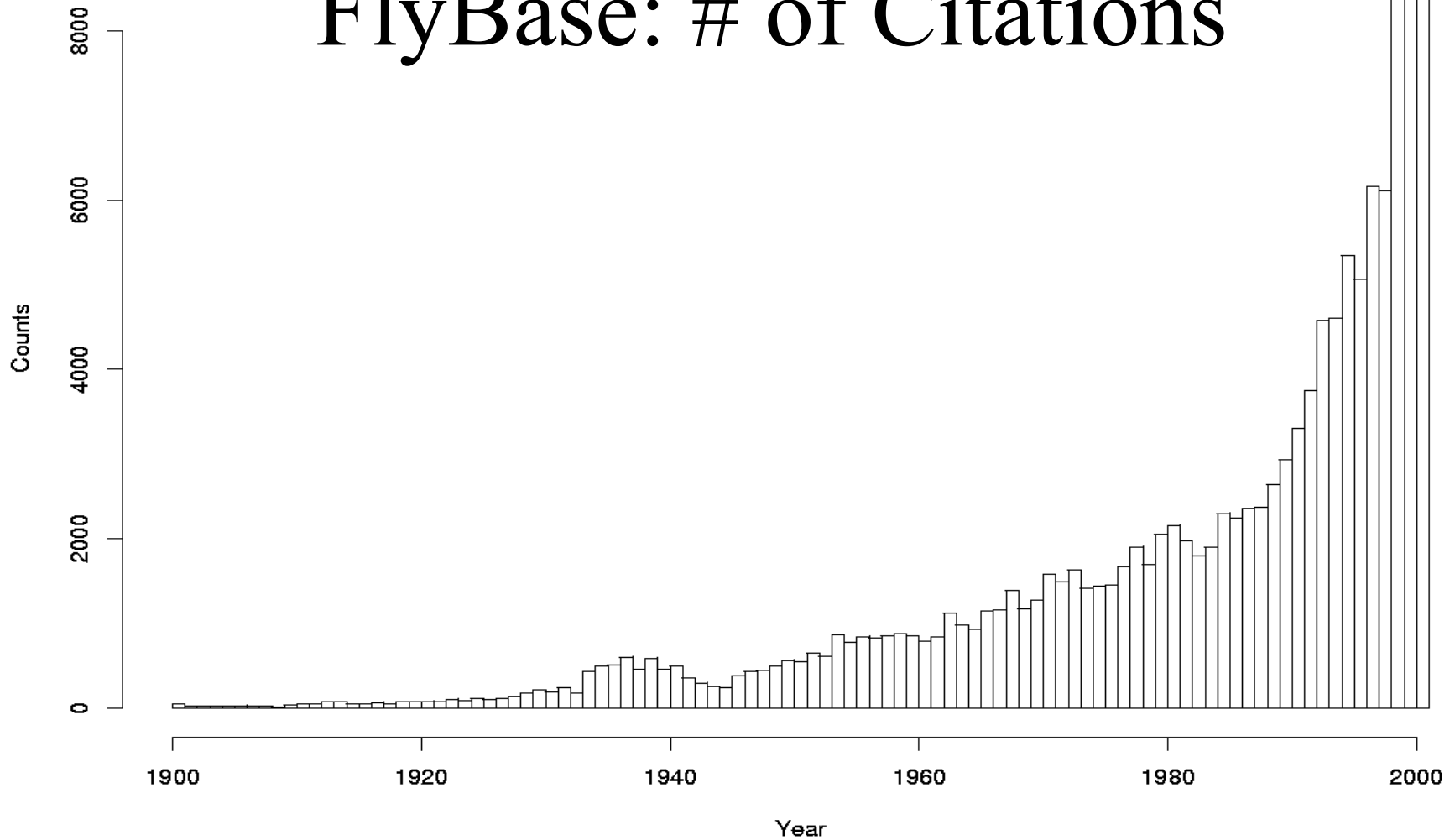
- Motivation
- ProFAL
- Case-Study
- Results
- Conclusions

Genbank

of sequences



FlyBase: # of Citations



Biologic Literature

- MEDLINE
 - “Beginning in 2002 over **2,000** completed references are added daily each Tuesday through Saturday, January through October”
 - “over **460,000** added last year”
 - Reading 10 articles per day, takes **112** years to read those articles





More Information About
 · Bio-ITWorld Conference & Expo
 · IDC Life Sciences Market Data
 · Buyer's Guide Submissions



Ness Technologies

Trusted Advisors
 and Innovative Solutions
 Since 1985

Bio-IT World

Technology for the Life Sciences

SITE SEARCH

Home > Archive > Mar 10, 2003 > Common Knowledge

- Site Search
- News
- Current Issue
- Archives
- Buyer's Guide
- Op-Ed/Newsletters
- Free Subscriptions
- Product Coverage
- Health-IT World
- Info Technology Resource Center
- Life Sciences Resource Center
- Career Center
- Events Calendar

Common Knowledge

Common Knowledge

Two heads (or more) are better than one, except when they don't share information. That's where knowledge management comes in.

BY SALVATORE SALAMONE

What's more difficult than finding a needle in a haystack? Defining knowledge management.

What's in a Name?

A somewhat different approach to finding answers in numerous databases of unstructured data is to use more intelligent search tools. More intelligent tools are needed because common text search tools that find things so well on the Web do not work as well with life science databases, largely because there is no universal, consistent way to denote genes, disease names, or other entities in the scientific literature.

Bio-IT World
Buyer's Guide
 >> For the Life Sciences

ADVERTISEMENT

Buyer's Guide
 is live

Av
 WI
 Ab
 Ed
 Ma
 Me
 Pr
 Co
 ID
 Site

knowledge management solutions for pharmaceuticals and





nature genetics

volume 28 no. 1 may 2001

Nature Publishing Group <http://genetics.nature.com>

Community watch

Microarray analysis yields a wealth of information that presents the geneticist with a vexatious challenge: where does one go with a cluster of genes defined by similar or identical expression levels in different tissues, or during adaptation to a new environment? How does one select out genes that warrant further analysis? Three

PubGene, the GO initiative and other efforts to establish controlled vocabularies are testimony to the power of *in silico* analysis and cast a new light on the commodity that is published text. The extent to which researchers are able to use published text—as a tool—depends on publishing strategy, and poses interesting questions for the publishers of on-line scientific literature. It is one that will become more interesting as literature-mining tools evolve.



The better use of scientific literature

A genetic paperchase

May 3rd 2001

From The Economist print edition

THERE is a difference between information and knowledge. With the completion of the Human Genome Project, biologists have a lot of information about what human genes are. Knowledge of how they work—and in particular of how their products interact to form the vast network of biochemical pathways referred to as “life”—is harder to assemble. But a group of biologists led by Eivind Hovig of the Norwegian Radium Hospital in Oslo has managed to find a way to speed the process up, by linking together the vast quantities of disconnected biochemical information that have already been published.

Select an option below to continue:

Subscribe

Get unlimited site access. Sign up for:

- a month, \$19.95 ([automatic renewal](#))
- a year, \$69 ([automatic renewal](#); save 70%)
- a print subscription (complimentary web access)

Subscribe

Pay Per View

Buy credits to view individual articles:

- 1 article-credit, \$2.95
- 5 article-credits, \$9.95
- Use credits already on my account

Go

Current Subscribers

Log in to view the article.

E-mail address

Password

Outline

- Motivation

- ProFAL

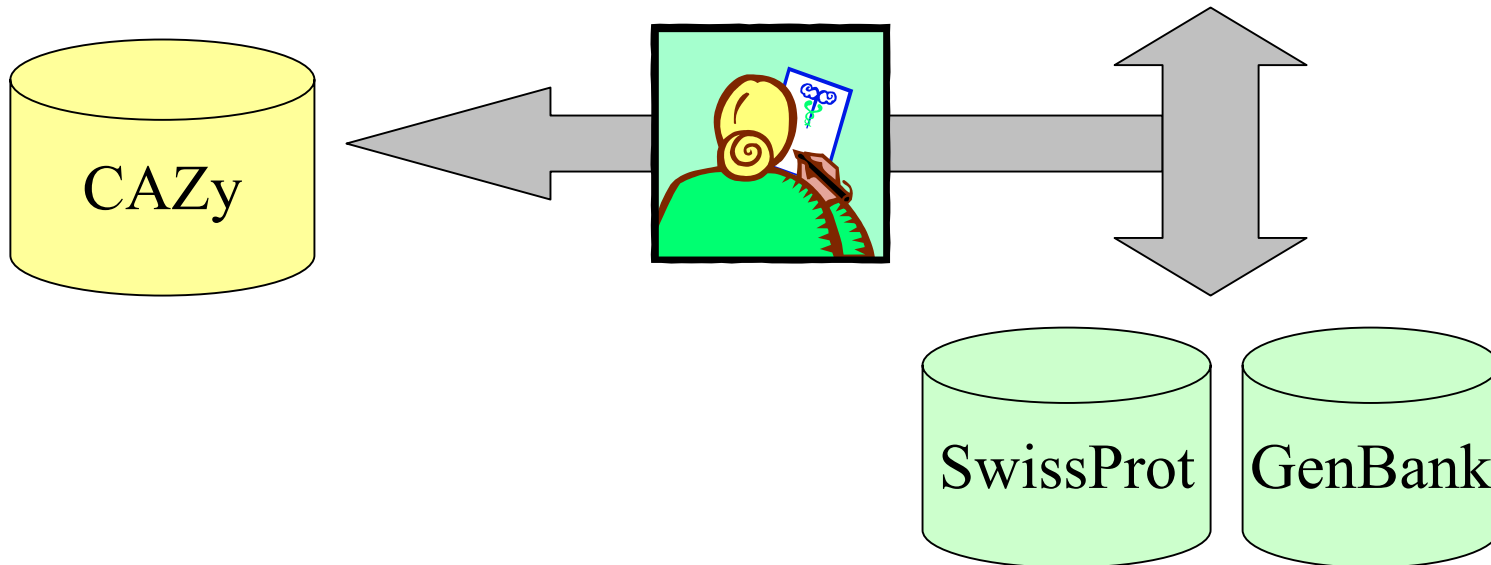
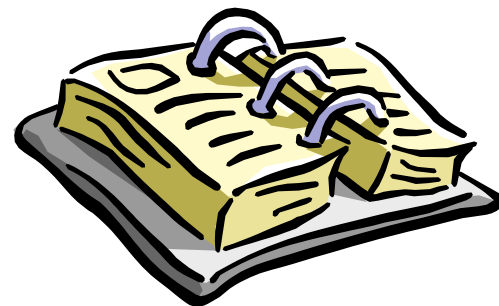
- Case-Study

- Results

- Conclusions

Problem

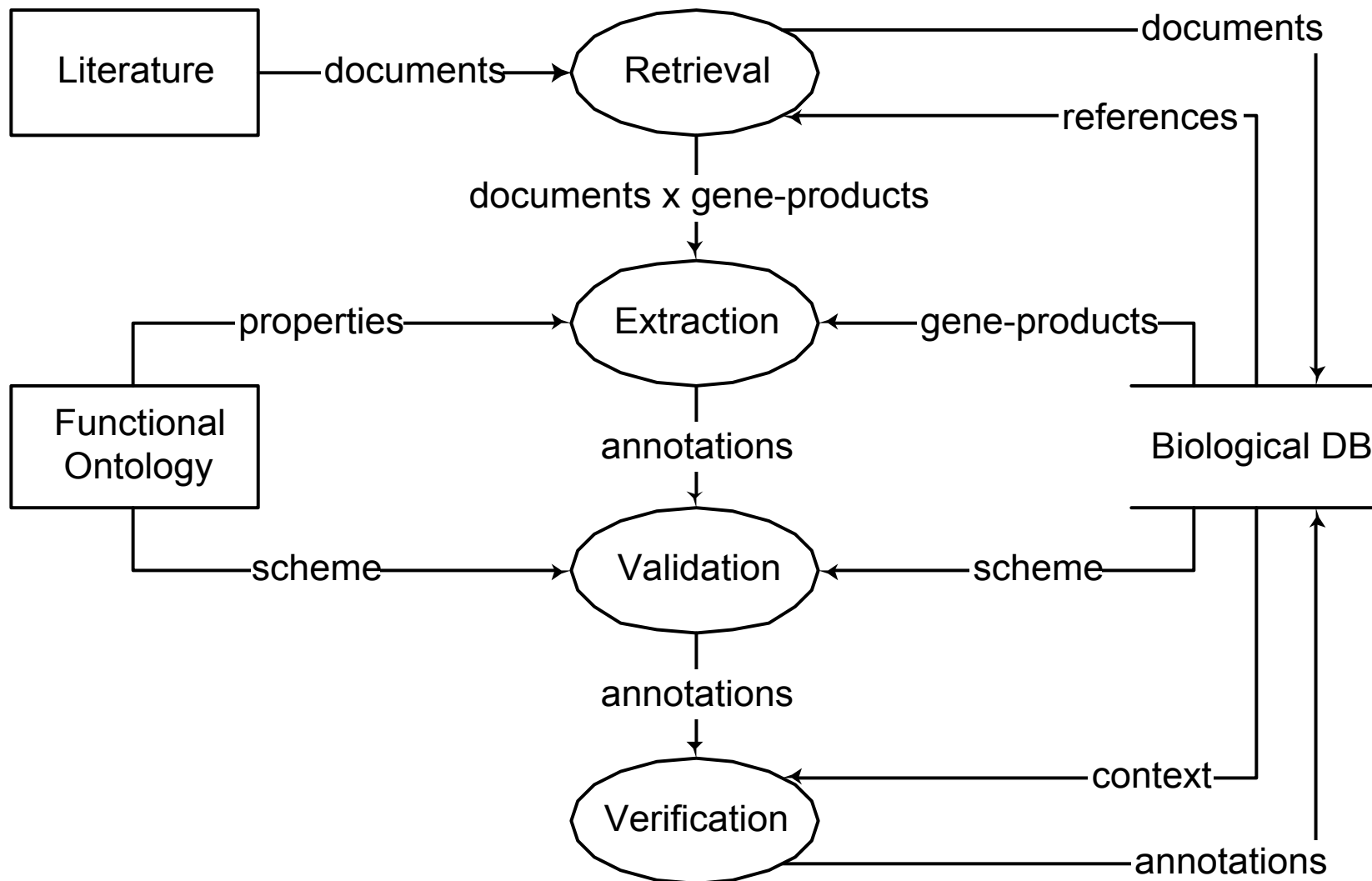
- Human curators read the literature to identify annotations and update their databases



ProFAL - Protein Functional Annotation Through Literature

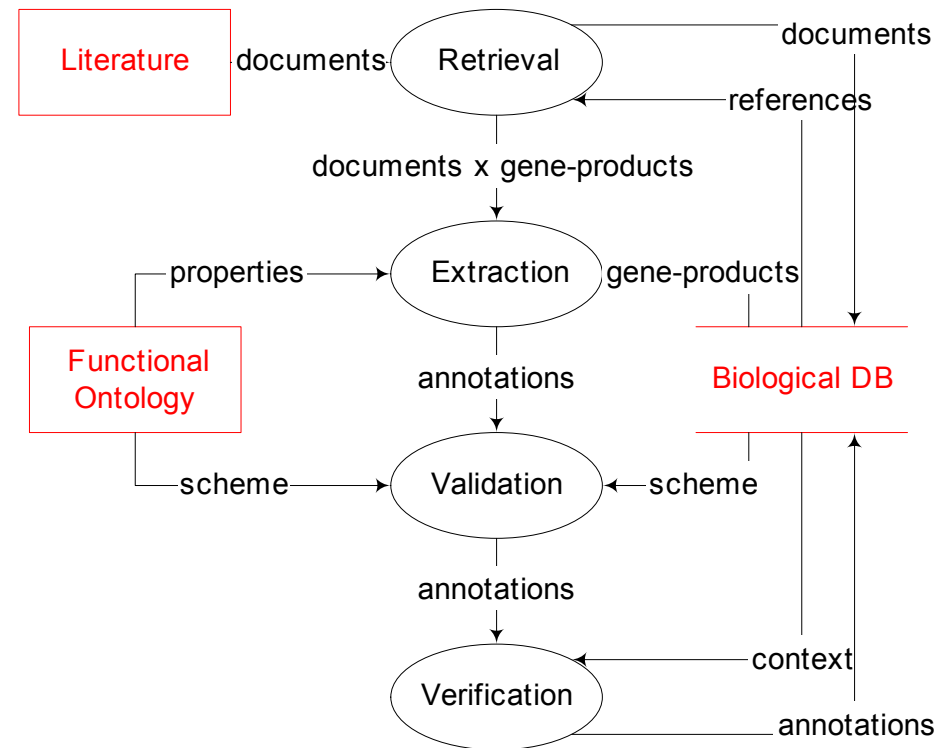
- Software tool for biologic database annotation and verification
 - Automatic retrieval and extraction of annotations from literature
 - Annotation validation method based on the correlation between structure and biologic function

ProFAL

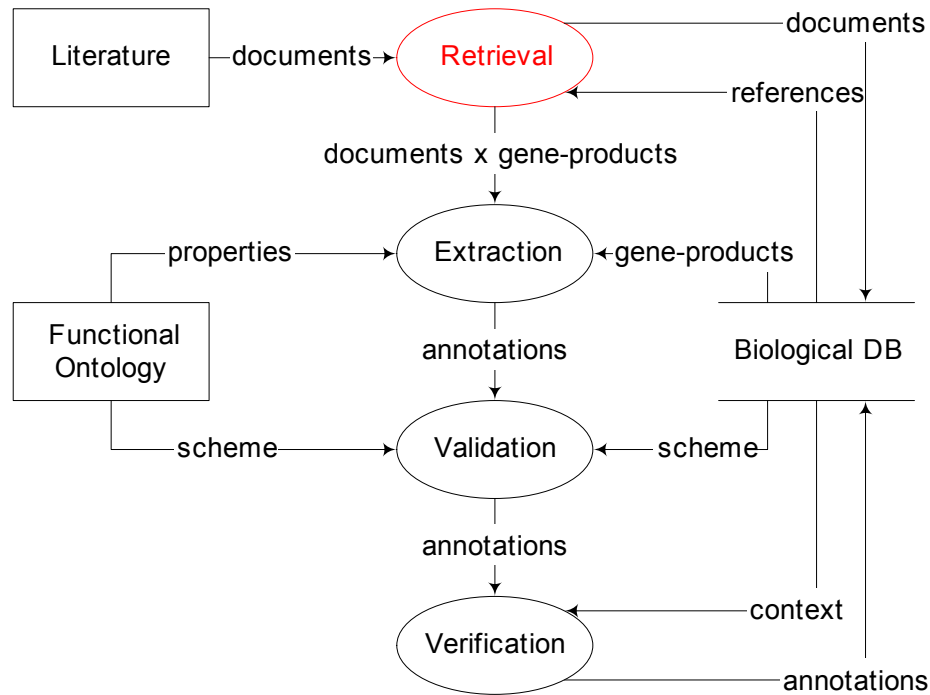


Information Sources

- **Biological DB**
 - Stores gene-products
 - Classified in families according to their structure
- **Functional Ontology**
 - A set of functional terms
 - Classified in a hierarchical taxonomy
- **Literature**
 - A collection of documents

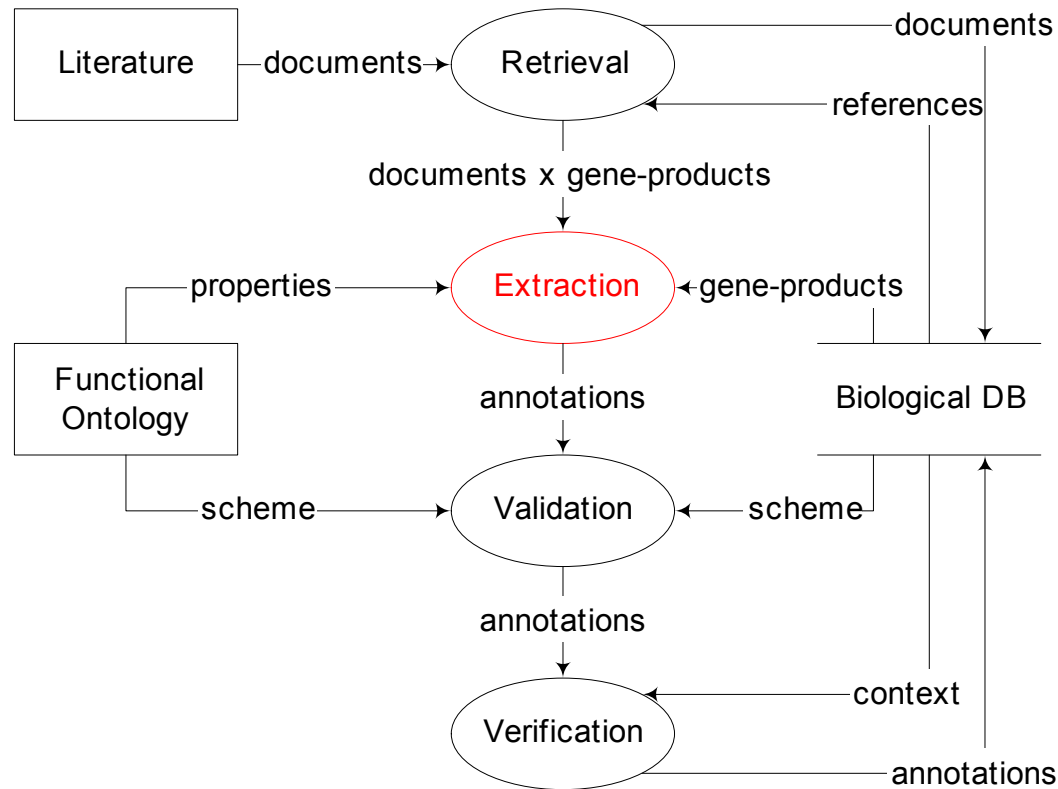


Processes: Retrieval



- Assign each gene-product to a set of documents
- Retrieve the document information

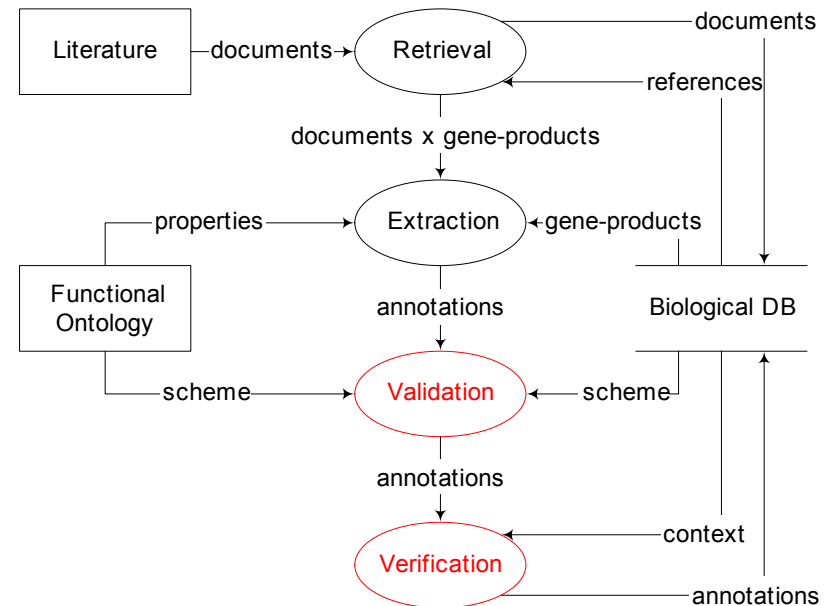
Extraction



- Annotate each gene-product with functional properties found in its assigned documents.

Validation & Verification

- Validation (automatic)
 - Use of heuristics about the domain
 - Quantitative classification of each annotation according to its reliability.
- Verification (computer assisted)
 - Analysis of annotations
 - Discarding of misannotations



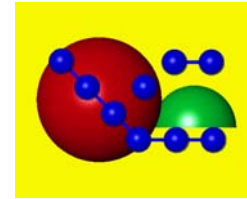
Outline

- Motivation
- ProFAL
- Case-Study
- Results
- Conclusions

Goals

- Validate the use of our annotation extraction and validation methods in a real environment
- Automate the annotation process on an existing biologic database

CAZy



- CAZy - Carbohydrate Active enZYmes
- Enzymes Classified in families
 - annotated with EC (Enzyme Commission) numbers
- Enzyme Sequences structure in modules (functional/structural segments of sequence)
 - Catalytic
 - carbohydrate Families

(<http://afmb.cnrs-mrs.fr/CAZY>)

Cazy Information Sources

- GenBank/GenPept,
- SwissProt,
- PDB
- Literature
- ...



CAZyModO - Glycoside Hydrolase Classification

Family GH48

Modular Organization

Go

CAZy Family Glycoside Hydrolase Family 48

Known Activities endoglucanase (EC [3.2.1.4](#)); cellobiohydrolase (EC [3.2.1.91](#)).

Mechanism Inverting

Catalytic Nucleophile/Base Not known

Catalytic Proton Donor Not known

3D Structure Status Available (see PDB)

Note formerly known as cellulase family L.

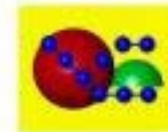
Relevant Links [InterPro](#); [PRINTS](#)

Statistics CAZyModO(12)

#ac	Protein	Organism	Modular Organisation	Status
154	cellulase CelA	Anaerocellum thermophilum	1 GH9 443 445CBM3612 613 681 682CBM3834 835 880 881CBM31034 1035 1080 1081 GH48 1111	
1497	CelA	Caldicellulosiruptor saccharolyticus	123 24 GH9 466 467CBM3642 643 700 701CBM3857 858 903 904CBM31060 1061 1112 1113 GH48 1120	
1591	cellobiohydrolase B	Cel48A Cellulomonas fimi	133 145 GH48 150 151EN3785 152 EN3786 153 EN3787 154 EN3788 989CBM21090	
10828	ORF CAC0911	Clostridium acetobutylicum ATCC 824	50 GH48 55 661DOC1726	
1716	CelF	Cel48A Clostridium cellulolyticum	129 13 GH48 160 665DOC1722	3D
1730	exoglucanase S	Clostridium cellulovorans	137 13 GH48 161 644DOC1703	
1733	CelD	Clostridium josui	129 13 GH48 150 651DOC1719	



CAZyModO



Family GH48

Modular Organization

Go

CAZy Family Glycoside Hydrolase Family 48

Known Activities endoglucanase (EC [3.2.1.4](#)); cellobiohydrolase (EC [3.2.1.91](#)).

Mechanism Inverting

Catalytic Nucleophile/Base Not known

Catalytic Proton Donor Not known

3D Structure Status Available (see PDB)

Note formerly known as cellulase family L.

Relevant Links [InterPro](#); [PRINTS](#)

Statistics CAZyModO(12)

Modular Organisation

Status

[1](#) [GH9](#) [443](#) [445CBM3612](#) [613](#) [681](#) [682CBM3834](#) [835](#) [880](#) [881CBM31034](#) [1035](#) [1080](#) [1080](#) [GH48](#) [1100](#)

[1](#) [23](#) [24](#) [GH9](#) [466](#) [467CBM3642](#) [643](#) [700](#) [701CBM3857](#) [858](#) [903](#) [904CBM31060](#) [1061](#) [1112](#) [1112](#) [GH48](#) [1147](#)

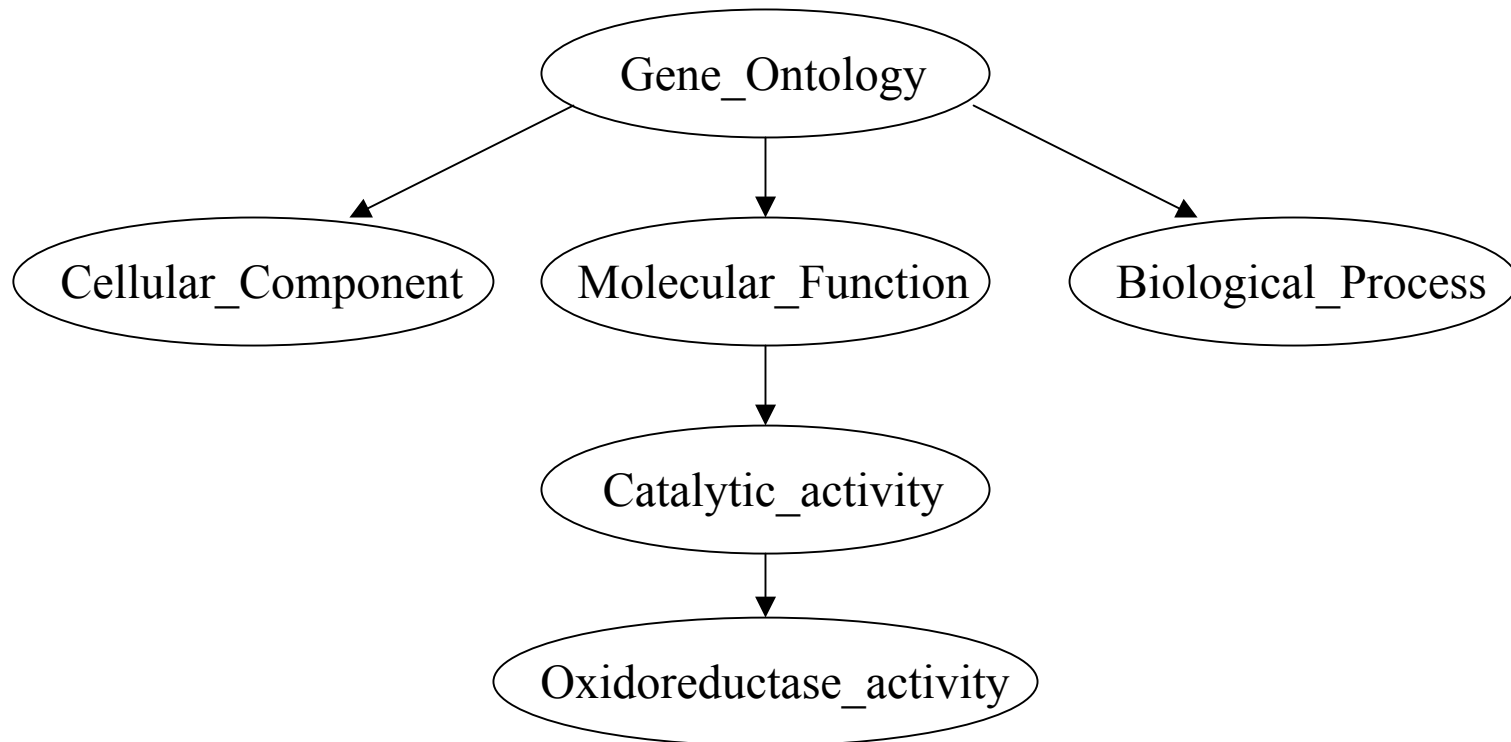
[1](#) [33](#) [34](#) [51](#) [GH48](#) [699](#) [699](#) [EN317](#) [699](#) [EN317](#) [699](#) [EN317](#) [989](#) [CBM21090](#)

Mechanism	Inverting
Catalytic Nucleophile/Base	Not known
Catalytic Proton Donor	Not known
3D Structure Status	Available (see 3D Structure)
Note	formerly known as Cel48A
Relevant Links	InterPro ; PRIN
Statistics	CAZyModO(1)

#ac	Protein	Organism	Modular Organisation
154	cellulase CelA	<i>Anaerocellum thermophilum</i>	1 GH9 443 445 CBM3612 613 681 682 CBM3834 835 880 881 CBM3
1497	CelA	<i>Caldicellulosiruptor saccharolyticus</i>	1 23 24 GH9 466 467 CBM3642 643 700 701 CBM3857 858 903 904 CB
1591	cellobiohydrolase B	<i>Cel48A Cellulomonas fimi</i>	1 33 34 53 GH48 609 700 FN3785 784 FN3884 801 FN397
10828	ORF CAC0911	<i>Clostridium acetobutylicum ATCC 824</i>	60 GH48 660 661 DOC1726
1716	CelF	<i>Cel48A Clostridium cellulolyticum</i>	1 20 GH48 660 665 DOC1722
1730	exoglucanase S	<i>Clostridium cellulovorans</i>	1 33 GH48 644 644 DOC1703
1733	CelD	<i>Clostridium josui</i>	1 29 GH48 650 651 DOC1719

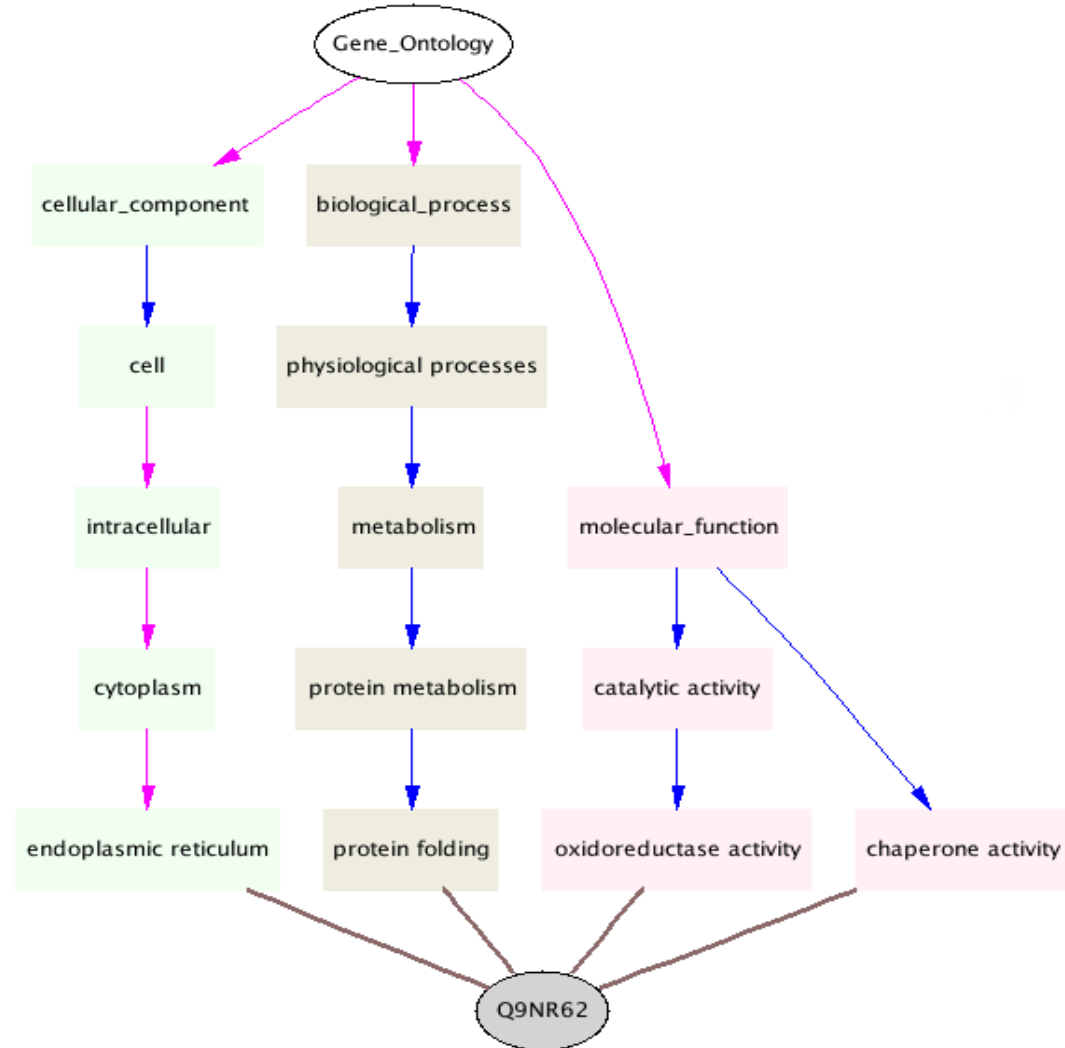
Gene Ontology

- Structure concepts and their relations:



Annotation

- Annotate objects related to the concepts:
- Example:
Protein name
Endoplasmic reticulum oxidoreductin 1-Lbeta



PubMed Abstracts

The screenshot displays the PubMed search interface. At the top, there are logos for NCBI, PubMed, and the National Library of Medicine (NLM). Below the logos is a navigation bar with tabs for PubMed, Nucleotide, Protein, Genome, Structure, PopSet, Taxonomy, OMIM, and Books. The search bar contains the text 'Carbohydrate Active enZymes' and has buttons for 'Go' and 'Clear'. Below the search bar are options for 'Limits', 'Preview/Index', 'History', 'Clipboard', and 'Details'. A second navigation bar includes 'Display' (set to 'Summary'), 'Sort', 'Save', 'Text', 'Clip Add', and 'Order'. Below this, it shows 'Show: 20' items, 'Items 1-20 of 14994', 'Page 1 of 750', and 'Select page: 1 2 3 4 5 6 7 8 9 10 »'. The main content area lists three abstracts, each with a checkbox, a link to the full article, and a 'Related Articles' link.

NCBI

PubMed

National Library of Medicine NLM

PubMed Nucleotide Protein Genome Structure PopSet Taxonomy OMIM Books

Search PubMed for Carbohydrate Active enZymes Go Clear

Limits Preview/Index History Clipboard Details

Display Summary Sort Save Text Clip Add Order

Show: 20 Items 1-20 of 14994 Page 1 of 750 Select page: 1 2 3 4 5 6 7 8 9 10 »

1: [Everard JD, Franceschi VR, Loescher WH](#) Related Articles
Mannose-6-Phosphate Reductase, a Key Enzyme in Photoassimilate Partitioning, Is Abundant and Located in the Cytosol of Photosynthetically Active Cells of Celery (*Apium graveolens* L.) Source Leaves.
 Plant Physiol. 1993 Jun;102(2):345-356.
 PMID: 12231825 [PubMed - as supplied by publisher]

2: [Bernassola F, Federici M, Corazzari M, Terrinoni A, Hribal ML, De Laurenzi V, Ranalli M, Massa O, Sesti G, McLean WH, Citro G, Barbetti F, Melino G](#) Related Articles
Role of transglutaminase 2 in glucose tolerance: knockout mice studies and a putative mutation in a MODY patient.
 FASEB J. 2002 Sep;16(11):1371-8.
 PMID: 12205028 [PubMed - indexed for MEDLINE]

3: [Keeshan K, Cotter TG, McKenna SL](#) Related Articles
High Bcr-Abl expression prevents the translocation of Bax and Bad to the mitochondrion.
 Leukemia. 2002 Sep;16(9):1725-34.
 PMID: 12200687 [PubMed - indexed for MEDLINE]

About Entrez

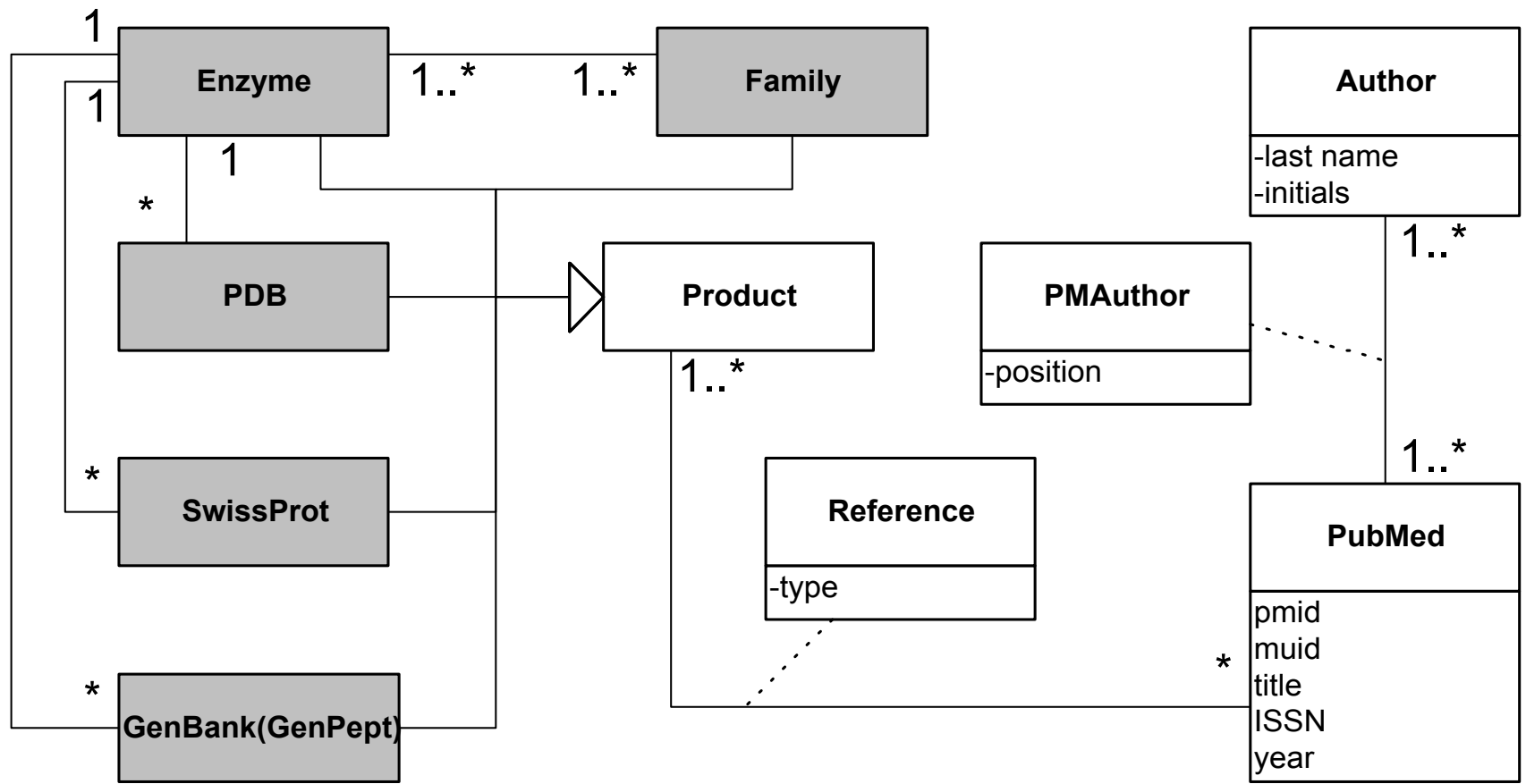
Text Version

Entrez PubMed
 Overview
 Help | FAQ
 Tutorial
 New/Noteworthy
 E-Utilities

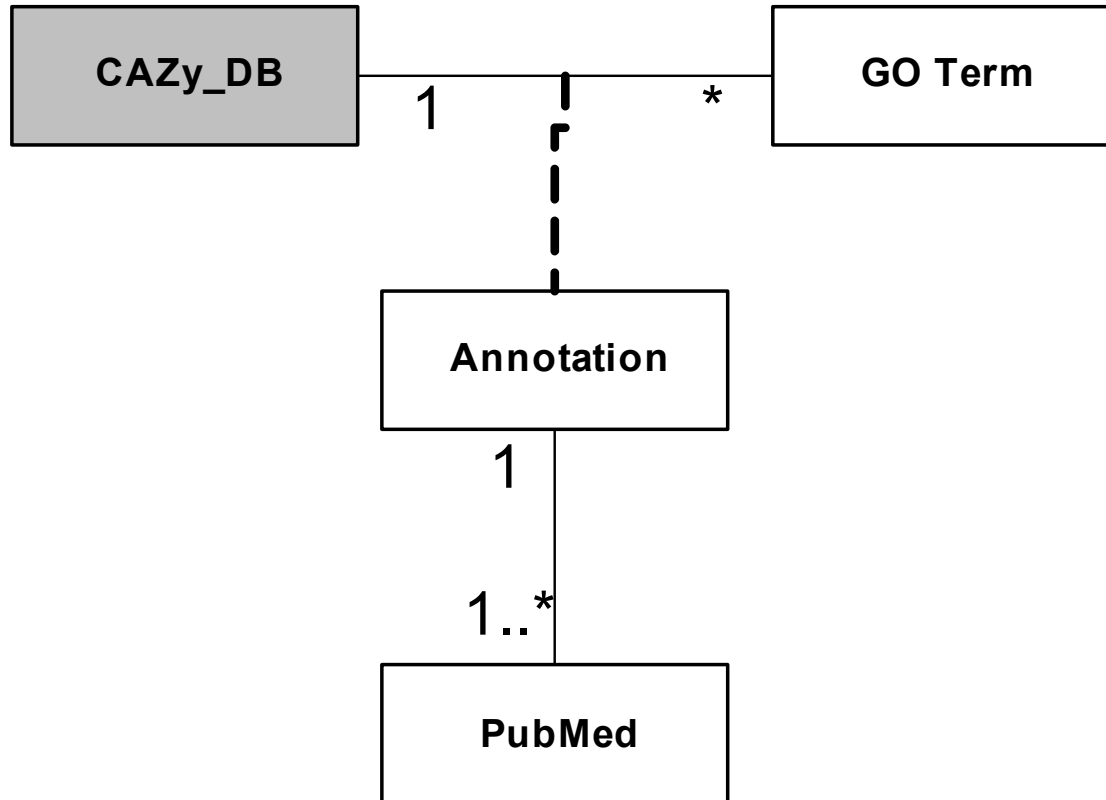
PubMed Services
 Journal Browser
 MeSH Browser
 Single Citation
 Matcher
 Batch Citation
 Matcher
 Clinical Queries
 LinkOut
 Cubby

Related Resources

Retrieval

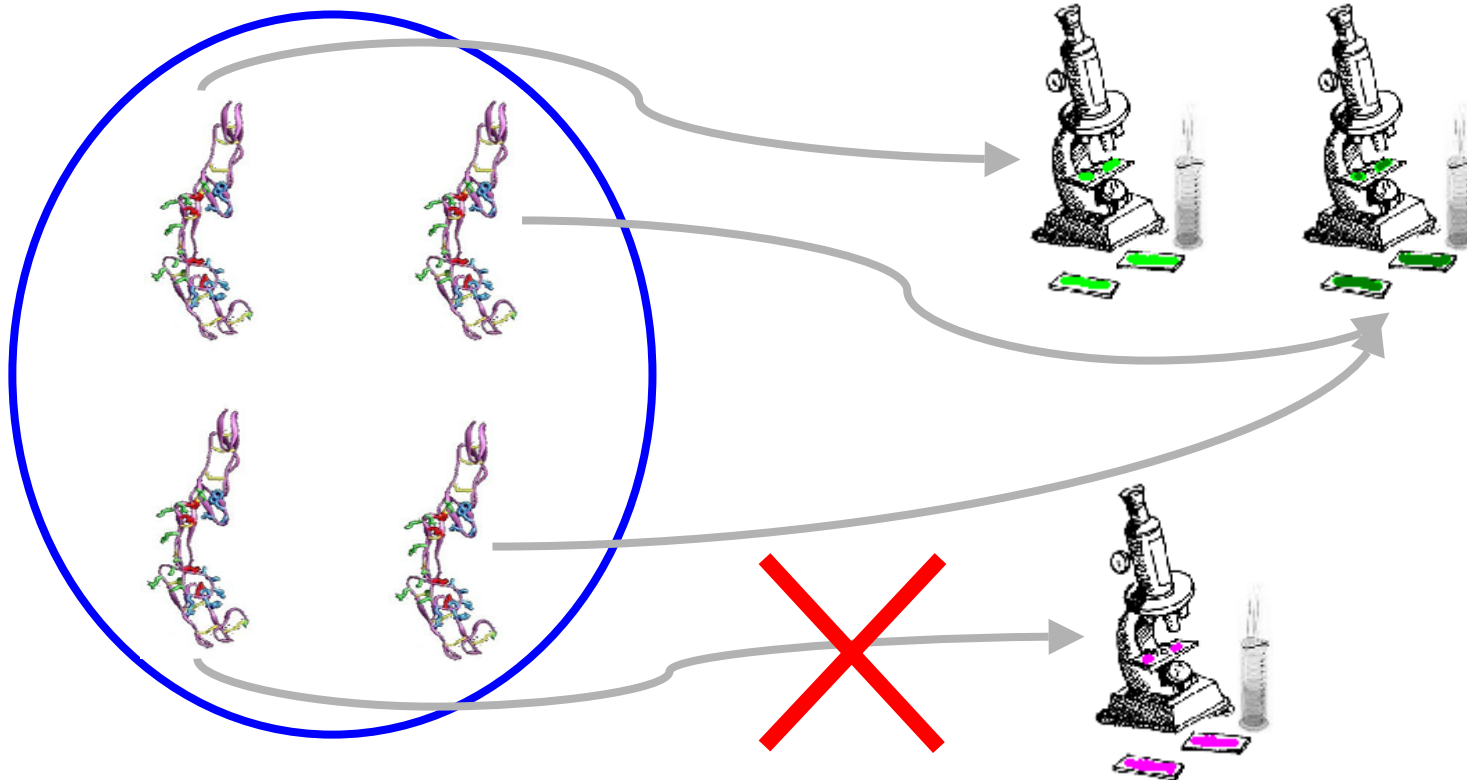


Extraction



Validation Method

- Correlation between function and structure



Verification

Publications								
PubMedID	MedlineID	Title	ISSN	Year	Classification	Note	#Authors	DB_ac
11435116	21345419	The kappa-carrageenase of <i>P. carrageenovora</i> features a tunnel-shaped active site: a novel insight in the evolution of Clan-B glycoside hydrolases.	0969-2126	2001	1	-3D-	7	
8112578	94156170	The gene encoding the kappa-carrageenase of <i>Alteromonas carrageenovora</i> is related to β -1,3-1,4-glucanases.	0378-1119	1994	1		3	
Options	Search All		<input type="text"/>	Insert PMID	<input type="text" value="11435116"/>	<input type="text" value="9"/>	Alter Classification	

Annotations					
TermsID	TermsName	Classification	Note	PubMedIDs	
GO:0016787	hydrolase	1		8112578	
GO:0008810	cellulase	1		11435116	
Options	<input type="text"/>	Insert Term	<input type="text" value="GO:0008810"/>	<input type="text" value="9"/>	Alter Classification

Outline

- Motivation
- ProFAL
- Case-Study
- Results
- Conclusions

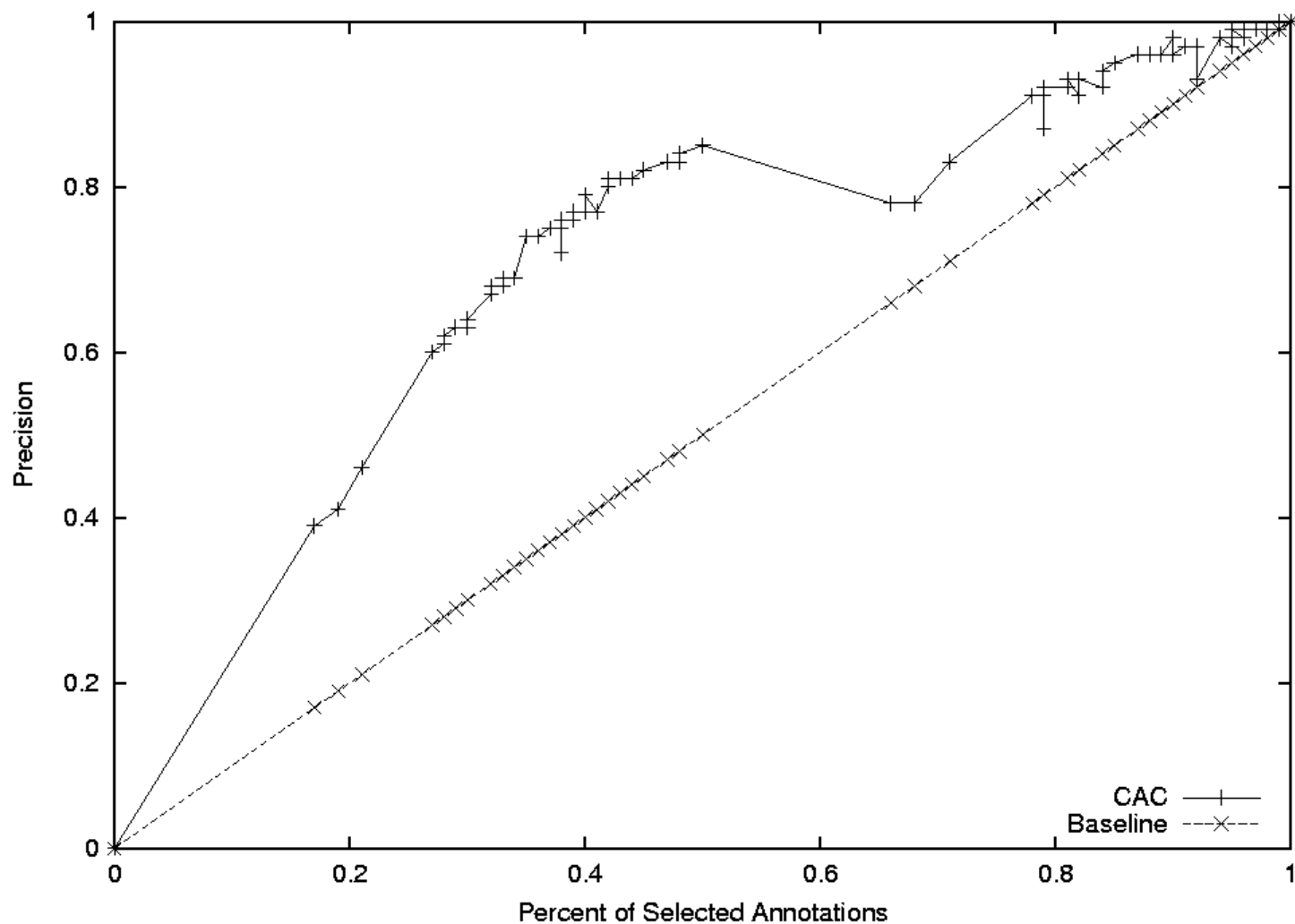
Retrieval

	References /Citations	Documents
GenBank(GenPept)	22849	4575
SwissProt	8998	4006
PDB	3561	785
Total	35408	6377

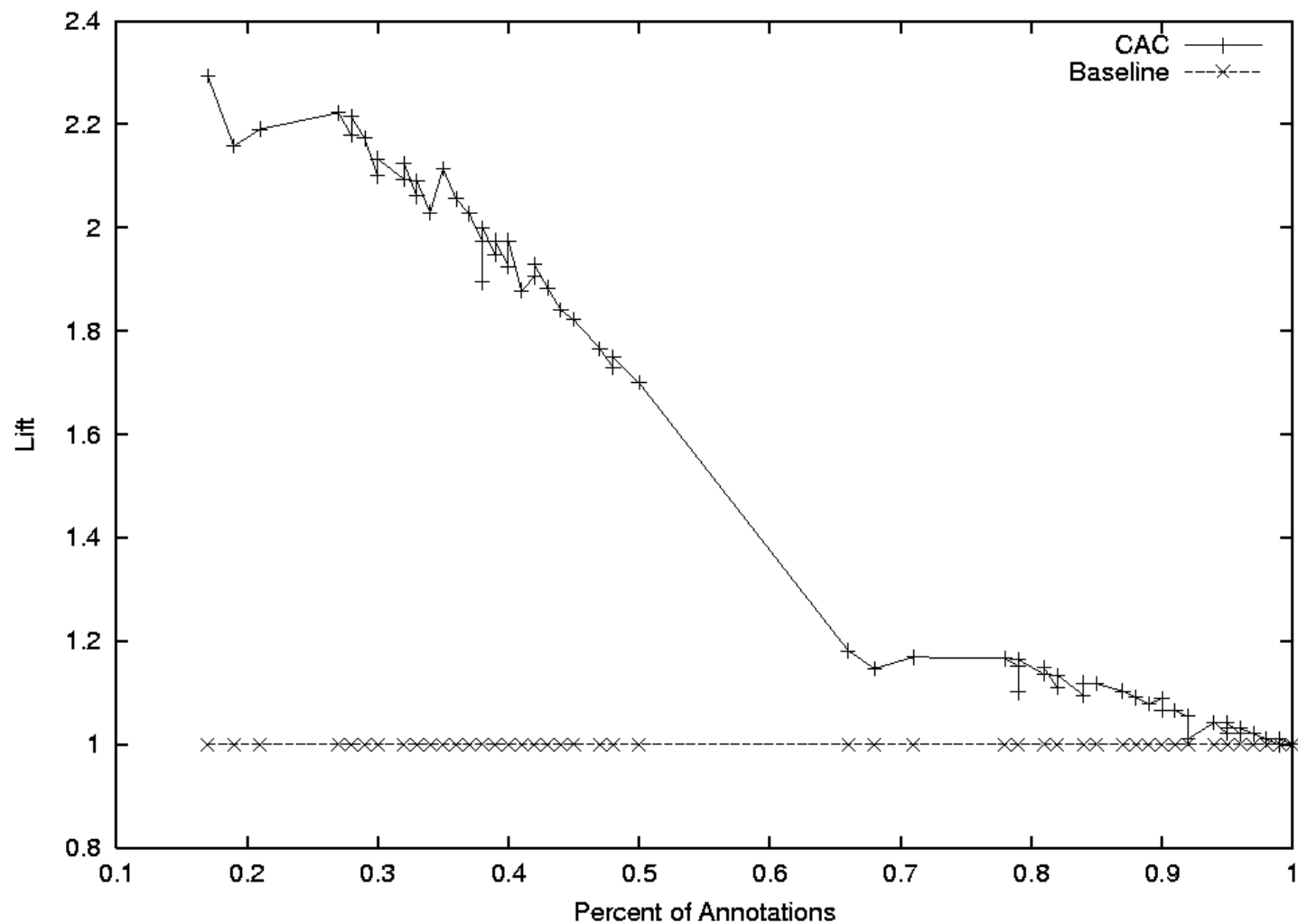
Extraction

- 13869 annotations
- 6918 enzymes
- 1342 terms
- Only 40% of enzymes were annotated
 - Lack of bibliographic references for 60% of enzymes
 - Few references in Cazy
 - Abstracts only
- Average of 2.2 annotations per document

Validation: Gain Chart



Validation: Lift Chart



Verification

- 173 annotations from 5 families
- Human verification
 - 95 correct and 78 incorrect annotations
 - Precision of 55%
 - Task completed in one hour
- Assuming a correct annotation per document:
 - The curator had to read 95 documents to identify 95 correct annotations
 - Assuming that reading a document takes 10 minutes, this would require about 16 hours without ProFAL.

Outline

- Motivation
- ProFAL
- Case-Study
- Results
- Conclusions



Conclusions

- ProFAL
 - Retrieves related Literature
 - Extracts annotations between gene-products and functional properties
 - Validates the annotations
 - Verifies the annotations
- Applied to CAZy
 - Improves the curation process done by humans
- Application to Arabidopsis gene expression data under development

More Information:

<http://xldb.di.fc.ul.pt/rebil/>