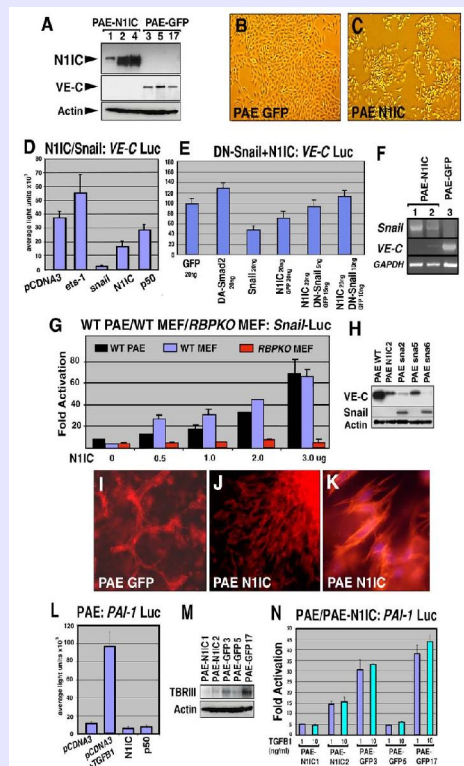# Text mining the Biomedical Literature

"We have here much  data, and we must proceed to lay out our campaign",

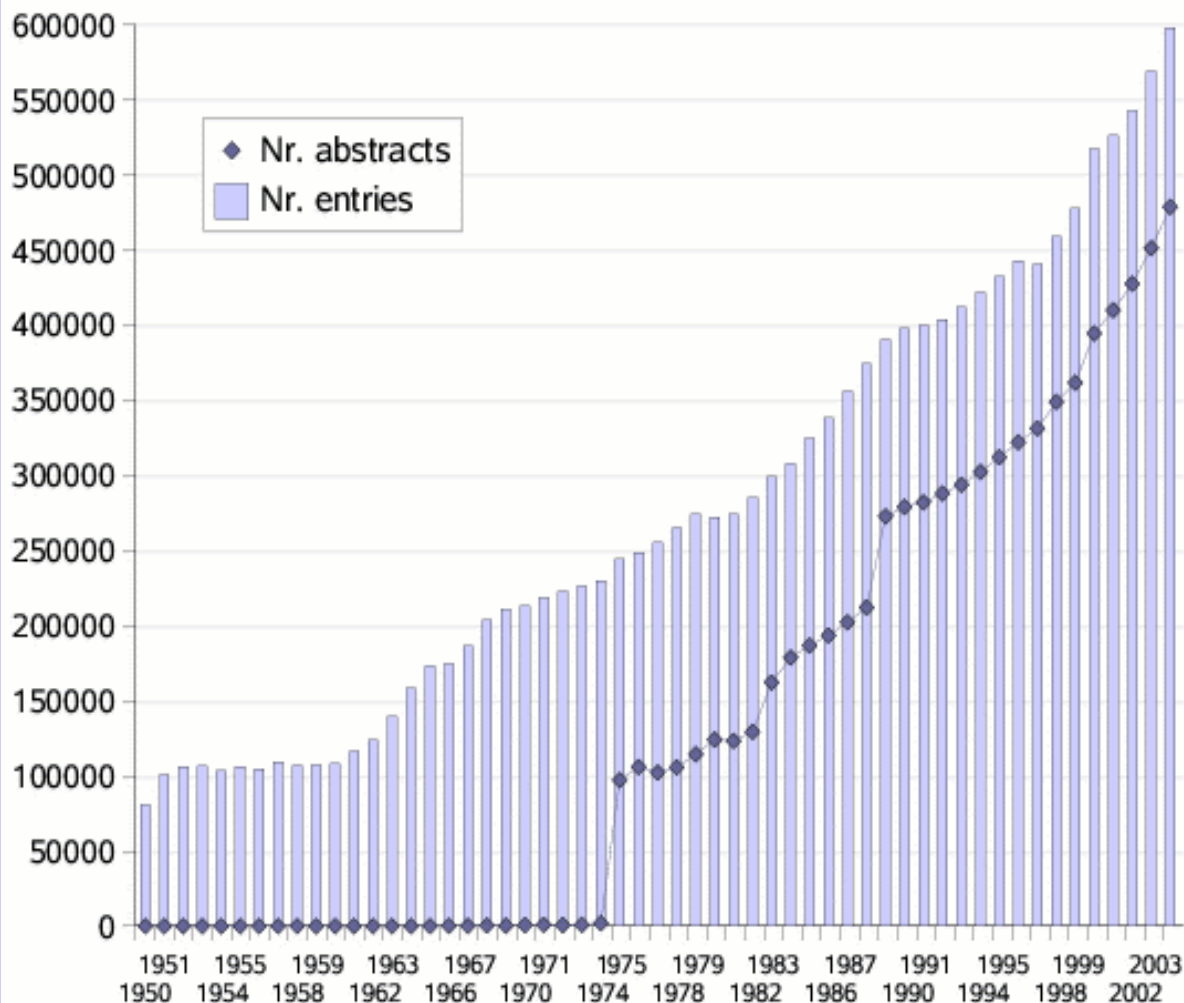Van Helsing in Bram Stockers Dracula

# Talk overview:

- The Biomedical literature
- Natural language processing (NLP)
- NLP in the Molecular Biology domain
- Text mining applications
- Evaluation of Text mining tools
- Conclusions and outlook
- Useful links, reviews and articles

# From experiments to scientific publications

## 1- Experiments

**Planning and carrying out experiments (lab work)**

## 2- Results

**Processing and interpretation of obtained results**



## 3- Scientific Peer-reviewed articles

**'Relevant' results are published in scientific journals**

**PDG** Protein Design Group

# Data in scientific articles

**Scientific Journals**



**Free Text**



**Title**

**Abstracts**

**Keywords**

**Text body**

**References**

**Tables**



**Figures**



**Journal-specific Information:**

• **Format**
• **Paper structure (sections)**
• **Article type**

**Biomedical literature characteristics**
– Heavy use of domain specific terminology (12% biochemistry related technical terms).
– Polysemic words (word sense disambiguation).
– Most words with low frequency (data sparseness).
– New names and terms created.
– Typographical variants
– Different writing styles (native languages)

**Text mining biomedical literature (2005)**

# PubMed/Medline database at NCBI

## PubMed growth



**Pub**Med

- Developed at the National Center for Biotechnology Information (NCBI).

– The core 'Textome'.

– repository of citation entries of scientific articles.

– PubMed titles and abstracts are primary data source for Bio–NLP.

– ~ 450,000 new abstracts/a

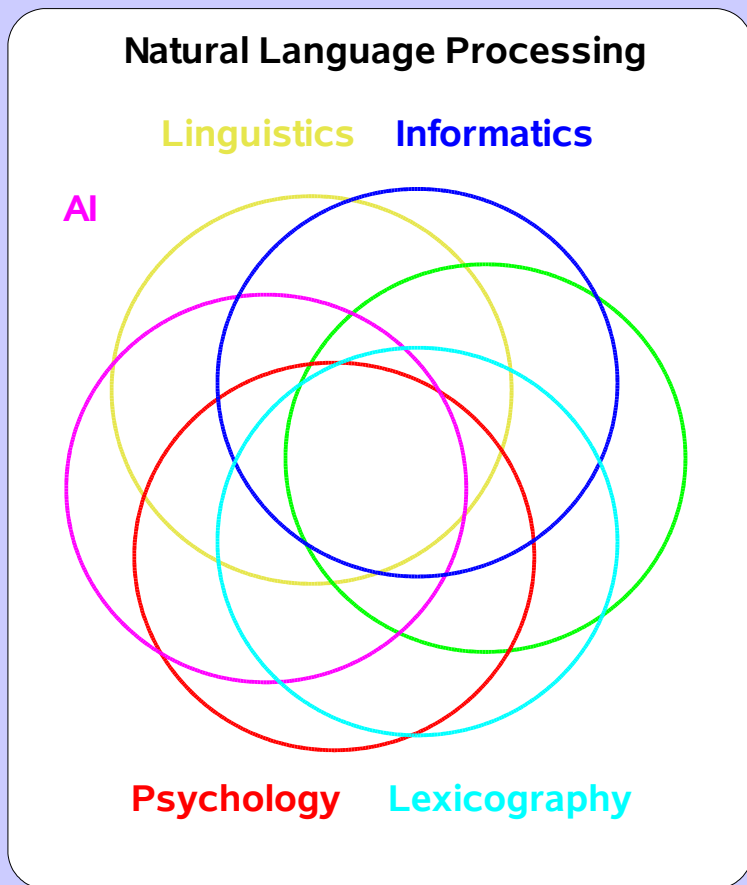– > 4,800 biomedical journals

– ENTREZ search engine

# PubMed online



**Text mining biomedical literature (2005)**

# Natural Language Processing (NLP) basics

**Natural Language Processing**

**Linguistics**  **Informatics**

**AI**

**Psychology**  **Lexicography**

**Domain, e.g. Biomedicine/ Molecular Biology**

➢ **Techniques that analyse, understand and generate language** (free text, speech).

➢ Linguistic tools, e.g. syntactic analyser and semantic classification.

➢ Multidisciplinary field.

➢ Strongly language dependent.

➢ Create computational models of language.

➢ Learn statistical properties of language.

➢ Methods: statistical analysis, machine learning, rule-based, pattern-matching, AI, etc...

➢ Domain dependent (biomedical) vs generic NLP.

# Major NLP tasks

- Information Retrieval (IR).
- Information extraction/Text mining (IE).
- Question Answering (QA).
- Natural Language Generation (NLG).

- Automatic summarisation.
- Machine translation.
- Text proofing.
- Speech recognition.
- Optical character recognition (OCR).

# Information Retrieval (IR)

➢ 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.

➢ First step: 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.

# Boolean query

- Based on **combination of terms** using Boolean operators.

- Basic **Boolean operators**: AND, OR and NOT.

- Queries matched against the terms in the inverted index file.

- Entrez – Boolean search in PubMed.

- Fast, easy to implement.

- **Search engines**: often stop word removal and case folding.

- Stop word removal : space saving speed improvement.

- Return a **unranked list**.

- Return large list of documents, many not relevant.

- Terms have different information content ->

  better weighted query.

# Zipf's law



> A small number of words occur very often

> Those high frequency words are often function words (e.g. prepositions)

> Most words with low frequency .

From: Rebholz-Schuhmann D, Kirsch H, Couto F (2005) Facts from Text—Is Text Mining Ready to Deliver? PLoS Biol 3(2): e65

# Commonly excluded stop words

| after | also | an | and |
|---|---|---|---|
| as | at | be | because |
| before | between | but | before |
| for | however | from | if |
| in | into | of | or |
| other | out | since | such |
| than | that | the | these |
| there | this | those | to |
| under | upon | when | where |
| whether | which | with | within |
| without | . | . | . |

# Vector space model

➢ Measure **similarity** between query and documents.

➢ Query can be a list of terms or whole documents.

➢ Documents and queries as **vectors of terms**.

➢ **Term weighting** according to their frequency:

  ➢ within the document

  ➢ within the document collection

➢ Widespread term weighting: tf x idf.

➢ Calculate similarity between

   those vectors.

➢ Cosine similarity.

➢ Return a ranked list.

➢ Example:  related article

   search in PubMed

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

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

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

**Text mining biomedical literature (2005)**

# PubMed online



**Text mining biomedical literature (2005)**

# eTBlast (1)



Science, May 14, 2004 issue. Under NetWatch, see the topic, "TOOLS: Just the Right Words".

# eTBlast (2)

# eTBlast (3)



**Text mining biomedical literature (2005)**

# eTBlast (4)



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

**Text mining biomedical literature (2005)**

# IR performance

➢ **Precision**: fraction of relevant documents retrieved

    divided by the total returned documents

➢ **Recall**: proportion of relevant documents returned

    divided by the total number of  relevant documents

➢ **F-score**: the harmonic mean of precision and recall

➢ Precision-recall curves

# Information Extraction and Text mining

➢ 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,..

# Stemming

Process of removing affixes of words transforming them to their corresponding morphological base form or root.

### Porter's Stemming Algorithm Online - Mozilla

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!

Done

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

### Porter's Stemming Algorit

File Edit View Go Bookmarks Tools Window

Search

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 |

# POS-tagging



Providing each word given a sentence with its corresponding part of speech label , e.g. whether it is a noun, verb, preposition, article, etc.

# Question Answering (QA)

➢ Humans formulate questions using natural language.

➢ Example: *What are the molecular functions of Glycogenin?*.

➢ QA: **automatic generation of answers** to queries in form

   NL expressions from document collections.

➢ Most systems limited to generic literature or newswire.

➢ QA difficult: heterogeneous, poorly formalised domain,

   new scientific terms

➢ Ad hoc retrieval task of the TREC Genomics Track 2005.

➢ Galitsky system (semantic skeletons (SSK), logical

   programming).

# **Natural Language Generation (NLG)**

➢ NLG: constructing automatically natural language texts.

➢ Display the content of databases: reports, error messages.

➢ Based on semantic input, providing computer-internal

representation of the information.

➢ Different degrees of complexity.

➢ Biology: modelling the domain language difficult.

➢ Simpathica/XSSYS trace analysis tool.

# Named entity recognition (NER)

➢ **Identification of entity types** in textual data.

➢ Semantic tagging.

➢ Example identification of company names and places

➢ Mainly identification of proper nouns.

➢ NER in Molecular Biology: identification of genes,

   proteins, chemical compounds, diseases,...

➢ Methods: ad-hoc rule based systems,

➢    ML techniques (HMM,SVM,...), statistical tools.

➢ Tools: GAPSCORE, ABNER, AbGene, NLProt

# ABNER



Burr Settles. "ABNER: A Biomedical Named Entity Recognizer."
http://www.cs.wisc.edu/~bsettles/abner/. 2004.

**Text mining biomedical literature (2005)**

# Basic NLP terms

**Corpus**: collection of documents.

**POS tagging**: labeling each word in

a sentence with its part of speech (verb,noun,..)

depending on its context.

**Word sense disambiguation**: assigning the

semantic class (meaning) to a given word

depending on its context.

# NLP in Biomedicine – Timeline

**AI / MACHINE LEARNING**
- Bayesian Classifier
- HMMs | CRFs | SVMs
- Neural Networks | MEMMs

**NLP**
- Shallow Parsing
- POS - Tagging
- Stemming

Pathology reports

Neighboring relationships

MEDLINE

UMLS

Protein/Gene NER tagging

Biology databases

Assessments
Applications
Methods
Data resources

**Before 1990** → **1990-1995** → **1995-2000** → **2000-2004**

Function prediction

Automatic Annotations

Protein sequence analysis

New generation of Visualization and Browsing systems

Protein interactions

BioCreative I corpus

Microarrays analysis

PubMed

Gene Ontology

GENIA corpus

Cellular localization

KDD cup

TREC I

BioCreative 1

JNLPBA Shared task

TREC II

LLL05

**Text mining biomedical literature (2005)**

# Text mining applications in biology

➢ NER: tagging biological entities.

➢ Automatic annotation: associating proteins to

   functional descriptions.

➢ Protein interactions: extracting interactions of

   proteins, genes and drugs.

➢ Microarray analysis: providing biological context

   through literature mining

➢ Protein localisation

➢ Improving sequence-based homology detection.

# Text mining applications in biology



**Text mining biomedical literature (2005)**

# Tagging Biological entities

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).

➢ 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).

# GAPSCORE (1)

Gene and Protein Search - Mozilla Firefox

Archivo   Editar   Ver   Ir   Marcadores   Herramientas   Ayuda

http://acronym.stanford.   Ir

Search for gene and protein names in some text.

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,

SEARCH

| | Gene or Protein Name | Quality (Score) |
|---|---|---|
| 1 | 75% VO2 | Good (0.70) |
| 2 | Glycogenin | Good (0.67) |
| 3 | Glycogenin | Good (0.67) |
| 4 | Glycogenin | Good (0.67) |
| 5 | elevated glycogenin | Good (0.67) |
| 6 | free (deglycosylated) glycogenin | Good (0.67) |
| 7 | glycogenin | Good (0.67) |
| 8 | glycogenin | Good (0.67) |
| 9 | glycogenin | Good (0.67) |
| 10 | glycogenin | Good (0.67) |

Terminado

➢ Scores words based on a statistical model of gene names

➢ Quantifies:
  ➢ Appearance
  ➢ Morphology
  ➢ Context.

➢ Online.

http://bionlp.stanford.edu/gapscore/

**Text mining biomedical literature (2005)**

# GAPSCORE (2)



> ➢ Choose cut-offs.

> ➢ Online.

> ➢ Based on Medline analysis

> ➢ Score new words using SVM

> ➢ Statistical analysis of PubMed words.

Chang JT, Schütze H, and Altman RB.
GAPSCORE: Finding Gene and Protein
Names One Word at a Time.
*Bioinformatics*. 2004 Jan 22;20(2):216-25.

# NLProt

| NAME | ORGANISM | TXT-POS | SCORE | METHOD | DB-ID(S) | |
|------|----------|---------|-------|--------|----------|---|
| Glycogenin | homo sapiens | 1 | 1.040 | SVM | GYG2 HUMAN | (86%) |
| glycogenin | homo sapiens | 96 | 0.856 | SVM | GYG2 HUMAN | (91%) |
| glycogenin | homo sapiens | 103 | 1.040 | SVM | GYG2 HUMAN | (91%) |
| Glycogenin | homo sapiens | 109 | 0.871 | SVM | GYG2 HUMAN | (86%) |
| glycogenin | homo sapiens | 138 | 0.980 | SVM | GYG2 HUMAN | (91%) |
| Glycogenin | homo sapiens | 157 | 0.971 | SVM | GYG2 HUMAN | (86%) |
| glycogenin | homo sapiens | 161 | 0.311 | SVM | GYG2 HUMAN | (91%) |
| glycogenin | homo sapiens | 214 | 0.819 | SVM | GYG2 HUMAN | (91%) |
| glycogenin | homo sapiens | 234 | 0.747 | SVM | GYG2 HUMAN | (91%) |

➢ Online (e-mail alert).

➢ Downloadable.

➢ SVM-based

➢ Pre-processing dictionary

➢ Rule-based filtering step

➢ PubMed words.

➢ Precision of 75%

➢ Recall of 76%

http://cubic.bioc.columbia.edu/services/nlprot/

Chang JT, Schutze H, Altman RB. GAPSCORE: finding gene and protein names one word at a time. *Bioinformatics*. 2004 Jan 22;20(2):216-25.

# ABNER

> A Biomedical Named
> Entity Recogniser

> Downloadable.

> CRF-based

> Trained on BioCreative
> and GENIA

> orthographic and
> contextual features

> Can be trained on
> new corpora

Burr Settles. "ABNER: A Biomedical Named Entity Recognizer."
http://www.cs.wisc.edu/~bsettles/abner/. 2004.

**Text mining biomedical literature (2005)**

# Extracting functional annotations

- **Manual annotation** extraction by database curators.
  - Scientific literature analysis.
  - Time-consuming & labour-intensive.
  - Example: Gene Ontology annotation (GOA).

- **Text mining** to assist annotation extraction:
  - Identification of annotation relevant sentences.
  - Identification of protein-term associations.

# Function extraction – applications

- ➢ Andrade et. (1997)

- ➢ iHOP

- ➢ Textpresso system

- ➢ Gene Information System (GIS)

- ➢ Medical Knowledge Explorer (MeKE)

- ➢ GO engine,...

# Andrade et al. (1997)

➢ Extracts sentences from PubMed which contain

functional information.

➢ Statistical analysis of the word frequencies.

➢ Analysis in protein families.

➢ Background frequencies of those words.
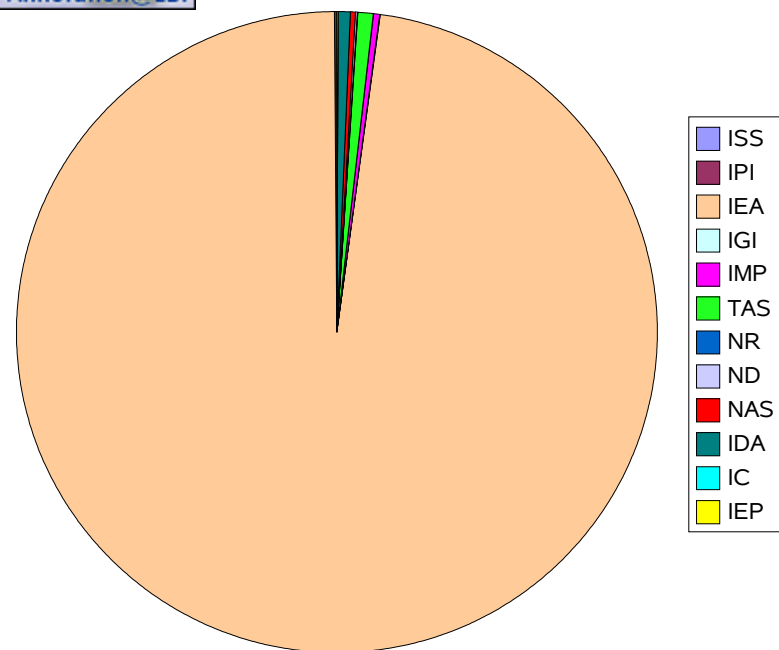
# GENE ONTOLOGY (GO)

➢ Ontology direct acyclic graph structure.

➢ Controlled vocabulary of concepts.

➢ Three main categories:

  ➢ Molecular Function

  ➢ Cellular Component

  ➢ Biological Process

➢ Describe relevant biological aspects of gene products

➢ Synonyms, links to external keywords.

➢ Currently most important source annotation terms.

http://www.geneontology.org/

# Gene Ontology Annotation

| Ev.C. | Annot | Perc. | |
|---|---|---|---|
| IEA | 6421817 | 0.97529 | Electronic/ |
| ISS | 19576 | 0.00297 | sequence- |
| NR | 2191 | 0.00033 | based |
| ND | 4433 | 0.00067 | annotation |
| IPI | 7130 | 0.00108 | |
| IGI | 3014 | 0.00046 | Experimental |
| IMP | 19072 | 0.00290 | evidence |
| IDA | 38862 | 0.00590 | |
| IEP | 1495 | 0.00023 | |
| IC | 831 | 0.00013 | |
| TAS | 49630 | 0.00754 | Curator |
| NAS | 16456 | 0.00250 | knowledge |

Legend: ISS, IPI, IEA, IGI, IMP, TAS, NR, ND, NAS, IDA, IC, IEP

TAS: Traceable Author Statement;  IDA: Inferred by direct assay;  IC: Inferred by curator ; ND:No data; IMP:Inferred from mutant phenotype;  IGI: Inferred from genetic interaction; 3.8) IPI :Inferred from physical interaction; ISS: Inferred from sequence similarity;  IEP: Inferred from expression pattern; NAS: Non traceable author statement;  IEA: Inferred by electronic annotation;  NR: Not recorded;

http://www.ebi.ac.uk/GOA/  04/22/05

**Text mining biomedical literature (2005)**

# Gene Ontology Growth



GO growth

- MF:Molecular
  Function
- CC: Cellular
  Component
- BP: Biological
  Process

# iHOP

➢ Protein centric: nucleates the literature around protein name.

➢ For a range of model organisms (e.g. Human, yeast,..)

➢ Hyperlinks proteins through co-occurrence

➢ Highlight direct associations between proteins and functional

   terms.

➢ Online, fast, easy to use.

Hoffmann R, Valencia A. A gene network for navigating the literature *Nat Gene*t. 2004 Jul;36(7):664.

# iHOP

# iHOP



**Text mining biomedical literature (2005)**

iHOP - Information Hyperlinked over Proteins [ GYG ] - Mozilla

File  Edit  View  Go  Bookmarks  Tools  Window  Help

http://www.pdg.cnb.uam.es/UniPub/iHOP/gs/88913.html?IN=1    Search

Home  Bookmarks  Yahoo  Google  MK Homepage  ORF  Zope on http://...  PubMed  Python  Zope  PyTut  OEAW  GeneDic  biocreative  GenomeNet

**iHOP**
information hyperlinked over proteins

| Symbol | Name | Synonyms | Organism |
| --- | --- | --- | --- |
| **GYG** | Glycogenin-1 | glycogenin, GYG1 | Homo sapiens |
| UniProt | P46976, Q9UNV0 | | |
| OMIM | 603942 | | |
| NCBI Gene | 2992 | | |
| NCBI RefSeq | NP_004121 | | |
| NCBI Accession | AAB00114, AAB09752, AAD31084 | | |

**Homologues of GYG ...** new

**Definitions for GYG** ...

**Enhanced PubMed/Google query ...** new

WARNING: Please keep in mind that gene detection is done automatically and can exhibit a certain error. Read more.

Find in this Page

Search Gene

Show overview new
Find in this Page

Filter and options
Gene Model

Developer's Zone
new
Help

Mutation of Tyr-196 in glycogenin-2 to a Phe residue abolished the ability of glycogenin-2 to self-glucosylate but not to **interact** with glycogenin-1.

Mutational analysis of the coding regions of the genes encoding protein kinase B-alpha and -beta, phosphoinositide-dependent protein kinase-1, phosphatase targeting to glycogen, protein phosphatase inhibitor-1, and glycogenin: lessons from a search for genetic variability of the insulin-**stimulated** glycogen synthesis pathway of skeletal muscle in NIDDM patients.

Effects of exercise on GLUT-4 and glycogenin gene expression in human skeletal muscle.

The third cDNA encoded a polypeptide of unknown function and was designated GNIP (glycogenin interacting protein).

GNIP, a novel protein that binds and activates glycogenin, the self-glucosylating initiator of glycogen biosynthesis.

Overall, GN-2 has 40-45% identity to muscle glycogenin but is 72% identical over a 200-residue segment thought to contain the catalytic domain.

Glycogenin-1 and glycogenin-2 interact with one another, based on in vitro interactions and co-immunoprecipitation from liver and cell extracts.

Mouse glycogenin-1 has a predicted molecular mass of 37¿ omitted¿399 Da, and the deduced amino acid sequence exhibited 87% homology with human glycogenin-1.

For the first time, we report that a single bout of exercise is sufficient to cause upregulation of GLUT-4 and glycogenin gene expression in human skeletal muscle.

Fasting plasma insulin concentrations, muscle creatine, glycogen and GLUT-4 protein content as well as GLUT-4, glycogen synthase-1 (GS-1) and glycogenin-1 (Gln-1) mRNA expression were determined.

In conclusion, the co-expression of glycogenin with GLUT3 might enable glycogen-storing cells to exchange glucose quite effectively according to prevailing metabolic demands of glycogen synthesis or degradation.

The discovery of a second human gene, GYG2, encoding a liver-specific isoform of glycogenin, the self-glucosylating initiator of glycogen biosynthesis, raised the possibility for differential controls of this protein in liver and muscle.

The present study investigated the expression of glycogenin, the protein primer for glycogen synthesis, and the high affinity glucose transporter isoform GLUT3 as a further potential regulator of cellular glycogen metabolism, in first trimester and term human placenta using immunohistochemistry and

Concept & Implementation
by Robert Hoffmann

Transferring data from www.pdg.cnb.uam.es...

# Text mining biomedical literature (2005)

# Textpresso

# Gene Information  System (GIS)

➢ Focus on 4 types of gene-related info:

  ➢ Biological function

  ➢ Associated disease

  ➢ Related genes

  ➢ Gene-gene relations

➢ Gene information screening

➢ Gene-gene relation extraction.

➢ Downloadable

http://iir.csie.ncku.edu.tw/~yuhc/gis/

# GIS



Chiang JH, Yu HC, Hsu HJ.GIS: a biomedical text-mining system for gene information discovery.*Bioinformatics.* 2004 Jan 1;20(1):120-1.

**PDG** Protein Design Group

# Keyword Annotation Tool (KAT)

➢ Extraction of mappings between related terms using a model of fuzzy associations

➢ Mesh terms/SwissProt keywords/GO terms

Perez AJ, Perez-Iratxeta C, Bork P, Thode G, Andrade MA.Gene annotation from scientific literature using mappings between keyword systems. Bioinformatics. 2004 Sep 1;20(13): 2084-91. Epub 2004 Apr 1.

**Text mining biomedical literature (2005)**

# Medical Knowledge Explorer (MeKE)

➢ Ontology-based text mining system.

➢ Methods of sentence alignment.

➢ Sentence classification methods.

➢ Flexible matching, stemming and indexing.

➢ Create new GO-term synonyms from text.

➢ Edit distance calculation

➢ Learn sentence motifs via sentence alignment

➢ Naïve Bayes sentence classifier

http://gen.csie.ncku.edu.tw/meke3/

Chiang JH, Yu HC.MeKE: discovering the functions of gene products from biomedical literature via sentence alignment. *Bioinformatics*. 2003 Jul 22;19(11):1417-22

# GO engine

- ➢ Computational platform for GO annotation.

- ➢ Correlation of text info with specific GO nodes.

- ➢ Combines: homology info, protein clustering and text analysis.

- ➢ Calculate frequency of association of terms to GO nodes.

Xie H, Wasserman A, Levine Z, Novik A, Grebinskiy V, Shoshan A, Mintz L. Large-scale protein annotation through gene ontology. Genome Res. 2002 May;12(5):785-94.

# Protein interactions

- Advances in experimental large scale protein

  interaction analysis

- Exp. Methods for protein interaction characterization:

  - protein arrays

  - mRNA expression microarrays

  - Yeast two-hybrid

  - Affinity purification with MS

  - X-ray, NMRFRET, chemical cross-linking,..

- Bioinformatics methods for protein characterization:

  - Genome-based

  - Sequence-based

# Protein interaction databases

| Database Name | Reference | URL |
|---|---|---|
| BIND | (Bader *et al*., 2003) | http://bind.ca |
| DIP | (Xenarios *et al*., 2002) | http://dip.doe-mbi.ucla.edu |
| GRID | (Breitkreutz *et al*. 2003) | http://biodata.mshri.on.ca/grid |
| HPID | (Han *et al.,* 2004) | http://www.hpid.org |
| HPRD | (Peri *et al*., 2004) | http://www.hprd.org |
| IntAct | (Hermjakob *et al*., 2004) | http:/www.ebi.ac.uk/intact |
| MINT | (Zanzoni *et al*., 2002) | http://cbm.bio.uniroma2.it/mint |
| STRING | (vonMering et al., 2003) | http://string.embl.de |
| ECID | (Juan *et al*., 2004) | http://www.pdg.cnb.uam.es/ECID |

# Text mining and Protein interactions

➢ Extract automatically those interactions from articles.

➢ NL used to characterise the nature of the interaction

and its directionality.

➢ Literature-derived interaction networks:

  ➢ power law distribution

  ➢ scale free topology

➢ Visualised using network graphs.

➢ Methods range from: simple occurrence, expert derived

word patterns (frames) to machine learning.

# PubGene

➢ Use the co-occurrence of protein and gene names.

➢ Assumption: co-occurrence imply biological relationship

➢ Indexing PubMed abstracts and titles with human proteins.

➢ Construction of  interaction networks.

➢ Build upon binary interactions between co-occurring proteins

Jenssen TK, Laegreid A, Komorowski J, Hovig E.A literature network of human genes for high-throughput analysis of gene expression.Nat Genet. 2001 May;28(1):21-8.

http://www.pubgene.org/

**iHOP:**

**Visualization**

**of protein**

**interactions**

**using network**

**graphs**

# SUISEKI

- ➢ Relationship between the co-occurring proteins using **frames**

- ➢ Frames: **textual patterns** used to express interactions

- ➢ Initial set of 14 interaction words based on domain knowledge.

- ➢ Examples: *activate, bind, suppress*

- ➢ Analysed the **order** of protein names within sentences.

- ➢ Take into account **distance** (off-set) between protein names.

- ➢ System effective for simple interaction types.

- ➢ Difficult cases: long sentences with complex

   grammatical structures

# SUISEKI interaction network

# iProLINK



➢ Mapped citations

➢ Annotation tagged literature corpora

➢ Data source for protein name ontology development

http://pir.georgetown.edu/iprolink/

Hu ZZ, Narayanaswamy M, Ravikumar KE, Vijay-Shanker K, Wu CH.Literature mining and database annotation of protein phosphorylation using a rule-based system.

*Bioinformatics*. 2005 Jun 1;21(11):2759-65. Epub 2005 Apr 6

# Chilibot

➢ NLP-based text mining approach.

➢ Content-rich relationship networks among biological

➢ Concepts, genes, proteins or drugs.

➢ Nature of the relationship: inhibitory, stimulative, neutral

and simple co-occurrence.

➢ Internet-based application with graphical visualisation

➢ Sentence as unit, POS tagging, shallow parsing and rules.

Chen H, Sharp BM.Content-rich biological network constructed by mining PubMed abstracts.BMC Bioinformatics. 2004 Oct 8;5(1): 147.

http://www.chilibot.net/

# Chilibot (2)



➢ Need registration.

➢ Hypothesis generation.

Chen H, Sharp BM.

Content-rich biological network constructed by mining PubMed abstracts.

BMC Bioinformatics. 2004 Oct 8;5(1):147.

http://www.chilibot.net/

- Based on SVM.
- Query protein or accession number.
- Assist the Biomolecular Interaction Network Database (BIND)

Donaldson I, Martin J, de Bruijn B, Wolting C, Lay V, Tuekam B, Zhang S, Baskin B, Bader GD, Michalickova K, Pawson T, Hogue CW.PreBIND and Textomy--mining the biomedical literature for protein-protein interactions using a support vector machine.*BMC Bioinformatics*. 2003 Mar 27;4(1):11.

http://www.blueprint.org/products/prebind

**Text mining biomedical literature (2005)**

# Microarray data analysis

➢ Co-ordinated expression of genes.

➢ Functional co-regulation within biological processes.

➢ Mine micro array data using the associated

biomedical literature.

➢ Characterise groups of genes extracting functional keywords.

➢ Score the coherence of gene clusters.

➢ Group genes based on their associated literature and

functional descriptions.

# GEISHA

➢ Text mining tool for microarray analysis.

➢ Analyse the correlation between:

  ➢ the increase of the level of expression patterns and

  ➢ the significance of functional information derived

    from the literature.

➢ Extract functional information from the literature linked

  to the microarray genes.

➢ Calculates statistical significance of terms from

  documents associated to genes of each cluster.

# Protein localization

- Protein activity -> specific cellular environments.

- Localisation determination:

    - Experimental techniques.

    - Bioinformatics techniques (PSORT).

    - Text mining.

- Nair and Rost: lexical information in annotation database records.

- Stapley et al: Use SVM to classify proteins according to their subcellular localisation, extracted from PubMed abstracts.

# NLP and sequence analysis: MedBlast

➢ Use NLP techniques to retrieve the related articles

for a given sequence (online).

➢ Related articles:

  ➢ those describing the query sequence (protein) or

  ➢ Its redundant sequences and close homologues

➢ Direct search with the sequence.

➢ Indirect search with gene symbols.

➢ Use Blast against GenBank.

➢ Use Eutilities toolset to retrieve documents

http://medblast.sibsnet.org/

# NLP and sequence analysis: SAWTED



Sequence similarity

the base for identifying

structure templates

for query sequence

Structure Assignment

With Text Description

Document comparison

algorithms

http://www.bmm.icnet.uk/~sawted/

**Text mining biomedical literature (2005)**

Use information contained in text descriptions of SwissProt annotations

identification of remote homologues

http://www.bmm.icnet.uk/~sawted/

**Text mining biomedical literature (2005)**

# Community wide evaluations

## BIOINFORMATICS

- CASP
- CAMDA
- CAPRI
- GASP
- GAW
- PTC

## BIO-NLP

- KDD cup
- BioCreative
- JNLPBA shared task
- TREC Genomics track
- LLL05 challenge

## NLP/IR/IE

- MUC
- TREC

CASP: Critical asessment 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

**Text mining biomedical literature (2005)**

# Overview: BioCreative tasks

**MITRE/ NCBI**

**BioCreative**

**PDG-CNB/ GOA-EBI**

**Task 1: NER (Protein Tagging)**

**Task 2: Automatic annotation**

**Sub-task 1: Gene mention finding**

**Sub-task 2: Normalized gene list**

**Sub-task 1: Annotation passage retrieval**

**Sub-task 2: Annotation prediction**

**Sub-task 3: Ad Hoc Retrieval**

NCBI: Nat. Cen. for Biotech. Inf.
GOA: Gene Ontology Annotation
EBI: European Bioinf. Institute
CNB: Centro Nacional de Biotecnologia

**Text mining biomedical literature (2005)**

# BioCreative – Why?

- Open evaluation to determine the state of the art.

- Compare the performance of different methods.

- Produce a gold standard training set.

- Monitor improvements in the field.

- Produce useful evaluation tools/metrics.

# BioCreative Task 1.1 summary

- Finding gene mentions in abstracts (NER).

- 15 teams, 3-4 submissions per team.

- Data and evaluation software provided by the NCBI.

- Performance: over 80% F-score (balanced precision and recall).

- Top scoring participants used some type of markov modelling (ME,HMM,CRF), SVM or manual rules.

# BioCreative Task 1.2 summary

- Gene identifier list task.

- 8 teams, 3 submissions per team.

- Given an abstract from a specific model organism (Fly, Mouse, yeast) create the list of unique gene identifiers.

- F-score: yeast 0.92, fly 0.82 and mouse 0.79.

- Difficulties: ambiguity, complex names, distinguish between multiple identifiers.

- Methods: matching against lexical resources (e.g. exhaustive matching) or task 1.1 type systems.

# BioCreative Task 2 description

- Automatic extraction and assignment of GO annotations for human proteins using full text articles.

- Based on triplets: protein – GO term – article passage.

- Task 2.1: Passage retrieval task, find the text passage which support a protein – GO term annotation.

- Task 2.2: text categorization task, predict protein – GO term associations and the corresponding text passage.

- Task 2.3*: ad hoc information retrieval, retrieve annotation relevant articles

# Data sets and evaluation strategy

- GO: Gene Ontology: controlled vocabulary (concepts) within

    an ontology (DAG), 3 categories, MF: Molecular function,

    BP: Biological Process and CC: Cellular Component.

- GO concepts used for annotation purposes: GOA.

- Training set: 803 GOA derived full text articles from JBC journal.

- Test set: 113 articles for task 2.1 and 99 for 2.2 and triplets.

- Triplets: GO-term - protein – article -> return passage (task 2.1).

- Evaluation by GOA annotators from the EBI.

- Manually evaluation of the predicted passages within its context

    in the paper using a highlighting tool.

- Evaluation types: High:correct, Generally: OK but to general for

    practical use and Low:wrong.

# BioCreative Task 2 data set

| Description | Training set | Test set 2.1 | Test set 2.2 |
|---|---|---|---|
| Full text articles | 803 | 113 | 99 |
| Total of GO annotation | 2317 | 1076 | 1227 |
| Nr of proteins in the GO annot | 939 | 138 | 138 |
| Nr GO terms used for annot | 776 | 580 | 544 |
| Average nr of annot/protein | 2.467 | 7.797 | 8.891 |
| Annotations of MF GO terms | 709 | 330 | 356 |
| Annotations of BP GO terms | 1061 | 544 | 701 |
| Annotations of CC GO terms | 547 | 182 | 170 |
| MF terms in the annotations | 343 | 173 | 179 |
| BP terms in the annotations | 339 | 334 | 314 |
| CC terms in the annotations | 94 | 57 | 51 |

# Data sets and evaluation strategy

- GO: Gene Ontology: controlled vocabulary (concepts) within

    an ontology (DAG), 3 categories, MF: Molecular function,

    BP: Biological Process and CC: Cellular Component.

- GO concepts used for annotation purposes: GOA.

- Training set: 803 GOA derived full text articles from JBC journal.

- Test set: 113 articles for task 2.1 and 99 for 2.2 and triplets.

- Triplets: GO-term - protein – article -> return passage (task 2.1).

# BioCreative Task 2.1 sample submission

```xml
<protein>
        <namefile>JBC_2001-2/bc4501042445.gml</namefile>
        <idTask>2.1</idTask>
        <participant>user14</participant>
        <nameProtein></nameProtein>
        <dbId>O15023</dbId>
        <sourceDb>Swiss-Prot</sourceDb>
        <goCode>
                <name>phosphatidylinositol binding</name>
                <code>0005545</code>
                <evidenceText>In addition, a single point mutation in the
FYVE finger motif at cysteine residue 753 (C753S) is sufficient to
abolish its endosomal association. Its endosomal localization is also
sensitive to the phosphatidylinositol 3-kinase inhibitor, wortmannin.
Using in vitro liposome binding assays, we demonstrate that Myc-tagged
endofin associates preferentially with phosphatidylinositol 3-
phosphate, whereas the C753S point mutant was unable to do so. We also
show that endofin co-localizes with SARA but that they are not
associated in a common complex because they failed to co-
immunoprecipitate in co-expressing cells.</evidenceText>
        </goCode>
<protein>
```

# BioCreative Task 2 participating systems

- 8 groups, max. 3 runs.

- Three strategies:

  (1) GO term centred, pattern matching, GO words

  (IC, word weight), recall centred

  (2) Machine learning techniques.

  (3) High precision, pattern matching and template

  extraction

- Tendency: sentence level, pattern matching, regular expressions and

  use of external resources (e.g. HuGO, UMLS), but:

- In general little overlap between the methods and the used resources.

# BioCreative Task 2.1 results



Task 2.1: TP vs precision

**TP: prediction evaluated as protein   and GO terms correct**

**Precision: TP / Total nr. of**
                              **evaluated**

**submissions**

**Teams:**
**1: Chiang et al.**
**2: Couto et al.**
**3: Ehrler et al.**
**4: Krallinger et al.**
**5: Krymolowski et al.**
**6: Ray et al.**
**7: Rice et al.**
**8: Verspoor et al.**

# BioCreative Task 2.1 examples

| | |
|---|---|
| Query_id | Q96PH1_0000910_11483596 |
| PMID: | 11483596 |
| UserId: | user20_1 |
| UserName: | Couto et al. |
| ProteinName: | NADPH oxidase 5 gamma |
| AccessionNr: | Q96PH1 |
| GO_term: | **cytokinesis** |
| GO_id: | 0000910 |
| EvalProtein: | high |
| EvalGO: | high |
| EvalAnnot: | highhigh |
| Len_GO: | 1 |
| GO_cat: | P |

**EvidenceText:** Thus, <GLOSREF RID="G8">NOX5</GLOSREF> might have a function in the early stages of spermatogenesis such as cell division, induction of apoptosis, or DNA compaction.

# BioCreative Task 2.2 results



Task 2.2: TP vs. precision

**TP: prediction evaluated as protein and GO terms correct**
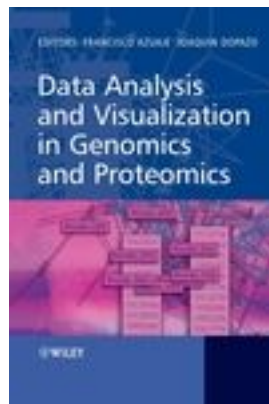
**Precision: TP / Total nr. of evaluated submissions**

1: Chiang et al.
2: Couto et al.
3: Ehrler et al.
4: Ray et al.
5: Rice et al.
6: Verspoor et al.

# Selected review references

R. Hoffmann, M. Krallinger, E. Andres, J. Tamames, C. Blaschke and A. Valencia. Text Mining for Metabolic Pathways, Signaling Cascades, and Protein Networks. Science STKE 283, pe21 (2005).

M. Krallinger, R. Alonso-Allende Erhadt and A. Valencia. Text-mining approaches in molecular biology and biomedicine. Drug Discovery Today 10, 439-445 (2005).

M. Krallinger and A. Valencia. Applications of Text Mining in Molecular Biology, from name recognition to Protein interaction maps. In Data Analysis and Visualization in Genomics and Proteomics, chapter 4, Wiley.

# Selected links

http://www.pdg.cnb.uam.es/martink/LINKS/bionlp_tools_links.htm

http://www.pdg.cnb.uam.es/martink/links.htm

# Acknowledgements

➢ To the audience for paying attention

➢ To Alfonso Valencia for his supervision, support and

  discussions.

➢ To Belen Bañeres for organisational aid.

➢ The the Protein Design group at CNB and

  especially Maria Padron for discussions and

  suggestions.