# Text Mining Biomedical Repositories

Ashish Vijay Tendulkar

School of Technology and Computer Science,
Tata Institute of Fundamental Research, Mumbai.
http://www.tifr.res.in/∼ashishvt/

Dec 20, 2011

# Acknowledgement

The tutorial slides were jointly prepared with Martin Krallinger of Spanish National Cancer Research Center, Madrid.

# Table of Contents

# Biology 101

## Biological Entities

- Genes, Genome
- mRNA, transcriptome
- Protein, Proteome
- Cell
- Tissues
- Organisms

# Biomedical Literature

## Generation and Form

- Biologists conduct experiments and generate heterogeneous data types such as sequence, structure, expression etc.
- A large fraction in form of Natural Language (NL) report of experimental findings in form of research papers, reports, patents, newswire articles
- Structured database records (in relatively less proportion than NL)(e.g. Genes/Protein sequence, structure etc.)

# Biomedical Literature

## Generation and Form

- ▶ Biologists conduct experiments and generate heterogeneous data types such as sequence, structure, expression etc.
- ▶ A large fraction in form of Natural Language (NL) report of experimental findings in form of research papers, reports, patents, newswire articles
- ▶ Structured database records (in relatively less proportion than NL)(e.g. Genes/Protein sequence, structure etc.)

## Importance

- ▶ Curation of structured databases (UniProt)
- ▶ Deriving functional annotations beyond what is present in DBs.
- ▶ Contextual information about experimental results and conditions(Cell lines, tissues, etc.)

# Literature and Scientific Discovery Process

## Biology

- 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

## Clinics

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

# Literature and Scientific Discovery Process

### Pharma

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

### Funding

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

### Publishing Groups

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

# Relevance of Literature in Bioinformatics

# Challenges in Exploring Biomedical Literature

Rapid growth of literature data poses following challenges:

- ▶ Efficient methods for extraction of information
- ▶ Effective ways of querying the information

# Curation of Biological Databases from Literature

## Classical Method: Manual Curation

- Trained human experts reads scientific literature and extracts information of interest
- Manual time consuming and labor intensive process
- Accurate through human inference and background knowledge
- Example DBs: Uniprot, GOA, SGD, MGI etc.

## Text Mining assisted Curation

- Retrieval of relevant literature from literature repositories
- Textual evidence and entity detection
- Revision and editing of manual records
- E.g. TextPresso, Rodriguez-Penagos et al (gene regulation), Grover el at (PPI), Chang et al (Pathways), Ongenaert et al (methylation), Shtatland (peptides), Miotto (allergen cross-reactivity).

# Overview of Current Literature Repositories

- e-Books: NCBI Bookshelf
- Citation of Biomedical Research Articles + Abstract: PubMed
- Full text research articles:
  - PubMed Central (PMC)
  - Highwire Press
  - BioMed Central

# PubMed

## Overview

- Developed by NCBI
- Citation entries of scientific articles of all biomedical sciences
- Each entry is characterized by a unique identifier, the PubMed identifier: PMID
- Often links to the full text articles are displayed

## Statistics

| | |
|---|---|
| No. of Citations | 16 million |
| No. of Indexed Journals | approx. 5000 |
| No. of English Articles | 12 million |
| No. of Articles with Abstracts | 7,000,000 |

# Importance of PubMed in Biomedical Text Mining

- Approximately 1 million entries 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

# Entrez



Figure: Source http://www.ncbi.nlm.nih.gov/pubmed/

# PubMed Search Results



Figure: Source http://www.ncbi.nlm.nih.gov/pubmed/

# PubMed XML Record



```xml
<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>
```

Figure: Source http://www.ncbi.nlm.nih.gov/pubmed/

# PubMed Query Translation



Figure: Source http://www.ncbi.nlm.nih.gov/pubmed/

# PubMedCentral

- Digital archive of full text life science journals
- Articles have a unique PMCID
- Allows Boolean query search
- Offers free full text articles
- Journal Publishing XML DTD, but also other widely used DTD in life science

# Example PubMedCentral Query



Figure: Source:http://www.ncbi.nlm.nih.gov/pmc/

# PubMed Journals



Figure: Source:http://www.ncbi.nlm.nih.gov/pmc/

# NCBI Bookshelf

- Collection of biomedical text books
- Allows boolean query searches
- Offers free full text articles
- Direct searching the books or from PubMed abstract

# Retrieving Electronic Literature Data from Web

- Get a local copy of some centralized literature repository (PubMed, PubMed Central, journals, etc): Leasing PubMed
- Use literature retrieval modules:
  - BioPython/BioPerl: Gazelle Z39.50 interface to PubMed
  - pubmed.pm by J. Smyser
  - Pubmed crawler written in Perl (http://pubcrawler.ie/)
  - eUtils: Entrez programming utilities
- Adopt web crawler, spiders or focussed crawler such as DataparkSearch, GNU Wget, Heritrix, Mutch.

# Preprocessing Scientific Articles

1. **Document Standardization**: variety of formats (ASCII, HTML, XML, PDF, scanned PDF, SGML), convert them into a common format and encoding.
2. XML /Extensible Markup language, standard way to insert tags onto a text to identify its parts
3. OCR (Optical Character Recognition), used to digitize older literature (PMC Back Issue Digitization initiative)
4. Recover article Structure and content using pdftotext, PDFLib,PDF Converter

# Preprocessing Scientific Articles

1. **Tokenization** break a stream of characters into words (tokens), e.g. white space, special chars. Each token is an instance of a type

2. **Stemming and lemmatization** standardize word tokens (e.g. Morphological analysis and Inflectional stemming, convert words to their corresponding root form)

3. Lexical analysis of the text with the objective of treating digits, hyphens, punctuation marks, and the case of letters

4. Elimination of stop-words

5. Selection of index terms

# Journal Specific Characteristics of Literature

- Journal/Article format
- Paper structure (Section types)
- Article types (Review, Clinical Study, etc.)
- Target audience of journal/article

## Processing of Full Text Articles

- Extract title, authors, abstract, text body, references
- Extract tables and tables legends
- Extract figures and figure legends

# Basic Features of Biomedical Literature Data

- Heavy use of domain specific terminology (12% biochemistry related technical terms). E.g. chemoattractant, fibroblasts, angiogenesis
- Polysemic words (Word Sense Disambiguation). For example, APC means either:
  - Argon Plasma Coagulation
  - Activated Protein C

  Teashirt means:
  - a type of cloth
  - tsh gene
- Heavy use of acronyms, e.g. Activated protein C (APC), or vascular endothelial growth factor (VEGF)
- Data sparseness: Many words occur with low frequency

# BioTerms Characteristics

- Novelty: New names and terms are frequently created. E.g. `This disorder maps to chromosome 7q11-21, and this locus was named CLAM.`[PMID:12771259]
- Typographical variants. E.g. TNF-Alpha and TNF Alpha
- Different writing styles
- 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.`

# Biomedical Corpora and Text Collections

- Medtag corpus includes 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 of 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.

# 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. E.g. noun, verb, adjective, adverb or preposition. E.g. [PMID 12700631]

| Token | POS |
|-------|-----|
| 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. |

# POS Taggers

- Programs to label words with POS
- 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 POS tagger for biomedical domain. MedPost: 97% accuracy in PubMed abstracts (86.8% general POS tagger)

# GENIA POS Tagger



Figure: Source: http://www-tsujii.is.s.u-tokyo.ac.jp/GENIA/tagger/

# GENIA POS Tagger Output



Figure: Source: http://www-tsujii.is.s.u-tokyo.ac.jp/GENIA/tagger/

# 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
- Problem: collapse two semantically different words, e.g: gallery and gall.

# Online Stemmer: SnowBall



Figure: Source: http://snowball.tartarus.org/demo.php

# 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).
- Output of Connexor shallow parser

| Token | Syntactic Role |
|-------|----------------|
| 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
  - verbal phrases (head is a verb, VP).

# Syntactic Features

# Syntactic Features

- Identification of subject-object relationships
- NP-VP-NP E.g.
  ```
  Overexpression of <gene>IMEl</gene> induced an
  <GO> early meiotic event (recombination) </GO> in
  rich medium, but later meiotic events did not
  occur (i.e., they detected [no spore formation])
  ```
  Subject: IMEl gene and object is GO term early meiotic event

# Semantic features

- Associations of words with their corresponding meaning in a given context.
- Semantics (meanings) of a word $\rightarrow$ 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.)

| Token | GO Symbols |
| --- | --- |
| Caspase-3 | GENE PRODUCT |
| was | |
| partially | |
| activated | INTERACTION VERB |
| by | |
| IFN-gamma | GENE PRODUCT |

# Contextual Features

- Words occurrence in textual context - association.
- Co-occurrence of Caspase-3 and INF-gamma in the same sentence indicates some relationship between them.
- Determine contextual similarity of proteins documents.
- Use for instance: list of words (bag of words)
- The statistical analysis of word frequencies or patterns
- Features are interrelated

Part 2: BioText Mining

# Key Technologies in BioText Mining

- Information Retrieval (IR)
- Information Extraction (IE)
- Text Classification
- Text Clustering

# Information Retrieval

## IR

- It is a process of recovery of those documents from a collection of documents which satisfy a given information demand.
- Information demand is posed in form of a query

Efficient Indexing is required to reduce vocabulary of terms and query formulation.

## Important Steps in Indexing Document Collection

- Tokenization
- Case folding
- Stemming
- Stop word removal

Query Types:
- Boolean queries

# Zipf's Law

# Boolean Queries

- Based on combination of terms using Boolean operators
- Basic Boolean operators: **AND**, **OR**, **NOT**
- Queries matched against the terms in the inverted index file
- E.g. Entrez - Boolean search in PubMed
- Fast and easy to implement

# Vector Space Model

- Measure similarity between query and documents
- Query may be a list of terms or even whole documents
- Represent document and query using a vector of terms.
- Each term $t$ is weighted according to its frequency in the document $d$ and in the whole document collection $D$.
- Calculate cosine similarity between query and document vector.
- Return ranked list of documents
- E.g. Related article search in PubMed

# eTBLAST



Figure: Source: http://invention.swmed.edu/etblast/index.shtml

# eTBLAST



Figure: Source: http://invention.swmed.edu/etblast/index.shtml

# eTBLAST



Figure: Source: http://invention.swmed.edu/etblast/index.shtml

# eTBLAST Results



Figure: Source: http://invention.swmed.edu/etblast/index.shtml

# eTBLAST



Figure: Source: http://invention.swmed.edu/etblast/index.shtml

# Text Similarity and Deja Vu



Figure: Source: http://spore.swmed.edu/dejavu/

# Text Similarity and Deja Vu



Figure: Source: http://spore.swmed.edu/dejavu/

# IR Evaluation

- 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

# 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.

# Text Clustering

- Clustering programs often choose first the documents that seem representative of the middle of each of the clusters (candidate centers of the clusters).
- Then it compares all the documents to these initial representatives.
- Each documents is assigned to the cluster it is most similar to.
- Similarity is based on how many words the documents have in common, and how strongly they are weighted.
- The topical terms of the clusters are chosen from words that represent the center of the cluster.
- The best clustering is one in which the average difference of the documents to their cluster centers smallest.
- Agglomerative clustering: first comparing every pair of documents, and finding the pair of documents which are most similar to each other.

# Text Classification

- Common problem in information science.
- Assignment of an electronic document to one or more categories, based on its contents (words).
- Supervised document classification where training examples of document classification are provided and the correct classification model is learnt based on one of the following techniques:
  - 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

# Information Extraction (IE)

IE refers to the automatic extraction of structured information such as entities, relationships between entities, and attributes describing entities from unstructured sources [1]

## Applications of IE

- ► Enterprise applications such as news tracking, customer care, data cleaning, classified Ads
- ► Personal Information Management such as Emails, Documents, Presentations
- ► Scientific applications such as BioIE
- ► Web oriented applications such as citation databases, opinion databases, community websites, comparison shopping, ad-placement

---

[1]Sarawagi, S. (2008) Information Extraction, FnT Databases, 1(3).

# IE Taxonomy

1. Structure extracted (entities, relationships, lists, tables, attributes etc.)
2. Unstructure source (short strings or documents, templatized or open ended.)
3. Input resources available for extraction (structured databases, labeled unstructured data, linguistic tags, dictionaries etc.)
4. Extraction method (rule based or statistical, manually coded or trained from examples)
5. Output of extraction (annotated unstructured text or a database)

# Key Challenges in IE

## Accuracy

- Diversity of clues such as orthographic features of words, POS, similarity with existing entries in database, presence of trigger words and so on

- Difficulty of detecting missed extractions: High recall models are desirable, but how do we ensure high recall without extensive labeled data?

- Increased complexity of the structures extracted: Extraction of longer entities where the boundary is not defined clearly. e.g. extracting name of restaurant from Blog.

# Key Challenges in IE

## Running Time

- Efficiently filtering right subset of documents that are likely to contain structured information of interest
- Efficiently locating portion of document containing relevant information of interest
- Efficient extraction of information

## Other Systems Issues

- Dynamically Changing Sources
- Data Integration
- Extraction Errors

Part 3: Applications of Text Mining

# Applications of BioText Mining

- Named entity recognition of biological entities (BioNER)
- Gene normalization
- Protein-Protein interaction
- Functional Analysis of genes and gene sets
- Extraction of gene-disease association
- Extraction of mutations and epigenetic characteristics
- Extraction of protein location information
- Building terminology resource for a specific domain
- Knowledge discovery and pathways

# Bio-Named Entity Recognition (BioNER)

## Objective

Identify biological entities in articles and to link them to entries in biological databases.

## Challenges

- more complex (synonyms, disambiguation, typographical variants, official symbols not used)
- Performance organism dependent

## Methods

- POS tagging,
- Rule-based,
- Flexible matching,
- Statistical and Machine Learning (naive Bayes, ME, SVM, CRF, HMM).

# BioNER Example

### Input: Abstract/Full Text Article

The distribution of myeloid lineage-associated cytokine receptors and lysosomal proteins was analyzed in human CD34+ cord blood cell (CB) subsets at different stages of myeloid commitment by reverse-transcriptase polymerase chain reaction (RT-PCR).

### Output: Input + Tagged Entities

The distribution of myeloid lineage-associated cytokine receptors and lysosomal proteins was analyzed in human CD34+ cord blood cell (CB) subsets at different stages of myeloid commitment by reverse-transcriptase polymerase chain reaction (RT-PCR).

Protein Cell Line

# BioNER Challenges

- Authors often do not use the official gene symbols
- Genes have often synonyms.
- Use of full gene names and/or gene symbols/acronyms
- Gene names - medical terms ambiguity
- Gene names - common English words ambiguity (fly)
- Alternative typographical variants
- 14% of genes display inter-species ambiguity (Chen, 2005).
- Ambiguity between protein names and their protein family names
- Identification of new gene names (novel genes)

# Tricky Issues in Gene Tagging

- The nightcap mutation caused severe defects in these cells [PMID:12399306].

- 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]

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

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

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

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

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

# GapScore

- Scores words based on a statistical model of gene names
- Quantifies: Appearance, Morphology, Context.
- Performance on Yapex corpus:

| Match Type | Precision | Recall | F1 Score |
|---|---|---|---|
| Partial Match | 81.5% | 83.3% | 82.5% |
| Exact Match | 56.7% | 58.5% | 57.6% |

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

# GapScore

**Welcome to our Gene and Protein Name Server!**

We have developed a method GAPSCORE that will scan text and identify the names of genes and proteins. This tool has many potential applications including: allowing users to search and index documents by genes of interest, analyzing the scientific literature for genes of interest, and automatically building knowledge bases from text.

We use a machine learning-based approach described in Chang JT, Schütze H, and Altman RB. *GAPSCORE: Finding Gene and Protein Names One Word at a Time*. In preparation.

Search for gene and protein names in some text:

```
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
```

[SEARCH]

You can access this functionality from computer programs using XML-RPC. We have provided bindings for Python, PERL, and Java. There is code and documentation at http://bionlp.stanford.edu/webservices.html.

# GapScore Output

| | Gene or Protein Name | Quality (Score) |
|---|---|---|
| 1 | human BRCA1 | Excellent (1.00) |
| 2 | human BRCA2 | Excellent (0.97) |
| 3 | human BRCA2 | Excellent (0.97) |
| 4 | human BRCA2 | Excellent (0.97) |
| 5 | Brca2 | Excellent (0.91) |
| 6 | Brca2 | Excellent (0.91) |
| 7 | Brca2 | Excellent (0.91) |
| 8 | murine Brca2 | Excellent (0.91) |
| 9 | murine Brca2 | Excellent (0.91) |
| 10 | murine Brca2 | Excellent (0.91) |
| 11 | murine Brca2 | Excellent (0.91) |
| 12 | murine Brca2 | Excellent (0.91) |
| 13 | Brca1 | Excellent (0.90) |
| 14 | RAD51 | Good (0.62) |
| 15 | Rad51 | Good (0.55) |
| 16 | protein encoded | Good (0.14) |
| 17 | dispensable | Poor (0.10) |
| 18 | similar carboxyl-terminal truncating mutations | Poor (0.09) |
| 19 | differences | Poor (0.09) |
| 20 | hypomorphic | Poor (0.09) |

# iHOP: Gene Name Tagging

# ABNER



- ABNER is a software tool for molecular biology text analysis.
- It uses linear chain conditional random fields (CRFs) with a variety of orthographic and contextual features.
- URL: http://pages.cs.wisc.edu/ bsettles/abner/

# Gene Normalization

## Objective

- Linking genes or gene products mentioned in the literature to biological databases.
- Key step in enabling accurate search of biological literature and linking database information to passages in research articles.

## Challenges

- Genes are often described rather than referred to by gene symbol
- One gene name may refer to different genes (often from different organism)
- Incomplete dictionaries of gene names

# Example

## Input: Abstract/Full Text Articles with Entites Tagged

The double-stranded (ds) RNA-activated protein kinase from human cells is a 68 kd protein (p68 kinase) induced by interferon. On activation by dsRNA in the presence of ATP, the kinase becomes autophosphorylated and can catalyze the phosphorylation of the alpha subunit of eIF2, which leads to an inhibition of the initiation of protein synthesis.

## Gene Normalized to Database Records

| Entity | Normalized Gene Name | UniProt ID |
|--------|---------------------|------------|
| RNA-activated protein kinase | E2AK2_HUMAN | P19525 |
| p68 kinase | E2AK2_HUMAN | P19525 |
| alpha subunit of eIF2 | IF2A_HUMAN | P05198 |

# iHOP System



Figure: Source: http://www.ihop-net.org/

# iHOP: Query to DB Records



Figure: Source: http://www.ihop-net.org/

# iHOP: Information about Gene

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 wall 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]

Figure: Source: http://www.ihop-net.org/

# Interaction Extraction

## Objective

Extract interaction information between biological entities from literature. For example, protein-protein interaction.

## Key Techniques

- Co-occurrence of bioentities within close vicinity
- Machine learning based methods (Relationship extraction)
- Linguistic methods (Dependency parsers, link parsers)
- Rule based

## Manually curated Interaction Databases

- MINT
- BioGRID
- IntAct

# iHOP: Interaction Information



Figure: Source: http://www.ihop-net.org/

# iHOP: 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 pathogenic mutation in **BRCA1** ☆ or **BRCA2** ☆. [2007]

**Figure:** Source: http://www.ihop-net.org/

# iHOP System: Gene Model/Graph



Figure: Source: http://www.ihop-net.org/

# PLAN2L: Plant Annotation to Literature

- Web tool for integrated text mining and literature-derived bio-entity relation extraction
- Provides following searches
  - Searching *Arabidopsis* bibliome
  - Searching for PPI
  - Searching for gene regulation association
  - Searching for location sentences
  - Searching for cell cycle association
  - Association retrieval

# PLAN2L: PPI Extraction



Figure: Source:http://zope.bioinfo.cnio.es/plan2l/plan2l.html

# EBIMed



Figure: Source:http://www.ebi.ac.uk/Rebholz-srv/ebimed/index.jsp

# EBIMed



Figure: Source:http://www.ebi.ac.uk/Rebholz-srv/ebimed/index.jsp

# EBIMed



Figure: Source:http://www.ebi.ac.uk/Rebholz-srv/ebimed/index.jsp

# InfoPubMed

- Info-PubMed provides information from Medline on protein-protein interactions.
- Given the name of a gene or protein, it shows a list of the names of other genes/proteins which co-occur in sentences from Medline, along with the frequency of co-occurrence.
- Uses information extraction techniques to identify interacting entities.
- URL: http://www-tsujii.is.s.u-tokyo.ac.jp/info-pubmed/

# InfoPubMed



Figure: Source:http://www-tsujii.is.s.u-tokyo.ac.jp/info-pubmed/

# PPI Finder: Extraction of Human PPI



Figure: Source:He et. al. (2009) PLoS ONE 4(2):e4554

# STITCH

## Salient Features

- A resource to explore known and predicted interactions of chemicals and proteins.
- Chemicals are linked to other chemicals and proteins by evidence derived from experiments, databases and the literature.
- Contains interactions for over 74,000 small molecules and over 2.5 million proteins in 630 organisms.

## Search Types

- Protein/chemical name
- Chemical structure
- Protein sequence
- Multiple names and sequences

STITCH URL: http://stitch.embl.de/

Figure: Source: Kuhn et. al. (2009) Nucleic Acid Res.:Database Issue

# Functional Analysis of Genes and Gene Sets

### Objective

Extract information about gene/protein function and gene sets.

### Key Techniques

- ▶ Rule based
- ▶ Dictionary based

# SciMiner

## Salient Features

- ▶ Web based literature mining and functional analysis tool
- ▶ Identifies gene and protein names via context specific analysis of MEDLINE abstract
- ▶ Accepts query in form of a list of PubMed IDs or Entrez style free text search
- ▶ Scans biomedical literature for gene/protein of user's interest
- ▶ Finds significant enrichment in
  - ▶ Target gene list
  - ▶ GO terms
  - ▶ Mesh Terms
  - ▶ Pathways
  - ▶ PPI Network

http://jdrf.neurology.med.umich.edu/SciMiner/

# Chilibot



Figure:  Source:http://www.chilibot.net/

# Chilibot

- The **BRCA2** homologue Brh2 nucleates **RAD51** filament formation at a dsDNA ssDNA junction. Ref: Nature, 2005

- Identification of **Rad51** regulation by **BRCA2** using Caenorhabditis elegans **BRCA2** and bimolecular fluorescence complementation analysis. Ref: Biochem Biophys Res Commun, 2007

- Human **BRCA2** interacts with the recombinase **RAD51** via eight BRC repeats. Ref: Proc Natl Acad Sci U S A, 2007

- It is known that **BRCA2** interacts directly with **RAD51** through a series of degenerative motifs known as the BRC repeats. Ref: Philos Trans R Soc Lond B Biol Sci, 2004

- **BRCA2** has important roles in **RAD51** focus formation and HRR of DNA double strand breaks ( DSBs ). Ref: Mol Cell Biol, 2005

- These results show that **BRCA2** repeats mimic the **RAD51** PM and imply analogous **RAD51** interactions with RAD52 and RAD54. Ref: EMBO J, 2003

- This modification blocks C terminal interactions between **BRCA2** and **RAD51**. Ref: Nature, 2005

- Here we have used the yeast two hybrid system to test for direct interaction between **BRCA2** or its effector **RAD51** and the FANCA, FANCC and FANCG proteins. Ref: Hum Mol Genet, 2003

Figure: Source:http://www.chilibot.net/

# Chilibot



Figure: Source:http://www.chilibot.net/

# GoPubmed



Figure: Source: http://www.gopubmed.org/

# Extraction of Gene-Disease Association

### Objective

Given a biomedical research article/abstract, extract gene disease associations.

### Tools

- Facta: Finding Associated Concepts with Text Analytics

# Facta Query



Figure: Source:http://text0.mib.man.ac.uk/software/facta/main.html

# Facta Relevant Concepts



Figure: Source:http://text0.mib.man.ac.uk/software/facta/main.html

# Facta Medline Search



Figure: Source:http://text0.mib.man.ac.uk/software/facta/main.html

# Extraction of Mutations and Epigenetic Characteristics

## Objective

Given a biomedical research article/abstract, extract mutation and epigenetic characteristics.

## Challenges

- Deluge of experimental data from high throuput screens such as microarray and RNAi
- Traditionally hundreds of genes are clustered via enrichment of GO terms.
- Analysis results in GO annotations and not always available for all genes in the model organism
- Deluge of literature data calls for high through put techniques to extract mutations and epigenetic characteristics

# GoGene

## Salient Features

- Extract co-occurrences of genes and ontology terms from literature
- Combines disease, compounds, techniques, and mutations information
- Claims to provide most recent facts about genes and rank them according to novelty and importance

| No. of Associations | 4,000,000 |
|---|---|
| No. of Model Organisms | 10 |
| No. of PubMed Articles | 18,000,000 |

## Query Types

PubMed, Entrez, Sequence

http://gopubmed2.biotec.tu-dresden.de/gogene/gogene/

# GoGene PubMed Query



Figure: Source:http://gopubmed2.biotec.tu-dresden.de/gogene/gogene/

# GoGene Entrez Query



Figure: Source:http://gopubmed2.biotec.tu-dresden.de/gogene/gogene/

# GoGene Sequence Query



Figure: Source:http://gopubmed2.biotec.tu-dresden.de/gogene/gogene/

# MeInfoText: Salient Features

- Extracts associated gene mythylation and cancer information from biomedical text and integrates it with biological pathways and protein protein interaction.
- DNA methylation, occurring predominantly in CpG islands, is an important epigenetic modification of the genome that is involved in mediating various cellular processes (Robertson, 2005).
- Abnormal methylation of DNA may result in increased transcription of oncogenes or silencing of tumor suppressor genes and is common in a variety of human cancer cells (Esteller, 2005).

# MeInfoText: Search Types

- Associations among gene, methylation and cancer
- Gene methylation associations
- Profile of gene methylation across human cancer types
- Gene methylation of a specific cancer type
- http://mit.lifescience.ntu.edu.tw/

# Extraction of Protein Location from Literature

### Objective

Given a biomedical research article/abstract, extract protein localization information

### Significance

The role of proteins in biochemical reactions depends on its location.

# Location extraction from PLAN2L



Figure: Source:http://zope.bioinfo.cnio.es/plan2l/plan2l.html

# Automatic Construction of Lexical Resources

- Automatic construct lexical resource such as dictionaries, thesauri from biomedical literature
- Rapidly expanding biomedical literature makes it hard to manually maintain dictionaries and thesauri
- High throughput and automated methods are required for creation of lexical resources
- Highly domain specific and virtually open terminology used in biomedical domain makes it challenging to extract appropriate lexical resources

# Termine

- Uses C-value term extraction and Acromine acronym recognition
- Statistical analysis assigns termhood to a candidate term using the following characteristics:
  - occurrence frequency of the candidate term
  - frequency of the candidate term as part of other longer candidate terms
  - number of these longer candidate terms
  - length of the candidate term
- URL: http://www.nactem.ac.uk/software/termine/

# Acronym Finder

- Abbreviation is a short form of a word or a phrase.
- Identification of correct abbreviation and its long form pair is crucial for IR and IE applications
- Fast rate of growth of biomedical literature makes it hard for thesauruses to keep track of all abbreviations
- Abbreviations are of two types:
  - Acronyms: Word formed by initial letters or letters of each successive parts or major parts of long form. E.g. CKB stands for Brain Creatin Kinase
  - Non-acronyms do not follow lexical patterns with long forms. E.g. 11p stands for the short arm of chromosome 11

## ADAM: Another Database of Abbreviations in MEDLINE

- URL: http://128.248.65.210/arrowsmith_uic/adam.html
- Input: Short form or long form
- Output: Corresponding long form or short form

# MedlineRanker

- Flexible ranking system for MEDLINE abstract
- Given an abstract related to a specific topic, MedlineRanker returns the most discriminative words in comparison with a random selection
- These words are used to score other abstracts.
- URL: http://cbdm.mdc-berlin.de/tools/medlineranker

Part 4: Issues and Challenges in Evaluation of Text Mining Systems

# Why Community Assessment?

- 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

# BioCreative Challenge

- ▶ Critical Assessment of Information Extraction systems in Biology
- ▶ Community challenge evaluation a community-wide effort for evaluating text mining and information extraction systems applied to the biological domain
- ▶ Increasing number of groups working in the area of text mining, new systems, publications
- ▶ Need of common standards or shared evaluation criteria to enable comparison
- ▶ Avoid the limitations of using private data sets: One system = one evaluation data set
- ▶ Promote development of systems which scale to real applications
- ▶ Community assessment of scientific progress: Monitor improvements
- ▶ Involve domain experts (end users) and biological database curators and domain experts
- ▶ Extraction of biologically relevant and useful information from

# BioCreative Challenge 1

# BioCreative Challenge 2

# BioCreative Challenge 2.5

# BioCreative Meta-Server (BCMS)

# BioCreative Meta-Server (BCMS)

# BioCreative Meta-Server (BCMS)

# BioCreative Conclusion

- ▶ Repeatability of experimental results
- ▶ Comparability of the experimental results
- ▶ Take into account potential user community: Biologists and Interaction databases
- ▶ Estimate how hard the task actually is and the quality of the training data: Inter-annotator agreement, e.g. kappa score (will be done as well, GB article)
- ▶ Evaluate also sub-aspects: bio-entity, functional term, relationships, sub-categories, organism source, sampling

Part 5: Practical Case Studies

# PLAN2L: Flowchart



Figure: Source:http://zope.bioinfo.cnio.es/plan2l/plan2l.html

# PLAN2L: Protein Normalization



Figure: Source:http://zope.bioinfo.cnio.es/plan2l/plan2l.html

# PLAN2L: Search Bibliome



Figure: Source:http://zope.bioinfo.cnio.es/plan2l/plan2l.html

# PLAN2L: Search PPI



Figure. Source:http://zope.bioinfo.cnio.es/plan2l/plan2l.html

# PLAN2L: Search Gene Regulation Association



Figure: Source:http://zope.bioinfo.cnio.es/plan2l/plan2l.html

# PLAN2L: Search Location Sentences



Figure: Source:http://zope.bioinfo.cnio.es/plan2l/plan2l.html

# PLAN2L: Search Cell Sycle Associations



Figure: Source:http://zope.bioinfo.cnio.es/plan2l/plan2l.html

# PLAN2L: Association between Bio-entities



Figure: Source:http://zope.bioinfo.cnio.es/plan2l/plan2l.html

# References

http://zope.bioinfo.cnio.es/teaching/
http://zope.bioinfo.cnio.es/bionlp_tools/
http://www.tifr.res.in/~ashishvt/biotextmining/