

Biomedical Text Mining

Linking genes to literature: text mining, information extraction, and retrieval applications for biology

Krallinger M. et. al. *Genome Biology* 2008, 9(Suppl 2):S8

Department of Computer and Information Science (DA)
Linköping universitet, Sweden

Disciplines

- **Information Retrieval (IR)**
 - finding the relevant documents
- **Entity Recognition (ER)**
 - identifying the entities
- **Information Extraction (IE)**
 - formalizing the facts
- **Text Mining (TM)**
 - finding nuggets in the literature
- **Integration**
 - combining text and biological data

Literature mining for the biologist:
from information retrieval to biological discovery
Jensen et al., *Nature Reviews Genetics*, 2006

Department of Computer and Information Science (DA)
Linköping universitet, Sweden

Ontologies and Text Mining

Text mining and ontologies in biomedicine: Making sense of raw text
Spasic I et al., *Briefings in bioinformatics*, 6(3):239-251, 2005

Department of Computer and Information Science (DA)
Linköping universitet, Sweden

Ontologies and Text Mining

Department of Computer and Information Science (DA)
Linköping universitet, Sweden

Outline

- Biomedical Text Mining
- Biomedical Ontologies
- TM systems using Ontologies
- Ontology motivated corpus construction

Department of Computer and Information Science (DA)
Linköping universitet, Sweden

The Gene Ontology (GO)

- The goal of the Gene Ontology Consortium (GOC) is to provide three ontologies of defined terms representing gene product properties.

34164 terms (100.0% defined) May 16, 2011 at 13:38 Pacific time

- 20764 biological_process
- 2831 cellular_component
- 9010 molecular_function

relations,

- is_a,
- part_of
- regulates (positively_regulates, negatively_regulates)

Department of Computer and Information Science (DA)
Linköping universitet, Sweden

Biomedical Language

- n Biomedical Language heavily use of domain specific terminology
 - q e.g. *chemoattractant*, *fibroblasts*, *endocytosis*, *exocytosis*
- n Short forms and abbreviations are often used
 - q e.g. *vascular endothelial growth factor (VEGF)*
- n Genes/proteins have often synonyms
 - q e.g. *Thermoactinomyces candidus* vs. *Thermoactinomyces vulgaris*
- n Orthographic variants
 - q e.g. *TNF α* , *TNF-alpha* and *TNF alpha* (without hyphen)

Department of Computer and Information Science (IDA)
Linköping universitet, Sweden

19

Biomedical Language

- n A name
 - q general English term
 - q may refer to a particular gene
 - q may include homologues of this gene in other organisms
 - q may denote an RNA, DNA, or the protein the gene encodes
 - q may be restricted to a specific splice variant

```

graph TD
  NF2 --> NF2_Protein[Neurofibromin 2 [protein]]
  NF2 --> NF2_Gene[Neurofibromin 2 gene [gene]]
  NF2_Protein --> NF2_Protein_1[Neurofibromin 2]
  NF2_Protein --> NF2_Protein_2[Neurofibromin 2 gene [gene]]
  
```

Department of Computer and Information Science (IDA)
Linköping universitet, Sweden

20

ER – Ontologies as terminologies

- n In practice, the distinction between ontological and terminological resources is somewhat arbitrary.
- n Map ontological terms to free text
 - q Features such as the position of a term in hierarchies and the semantic categorization of biomedical concepts can help disambiguate polysemous terms
 - q Coverage
 - q Term variation and management issues

Department of Computer and Information Science (IDA)
Linköping universitet, Sweden

21

IR – Search PubMed with MeSH

- n Search Indexed for MEDLINE citations (90% of the PubMed database) using MeSH terms
 - q MeSH term represents the major focus of the article
 - q Assigned by professional human indexers at the NLM
- n Limit searches to citations with the MeSH term
- n Broaden/Narrow a search with the MeSH hierarchy
 - q A search automatically include all articles which focus not only on the query term, but also focus on narrower terms
- n Use subheadings to build complex and focused search strategies
 - q Combine MeSH Terms

Department of Computer and Information Science (IDA)
Linköping universitet, Sweden

22

IR - GoPubMed

- n It uses GO and MeSH to index search results,
 - q Automatically map all ontology terms in GO and MeSH to the PubMed database.
 - q Categorize the search results
 - q Identify relevant terms
 - q Summarize trends for a topic

GoPubMed: exploring PubMed with the Gene Ontology.
Doms A, Schroeder M. *Nucleic acids research*, 33: W783-W786, 2005.

Department of Computer and Information Science (IDA)
Linköping universitet, Sweden

23

IR - GoPubMed

Department of Computer and Information Science (IDA)
Linköping universitet, Sweden

24

IE/TM - PASTA

- PASTA Templates
 - 12 domain classes, e.g. *protein*, *residue*, *region*.
 - Object-oriented Template
 - template object
 - a specific entity
 - a relation between objects, e.g. *in_protein*, *in_species*
 - Scenario, e.g. a metabolic reaction
 - Slot filler: filling with information extracted from the text
 - Corpus - 1513 Medline abstract relevant to the study of protein structure.

```

@RESIDUE 134: ..
NAME: SERINE
NO: 67
PDB_FUNCTION: "catalytic"
"reaction-binding"
"active-site"
SITE_OBJECT: "residue"
SITE_PROPERTY: "site_location"
RESIDUE: "134"
INTERACTION: "not specified"
..
-IE PROTEIN: ..
RESIDUE: @RESIDUE-134:
PROTEIN: @PROTEIN-24
..
+IE PROTEIN: ..
NAME: "trypsinogen lipase"
RESIDUE: @RESIDUE-134:
RESIDUE: @RESIDUE-134:

```

Fig. 1. PASTA template examples

given the sentence *Ser154, Tyr67 and Asp177 are found at the active site.*

```

residue(e1), name(e1, "Ser154"),
residue(e2), name(e2, "Tyr67"),
residue(e3), name(e3, "Asp177"),
set(e4), set_member(e4, e1),
set_member(e4, e2), set_member(e4, e3),
find(e5), job(e5, e4),
name_of(e5, e1, e4, e5)

```

Protein structures and information extraction from biological texts: the PASTA system.
R. Guizauskas, et al. *Bioinformatics*, 19(1): 135-143, 2003.
31

IE/TM - PASTA

- Discourse Processing
 - extract information from multiple sentences
 - make inferences using a limited predefined *domain ontology*.

S1: *The three-dimensional structure of Endo H has been determined ...*
 S5: *A shallow curved cleft runs across the surface of the molecule from ...*
 S6: *This cleft contains the putative catalytic residue Asp130 ...*

```

S1: protein(e1), name(e1, "Endo H")
S5: cleft(e23), molecule(e25)
    locate_in(e23, e25)
    locate_in(e23, e1)
S6: cleft(e52),
    residue(e61), name(e61, "Asp130")
    contain(e52, e61)
    locate_in(e61, e52)
    locate_in(e61, e1)

```

Ontological Domain Knowledge

i.e. Asp130 is in EndoH.

Department of Computer and Information Science (IDA)
Linköping universitet, Sweden
32

IE/TM - GenIE

Fig. 1. Linear IE architecture.

Fig. 2. GenIE's system architecture.

- GenIE (Genome Information Extraction)
 - A general IE framework.

Ontology-driven discourse analysis for information extraction
Philipp C, et al. *Data & Knowledge Engineering* 55(1): 59-83, 2005.

Department of Computer and Information Science (IDA)
Linköping universitet, Sweden
33

IE/TM - GenIE

Fig. 3. GenIE's system architecture.

- The lexical and knowledge sources
 - A lexicon of gene names for *Saccharomyces cerevisiae*
 - Semantic lexicon - 50 different linguistic variations of the term *yeast*, or *Saccharomyces cerevisiae* in 9000 Medline abstracts
 - POS Tagger - TreeTagger was trained on a manually annotated training corpus (600 SWISS-PROT function slots)
 - Sub-categorization frames for verbs
 - How to acquire the possible sub-categorization frames for all the relevant verbs?
 - binding* events: *bind*, *binds*, *binding*
 - The frames are automatically acquired from a domain-specific corpus (Swiss-Prot function slots)

bind(subj(agent):protein,obj(patient):protein domain)

Department of Computer and Information Science (IDA)
Linköping universitet, Sweden
34

IE/TM - GenIE

- The lexical and knowledge sources
 - An ontology of biochemical events
 - Classification of biochemical events
 - A taxonomy of biochemical events
 - Relations between biochemical events
 - Ontology-driven approach to discourse analysis

Fig. 2. GenIE's system architecture.

Fig. 7. Top level of the ontology O_{Bio} .

Department of Computer and Information Science (IDA)
Linköping universitet, Sweden
35

Outline

- Biomedical Text Mining
 - Biomedical Ontologies
 - TM systems using Ontologies
 - Ontology motivated corpus management

Department of Computer and Information Science (IDA)
Linköping universitet, Sweden
36

Semantic Role Labeling

- Semantic Role Labeling (SRL) is a process that, for each predicate in a sentence, indicates what semantic relations hold among the predicate and other sentence constituents that express the participants in the event.



- It is believed to play a key role in Information Extraction, Question Answering and Summarization.
- Large corpora annotated with semantic roles
 - FrameNet
 - PropBank

Department of Computer and Information Science (DA)
Linköping universitet, Sweden

37

PropBank

- Penn TreeBank → PropBank
 - Add a semantic layer on Penn TreeBank
 - Define a set of semantic roles for each verb
 - VerbNet project maps PropBank verb types to their corresponding Levin classes

hit.01 "strike"

- A0: agent, hitter; A1: thing hit;
- A2: instrument, thing hit by or with

[_{A0} Kristina] hit [_{A1} Scott] [_{A2} with a baseball] yesterday.

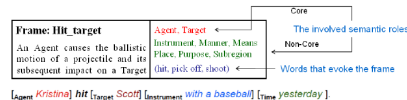


Department of Computer and Information Science (DA)
Linköping universitet, Sweden

38

FrameNet

- Sentences from the British National Corpus
- Method of building FrameNet
 - collects and analyzes the corpus attestations of target words with semantic overlapping.
 - The attestations are divided into semantic groups, and then these small groups are combined into frames.



Department of Computer and Information Science (DA)
Linköping universitet, Sweden

39

FrameNet vs. PropBank

- FrameNet includes semantic analysis of all major parts of speech, not only verb.
- PropBank makes reference to specific tree nodes of TreeBank's syntactic parses of the corpus data.
- In FrameNet, different lexical units within the same Frame will have consistent uses of semantic roles.

Domain-Specific Corpus

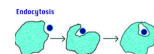
- As with other technologies in natural language processing (NLP), researchers have experienced the difficulties of adapting SRL systems to a new domain, different than the domain used to develop and train the system.
- Biomedical text considerably differs from the text in these corpus, both in the style of the written text and the predicates involved.

Department of Computer and Information Science (DA)
Linköping universitet, Sweden

41

Biomedical Language

- Predicates in biomedical text often prefers nominalizations, gerunds and relational nouns
 - e.g. interaction, association, binding, transcription
- Domain specific predicates are absent from both the FrameNet and PropBank data
 - e.g. endocytosis, exocytosis and translocate
- Predicates have been used in biomedical documents with different semantic senses and require different number of semantic roles compared to FrameNet and PropBank data
 - e.g. block, generate and transform.



Department of Computer and Information Science (DA)
Linköping universitet, Sweden

42

Difficulties of building frame lexicon

- 1. How to discover and define semantic frames together with associated semantic roles within the domain?
- 2. How to collect and group domain-specific predicates to each semantic frame?
- 3. How to select example sentences from publication databases, such as the PubMed/MEDLINE database containing over 20 million articles?

43

Biological procedure ontology of GO

"the directed movement of proteins into, out of or within a cell, or between cells, by means of some agent such as a transporter or pore".

- GO:0008150 : biological_process (353603 gene products)
- GO:0051234 : establishment of localization (40515 gene products)
- GO:0045194 : establishment of protein localization (8102 gene products)
- GO:0015031 : **protein transport** (7366 gene products)
- GO:0071828 : apolipoprotein E recycling (0 gene products)
- GO:0043213 : bacteriocin transport (11 gene products)
- GO:0072321 : chaperone-mediated protein transport (4 gene products)
- GO:0034436 : glycoprotein transport (3 gene products)
- GO:0006886 : intracellular protein transport (4290 gene products)
- GO:0033571 : leuciferin transport (0 gene products)
- GO:0042953 : lipoprotein transport (58 gene products)
- GO:0051224 : negative regulation of protein transport (240 gene products)
- GO:0051222 : positive regulation of protein transport (124 gene products)
- GO:0017038 : protein import (1420 gene products)
- GO:0009306 : protein secretion (1500 gene products)
- GO:0071806 : protein transmembrane transport (504 gene products)
- GO:0072322 : protein transport across periplasmic space (0 gene products)
- GO:0071693 : protein transport within extracellular region (0 gene products)
- GO:0051223 : regulation of protein transport (795 gene products)
- GO:0033572 : transferin transport (12 gene products)
- GO:0051808 : translocation of peptides or proteins into other organism involved in symbiotic interaction (0 gene products)
- GO:0006810 : transport (40099 gene products)
- GO:0015031 : **protein transport** (7366 gene products)
- GO:0071828 : apolipoprotein E recycling (0 gene products)
- GO:0043213 : bacteriocin transport (11 gene products)

177 descendant classes.
581 class names and synonyms

44

Underlying compositional structures in ontological terms

9 direct subclasses of protein transport

The possible predicates *translocation, import, recycling, secretion* and *transport*

The more complex expressions, e.g. "*translocation of peptides or proteins into other organism involved in symbiotic interaction*" (GO:0051808), express participants involved in the event, i.e. the **entity** (*peptides or proteins*), **destination** (*into other organism*) and **condition** (*involved in symbiotic interaction*) of the event.

45

Aspects of the method

- 1. The structure and semantics of domain knowledge in ontologies constrain the frame semantics analysis, i.e. decide the coverage of semantic frames and the relations between them;
- 2. Ontological terms can comprehensively describe the characteristics of events/scenarios in the domain, so domain-specific semantic roles can be determined based on terms;
- 3. Ontological terms provide a list of domain specific predicates, so the semantic sense of the predicates in the domain are determined;
- 4. The collection and selection of example sentences can be based on knowledge-based search engine for biomedical text.

46

"Protein Transport" Frame

TE	definition
Transport-Entity (TE)	Person or protein complex which is undergoing the motion event into, out of or within a cell, or between cells, or within a multicellular organism.
Transport-Origin (TO)	The organelle, cell, tissue or gland from which the Transport-Entity moves to a different location.
Transport-Destination (TD)	The organelle, cell, tissue or gland to which the Transport-Entity moves from a different location.
Transport-Condition (TC)	The event, substance, organelle or chemical environment which positively or negatively directly influences or is influenced by, the motion event. The substance/organelle does not necessarily move with the Transport-Entity.
Transport-Location (TL)	The organelle, cell, tissue or gland where the motion event takes place when the origin and the destination are the same or when origin or destination is not specified.
Transport-Path (TP)	The substance or organelle which helps the entity to move from the Transport-Origin to the Transport-Destination, sometimes by connecting the two locations, without itself undergoing translocation.
Transport-Transporter (TT)	The substance, organelle or cell which is the motion event, that moves along with the Transport-Entity, taking it from the Transport-Origin to the Transport-Destination.
Transport-Direction	The direction in which the motion event is taking place with re-

47

The Frame Lexicon

- 1. First work on ontology driven corpus management
- 2. Released the corpus covering "*protein transport*" event, <http://www.ida.liu.se/~hetan/bio-onto-frame-corpus/>
- 3. We aim to extend the corpus to cover other biological events.
 - 1. GO ontologies, other ontologies (e.g. pathway ontologies)
- 4. The identification of frames and the relations between frames are needed to be investigated.
- 5. We will study the definition of Semantic Type (ST) in the domain corpus and their mappings to classes in top domain ontologies
 - 1. ST: "Sentient" defined for the semantic role "Cognizer" in the frame "Cogitation".

48

Thanks



Department of Computer and Information Science (IDA)
Linköping universitet, Sweden