

Introduction to Text Mining

—Tutorial at EDBT'06—

René Witte

Faculty of Informatics
Institute for Program Structures and Data Organization (IPD)
Universität Karlsruhe, Germany
<http://rene-witte.net>

27.03.2006

Lack of Information?

The screenshot shows a Konqueror browser window titled "text mining - Google Search - Konqueror". The address bar contains the URL "http://www.google.com/search?hl=en&lr=&q=text+mining&btnG=Search". The search results page displays the Google logo, a search bar with "text mining" entered, and a "Search" button. Below the search bar, it indicates "Results 1 - 10 of about 55,900,000 for text mining. (0.11 seconds)".

The search results are categorized into "Web" and "Sponsored Links".

Web Results:

- Powerful Text Mining?** (Sponsored Link)
www.semantic-knowledge.com Use a Desktop Search Engine, with Semantics and Text Analysis
- Marti Hearst: What Is Text Mining?**
What is **text mining**? What are its potential applications and limitations? **Text Mining** is the discovery by computer of new, previously unknown information, ...
www.sims.berkeley.edu/~hearst/text-mining.html - 9k - [Cached](#) - [Similar pages](#)
- text-mining.org**
Provide a web home for people interested in **text mining** related technologies, with a mailing list and a resources section.
www.text-mining.org/ - 30k - 16 Mar 2006 - [Cached](#) - [Similar pages](#)
- Data Mining, Text Mining and Web Mining Software**
Data, **text**, and web **mining** software. PolyAnalyst includes in-place **mining**, strong Microsoft integration.
www.megaputer.com/ - 21k - [Cached](#) - [Similar pages](#)

Sponsored Links:

- Google Search Appliance**
Make Your Intranet More Useful Search with the Power of Google
www.google.de/appliance
- Text Mining Solutions**
For Unstructured Multilingual Data. Foreign Language **Text** Analysis.
www.basistech.com
- Text Data Mining**
Fact-based data discovery and **text mining**: Insightful InFact
www.insightful.com
- Data Matching Solutions**
Flexible solution to your DQ issues Deployable in weeks not months
www.datactics.com

Lack of Information?

The screenshot shows a web browser window titled "text mining tools German - Google Search - Konqueror". The address bar contains the search URL: `http://www.google.com/search?hl=en&lr=&q=text+mining+tools+German&btnG=Search`. The search results are displayed under the heading "Web" and show "Results 1 - 10 of about 20,400,000 for text mining tools German. (0.31 seconds)".

The first search result is for "txtkit - Visual Text Mining Tool | Homepage". The snippet reads: "txtkit is an Open Source visual text mining tool for exploring large amounts of ... Our servers offer four sources at the moment: Hans Ulrich Reck (german ... www.txtkit.sw.ofcd.com/ - 13k - [Cached](#) - [Similar pages](#)".

The second search result is for "txtkit - Visual Text Mining Tool | Download". The snippet reads: "txtkit - Visual Text Mining Tool. Shell Interface · Visual Bot · Architecture; Download ... HUReck_Collection.de.txtkit (Hans Ulrich Reck / german) ... www.txtkit.sw.ofcd.com/obj/download - 14k - [Cached](#) - [Similar pages](#)".

The third search result is for "SRA International - Text Mining Solutions". The snippet reads: "We provide text mining tool recommendations and recommendations for other ... We then apply the text mining tools and techniques and gather feedback on how ... www.sra.com/services/index.asp?id=172 - 19k - [Cached](#) - [Similar pages](#)".

On the right side of the page, there are "Sponsored Links". The first is "Text Mining Software" with the text "Content Analysis & Text Mining Tool Download a free trial version now ! www.provallsresearch.com". The second is "Text Mining Tools" with the text "Find Solutions for your Business Free Reports, Info & Registration! www.KnowledgeStorm.com".

The status bar at the bottom of the browser window indicates "Page loaded."

Tutorial Overview

Today's Tutorial contains...

Introduction: Motivation, definitions, applications

Foundations: Theoretical background in Computational Linguistics

Technology: Technological foundations for building Text Mining systems

Applications: In-depth description of two application areas (summarization, biology) and overview on two others (question-answering, opinion mining)

Conclusions: the end.

Each part contains some references for further study.

Part I

Introduction

3 Introduction
• Motivation

4 Definitions
• Text Mining

5 Applications
• Domains

Information Overload

Too much (textual) information

- We now have electronic books, documents, web pages, emails, blogs, news, chats, memos, research papers, ...
- ... all of it immediately accessible, thanks to databases and Information Retrieval (IR)
- An estimated 80–85% of all data stored in databases are natural language texts
- But humans did not scale so well...

This results in the common perception of **Information Overload**.



Example: The BioTech Industry

Access to information is a serious problem

- 80% of biological knowledge is **only** in reasearch papers
- finding the information you need is prohibitively expensive

Humans do not scale well

- if you read 60 research papers/week...
- ...and 10% of those are interesting...
- ...a scientist manages 6/week, or 300/year

This is not good enough

- MedLine adds more than 10 000 abstracts each *month*!
- Chemical Abstracts Registry (CAS) registers 4000 entities **each day**, 2.5 million in 2004 alone

[cf. Talk by Robin McEntire of GlaxoSmithKline at KBB'05]



Example: The BioTech Industry

Access to information is a serious problem

- 80% of biological knowledge is **only** in reasearch papers
- finding the information you need is prohibitively expensive

Humans do not scale well

- if you read 60 research papers/week...
- ...and 10% of those are interesting...
- ...a scientist manages 6/week, or 300/year

This is not good enough

- MedLine adds more than 10 000 abstracts each *month*!
- Chemical Abstracts Registry (CAS) registers 4000 entities **each day**, 2.5 million in 2004 alone

[cf. Talk by Robin McEntire of GlaxoSmithKline at KBB'05]



Example: The BioTech Industry

Access to information is a serious problem

- 80% of biological knowledge is **only** in research papers
- finding the information you need is prohibitively expensive

Humans do not scale well

- if you read 60 research papers/week...
- ...and 10% of those are interesting...
- ...a scientist manages 6/week, or 300/year

This is not good enough

- MedLine adds more than 10 000 abstracts each *month*!
- Chemical Abstracts Registry (CAS) registers 4000 entities **each day**, 2.5 million in 2004 alone

[cf. Talk by Robin McEntire of GlaxoSmithKline at KBB'05]

Definitions

One usually distinguishes

- Information Retrieval
- Information Extraction
- Text Mining

Text Mining (Def. *Wikipedia*)

Text mining, also known as intelligent text analysis, text data mining or knowledge-discovery in text (KDT), refers generally to the process of extracting interesting and non-trivial information and knowledge from unstructured text. Text mining is a young interdisciplinary field which draws on information retrieval, data mining, machine learning, statistics and computational linguistics. As most information (over 80%) is stored as text, text mining is believed to have a high commercial potential value.

Definitions

One usually distinguishes

- Information Retrieval
- Information Extraction
- Text Mining

Text Mining (Def. *Wikipedia*)

Text mining, also known as intelligent text analysis, text data mining or knowledge-discovery in text (KDT), refers generally to the process of extracting interesting and non-trivial information and knowledge from unstructured text. Text mining is a young interdisciplinary field which draws on information retrieval, data mining, machine learning, statistics and computational linguistics. As most information (over 80%) is stored as text, text mining is believed to have a high commercial potential value.

Definitions

One usually distinguishes

- Information Retrieval
- Information Extraction
- Text Mining

Text Mining (Def. *Wikipedia*)

Text mining, also known as intelligent text analysis, text data mining or knowledge-discovery in text (KDT), refers generally to the process of extracting interesting and non-trivial information and knowledge from unstructured text. Text mining is a young interdisciplinary field which draws on information retrieval, data mining, machine learning, statistics and computational linguistics. As most information (over 80%) is stored as text, text mining is believed to have a high commercial potential value.

Definitions

One usually distinguishes

- Information Retrieval
- Information Extraction
- Text Mining

Text Mining (Def. *Wikipedia*)

Text mining, also known as intelligent text analysis, text data mining or knowledge-discovery in text (KDT), refers generally to the process of extracting interesting and non-trivial information and knowledge from unstructured text. Text mining is a young interdisciplinary field which draws on information retrieval, data mining, machine learning, statistics and computational linguistics. As most information (over 80%) is stored as text, text mining is believed to have a high commercial potential value.

What to mine?

Emails, Instant Messages, Blogs, ...

Look for:

- Entities (Persons, Companies, Organizations, ...)
- Events (Inventions, Offers, Attacks, ...)

Biggest existing system: **ECHELON** (UKUSA)

What to mine? (II)

News: Newspaper articles, Newswires, ...

Similar to last, but additionally:

- collections of articles (e.g., from different agencies, describing the same event)
- contrastive summaries (e.g., event described by U.S. newspaper vs. Arabic newspaper)
- also needs temporal analysis
- main problems: cross-language and cross-document analysis

Many publicly accessible systems, e.g. *Google News* or *Newsblaster*.

What to mine? (III)

(Scientific) Books, Papers, ...

- detect new trends in research
- automatic curation of research results in Bioinformatics

need to deal with highly specific language

Software Requirement Specifications, Documentation, ...

- extract requirements from software specification
- detect conflicts between source code and its documentation

Web Mining

extract and analyse information from web sites

- mine companies' web pages (detect new products & trends)
- mine Intranets (gather knowledge, find "illegal" content, ...)

problems: not simply plain text, also hyperlinks and hidden information ("deep web")

What to mine? (III)

(Scientific) Books, Papers, ...

- detect new trends in research
- automatic curation of research results in Bioinformatics

need to deal with highly specific language

Software Requirement Specifications, Documentation, ...

- extract requirements from software specification
- detect conflicts between source code and its documentation

Web Mining

extract and analyse information from web sites

- mine companies' web pages (detect new products & trends)
- mine Intranets (gather knowledge, find "illegal" content, ...)

problems: not simply plain text, also hyperlinks and hidden information ("deep web")

What to mine? (III)

(Scientific) Books, Papers, ...

- detect new trends in research
- automatic curation of research results in Bioinformatics

need to deal with highly specific language

Software Requirement Specifications, Documentation, ...

- extract requirements from software specification
- detect conflicts between source code and its documentation

Web Mining

extract and analyse information from web sites

- mine companies' web pages (detect new products & trends)
- mine Intranets (gather knowledge, find "illegal" content, ...)

problems: not simply plain text, also hyperlinks and hidden information ("deep web")

Typical Text Mining Tasks

Classification and Clustering

- Email Spam-Detection, Classification (Orders, Offers, ...)
- Clustering of large document sets (vivisimo.com)
- Creation of topic maps (www.leximancer.com)

Web Mining

- Trend Mining, Opinion Mining, Novelty Detection
- Ontology Creation, Entity Tracking, Information Extraction

“Classical” NLP Tasks

- Machine Translation (MT)
- Automatic Summarization
- Question-Answering (QA)

Typical Text Mining Tasks

Classification and Clustering

- Email Spam-Detection, Classification (Orders, Offers, ...)
- Clustering of large document sets (vivisimo.com)
- Creation of topic maps (www.leximancer.com)

Web Mining

- Trend Mining, Opinion Mining, Novelty Detection
- Ontology Creation, Entity Tracking, Information Extraction

“Classical” NLP Tasks

- Machine Translation (MT)
- Automatic Summarization
- Question-Answering (QA)

Typical Text Mining Tasks

Classification and Clustering

- Email Spam-Detection, Classification (Orders, Offers, ...)
- Clustering of large document sets (vivisimo.com)
- Creation of topic maps (www.leximancer.com)

Web Mining

- Trend Mining, Opinion Mining, Novelty Detection
- Ontology Creation, Entity Tracking, Information Extraction

“Classical” NLP Tasks

- Machine Translation (MT)
- Automatic Summarization
- Question-Answering (QA)

Information Overload, Part II

Can't you just summarize this for me?

Create “intelligent assistants” that retrieve, process, and condense information for you.

We already have: Information Retrieval

We need: Technologies to process the retrieved information

One example is **Automatic Summarization** to condense a single document or a set of documents.

For example...

Mrs. Coolidge: *What did the preacher discuss in his sermon?*

President Coolidge: *Sin.*

Mrs. Coolidge: *What did he say?*

President Coolidge: *He said he was against it.*

Information Overload, Part II

Can't you just summarize this for me?

Create “intelligent assistants” that retrieve, process, and condense information for you.

We already have: Information Retrieval

We need: Technologies to process the retrieved information

One example is **Automatic Summarization** to condense a single document or a set of documents.

For example. . .

Mrs. Coolidge: What did the preacher discuss in his sermon?

President Coolidge: Sin.

Mrs. Coolidge: What did he say?

President Coolidge: He said he was against it.

Automatic Summarization

Example source (newspaper article)

HOUSTON – The Hubble Space Telescope got smarter and better able to point at distant astronomical targets on Thursday as spacewalking astronauts replaced two major pieces of the observatory's gear. On the second spacewalk of the shuttle Discovery's Hubble repair mission, the astronauts, C. Michael Foale and Claude Nicollier, swapped out the observatory's central computer and one of its fine guidance sensors, a precision pointing device. The spacewalkers ventured into Discovery's cargo bay, where Hubble towers almost four stories above, at 2:06 p.m. EST, about 45 minutes earlier than scheduled, to get a jump on their busy day of replacing some of the telescope's most important components. ...

Summary (10 words)

Space News: [the shuttle Discovery's Hubble repair mission, the observatory's central computer]

Automatic Summarization

Example source (newspaper article)

HOUSTON – The Hubble Space Telescope got smarter and better able to point at distant astronomical targets on Thursday as spacewalking astronauts replaced two major pieces of the observatory's gear. On the second spacewalk of the shuttle Discovery's Hubble repair mission, the astronauts, C. Michael Foale and Claude Nicollier, swapped out the observatory's central computer and one of its fine guidance sensors, a precision pointing device. The spacewalkers ventured into Discovery's cargo bay, where Hubble towers almost four stories above, at 2:06 p.m. EST, about 45 minutes earlier than scheduled, to get a jump on their busy day of replacing some of the telescope's most important components. ...

Summary (10 words)

Space News: [the shuttle Discovery's Hubble repair mission, the observatory's central computer]

Dealing with Text in Natural Languages

Problem

How can I automatically create a summary from a text written in natural language?

Solution: Natural Language Processing (NLP)

Current trends in NLP:

- deal with “real-world” texts, not just limited examples
- requires robust, fault-tolerant algorithms (e.g., partial parsing)
- shift from rule-based approaches to statistical methods and machine learning
- focus on “knowledge-poor” techniques, as even shallow semantics is quite tough to obtain

Dealing with Text in Natural Languages

Problem

How can I automatically create a summary from a text written in natural language?

Solution: Natural Language Processing (NLP)

Current trends in NLP:

- deal with “real-world” texts, not just limited examples
- requires robust, fault-tolerant algorithms (e.g., partial parsing)
- shift from rule-based approaches to statistical methods and machine learning
- focus on “knowledge-poor” techniques, as even shallow semantics is quite tough to obtain

Part II

Foundations

6 Introduction

7 Computational Linguistics

- Introduction
- Ambiguity
- Rule-based vs. Statistical NLP
- Preprocessing and Tokenisation
- Sentence Splitting
- Morphology
- Part-of-Speech (POS) Tagging
- Chunking and Parsing
- Semantics
- Pragmatics: Co-reference resolution

8 Performance Evaluation

- Evaluation Measures
- Accuracy and Error
- Precision and Recall
- F-Measure and Inter-Annotator Agreement
- More complex evaluations

9 Literature

Take your PP-Attachement out of my Garden Path!

Understanding Computational Linguists

Text Mining is concerned with processing documents written in natural language:

- this is the domain of *Computational Linguistics (CL)* and *Natural Language Processing (NLP)*
- practical application, with more of an engineering perspective, also called *Language Technology (LT)*
- Text Mining (TM) is concerned with concrete practical applications (compare: “Information Systems” and “Databases”)

Hence, we need to review some concepts, terminology, and foundations from these areas.



Computational Linguistics 101

Classical Categorization

To deal with the complexity of natural language, it is typically regarded on several levels (cf. Jurafsky & Martin):

Phonology the study of linguistic sounds

Morphology the study of meaningful components of words

Syntax the study of structural relationships between words

Semantics the study of meaning

Pragmatics the study of how language is used to accomplish goals

Discourse the study of larger linguistic units

Importance for Text Mining

- *Phonology* only concerns spoken language
- *Discourse*, *Pragmatics*, and even *Semantics* is still rarely used

Ambiguity

Ambiguity appears on every analysis level

The classical examples:

- *He saw the man with the telescope.*
- *Time flies like an arrow. Fruit flies like a banana.*

And those are simple. . .

This does not get better with real-world sentences:

- *The board approved [its acquisition] [by Royal Trustco. Ltd.] [of Toronto] [for \$27 a share] [at its monthly meeting].*

(cf. Manning & Schütze)

Ambiguity

Ambiguity appears on every analysis level

The classical examples:

- *He saw the man with the telescope.*
- *Time flies like an arrow. Fruit flies like a banana.*

And those are simple. . .

This does not get better with real-world sentences:

- *The board approved [its acquisition] [by Royal Trustco. Ltd.] [of Toronto] [for \$27 a share] [at its monthly meeting].*

(cf. Manning & Schütze)

Tokenization (II)

What is a word?

Unfortunately, even tokenization can be difficult:

- Is *John's sick* one token or two?
If one → problems in parsing (where's the verb?)
If two → what do we do with *John's house*?
- What to do with hyphens?
E.g., *database* vs. *data-base* vs. *data base*
- what to do with "C++", "A/C", ":-)", "..."?

Even worse...

- Some languages don't use whitespace (e.g., Chinese)
→ need to run a *word segmentation* first
- Heavy compounding e.g. in German, decomposition necessary
"Rinderbraten" (roast beef) → *Rind|erbraten?*
Rind|erb|raten? *Rinder|braten?*

Tokenization (II)

What is a word?

Unfortunately, even tokenization can be difficult:

- Is *John's sick* one token or two?
If one → problems in parsing (where's the verb?)
If two → what do we do with *John's house*?
- What to do with hyphens?
E.g., *database* vs. *data-base* vs. *data base*
- what to do with "C++", "A/C", ":-)", "..."?

Even worse. . .

- Some languages don't use whitespace (e.g., Chinese)
→ need to run a *word segmentation* first
- Heavy compounding e.g. in German, decomposition necessary
"Rinderbraten" (roast beef) → *Rind|erbraten?*
Rind|erb|raten? *Rinder|braten?*

Tokenization (III)

The good, the bad, and the ...

Tokenization can become even more difficult in specific domains.

Software Documents

Documents include lots of source code snippets:

- `package java.util.*`
- *The range-view operation, `subList(int fromIndex, int toIndex)`, returns a List view of the portion of this list whose indices range from `fromIndex`, inclusive, to `toIndex`, exclusive.*

Need to deal with URLs, methods, class names, etc.

Tokenization (III)

The good, the bad, and the ...

Tokenization can become even more difficult in specific domains.

Software Documents

Documents include lots of source code snippets:

- `package java.util.*`
- *The range-view operation, `subList(int fromIndex, int toIndex)`, returns a List view of the portion of this list whose indices range from `fromIndex`, inclusive, to `toIndex`, exclusive.*

Need to deal with URLs, methods, class names, etc.

Tokenization (IV)

Biological Documents

Highly complex expressions, chemical formulas, etc.:

- *1,4- β -xylanase II from Trichoderma reesei*
- *When N-formyl-L-methionyl-L-leucyl-L-phenylalanine (fMLP) was injected. . .*
- *Technetium-99m-CDO-MeB [Bis[1,2-cyclohexanedione-dioximato(1-)-O]-[1,2-cyclohexanedione dioximato(2-)-O]methyl-borato(2-)-N,N',N'',N''',N''''',N'''''')-chlorotechnetium) belongs to a family of compounds. . .*

Sentence Splitting

Mark Sentence Boundaries

Detects sentence units. Easy case:

- often, sentences end with “.”, “!”, or “?”

Hard (or annoying) cases:

- difficult when a “.” do not indicate an EOS:
“MR. X”, “3.14”, “Y Corp.”, ...
- we can detect common abbreviations (“U.S.”), but what if a sentence ends with one?
“...announced today by the U.S. The...”
- Sentences can be *nested* (e.g., within quotes)

Correct sentence boundary is important

for many downstream analysis tasks:

- POS-Taggers maximize probabilities of tags within a sentence
- Summarization systems rely on correct detection of sentence

Sentence Splitting

Mark Sentence Boundaries

Detects sentence units. Easy case:

- often, sentences end with “.”, “!”, or “?”

Hard (or annoying) cases:

- difficult when a “.” do not indicate an EOS:
“MR. X”, “3.14”, “Y Corp.”, ...
- we can detect common abbreviations (“U.S.”), but what if a sentence ends with one?
“...announced today by the U.S. The...”
- Sentences can be *nested* (e.g., within quotes)

Correct sentence boundary is important

for many downstream analysis tasks:

- POS-Taggers maximize probabilities of tags within a sentence
- Summarization systems rely on correct detection of sentence

Morphological Analysis

Morphological Variants

Words are changed through a morphological process called *inflection*:

- typically indicates changes in case, gender, number, tense, etc.
- example *car* → *cars*, *give* → *gives*, *gave*, *given*

Goal: “normalize” words

Stemming and Lemmatization

Two main approaches to normalization:

Stemming reduce words to a *base form*

Lemmatization reduce words to their *lemma*

Main difference: stemming just finds **any** base form, which doesn't even need to be a word in the language! Lemmatization find the actual *root* of a word, but requires morphological analysis.

Morphological Analysis

Morphological Variants

Words are changed through a morphological process called *inflection*:

- typically indicates changes in case, gender, number, tense, etc.
- example *car* → *cars*, *give* → *gives*, *gave*, *given*

Goal: “normalize” words

Stemming and Lemmatization

Two main approaches to normalization:

Stemming reduce words to a *base form*

Lemmatization reduce words to their *lemma*

Main difference: stemming just finds **any** base form, which doesn't even need to be a word in the language! Lemmatization find the actual *root* of a word, but requires morphological analysis.

Stemming vs. Lemmatization

Stemming

Commonly used in Information Retrieval:

- Can be achieved with rule-based algorithms, usually based on suffix-stripping
- Standard algorithm for English: the *Porter* stemmer
- Advantages: simple & fast
- Disadvantages:
 - Rules are language-dependent
 - Can create words that do not exist in the language, e.g., *computers* → *comput*
 - Often reduces different words to the same stem, e.g.,
army, arm → *arm*
stocks, stockings → *stock*
- Stemming for German: German stemmer in the full-text search engine *Lucene*, *Snowball* stemmer with German rule file

Stemming vs. Lemmatization, Part II

Lemmatization

Lemmatization is the process of deriving the base form, or *lemma*, of a word from one of its inflected forms. This requires a morphological analysis, which in turn typically requires a *lexicon*.

- Advantages:
 - identifies the *lemma* (root form), which is an actual word
 - less errors than in stemming
- Disadvantages:
 - more complex than stemming, slower
 - requires additional language-dependent resources
- While stemming is good enough for Information Retrieval, Text Mining often requires lemmatization
 - Semantics is more important (we need to distinguish an *army* and an *arm*!)
 - Errors in low-level components can multiply when running downstream

Lemmatization Example

Lemmatization in German

Lemmatization for a morphologically complex language like German is complicated

- Cannot be solved through a rule-based algorithm

Kinder → *Kind* *Vorlesungen* → *Vorlesung* *Länder* → *Land*
Leiter → **Leit* *Leben* → **Leb* *Affären* → **Affare*

- An accurate lemmatization for German requires a lexicon
 - For each word, all inflected forms or morphological rules

The *Durm* German Lemmatizer

A self-learning context-aware lemmatization system for German that can create (and correct) a lexicon by processing German documents:

```
Menschen  Sg  Masc  Akk  Mensch  6  4/11/2005  15:8:16
           4/11/2005  15:10:11  116  unlocked
```

Part-of-Speech (POS) Tagging

Where are we now?

So far, we splitted texts into *tokens* and *sentences* and performed some *normalization*.

- Still a long way to go to an *understanding* of natural language. . .

Typical approach in NLP: deal with the complexity of language by applying intermediate processing steps to acquire more and more structure. Next stop: *POS-Tagging*.

POS-Tagging

A statistical POS Tagger scans tokens and assigns **POS Tags**.
A black cat plays. . . → *A/DT black/JJ cat/NN plays/VB. . .*

- relies on different word order probabilities
- needs a manually tagged corpus for machine learning

Note: *this is not parsing!*

Part-of-Speech (POS) Tagging

Where are we now?

So far, we splitted texts into *tokens* and *sentences* and performed some *normalization*.

- Still a long way to go to an *understanding* of natural language. . .

Typical approach in NLP: deal with the complexity of language by applying intermediate processing steps to acquire more and more structure. Next stop: *POS-Tagging*.

POS-Tagging

A statistical POS Tagger scans tokens and assigns **POS Tags**.
A black cat plays. . . → *A/DT black/JJ cat/NN plays/VB. . .*

- relies on different word order probabilities
- needs a manually tagged corpus for machine learning

Note: *this is not parsing!*

Part-of-Speech (POS) Tagging (II)

Tagsets

A **tagset** defines the tags to assign to words. Main POS classes are:

Noun refers to entities like people, places, things or ideas

Adjective describes the properties of nouns or pronouns

Verb describes actions, activities and states

Adverb describes a verb, an adjective or another adverb

Pronoun word that can take the place of a noun

Determiner describes the particular reference of a noun

Preposition expresses spatial or time relationships

Note: real tagsets have from 45 (Penn Treebank) to 146 tags (C7).

POS Tagging Algorithms

Fundamentals

POS-Tagging generally requires:

Training phase where a **manually annotated** corpus is processed by a machine learning algorithm; and a

Tagging algorithm that processes texts using learned parameters.

Performance is generally good (around 96%) when staying in the same domain.

Algorithms used in POS-Tagging

There is a multitude of approaches, commonly used are:

- Decision Trees
- Hidden Markov Models (HMMs)
- Support Vector Machines (SVM)
- Transformation-based Taggers (e.g., the **Brill** tagger)

POS Tagging Algorithms

Fundamentals

POS-Tagging generally requires:

Training phase where a **manually annotated** corpus is processed by a machine learning algorithm; and a

Tagging algorithm that processes texts using learned parameters.

Performance is generally good (around 96%) when staying in the same domain.

Algorithms used in POS-Tagging

There is a multitude of approaches, commonly used are:

- Decision Trees
- Hidden Markov Models (HMMs)
- Support Vector Machines (SVM)
- Transformation-based Taggers (e.g., the **Brill** tagger)

Syntax: Chunking and Parsing

Finding Syntactic Structures

We can now start a **syntactic analysis** of a sentence using:

Parsing producing a *parse tree* for a sentence using a parser, a grammar, and a lexicon

Chunking finding syntactic constituents like *Noun Phrases (NPs)* or *Verb Groups (VGs)* within a sentence

Chunking vs. Parsing

Producing a *full parse tree* often fails due to grammatical inaccuracies, novel words, bad tokenization, wrong sentence splits, errors in POS tagging, ...

Hence, *chunking* and *partial parsing* are more commonly used.

Syntax: Chunking and Parsing

Finding Syntactic Structures

We can now start a **syntactic analysis** of a sentence using:

Parsing producing a *parse tree* for a sentence using a parser, a grammar, and a lexicon

Chunking finding syntactic constituents like *Noun Phrases (NPs)* or *Verb Groups (VGs)* within a sentence

Chunking vs. Parsing

Producing a *full parse tree* often fails due to grammatical inaccuracies, novel words, bad tokenization, wrong sentence splits, errors in POS tagging, ...

Hence, *chunking* and *partial parsing* are more commonly used.

Noun Phrase Chunking

NP Chunker

Recognition of noun phrases through context-free grammar with Earley-type chart parser

Grammar Excerpt

```
(NP      (DET MOD HEAD))  
(MOD     (MOD-ingredients)  
         (MOD-ingredients MOD))  
(      )  
(HEAD   (NN) ...)
```

Example

Noun Phrase Chunking

NP Chunker

Recognition of noun phrases through context-free grammar with Earley-type chart parser

Example

Grammar Excerpt

```
(NP      (DET MOD HEAD))
(MOD     (MOD-ingredients)
         (MOD-ingredients MOD)
         ())
(HEAD    (NN) ...)
```

Noun Phrase Chunking

NP Chunker

Recognition of noun phrases through context-free grammar with Earley-type chart parser

Grammar Excerpt

```
(NP      (DET MOD HEAD))
(MOD     (MOD-ingredients)
         (MOD-ingredients MOD))
         ( ))
(HEAD    (NN) ...)
```

Example

"I couldn't believe what I saw," said McNeill, who also discovered bomb-making instructions and detailed maps of U.S. landmarks in the cave. "On top of all the destruction these people had already unleashed, plans were underway to harass the American people with a merciless assault of offers for everything from discounts on home DSL lines to pre-approved, low-interest credit cards."

For all the evidence collected by the CIA, the "smoking gun" in the investigation may turn out to be an alleged Osama bin Laden motivational videotape, currently in the possession of CNN. The controversial tape, which has never aired on the cable network, is rumored to feature bin Laden urging his followers to think positive and believe in the quality of the product they are pitching, closing on the grim slogan "Smile And Dial."

type	Set	Start	End	Features
P	Default	3582	3596	{DET="", MOD="", HEAD="Guantanamo Bay"}
P	Default	776	791	{DET="the", MOD="dinner", HEAD="hour"}
P	Default	2259	2262	{DET="", MOD="", HEAD="out"}
P	Default	1806	1807	{DET="", MOD="", HEAD="I"}
P	Default	3849	3852	{DET="", MOD="", HEAD="one"}
P	Default	987	996	{DET="The", MOD="", HEAD="video"}
P	Default	1487	1494	{DET="", MOD="", HEAD="McNeill"}
P	Default	2280	2318	{DET="", MOD="Osama bin Laden motivational", HEAD="videotape"}
P	Default	894	910	{DET="", MOD="money", HEAD="laundering"}

Chunking vs. Parsing, Round 2

What can we do with chunks?

(NP) chunks are very useful in finding **named entities** (NEs), e.g., *Persons, Companies, Locations, Patents, Organisms, ...*

But additional methods are needed for finding **relations**:

- *Who* invented *X*?
- *What* company created product *Y* that is doomed to fail?
- *Which* organism is this protein coming from?

Parse trees can help in determining these relationships

Parsing Challenges

Parsing is hard due to many kinds of ambiguities:

PP-Attachement which NP takes the PP? Compare:

He ate spaghetti with a fork.

He ate spaghetti with tomato sauce.

NP Bracketing *plastic cat food can cover*

Chunking vs. Parsing, Round 2

What can we do with chunks?

(NP) chunks are very useful in finding **named entities** (NEs), e.g., *Persons, Companies, Locations, Patents, Organisms, ...*

But additional methods are needed for finding **relations**:

- *Who* invented *X*?
- *What* company created product *Y* that is doomed to fail?
- *Which* organism is this protein coming from?

Parse trees can help in determining these relationships

Parsing Challenges

Parsing is hard due to many kinds of ambiguities:

PP-Attachement which NP takes the PP? Compare:

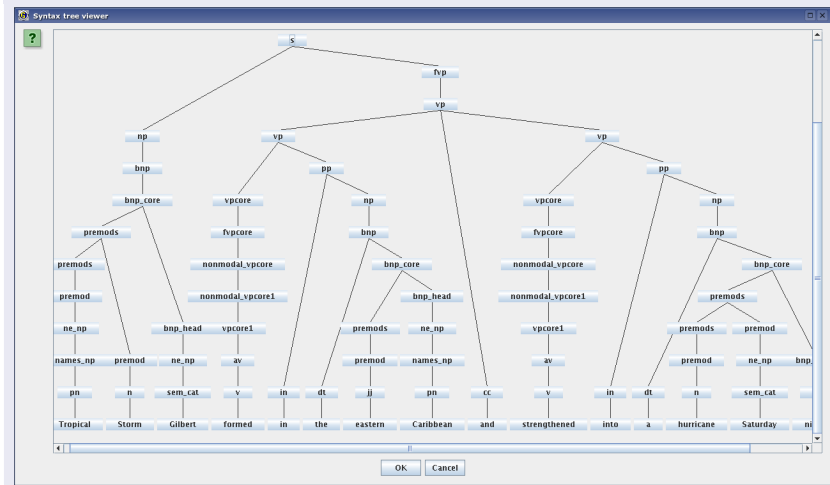
He ate spaghetti with a fork.

He ate spaghetti with tomato sauce.

NP Bracketing *plastic cat food can cover*

Parsing: Example

Example of a (partial) parser output using SUPPLE



Semantics

Moving on...

Now that we have syntactic information, we can start to address the *meaning* of words.

WordNets

A **WordNet** is a semantic network encoding the words of a single (or multiple) language(s) using:

Synsets encoding the *meanings* for each word (e.g., *bank*)

Relations synonymy, antonymy, hypernymy, hyponymy, holonymy, meronymy, homonymy, troponymy, ...

The English WordNet currently encodes 147249 words (v2.1) and is freely available.

Example

Use WordNet to find out whether **tea** is something we can drink.

Semantics

Moving on...

Now that we have syntactic information, we can start to address the *meaning* of words.

WordNets

A **WordNet** is a semantic network encoding the words of a single (or multiple) language(s) using:

Synsets encoding the *meanings* for each word (e.g., *bank*)

Relations synonymy, antonymy, hypernymy, hyponymy, holonymy, meronymy, homonymy, troponymy, ...

The English WordNet currently encodes 147249 words (v2.1) and is freely available.

Example

Use WordNet to find out whether **tea** is something we can drink.

Semantics

Moving on...

Now that we have syntactic information, we can start to address the *meaning* of words.

WordNets

A **WordNet** is a semantic network encoding the words of a single (or multiple) language(s) using:

Synsets encoding the *meanings* for each word (e.g., *bank*)

Relations synonymy, antonymy, hypernymy, hyponymy, holonymy, meronymy, homonymy, troponymy, ...

The English WordNet currently encodes 147249 words (v2.1) and is freely available.

Example

Use WordNet to find out whether **tea** is something we can drink.

WordNet Example

Lookup for "tea"

WordNet Search - 2.1 - Konqueror

Location Edit View Go Bookmarks Tools Settings Window Help

Location: <http://wordnet.princeton.edu/perl/webwn>

WordNet Search - 2.1

[Return to WordNet Home](#)

[Glossary](#) - [Help](#)

SEARCH DISPLAY OPTIONS: (Select option to change)

Enter a word to search for:

KEY: "S:" = Show Synset (semantic) relations, "W:" = Show Word (lexical) relations

Noun

- **S: (n) tea** (a beverage made by steeping tea leaves in water) *"iced tea is a cooling drink"*
- **S: (n) tea**, [afternoon tea](#), [teatime](#) (a light midafternoon meal of tea and sandwiches or cakes) *"an Englishman would interrupt a war to have his afternoon tea"*
- **S: (n) tea**, [tea leaf](#) (dried leaves of the tea shrub; used to make tea) *"the store shelves held many different kinds of tea"; "they threw the tea into Boston harbor"*
- **S: (n) tea** (a reception or party at which tea is served) *"we met at the Dean's tea for newcomers"*
- **S: (n) tea**, [Camellia sinensis](#) (a tropical evergreen shrub or small tree extensively cultivated in

WordNet Example (II)

Hypernyms of "tea", Sense 2

WordNet Search - 2.1 - Konqueror

Location: wn7o2=&o0=1&o7=&o5=&o1=1&o6=&o4=&o3=&r=1&s=tea&i=2&h=0100000#c

- **S: (n) tea** (a beverage made by steeping tea leaves in water) *"iced tea is a cooling drink"*
- **S: (n) tea**, [afternoon tea](#), [teatime](#) (a light midafternoon meal of tea and sandwiches or cakes) *"an Englishman would interrupt a war to have his afternoon tea"*
 - ◊ [direct hypernym](#) / [inherited hypernym](#) / [sister term](#)
 - **S: (n) meal**, [repast](#) (the food served and eaten at one time)
 - **S: (n) nutriment**, [nourishment](#), [nutrition](#), [sustenance](#), [aliment](#), [alimentation](#), [victuals](#) (a source of materials to nourish the body)
 - **S: (n) food**, [nutrient](#) (any substance that can be metabolized by an organism to give energy and build tissue)
 - **S: (n) substance**, [matter](#) (that which has mass and occupies space) *"an atom is the smallest indivisible unit of matter"*
 - **S: (n) physical entity** (an entity that has physical existence)
 - **S: (n) entity** (that which is perceived or known or inferred to have its own distinct existence (living or nonliving))
 - ◊ [domain region](#)
- **S: (n) tea**, [tea leaf](#) (dried leaves of the tea shrub; used to make tea) *"the store shelves held many different kinds of tea"; "they threw the tea into Boston harbor"*
- **S: (n) tea** (a reception or party at which tea is served) *"we met at the Dean's tea for newcomers"*
- **S: (n) tea**, [Camellia sinensis](#) (a tropical evergreen shrub or small tree extensively cultivated in e.g. China and Japan and India; source of tea leaves) *"tea has fragrant white flowers"*

Logical Forms and Predicate-Argument Structures

Transforming Text into Logical Units

Suppose we found the correct sense for each word. We can now transform the text into a formal representation, e.g., first-order predicate logic or description logics.

- knowledge is encoded independently from the textual description (e.g., “X bought A” and “A was acquired by X” both encode the same information)
- with this, **formal reasoning** becomes possible

Predicate-Argument Structures

Convert text into logical structures using predicates:

- $company(x_1) \wedge company(x_2) \wedge buy-act(x_1, x_2)$

PA structures can be derived from parse and additionally incorporate semantic information (e.g., using WordNet).

Logical Forms and Predicate-Argument Structures

Transforming Text into Logical Units

Suppose we found the correct sense for each word. We can now transform the text into a formal representation, e.g., first-order predicate logic or description logics.

- knowledge is encoded independently from the textual description (e.g., “X bought A” and “A was acquired by X” both encode the same information)
- with this, **formal reasoning** becomes possible

Predicate-Argument Structures

Convert text into logical structures using predicates:

- $company(x_1) \wedge company(x_2) \wedge buy-act(x_1, x_2)$

PA structures can be derived from parse and additionally incorporate semantic information (e.g., using WordNet).

Pragmatics: Coreference Resolution

Problem

Entities in natural language texts are not identified with convenient unique IDs, but rather with constantly changing descriptions.

Example: *Mr. Bush, The president, he, George W., ...*

Solution

Automatic detection and collection of all textual descriptors that refer to the same entity within a *coreference chain*.

- can be used to find information about an entity, even when referenced by a different name
- important for many higher-level text analysis tasks

Coreference Resolution Algorithms

Pronominal coreferences can be detected quite reliably (also called *Anaphora Resolution*). Full (nominal) coreference resolution is hard.

Pragmatics: Coreference Resolution

Problem

Entities in natural language texts are not identified with convenient unique IDs, but rather with constantly changing descriptions.

Example: *Mr. Bush, The president, he, George W., ...*

Solution

Automatic detection and collection of all textual descriptors that refer to the same entity within a *coreference chain*.

- can be used to find information about an entity, even when referenced by a different name
- important for many higher-level text analysis tasks

Coreference Resolution Algorithms

Pronominal coreferences can be detected quite reliably (also called *Anaphora Resolution*). Full (nominal) coreference resolution is hard.

Pragmatics: Coreference Resolution

Problem

Entities in natural language texts are not identified with convenient unique IDs, but rather with constantly changing descriptions.

Example: *Mr. Bush, The president, he, George W., ...*

Solution

Automatic detection and collection of all textual descriptors that refer to the same entity within a *coreference chain*.

- can be used to find information about an entity, even when referenced by a different name
- important for many higher-level text analysis tasks

Coreference Resolution Algorithms

Pronominal coreferences can be detected quite reliably (also called *Anaphora Resolution*). Full (nominal) coreference resolution is hard.

Evaluation of NLP Systems

General Approach

The results of a system are compared to a manually created *gold standard* using various metrics.

Main Challenges

Manually annotating large amounts of texts for specific linguistic phenomena is **very** time-consuming (thus expensive):

- test set needs to be different from training set
- for some tasks, two or more annotations of the same data are needed (to measure *inter-annotator agreement*)

Annotated Corpora

For some tasks (e.g., POS tagging), annotated corpora are (freely) available.

Evaluation of NLP Systems

General Approach

The results of a system are compared to a manually created *gold standard* using various metrics.

Main Challenges

Manually annotating large amounts of texts for specific linguistic phenomena is **very** time-consuming (thus expensive):

- test set needs to be different from training set
- for some tasks, two or more annotations of the same data are needed (to measure *inter-annotator agreement*)

Annotated Corpora

For some tasks (e.g., POS tagging), annotated corpora are (freely) available.

Evaluation of NLP Systems

General Approach

The results of a system are compared to a manually created *gold standard* using various metrics.

Main Challenges

Manually annotating large amounts of texts for specific linguistic phenomena is **very** time-consuming (thus expensive):

- test set needs to be different from training set
- for some tasks, two or more annotations of the same data are needed (to measure *inter-annotator agreement*)

Annotated Corpora

For some tasks (e.g., POS tagging), annotated corpora are (freely) available.

Evaluation Measures

Accuracy and Error

Simplest measure are *accuracy* (percentage of correct results) and *error* (percentage of wrong results).

- not often used, as they are very insensitive to the interesting numbers
- reason is the usually large number of non-relevant and non-selected entities that is “hiding” all other numbers
- in other words, accuracy only reacts to real errors, and doesn't show how many correct results have been found as such

Precision and Recall

Precision

Like in Information Retrieval, *Precision* show the percentage of correct results within an answer:

$$\text{Precision} = \frac{\text{Correct} + \frac{1}{2}\text{Partial}}{\text{Correct} + \text{Spurious} + \frac{1}{2}\text{Partial}}$$

Recall

And *Recall* the percentage of the correct system results over all correct results:

$$\text{Recall} = \frac{\text{Correct} + \frac{1}{2}\text{Partial}}{\text{Correct} + \text{Missing} + \frac{1}{2}\text{Partial}}$$

Tradeoff

Note that you can always get 100% Precision by selecting nothing and 100% Recall by selecting everything. However, in NLP there is often no clear trade-off between the two.

Precision and Recall

Precision

Like in Information Retrieval, *Precision* show the percentage of correct results within an answer:

$$\text{Precision} = \frac{\text{Correct} + \frac{1}{2}\text{Partial}}{\text{Correct} + \text{Spurious} + \frac{1}{2}\text{Partial}}$$

Recall

And *Recall* the percentage of the correct system results over all correct results:

$$\text{Recall} = \frac{\text{Correct} + \frac{1}{2}\text{Partial}}{\text{Correct} + \text{Missing} + \frac{1}{2}\text{Partial}}$$

Tradeoff

Note that you can always get 100% Precision by selecting nothing and 100% Recall by selecting everything. However, in NLP there is often no clear trade-off between the two.

Precision and Recall

Precision

Like in Information Retrieval, *Precision* show the percentage of correct results within an answer:

$$\text{Precision} = \frac{\text{Correct} + \frac{1}{2}\text{Partial}}{\text{Correct} + \text{Spurious} + \frac{1}{2}\text{Partial}}$$

Recall

And *Recall* the percentage of the correct system results over all correct results:

$$\text{Recall} = \frac{\text{Correct} + \frac{1}{2}\text{Partial}}{\text{Correct} + \text{Missing} + \frac{1}{2}\text{Partial}}$$

Tradeoff

Note that you can always get 100% Precision by selecting nothing and 100% Recall by selecting everything. However, in NLP there is often no clear trade-off between the two.

F-Measure and IAA

Combining Precision and Recall

Often a combined measure of Precision and Recall is helpful. This can be done using the *F-Measure* (equal weight for $\beta = 1$):

$$\text{F-measure} = \frac{(\beta^2 + 1)P \cdot R}{(\beta^2 R) + P}$$

Measuring Inter-Annotator Agreement

There are many measures for computing IAA (Cohen's Kappa, prevalence, bias, ...), depending on the concrete task. One way to obtain the IAA is to compute P , R , and F values between two humans and averaging the results of $P(H_1)$ vs. $P(H_2)$ and $P(H_2)$ vs. $P(H_1)$.

In essence, IAA shows how *hard* a task is: if humans cannot agree on the correct result in more than 90% of all cases, don't expect your system to be better!

F-Measure and IAA

Combining Precision and Recall

Often a combined measure of Precision and Recall is helpful. This can be done using the *F-Measure* (equal weight for $\beta = 1$):

$$\text{F-measure} = \frac{(\beta^2 + 1)P \cdot R}{(\beta^2 R) + P}$$

Measuring Inter-Annotator Agreement

There are many measures for computing IAA (Cohen's Kappa, prevalence, bias, ...), depending on the concrete task. One way to obtain the IAA is to compute P , R , and F values between two humans and averaging the results of $P(H_1)$ vs. $P(H_2)$ and $P(H_2)$ vs. $P(H_1)$.

In essence, IAA shows how *hard* a task is: if humans cannot agree on the correct result in more than 90% of all cases, don't expect your system to be better!

Evaluation Example

Evaluation of a Noun Phrase (NP) Chunker

Annotation Diff Tool

Key: Annotation Set: Annotation Type: F-Measure Weight:

Response: [Default set] Features: All Some None

Start	End	Key	Features	Start	End	Response	Feat
1503	1520	10 inches of rain	0	1503	1512	10 inches	(HEAD_END=1512, HEAD_START=1506, HEAD...
805	823	the larger systems	0	805	826	the larger systems we	(DET=the, MOD=larger systems, HEAD=we, H...
1570	1580	last night	0	1564	1580	Haiti last night	(HEAD_END=1580, HEAD_START=1575, HEAD...
2109	2132	Guayanilla, Puerto Rico	0	2109	2119	Guayanilla	(HEAD=Guayanilla, HEAD_START=2109, HEAD...
187	209	100-mile-an-hour winds	0	191	209	mile-an-hour winds	(HEAD_END=209, HEAD_START=204, HEAD=v...
824	826	we	0				
676	681	south	0				
697	717	the Caribbean island	0				
1564	1569	Haiti	0				
1855	1866	a hurricane	0				
172	181	yesterday	0				
342	351	the roofs	0				
				1086	1093	Jamaica	(HEAD=Jamaica, HEAD_START=1086, HEAD_E...
				80	101	LATAM SANTO DOMINGO	(HEAD_END=101, HEAD_START=88, HEAD=SA...
				2121	2132	Puerto Rico	(HEAD=Puerto Rico, HEAD_START=2121, HEA...

Correct: 92 Recall Precision F-Measure

Partially Correct: 20 Strict: 0.7731 0.7797 0.7764

Missing: 7 Lenient: 0.9412 0.9492 0.9451

False Positives: 6 Average: 0.8571 0.8644 0.8608

More Complex Metrics

OK, but...

...how do I define precision and recall for more complex tasks?

- Parsing Sentences (need to compare *parse trees*)
- Coreference Chains (need to compare *graphs*)
- Automatic Summaries (need to compare *whole texts*)

Parser Evaluation: The PARSEVAL Measure

A classical measure for parser evaluation is *PARSEVAL*. Compare a gold-standard parse tree to a system's one by segmenting it into its constituents (brackets). Then:

Precision is the number of brackets appearing the gold standard;

Recall measures how many of the gold standard's brackets are in the parse

Crossing Brackets measures how many brackets are crossing on average

More Complex Metrics

OK, but...

...how do I define precision and recall for more complex tasks?

- Parsing Sentences (need to compare *parse trees*)
- Coreference Chains (need to compare *graphs*)
- Automatic Summaries (need to compare *whole texts*)

Parser Evaluation: The *PARSEVAL* Measure

A classical measure for parser evaluation is *PARSEVAL*. Compare a gold-standard parse tree to a system's one by segmenting it into its constituents (brackets). Then:

Precision is the number of brackets appearing the gold standard;

Recall measures how many of the gold standard's brackets are in the parse

Crossing Brackets measures how many brackets are crossing on average

Evaluation: Summary

Some remarks

- Evaluation is often *very* expensive due to the large amount of time needed for manually annotating documents
- For some tasks (e.g., automatic summarization) the evaluation can be (almost) as difficult as the task itself
- Development of metrics for certain tasks, as well as the evaluation of evaluation metrics, is another branch of research
- Due to the high costs involved, and in order to ensure comparability of the results, the NLP community organises various *competitions* where system developers participate in solving prescribed tasks on the same data, using the same evaluation metrics. Examples are MUC, TREC, DUC, BioCreAtIvE, ...

Recommended Literature

NLP Foundations

- Daniel Jurafsky and James H. Martin, *Speech and Language Processing*, Prentice Hall, 2000
- Christopher D. Manning and Hinrich Schütze, *Foundations of Statistical Natural Language Processing*, MIT Press, 1999.

Online

- Statistical natural language processing and corpus-based computational linguistics: An annotated list of resources
<http://www-nlp.stanford.edu/links/statnlp.html>

Major Conferences

ACL, NAACL, EAACL, COLING, HLT, EMNLP, LREC, ANLP, NLDB, ...

Part III

Technology

- 10 Technology
 - Toolkits and Frameworks
- 11 GATE
 - GATE Overview
 - JAPE Transducers
- 12 Example: Information Extraction with ANNIE
 - The Task
 - Step 1: Tokenization
 - Step 2: Gazetteering
 - Step 3: Sentence Splitting
 - Step 4: Part-of-Speech (POS) Tagging
 - Step 5: Named Entity (NE) Detection
 - Step 6: Coreference Resolution
- 13 Other Resources
 - More GATE Plugins
 - SUPPLE
 - MuNPE_x
 - The Durm German Lemmatizer
- 14 References

So you want to build a Text Mining system...

Requirements

A TM system requires a large amount of infrastructure work:

- Document handling, in various formats (plain text, HTML, XML, PDF, ...), from various sources (files, DBs, email, ...)
- Annotation handling (stand-off markup)
- Component implementations for standard tasks, like Tokenizers, Sentence Splitters, Part-of-Speech (POS) Taggers, Finite-State Transducers, Full Parsers, Classifiers, Noun Phrase Chunkers, Lemmatizers, Entity Taggers, Coreference Resolution Engines, Summarizers, ...

As well as *resources* for concrete tasks and languages:

- Lexicons, WordNets
- Grammar files and Language models
- etc.

Existing Resources

Fortunately, you don't have to start from scratch

Many (open source) tools and resources are available:

Tools: programs performing a single task, like classifiers, parsers, or NP chunkers

Frameworks: integrating architectures for combining and controlling all components and resources of an NLP system

Resources: for various languages, like lexicons, wordnets, or grammars

GATE and UIMA

Major Frameworks

Two important frameworks are:

- GATE (*General Architecture of Text Engineering*), under development since 1995 at University of Sheffield, UK
- UIMA (*Unstructured Information Management Architecture*), developed by IBM

Both frameworks are open source (GATE: LGPL, UIMA: CPL)

In the following, we will focus on GATE only.

General Architecture for Text Engineering (GATE)

GATE features

GATE (*General Architecture for Text Engineering*) is a component framework for the development of NLP applications.

Rich Infrastructure: XML Parser, Corpus management, Unicode handling, Document Annotation Model, Finite State Transducer (JAPE Grammar), etc.

Standard Components: Tokeniser, Part-of-Speech (POS) Tagger, Sentence Splitter, etc.

Set of NLP tools: Information Retrieval (IR), Machine Learning, Database access, Ontology editor, Evaluation tool, etc.

Clean Framework: Java Beans component model; Other tools can easily be integrated into GATE via *Wrappers*

- File Options Tools Help
- Gate
 - Applications
 - Language Resources
 - Processing Resources
 - CLaC Verb Grouper
 - CLaC AnnotationSetTransfer-NP
 - CLaC Auxiliaries Transducer
 - CLaC DE-NPE
 - CLaC StupidStemmer
 - CLaC DE TreeTagger
 - CLaC DE-Gazetteer
 - CLaC DE AAMarker
 - ANNIE OrthoMatcher_00080
 - CLaC Summarizer_00091
 - CLaC Classifier_00046
 - CLaC FuzzyCoreferencer_0004C
 - CLaC Noun Phrase Extractor
 - CLaC AnnotationSetTransfer_0005D
 - ANNIE VP Chunker_0005A
 - ANNIE OrthoMatcher_00057
 - CLaC NE Transducer
 - ANNIE POS Tagger_00052
 - CLaC HeadPhrase Remover
 - CLaC SentenceSplitter_0008D
 - CLaC Gazetteer
 - CLaC NumberCombiner
 - CLaC AAMarker_00049
 - ANNIE English Tokeniser_00041
 - CLaC DocumentResetter_0003E
 - Data stores

Messages AI CLaC DUCTape Summarizer

Text Annotations Annotation Sets Coreference Print

conference - many of these cases have occurred, unfortunately, during the dinner hour."

Added Fener: "He had known about al-Qaeda's practice of raising money through drug trafficking and money laundering, but it seems the full scope of their depravity had barely been imagined."

The video is not the only evidence of telemarketing activity within al-Qaeda. According to Fener, CIA agents tracking the terrorist organization over the past 12 months made steady progress infiltrating its communications network, eventually gaining access to transmissions to and from al-Qaeda operatives. These transmissions included a number of telemarketing "cold calls" to randomly chosen U.S. citizens.

Last December, during a sweep of caves near the Afghan-Pakistani border, Maj. Gen. Dan K. McNeill, leader of U.S. forces in Afghanistan, unearthed further evidence corroborating the phone-solicitation theory. Inside one cave, McNeill and his troops found a bank of empty cubicles with individual phone lines, a bullhorn, and 10 desktop bells, commonly rung in the event of a "sale."

"I couldn't believe what I saw," said McNeill, who also discovered bomb-making instructions and detailed maps of U.S. landmarks in the cave. "On top of all the destruction these people had already unleashed, plans were underway to harass the American people with a merciless assault of offers for everything from discounts on home DSL lines to pre-approved, low-interest credit cards."

For all the evidence collected by the CIA, the "smoking gun" in the investigation may turn out to be an alleged Osama bin Laden motivational videotape, currently in the possession of CNN. The controversial tape, which has never aired on the cable network, is rumored to feature bin Laden urging his followers to think positive and believe in the quality of the product they are pitching, closing on the grim slogan "Smile. And Dial."

type	Set	Start	End	Features
P	Default	3582	3596	{DET="", MOD="", HEAD="Guantanamo Bay"}
P	Default	776	791	{DET="the", MOD="dinner", HEAD="hour"}
P	Default	2259	2262	{DET="", MOD="", HEAD="out"}
P	Default	1806	1807	{DET="", MOD="", HEAD="I"}
P	Default	3849	3852	{DET="", MOD="", HEAD="one"}
P	Default	987	996	{DET="The", MOD="", HEAD="video"}
P	Default	1487	1494	{DET="", MOD="", HEAD="McNeill"}
P	Default	2280	2318	{DET="", MOD="Osama bin Laden motivational", HEAD="videotape"}
P	Default	894	910	{DET="", MOD="money", HEAD="laundering"}
P	Default	4028	4034	{DET="", MOD="", HEAD="action"}
P	Default	3382	3401	{DET="", MOD="bored", HEAD="receptionists"}
P	Default	3046	3058	{DET="", MOD="Sears", HEAD="charge"}

- caves (3)
- its communications network
- me (2)
- money (3)
- phone lines (2)
- the \$3,000 (2)
- the CIA (2)
- the only evidence (4)
- these people (4)
- transmissions (2)
- Default annotations
 - Chain
 - Classification
 - Date
 - firstPerson
 - JobTitle
 - Location
 - Lookup
 - Money
 - NP
 - Organization
 - Person
 - Sentence
 - SpaceToken
 - Split
 - Title
 - Token
 - Unknown
 - Vc
- Original markups annotations
 - Paragraph
- Summary annotations
 - Selected Chains
 - Summary
- ToBeParsed annotations

GATE Concepts

A *Processing Pipeline* holds the required components

Component-based applications, assembled at run-time:

Messages * CLaC DUCTape Summarizer * Profiler

Loaded Processing resources

Name	Type
ANNIE OrthoMatcher_00080	ANNIE OrthoMatch
CLaC AnnotationSetTransfer-NP	CLaC AnnotationS
CLaC Auxiliaries Transducer	Jape Transducer
CLaC DE AAMarker	CLaC AAMarker
CLaC DE TreeTagger	CLaC TreeTagger
CLaC DE-Gazetteer	ANNIE Gazetteer
CLaC DE-NPE	CLaC EarleyParser
CLaC StupidStemmer	Jape Transducer
CLaC Verb Grouper	CLaC EarleyParser

Selected Processing resources

Name	Type
CLaC DocumentResetter_0003E	CLaC DocumentResetter
ANNIE English Tokeniser_00041	ANNIE English Tokeniser
CLaC AAMarker_00049	CLaC AAMarker
CLaC NumberCombiner	Jape Transducer
CLaC Gazetteer	ANNIE Gazetteer
CLaC SentenceSplitter_0008D	CLaC SentenceSplitter
CLaC HeadPhrase Remover	CLaC HeadPhrase Remover for DUC 20
ANNIE POS Tagger_00052	ANNIE POS Tagger
CLaC NE Transducer	ANNIE NE Transducer
ANNIE OrthoMatcher_00057	ANNIE OrthoMatcher
ANNIE VP Chunker_0005A	ANNIE VP Chunker
CLaC AnnotationSetTransfer_0005D	CLaC AnnotationSetTransfer
CLaC Noun Phrase Extractor	CLaC EarleyParser
CLaC FuzzyCoreferencer_0004C	CLaC FuzzyCoreferencer
CLaC Classifier_00046	CLaC Classifier
CLaC Summarizer_00091	CLaC Summarizer

Results are exchanged between the components through document *annotations*.

Finite-State Language Processing with GATE

JAPE Transducers

JAPE (*Java Annotation Patterns Engine*) is a component to build finite-state **transducers** running over annotations from grammars.

- this is an application of *finite-state language processing*
- Transducers are basically (non-deterministic) finite-state machines, running over a graph data structure
- expressiveness of JAPE grammars corresponds to regular expressions
- basic format of a JAPE rule: LHS:RHS
left-hand side matches annotations in documents, right-hand side adds annotations
- Java code can be included on the RHS, allowing computations that cannot be expressed in JAPE alone

Example for a JAPE grammar rule

Finding IP Addresses

```
// IP Address Rules
Rule: IPaddress1
(
    {Token.kind == number}
    {Token.string == "."}
    {Token.kind == number}
    {Token.string == "."}
    {Token.kind == number}
    {Token.string == "."}
    {Token.kind == number}
):ipAddress -->
:ipAddress.Ip = {kind = "ipAddress", rule = "IPaddress1"}
```

Results

- matches e.g. 141.3.49.133.
- for each detected address an **annotation** is added to the document at the matching start- and end-positions

A Nearly-New Information Extraction System (ANNIE)

Task: Find all *Persons* mentioned in a document

- A simple “search” function doesn’t help here
- What we need is *Information Extraction* (IE), particularly *Named Entity (NE) Detection* (entity-type Person)

ANNIE

GATE includes an example application, ANNIE, which can solve this task.

- developed for the news domain (newspapers, newswires), but can be adapted to other domains
- good starting point to practice NLP, IE, and TM

A Nearly-New Information Extraction System (ANNIE)

Task: Find all *Persons* mentioned in a document

- A simple “search” function doesn’t help here
- What we need is *Information Extraction* (IE), particularly *Named Entity (NE) Detection* (entity-type Person)

ANNIE

GATE includes an example application, ANNIE, which can solve this task.

- developed for the news domain (newspapers, newswires), but can be adapted to other domains
- good starting point to practice NLP, IE, and TM

Persons detected by ANNIE

GATE 3.1-beta1 build 2233

File Options Tools Help

GATE

- Applications
 - ANNIE_00016
- Language Resources
 - GATE document_00023
 - GATE corpus_00022
- Processing Resources
 - ANNIE OrthoMatcher_00021
 - ANNIE NE Transducer_00020
 - ANNIE POS Tagger_0001F
 - ANNIE Sentence Splitter_0001C
 - ANNIE Gazetteer_0001B
 - ANNIE English Tokeniser_00018
 - Document Reset PR_00017
- Data stores

MatchesAnnots {null={}}

MimeType text/plain

docNewLineType LF

gate.SourceURL file:/home/witte

ANNIE_00016 run in 0.716 seconds

Messages GATE corpus_00022 ANNIE_00016 GATE document_00023

Annotation Sets Annotations Co-reference Editor Text

stringently than the common kind of asbestos, chrysotile, found in most schools and other buildings, Dr. Talcott said. The U.S. is one of the few industrialized nations that doesn't have a higher standard of regulation for the smooth, needle-like fibers such as crocidolite that are classified as amphiboles, according to Brooke T. Mossman, a professor of pathology at the University of Vermont. Mossman explained. In July, the Environ on virtually all use uses of cancer-causi

Person

Type	Set	Start	End	
Person	1017	1021		
Person	1103	1119		
Person	1162	1173		(gender= null, rule= PersonFinal, rule1= PersonTitle)
Person	2043	2054		(gender= null, rule= PersonFinal, rule1= PersonTitle)
Person	2609	2620		(gender= null, rule= PersonFinal, rule1= PersonTitle)
Person	2838	2848		(gender= null, rule= PersonFinal, rule1= PersonFull)
Person	3006	3017		(gender= null, rule= PersonFinal, rule1= PersonTitle)
Person	3268	3272		(gender= male, rule= PersonFinal, rule1= GazPersonFirst)
Person	3815	3831		(gender= male, rule= PersonFinal, rule1= PersonFull)

9 Annotations (1 selected)

Document Editor Initialisation Parameters

- Date
- FirstPerson
- JobTitle
- Location
- Lookup
- Organization
- Person
- Sentence
- SpaceToken
- Split
- Title
- Token
- Unknown
- Original markups
 - paragraph

New

Step 1: Tokenization

Tokenization Component

Tokenization is performed in two steps:

- a generic Unicode Tokeniser is fed with tokenisation rules for English
- afterwards, a grammer changes some of these tokens for later processing: e.g., “don’t” results in three tokens: “don”, “'”, and “t”. This is converted into two tokens, “do” and “n’t” for downstream components

For each detected token, a corresponding `Token` annotation is added to the document.

Step 1: Tokenization (Example)

Example Tokenisation Rules

```
#numbers#
```

```
// a number is any combination of digits
"DECIMAL_DIGIT_NUMBER"+ >Token;kind=number;
```

```
#whitespace#
```

```
(SPACE_SEPARATOR) >SpaceToken;kind=space;
(CONTROL) >SpaceToken;kind=control;
```

Example Output

Type	Set	Start	End	Features
Token		158	163	{kind=word, length=5, orth=lowercase, string=years}
SpaceToken		163	164	{kind=space, length=1, string= }
Token		164	167	{kind=word, length=3, orth=lowercase, string=ago}
Token		167	168	{kind=punctuation, length=1, string=,}
SpaceToken		168	169	{kind=space, length=1, string= }
Token		169	180	{kind=word, length=11, orth=lowercase, string=researchers}
SpaceToken		180	181	{kind=space, length=1, string= }

1417 Annotations (0 selected)

Step 2: Gazetteering

Gazetteer Component

The *Gazetteer* uses structured plain text lists to annotate words with a `major_type` and `minor_type`

- each lists represents a concept or type, e.g., female first names, mountains, countries, male titles, streets, festivals, dates, planets, organizations, cities, ...
- ambiguities are not resolved at this step—e.g., a string can be annotated both as *female first name* and *city*
- GATE provides several different Gazetteer implementation: Simple Gazetteer, HashGazetteer, FlexibleGazetteer, OntoGazetteer, ...
- Gazetteer lists can be (a) created by hand, (b) derived from databases, (c) “learned” through patterns, e.g., from web sites

Step 2: Gazetteering (Example)

Gazetteer Definition

Connecting lists with major/minor types:

```
organization.lst:organization
organization_nouns.lst:organization_noun
person_ambig.lst:person_first:ambig
person_ending.lst:person_ending
person_female.lst:person_first:female
person_female_cap.lst:person_first:female
person_female_lower.lst:person_first:female
person_full.lst:person_full
```

Example List

Person_female.lst:

Acantha
Acenith
Achala
Achava
Achсах
Ada
Adah
Adalgisa

9.8 billion Kent with the filters were sold, the company said.
Among 33 men who with the substance, 28 have died -- more than three times the expected number. Four of the five surviving workers have asbestos-related diseases including three

Type	Set	Start	End	Features
Lookup		1480	1490	{majorType=stop}
Lookup		1509	1516	{majorType=number}
Lookup		1517	1521	{majorType=person_first, minorType=male}
Lookup		1517	1521	{majorType=location, minorType=region}
Lookup		1565	1572	{majorType=organization_noun}

94 Annotations (1selected)

Step 3: Sentence Splitting

Task: Split Stream of Tokens into Sentences

Sentences are important units in texts

- Correct detection important for downstream components, e.g., the POS-Tagger

Precise splitting can be annoyingly hard:

- a “.” (dot) often does **not** indicate an EOS
- Abbreviations “*The U.S. government*”, but: “*... announced by the U.S.*”
- Ambiguous boundaries “!” , “;” , “:” , nested sentences (e.g., inside quotations) etc.
- Formatting detection (headlines, footnotes, tables, ...)

ANNIE Sentence Splitter

Uses grammar rules and abbreviation lists to detect sentence boundaries.

Step 4: Part-of-Speech (POS) Tagging

Producing POS Annotations

POS-Tagging assigns a part-of-speech-tag (POS tag) to each Token.

- GATE includes the Hepple tagger for English, which is a modified version of the Brill tagger

Example output

inc., the unit of New York-based Loews Corp. that makes Kent cigarettes, stopped using crocidolite in its Micronite cigarette filters in 1956.

Although preliminary findings were reported more than a year ago, the latest results appear in today's New England Journal of

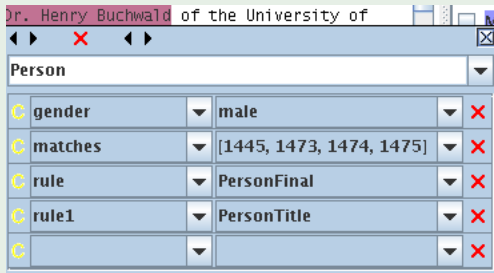
Type	Set	Start	End	Features
Token		485	494	{category=NN, kind=word, length=9, orth=upperlit
Token		495	504	{category=NN, kind=word, length=9, orth=lowercas
Token		505	512	{category=NNS, kind=word, length=7, orth=lowerca
Token		513	515	{category=IN, kind=word, length=2, orth=lowercase
Token		516	520	{category=CD, kind=number, length=4, string=1956

Step 5: Named Entity (NE) Detection

Transducer-based NE Detection

Using all the information obtained in the previous steps (Tokens, Gazetteer lookups, POS tags), ANNIE now runs a sequence of JAPE-Transducers to detect Named Entities (NE)s.

Example for a detected *Person*



The screenshot shows a window with the text "Dr. Henry Buchwald of the University of". Below the text, a dropdown menu is set to "Person". A table displays the following details for the detected entity:

gender	male	X
matches	[1445, 1473, 1474, 1475]	X
rule	PersonFinal	X
rule1	PersonTitle	X
		X

We can now look at the grammar rules that found this person.

Entity Detection: Finding Persons

Strategy

A JAPE grammar rule combines information obtained from POS-tags with Gazetteer lookup information

- although the last name in the example is not in any list, it can be found based on its POS tag and an additional first name/last name rule (not shown)
- many additional rules for other Person patterns, as well as Organizations, Dates, Addresses, ...

Persons with Titles

```
Rule:   PersonTitle
Priority: 35
(
  {Token.category == DT}|
  {Token.category == PRP}|
  {Token.category == RB}
)?
(
  (TITLE)+
  ((FIRSTNAME | FIRSTNAMEAMBIG
   | INITIALS2)
  )?
  (PREFIX)*
  (UPPER)
  (PERSONENDING)?
)
:person --> ...
```

Step 6: Coreference Resolution

Finding Coreferences

Remember the problem of coreference resolution:

- need to find all instances of an entity in a text,
- even when referred to by different textual descriptors

Coreference resolution in ANNIE

GATE provides two components for performing a restricted subset of coreference resolution:

Pronominal Coreferences finds anaphors (e.g., “he” referring to a previously mentioned person) and also some cataphors (e.g., “Before *he* bought the car, *John*..”)

Nominal Coreferences a number of JAPE rules match entities based on orthographic features, e.g., a person “John Smith” will be matched with “Mr. Smith”

Coreference Resolution Example

The screenshot displays the GATE 3.1-beta1 build 2233 interface. The main window shows a text document with several sentences. The text is as follows:

Iraqi President Saddam Hussein seemed determined to solve his financial problems and fulfill territorial ambitions by dethroning the government of neighboring Kuwait.

The invasion, unprecedented in modern Arab history, reflected the brutality Saddam has used to crush all opposition at home since coming to power in 1979.

The sentence for criticizing Saddam is death, and the president considered Kuwait's oil overproduction and the subsequent fall in the price of oil a personal affront to Iraq.

"Iraqis will not forget the saying that cutting necks is better than cutting means of living. Oh God Almighty, be witness that we have warned them," he said last month after prices tumbled below

The interface shows the following components:

- Annotations:** The text is annotated with various entities and relations. For example, "Saddam Hussein" is highlighted in red, "Kuwait" in yellow, "Iraq" in blue, "Saudi Arabia" in yellow, "Iran" in blue, "United Arab Emirates" in pink, "Israel" in blue, "President Saddam Hussein" in red, "Gulf Cooperation Council" in green, "AP News" in blue, "Abdel Nassar" in blue, and "President Hosni Mubarak" in pink.
- Co-reference Editor:** A list of entities is shown on the right, with checkboxes indicating their coreference status. The checked entity is "President Saddam Hussein".
- Document Editor:** The text is displayed in a scrollable area with various annotations.
- Initialisation Parameters:** A section at the bottom of the window.

The status bar at the bottom indicates: ANNIE_00016 run in 1.87 seconds

GATE Plugins

More GATE Plugins

GATE comes with a number of other language plugins, which are either implemented directly for GATE, or use *wrappers* to access external resources:

Verb Grouper: a JAPE grammar to analyse verb groups (VGs)

SUPPLE Parser: a Prolog-based parser for (partial) parsing that can create logical forms

Chemistry Tagger: component to find chemistry items (formulas, elements etc.)

Web Crawler: wrapper for the Websphinx crawler to construct a corpus from the Web

Kea Wrapper: for the *Kea* keyphrase detector

Ontology tools: for using (Jena) ontologies in pipelines, e.g., with the OntoGazetteer and Ontology-aware JAPE transducer

GATE Plugins

Plugin Management Console

Name	URL	Load now	Load always	Delete	CREOLE resources in directory
Machine_Learning	file:/usr/local/clactools/GATE3_CVS/gate/plugins/Machine_Learning/	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	
Kea	file:/usr/local/clactools/GATE3_CVS/gate/plugins/Kea/	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	
rasp	file:/usr/local/clactools/GATE3_CVS/gate/plugins/rasp/	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	
uima	file:/usr/local/clactools/GATE3_CVS/gate/plugins/uima/	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	
ANNIE	file:/usr/local/clactools/GATE3_CVS/gate/plugins/ANNIE/	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	
Tools	file:/usr/local/clactools/GATE3_CVS/gate/plugins/Tools/	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	
crawl	file:/usr/local/clactools/GATE3_CVS/gate/plugins/crawl/	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	
hindi	file:/usr/local/clactools/GATE3_CVS/gate/plugins/hindi/	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	
NLG_Tools	file:/usr/local/clactools/GATE3_CVS/gate/plugins/NLG_Tools/	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	
SUPPLE	file:/usr/local/clactools/GATE3_CVS/gate/plugins/SUPPLE/	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	
Montreal_Transducer	file:/usr/local/clactools/GATE3_CVS/gate/plugins/Montreal_Transducer/	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	
WordNet	file:/usr/local/clactools/GATE3_CVS/gate/plugins/WordNet/	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	
Buchart	file:/usr/local/clactools/GATE3_CVS/gate/plugins/Buchart/	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	
spanish	file:/usr/local/clactools/GATE3_CVS/gate/plugins/spanish/	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	
Information_Retrieval	file:/usr/local/clactools/GATE3_CVS/gate/plugins/Information_Retrieval/	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	
arabic	file:/usr/local/clactools/GATE3_CVS/gate/plugins/arabic/	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	
google	file:/usr/local/clactools/GATE3_CVS/gate/plugins/google/	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	
romanian	file:/usr/local/clactools/GATE3_CVS/gate/plugins/romanian/	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	
Obsolete	file:/usr/local/clactools/GATE3_CVS/gate/plugins/Obsolete/	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	
Jape_Compiler	file:/usr/local/clactools/GATE3_CVS/gate/plugins/Jape_Compiler/	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	
TreeTagger	file:/usr/local/clactools/GATE3_CVS/gate/plugins/TreeTagger/	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	
Chemistry_Tagger	file:/usr/local/clactools/GATE3_CVS/gate/plugins/Chemistry_Tagger/	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	
Minipar	file:/usr/local/clactools/GATE3_CVS/gate/plugins/Minipar/	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	
Stemmer	file:/usr/local/clactools/GATE3_CVS/gate/plugins/Stemmer/	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	
Minorthird	file:/usr/local/clactools/GATE3_CVS/gate/plugins/Minorthird/	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	
probability	file:/usr/local/clactools/GATE3_CVS/gate/plugins/probability/	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	
cebuano	file:/usr/local/clactools/GATE3_CVS/gate/plugins/cebuano/	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	
Ontology_Tools	file:/usr/local/clactools/GATE3_CVS/gate/plugins/Ontology_Tools/	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	
NP_Chunking	file:/usr/local/clactools/GATE3_CVS/gate/plugins/NP_Chunking/	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	
build	file:/home/witte/Repository/durm/Components/FuzzyCoreferencer/build/	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	

You can also

SUPPLE Parser

Bottom-up Parser for English

Constructs (partial) syntax trees and logical forms for English sentences. Implemented in Prolog.

determined to solve his financial problems and fulfill territorial ambitions by dethroning the government of neighboring Kuwait.

The invasion, unprecedented in modern Arab history, reflected the brutality Saddam has used to crush all opposition at home since coming to power in 1979.

The sentence for criticizing Saddam is death, and the president considered Kuwait's oil overproduction and the subsequent fall in the price of oil a personal affront to Iraq.

"Iraqis will not forget the saying that cutting necks is better than cutting means of living. Oh God Almighty, be witness that we have warned them," he said last month after prices tumbled below

Type	Set	Start	End
{qlf=[name(e60, 'Saddam'), person(e60), '_'(e61), qual(e60, e61), number(e61, sing), realisati			
{qlf=[at(e65, e66), name(e66, home), '_'(e66), realisation(e66, offsets(630, 634))]			
{qlf=[since(e67, e68), '_'(e68), time(e68, none), aspect(e68, simple), voice(e68, active), incom			
{qlf=['_'(e71), number(e71, sing), det(e71, 'The'), realisation(e71, offsets(669, 681))]			
{qlf=[for(e72, e73), name(e73, 'Saddam'), person(e73), adj(e73, criticizing), realisation(e73, o			
{qlf=['_'(e74), time(e74, present), aspect(e74, simple), voice(e74, active), lobj(e74, e75), '_'(e7			
{qlf=[name(e77, president), person(e77), det(e77, the), realisation(e77, offsets(719, 732)), '_'			
{qlf=[in(e82, e83), '_'(e83), number(e83, sing), det(e83, the), realisation(e83, offsets(799, 808			
{qlf=['_'(e85), number(e85, sing), adj(e85, personal), det(e85, a), realisation(e85, offsets(816,			

Multi-lingual Noun Phrase Chunker

MuNPEX

MuNPEX is an open-source multi-lingual noun phrase (NP) chunker implemented in JAPE. Currently supported are English, German, French, and Spanish (in beta).

The screenshot shows the GATE 2.2 build 1350 interface. The main window displays a text document with several noun phrases highlighted in yellow. The text is in German and discusses the purpose of plastering walls. The annotations are represented by yellow boxes around the text. Below the text, a table lists the features for each annotation. The 'Edit Annotation' dialog box is open, showing the details for the selected annotation.

Default annotations:

- AbbrNAcro
- FirstPerson
- Lookup
- NP

End	Features
714	{HEAD_END=714, HEAD_START=705, HEAD=Mauerwerk, MOD=schlecht aussehendem}
733	{HEAD_END=733, HEAD_START=727, HEAD=Räumen, MOD=inneren}
756	{HEAD_END=756, HEAD_START=752, HEAD=Putz, DET=den}
786	{HEAD_END=786, HEAD_START=782, HEAD=Wand, MOD=glatte}
805	{HEAD=Deckenflächen, HEAD_START=792, HEAD_END=805}
848	{HEAD_END=848, HEAD_START=835, HEAD=Ausschmückung, MOD=weiteren}
883	{HEAD=Teil, HEAD_START=879, HEAD_END=883}
887	{HEAD=III, HEAD_START=884, HEAD_END=887}
893	{HEAD=Band, HEAD_START=889, HEAD_END=893}

Edit Annotation dialog box:

Annotation type: NP

Feature	Value
HEAD	Mauerwerk
HEAD_END	714
HEAD_START	705
MOD	schlecht aussehendem

New feature name: _____ New feature value: _____

The *Durm* German Lemmatizer

An Open Source Lemmatizer for German

Annotation Sets Annotations Co-reference Editor Text

Der Bundeskanzler wird vom Bundestag in geheimer Wahl ohne
Aussprache gewählt. Zunächst erfolgt ein Vorschlag des
Bundespräsidenten, der hinsichtlich der Person, die er

◀ ▶ ✖ ▶▶
 Mehr stark
 Bundestagswahl
 wird. Dies ist
 siegreichen
 n mit den
 d der Gewählte
 r Kandidat vom
 Bundestag gewählt worden. Wählt der Bundestag den
 Vorgeschlagenen nicht, so hat der Bundestag vierzehn Tage
 Zeit, nach Vorschlägen aus seiner Mitte einen Bundeskanzler
 mit den Stimmen der Mehrheit seiner Mitglieder (absolute
 Mehrheit) zu wählen. Gelingt es dem Bundestag nicht, in
 dieser Zeit eine Person zu wählen, so findet nach Ablauf der
 Frist unverzüglich ein neuer Wahlgang statt, in dem gewählt
 ist, wer die meisten Stimmen erhält. Ist diese Mehrheit
 zugleich eine absolute Mehrheit, so muss der Bundespräsident
 den Gewählten binnen sieben Tagen ernennen. Konnte der
 Gewählte nur eine relative Mehrheit auf sich vereinen, so

Type	Set	Start	End	
Lemma		17183	17196	{Lemma=Bundeskanzler}
Lemma		17206	17215	{Lemma=Bundestag}
Lemma		17228	17232	{Lemma=Wahl}
Lemma		17279	17288	{Lemma=Vorschlag}
Lemma		17293	17310	{Lemma=Bundespräsident}
Lemma		17439	17444	{Lemma=Abend}

- DATE
- DE-Morph
- DEFAULT_TOKEN
- Exception
- LOC
- Lemma
- Lookup
- NMB
- NP
- ORG
- PER
- PRC
- PROF
- Sentence
- SpaceToken
- Split
- TIME
- Token
- tempNP

References

Frameworks

The GATE (*General Architecture for Text Engineering*) System:

- <http://gate.ac.uk>
- <http://sourceforge.net/projects/gate>
- User's Guide: <http://gate.ac.uk/sale/tao/>

IBM's UIMA (*Unstructured Information Management Architecture*):

- <http://www.research.ibm.com/UIMA/>
- <http://sourceforge.net/projects/uima-framework/>

Other Resources

- WordNet: <http://wordnet.princeton.edu/>
- MuNPEX: <http://www.ipd.uka.de/~durm/tm/munpex/>

Part IV

Applications

15 Introduction

- Applications

16 Summarization

- Introduction
- Example System: NewsBlaster
- Document Understanding Conference (DUC)
- Example System: ERSS
- Evaluation
- Summarization: Summary

17 Opinion Mining

18 Question-Answering (QA)

19 Text Mining in Biology and Biomedicine

- Introduction
- The BioRAT System
- Mutation Miner

20 References

- References

Text Mining Applications

Bringing it all together...

We now look at some actual Text Mining applications:

Automatic Summarization: of single and multiple documents

Opinion Mining: extracting *opinions* by consumers regarding companies and their products

Question-Answering: answering factual questions

Text Mining in Biology: the BioRAT and MutationMiner systems

For **Summarization** and **Biology**, we'll look into some systems in detail.

15 Introduction

16 Summarization

- Introduction
- Example System: NewsBlaster
- Document Understanding Conference (DUC)
- Example System: ERSS
- Evaluation
- Summarization: Summary

17 Opinion Mining

18 Question-Answering (QA)

19 Text Mining in Biology and Biomedicine

20 References

Automatic Summarization

Definition

A summary text is a condensed derivative of a source text, reducing content by selection and/or generalisation on what is important.

Note

Distinguish between:

- *abstracting*-based summaries, and
- *extracting*-based summaries.

Automatically created summaries are (almost) exclusively text extracts.

The Challenge

to identify the informative segments at the expense of the rest

The NewsBlaster System (Columbia U.)

The screenshot shows a web browser window titled "Columbia Newsblaster: Summarizing All the News on the Web (03/18/2006 - 03/21/2006) - Konqueror". The address bar shows "http://newsblaster.cs.columbia.edu/". The page header includes the date "Tuesday, March 21, 2006" and "Articles from 03/18/2006 to 03/21/2006". The main content area features a search bar, a navigation menu with categories like "U.S.", "World", "Finance", "Sci/Tech", "Entertainment", and "Sports", and a featured article titled "Three years after Iraq invasion leader says civil war has begun". The article includes a photo of George W. Bush and a list of other stories about Iraq, war, and Bush.

Columbia Newsblaster
Summarizing all the news on the Web

Tuesday, March 21, 2006
Articles from 03/18/2006 to 03/21/2006
Last update: 8:58 AM EST

Search for:

Offline summarization (na) ▾
Go

U.S.
[World](#)
[Finance](#)
[Sci/Tech](#)
[Entertainment](#)
[Sports](#)


[View Today's Images](#)

[View Archive](#)

[About Newsblaster](#)
[About today's run](#)
[Newsblaster in Press](#)
[Academic Papers](#)

Three years after Iraq invasion leader says civil war has begun (World, 21 articles) [UPDATE]

WASHINGTON - US President George W Bush and his senior advisers sought to mark Sunday's third anniversary of the Iraq war with declarations of progress, but found themselves embroiled in renewed debate about whether Iraq has fallen into civil war. Bush gave a blunt defense of the American strategy in Iraq today, while acknowledging that ordinary Iraqis had been left exposed to terrorism during the war's earlier stages. Washington On the third anniversary of a war that they once expected to be over by now, Bush and senior officials contended that their strategy is working despite the escalating sectarian violence in Iraq. Polls have shown American support for the Iraq war dropping since the bombing last month of a Shiite shrine in Samarra led to widespread communal violence. LONDON - Iraq is in a state of civil war and is nearing the point of no return when the country's sectarian violence will spill over throughout the Middle East, Iraqi Prime Minister Ayad Allawi said on Sunday. (CBS/AP) Iraq is in the middle of a civil war, Allawi said in an interview with the British Broadcasting Corp. aired on Sunday.



Other stories about Iraq, war and Bush:

- [Bush marks Iraq date, omits using 'war' word](#) (12 articles)
- [Bush Asks U.S. to Look Past Iraq Bloodshed](#) (9 articles)

A Multi-Document Summary generated by NewsBlaster

The screenshot shows a web browser window titled "Columbia Newsblaster: Three years after Iraq invasion leader says civil war has begun - Conqueror". The address bar shows the URL: <http://newsblaster.cs.columbia.edu/summaries/2006-03-21-08-58-15-021.html>. The page header includes the "Columbia Newsblaster" logo, the date "Tuesday, March 21, 2006", and the text "Articles from 03/18/2006 to 03/21/2006" and "Last update: 8:58 AM EST".

The main content area features a search bar with the text "Search for:" and a dropdown menu for "Offline summarization (na)". Below the search bar is a navigation menu with links for "U.S.", "World", "Finance", "Sci/Tech", "Entertainment", and "Sports". There are also links for "View Today's Images", "View Archive", "About Newsblaster", "About today's run", "Newsblaster in Press", and "Academic Papers".

The article title is "Three years after Iraq invasion leader says civil war has begun". The sub-headline is "Summary from multiple countries, from articles in English" and it is marked as "[UPDATED] (see summary with new information since yesterday)".

The article text begins with: "WASHINGTON - US President George W Bush and his senior advisers sought to mark Sunday's third anniversary of the Iraq war with declarations of progress, but found themselves embroiled in renewed debate about whether Iraq has fallen into civil war. (article 13) Bush gave a blunt defense of the American strategy in Iraq today, while acknowledging that ordinary Iraqis had been left exposed to terrorism during the war's earlier stages. (article 11) Washington On the third anniversary of a war that they once expected to be over by now, Bush and senior officials contended that their strategy is working despite the escalating sectarian violence in Iraq. (article 10) Polls have shown American support for the Iraq war dropping since the bombing last month of a Shiite shrine in Samarra led to widespread communal violence. (article 15) LONDON - Iraq is in a state of civil war and is nearing the point of no return when the country's sectarian violence will spill over throughout the Middle East, Iraqi Prime Minister Ayad Allawi said on Sunday. (article 9) (CBS/AP) Iraq is in the middle of a civil war, Allawi said in an interview with the British Broadcasting Corp. aired on Sunday. (article 8)"

On the right side of the article, there is a vertical stack of three images: the top image shows a man in a suit (likely Ayad Allawi), the middle image shows a group of people in a public setting, and the bottom image shows soldiers in military gear.

NewsBlaster: Article Classification

Columbia Newsblaster: Summarizing All the News on the Web (03/18/2006 - 03/21/2006) - Konqueror

Location Edit View Go Bookmarks Tools Settings Window Help

Location: <http://newsblaster.cs.columbia.edu/>

U.S.

- [Labour secret loans to be debated](#) (13 articles)
- [Agent: FBI bosses hindered Moussaoui probe](#) (10 articles)
- [Dark Side Of The Mesa](#) (8 articles)
- [The Seattle Times: Chicago's got the headquarters, but Seattle's still Jet City, USA](#) (8 articles)

[See all 26 U.S. stories...](#)

World

- [French unions call for general strike over job law](#) (13 articles) [UPDATE]
- [Australia begins cyclone clean-up](#) (11 articles) [UPDATE]
- [Russia said to still object to UN Iran statement](#) (9 articles) [UPDATE]
- [Putin visits China to boost ties](#) (6 articles)

[See all 17 World stories...](#)

Finance

- [FCC near deciding Verizon's broadband request](#) (6 articles)
- [The Seattle Times: Microsoft's new Vista should reach store shelves by holidays](#) (5 articles)
- [Dell to double its staff in India to 20,000 by 2009](#) (5 articles)
- [Aviva courting Prudential; rules out hostile bid](#) (5 articles)

[See all 8 Finance stories...](#)

Science/Technology

- [Breast Asymmetry Linked to Increased Cancer Risk](#) (7 articles)

Entertainment

- [Fashion designer Oleg Cassini dies at 92](#) (5 articles)

Sports

- [Bonds probe? Selig won't say](#) (13 articles)
- [Live: Commonwealth Games](#) (12 articles)
- [Connecticut, Villanova Survive Comebacks](#) (10 articles)
- [Big Dance man Wright stuffs Pitt's Krauser](#) (9 articles)

NewsBlaster: Tracking Events over Time

Columbia Newsblaster
Summarizing all the news on the Web

Tuesday, March 21, 2006
Articles from 03/18/2006 to 03/21/2006
Last update: 8:58 AM EST

Search for:

Offline summarization (ma) ▾
Go

Three years after Iraq invasion leader says civil war has begun
Summary from multiple countries, from articles in English
[UPDATED] (see summary with new information since yesterday)

Tracking this event across days (click [here](#) to return)

3/19/2006	3/20/2006	3/21/2006
Journalist's Alleged Killers Held in Iraq [recenter]	///////	
As Iraq War Heads Into 4th Year, Bush Pledges 'Complete Victory' [recenter]	----- /////////	Top Iraqi Leaders Agree to Form a Policy Council [recenter]
Global anti-war protesters rally to mark 3 years in Iraq [recenter]	-----	Three years after Iraq invasion leader says civil war has begun [recenter]
	Three years after Iraq invasion leader says civil war has begun [recenter]	///////

Research in Automatic Summarization

The Challenge

- Various summarization systems produce different kinds of summaries, from different data, for different purposes, using different evaluations
- Impossible to measure (scientific) progress

Document Understanding Conference (DUC)

The solution: hold a *competition*

- Started in 2001
- Organized by U.S. National Institute of Standardization and Technology (NIST)
- Forum to compare summarization systems
- For all systems the same tasks, data, and evaluation methods

Research in Automatic Summarization

The Challenge

- Various summarization systems produce different kinds of summaries, from different data, for different purposes, using different evaluations
- Impossible to measure (scientific) progress

Document Understanding Conference (DUC)

The solution: hold a *competition*

- Started in 2001
- Organized by U.S. National Institute of Standardization and Technology (NIST)
- Forum to compare summarization systems
- For all systems the same tasks, data, and evaluation methods

Document Understanding Conference (DUC)

Data

- newspaper and newswire articles (AP, NYT, XIE, ...)
- topical clusters of various length (2004: 10, 2005: 25–50, 2006: 25)

Tasks

In 2004:

- short summaries of single articles (10 words)
- summaries of single articles (100 words)
- multi-document summaries of a 10-document cluster
- cross-language summaries (machine translated Arabic)
- summaries focused by a question “Who is X?”

In 2005–2006:

- Focused multi-document summaries for a given *context*

Summarization System ERSS (CLaC/IPD)

Main processing steps

Preprocessing Tokenizer, Sentence Splitter, POS Tagger, ...

MuNPEx noun phrase chunker (JAPE-based)

FCR fuzzy coreference resolution algorithm

Classy naive Bayesian classifier for multi-dimensional text categorization

Summarizer summarization framework with individual strategies

Implementation based on the GATE architecture.

Summarization System ERSS (CLaC/IPD)

Main processing steps

Preprocessing Tokenizer, Sentence Splitter, POS Tagger, ...

MuNPE_x noun phrase chunker (JAPE-based)

FCR fuzzy coreference resolution algorithm

Classy naive Bayesian classifier for multi-dimensional text categorization

Summarizer summarization framework with individual strategies

Implementation based on the GATE architecture.

Summarization System ERSS (CLaC/IPD)

Main processing steps

Preprocessing Tokenizer, Sentence Splitter, POS Tagger, ...

MuNPEx noun phrase chunker (JAPE-based)

FCR fuzzy coreference resolution algorithm

Classy naive Bayesian classifier for multi-dimensional text categorization

Summarizer summarization framework with individual strategies

Implementation based on the GATE architecture.

Summarization System ERSS (CLaC/IPD)

Main processing steps

Preprocessing Tokenizer, Sentence Splitter, POS Tagger, ...

MuNPEx noun phrase chunker (JAPE-based)

FCR fuzzy coreference resolution algorithm

Classy naive Bayesian classifier for multi-dimensional text categorization

Summarizer summarization framework with individual strategies

Implementation based on the GATE architecture.

Summarization System ERSS (CLaC/IPD)

Main processing steps

Preprocessing Tokenizer, Sentence Splitter, POS Tagger, ...

MuNPEx noun phrase chunker (JAPE-based)

FCR fuzzy coreference resolution algorithm

Classy naive Bayesian classifier for multi-dimensional text categorization

Summarizer summarization framework with individual strategies

Implementation based on the GATE architecture.

Summarization System ERSS (CLaC/IPD)

Main processing steps

Preprocessing Tokenizer, Sentence Splitter, POS Tagger, ...

MuNPEx noun phrase chunker (JAPE-based)

FCR fuzzy coreference resolution algorithm

Classy naive Bayesian classifier for multi-dimensional text categorization

Summarizer summarization framework with individual strategies

Implementation based on the GATE architecture.

ERSS: Preprocessing Steps

Basic Preprocessing

Tokenization, Sentence Splitting, POS Tagging, ...

Number Interpreter

Locates number expressions and assigns numerical values, e.g.,
"two" → 2.

Abbreviation & Acronym Detector

Scans tokens for acronyms ("GM", "IBM", ...) and abbreviations (e.g., "e.g.", "Fig.", ...) and adds the full text.

Gazetteer

Scans input tokens and adds *type* information based on a number of word lists: *city, company, currency, festival, mountain, person_female, planet, region, street, timezone, title, water, ...*

ERSS: Preprocessing Steps

Basic Preprocessing

Tokenization, Sentence Splitting, POS Tagging, ...

Number Interpreter

Locates number expressions and assigns numerical values, e.g.,
“two” → 2.

Abbreviation & Acronym Detector

Scans tokens for acronyms (“GM”, “IBM”, ...) and abbreviations (e.g., “e.g.”, “Fig.”, ...) and adds the full text.

Gazetteer

Scans input tokens and adds *type* information based on a number of word lists: *city, company, currency, festival, mountain, person_female, planet, region, street, timezone, title, water, ...*

ERSS: Preprocessing Steps

Basic Preprocessing

Tokenization, Sentence Splitting, POS Tagging, ...

Number Interpreter

Locates number expressions and assigns numerical values, e.g.,
“two” → 2.

Abbreviation & Acronym Detector

Scans tokens for acronyms (“GM”, “IBM”, ...) and abbreviations (e.g., “e.g.”, “Fig.”, ...) and adds the full text.

Gazetteer

Scans input tokens and adds *type* information based on a number of word lists: *city, company, currency, festival, mountain, person_female, planet, region, street, timezone, title, water, ...*

ERSS: Preprocessing Steps

Basic Preprocessing

Tokenization, Sentence Splitting, POS Tagging, ...

Number Interpreter

Locates number expressions and assigns numerical values, e.g.,
“two” → 2.

Abbreviation & Acronym Detector

Scans tokens for acronyms (“GM”, “IBM”, ...) and abbreviations (e.g., “e.g.”, “Fig.”, ...) and adds the full text.

Gazetteer

Scans input tokens and adds *type* information based on a number of word lists: *city, company, currency, festival, mountain, person_female, planet, region, street, timezone, title, water, ...*

Preprocessing Steps (II)

Named Entity (NE) Recognition

Scans a sequence of (annotated) tokens with JAPE grammars and adds NE information: *Date*, *Person*, *Organization*, ...

Example: Tokens "10", "o", "'", "clock" → *Date::TimeOClock*

JAPE Grammars

- Regular-expression based grammars
- used to generate *finite state Transducers* (non-deterministic finite state machines)

Example Grammar

```
Rule: TimeOClock // ten o'clock
({Lookup.minorType == hour}
 {Token.string == "o"}
 {Token.string == "'"}
 {Token.string == "clock"}
):time
-->
:time.TempTime = {kind = "positive",
                  rule = "TimeOClock"}
```

Preprocessing Steps (II)

Named Entity (NE) Recognition

Scans a sequence of (annotated) tokens with JAPE grammars and adds NE information: *Date*, *Person*, *Organization*, ...

Example: Tokens "10", "o", "'", "clock" → *Date::TimeOClock*

JAPE Grammars

- Regular-expression based grammars
- used to generate *finite state Transducers* (non-deterministic finite state machines)

Example Grammar

```
Rule: TimeOClock // ten o'clock
({Lookup.minorType == hour}
 {Token.string == "o"}
 {Token.string == "'"}
 {Token.string == "clock"}
):time
-->
:time.TempTime = {kind = "positive",
                  rule = "TimeOClock"}
```


Preprocessing Steps (II)

Named Entity (NE) Recognition

Scans a sequence of (annotated) tokens with JAPE grammars and adds NE information: *Date*, *Person*, *Organization*, ...

Example: Tokens "10", "o", "'", "clock" → *Date::TimeOClock*

JAPE Grammars

- Regular-expression based grammars
- used to generate *finite state Transducers* (non-deterministic finite state machines)

Example Grammar

```
Rule: TimeOClock // ten o'clock
({Lookup.minorType == hour}
 {Token.string == "o"}
 {Token.string == "'"}
 {Token.string == "clock"}
):time
-->
:time.TempTime = {kind = "positive",
                  rule = "TimeOClock"}
```

Fuzzy Coreference Resolution

Coreference Resolution

Input to a coreference resolution algorithm is a set of noun phrases (NPs). Example: *Mr. Bush* $\overset{?}{\longleftrightarrow}$ *the president* $\overset{?}{\longleftrightarrow}$ *he*

Fuzzy Representation of Coreference

Core idea: coreference between noun phrases is almost never “100% certain”

- fuzzy model: represent certainty of coreference *explicitly* with a membership degree
- formally: represent fuzzy chain \mathcal{C} with a fuzzy set $\mu_{\mathcal{C}}$, mapping the domain of all NPs in a text to the $[0,1]$ -interval
- then, each noun phrase np_i has a corresponding membership degree $\mu_{\mathcal{C}}(np_i)$, indicating how certain this NP is a member of chain \mathcal{C}

Fuzzy Coreference Resolution

Coreference Resolution

Input to a coreference resolution algorithm is a set of noun phrases (NPs). Example: *Mr. Bush* $\overset{?}{\longleftrightarrow}$ *the president* $\overset{?}{\longleftrightarrow}$ *he*

Fuzzy Representation of Coreference

Core idea: coreference between noun phrases is almost never “100% certain”

- fuzzy model: represent certainty of coreference *explicitly* with a membership degree
- formally: represent fuzzy chain \mathcal{C} with a fuzzy set $\mu_{\mathcal{C}}$, mapping the domain of all NPs in a text to the $[0,1]$ -interval
- then, each noun phrase np_i has a corresponding membership degree $\mu_{\mathcal{C}}(np_i)$, indicating how certain this NP is a member of chain \mathcal{C}

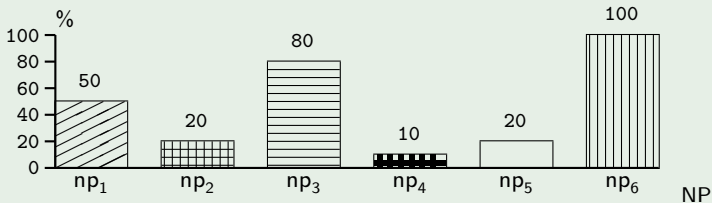
Fuzzy Coreference Resolution

Fuzzy Coreference Chain

Fuzzy set $\mu_C : NP \rightarrow [0, 1]$

Example

Fuzzy Coreference Chain \mathcal{C}



Fuzzy Coreference Chains

Properties of fuzzy chains

- each chain holds *all* noun phrases in a text
- i.e., each NP is a member of every chain (but with very different certainties)
- we don't have to reject inconsistencies right away — they can be reconciled later through suitable fuzzy operators
- also, there is no arbitrary boundary for discriminating between “corefering” and “not corefering”
- thus, in this step we don't lose information we might need later

Fuzzy Clustering

How can we *build* fuzzy chains?

- Use knowledge-poor heuristics to check for coreference between NP pairs
- Examples: Substring, Synonym/Hypernym, Pronoun, CommonHead, Acronym...
- Fuzzy heuristic: return a *degree* of coreference $\in [0, 1]$

Creating Chains by Clustering

Idea: initially, each NP represents one chain (where it is its *medoid*). Then:

- apply a single-link hierarchical clustering strategy,
- using the fuzzy degree as an (inverse) distance measure

This results in NP clusters, which can be converted into coreference chains.

Fuzzy Clustering

How can we *build* fuzzy chains?

- Use knowledge-poor heuristics to check for coreference between NP pairs
- Examples: Substring, Synonym/Hypernym, Pronoun, CommonHead, Acronym. . .
- Fuzzy heuristic: return a *degree* of coreference $\in [0, 1]$

Creating Chains by Clustering

Idea: initially, each NP represents one chain (where it is its *medoid*). Then:

- apply a single-link hierarchical clustering strategy,
- using the fuzzy degree as an (inverse) distance measure

This results in NP clusters, which can be converted into coreference chains.

Designing Fuzzy Heuristics

Fuzzy Heuristics

How can we compute a coreference degree $\mu_{(np_j, np_k)}^{\mathcal{H}_i}$?

Fuzzy Substring Heuristic: (character n-gram match) return coreference degree of 1.0 if two NP string are identical, 0.0 if they share no substring. Otherwise, select longest matching substring and set coreference degree to its percentage of first NP.

Fuzzy Synonym/Hypernym Heuristic: Synonyms (determined through *WordNet*) receive a coreference degree of 1.0. If two NPs are hypernyms, set the coreference degree depending on distance in the hierarchy (i.e., longer paths result in lower certainty degrees).

Designing Fuzzy Heuristics

Fuzzy Heuristics

How can we compute a coreference degree $\mu_{(np_j, np_k)}^{\mathcal{H}_i}$?

Fuzzy Substring Heuristic: (character n-gram match) return coreference degree of 1.0 if two NP string are identical, 0.0 if they share no substring. Otherwise, select longest matching substring and set coreference degree to its percentage of first NP.

Fuzzy Synonym/Hypernym Heuristic: Synonyms (determined through *WordNet*) receive a coreference degree of 1.0. If two NPs are hypernyms, set the coreference degree depending on distance in the hierarchy (i.e., longer paths result in lower certainty degrees).

Designing Fuzzy Heuristics

Fuzzy Heuristics

How can we compute a coreference degree $\mu_{(np_j, np_k)}^{\mathcal{H}_i}$?

Fuzzy Substring Heuristic: (character n-gram match) return coreference degree of 1.0 if two NP string are identical, 0.0 if they share no substring. Otherwise, select longest matching substring and set coreference degree to its percentage of first NP.

Fuzzy Synonym/Hypernym Heuristic: Synonyms (determined through *WordNet*) receive a coreference degree of 1.0. If two NPs are hypernyms, set the coreference degree depending on distance in the hierarchy (i.e., longer paths result in lower certainty degrees).

toward the Dominican Republic Sunday, and the Civil Defense alerted its heavily populated south coast to prepare for high winds, heavy rains and high seas.

The storm was approaching from the southeast with sustained winds of 75 mph gusting to 92 mph.

“There is no need for alarm,” Civil Defense Director Eugenio Cabral said in a television alert shortly before midnight Saturday.

Cabral said residents of the province of Barahona should closely follow Gilbert’s movement. An estimated 100,000 people live in the province, including 70,000 in the city of Barahona, about 125 miles west of Santo Domingo.

Tropical Storm Gilbert formed in the eastern Caribbean and strengthened into a hurricane Saturday night. The National Hurricane Center in Miami reported its position at 2 a.m. Sunday at latitude 16.1 north, longitude 67.5 west, about 140 miles south of Ponce, Puerto Rico, and 200 miles southeast of Santo Domingo.

The National Weather Service in San Juan, Puerto Rico, said Gilbert was moving westward at 15 mph with a “broad area of cloudiness and heavy weather” rotating around the center of the storm.

The weather service issued a flash flood watch for Puerto Rico and the Virgin Islands until at least 6 p.m. Sunday.

Strong winds associated with the Gilbert brought coastal flooding, strong southeast winds and up to 12 feet of rain to Puerto Rico’s south coast. There were no reports of casualties.

San Juan, on the north coast, had heavy rains and gusts Saturday, but they subsided during the night.

On Saturday, Hurricane Florence was downgraded to a tropical

- Lookup
- Measurement
- Number
- Organization
- Person
- Sentence
- SpaceToken
- Split
- Title
- Token

NP Coreference annotations

- (02) Barahona
- (02) Ponce, Puerto Rico
- (02) Santo Domingo
- (02) residents
- (02) the province
- (03) Saturday
- (03) heavy rains
- (03) south coast
- (04) Hurricane Gilbert
- (04) The storm
- (04) winds

Original markups annotations

- BYLINE
- DATELINE
- DOC

Summarizer

ERSS (Experimental Resolution System Summarizer)

A Summary should contain the **most important entities** within a text. Assumption: these are also mentioned more often, hence result in longer **coreference chains**.

Summarization Algorithm (Single Documents)

- 1 Rank coreference chains by size (and other features)
- 2 For each chain: select highest-ranking NP/Sentence
- 3 extract NP (short summary)
or complete sentence (long summary)
- 4 continue with next-longest chain until length limit has been reached

Summarizer

ERSS (Experimental Resolution System Summarizer)

A Summary should contain the **most important entities** within a text. Assumption: these are also mentioned more often, hence result in longer **coreference chains**.

Summarization Algorithm (Single Documents)

- 1 Rank coreference chains by size (and other features)
- 2 For each chain: select highest-ranking NP/Sentence
- 3 extract NP (short summary)
or complete sentence (long summary)
- 4 continue with next-longest chain until length limit has been reached

ERSS: Keyword-style Summary Examples

Automatically created 10-word-summaries

Can you guess the text's topic?

Space News: [the shuttle Discovery's Hubble repair mission,
the observatory's central computer]

People & Politics: [Lewinsky, President Bill Clinton, her testimony,
the White House scandal]

Business & Economics: [PAL, the company's stock, a management-
proposed recovery plan, the laid-off workers]

(from DUC 2003)

ERSS: Single-Document Summary Example

Automatically created 100-word summary (from DUC 2004)

President Yoweri Museveni insists they will remain there until Ugandan security is guaranteed, despite Congolese President Laurent Kabila's protests that Uganda is backing Congolese rebels attempting to topple him. After a day of fighting, Congolese rebels said Sunday they had entered Kindu, the strategic town and airbase in eastern Congo used by the government to halt their advances. The rebels accuse Kabila of betraying the eight-month rebellion that brought him to power in May 1997 through mismanagement and creating divisions among Congo's 400 tribes. A day after shooting down a jetliner carrying 40 people, rebels clashed with government troops near a strategic airstrip in eastern Congo on Sunday.

Summarizer (II): more complicated summaries

Multi-Document Summaries

Many tasks in DUC require summaries of *multiple* documents:

- cross-document summary
- focused summary
- context-based summary (DUC 2005, 2006)

Solution

Additionally build *cross-document coreference chains* and summarize using a *fuzzy cluster graph algorithm*.

- For focused and context-based summaries, only use those chains that connect the question(s) with the documents (even if they have a lower rank)

Summarizer (II): more complicated summaries

Multi-Document Summaries

Many tasks in DUC require summaries of *multiple* documents:

- cross-document summary
- focused summary
- context-based summary (DUC 2005, 2006)

Solution

Additionally build *cross-document coreference chains* and summarize using a *fuzzy cluster graph algorithm*.

- For focused and context-based summaries, only use those chains that connect the question(s) with the documents (even if they have a lower rank)

Example for a Focused Summary generated by ERSS

"Who is Stephen Hawking?"

Hawking, 56, is the Lucasian Professor of Mathematics at Cambridge, a post once held by Sir Isaac Newton. Hawking, 56, suffers from Lou Gehrig's Disease, which affects his motor skills, and speaks by touching a computer screen that translates his words through an electronic synthesizers. Stephen Hawking, the Cambridge University physicist, is renowned for his brains. Hawking, a professor of physics and mathematics at Cambridge University in England, has gained immense celebrity, written a best-selling book, fathered three children, and done a huge amount for the public image of disability. Hawking, Mr. Big Bang Theory, has devoted his life to solving the mystery of how the universe started and where it's headed.

Example for a context-based summary (Excerpt)

Question

What countries are or have been involved in land or water boundary disputes with each other over oil resources or exploration? How have disputes been resolved, or towards what kind of resolution are the countries moving? What other factors affect the disputes?

System summary (first ~70 words of 250 total)

The ministers of Asean - grouping Brunei, Indonesia, Malaysia, the Philippines, Singapore and Thailand - raised the Spratlys issue at a meeting yesterday with Qian Qichen, their Chinese counterpart. The meeting takes place against a backdrop of the continuing territorial disputes involving three Asean members - China, Vietnam and Taiwan - over the Spratley Islands in the South China Sea, a quarrel which could deteriorate shortly with the expected start of oil exploration in the area...

Example for a context-based summary (Excerpt)

Question

What **countries** are or have been involved in land or water boundary disputes with each other over **oil resources or exploration**? How have disputes been resolved, or towards what kind of resolution are the countries moving? What other factors affect the disputes?

System summary (first ~70 words of 250 total)

The ministers of Asean - grouping **Brunei, Indonesia, Malaysia, the Philippines, Singapore and Thailand** - raised the Spratlys issue at a meeting yesterday with Qian Qichen, their Chinese counterpart. The meeting takes place against a backdrop of the continuing territorial disputes involving three Asean members - **China, Vietnam and Taiwan** - over the Spratley Islands in the South China Sea, a quarrel which could deteriorate shortly with the expected start of **oil exploration** in the area...

How can we evaluate summaries?

Problem

A summary is not **right** or **wrong**. Hard to find criterias.

Intrinsic

- Compare with model summaries
- Compare with source text
- Look solely at summary

Extrinsic

- Regarding external task
- Example: Use summary to cook a meal

Manual

- Subjective view
- High costs (40 systems X 50 clusters X 2 assessors = 4000 summaries)

Automatic

- High availability (during development)
- Repeatable and fast

How can we evaluate summaries?

Problem

A summary is not **right** or **wrong**. Hard to find criterias.

Intrinsic

- Compare with model summaries
- Compare with source text
- Look solely at summary

Extrinsic

- Regarding external task
- Example: Use summary to cook a meal

Manual

- Subjective view
- High costs (40 systems X 50 clusters X 2 assessors = 4000 summaries)

Automatic

- High availability (during development)
- Repeatable and fast

How can we evaluate summaries?

Problem

A summary is not **right** or **wrong**. Hard to find criterias.

Intrinsic

- Compare with model summaries
- Compare with source text
- Look solely at summary

Extrinsic

- Regarding external task
- Example: Use summary to cook a meal

Manual

- Subjective view
- High costs (40 systems X 50 clusters X 2 assessors = 4000 summaries)

Automatic

- High availability (during development)
- Repeatable and fast

Manual Measures

Summary Evaluation Environment: Linguistic quality

- Grammaticality
- Non-redundancy
- Referential clarity
- Focus
- Structure & Coherence

Responsiveness (2005)

- Pseudo-extrinsic
- How well was the question answered
- Form & Content
- In relation to the other systems' summaries

Manual Measures

Summary Evaluation Environment: Linguistic quality

- Grammaticality
- Non-redundancy
- Referential clarity
- Focus
- Structure & Coherence

Responsiveness (2005)

- Pseudo-extrinsic
- How well was the question answered
- Form & Content
- In relation to the other systems' summaries

Manual Measures: SEE – Quality evaluation

The screenshot shows the SEE software interface. At the top, there is a menu bar with 'File', 'Options', and 'Help'. Below the menu bar, there are two input fields for file paths: 'Peer Summary Path' and 'Model Summary Path', both pointing to 'E:\Dokumente und Einstellungen\axis58\Desktop\D132.M.100.'. To the right of each field is a button labeled 'Open Peer Summary' and 'Open Model Summary' respectively. The main area is divided into two panes: 'Peer Summary' on the left and 'Model Summary' on the right. The 'Peer Summary' pane contains a text snippet with two numbered references: [1] Treasury Secretary Robert Rubin arrived in Malaysia. Sunday for a two-day visit to discuss the regional economic situation, the U.S. Embassy said. [2] So it comes as something of a surprise that Rubin, now treasury secretary, may have missed out on what would probably have created the biggest windfall of his life: Robert Rubin resigns as... The 'Model Summary' pane is currently empty. Below the panes is a tabbed interface with tabs for 'Quality Judgment 1', 'Quality Judgment 2', 'Content', 'Unmarked Peer Units', 'Auto Evaluation', and 'Results'. The 'Quality Judgment 2' tab is active, displaying two questions: Q2. If you were editing the summary to make it more concise and to the point, how much useless, confusing repetitive text would you remove from the existing summary? and Q3. To what degree does the summary say the same thing over again? Each question has five radio button options: None, A little, Some, A lot, and Most of the text. At the bottom of the window, there is a status bar showing '0 of 12 quality questions judged' and '0/6 model units judged'.

SEE - *Untitled

File Options Help

Peer Summary Path E:\Dokumente und Einstellungen\axis58\Desktop\D132.M.100. Open Peer Summary

Model Summary Path E:\Dokumente und Einstellungen\axis58\Desktop\D132.M.100. Open Model Summary

Peer Summary

[1] [Treasury Secretary Robert Rubin arrived in Malaysia. Sunday for a two-day visit to discuss the regional economic situation, the U.S. Embassy said.](#) [2] [So it comes as something of a surprise that Rubin, now treasury secretary, may have missed out on what would probably have created the biggest windfall of his life: Robert Rubin resigns as...](#)

Model Summary

Quality Judgment 1 Quality Judgment 2 Content Unmarked Peer Units Auto Evaluation Results

Q2. If you were editing the summary to make it more concise and to the point, how much useless, confusing repetitive text would you remove from the existing summary?

- None
- A little
- Some
- A lot
- Most of the text

Q3. To what degree does the summary say the same thing over again?

0 of 12 quality questions judged 0/6 model units judged

Automatic Measures: ROUGE

ROUGE (Recall-Oriented Understudy for Gisting Evaluation)

measures n-gram overlap between a peer and a set of reference summaries.

Definition

$$ROUGE_n = \frac{\sum_{C \in ModelUnits} \sum_{n\text{-gram} \in C} Count_{match}(n\text{-gram})}{\sum_{C \in ModelUnits} \sum_{n\text{-gram} \in C} Count(n\text{-gram})}$$

$ROUGE_{SU4} = ROUGE_2$ with skip of max. 4 words between two 2-grams

$ROUGE_2 / ROUGE_{SU4}$

S1 police killed the gunman

S2 police stopped the gunman

Evaluation: ERSS Results

DUC 2004

26 systems from 25 different groups, both industry and academic.
Evaluation performed by NIST (see <http://duc.nist.gov>).

ROUGE Results

Task 2: Cross-Document Common Topic Summaries

- *Best: 0.38, Worst: 0.24, Average: 0.34, **ERSS: 0.36***
- ERSS statistically indistinguishable from top system within a 0.05 confidence level

Task 5: Focused Summaries

- *Best: 0.35, Worst: 0.26, Average: 0.31, **ERSS: 0.33***
- same as above

Similar results for all other tasks.

Automatic Measures: Pyramids & Basic Elements

Driving force

- Scores of systems are not distinguishable.
- Only exact matches count. Abstractions are ignored.

Pyramids

- Comparing content units (not n-grams) of peer and models.
- Chunks occurring in more models get higher points.
- Needs manual annotation of peers **and** models.

Basic Elements

- Peer and Model summaries are parsed, extracting general relations between words of a sentence.
- Compute overlap of extracted **Head-Modifier-Relation-Triples** between peer and models.

⇒ Peers don't have to be annotated by hand!

Automatic Measures: Pyramids & Basic Elements

Driving force

- Scores of systems are not distinguishable.
- Only exact matches count. Abstractions are ignored.

Pyramids

- Comparing content units (not n-grams) of peer and models.
- Chunks occurring in more models get higher points.
- Needs manual annotation of peers **and** models.

Basic Elements

- Peer and Model summaries are parsed, extracting general relations between words of a sentence.
- Compute overlap of extracted **Head-Modifier-Relation-Triples** between peer and models.

⇒ Peers don't have to be annotated by hand!

Automatic Measures: Pyramids & Basic Elements

Driving force

- Scores of systems are not distinguishable.
- Only exact matches count. Abstractions are ignored.

Pyramids

- Comparing content units (not n-grams) of peer and models.
- Chunks occurring in more models get higher points.
- Needs manual annotation of peers **and** models.

Basic Elements

- Peer and Model summaries are parsed, extracting general relations between words of a sentence.
- Compute overlap of extracted **Head-Modifier-Relation**-Triples between peer and models.

⇒ Peers don't have to be annotated by hand!

Automatic Measures: Pyramids – GUI

DucView v. 1.2 - Annotating Peer

File Edit Options Help

the release of carcinogenic dioxins and toxic ashes. The two most commonly occurring hazards in confined spaces, according to Moran, are oxygen deficiency and potentially explosive conditions created by the presence of methane gas. After the gold is extracted, residual cyanide would be chemically neutralized. In 1976, cries lamenting damage to the desert reached a peak, and Congress – sensitive to pressure from the politically potent environmental community – passed the Federal Land Management Policy Act. They feared exposure to rodent poison containing arsenic and cyanide, which is easily absorbed through the skin, and swimming pool supplies containing chlorine, whose fumes can cause lung damage, said Pat Askren, chief of Fillmore’s volunteer fire department.

Sodium cyanide spilled into Little Fork Creek, about 2,000 feet from the reservoir.

The reservoir held 13 million to 14 million gallons of sodium cyanide, Berry said.

In addition, some scientists worry about damage to the liver from having to break down

the waste from this process is discharged into self-contained ponds. Environmental problems arise during times of heavy rainfall when ponds overflow and run into natural streams. Environmentalists say the runoff could pollute drinking water and endanger salmon fisheries. In addition, birds and other wildlife are endangered when they see the open blue cyanide ponds and stop to drink from them. Another industrial use of cyanide, in the form of hydrogen cyanide gas, is in the plating industry. There is a reported incident of five workers in Indiana dying of asphyxiation when working in an enclosed space in the presence of the cyanide gas. A plating company in Hollywood, California was charged with dumping cyanide into the sewer system and with reckless storage of chemicals. Another plating company in Burbank, California was closed by the EPA for reckless storage of chemicals including hydrogen cyanide. The Japanese use hydrogen cyanide to manufacture tryptophan, an amino acid used as a nutritional supplement. An unusual use of cyanide is to assist Cameroon villagers to gather honey from hives in tall trees. Climbers stun the bees with smoking leaves and a cyanide compound.

D366.M.250.IE

Gold mining operations make extensive use of cyanide. A process known as heap-leach or heap mining is used to extract gold from low-grade ore. The ore is spread on impermeable plastic pads then sprinkled with a weak cyanide-water mixture. Unfortunately, in some cases runoff from this cyanide process has discharged into ponds with the result that birds and animals were killed.

Measures:

- Leaks contaminate waterways and groundwater
 - There were more spills of water laced with cyanide and heavy metals
 - Sodium cyanide spilled into Little Fork Creek, about 2,000 feet from the reservoir
- Cyanide is used in mining to recover metal from low-yielding ore
 - the cyanide solution used to leach the ore
- Cyanide is used in the metal plating industry
- Sodium cyanide use by fisherman decimates fish
- and a weak cyanide solution is poured over it to pull the gold from the rock
 - After the gold is extracted, residual cyanide
- animals died by thousands from drinking at cyanide-laced holding ponds
- mining industry uses method known as heap leaching
- Philippine fishermen use cyanide in fishing
- The death of hundreds of birds has been caused by drinking from holding ponds
- waste reservoir mists affected people’s breathing
- A film recovery worker died from inhaling cyanide fumes
- and when a plating company dumped cyanide-laced waste water into the L
- Another use of cyanide is in manufacturing the nutritional supplement tryptophan
- Cameroon honey gatherers use cyanide to stun bees
- Cyanide fumes killed five workers cleaning one tank
- Cyanide is used extensively in gold mining
- Cyanide is used to extract silver from used x-ray films
- Leak onto Alamosa river corroded irrigation equipment in the valley
- leaks from waste ponds poison the fish
- Mine waste-water poses an environmental hazard
 - Local groups have expressed fears of environmental damage from a possible spill
- Ore is placed on a plastic pad
- The waste from this process is discharged into self-contained ponds
- cyanide is used in gold, silver, and copper mining
- Cyanide-leaching recovers silver from photographic film;
- heap leach is a cost effective method
- In Summitville Colorado, pollution from a mountain mine killed fish for 17 years
- leaks from waste ponds are contaminating water
- A broken dam spilled cyanide-contaminated water into a creek
- a mine leaked contaminated water into the Alamosa River
- a number of dangerous situations may occur
- A pioneering “closed-loop” cyanide leaching system eliminates open ponds
- animals are killed by leak of cyanide into water bodies
- Cameroon villagers use cyanide to gather honey from hives in tall trees
- Cyanide fishing is linked to the destruction of area reefs
- Cyanide is used in electroplating aircraft parts
- Cyanide is used in pesticides
- Cyanide is used in rodent poison
 - rodent poison containing arsenic and cyanide

Automatic Measures: Basic Elements

“Law enforcement officers from nine African countries are meeting in Nairobi this week to create a regional task force to fight international crime syndicates dealing in ivory, rhino horn, diamonds, arms, and drugs.”

officers—enforcement—nn		syndicates—intern.—mod
officers—countries—from		meeting—officers—subj
nairobi—create—rel	create—week—subj	force—regional—mod
diamonds—arms—conj	countries—nine—nn	force—fight—rel
fight—force—subj	create—force—obj	syndicates—crime—nn
fight—syndicates—obj	horn—rhino—nn	ivory—horn—conj
horn—diamonds—conj	force—task—nn	arms—and—punc
arms—drugs—conj	meeting—nairobi—in	countries—african—nn

Basic Elements (head—modifier—relation) of the sentence shown on top

Summarization: Summary

Some Conclusions...

- Systems score very close to each other, partly due to the automatic ROUGE measure
- Automatic summaries still have a long way to go regarding style, coherence, and capabilities for abstraction
- Evaluation (almost) as difficult as the actual task

The Future?

Still, context-based summarization is promising:

- Do you **really** want to spent hours with Google? Scenario:
 - When writing a report/paper/memo on a certain topic,
 - a system will permanently scan your context,
 - retrieve documents pertaining to your topic,
 - and propose (hopefully relevant) information by itself
- Prediction: This will eventually find its way into Email clients, Word processors, Web browsers, etc.

[cf. Witte 2004 (IIWeb), Witte et al. 2005 (Semantic Desktop)]

Summarization: Summary

Some Conclusions...

- Systems score very close to each other, partly due to the automatic ROUGE measure
- Automatic summaries still have a long way to go regarding style, coherence, and capabilities for abstraction
- Evaluation (almost) as difficult as the actual task

The Future?

Still, context-based summarization is promising:

- Do you **really** want to spent hours with Google? Scenario:
 - When writing a report/paper/memo on a certain topic,
 - a system will permanently scan your context,
 - retrieve documents pertaining to your topic,
 - and propose (hopefully relevant) information by itself
- Prediction: This will eventually find its way into Email clients, Word processors, Web browsers, etc.

[cf. Witte 2004 (IIWeb), Witte et al. 2005 (Semantic Desktop)]

Opinion Mining

Motivation

Nowadays, there are countless websites containing huge amounts of product reviews written by consumers:

- E.g., Amazon.com, Epinions.com

But, like always, now there's too much information:

- You do not really want to spend more time on reading the reviews for a book than the book itself
- For a company, it is difficult to track all opinions regarding its product published on websites

Solution: use Text Mining to process and summarize opinions.

Opinion Mining: General Approach

Processing Steps

Detect Product Features: discussed in the review

Detect Opinions: regarding these features

Determine Polarity: of these opinions (positive? negative?)

Rank opinions: based on their strength (compare “so-so” vs. “desaster”)

[cf. Popescu & Etzioni, HLT/EMNLP 2005]

Solution?

- Use *NE Detection* and *NP Chunking* to identify features
- Find opinions either within the NPs “*a very high resolution*”, or within adjacent constituents using parsing
- Match opinions (using stemming or lemmatization) against a lexicon containing polarity information
- Sort and rank opinions based on the number of reviews and strength

Opinion Mining: General Approach

Processing Steps

Detect Product Features: discussed in the review

Detect Opinions: regarding these features

Determine Polarity: of these opinions (positive? negative?)

Rank opinions: based on their strength (compare “so-so” vs. “desaster”)

[cf. Popescu & Etzioni, HLT/EMNLP 2005]

Solution?

- Use *NE Detection* and *NP Chunking* to identify features
- Find opinions either within the NPs “*a very high resolution*”, or within adjacent constituents using parsing
- Match opinions (using stemming or lemmatization) against a lexicon containing polarity information
- Sort and rank opinions based on the number of reviews and strength

Question-Answering (QA)

Answering Factual Questions

A task somewhat related to automatic summarization is answering (factual) questions posed in natural languages.

Examples

From TREC-9 (2000):

- *Who invented the paper clip?*
- *Where is the Danube?*
- *How many years ago did the ship Titanic sink?*

The TREC Competition

The *Text REtrieval Conference* (TREC), also organized by NIST, includes a QA track.

QA Systems

Typical Approach in QA

Most QA systems roughly follow a three-step process:

Retrieval Step: find documents from a set that might be relevant for the question

Answer Detection Step: process retrieved documents to find possible answers

Reply Formulation Step: create an answer in the required format (single NP, full sentence etc.)

How to find the answer?

Again, a multitude of approaches:

Syntactic: find matching patterns or parse (sub-)trees (with some transformations) in both Q and A

Semantic: transform both Q and A into a logical form and use inference to check consistency

Google: plug the question into Google and select the answer with a syntactic strategy...

QA Systems

Typical Approach in QA

Most QA systems roughly follow a three-step process:

Retrieval Step: find documents from a set that might be relevant for the question

Answer Detection Step: process retrieved documents to find possible answers

Reply Formulation Step: create an answer in the required format (single NP, full sentence etc.)

How to find the answer?

Again, a multitude of approaches:

Syntactic: find matching patterns or parse (sub-)trees (with some transformations) in both Q and A

Semantic: transform both Q and A into a logical form and use inference to check consistency

Google: plug the question into Google and select the answer with a syntactic strategy...

Google does some QA. . .

Ask Google: When was Julius Caesar born?

The screenshot shows a web browser window titled "When was Julius Caesar born? - Google Search - Konqueror". The address bar contains the search query: "G=When+was+Julius+Caesar+born%3F&btnG=Search". The search results are displayed under the heading "Web" and show "Results 1 - 10 of about 1,750,000 for When was Julius Caesar born?. (0.07 seconds)".

The top result is for "Julius Caesar" with the sub-heading "Date of Birth: 101 BC" and a link to "http://www.who2.com/juliuscaesar.html".

Below this, there is a section for "Book results for When was Julius Caesar born?" listing three books:

- [The Conquest of Gaul](#) - by **Julius Caesar** - 272 pages
- [The Civil War](#) - by **Julius Caesar** - 368 pages
- [Max Notes - Julius Caesar](#) - by William Shakespeare, Joseph E Scalia - 96 pages

15 Introduction

16 Summarization

17 Opinion Mining

18 Question-Answering (QA)

19 Text Mining in Biology and Biomedicine

- Introduction
- The BioRAT System
- Mutation Miner

20 References

Text Mining in the Biological Domain

Biological Research

Like in other disciplines, researchers and practitioners in biology

- need up-to-date information
- but have too much literature to cope with

Particular to Biology

- biological databases containing results of experiments
- manually curated databases
- central repositories for literature (PubMed/Medline/Entrez)

General Idea of our Work

Support researchers in biology, by information extraction (automatic curation support) and combining NLP results with databases and end user's tools

Text Mining in the Biological Domain

Biological Research

Like in other disciplines, researchers and practitioners in biology

- need up-to-date information
- but have too much literature to cope with

Particular to Biology

- biological databases containing results of experiments
- manually curated databases
- central repositories for literature (PubMed/Medline/Entrez)

General Idea of our Work

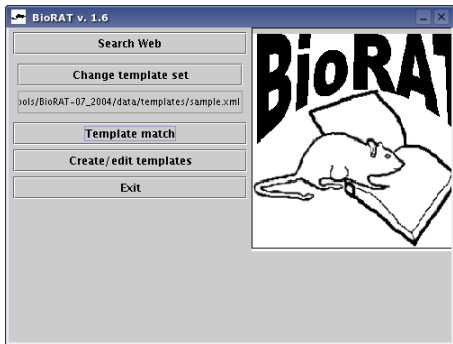
Support researchers in biology, by information extraction (automatic curation support) and combining NLP results with databases and end user's tools

The BioRAT System

BioRAT

BioRAT is a *search engine and information extraction tool for biological research*

- developed at University College London (UCL) in cooperation with GlaxoSmithKline



BioRAT provides

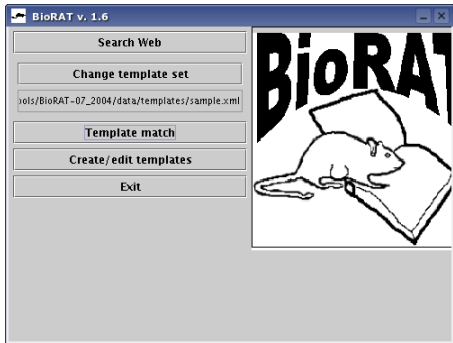
- a web spidering/information retrieval engine
- an information extraction system based on GATE
- a “template” design tool for IE patterns

The BioRAT System

BioRAT

BioRAT is a search engine and information extraction tool for biological research

- developed at University College London (UCL) in cooperation with GlaxoSmithKline



BioRAT provides

- a web spidering/information retrieval engine
- an information extraction system based on GATE
- a “template” design tool for IE patterns

BioRAT: Information Retrieval

The screenshot shows the BioRAT Web search interface. At the top, there is a menu bar with 'File', 'Search', and 'Help'. Below the menu bar, the search query is 'Xylanase'. There are several buttons: 'Search PubMed', 'Search Google', 'Clear screen', 'Copy abstract', and 'Save abstract'. To the right of these buttons, there are input fields for 'Num. to get:' (set to 10) and 'Limit search:' (set to 'No limit'), along with an 'Exit' button. A checkbox for 'Batch mode?' is also present.

The search results are displayed in a list on the left side of the window. Each result includes the date, author(s), and PMID, followed by a link to the abstract and, where available, a link to the full text. The results are as follows:

- (2004 Dec) Author(s): PMID 15576784 [Abstract](#) (no paper available)
- (2004 Dec 3) Author(s): PMID 15575727 [Abstract](#) (no paper available)
- (2004 Oct) Author(s): PMID 15564518 [Abstract](#) [Full text from pcp.oupjournals.org](#)
- (2004 Nov) Author(s): PMID 15556384 [Abstract](#) (no paper available)
- (2004 Nov 17) Author(s): PMID 15537325 [Abstract](#) [Full text from dx.doi.org](#)
- (2004 Jul-Aug) Author(s): PMID 15535463 [Abstract](#) (no paper available)
- (2004 Nov) Author(s): PMID 15531785 [Abstract](#) [Full text from www.humanapress.com](#)
- (2004 Aug) Author(s): PMID 15527073 [Abstract](#) (no paper available)

The full text of the first result is displayed on the right side of the window. It is a scientific paper titled 'Isolation and Expression of the xynB Gene and Its Product, XynB, a Consistent Component of the Clostridium cellulovorans Cellulosome'. The abstract text is as follows:

1: J Bacteriol. 2004 Dec;186(24):8347-55. Isolation and Expression of the xynB Gene and Its Product, XynB, a Consistent Component of the Clostridium cellulovorans Cellulosome. Han SO, Yukawa H, Inui M, Doi RH. Section of Molecular and Cellular Biology, University of California, Davis, CA 95616. rhdoi@ucdavis.edu. The nucleotide sequence of the Clostridium cellulovorans xynB gene, which encodes the XynB xylanase, consists of 1,821 bp and encodes a protein of 607 amino acids with a molecular weight of 65,976. XynB contains a typical N-terminal signal peptide of 29 amino acid residues, followed by a 147-amino-acid sequence that is homologous to the family 4-9 (subfamily 9 in family 4) carbohydrate-binding domain. Downstream of this domain is a family 10 catalytic domain of glycosyl hydrolase. The C terminus separated from the catalytic domain by a short linker sequence contains a dockerin domain responsible for cellulosome assembly. The XynB sequence from mass spectrometry and N-terminal amino acid sequence analyses agreed with that deduced from the nucleotide sequence. XynB was highly active toward xylan, but not active toward carboxymethyl cellulose. The enzyme was optimally active at 40 degrees C and pH 5.0. Northern hybridizations revealed that xynB is transcribed as a monocistronic 1.9-kb mRNA. RNA ligase-mediated rapid amplification of 5' cDNA ends by PCR (RLM-5'RACE PCR) analysis of C. cellulovorans RNA identified a single transcriptional start site of xynB located 47 bp upstream from the first nucleotide of the translation initiation codon. Alignment of the xynB promoter region provided evidence for highly conserved sequences that exhibited strong similarity to the sigma(A) consensus promoter sequences of gram-positive bacteria. Expression of xynB mRNA increased from early to middle exponential phase and decreased during the early stationary phase when the cells were grown on cellobiose. No alternative promoter was observed by RLM-5'RACE PCR and reverse transcriptase PCR analyses during expression. The analysis of the products from xylan hydrolysis by thin-layer chromatography indicated its endoxylanase activity. The results suggest that XynB is a consistent and major cellulosomal enzyme during growth on cellulose or xylan. PMID: 15576784 [PubMed - in process]

At the bottom of the window, the status bar indicates: 'Status: Finished searching PubMed (Retrieved 10 matches out of 1350)'.

BioRAT: Information Retrieval

BioRAT Web search

File Search Help

Query: Xylanase

Search PubMed Search Google Clear screen Copy abstract Save abstract

Num. to get: 10 Limit search: No limit Exit

Batch mode?

... Alkaline-active xylanase produced by an alkaliphilic *Bacillus* sp isolated from kraft pulp ... 1 ml⁻¹) of xylanase when cultivated in alkaline medium at pH 9. ... [Download paper](#) from <http://www.fpl.fs.fed.us/documents/pdf/1995/yang95a.pdf>

Xylanase User Guide Solutions for Crystal Growth Overview The Xylanase crystals are dissolved by adding phosphate buffer and glycerol to the final crystal cake. ... [Download paper](#) from <http://www.hamptonresearch.com/support/guides/7104G.pdf>

... The consequence can be that an added xylanase, for a certain application, is not as effective as desired. ... Xylanase Inhibitors Occurring in Cereals ... [Download paper](#) from <http://www.voeding.tno.nl/Common/PDF/voe284e.pdf>

... 1 F001 0112 Xylanase and cellulase: Structure, Mechanism and Applications (1993-1996) Hemicellulases and Cellulases: Structure, Mechanism and Applications (... [Download paper](#) from http://www.a-b.tugraz.at/data/Zentrum/PDFs_SFB/112%20Endbericht.pdf

21 Effect of xylanase addition in feed containing either

Status: Google returned 10 items (from a total of 1610)

1: J Bacteriol. 2004 Dec;186(24):8347-55. Isolation and Expression of the xynB Gene and Its Product, XynB, a Consistent Component of the Clostridium cellulovorans Cellulosome. Han SO, Yukawa H, Inui M, Doi RH. Section of Molecular and Cellular Biology, University of California, Davis, CA 95616. rhdoi@ucdavis.edu. The nucleotide sequence of the Clostridium cellulovorans xynB gene, which encodes the XynB xylanase, consists of 1,821 bp and encodes a protein of 607 amino acids with a molecular weight of 65,976. XynB contains a typical N-terminal signal peptide of 29 amino acid residues, followed by a 147-amino-acid sequence that is homologous to the family 4-9 (subfamily 9 in family 4) carbohydrate-binding domain. Downstream of this domain is a family 10 catalytic domain of glycosyl hydrolase. The C terminus separated from the catalytic domain by a short linker sequence contains a dockerin domain responsible for cellulosome assembly. The XynB sequence from mass spectrometry and N-terminal amino acid sequence analyses agreed with that deduced from the nucleotide sequence. XynB was highly active toward xylan, but not active toward carboxymethyl cellulose. The enzyme was optimally active at 40 degrees C and pH 5.0. Northern hybridizations revealed that xynB is transcribed as a monocistronic 1.9-kb mRNA. RNA ligase-mediated rapid amplification of 5' cDNA ends by PCR (RLM-5'RACE PCR) analysis of C. cellulovorans RNA identified a single transcriptional start site of xynB located 47 bp upstream from the first nucleotide of the translation initiation codon. Alignment of the xynB promoter region provided evidence for highly conserved sequences that exhibited strong similarity to the sigma(A) consensus promoter sequences of gram-positive bacteria. Expression of xynB mRNA increased from early to middle exponential phase and decreased during the early stationary phase when the cells were grown on cellobiose. No alternative promoter was observed by RLM-5'RACE PCR and reverse transcriptase PCR analyses during expression. The analysis of the products from xylan hydrolysis by thin-layer chromatography indicated its endoxylanase activity. The results suggest that XynB is a consistent and major cellulosomal enzyme during growth on cellulose or xylan. PMID: 15576784 [PubMed - in process]

BioRAT: Information Extraction

Template-based Extraction (actually regular expressions)

- Preprocessing provides *Tokens* and *POS tags*
- Gazetteering step uses lists derived from SwissProt and MeSH to annotate *entities* (genes, proteins, drugs, procedures, ...)
- Templates (JAPE grammars) define patterns for extraction

Templates

- *Sample*: find pattern
<noun> <prep>
<drug/chemical>
- *DIP*: find
protein-protein
interactions

Example Grammar

```
Rule: sample1
Priority: 1000
(
  ({Token.category == NN}):block0
  ({Token.category == IN}):block1
  ({Lookup.majorType
    == "chemicals_and_drugs"}):block2
) --> (add result...)
```

BioRAT: Information Extraction

Template-based Extraction (actually regular expressions)

- Preprocessing provides *Tokens* and *POS tags*
- Gazetteering step uses lists derived from SwissProt and MeSH to annotate *entities* (genes, proteins, drugs, procedures, ...)
- Templates (JAPE grammars) define patterns for extraction

Templates

- *Sample*: find pattern
<noun> <prep>
<drug/chemical>
- *DIP*: find
protein-protein
interactions

Example Grammar

```
Rule: sample1
Priority: 1000
(
  ({Token.category == NN}):block0
  ({Token.category == IN}):block1
  ({Lookup.majorType
    == "chemicals_and_drugs"}):block2
) --> (add result...)
```

BioRAT: Information Extraction

Template-based Extraction (actually regular expressions)

- Preprocessing provides *Tokens* and *POS tags*
- Gazetteering step uses lists derived from SwissProt and MeSH to annotate *entities* (genes, proteins, drugs, procedures, ...)
- Templates (JAPE grammars) define patterns for extraction

Templates

- *Sample*: find pattern
<noun> <prep>
<drug/chemical>
- *DIP*: find
protein-protein
interactions

Example Grammar

```
Rule: sample1
Priority: 1000
(
  ({Token.category == NN}):block0
  ({Token.category == IN}):block1
  ({Lookup.majorType
    == "chemicals_and_drugs"}):block2
) --> (add result...)
```

BioRAT: Extraction Results

Text file	Rule	Position	Length	block0	block1	block2	Context
...0.0: Op1/BioCase /Xylanase3.txt	sample1	6.61%	3	addition	of	disulfide	In practice, thermal stability of xylanases has been improved by addition of disulfide bridges (1115), engineering polar side chains into a protein surface (16), increasing aromatic interactions (17), and other amino acid substitutions (11 21).
...0.0: Op1/BioCase /Xylanase3.txt	sample1	9.61%	3	characterization	of	proteins	The advent of electrospray ionization (ESI) has opened up mass spectrometry (MS) as one of the most powerful analytical techniques for structural and functional characterization of proteins (28).
...0.0: Op1/BioCase /Xylanase3.txt	sample1	10.57%	3	configuration	for	protein	ESI in combination with a Fourier transform ion cyclotron resonance (FT ICR) mass analyzer (30, 31) represents the most powerful instrumental configuration for protein analyses.
...0.0: Op1/BioCase /Xylanase3.txt	sample1	28.96%	3	presence	of	acetonitrile	In contrast, clear changes appeared in the presence of acetonitrile.
...0.0: Op1/BioCase /Xylanase3.txt	sample1	38.9%	3	number	of	disulfide	Hence, to verify the number of disulfide bridges (i.e.
...0.0: Op1/BioCase /Xylanase3.txt	sample1	41.52%	3	number	of	disulfide	disulfide reduced), and 2red (two disulfides reduced, in DB1 only) were used to assign the correct number of disulfide bridges, i.e.
...0.0: Op1/BioCase /Xylanase3.txt	sample1	46.73%	3	ratio	in	solution	D(t)) Htotal aHfaste kfastt + Hslowekslowt) in which Htotal is a total number of exchangeable hydrogens in TRX II (341 based on the sequence), a is a D/H ratio in solution (0.95), Hfast and Hslow are the numbers for fast and slow exchanging hydrogens, and kfast and kslow are the corresponding pseudo first order rate constants.
...0.0: Op1/BioCase /Xylanase3.txt	sample1	60.64%	3	number	of	disulfides	Standard deviation calculations between the theoretical and the experimental isotopic distributions confirmed the expected number of disulfides in DS2, DS5, and DB1 mutants (i.e.

Close

Status:

BioRAT: Template Design Tool

BioRAT Template Designer

1 : Br J Clin Pharmacol . 2004 Apr ; 57 (4) : 522 - 524 .
 Overanticoagulation associated with combined use of lactulose and coumarin anticoagulants . Visser LE , Penning Van Beest FJ , Wilson JH , Vulto AG , Kasbergen AA , De Smet PA , Hofman A , Stricker BH . Pharmaco epidemiology Unit , Departments of Internal Medicine and **Epidemiology** & amp ; Biostatistics , Erasmus MC , Rotterdam , The Netherlands . Some medical textbooks on drug interactions take note of the potential interaction between laxatives and coumarin anticoagulants , but epidemiological evidence that this interaction is of practical importance is lacking . We conducted a follow up study in a large population based cohort to investigate which laxatives are associated with overanticoagulation during therapy with coumarins . Of the 1124 patients in the cohort , 351 developed an International Normalized Ratio & amp ; gt ; 6 . 0 . The only laxative with a moderate but significantly increased relative risk of overanticoagulation was lactulose (relative risk 3 . 4 , 95 % confidence interval 2 . 2 , 5 . 3) . In view of the widespread use of lactulose , especially among the elderly , awareness of this potential drug interaction is required . PMID : 15025752 [PubMed as supplied by publisher]

<MACRO: VERB> <LOOKUP: swissprot proteins , proteins> <LITERAL: follow up>

Status: Found 3 matches

Any: **NOUN** Or **Optional**

Word:

Part of spe...: **NNP** "noun, proper, singular"

Stem:

Gazetteer matches: <biological sciences , health occupations>

BioRAT: Some Observations

BioRAT Performance

- Authors report 39% recall and 48% precision on the DIP task
- Comparable to the SUISEKI system (Blaschke et al.), which is statistics-based

System Design

More interestingly,

- BioRAT is rather “low” on NLP knowledge,
- yet surprisingly useful for Biologists

Interesting pattern:

- NLP is “just another” system component
- Users (Biologists) are empowered: no need for computational linguists to add/modify/remove grammar rules

BioRAT: Some Observations

BioRAT Performance

- Authors report 39% recall and 48% precision on the DIP task
- Comparable to the SUISEKI system (Blaschke et al.), which is statistics-based

System Design

More interestingly,

- BioRAT is rather “low” on NLP knowledge,
- yet surprisingly useful for Biologists

Interesting pattern:

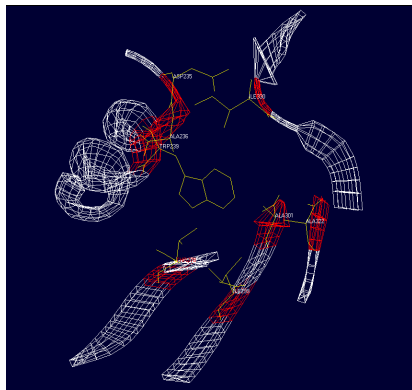
- NLP is “just another” system component
- Users (Biologists) are empowered: no need for computational linguists to add/modify/remove grammar rules

MutationMiner: Motivation

Challenge

Support Bio-Engineers designing proteins:

- need up-to-date, relevant information from research literature
- need for automated updates
- need for integration with structural biology tools



MutationMiner: Background

Existing Resources

Protein Mutant Database

*Center for Information Biology and DNA Data Bank of Japan
National Institute of Genetics*

- 1999: authors quote 3-year backlog of unprocessed publications
- Funding for manual curation limited / declining
- Manual data submission is slow and incomplete
- Sequence and structure databases expanding
- New techniques: Directed Evolution
- New alignment algorithms: e.g. Fugue, Muscle

Protein Mutant Database

Example PMD Entry (manually curated)

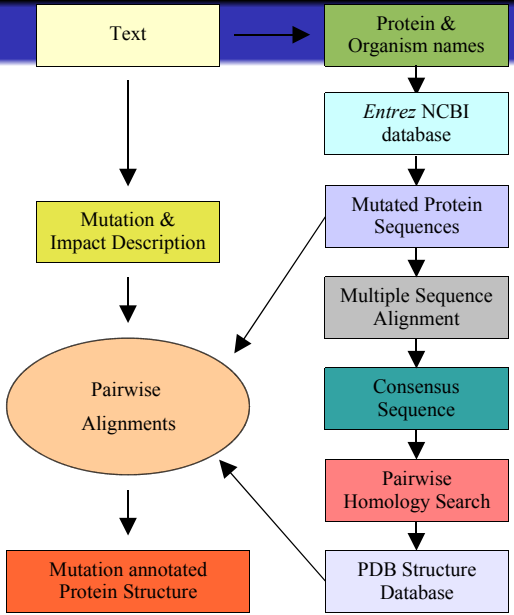
ENTRY A931290 - Artificial 1921240
 AUTHORS Lee Y.-E., Lowe S.E., Henrissat B. & Zeikus J.G.
 JOURNAL J.Bacteriol. (1993) 175(18), 5890-5898 [LINK-TO-MEDLINE]
 TITLE Characterization of the active site and thermostability
 regions of endoxylanase from *Thermoanaerobacterium saccharolyticum*
 CROSS-REFERENCE A48490 [LINK TO PIR "A48490"] No PDB-LINK for "A48490"
 PROTEIN Endoxylanase (endo-1,4-beta-xylanase) #EC3.2.1.8
 SOURCE *Thermoanaerobacterium saccharolyticum*
 N-TERMINAL MMKNN
 EXPRESSION-SYSTEM *Escherichia coli*
 CHANGE Asp 537 Asn FUNCTION Endoxylanase activity [0]
 CHANGE Glu 541 Gln FUNCTION Endoxylanase activity [=]
 CHANGE His 572 Asn FUNCTION Endoxylanase activity [=]
 CHANGE Glu 600 Gln FUNCTION Endoxylanase activity [0]
 CHANGE Asp 602 Asn FUNCTION Endoxylanase activity [0]

MutationMiner: Goal

Aim

Develop a system to

- extract annotations regarding mutations from full-text papers; and
- legitimately link them to protein structure visualizations



MutationMiner NLP: Input

Input documents are typically in HTML, XML, or PDF formats:

9556

Biochemistry 2004, 43, 9556–9566

Characterization of Mutant Xylanases Using Fourier Transform Ion Cyclotron Resonance Mass Spectrometry: Stabilizing Contributions of Disulfide Bridges and N-Terminal Extensions[†]

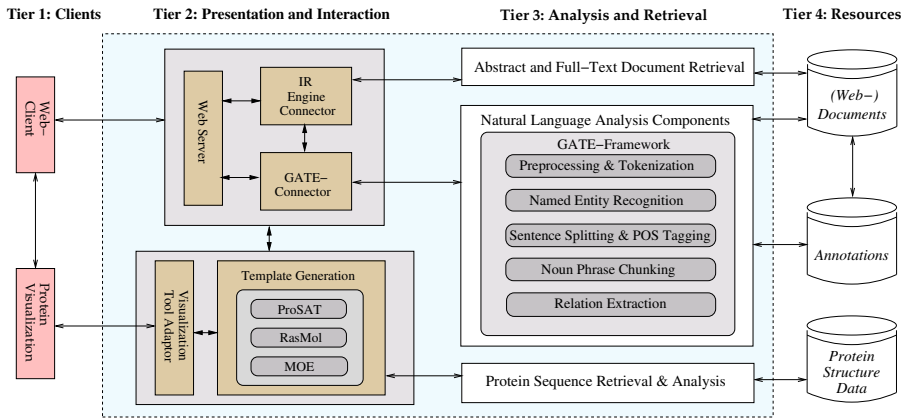
Janne Jänis,[‡] Ossi Turunen,[§] Matti Leisola,[§] Peter J. Derrick,^{||} Juha Rouvinen,[‡] and Pirjo Vainiotalo^{*‡}

Department of Chemistry, University of Joensuu, P.O. Box 111, FI-80101 Joensuu, Finland, Laboratory of Bioprocess Engineering, Helsinki University of Technology, P.O. Box 6100, FI-02015 HUT, Finland, and Department of Chemistry, Institute of Mass Spectrometry, University of Warwick, Coventry CV4 7AL, United Kingdom

Received February 27, 2004; Revised Manuscript Received May 17, 2004

ABSTRACT: Structural properties and thermal stability of *Trichoderma reesei* endo-1,4- β -xylanase II (TRX II) and its three recombinant mutants were characterized using electrospray ionization Fourier transform ion cyclotron resonance (ESI FT-ICR) mass spectrometry and hydrogen/deuterium (H/D) exchange reactions. TRX II has been previously stabilized by a disulfide bridge C110–C154 and other site-directed mutations (TRX II mutants DS2 and DS5). Very recently, a highly thermostable mutant was introduced by combining mutations of DS5 with an N-terminal disulfide bridge C2–C28 (mutant DB1). Accurate

MutationMiner Architecture



MutationMiner: NLP Subsystem

NLP Steps

Tokenization split input into tokens

Gazetteering using lists derived from Swissprot and MeSH

Named Entity recognition find proteins, mutations, organisms

Sentence splitting sentence boundary detection

POS tagging add part-of-speech tags

NP Chunking e.g. *the/DET catalytic/MOD activity/HEAD*

Relation detection find protein-organism and protein-mutation relations

“Wild-type and mutated xylanase II proteins (termed E210D and E210S) were expressed in S. cerevisiae grown in liquid culture.”

- Gate
 - Applications
 - Mutation Miner
 - Language Resources
 - GATE document_00044
 - xylanase
 - GATE corpus_00048
 - Processing Resources
 - ANNIE VP Chunker_0005F
 - maxNP Transducer
 - Bio-Transducer
 - CLaC POS Tagger
 - CLaC Sentence Splitter
 - Bio Gazetteer (ignore case, whole v
 - CLaC Gazetteer
 - StemmerPlus
 - Tokeniser
 - Document Reset PR
 - Data stores

Text Annotations Annotation Sets Print

To determine whether enzymatic activity is necessary for elicitor activity, we used site-directed mutagenesis to reduce the catalytic activity of xylanase II from *Trichoderma reesei*.

Compared with the wild-type form of xManase II, E210D had >100-fold and E210S 1,000-fold lower enzymatic activity.

Type	Set	Start	End	Features
Lookup	Default	113	119	(majorType=protein_expression, minorType=lower_case)
NP	Default	120	142	(HEAD_START=134, HEAD_END=142)
NP	Default	120	142	(DET=the, MOD=catalytic, HEAD=activity, HEAD_START=134, HEAD_END=142)
Lookup	Default	124	133	(majorType=swissprot_proteins, minorType=proteins)
Lookup	Default	134	142	(majorType=softtox_enzyme_change, minorType=extra_list)
Lookup	Default	146	157	(majorType=swissprot_proteins, minorType=proteins)
NP	Default	146	157	(HEAD_START=155, HEAD_END=157)
NP	Default	146	157	(HEAD_END=157, HEAD_START=155, HEAD=II, MOD=xylanase)
Protein	Default	146	157	()
Lookup	Default	155	157	(majorType=swissprot_proteins, minorType=proteins)
Lookup	Default	155	157	(majorType=swissprot_proteins, minorType=genes)
Lookup	Default	163	174	(majorType=organisms, minorType=algae_and_fungi)
NP	Default	163	181	(HEAD_END=181, HEAD_START=175, HEAD=reesei, MOD=Trichoderma)
NP	Default	163	181	(HEAD_START=175, HEAD_END=181)
Organism	Default	163	181	()
NP	Default	198	216	(DET=the, MOD=wild-type, HEAD=form, HEAD_START=212, HEAD_END=216)
NP	Default	198	216	(HEAD_START=212, HEAD_END=216)
Lookup	Default	220	231	(majorType=swissprot_proteins, minorType=proteins)
NP	Default	220	231	(HEAD_END=231, HEAD_START=229, HEAD=II, MOD=xylanase)
NP	Default	220	231	(HEAD_START=229, HEAD_END=231)

Annotations Editor Features Editor Initialisation Parameters

- Default annotations
- Lookup
 - Mutation
 - NP
 - Organism
 - Prot-Mut
 - Prot-Org
 - Protein
 - Sentence
 - SpaceToken
 - Split
 - Token
 - YG
 - maxNP
 - tempoNP
- Original markups annotations
- paragraph

Removes this resource from the system

MutationMiner: Further Processing

Results

- Results are information about *Proteins*, *Organisms*, and *Mutations*, along with *context information*

Next Step

- These results could already be used to (semi-)automatically curate PMD entries
- But remember the original goal: integrate results into end user's tools
- Needs data that can be further processed by bioinformatics tools

Thus, we need to find the corresponding real-world entities in biological databases: *amino acid sequences*

MutationMiner: Further Processing

Results

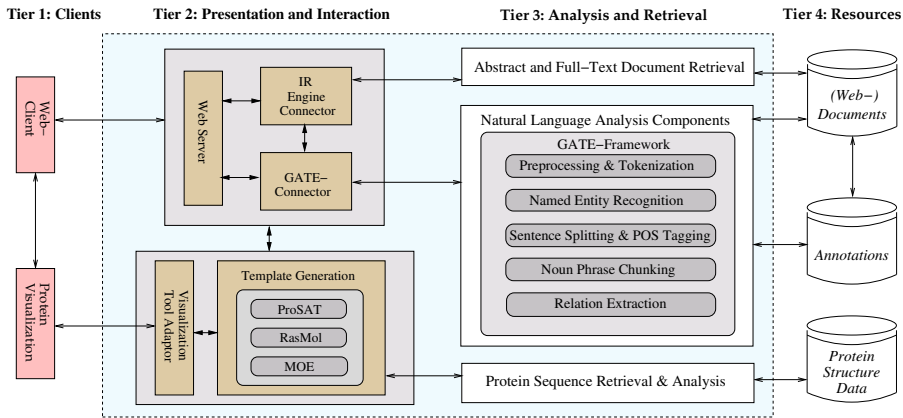
- Results are information about *Proteins*, *Organisms*, and *Mutations*, along with *context information*

Next Step

- These results could already be used to (semi-)automatically curate PMD entries
- But remember the original goal: integrate results into end user's tools
- Needs data that can be further processed by bioinformatics tools

Thus, we need to find the corresponding real-world entities in biological databases: *amino acid sequences*

MutationMiner Architecture



MutationMiner: Sequence Retrieval

Sequence Retrieval

- Retrieval of FASTA formatted sequences for protein accessions obtained by NLP analysis of texts
- Obtained through querying *Entrez* NCBI database (E-fetch)

The screenshot shows the NCBI Entrez Protein search interface. The search query is "Xylanase C Aspergillus Kawachii". The results are displayed in FASTA format. The sequence is annotated with three domains: Esterase (green), CBD_IV (blue), and Glyco_10 (red).

Search: Protein for Xylanase C Aspergillus Kawachii

Display: FASTA Show: 20 Send to: Text

Items 1 - 4 of 4

1: [P33557](#)
Endo-1,4-beta-xylanase 3 precursor (Xylanase 3) (1,4-beta-D-xylan xylanohydrolase 3) (Xylanase C)
gi|465492|sp|P33557|XYN3_ASPKA[465492]

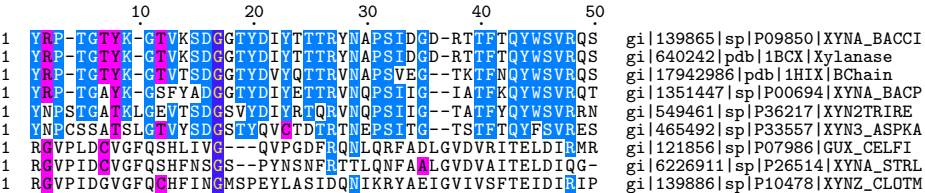
```
>gi|465492|sp|P33557|XYN3_ASPKA Endo-1,4-beta-xylanase 3 precursor (Xylanase 3)
MKVTAASAGLLGHAFAPVPQPVLVRSAGINYPVQNYNGNLADFTYDESAGTFMSYWEDGVSSDFVVLG
WTTGSSNAISYSAEYSASGSSSYLAVYGVVNPQAEYIIVEDYGDYNPCSSATS LGTVYSDGSTYQVCTD
TRTNEPSITGTSTFTQYFSVRESTRTSGTVTVANHFNFWAQHGFGNSDFNYQVMAVEAWSGAGSASVTIS
S
```

1 100 200 300 400 500 600 700 800 837

Esterase CBD_IV Glyco_10

MutationMiner: Sequence Analysis

CLUSTAL W (1.82) multiple sequence alignment



- sequence analyzed and sliced in regions using CDD (*conserved domain database*) search tools
- iterative removal of outlying sequences through statistical scoring using *Alistat*
- generation of a consensus sequence using a HMM (HMMER)
- locate NLP-extracted mutations on sequence

Sequence Analysis Results

Amino Acid Sequence Analysis

- We now have a set of filtered sequences, describing proteins and their mutations
- Still not a very intuitive presentation of results

Suitable visualization needed!

3D-Structure Visualization

- Idea: map mutations of proteins directly to a 3D-visualization of their structural representation
- However, for this we need to find a 3D-model (homolog)

Solution: access Protein Data Bank (PDB) using BLAST for a suitable 3D-model and map NLP results onto this structure

Sequence Analysis Results

Amino Acid Sequence Analysis

- We now have a set of filtered sequences, describing proteins and their mutations
- Still not a very intuitive presentation of results

Suitable visualization needed!

3D-Structure Visualization

- Idea: map mutations of proteins directly to a 3D-visualization of their structural representation
- However, for this we need to find a 3D-model (homolog)

Solution: access Protein Data Bank (PDB) using BLAST for a suitable 3D-model and map NLP results onto this structure

MutationMiner: PDB Structure Retrieval

Title Crystallographic Analyses Of Family 11 Endo-1,4-Xylanase Xyl1
Classification Hydrolase
Compound Mol_Id: 1; Molecule: Endo-1,4-Xylanase; Chain: A, B; Ec: 3.2.1.8;
Exp. Method X-ray Diffraction

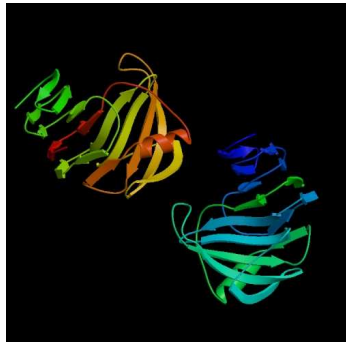
JRNL TITL 2 ENDO-[BETA]-1,4-XYLANASE XYL1 FROM STREPTOMYCES SP. S38

JRNL REF ACTA CRYSTALLOGR.,SECT.D V. 57 1813 2001

JRNL REFN ASTM ABCRE6 DK ISSN 0907-4449

```

...
DBREF  1HIX A    1   190  TREMBL  Q59962   Q59962
DBREF  1HIX B    1   190  TREMBL  Q59962   Q59962
...
ATOM  1 N      ILE A  4  48.459  19.245  17.075  1.00  24.52  N
ATOM  2 CA     ILE A  4  47.132  19.306  17.680  1.00  50.98  C
ATOM  3 C      ILE A  4  47.116  18.686  19.079  1.00  49.94  C
ATOM  4 O      ILE A  4  48.009  17.936  19.465  1.00  70.83  O
ATOM  5 CB     ILE A  4  46.042  18.612  16.837  1.00  50.51  C
ATOM  6 CG1    ILE A  4  46.419  17.217  16.338  1.00  51.09  C
ATOM  7 CG2    ILE A  4  45.613  19.514  15.687  1.00  54.39  C
ATOM  8 CD1    ILE A  4  46.397  17.045  14.836  1.00  46.72  C
ATOM  9 N      THR A  5  46.077  19.024  19.828  1.00  40.65  N
...
MASTER  321  0  0  2  28  0  0  9  3077  2  0  30
END
    
```



MutationMiner: Visualization

Visualization Tools

- ProSAT is a tool to map SwissProt sequence features and Prosite patterns on to a 3D structure of a protein.
- We are now able to upload the 3D structure together with our textual annotations for rendering using a Webmol interface

ProSAT - Protein Structure Annotation Tool

1ref

MOL_ID: 1; MOLECULE: ENDO-1,4-BETA-XYLANASE II; CHAIN: A, B; SYNONYM: X
links: [rcsb](#) or [pqs](#)

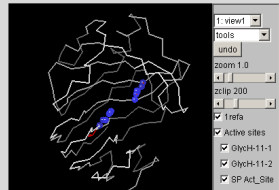
For more in-depth visualisation, choose from:

- [kinemage](#)
- [rasmol](#)
- [chime](#)
- [webmol](#)

The associated data:

- [pdb/pqs file](#)
- [feature residues](#)

Try another pdb entry:



MOL_ID: 1; MOLECULE: ENDO-1,4-BETA-XYLANASE; CHAIN: A, B;

Java Applet Window

Glycosyl hydrolases family 11 active site signature 1.

Site

Motif

Glycosyl hydrolases family 11 active site signature 2.

Site

Motif

Prosite Abundant

N-glycosylation site.

Protein kinase C phosphorylation site > N-glycosylation site.

empty

N-myristoylation site.

11377763: 549461

Table 1 Combinations of XYNII mutations built on the disulfide bridge (more...)

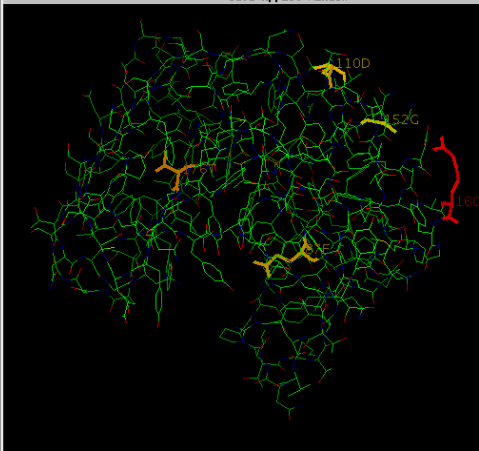
The mutations N11D and N38E did not have any significant effect: (more...)

The combination of the disulfide bridge (110 154) with mutations (more...)

RESET COLORS

WebmolEML

Java Applet Window

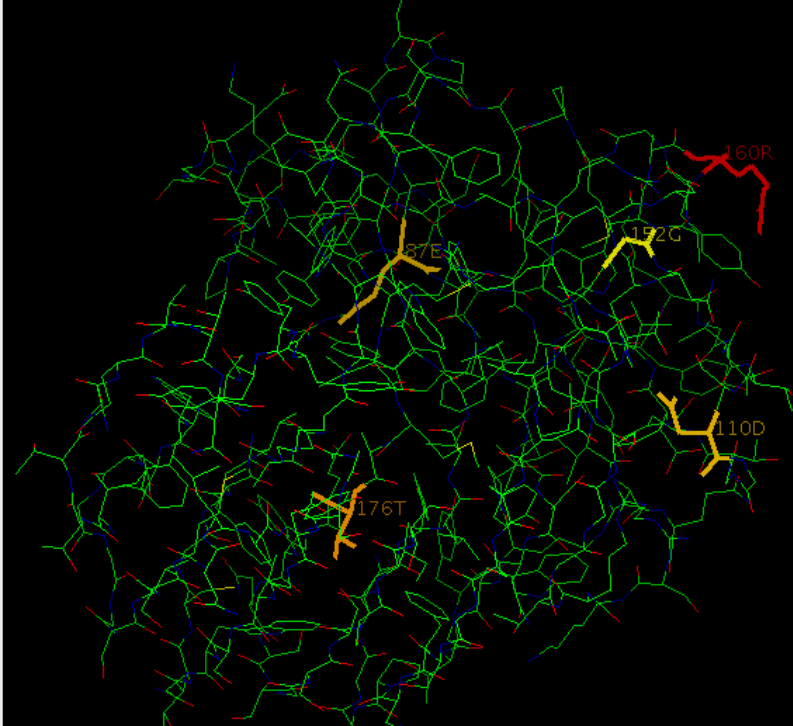


110D

152G

160R

WebMol Open Print C'n'P ResetSlab Center Control Info



Color

Surface

Labels

HOH

Stereo

< >

Rock

Select

Focus

Msure

DMat

Rama

Trace

Implementation and Evaluation

Implementation

NLP subsystem implemented using the GATE architecture

Testing Corpus

First evaluation performed on research literature concerning the *Xylanase* protein family (20 papers)

NLP subsystem partial evaluation results

	Abstract only		Full paper	
	Protein/Organism	Mutations	Protein/Organism	Mutations
Precision	0.88	1.00	0.91	0.84
Recall	0.71	0.85	0.46	0.97
F-Measure	0.79	0.92	0.61	0.90

Implementation and Evaluation

Implementation

NLP subsystem implemented using the GATE architecture

Testing Corpus

First evaluation performed on research literature concerning the *Xylanase* protein family (20 papers)

NLP subsystem partial evaluation results

	Abstract only		Full paper	
	Protein/Organism	Mutations	Protein/Organism	Mutations
Precision	0.88	1.00	0.91	0.84
Recall	0.71	0.85	0.46	0.97
F-Measure	0.79	0.92	0.61	0.90

- Gate
 - Applications
 - Mutation Miner
 - Language Resources
 - GATE document_00044
 - xylanase
 - GATE corpus_00048
 - Processing Resources
 - ANNIE VP Chunker_0005F
 - maxNP Transducer
 - Bio-Transducer
 - CLaC POS Tagger
 - CLaC Sentence Splitter
 - Bio Gazetteer (ignore case, whole v
 - CLaC Gazetteer
 - StemmerPlus
 - Tokeniser
 - Document Reset PR
 - Data stores

Text Annotations Annotation Sets Print

To determine whether enzymatic activity is necessary for elicitor activity, we used site-directed mutagenesis to reduce the catalytic activity of xylanase II from *Trichoderma reesei*.

Compared with the wild-type form of xManase II, E210D had >100-fold and E210S 1,000-fold lower enzymatic activity.

Type	Set	Start	End	Features
Lookup	Default	113	119	(majorType=protein_expression, minorType=lower_case)
NP	Default	120	142	(HEAD_START=134, HEAD_END=142)
NP	Default	120	142	(DET=the, MOD=catalytic, HEAD=activity, HEAD_START=134, HEAD_END=142)
Lookup	Default	124	133	(majorType=swissprot_proteins, minorType=proteins)
Lookup	Default	134	142	(majorType=softtox_enzyme_change, minorType=extra_list)
Lookup	Default	146	157	(majorType=swissprot_proteins, minorType=proteins)
NP	Default	146	157	(HEAD_START=155, HEAD_END=157)
NP	Default	146	157	(HEAD_END=157, HEAD_START=155, HEAD=II, MOD=xylanase)
Protein	Default	146	157	()
Lookup	Default	155	157	(majorType=swissprot_proteins, minorType=proteins)
Lookup	Default	155	157	(majorType=swissprot_proteins, minorType=genes)
Lookup	Default	163	174	(majorType=organisms, minorType=algae_and_fungi)
NP	Default	163	181	(HEAD_END=181, HEAD_START=175, HEAD=reesei, MOD=Trichoderma)
NP	Default	163	181	(HEAD_START=175, HEAD_END=181)
Organism	Default	163	181	()
NP	Default	198	216	(DET=the, MOD=wild-type, HEAD=form, HEAD_START=212, HEAD_END=216)
NP	Default	198	216	(HEAD_START=212, HEAD_END=216)
Lookup	Default	220	231	(majorType=swissprot_proteins, minorType=proteins)
NP	Default	220	231	(HEAD_END=231, HEAD_START=229, HEAD=II, MOD=xylanase)
NP	Default	220	231	(HEAD_START=229, HEAD_END=231)

Annotations Editor Features Editor Initialisation Parameters

- Default annotations
- Lookup
 - Mutation
 - NP
 - Organism
 - Prot-Mut
 - Prot-Org
 - Protein
 - Sentence
 - SpaceToken
 - Split
 - Token
 - YG
 - maxNP
 - tempoNP
- Original markups annotations
- paragraph

Removes this resource from the system

MutationMiner: Conclusions and Ongoing Work

Conclusions

- Integration of bio-NLP, bio-DBs, and bioinformatics tools is very promising and has high practical relevance

Current Work

- Analysis of *Dehalogenase*, *Biphenyl Dioxygenase*, and *Subtilisin*
- Application to human health related scenarios, like *BrcA* protein, which is involved in breast cancer

Future Work

- Extrinsic evaluation with domain specialists, both protein researchers and industry practitioners

MutationMiner: Conclusions and Ongoing Work

Conclusions

- Integration of bio-NLP, bio-DBs, and bioinformatics tools is very promising and has high practical relevance

Current Work

- Analysis of *Dehalogenase*, *Biphenyl Dioxygenase*, and *Subtilisin*
- Application to human health related scenarios, like *BrcA* protein, which is involved in breast cancer

Future Work

- Extrinsic evaluation with domain specialists, both protein researchers and industry practitioners

MutationMiner: Conclusions and Ongoing Work

Conclusions

- Integration of bio-NLP, bio-DBs, and bioinformatics tools is very promising and has high practical relevance

Current Work

- Analysis of *Dehalogenase*, *Biphenyl Dioxygenase*, and *Subtilisin*
- Application to human health related scenarios, like *BrcA* protein, which is involved in breast cancer

Future Work

- Extrinsic evaluation with domain specialists, both protein researchers and industry practitioners

References

General Text Mining Books

- Weiss, Indurkha, Zhang, Damerau, *Text Mining: Predictive Methods for Analyzing Unstructured Information*, Springer, 2005.
- Sophia Ananiadou and John McNaught (Eds.), *Text Mining for Biology and Biomedicine*, Artech House, 2006.
- Inderjeet Mani, *Automatic Summarization*, John Benjamins B.V., 2001

Papers

- Cunningham, *Information Extraction, Automatic*, Encyclopedia of Language and Linguistics, 2005.
<http://gate.ac.uk/sale/e112/ie/>

References (II)

Conferences

- Automatic Summarization: Document Understanding Conference (DUC), <http://duc.nist.gov>
- Text REtrieval Conference: <http://trec.nist.gov/>
- Exploring Attitude and Affect in Text: Theories and Applications, AAI Spring Symposium, 2004, AAI Press, Technical Report SS-04-07

Related Literature

- Zhong & Liu (Eds.), *Intelligent Technologies for Information Analysis*, Springer, 2004
- Peter Morville, *Ambient Findability*, O'Reilly, 2006

And of course...

- <http://rene-witte.net>

Part V

Conclusions

21 Conclusions

Some conclusions

The present

- A complete semantic understanding of natural language is currently (and in the mid-term future) impossible
- However, the available language technologies can already provide more information than simple keyword-indexing
- Implementations are not mainstream yet, but this is changing as more open-source systems become available

The future

- Prediction: within 5 years, Text Mining will start to enter mainstream, similarly to the *Google* effect
- Requires more interdisciplinary work: computational linguists, information system engineers, domain experts

Some conclusions

The present

- A complete semantic understanding of natural language is currently (and in the mid-term future) impossible
- However, the available language technologies can already provide more information than simple keyword-indexing
- Implementations are not mainstream yet, but this is changing as more open-source systems become available

The future

- Prediction: within 5 years, Text Mining will start to enter mainstream, similarly to the *Google* effect
- Requires more interdisciplinary work: computational linguists, information system engineers, domain experts