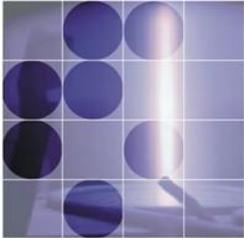


How to Represent Text? ...from Characters to Logic

Marko Grobelnik (marko.grobelnik@ijs.si)

Jozef Stefan Institute (<http://www.ijs.si/>)

Slovenia, Europe



Outline

- ▶ Some Initial thoughts
 - ...quick example why representation matters
- ▶ How we represent Text?
 - Big picture
 - Levels of representation:
 - Lexical
 - Syntactic
 - Semantic
- ▶ Further references
 - ...events, books, videos

Some Initial thoughts

- ▶ Two major steps in Machine learning:
 - Choosing representation (feature engineering)
 - Modeling (statistics+optimization)
- ▶ ...typically people do modeling well and often ignore data representation
- ▶ But...
 - ...**good representation with bad algorithm** gives typically better results than **good algorithm with bad representation**

Quick example:

Why representation matters?

- ▶ Selection of kernels when using SVM (Support Vector Machine) equals selecting the right representation for data

Easy decision problems require simple data representation

The screenshot shows a window titled "SVM - demo" with a 2D plot. The plot contains two classes of data points: blue 'x' marks and red 'x' marks. A diagonal decision boundary, shaded in gray, separates the two classes. An arrow points to this boundary with the text "Linear separation is enough to separate the data".

Point class

PLUS MINUS

Clear

Learning properties

Cost: 100

Linear kernel

Polynomial kernel

Degree: 2

Gaussian kernel

Sigma: 10

Do the right stuff!

Made by:
Blaz Fortuna
<http://kt.ijs.si/blazf/index.html>

**Jozef Stefan
Institute,
Slovenia**
<http://www.ijs.si/>

Harder decision problems require better data representation

The screenshot shows a software window titled "SVM - demo". The main area is a 2D plot with a gray background. On the left side, there are several blue 'x' marks representing one class of data points. On the right side, there are several red 'x' marks representing the other class. A complex, curved decision boundary, shaded in light gray, separates the two classes. An arrow points from the text "Polynomial separation of data" to this boundary.

Point class

PLUS MINUS

Clear

Learning properties

Cost:

Linear kernel

Polynomial kernel

Degree:

Gaussian kernel

Sigma:

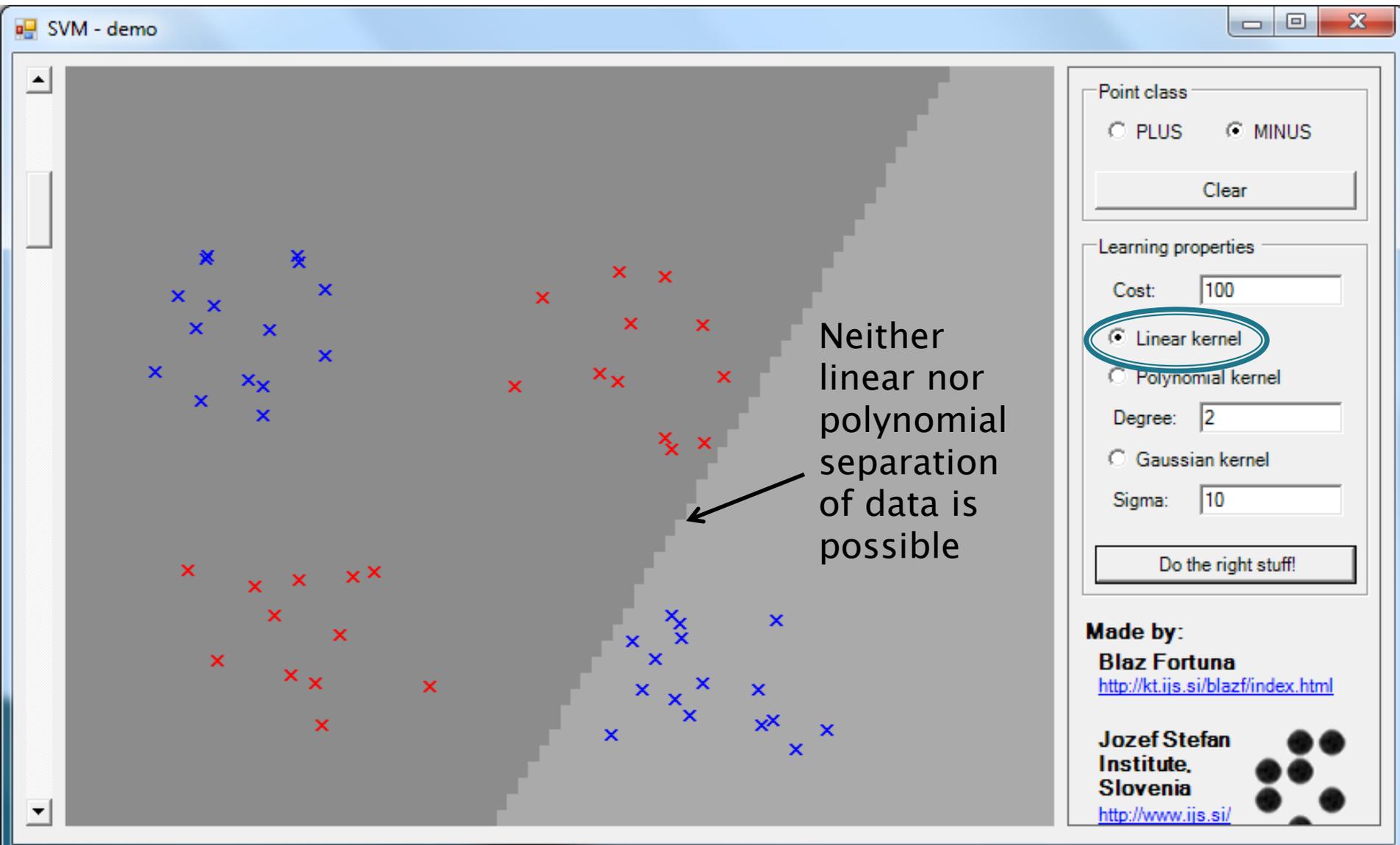
Do the right stuff!

Made by:
Blaz Fortuna
<http://kt.ijs.si/blazf/index.html>

**Jozef Stefan
Institute,
Slovenia**
<http://www.ijs.si/>

Polynomial separation of data

Hard decision problems require sophisticated data representation



The image shows a software window titled "SVM - demo". The main area is a 2D plot with a gray background. It contains two sets of data points: blue 'x' marks and red 'x' marks. A diagonal decision boundary separates the space. The blue points are mostly in the lower-left region, and the red points are mostly in the upper-right region. An arrow points to the decision boundary with the text "Neither linear nor polynomial separation of data is possible".

On the right side, there is a control panel with the following settings:

- Point class:** PLUS, MINUS
- Clear:** [button]
- Learning properties:**
 - Cost:** 100
 - Linear kernel
 - Polynomial kernel
 - Degree:** 2
 - Gaussian kernel
 - Sigma:** 10
- Do the right stuff!** [button]

Made by:
Blaz Fortuna
<http://kt.ijs.si/blazf/index.html>

Jozef Stefan Institute, Slovenia
<http://www.ijs.si/>

Hard decision problems require sophisticated data representation

SVM - demo

Gaussian kernel separates the data

Point class

PLUS MINUS

Clear

Learning properties

Cost: 100

Linear kernel

Polynomial kernel

Degree: 2

Gaussian kernel

Sigma: 10

Do the right stuff!

Made by:
Blaz Fortuna
<http://kt.ijs.si/blazf/index.html>

Jozef Stefan Institute, Slovenia
<http://www.ijs.si/>

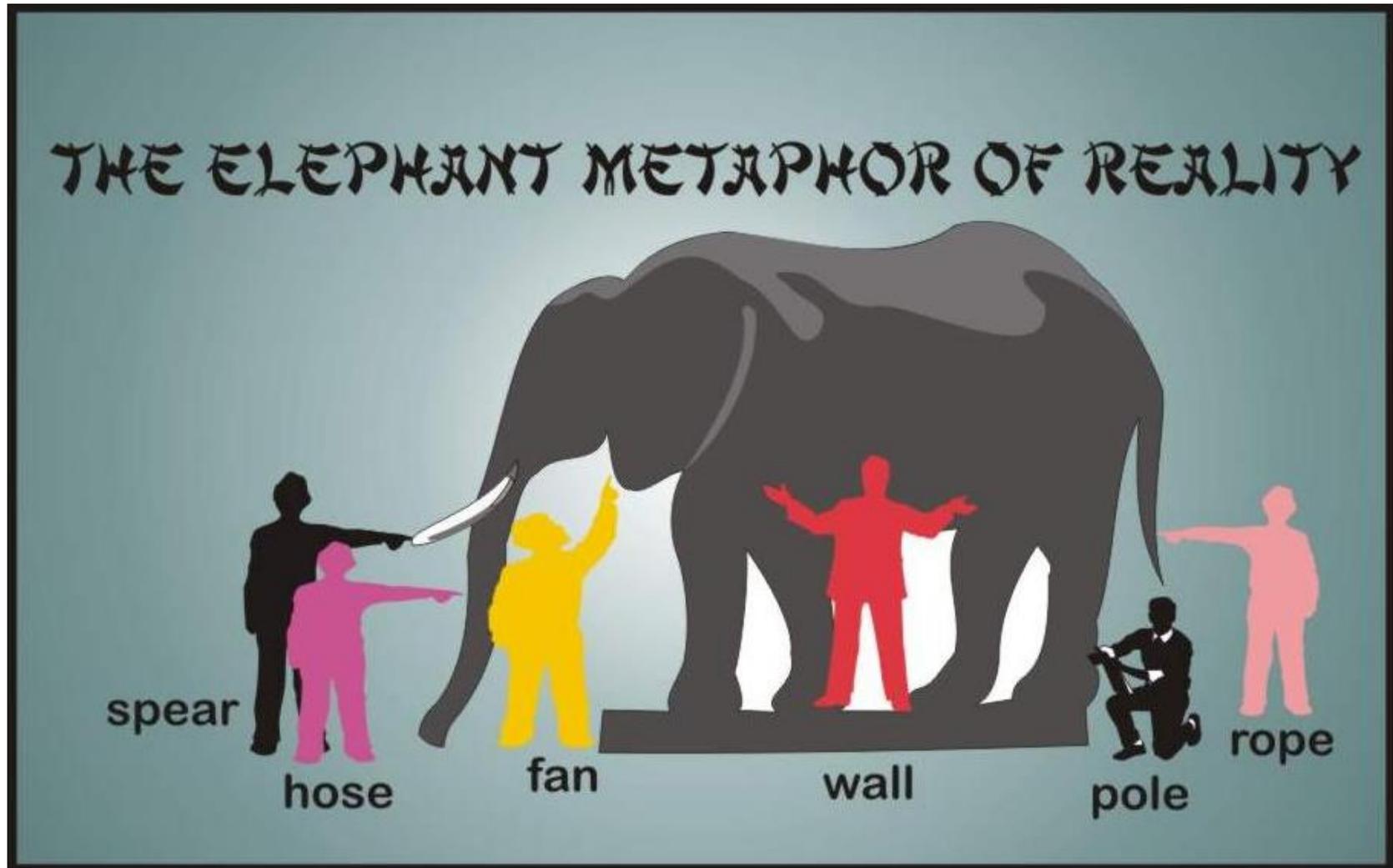
How we represent Text?

The Big picture

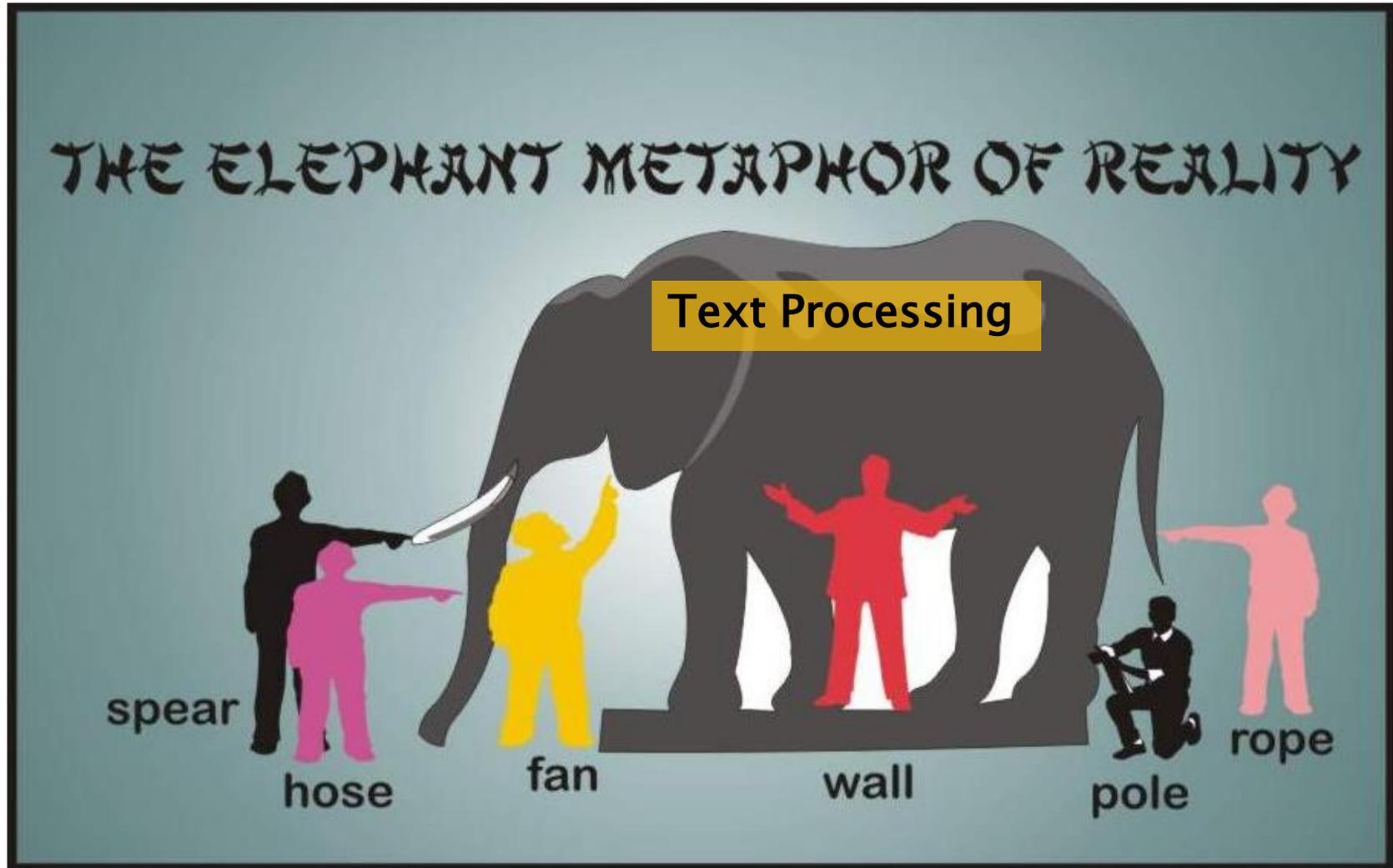
Key paradigms when dealing with data

- ▶ Three key scientific paradigms
 - **Top-down approaches (model driven)**
 - (Traditional NLP, KRR, Semantic Web)
 - **Bottom-up approaches (data driven)**
 - (Machine Learning, Data Mining)
 - **Collaborative approaches (socially driven)**
 - (Web2.0, Social Computing)

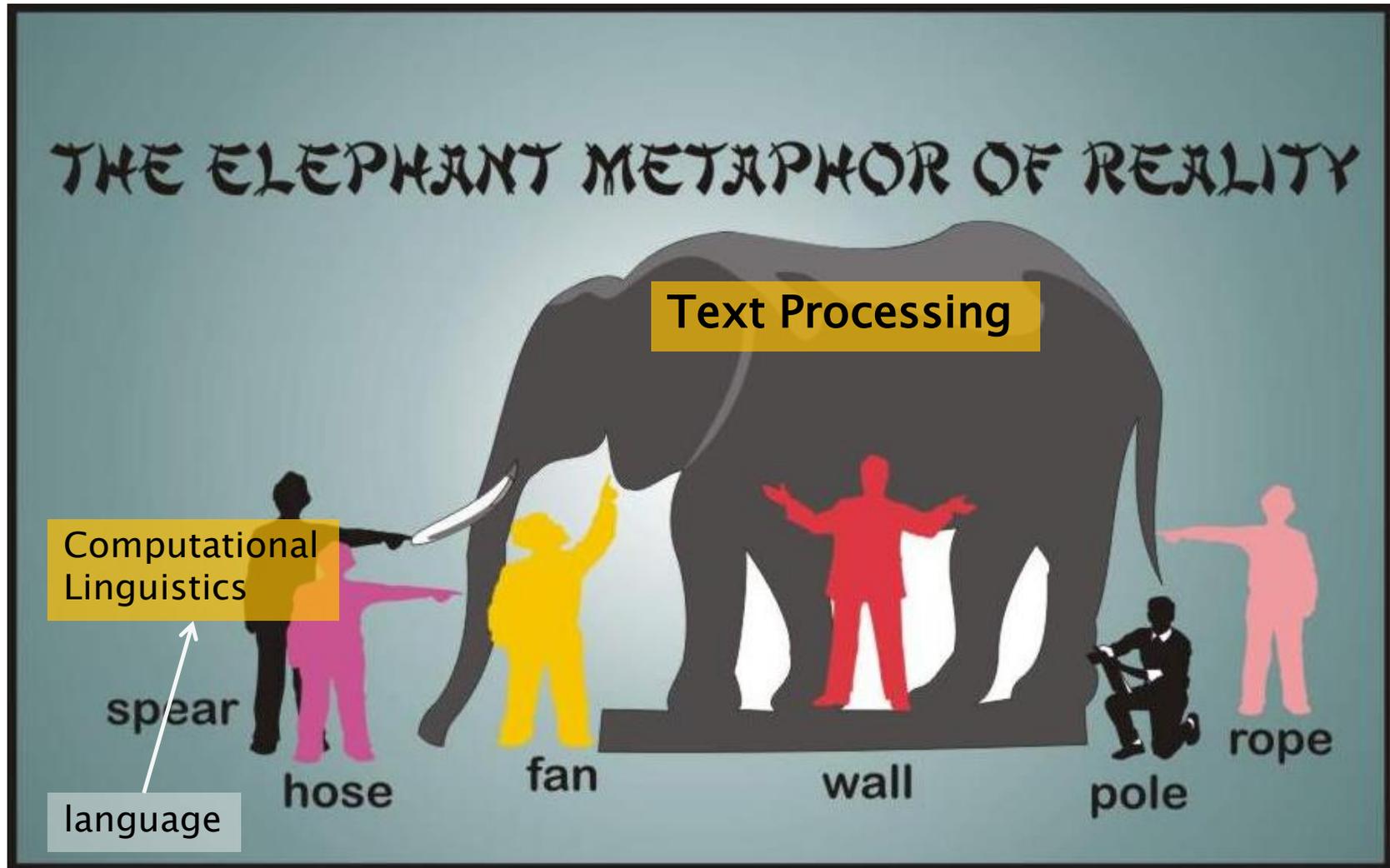
How different research areas approach text?



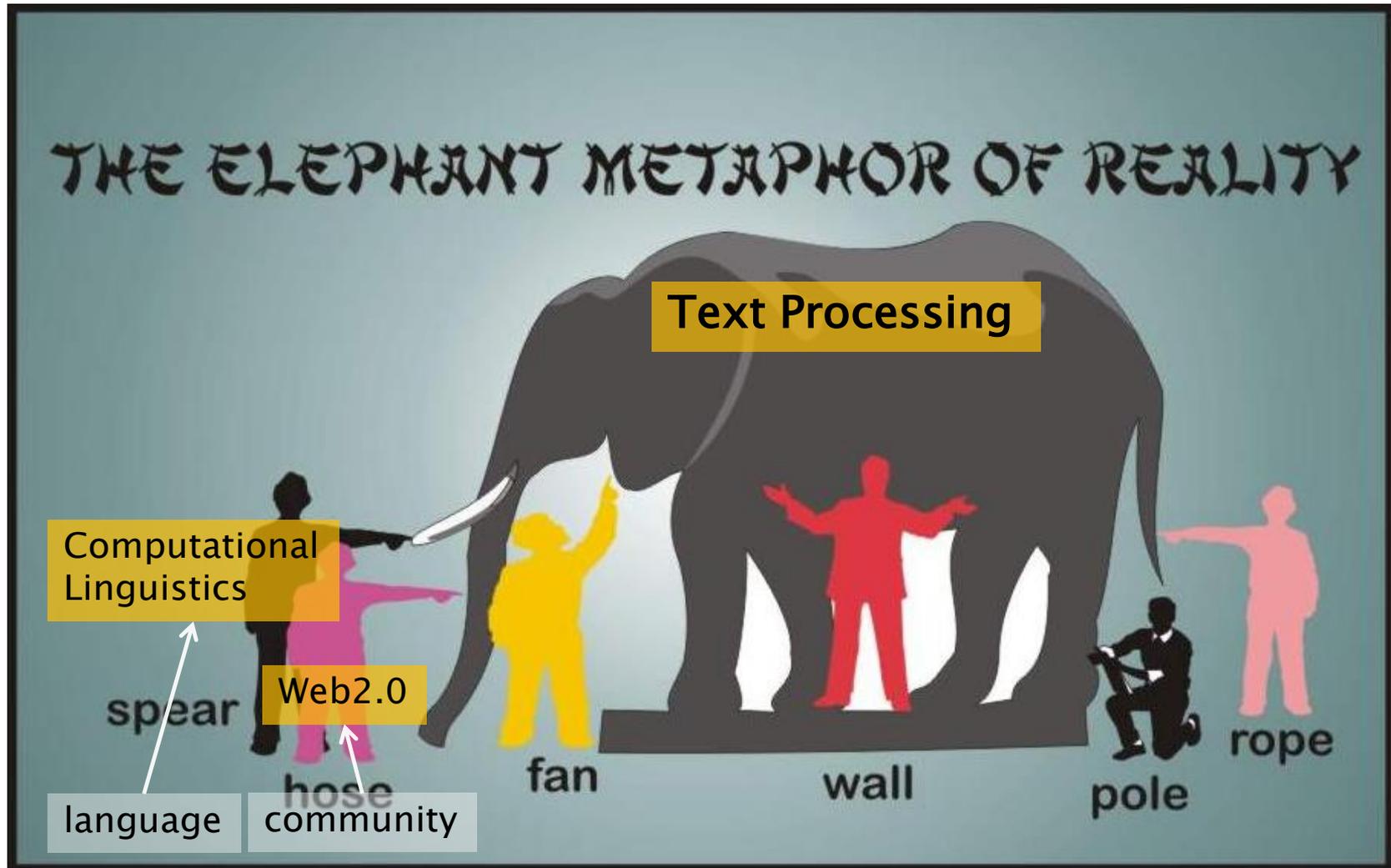
How different research areas approach text?



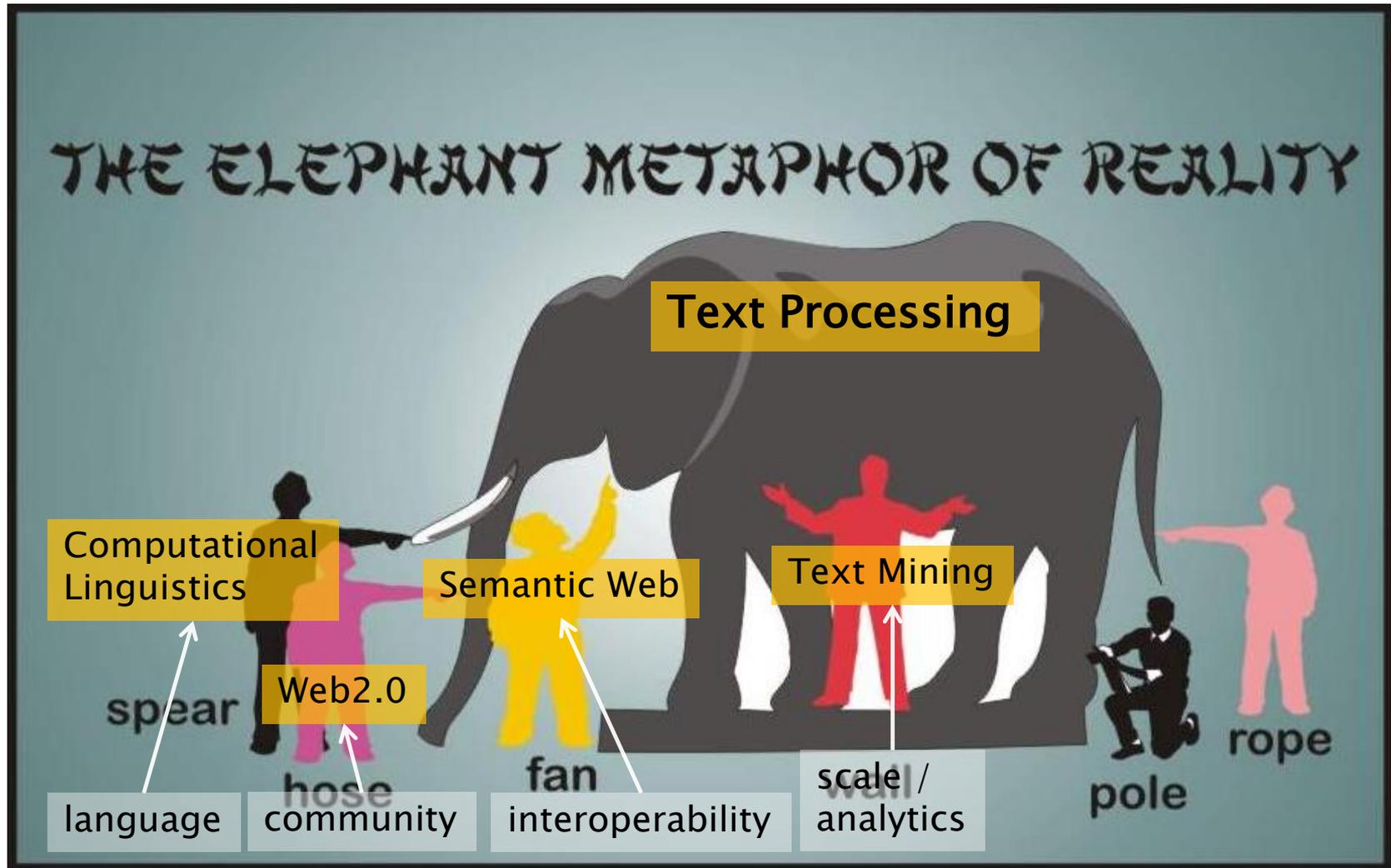
How different research areas approach text?



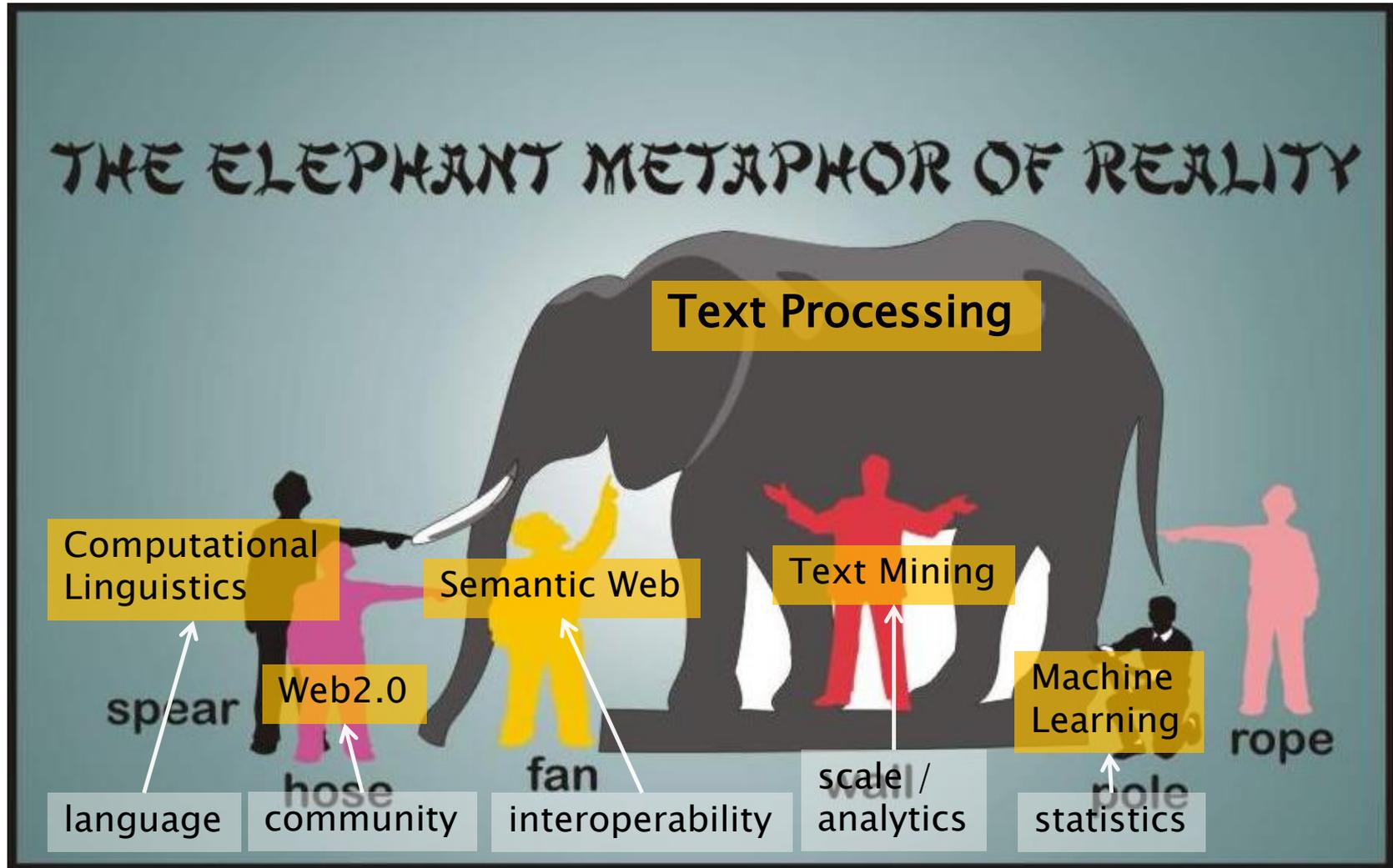
How different research areas approach text?



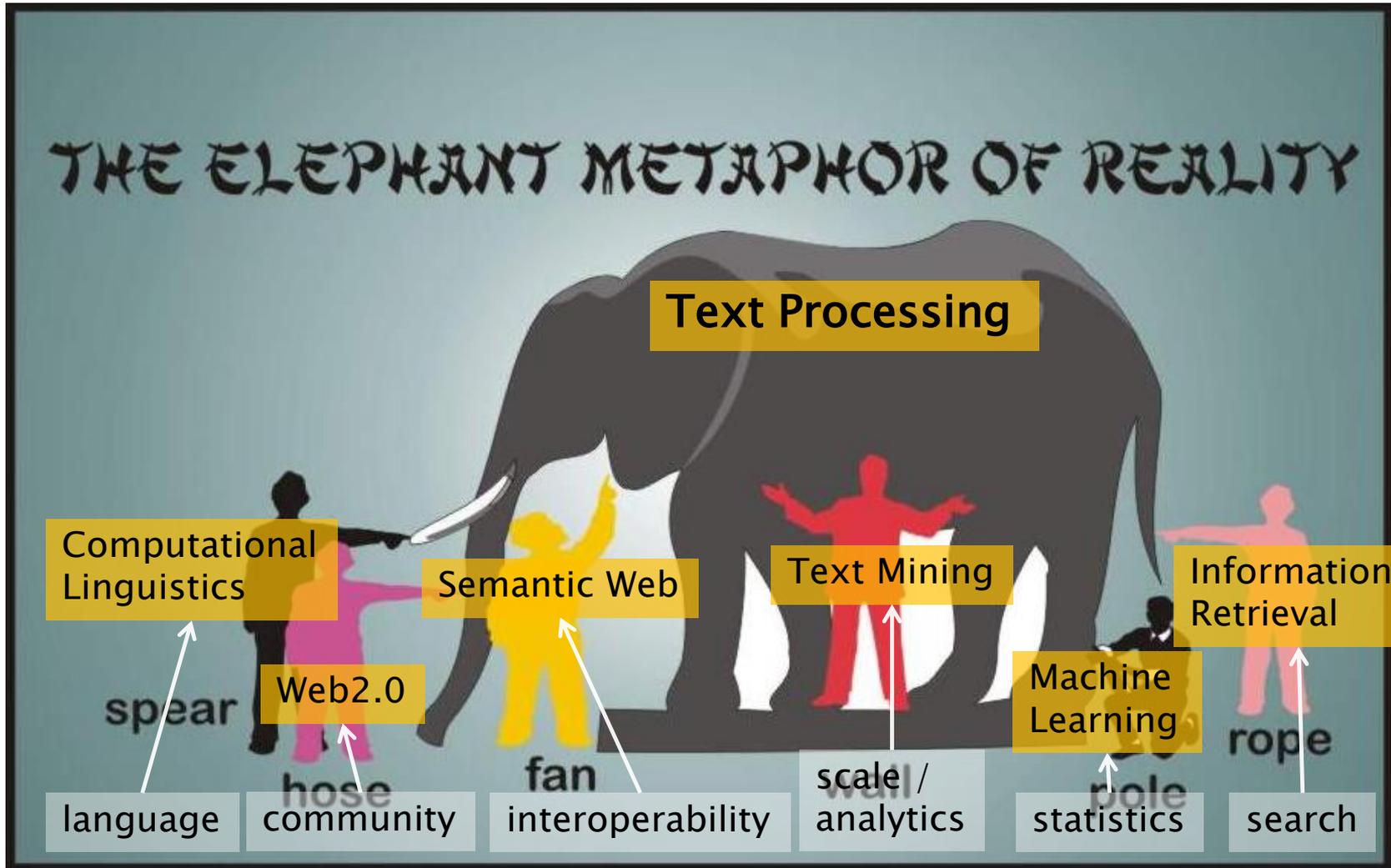
How different research areas approach text?



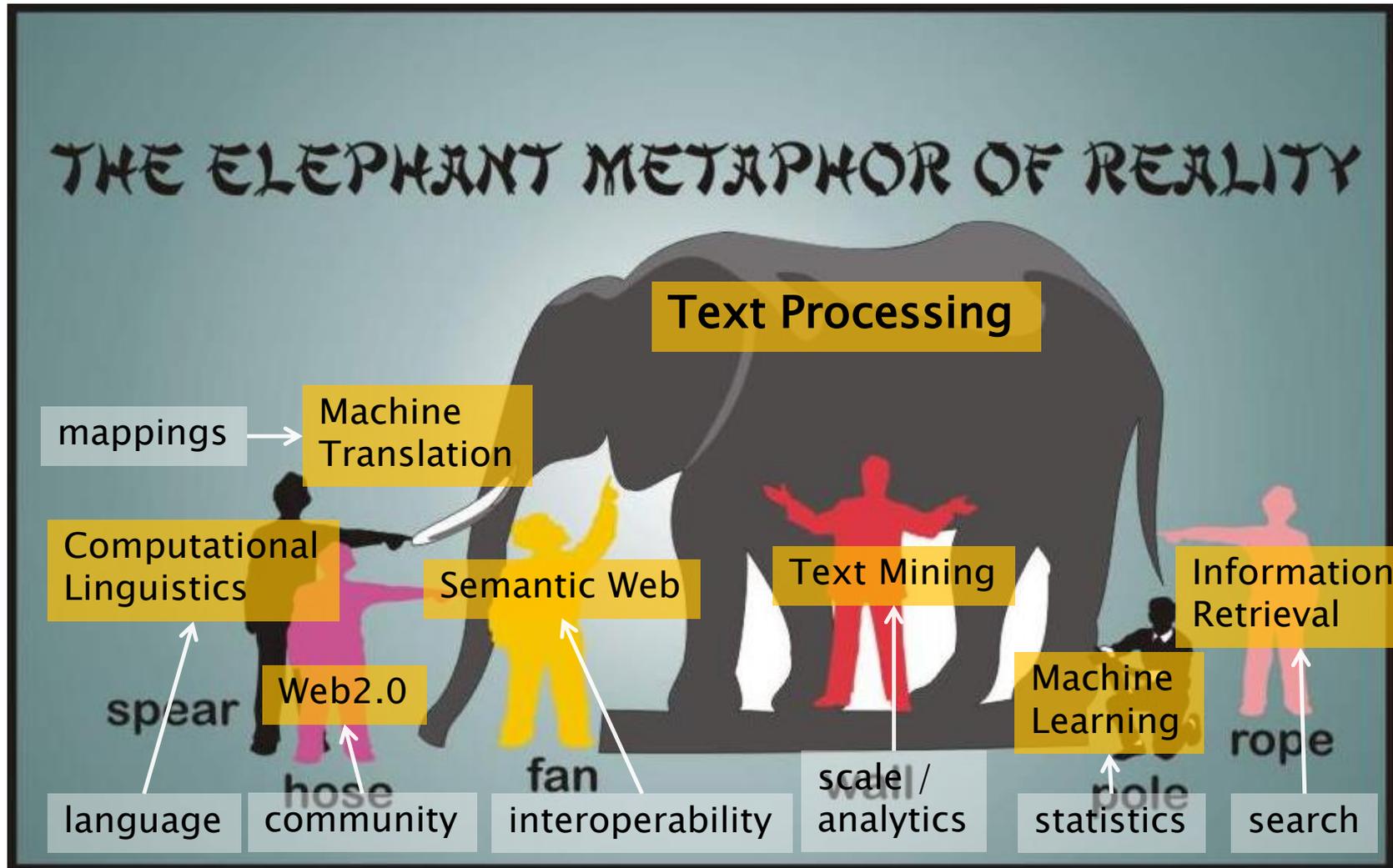
How different research areas approach text?



How different research areas approach text?

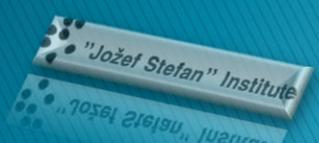


How different research areas approach text?



How do we represent text?

Levels of representation



Levels of text representations

- ▶ Character (character n-grams and sequences)
- ▶ Words (stop-words, stemming, lemmatization)
- ▶ Phrases (word n-grams, proximity features)
- ▶ Part-of-speech tags
- ▶ Taxonomies / thesauri

Lexical

-
- ▶ Vector-space model
 - ▶ Language models
 - ▶ Full-parsing
 - ▶ Cross-modality

Syntactic

-
- ▶ Collaborative tagging / Web2.0
 - ▶ Linked Data
 - ▶ Templates / Frames
 - ▶ Ontologies / First order theories

Semantic

Levels of text representations

Language
identification,
Copy detection

- ▶ Character (character n-grams and sequences)
- ▶ Words (stop-words, stemming, lemmatization)
- ▶ Phrases (word n-grams, proximity features)
- ▶ Part-of-speech tags
- ▶ Taxonomies / thesauri

Lexical

-
- ▶ Vector-space model
 - ▶ Language models
 - ▶ Full-parsing
 - ▶ Cross-modality

Syntactic

-
- ▶ Collaborative tagging / Web2.0
 - ▶ Linked Data
 - ▶ Templates / Frames
 - ▶ Ontologies / First order theories

Semantic

Levels of text representations

Language identification, Copy detection

- ▶ Character (character n-grams and sequences)
- ▶ Words (stop-words, stemming, lemmatization)
- ▶ Phrases (word n-grams, proximity features)
- ▶ Part-of-speech tags
- ▶ Taxonomies / thesauri

Named entity extraction (names of people, places, organizations)

-
- ▶ Vector-space model
 - ▶ Language models
 - ▶ Full-parsing
 - ▶ Cross-modality

Syntactic

-
- ▶ Collaborative tagging / Web2.0
 - ▶ Linked Data
 - ▶ Templates / Frames
 - ▶ Ontologies / First order theories

Semantic

Levels of text representations

- ▶ Character (character n-grams and sequences)
- ▶ Words (stop-words, stemming, lemmatization)
- ▶ Phrases (word n-grams, proximity features)
- ▶ Part-of-speech tags
- ▶ Taxonomies / thesauri

- ▶ Vector-space model
- ▶ Language models
- ▶ Full-parsing
- ▶ Cross-modality

Language identification,
Copy detection

Named entity extraction
(names of people, places,
organizations)

Text categorization,
Clustering, Search,
Summarization, ...

Syntactic

- ▶ Collaborative tagging / Web2.0
- ▶ Linked Data
- ▶ Templates / Frames
- ▶ Ontologies / First order theories

Semantic

Levels of text representations

- ▶ Character (character n-grams and sequences)
- ▶ Words (stop-words, stemming, lemmatization)
- ▶ Phrases (word n-grams, proximity features)
- ▶ Part-of-speech tags
- ▶ Taxonomies / thesauri

- ▶ Vector-space model
- ▶ Language models
- ▶ Full-parsing
- ▶ Cross-modality

Language identification,
Copy detection

Named entity extraction
(names of people, places,
organizations)

Text categorization,
Clustering, Search,
Summarization, ...

Spam filtering,
Machine translation

Syntactic

-
- ▶ Collaborative tagging / Web2.0
 - ▶ Linked Data
 - ▶ Templates / Frames
 - ▶ Ontologies / First order theories

Semantic

Levels of text representations

- ▶ Character (character n-grams and sequences)
- ▶ Words (stop-words, stemming, lemmatization)
- ▶ Phrases (word n-grams, proximity features)
- ▶ Part-of-speech tags
- ▶ Taxonomies / thesauri

- ▶ Vector-space model
- ▶ Language models
- ▶ Full-parsing
- ▶ Cross-modality

- ▶ Collaborative tagging / Web2.0
- ▶ Linked Data
- ▶ Templates / Frames
- ▶ Ontologies / First order theories

Language identification,
Copy detection

Named entity extraction
(names of people, places,
organizations)

Text categorization,
Clustering, Search,
Summarization, ...

Spam filtering,
Machine translation

Multilingual search,
Associating text with
images, ...

Syntactic

Semantic

Levels of text representations

- ▶ Character (character n-grams and sequences)
- ▶ Words (stop-words, stemming, lemmatization)
- ▶ Phrases (word n-grams, proximity features)
- ▶ Part-of-speech tags
- ▶ Taxonomies / thesauri

- ▶ Vector-space model
- ▶ Language models
- ▶ Full-parsing
- ▶ Cross-modality

- ▶ Collaborative tagging / Web2.0
- ▶ Linked Data
- ▶ Templates / Frames
- ▶ Ontologies / First order theories

Language identification,
Copy detection

Named entity extraction
(names of people, places,
organizations)

Text categorization,
Clustering, Search,
Summarization, ...

Spam filtering,
Machine translation

Multilingual search,
Associating text with
images, ...

Unifying
semantics
of data

Syntactic

Semantic

Levels of text representations

- ▶ Character (character n-grams and sequences)
- ▶ Words (stop-words, stemming, lemmatization)
- ▶ Phrases (word n-grams, proximity features)
- ▶ Part-of-speech tags
- ▶ Taxonomies / thesauri

- ▶ Vector-space model
- ▶ Language models
- ▶ Full-parsing
- ▶ Cross-modality

- ▶ Collaborative tagging / Web2.0
- ▶ Linked Data
- ▶ Templates / Frames
- ▶ Ontologies / First order theories

Language identification,
Copy detection

Named entity extraction
(names of people, places,
organizations)

Text categorization,
Clustering, Search,
Summarization, ...

Spam filtering,
Machine translation

Multilingual search,
Associating text with
images, ...

Data integration

Unifying
semantics
of data

Syntactic

Semantic

Levels of text representations

- ▶ Character (character n-grams and sequences)
- ▶ Words (stop-words, stemming, lemmatization)
- ▶ Phrases (word n-grams, proximity features)
- ▶ Part-of-speech tags
- ▶ Taxonomies / thesauri

- ▶ Vector-space model
- ▶ Language models
- ▶ Full-parsing
- ▶ Cross-modality

- ▶ Collaborative tagging / Web2.0
- ▶ Linked Data
- ▶ Templates / Frames
- ▶ Ontologies / First order theories

Language identification,
Copy detection

Named entity extraction
(names of people, places, organizations)

Text categorization,
Clustering, Search,
Summarization, ...

Spam filtering,
Machine translation

Syntactic

Multilingual search,
Associating text with
images, ...

Data integration

Unifying semantics
of data

Semantic

Reasoning,
Semantic Search



Levels of text representations

- ▶ Character (character n-grams and sequences)
- ▶ Words (stop-words, stemming, lemmatization)
- ▶ Phrases (word n-grams, proximity features)
- ▶ Part-of-speech tags
- ▶ Taxonomies / thesauri

Lexical

-
- ▶ Vector-space model
 - ▶ Language models
 - ▶ Full-parsing
 - ▶ Cross-modality

Syntactic

-
- ▶ Collaborative tagging / Web2.0
 - ▶ Linked Data
 - ▶ Templates / Frames
 - ▶ Ontologies / First order theories

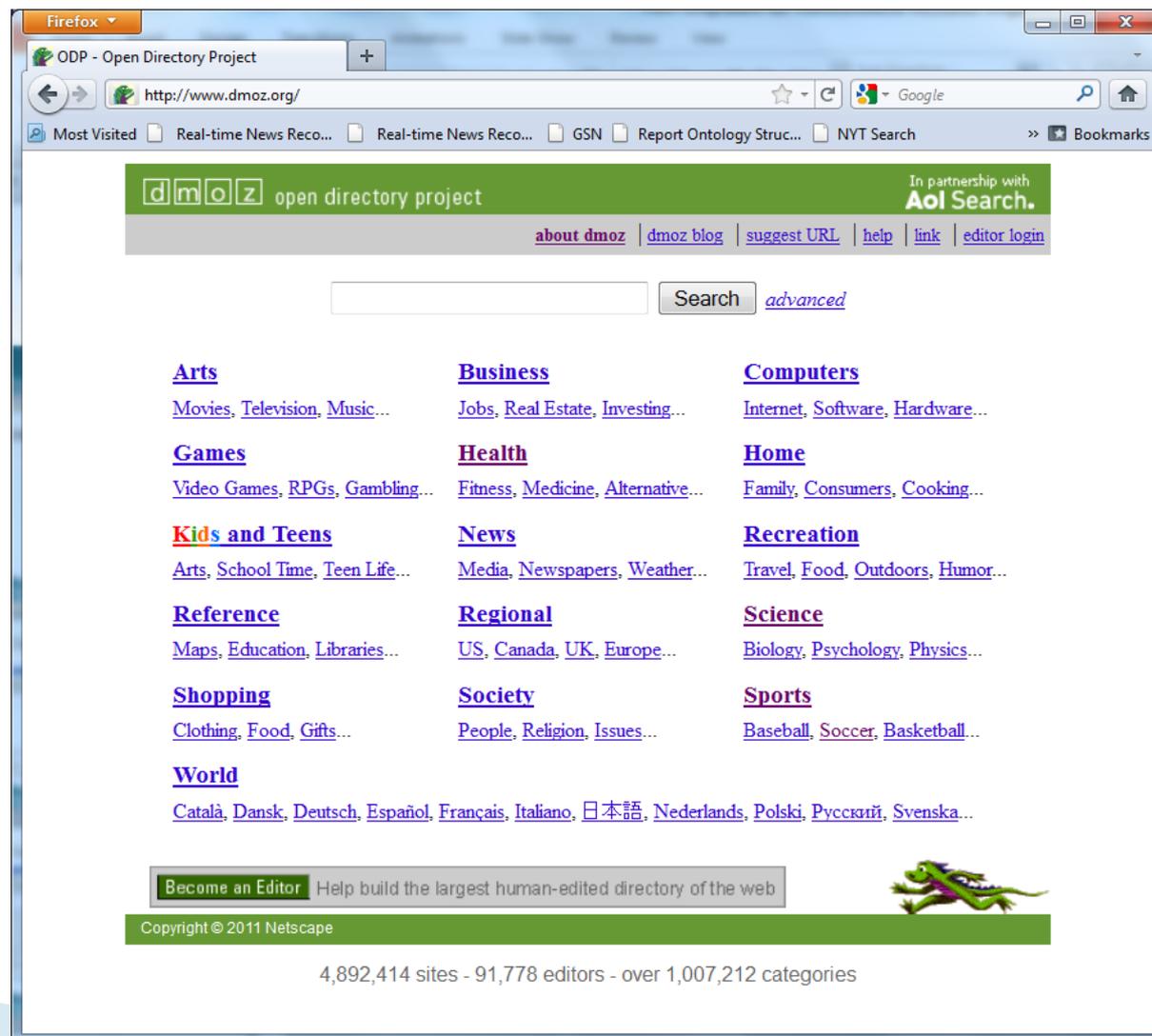
Semantic

Document Classification into large taxonomies

DMoz (Open Directory Project)

<http://dmoz.org>

- ▶ Largest handcrafted taxonomy on the Web
 - 4,892,414 sites
 - 91,778 editors
 - over 1,007,212 categories
- ▶ Data available for download
 - <http://www.dmoz.org/rdf.html>



The screenshot shows the DMOZ website in a Firefox browser window. The address bar displays "http://www.dmoz.org/". The page header includes the DMOZ logo and the text "open directory project". A search bar is present with a "Search" button and a link to "advanced". The main content area is organized into a grid of category links, including Arts, Business, Computers, Games, Health, Home, Kids and Teens, News, Recreation, Reference, Regional, Science, Shopping, Society, Sports, and World. Each category link is followed by a list of sub-topics. At the bottom of the page, there is a "Become an Editor" button and a copyright notice for 2011 Netscape. A small green dragon logo is visible in the bottom right corner.

4,892,414 sites - 91,778 editors - over 1,007,212 categories

Classification of a query into DMoz

Mozilla Firefox

File Edit View History Bookmarks Tools Help

http://localhost:8080/

Most Visited Getting Started Latest Headlines Customize Links Windows Marketplace

Classification into DMoz Top_Science

URL:

Text:

Categories:

Keyword Treshold (0=Non, 1=All):

Context:

XML Format:

Plain-Text Mime-Type:

Done



Classification - Mozilla Firefox

File Edit View History Bookmarks Tools Help

http://localhost:8080/Classify

Most Visited Getting Started Latest Headlines Customize Links Windows Marketplace

Results of Classification into DMoz Top_Science

Keywords:

- Science (0.183)
- Math (0.183)
- Number_Theory (0.182)
- Diophantine_Equations (0.154)
- Fermat's_Last_Theorem (0.115)

Categories:

1	<input type="checkbox"/>	0.504	Top/Science/Math/Number_Theory/Diophantine_Equations/Fermat's_Last_Theorem
2	<input type="checkbox"/>	0.341	Top/Science/Math/Number_Theory/Diophantine_Equations
3	<input type="checkbox"/>	0.117	Top/Science/Math/Number_Theory/Tables
4	<input type="checkbox"/>	0.113	Top/Science/Math/Number_Theory/History
5	<input type="checkbox"/>	0.074	Top/Science/Math/Number_Theory/Factoring
6	<input type="checkbox"/>	0.064	Top/Science/Math/Number_Theory/Factoring/Tables
7	<input type="checkbox"/>	0.053	Top/Science/Math/Number_Theory
8	<input type="checkbox"/>	0.053	Top/Science/Math/Number_Theory/Prime_Numbers/Primality_Tests/Pseudoprimes
9	<input type="checkbox"/>	0.037	Top/Science/Math/History/People
10	<input type="checkbox"/>	0.035	Top/Science/Math/Number_Theory/Prime_Numbers/Mersenne
11	<input type="checkbox"/>	0.027	Top/Science/Math/Number_Theory/Publications/Books

Done

Classification of a document into DMoz

Classification into DMoz Top_Science

URL:

Text:
international class scientific research. With this in mind, the in-house research has been reinforced by building strong links to universities, other research institutions and industry. The Institute is closely connected with the Slovenian universities, where many scientists who initially developed their research talents at the Institute have been appointed to teaching posts, while still retaining their research positions or research teams at the Institute. Since 1985 more than 800 postgraduate students have gained their MSc. and Ph.D. degrees at the Institute. Close contacts are also maintained with secondary schools, providing work practice on research projects in natural sciences and organising regular visits to the laboratories.

Categories: 25

Keyword Treshold (0=Non, 1=All): 0.75

Context:

XML Format:

Plain-Text Mime-Type:

Done

Classification - Mozilla Firefox

http://localhost:8080/Classify

Classification

Results of Classification into DMoz Top_Science

Keywords:

- Science (0.085)
- Institutions (0.069)
- Research_Centers (0.048)
- Social_Sciences (0.043)
- Institutes (0.041)
- Biology (0.039)
- Oceanography (0.038)
- Earth_Sciences (0.038)
- Research_Institutes (0.038)
- Regional (0.038)
- Research (0.036)
- Europe (0.034)
- Organizations (0.031)
- Math (0.031)
- Ecology (0.026)
- Environment (0.022)
- Associations (0.020)
- Economics (0.020)
- Agriculture (0.019)
- Greece (0.019)
- Plasma (0.018)

Categories:

1	<input type="checkbox"/>	0.387	Top_Science/Institutions/Research_Institutes
2	<input type="checkbox"/>	0.344	Top_Science/Biology/Institutions/Research_Centers
3	<input type="checkbox"/>	0.313	Top_Science/Institutions/Regional
4	<input type="checkbox"/>	0.309	Top_Science/Environment/Organizations/Research_Institutes
5	<input type="checkbox"/>	0.293	Top_Science/Math/Research/Institutes
6	<input type="checkbox"/>	0.289	Top_Science/Institutions/Associations
7	<input type="checkbox"/>	0.281	Top_Science/Math/Research
8	<input type="checkbox"/>	0.278	Top_Science/Social_Sciences/Economics/Institutes
9	<input type="checkbox"/>	0.276	Top_Science/Agriculture/Research_Centers
10	<input type="checkbox"/>	0.270	Top_Science/Institutions/Regional/Europe

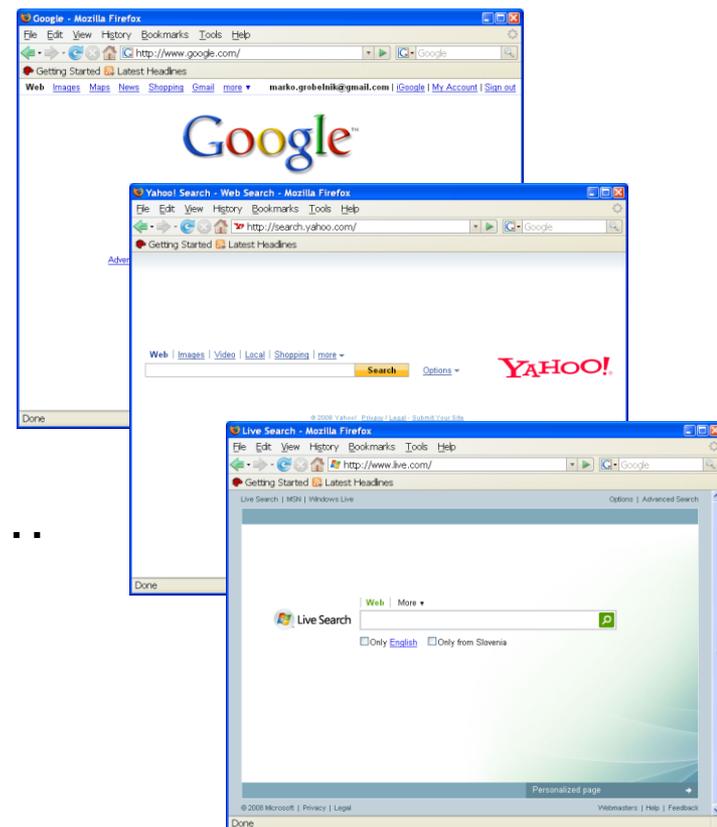
Done

Visual & Contextual Search

(WWW2008)

Contextualized search

- ▶ What is the most common tasks where we manipulate text in everyday life?
 - “Internet search”!
- ▶ ...but – how smart is search technology today?
 - ...not too smart!
 - It is sophisticated, but not smart...



Example: searching for “Jaguar”

- ▶ Query “jaguar” has many meanings...
- ▶ ...but the first page of search engines doesn't provide us with many answers

A screenshot of a Mozilla Firefox browser window showing a Google search for "jaguar". The browser title is "jaguar - Google Search - Mozilla Firefox". The address bar shows the URL "http://www.google.com/search?hl=en&q=jaguar". The search bar contains the word "jaguar" and a "Search" button. Below the search bar, there are navigation links for "Web", "Images", "Maps", "News", "Shopping", and "Gmail". The search results are displayed under the heading "Web" and show the first 10 results for "jaguar" (definition) in 0.05 seconds. The results include:

- Jaguar**: Official worldwide web site of Jaguar Cars. Directs users to pages tailored to country-specific markets and model-specific websites. www.jaguar.com/ - Similar pages - Note this
- Jaguar UK - Jaguar Cars**: Jaguar XF. TEST DRIVE. Brochure. Dealer. eNewsletter. SEARCH SITEMAP COMPANY Privacy Policy Accessibility Statement Contact Us TERMS & CONDITIONS ... www.jaguar.com/uk/ - 17k - Cached - Similar pages - Note this
- Jaguar US - Home**: Jaguar USA official website. ... Build Your Jaguar. Request Brochure. Get Email Updates. Locate a Dealer. Search Your Profile Site Map Contact Us Privacy ... www.jaguarusa.com/ - 20k - Cached - Similar pages - Note this
- Jaguar - Wikipedia, the free encyclopedia**: The jaguar (Panthera onca, pronounced /ˈdʒæɡjuːr/ in British English, or /ˈdʒæɡwɑːr/ in American English) is a New World mammal of the Felidae family and one ... en.wikipedia.org/wiki/Jaguar - 153k - Cached - Similar pages - Note this
- Jaguar Cars**: English · Français. www.jaguar.ca/ - 4k - Cached - Similar pages - Note this

The browser status bar at the bottom shows "Done".

Example: searching for “Jaguar”

- ▶ Query “jaguar” has many meanings...
- ▶ ...but the first page of search engines doesn't provide us with many answers
- ▶ ...there are 84M more results

jaguar - Google Search - Mozilla Firefox

File Edit View History Bookmarks Tools Help

http://www.google.com/search?hl=en&q=jaguar

Getting Started Latest Headlines

Web Images Maps News Shopping Gmail more

marko.grobelnik@gmail.com | Web History | My Account | Sign out

Google jaguar Search Advanced Search Preferences

Web Images Results 1 - 10 of about 84,200,000 for jaguar [definition]. (0.05 seconds)

Jaguar
Official worldwide web site of **Jaguar** Cars. Directs users to pages tailored to country-specific markets and model-specific websites.
www.jaguar.com/ - Similar pages - Note this

Jaguar UK - Jaguar Cars
Jaguar XF. TEST DRIVE. Brochure. Dealer. eNewsletter. SEARCH SITEMAP COMPANY
Privacy Policy Accessibility Statement Contact Us TERMS & CONDITIONS ...
www.jaguar.com/uk/ - 17k - Cached - Similar pages - Note this
More results from www.jaguar.com »

Jaguar US - Home
Jaguar USA official website. ... Build Your **Jaguar**. Request Brochure. Get Email Updates.
Locate a Dealer. Search Your Profile Site Map Contact Us Privacy ...
www.jaguarusa.com/ - 20k - Cached - Similar pages - Note this

Jaguar - Wikipedia, the free encyclopedia
The **jaguar** (*Panthera onca*, pronounced /ˈdʒæɡjuːr/ in British English, or /ˈdʒæɡwɑːr/ in American English) is a New World mammal of the Felidae family and one ...
en.wikipedia.org/wiki/Jaguar - 153k - Cached - Similar pages - Note this

Jaguar Cars
English · Français.
www.jaguar.ca/ - 4k - Cached - Similar pages - Note this

Done

Context sensitive search with <http://searchpoint.ijs.si>

The screenshot shows a search engine interface with the following elements:

- Search Bar:** Contains the query "jaguar".
- Search Buttons:** "Search via topics", "Search via query to ontology", and "Search via hits to ontology".
- Search Results:** A list of results for "jaguar", including:
 - (9) [Jaguar](#): General information and facts from Big Cats Online.
 - (59) [Jaguar, Jaguar Profile, Facts, Information, Photos, Pictures ...](#): Get jaguar profile, facts, information, photos, pictures, sounds, habitats, reports, news, and more from National Geographic.
 - (8) [Jaguar - Wikipedia, the free encyclopedia](#): The jaguar (Panthera onca) is a New World mammal of the Felidae family and one of four "big cats" in the Panthera genus, along with the tiger, ...
 - (11) [Jaguar](#): Jaguar Facts, Jaguar Photos and Jaguars in the news at the world's largest big cat rescue and sanctuary.
 - (1) [Jaguar](#): Official worldwide web site of Jaguar Cars. Directs users to pages tailored to country-specific markets.
 - (32) [Jaguar](#): Contains extensive information about the Jaguar. Information includes habitat, body size, and life span.
 - (2) [Jaguar UK - Jaguar Cars](#): Jaguar & Ownership. Highlights. Gallery. Models & Pricing. Design Your XK. TEST DRIVE. Brochure. Dealer. eNewsletter ...
 - (17) [Jaguar Enthusiasts' Club](#): World's largest audited membership. UK-based, JEC's site has extensive resources available for the enthusiast, including information about their Sections, ...
 - (20) [San Diego Zoo's Animal Bytes: Jaguar](#): Get fun and interesting jaguar facts in an easy-to-read style from the San Diego Zoo's Animal

On the right side of the interface, there is a **Conceptual map** showing a network of related terms. A red dot is placed on the "Mammalia" node, and a black arrow points from the search bar to this node. Other nodes include "Vehicles", "Shopping", "Sports", "Games", "Console Platforms", "Aviation", "Society", "Recreation", "Enthusiasts", "Aircraft", "Models", "Makes and", "Parts and Accessories", and "NFL".

Query

Conceptual map

Search Point

Dynamic contextual ranking based on the search point



News reporting bias

(Fortuna, Galleguillos, Cristianini 2008)

News Reporting Bias example

UK SOLDIERS CLEARED IN IRAQI DEATH – SEVEN BRITISH SOLDIERS WERE ACQUITTED ON THURSDAY OF CHARGES OF BEATING AN INNOCENT IRAQI TEENAGER TO DEATH WITH RIFLE BUTTS. A JUDGE AT A SPECIALLY CONVENED MILITARY COURT IN EASTERN ENGLAND ORDERED THE ADJUDICATING PANEL TO RETURN 'NOT GUILTY' VERDICTS AGAINST THE SEVEN BECAUSE HE DID NOT BELIEVE THERE WAS SUFFICIENT EVIDENCE AGAINST THEM, THE MINISTRY OF DEFENCE SAID. . . .

BRITISH MURDERERS IN IRAQ ACQUITTED – THE JUDGE AT A COURT-MARTIAL ON THURSDAY DISMISSED MURDER CHARGES AGAINST SEVEN SOLDIERS, FROM THE 3RD BATTALION, THE PARACHUTE REGIMENT, WHO'RE ACCUSED OF MURDERING IRAQI TEENAGER; CLAIMING THERE'S INSUFFICIENT EVIDENCE TO SECURE A CONVICTION, THE ASSOCIATED PRESS REPORTED THURSDAY. . . .

Experimental setup

- ▶ Time period: March 31st 2005 – April 14th 2006
- ▶ Size of collections:

Source	No. of news
Al Jazeera	2142
CNN	6840
Detroit News	2929
International Herald Tribune	9641

- ▶ Number of discovered matches:

	AJ	CNN	DN	IHT
AJ	–	816	447	834
CNN	816	–	1103	2437
DN	447	1103	–	895
IHT	834	2437	895	–

Prediction of news source

- ▶ **The task:** given a pair of news articles describing the same event, can we predict the news source for each?
- ▶ In this experiment we focused on CNN and Al Jazeera.
- ▶ SVM linear classifier was used for prediction
 - Evaluation was done using 10-fold cross-validation
 - Significance of results was tested against random matches

Detecting News Reporting Bias

- ▶ The task:
 - Given a news story, are we able to say from which news source it came?



Detecting News Reporting Bias

- ▶ The task:
 - Given a news story, are we able to say from which news source it came?
- ▶ We compared **CNN** and **Aljazeera** reports about the same events from the war in Iraq
 - ...300 aligned articles describing the same story from both sources



Detecting News Reporting Bias

- ▶ The task:
 - Given a news story, are we able to say from which news source it came?
- ▶ We compared **CNN** and **Aljazeera** reports about the same events from the war in Iraq
 - ...300 aligned articles describing the same story from both sources
- ▶ The same topics are expressed in both sources with the following keywords:
 - CNN with:
 - **Insurgents**, Troops, Baghdad, Iran, **Militant**, Police, **Suicide**, **Terrorist**, United, National, Hussein, **Alleged**, Israeli, Syria, Terrorism...
 - Aljazeera with:
 - Attacks, Claims, **Rebels**, Withdrawing, Report, **Fighters**, President, **Resistance**, Occupation, Injured, Army, Demanded, Hit, Muslim, ...

News Visualization

Topic landscape of the query “Clinton” from Reuters news 1996–1997

The screenshot displays the News Analyser interface. On the left, a search bar contains the query 'clinton'. Below it is a list of search results with columns for 'Date' and 'Title'. A tooltip is visible over the topic map, showing a list of terms: USA: U.S. will attend ..., #documents = 245, NATO, PALESTINIAN, ISRA, PEAC, ISRAEL, NETANYAHU, YELTSIN, ARAFAT, RUSSIA, SUMMIT. The topic map itself is a word cloud where the size of each word represents its frequency. A selected story is shown at the bottom right.

Query

Search Results

Topic Map

Selected group of news

Selected story

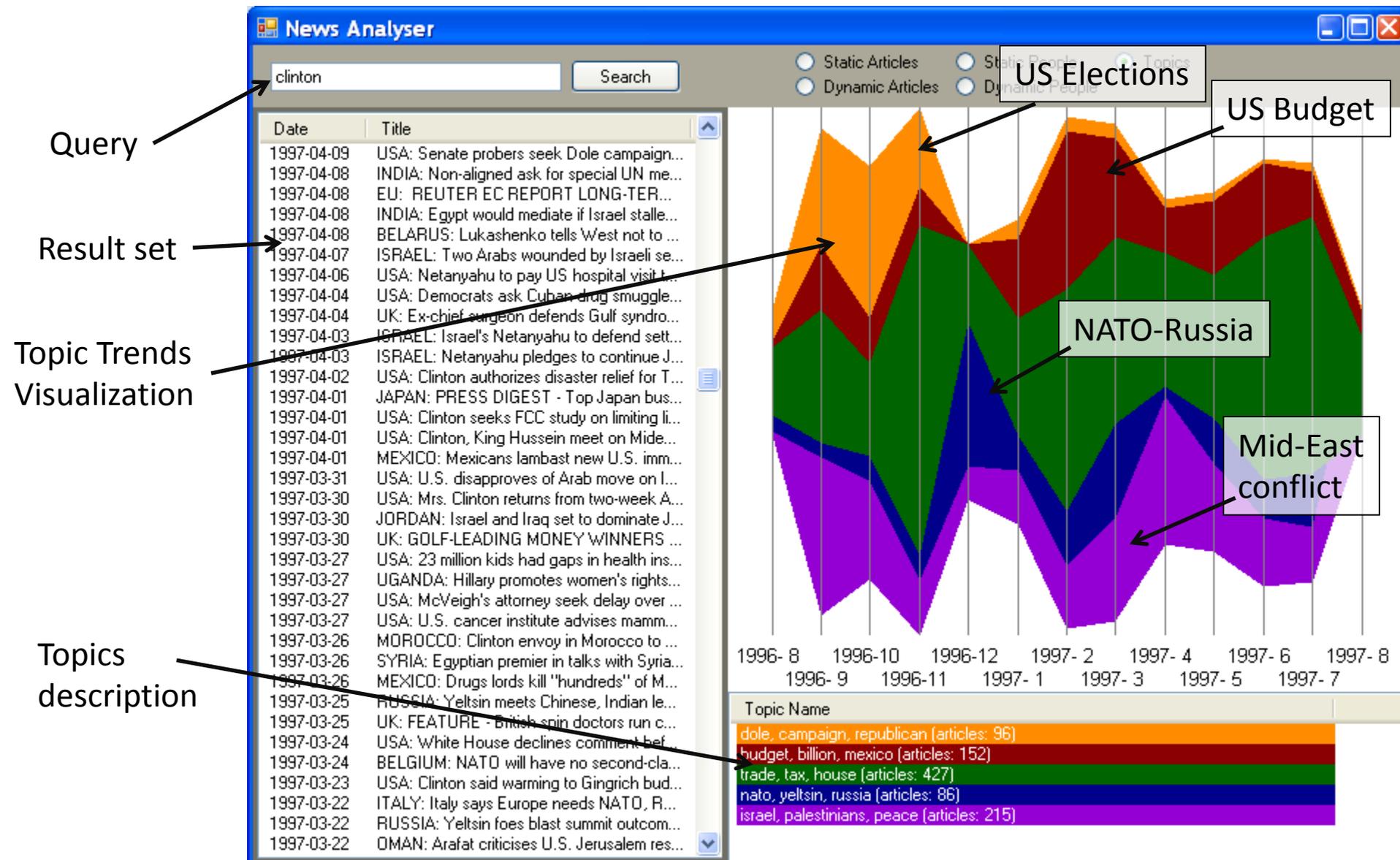
Date	Title
1997-04-09	USA: Senate probers seek Dole campaign...
1997-04-08	INDIA: Non-aligned ask for special UN me...
1997-04-08	EU: REUTER EC REPORT LONG-TER...
1997-04-08	INDIA: Egypt would mediate if Israel stalle...
1997-04-08	BELARUS: Lukashenko tells West not to...
1997-04-07	ISRAEL: Two Arabs wounded by Israeli se...
1997-04-06	USA: Netanyahu to pay US hospital visit t...
1997-04-04	USA: Democrats ask Cuban drug smuggle...
1997-04-04	UK: Ex-chief surgeon defends Gulf syndro...
1997-04-03	ISRAEL: Israel's Netanyahu to defend sett...
1997-04-03	ISRAEL: Netanyahu pledges to continue J...
1997-04-02	USA: Clinton authorizes disaster relief for T...
1997-04-01	JAPAN: PRESS DIGEST - Top Japan bus...
1997-04-01	USA: Clinton seeks FCC study on limiting li...
1997-04-01	USA: Clinton, King Hussein meet on Mide...
1997-04-01	MEXICO: Mexicans lambast new U.S. imm...
1997-03-31	USA: U.S. disapproves of Arab move on l...
1997-03-30	USA: Mrs. Clinton returns from two week A...
1997-03-30	JORDAN: Israel and Iraq set to dominate J...
1997-03-30	UK: GOLF-LEADING MONEY WINNERS ...
1997-03-27	USA: 23 million kids had gaps in health ins...
1997-03-27	UGANDA: Hillary promotes women's rights...
1997-03-27	USA: McVeigh's attorney seek delay over ...
1997-03-27	USA: U.S. cancer institute advises mamm...
1997-03-26	MOROCCO: Clinton envoy in Morocco to ...
1997-03-26	SYRIA: Egyptian premier in talks with Syria...
1997-03-26	MEXICO: Drugs lords kill "hundreds" of M...
1997-03-25	RUSSIA: Yeltsin meets Chinese, Indian le...
1997-03-25	UK: FEATURE - British spin doctors run c...
1997-03-24	USA: White House declines comment bef...
1997-03-24	BELGIUM: NATO will have no second-cla...
1997-03-23	USA: Clinton said warning to Gingrich bud...
1997-03-22	ITALY: Italy says Europe needs NATO, R...
1997-03-22	RUSSIA: Yeltsin foes blast summit outcom...
1997-03-22	OMAN: Arafat criticises U.S. Jerusalem res...

Static Articles **Static People** **Topics**
Dynamic Articles **Dynamic People**

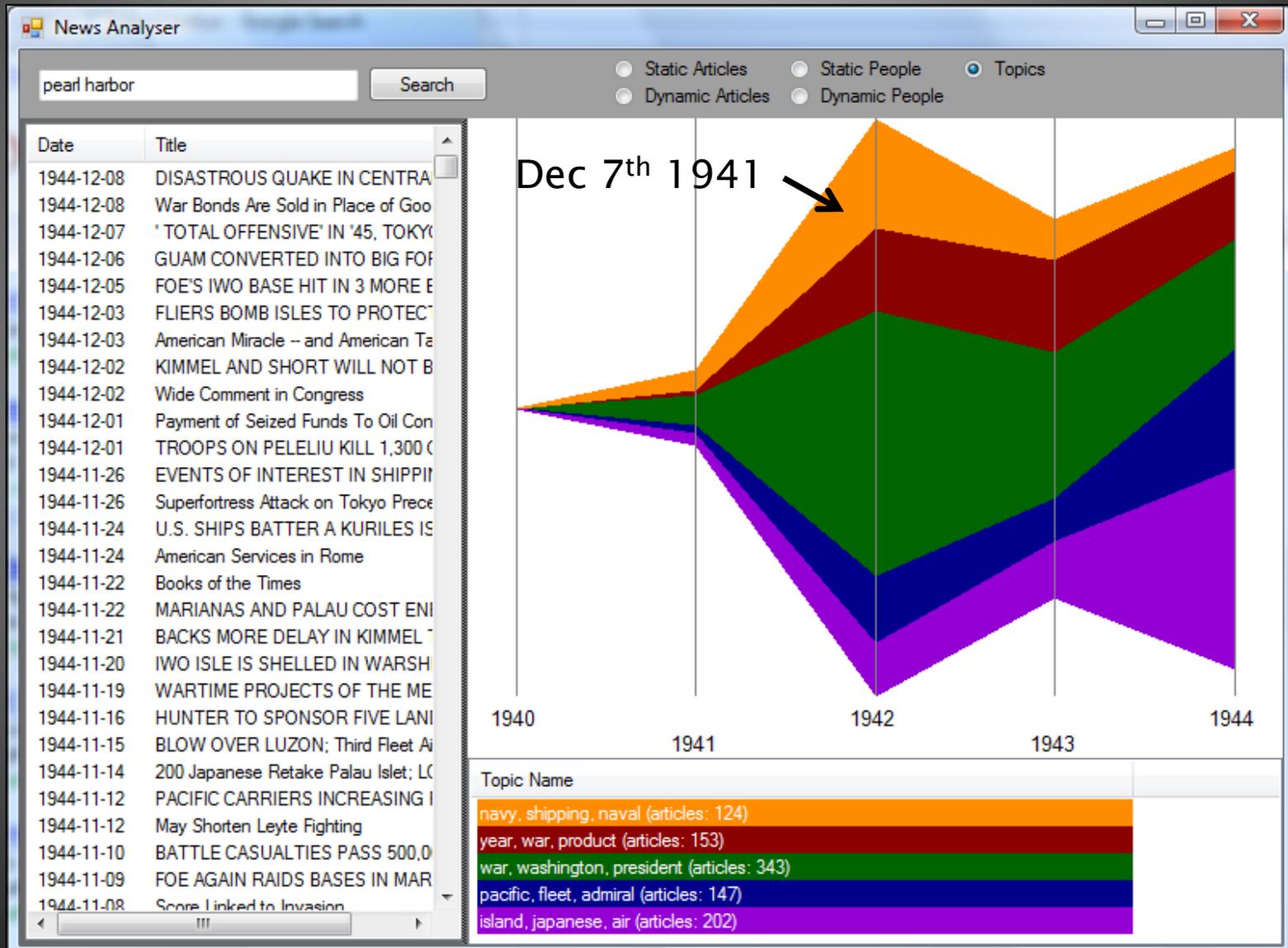
USA: U.S. will attend ...
#documents = 245
NATO, PALESTINIAN,
ISRA, PEAC, ISRAEL,
NETANYAHU,
YELTSIN, ARAFAT,
RUSSIA, SUMMIT

ISRAEL: Israel's Netanyahu to defend settlements in U.S.
Israel's Netanyahu to defend settlements in U.S.
Israel's leader Benjamin Netanyahu will insist when he meets U.S. President Bill Clinton on Monday that the Jewish state has a right to go on building Jewish settlements, Israel said on Thursday.

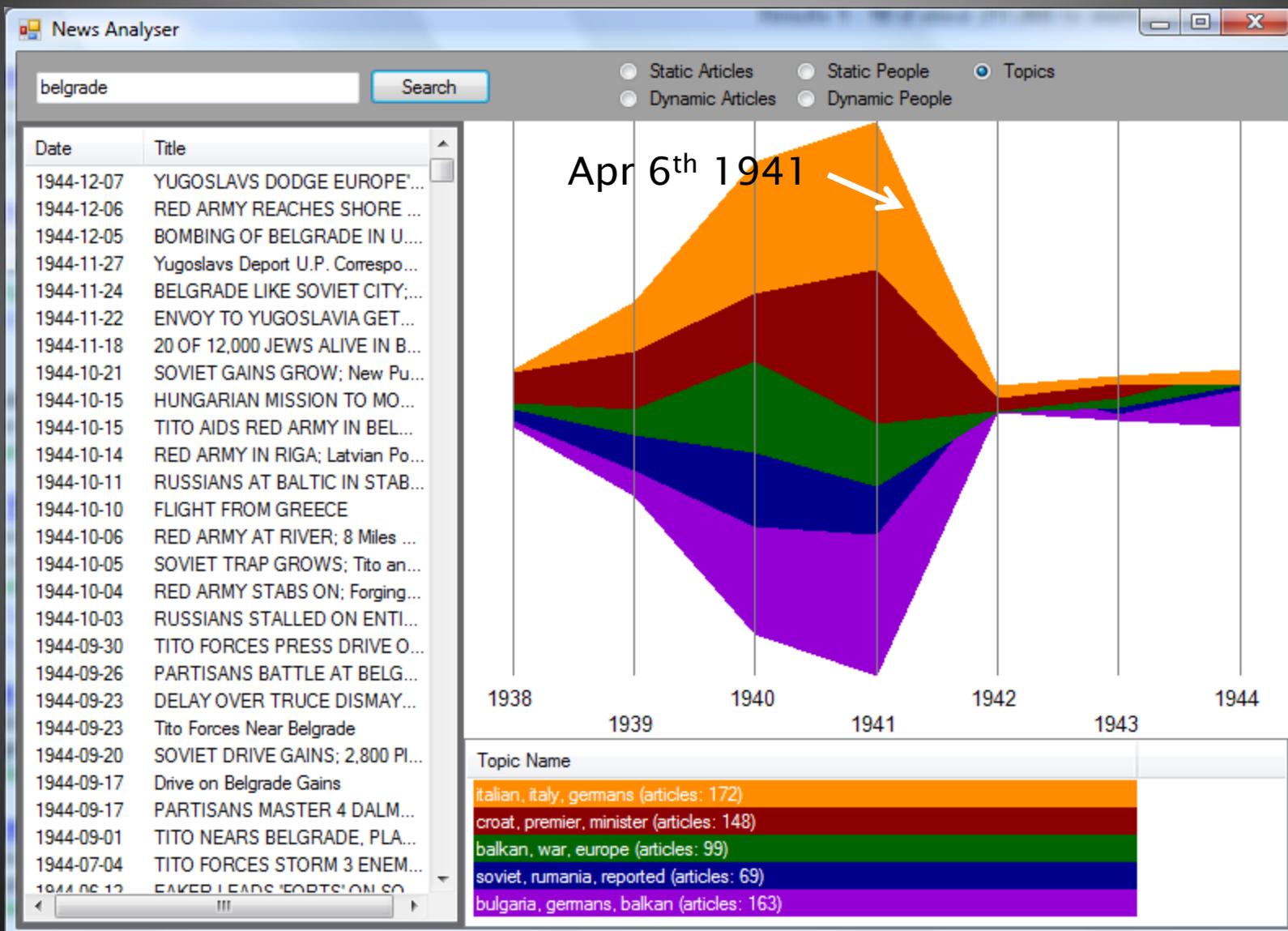
Topic Trends Tracking of the documents including "Clinton"



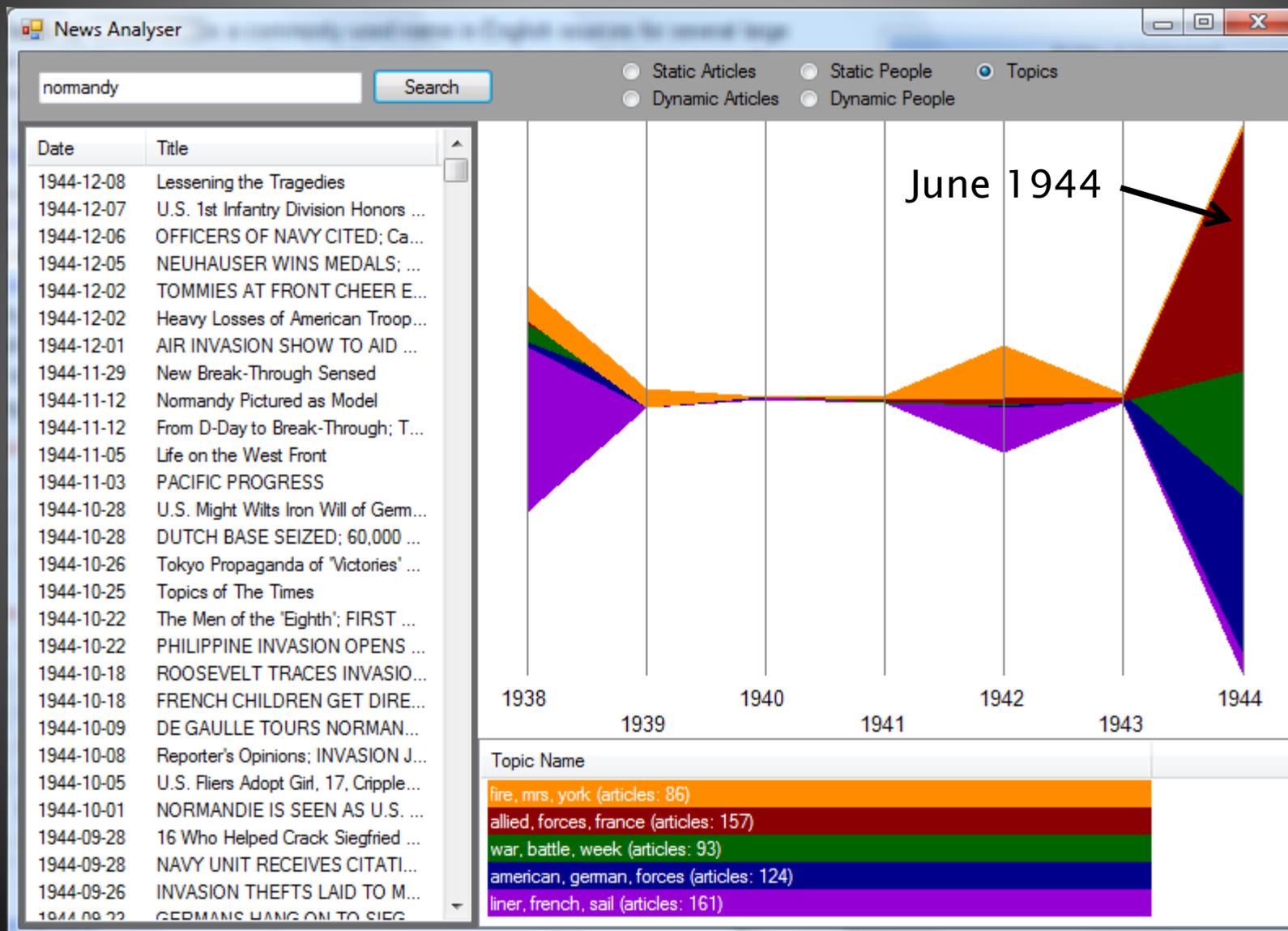
WW2 query “Pearl Harbor” into NYTimes archive



WW2 query “Belgrade” into NYTimes archive

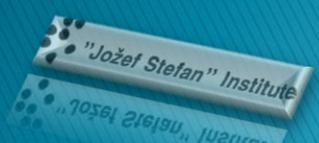


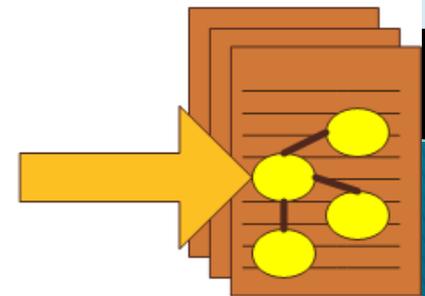
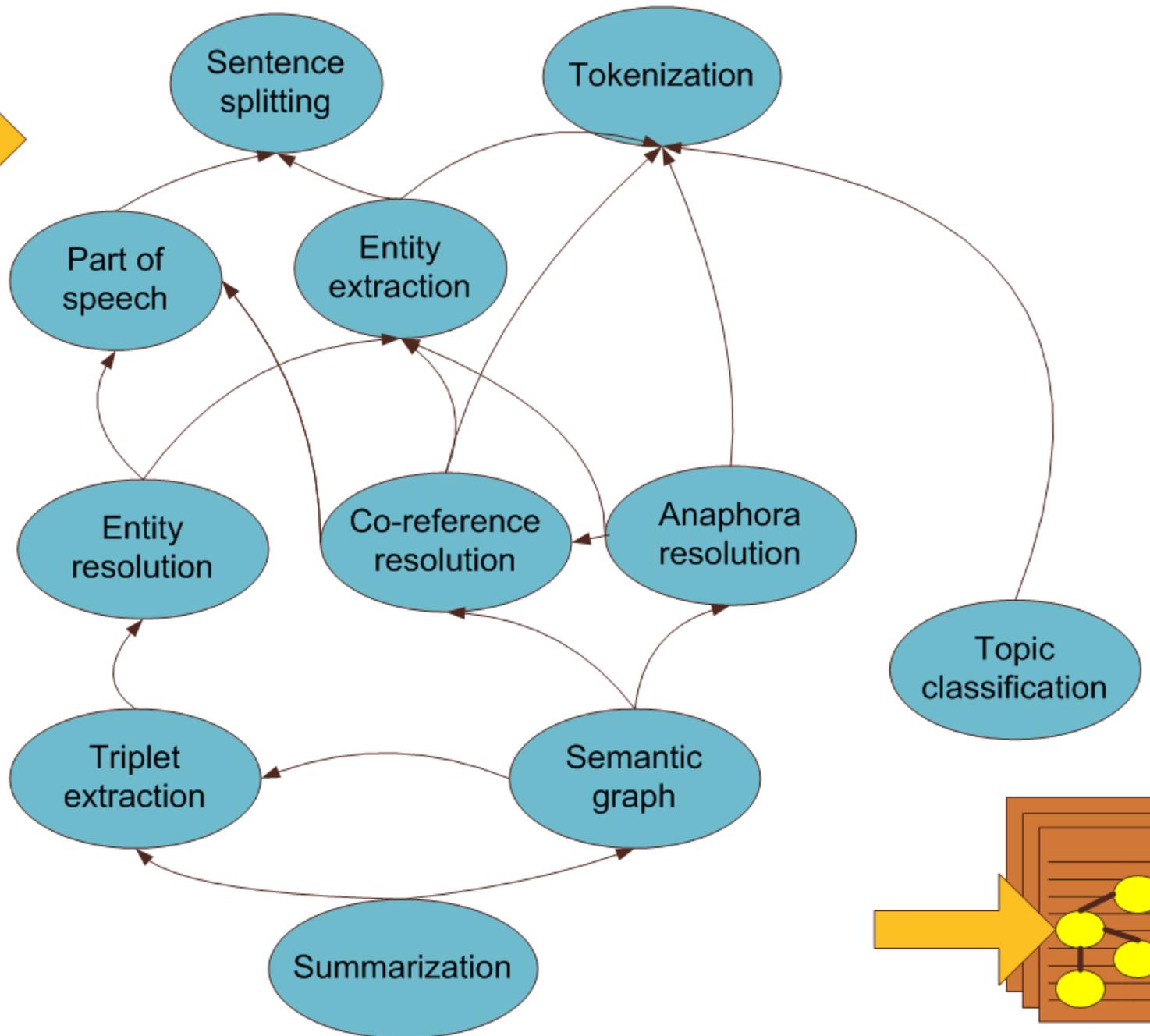
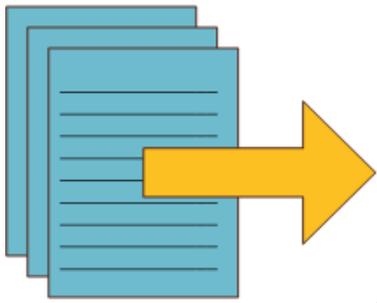
WW2 query “Normandy” into NYTimes archive



Text Enrichment

WWW2009–SemSearch





Knowledge based summarization

AAAI 2005

Detailed Summarization Procedure

Linguistic analysis of the text

- Deep parsing of sentences

Refinement of the text parse

- Named-entity consolidation

Determine that 'George Bush' = 'Bush' = 'U.S. president'

- Anaphora resolution

Link pronouns with name-entities

Extract **Subject-Predicate-Object** triples

Compose a **graph** from triples

Describe each triple with a set of features for learning

Learn a model to classify triples into the summary

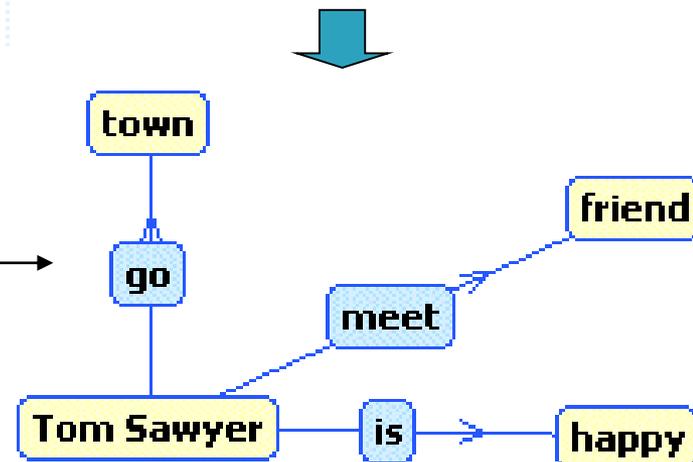
Generate a **summary graph**

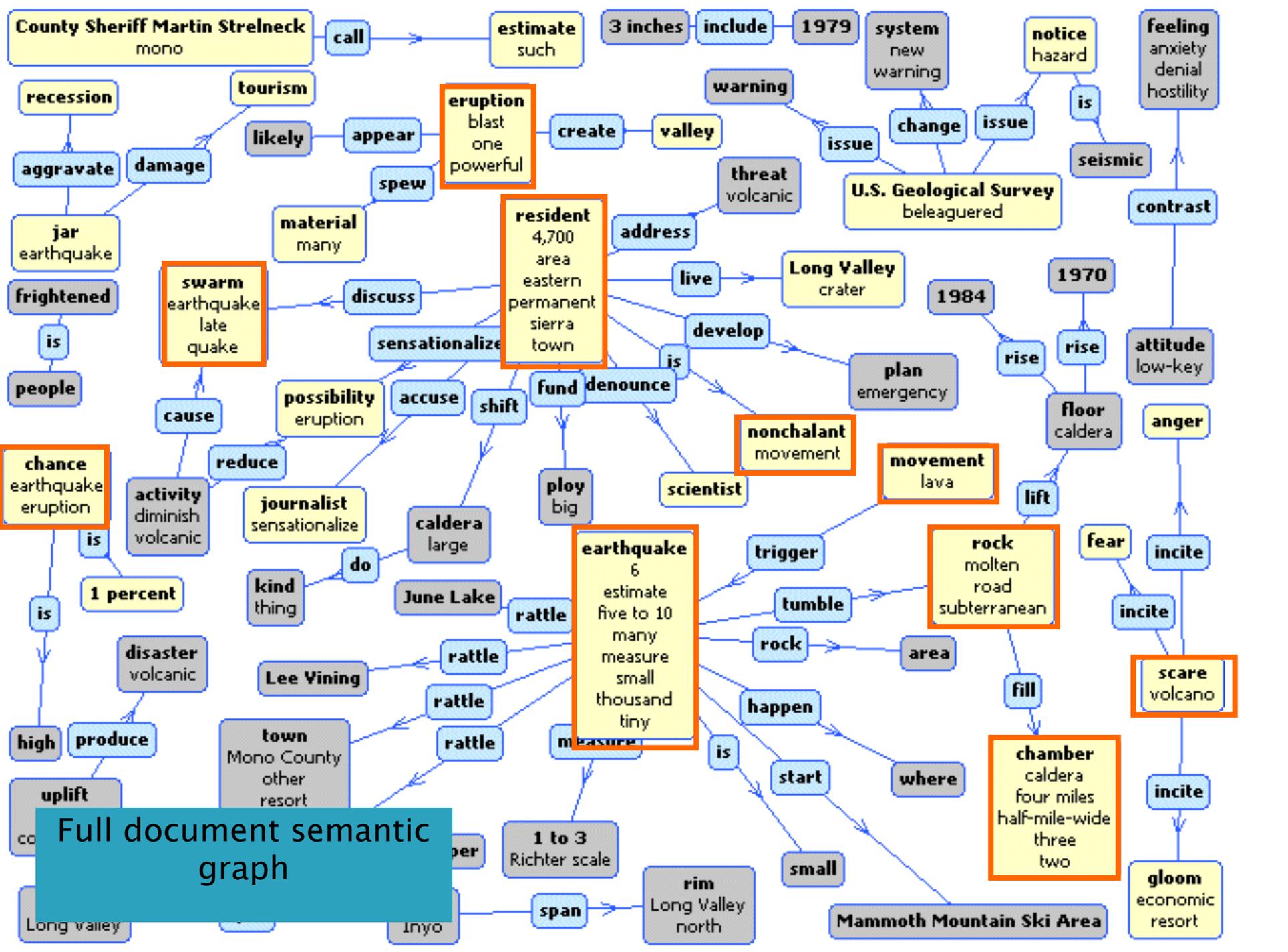
Use summary graph to generate textual document summary

Tom Sawyer went to town. He met a friend. Tom was happy. ...

Tom Sawyer went to town. He [**Tom Sawyer**] met a friend. Tom [**Tom Sawyer**] was happy. ...

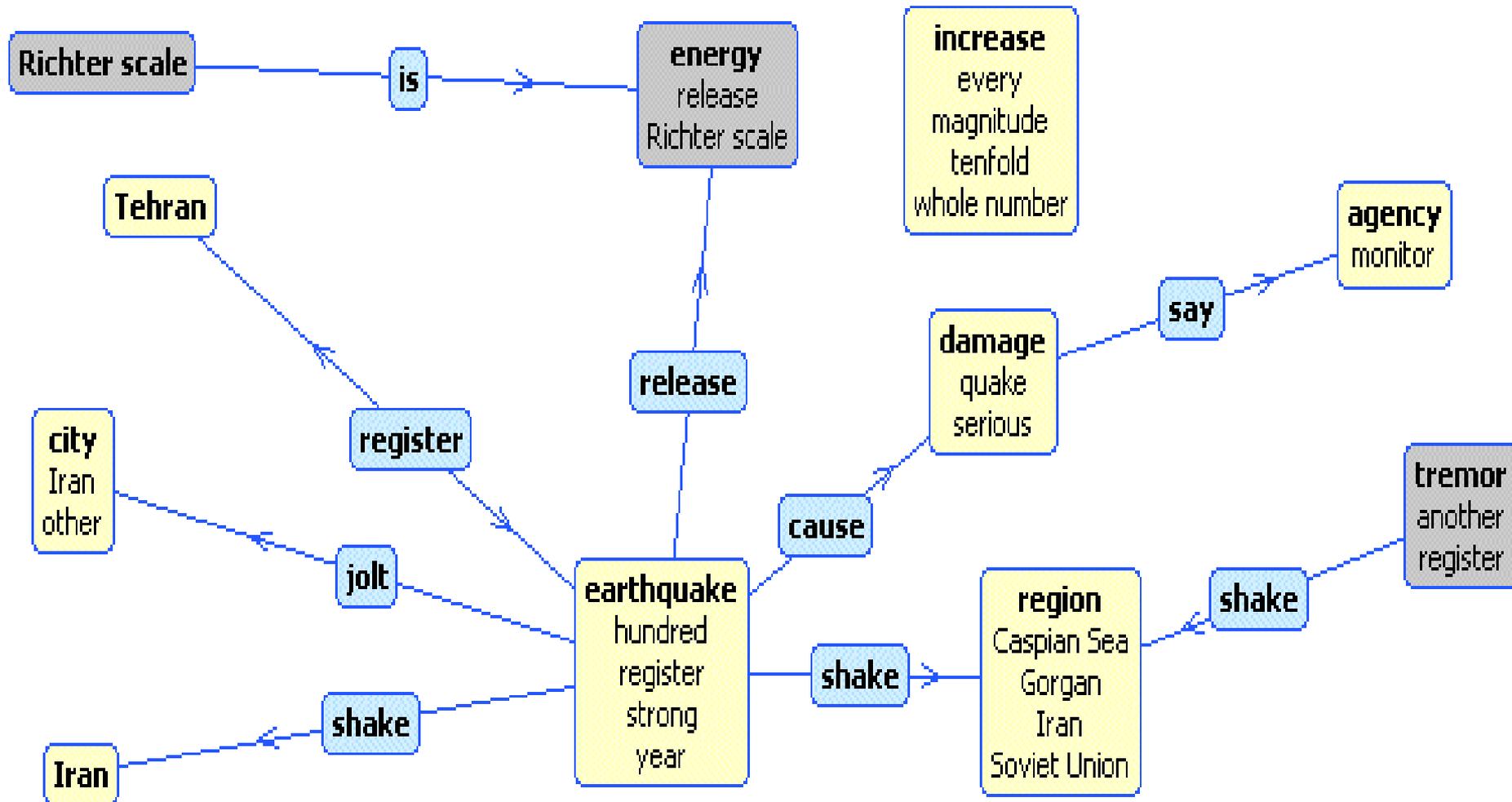
Tom \leftarrow go \rightarrow town
Tom \leftarrow meet \rightarrow friend
Tom \leftarrow is \rightarrow happy

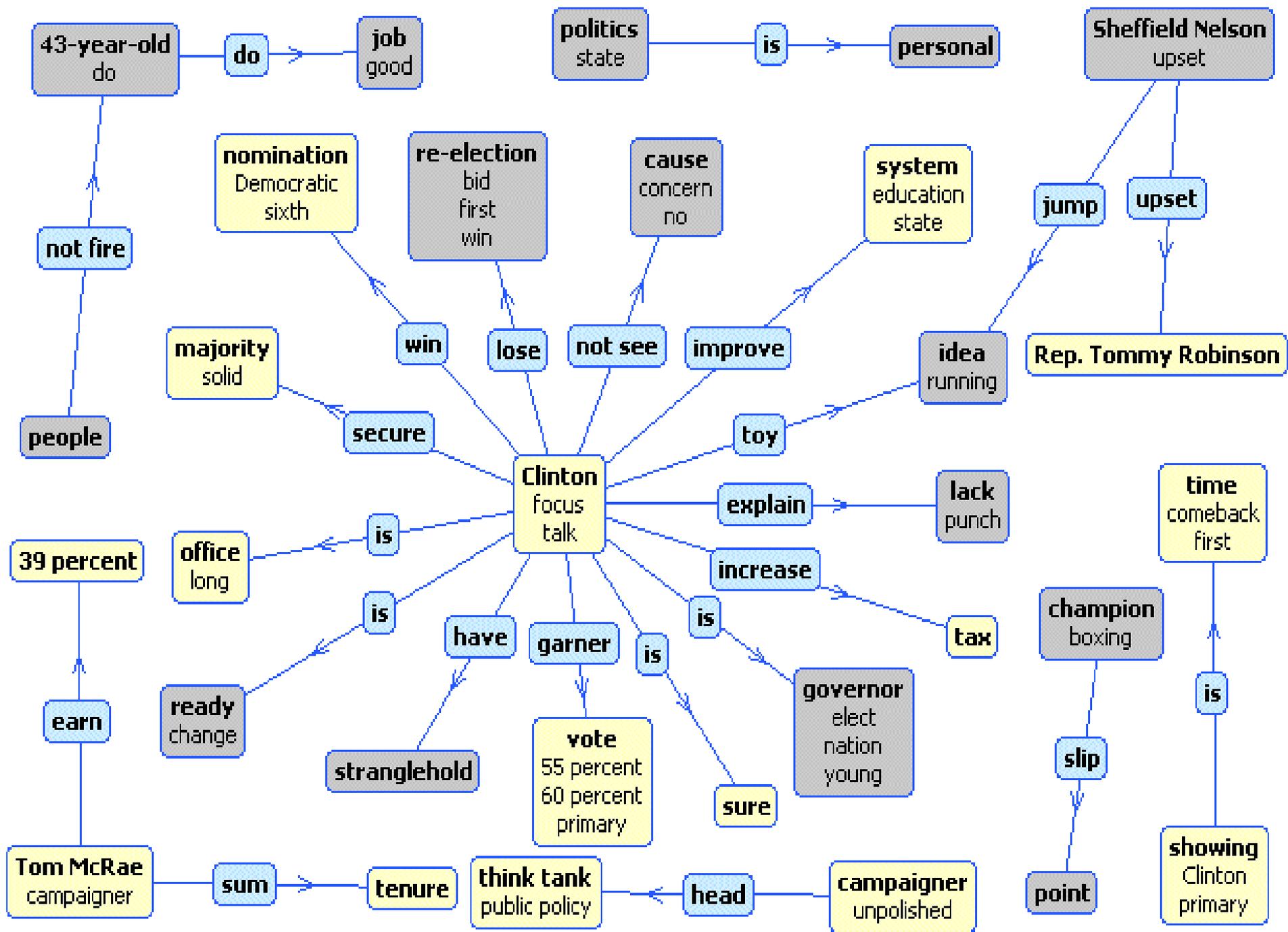




Full document semantic graph

More examples





Cyc Knowledge Base and Reasoning

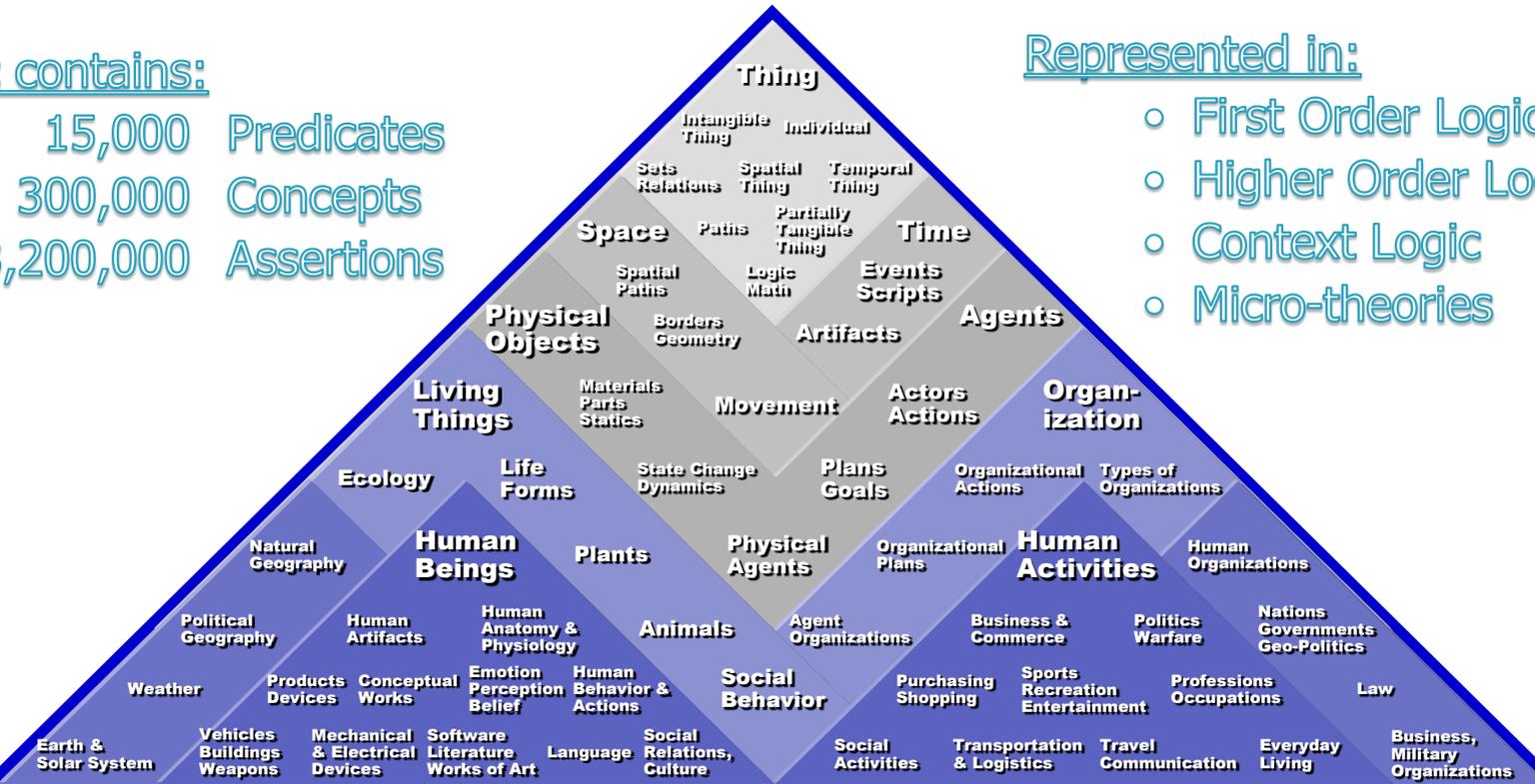
The Cyc Ontology

Cyc contains:

15,000 Predicates
 300,000 Concepts
 3,200,000 Assertions

Represented in:

- First Order Logic
- Higher Order Logic
- Context Logic
- Micro-theories



General Knowledge about Various Domains

Specific data, facts, and observations

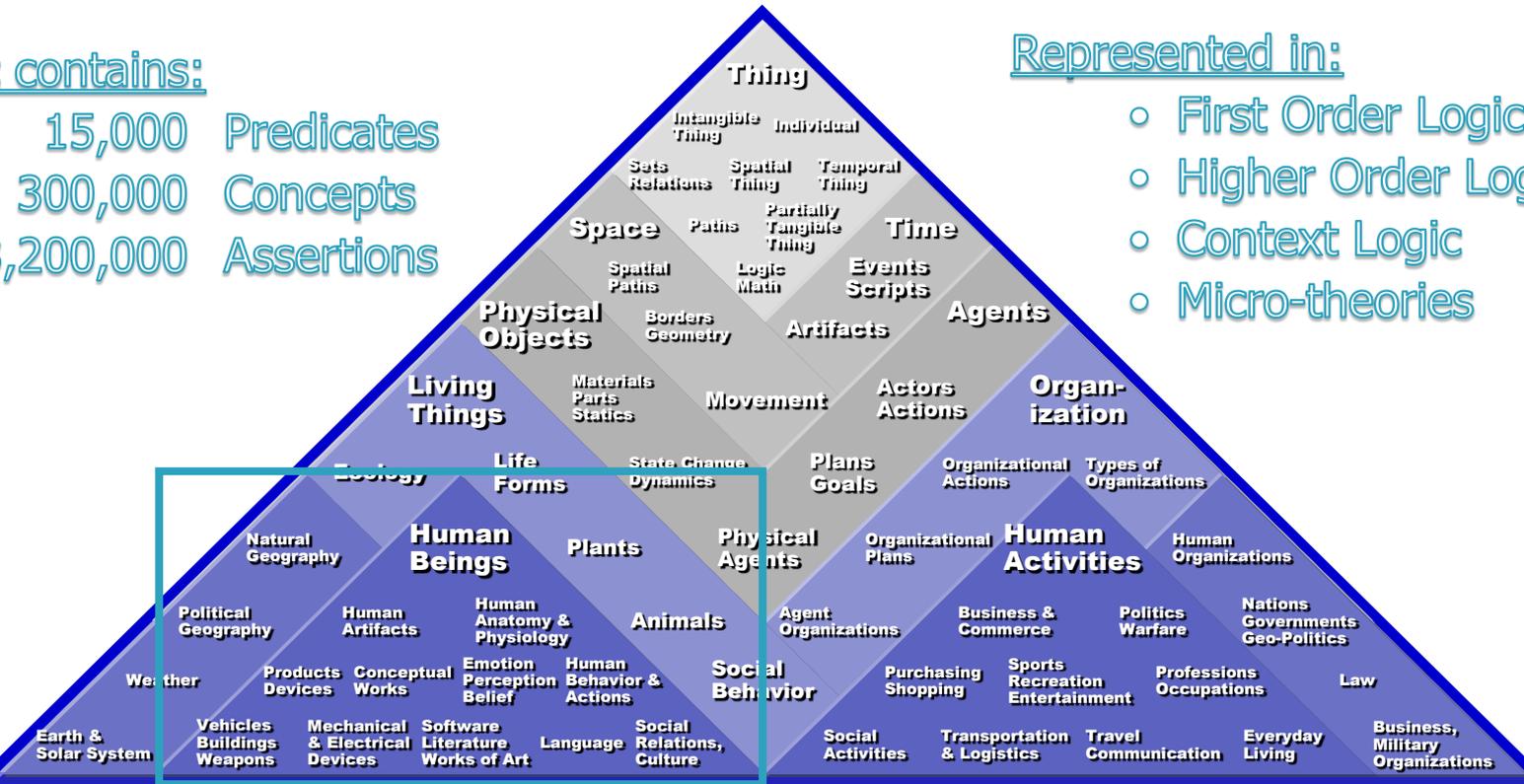
The Cyc Ontology

Cyc contains:

15,000 Predicates
 300,000 Concepts
 3,200,000 Assertions

Represented in:

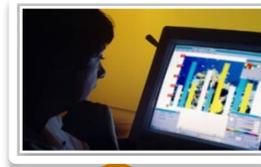
- First Order Logic
- Higher Order Logic
- Context Logic
- Micro-theories



General Knowledge about Various Domains

Specific data, facts, and observations

Knowledge Users



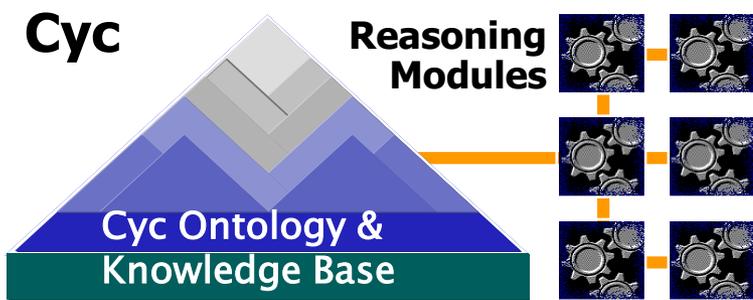
Knowledge Authors



Knowledge Entry Tools

User Interface (with Natural Language Dialog)

Other Applications



Interface to External Data Sources

External Data Sources



Cyc High-level Architecture



Cyc KB Extended w/Domain Knowledge

Thing

Intangible
Thing Individual

General Knowledge about Terrorism:

Terrorist groups are capable of directing assassinations:

(implies

(isa ?GROUP TerroristGroup)

(behaviorCapable ?GROUP AssassinatingSomeone directingAgent))

...

If a terrorist group considers an agent an enemy, that agent is vulnerable to an attack by that group:

(implies

(and

(isa ?GROUP TerroristGroup)

(considersAsEnemy ?GROUP ?TARGET))

(vulnerableTo ?GROUP ?TARGET TerroristAttack))

Solar System

Buildings
Weapons

& Electrical
Devices

Literature
Works of Art

Language
Relations,
Culture

Activities

& Logistics

Communication

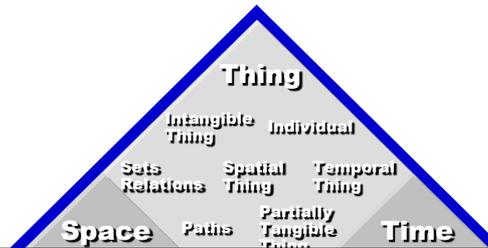
Living

Military
Organizations

General Knowledge about Terrorism

**Specific data, facts, and observations
about terrorist groups and activities**

Cyc KB Extended w/ Domain Knowledge



Specific Facts about Al Qaida:

(basedInRegion AlQaida Afghanistan) Al-Qaida is based in Afghanistan.

(hasBeliefSystems AlQaida IslamicFundamentalistBeliefs) Al-Qaida has Islamic fundamentalist beliefs.

(hasLeaders AlQaida OsamaBinLaden) Al-Qaida is led by Osama bin Laden.

...

(affiliatedWith AlQaida AlQudsMosqueOrganization) Al-Qaida is affiliated with the Al Quds Mosque.

(affiliatedWith AlQaida SudaneseIntelligenceService) Al-Qaida is affiliated with the Sudanese Intell Service

...

(sponsors AlQaida HarakatUIAnsar) Al-Qaida sponsors Harakat ul-Ansar.

(sponsors AlQaida LaskarJihad) Al-Qaida sponsors Laskar Jihad.

...

(performedBy EmbassyBombingInNairobi AlQaida) Al-Qaida bombed the Embassy in Nairobi.

(performedBy EmbassyBombingInTanzania AlQaida) Al-Qaida bombed the Embassy in Tanzania.

General Knowledge about Terrorism

Specific data, facts, and observations about terrorist groups and activities

Example of automatic translating text into Cyc Logic

Source: "Galileo Galilei was an Italian physicist and astronomer."

Learn Logic: (`and (isa GalileoGalilei ItalianPerson)`
`(isa GalileoGalilei Physicist)`
`(isa GalileoGalilei Astronomer)`)

Fact: Galileo was an Italian, a physicist, and an astronomer.

Source: "Galileo was born in Pisa on February 15, 1564."

Learn Logic: (`and (birthDate GalileoGalilei`
`(DayFn 15`
`(MonthFn February`
`(YearFn 1564)))`
`(birthPlace GalileoGalilei CityOfPisaItaly)`)

Fact: Galileo was born on February 15, 1564 and he was born in Pisa.

Source: "Albert Einstein was born in 1879 in Ulm, Germany."

Learn Logic: (`birthDate AlbertEinstein (YearFn 1879)`)

Fact: Albert Einstein was born in 1879.



Cyc's front-end: "Cyc Analytic Environment" – querying (1 / 2)

Task Info Document Search Concepts Related-to Query Creator Queries

Find Stop

WHO had a motive for the assassination of Hariri.

Continue
Save
New Tab
Reset

5 answers
Timed out

Allow speculation?

Answers (5)

Answer	Speculation Level	Sources
Bashar al-Assad	No Speculation	
Syria	Mildly Speculative	
al Qaeda	Moderately Speculative	
United States, the	No Speculation	2
Israel	No Speculation	2

Justify Fact Sheet Visualize Visualize All

Status: **Finished** Message: No appropriate visualizations found

Text query

Who has a motive for the assassination of Rafik Hariri?

Query (semi) automatically translated in the First Order Logic

Answers to the query

Cyc's front-end: "Cyc Analytic Environment" – justification (2/2)

Task Info Document Search Concepts Related-to Query Creator Queries Justification Justification Justification

Proof 1 Save... Copy

▼ **Query:** Who or what had a motive for the assassination of Hariri? ← **Query & Answer**
Answer: al Qaeda

▼ **Because:**

Since 2000, Lebanon has been responsible for according with Lebanese economic reform.  1

February 14, 2005 was the date of the assassination of Hariri.  2

Rafik Hariri was killed during the assassination of Hariri.  2 ← **Justification**
Rafik Hariri is an advocate of Lebanese economic reform.

Al Qaeda opposes Lebanese economic reform.

▼ **Detailed Justification:**
▶ Al Qaeda had a motive for the assassination of Hariri.

▼ **External Sources:**

1  Gary C. Gambill, "Dossier: Rafiq Hariri", *United States Committee for a Free Lebanon*, July 2001, http://www.meib.org/articles/0107_id1.htm.

2  "Huge blast kills Lebanese ex-PM", *the Cable News Network*, February 14, 2005, <http://www.cnn.com/2005/WORLD/meast/02/14/beirut.explosion.1910/>.

▼ Options
▼ Options

Sources for Reasoning and Justification ←

Cyc's front-end: "Cyc Analytic Environment" – justification (2/2)

Task Info Document Search Concepts Related-to Query Creator Queries Justification Justification Justification

Proof 1 Save... Copy

▼ Query: Who or what had a motive for [the assassination of Hariri](#)?
 Answer: [al Qaeda](#)
 Because:

← Query & Answer

Since 2000, [Lebanon](#) has been responsible for according with [Lebanese economic reform](#).  1

February 14, 2005 was the date of [the assassination of Hariri](#).  2

← Justification

[Rafik Hariri](#) was killed during [the assassination of Hariri](#).  2
[Rafik Hariri](#) is an advocate of [Lebanese economic reform](#).

[Al Qaeda](#) opposes [Lebanese economic reform](#).

▼ Detailed Justification:
 ▶ [Al Qaeda](#) had a motive for [the assassination of Hariri](#).

▼ External Sources:
 Gary C. Gambill, "Dossier: Rafiq Hariri", *United States Committee for a Free Lebanon*, July 2001, http://www.meib.org/articles/0107_id1.htm.
 "Huge blast kills Lebanese ex-PM", *the Cable News Network*, February 14, 2005, <http://www.cnn.com/2005/WORLD/meast/02/14/beirut.explosion.1910/>.

← Sources for Reasoning and Justification

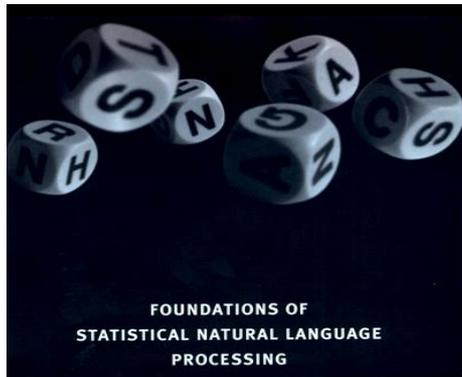
▶ If

- some intelligent agent opposes some policy,
- and some other intelligent agent *VICTIM* is an advocate of that policy,
- and some other intelligent agent *ADOPTER* is responsible for according with the policy,
- and it is adopted by *ADOPTER* in any some *ADOPT-TYPE*,
- and some *ACT* prevents *VICTIM* from playing the role "key participant" in any *ADOPT-TYPE*,

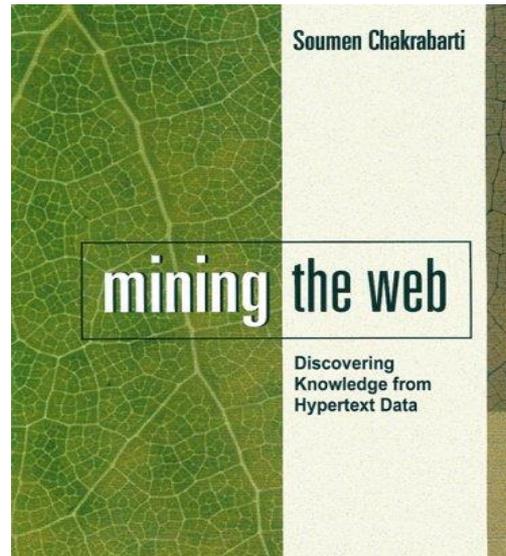
then that intelligent agent had a motive for *ACT*.

Further references...

References to some Text-Mining books



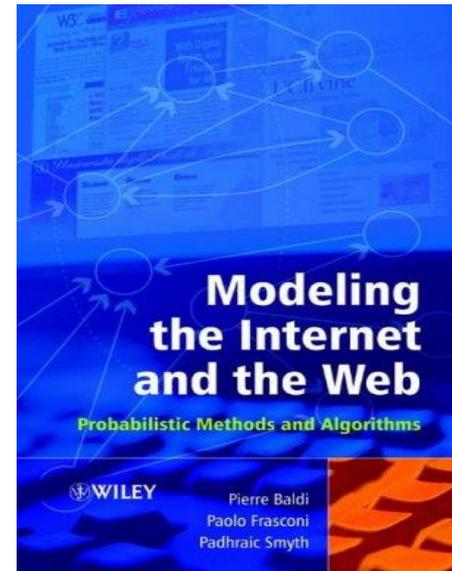
FOUNDATIONS OF
STATISTICAL NATURAL LANGUAGE
PROCESSING
CHRISTOPHER D. MANNING AND
HINRICH SCHÜTZE



Soumen Chakrabarti

mining the web

Discovering
Knowledge from
Hypertext Data

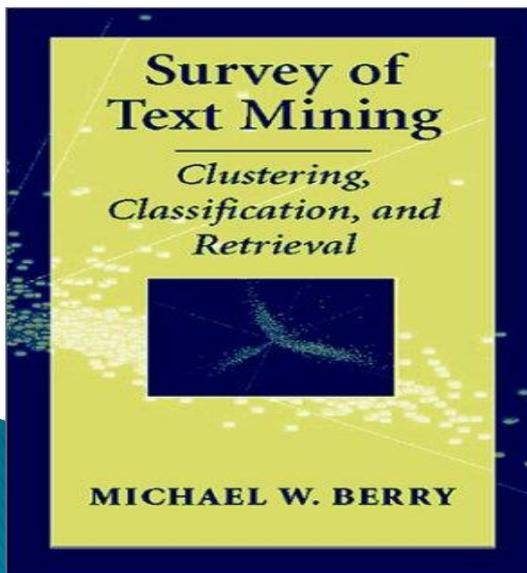


Modeling
the Internet
and the Web

Probabilistic Methods and Algorithms

WILEY

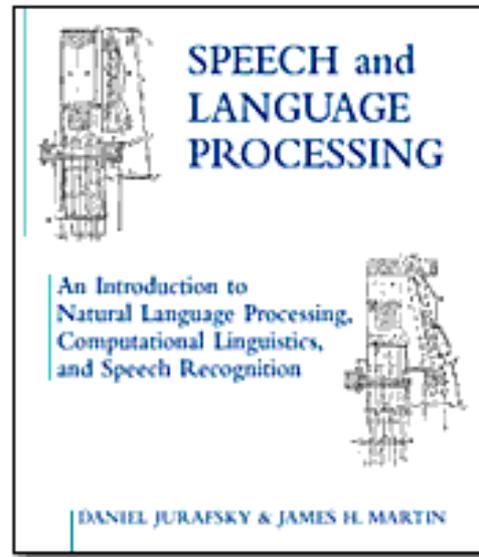
Pierre Baldi
Paolo Frasconi
Padhraic Smyth



Survey of
Text Mining

Clustering,
Classification, and
Retrieval

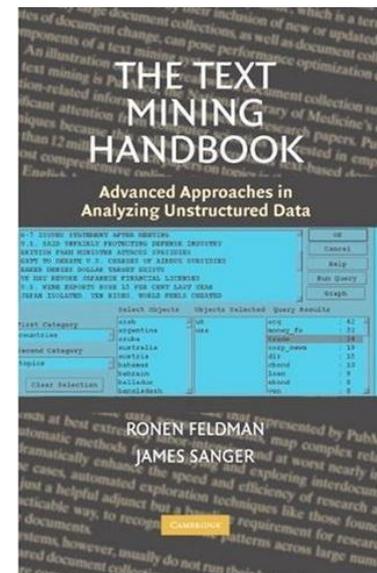
MICHAEL W. BERRY



SPEECH and
LANGUAGE
PROCESSING

An Introduction to
Natural Language Processing,
Computational Linguistics,
and Speech Recognition

DANIEL JURAFSKY & JAMES H. MARTIN

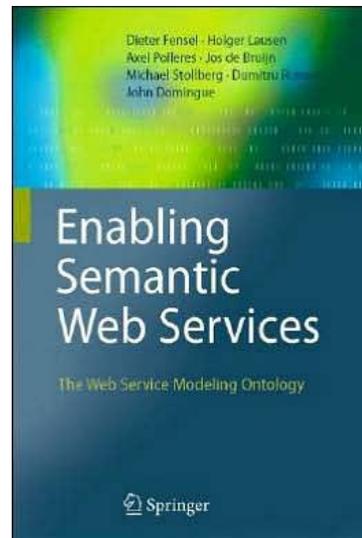
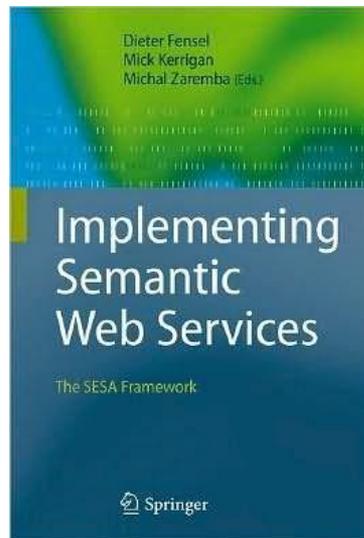
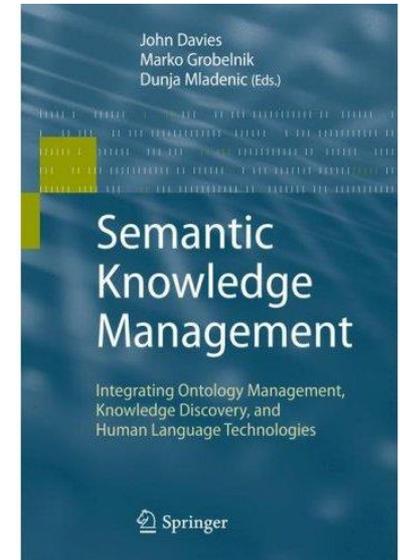
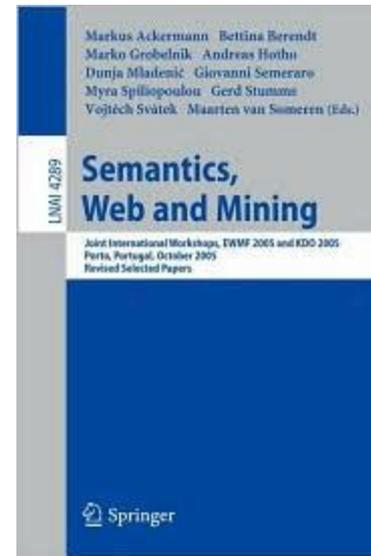
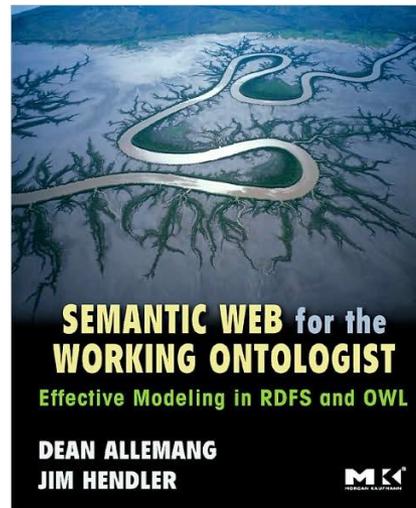
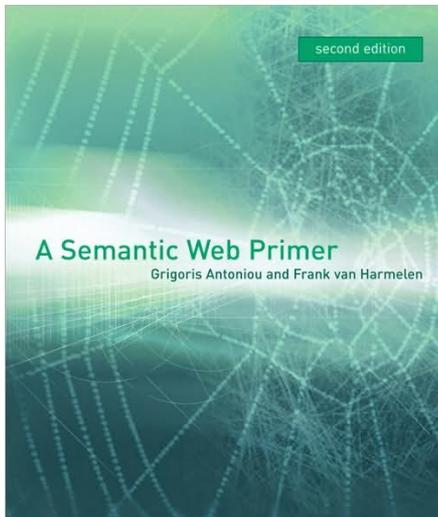


THE TEXT
MINING
HANDBOOK

Advanced Approaches in
Analyzing Unstructured Data

RONEN FELDMAN
JAMES SANGER

Books on Semantic Technologies



References to the main conferences

- ▶ **Information Retrieval:**
 - SIGIR, ECIR
- ▶ **Machine Learning/Data Mining:**
 - ICML, ECML/PKDD, KDD, ICDM, SDM
- ▶ **Computational Linguistics:**
 - ACL, EACL, NAACL
- ▶ **Semantic Web:**
 - ISWC, ESWC, ASWC

Videos on Text and Semantic Technologies

http://videolectures.net/Top/Computer_Science/Semantic_Web/

- ▶ Recorded tutorials, workshops, conferences, summer schools available from <http://videolectures.net>

The screenshot shows the VideoLectures website in a Mozilla Firefox browser window. The browser title is "VideoLectures - exchange ideas & share knowledge - Mozilla Firefox". The address bar shows the URL "http://videolectures.net/". The page content includes:

- Navigation:** HOME, MOST POPULAR, LATEST LECTURES, CATEGORIES, EVENTS, PEOPLE, INTERVIEWS, TUTORIALS, ABOUT US
- FEATURED LECTURES:** Five video thumbnails with titles and view counts:
 - POWERSET - Natural Language and the Semantic Web (1461 views, 01:09:28)
 - Leonardo: Goal assistance with divergent beliefs (346 views, 00:03:31)
 - Bayesian Kernel Methods (164 views, 04:32:50)
 - Functional Analysis in Data Modelling (90 views, 01:06:03)
 - Sparse Log Gaussian Processes via MCMC for Spatial Epidemiology (42 views, 00:05:51)
- RECENT EVENTS:** ISWC '08 - Karlsruhe (7th International Semantic Web Conference), 18.085 Computational Science and Engineering I (MIT 18.085 Computational Science and Engineering I - Fall 2007), ESTC '08 - Vienna (2nd European Semantic Technology Conference).
- NEWS:** MIT OpenCourseWare Collection (MIT OPENCOURSEWARE), Cambridge University Engineering Department - Machine Learning seminars (UNIVERSITY OF CAMBRIDGE), Carnegie Mellon Machine Learning Lunch seminar (Carnegie Mellon).
- CATEGORIES:** Architecture (2), Arts (24), Biology (38), Business (70), Chemistry (12), Computers (16), Computer Science (1375), Economics (9), Education (4), Environment (12), Events (35), History (2), Law (12), Mathematics (73), Medicine (29), Philosophy (7), Physics (15), Psychology (24), Science (41), Society (33), Technology (8).
- FEATURED:** Universal Access to All Knowledge - Archive.org (172 views, 00:26:53).