

---

# Relation Extraction

Feiyu Xu



# Outline

---

- ⊙ Introduction to relation extraction
- ⊙ Brief history of information extraction
- ⊙ Machine learning approaches to relation extraction
  - ⊙ State-of-the-art
  - ⊙ DARE System (<http://dare.dfki.de>)
- ⊙ Big Text Data Analytics
- ⊙ References

- ③ In a single day on the Internet, PayPal processes over \$315 million in transactions, Facebook takes in about 350 million uploaded images, and Twitter blasts out 400 million tweets from its users
- ③ From the perspectives of data content processing and data mining, language data belongs to unstructured data like images or videos because of the complexities of their interne structures.
- ③ Technologies such as information retrieval and text analytics are needed to ease human users to access the relevant information in a realistic time from a big amount of data



## Information Extraction is ...

---

a technology that is futuristic from the user's point of view in the current information-driven world.

Rather than indicating which documents need to be read by a user, it extracts pieces of information that are salient to the user's needs ...

provided by NIST:  
[[http://www-nlpir.nist.gov/related\\_projects/muc/](http://www-nlpir.nist.gov/related_projects/muc/)]



# Information Extraction: A Pragmatic Approach

---

- ③ Identify the types of entities that are relevant to a particular task
- ③ Identify the range of facts that one is interested in for those entities
- ③ Ignore everything else

[Appelt, 2003]



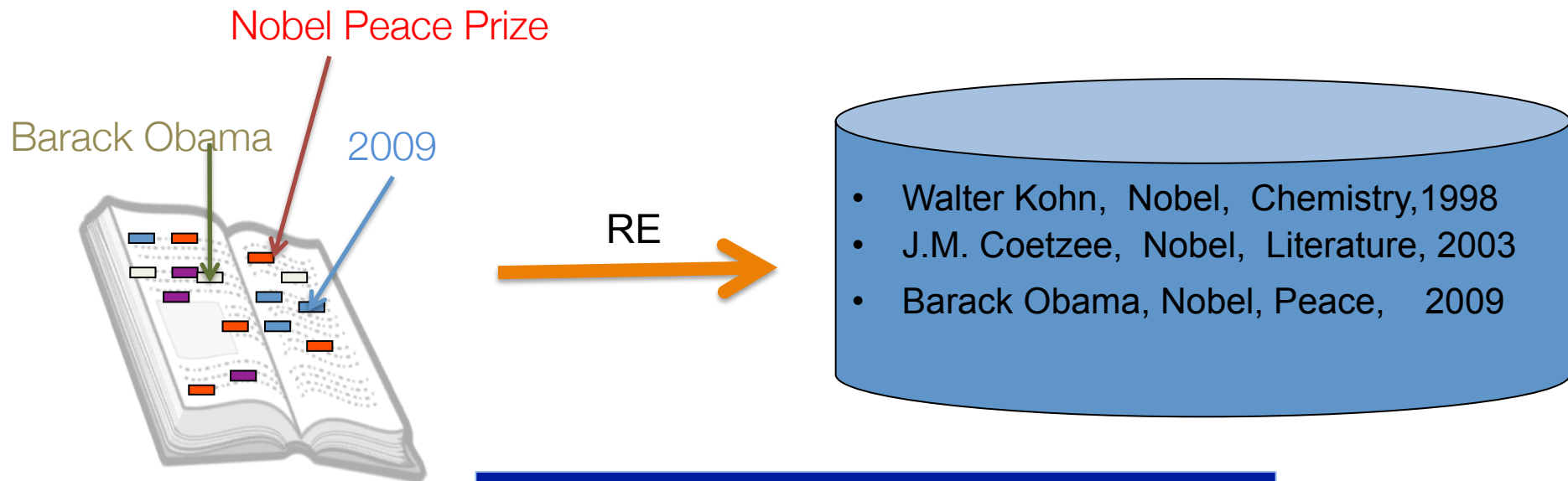
# Types of Information Extraction in LT

---

- ⊙ Topic Extraction
- ⊙ Term Extraction
- ⊙ Named Entity Extraction
- ⊙ Binary Relation Extraction
- ⊙ N-ary Relation Extraction
- ⊙ Event Extraction
- ⊙ Answer Extraction
- ⊙ Opinion Extraction
- ⊙ Sentiment Extraction



- Given an **unstructured** text, a relation extraction (RE) tool should be able to automatically recognize and extract relations among the relevant entities or concepts that are salient to the user's needs



## Linguistic Patterns:

- *<prize>* be awarded to *<person>*
- *<person>* win *<prize>* in *<year>*
- .....

# Types of Information Extraction in LT

---

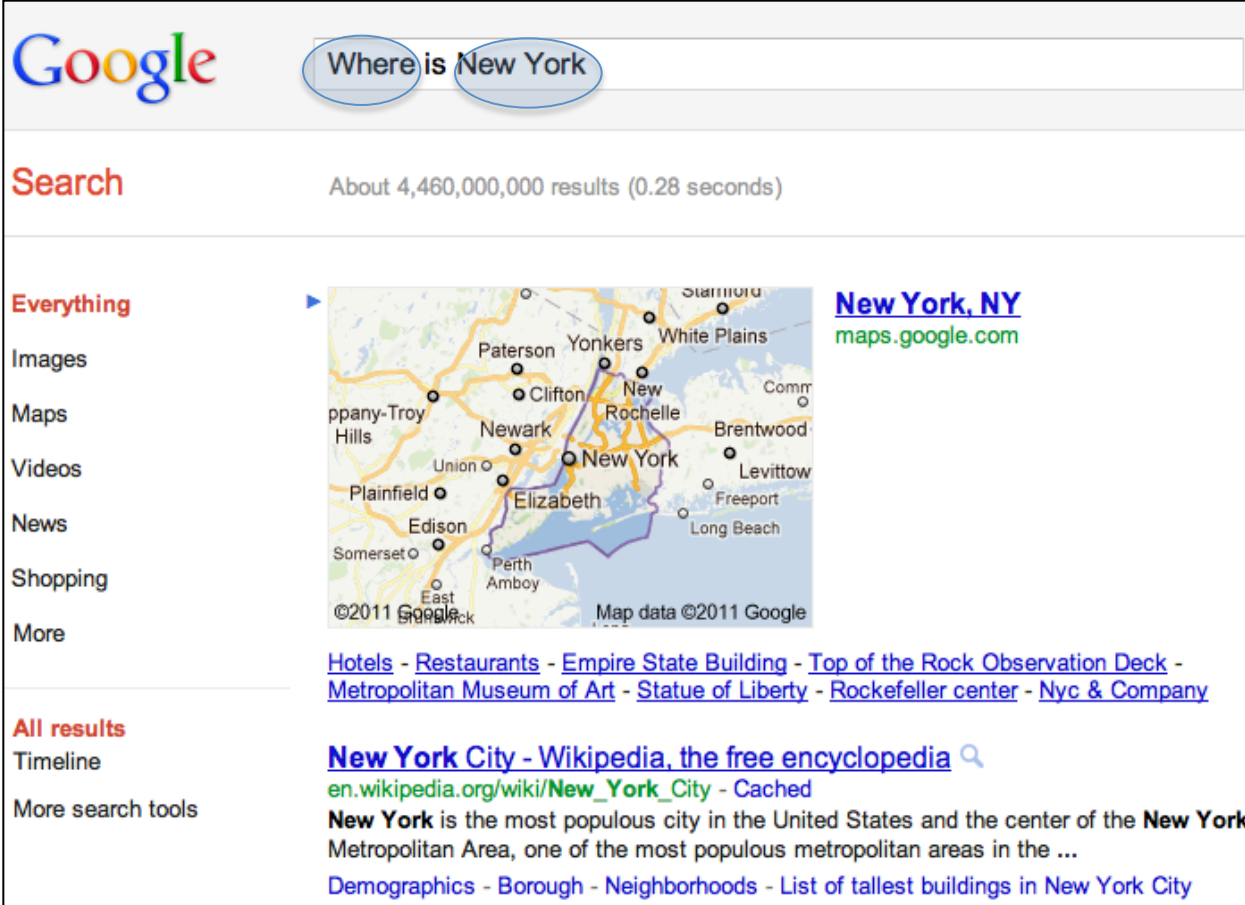
- ⊙ Topic Extraction
- ⊙ Term Extraction
- ⊙ Named Entity Extraction
- ⊙ Binary Relation Extraction
- ⊙ N-ary Relation Extraction
- ⊙ Event Extraction
- ⊙ Answer Extraction
- ⊙ Opinion Extraction
- ⊙ Sentiment Extraction

**Types of Relation Extraction**

# General application task 1:

## ☆ Information access for information finder

mapping unstructured textual queries of users to more structured formal query for search and answer engines



The image shows a screenshot of a Google search interface. At the top left is the Google logo. To its right is a search bar containing the text "Where is New York". Below the search bar, the word "Search" is displayed in red, followed by the text "About 4,460,000,000 results (0.28 seconds)". On the left side, there is a vertical menu with the following items: "Everything" (selected), "Images", "Maps", "Videos", "News", "Shopping", and "More". The main content area features a map of New York City and its surrounding areas, including labels for Paterson, Yonkers, White Plains, Newark, Rochelle, New York, Elizabeth, Union, Plainfield, Edison, Perth Amboy, and Long Beach. To the right of the map is a link for "New York, NY" with the URL "maps.google.com". Below the map, there is a list of search results, including "Hotels - Restaurants - Empire State Building - Top of the Rock Observation Deck - Metropolitan Museum of Art - Statue of Liberty - Rockefeller center - Nyc & Company" and "New York City - Wikipedia, the free encyclopedia" with a magnifying glass icon. The Wikipedia result includes the URL "en.wikipedia.org/wiki/New\_York\_City - Cached" and a snippet of text: "New York is the most populous city in the United States and the center of the New York Metropolitan Area, one of the most populous metropolitan areas in the ...". At the bottom, there are more search results: "Demographics - Borough - Neighborhoods - List of tallest buildings in New York City".

## General application task 2:

### ☆ Information acquisition for information provider

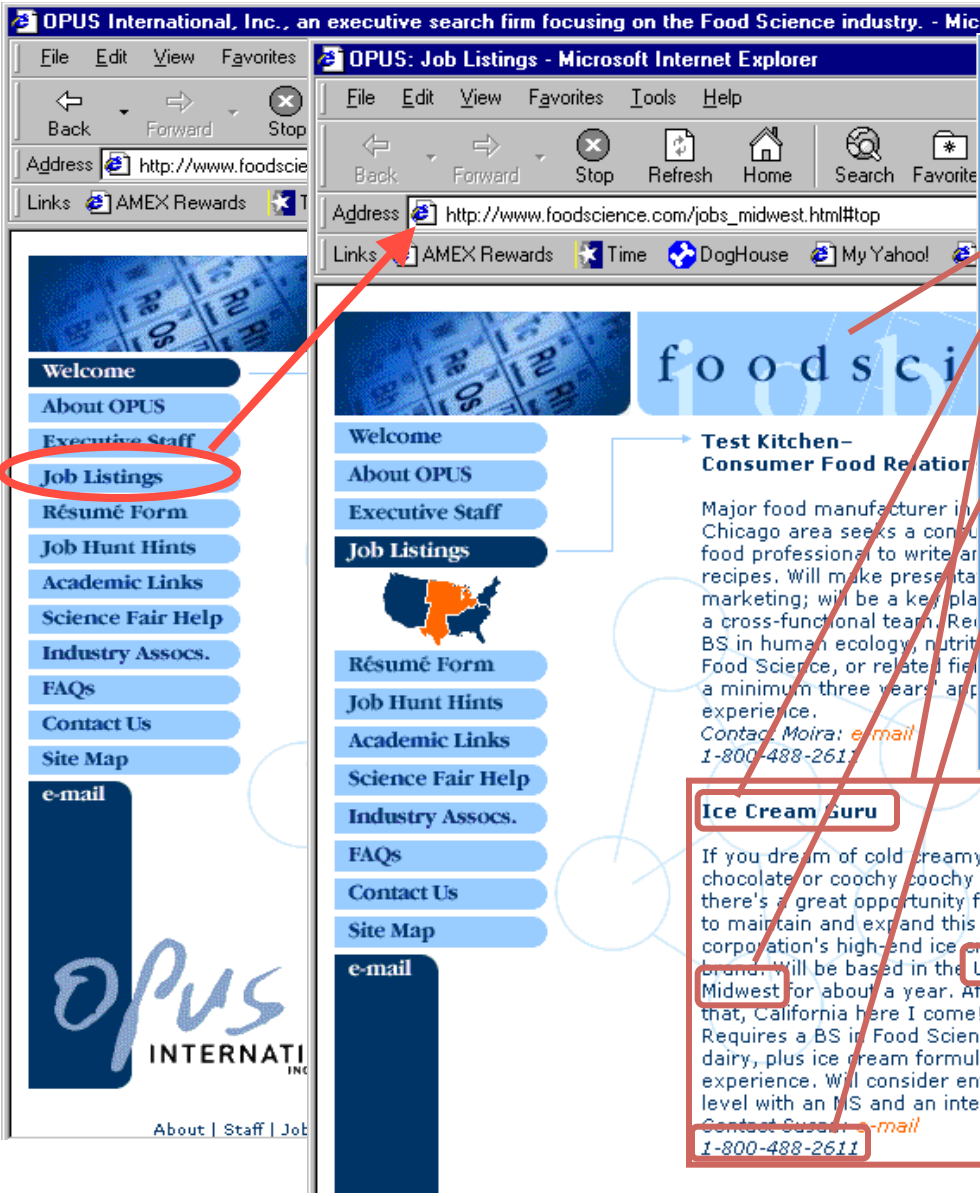
extract structured information from big amount free texts to construct knowledge bases



RE



# Extracting Job Openings from the Web: Semi-Structured Data



**foodscience.com-Job2**

JobTitle: Ice Cream Guru  
Employer: foodscience.com  
JobCategory: Travel/Hospitality  
JobFunction: Food Services  
JobLocation: Upper Midwest  
Contact Phone: 800-488-2611  
DateExtracted: January 8, 2001  
Source: www.foodscience.com/jobs\_midwest.html  
OtherCompanyJobs: foodscience.com-Job1

**Ice Cream Guru**

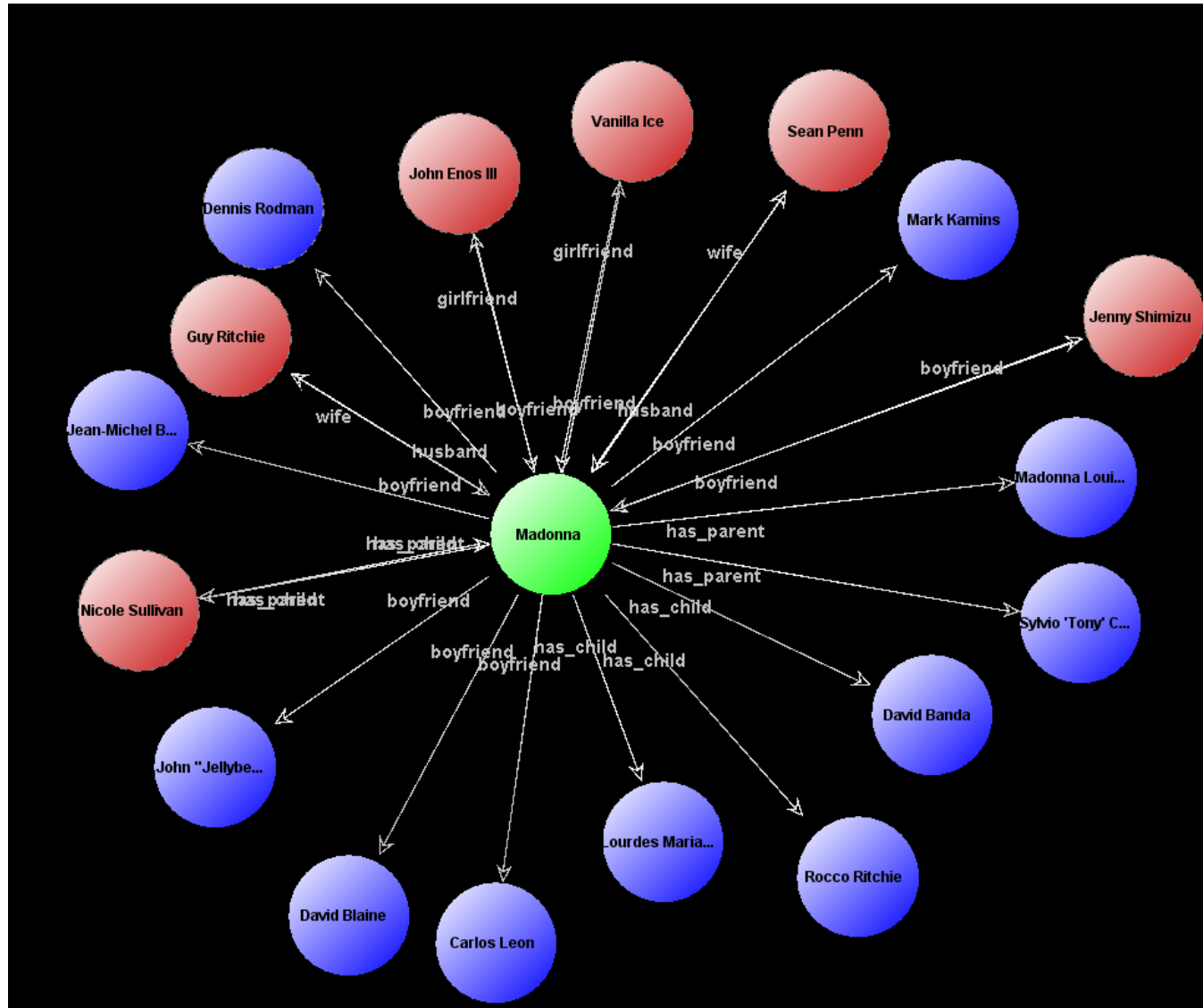
If you dream of cold creamy chocolate or coochy coochy cookie, there's a great opportunity for you to maintain and expand this major corporation's high-end ice cream brand. Will be based in the Upper Midwest for about a year. After that, California here I come! Requires a BS in Food Science or dairy, plus ice cream formulation experience. Will consider entry level with an MS and an internship.

Contact Susan [e-mail](#)  
1-800-488-2611

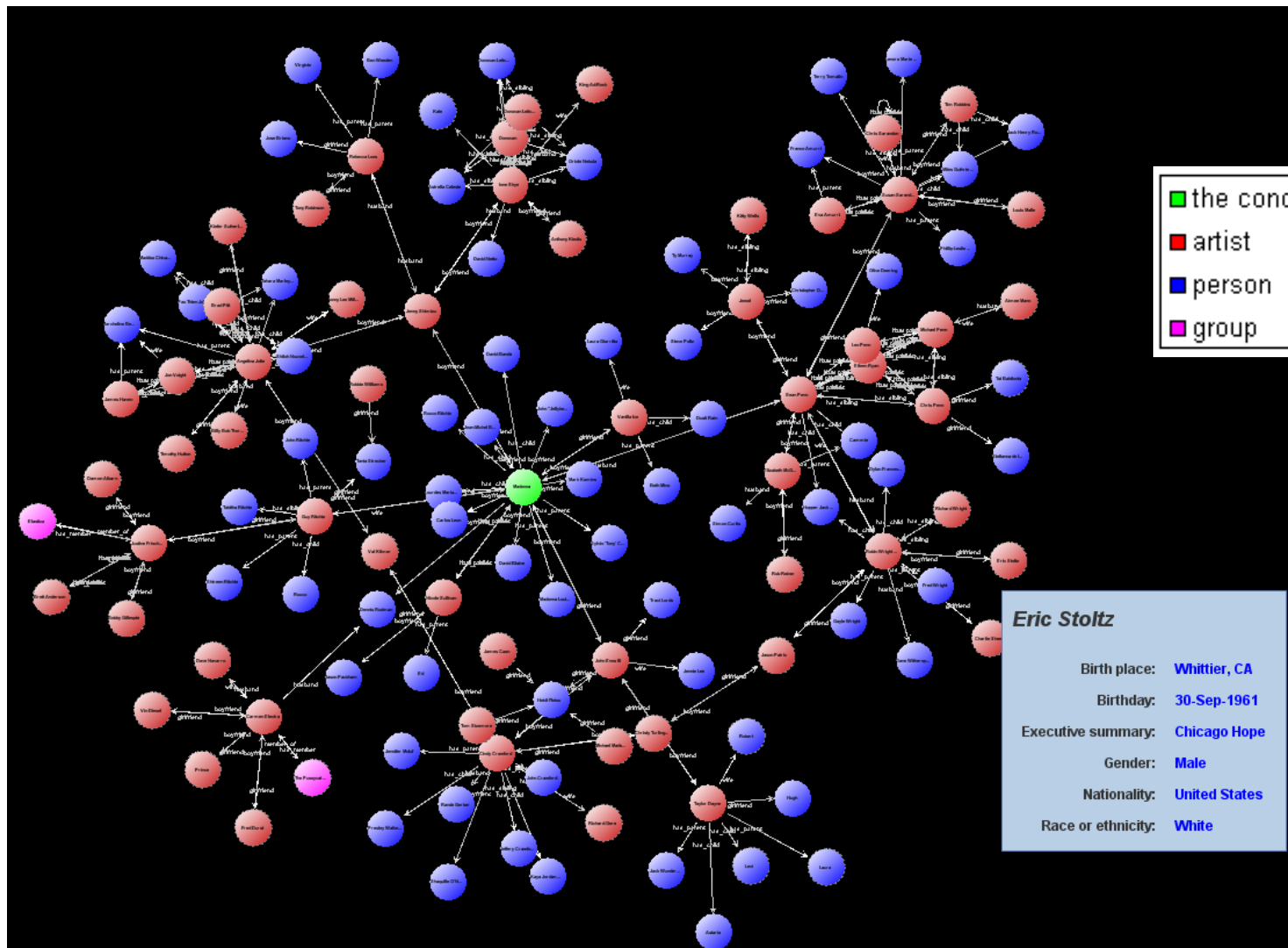


# Acquisition of Social Network of Pop Stars from Web

Social Network of "Madonna" (Depth = 1)

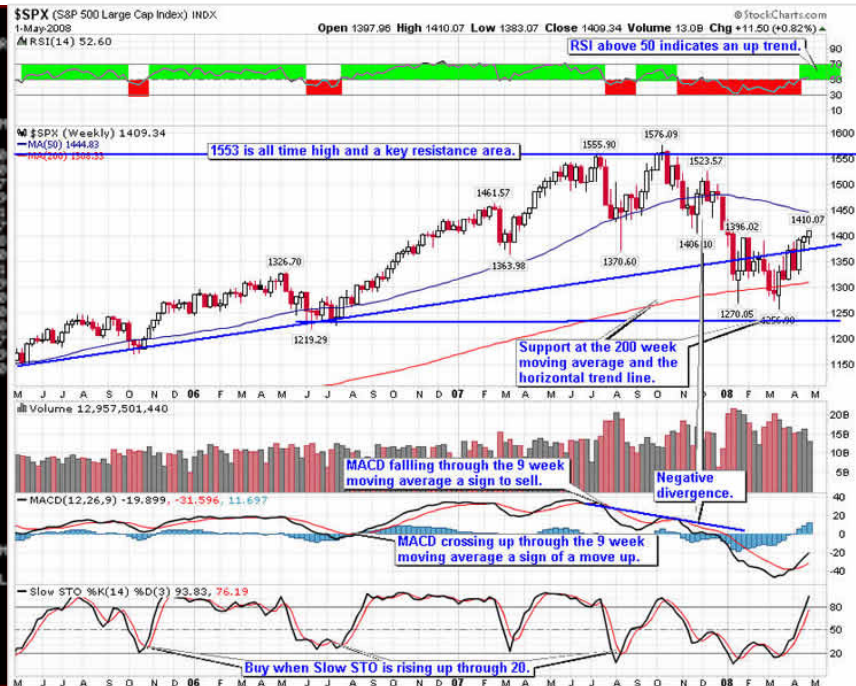
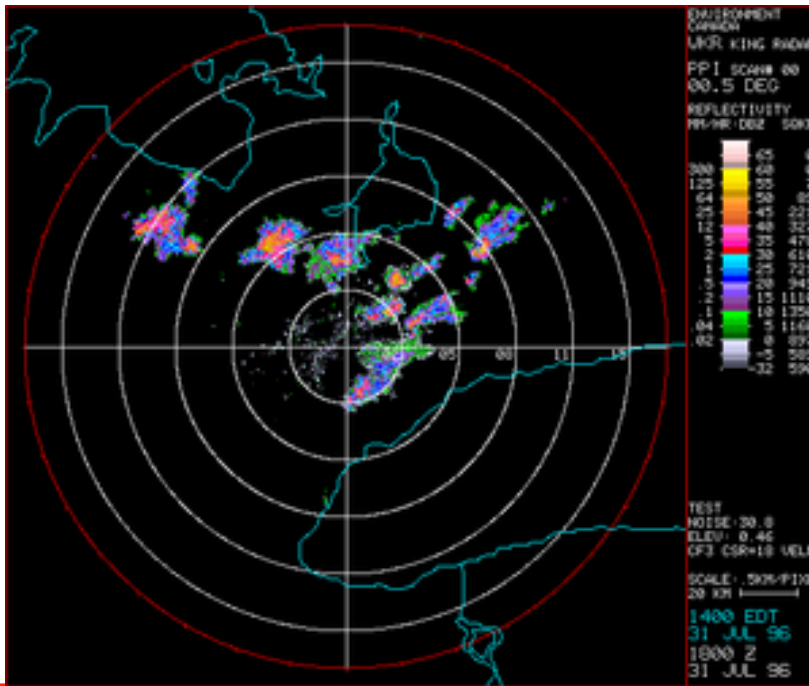


# Social Network (Depth=3)



# General application task 3: Big Data Analytics

- ⊙ Enabling the linking between structured and unstructured data
  - ⊙ Large-scale information monitoring
  - ⊙ Analytics: analyses of areas, markets, trends
  - ⊙ Watch: Scanning for relevant new developments



# Example: Network of Innovation Keyplayers

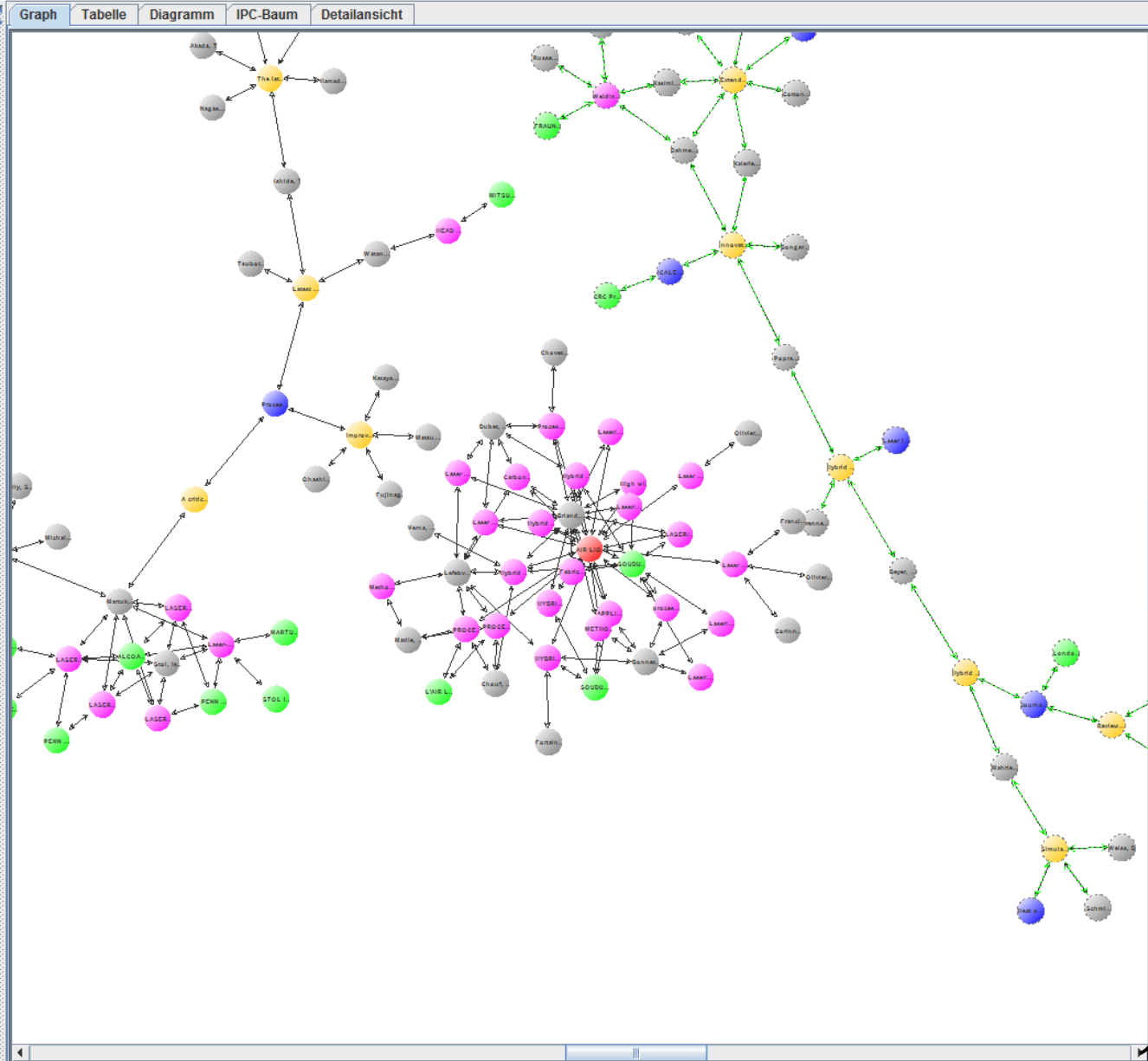
http://techwatchtool.dfki.de/thyssenkruppsteel/concept?queryId=33&queryName=hybrid+beam+welding&command=show

Ergebnisse für "hybrid beam v

- Patent (62)
- Publikation (43)
- Person (207)
- Organisation (52)
- AIR LIQUIDE
- ALCOA INC
- ALCOA INC.
- AMERICAN WELDING SOC
- AT-FACHVERLAG
- Aedermannsdorf, Switzerland: Tran
- American Welding Society
- BOEING CO

## AIR LIQUIDE

- All Patents in Database:
- [METHOD FOR LASER WELDING USING A NOZZLE CAPABLE OF STABILISING THE KEYHOLE](#) (2009-09-03)
  - [Laser Beam Welding Method with a Metal Vapour Capillary Formation Control](#) (2009-05-28)
  - [Laser beam welding method with a metal vapour capillary formation control](#) (2008-10-22)
  - [Process for laser-ARC hybrid welding aluminized metal workpieces](#) (2008-01-17)
  - [Laser arc hybrid welding method for surface coated metal parts, the surface coating containing aluminium](#) (2008-01-16)
  - [Laser/MIG hybrid welding process with a high wire speed](#) (2007-10-30)
  - [Laser or hybrid arc-laser welding with formation of a plasma on the back side.](#) (2006-11-29)
  - [LASER/MIG HYBRID WELDING PROCESS WITH A HIGH WIRE SPEED](#) (2006-05-19)
  - [Hybrid laser-Metal in Gas welding with a elevated welding and filler wire supply speeds and a high](#)



Legende

- Auswahl
- Patent
- Publikation
- Person
- Organisation
- Journal

Ansicht

Vergrößern (+)  
Verkleinern (-)  
Normalgröße  
Größe an Fenster a...

Graph anpassen

Fast Organic  
Circle

# Example in Opinion Mining

Mitten in der Euro-Krise geht **Altkanzler Helmut Kohl** mit **Angela Merkel** äußerst hart ins Gericht

-- Welt online, 25.08.2011

Opinion Holder

Opinion Target

Polarity



---

# A Brief History of IE



## Message Understanding Conference (MUC-6)

---

- ③ U.S. Government sponsored conferences with the intention to coordinate multiple research groups seeking to improve IE and IR technologies (since 1987)
- ③ defined several generic types of information extraction tasks (MUC Competition)
- ③ MUC 1-2 focused on automated analysis of military messages containing textual information
- ③ MUC 3-7 focused on information extraction from newswire articles
  - ③ terrorist events
  - ③ international joint-ventures
  - ③ management succession event

# Evaluation of IE systems in MUC

---

- ③ Participants receive description of the scenario along with the annotated training corpus in order to adapt their systems to the new scenario (1 to 6 months)
- ③ Participants receive new set of documents (test corpus) and use their systems to extract information from these documents and return the results to the conference organizer
- ③ The results are compared to the manually filled set of templates (answer key)

- ◎ precision and recall measures were adopted from the information retrieval research community

$$\textit{recall} = \frac{N_{\textit{correct}}}{N_{\textit{key}}} \qquad \textit{precision} = \frac{N_{\textit{correct}}}{N_{\textit{correct}} + N_{\textit{incorrect}}}$$

$$F = \frac{(\beta^2 + 1) \times \textit{precision} \times \textit{recall}}{\beta^2 \times \textit{precision} + \textit{recall}}$$

- ◎ Sometimes an F-measure is used as a combined recall-precision score

# Development Steps within IE Communities

---

- ③ from attempts to use the methods of full text understanding to shallow text processing;
- ③ from pure knowledge-based hand-coded systems to (semi-) automatic systems using machine learning methods;
- ③ from complex domain-dependent event extraction to standardized domain-independent elementary entity identification, simple semantic relation and event extraction

# New Evaluation Conferences

---

- ◎ New terms closely related to IE research:
  - ◎ ACE (Automatic Content Analysis), Text Analytics, Text Mining
- ◎ ACE: Automatic Content Extraction (ACE) Evaluation
  - ◎ <http://www.itl.nist.gov/iad/mig/tests/ace/>
- ◎ TAC: Text Analytics Conference
  - ◎ <http://www.nist.gov/tac/>
  - ◎ [Knowledge Base Population \(KBP\)](#)
  - ◎ [Biomedical Summarization \(BiomedSumm\)](#)



# Components of an IE Semantic Model (1)

---

## ⊙ Entities

- ⊙ Individuals in the world that are mentioned in a text

## ⊙ Simple entities

- ⊙ singular objects

## ⊙ Collective entities

- ⊙ sets of objects of the same type where the set is explicitly mentioned in the text

## ⊙ Relations

- ⊙ Properties that hold of tuples of entities

## ⊙ Complex Relations

- ⊙ Relations that hold among entities and relations

## ⊙ Attributes

- ⊙ one place relations are attributes or individual properties
-

## Components of an IE Semantic Model (2)

---

- ◎ Temporal points and intervals
- ◎ Relations may be timeless or bound to time intervals
- ◎ Events
  - ◎ A particular kind of simple or complex relation among entities involving a change in relation state at the end of a time interval.

# Relations in Time

---

- ⊙ time-dependent attribute:  $\text{age}(x)$
- ⊙ timeless two-place relation:  $\text{father}(x, y)$
- ⊙ time-dependent two-place relation:  $\text{boss}(x, y)$

# IE as a Semantic Analysis: Relating Language to the Model

---

## ⊙ Linguistic Mention

- ⊙ A particular linguistic phrase
- ⊙ Denotes a particular entity, relation, or event
  - A noun phrase, name, or possessive pronoun
  - A verb, nominalization, compound nominal, or other linguistic construct relating other linguistic mentions

## ⊙ Linguistic Entity

- ⊙ Equivalence class of mentions with same meaning
  - Coreferring noun phrases
  - Relations and events derived from different mentions, but conveying the same meaning

[Appelt, 2003]



# The Basic Semantic Tasks of an IE System

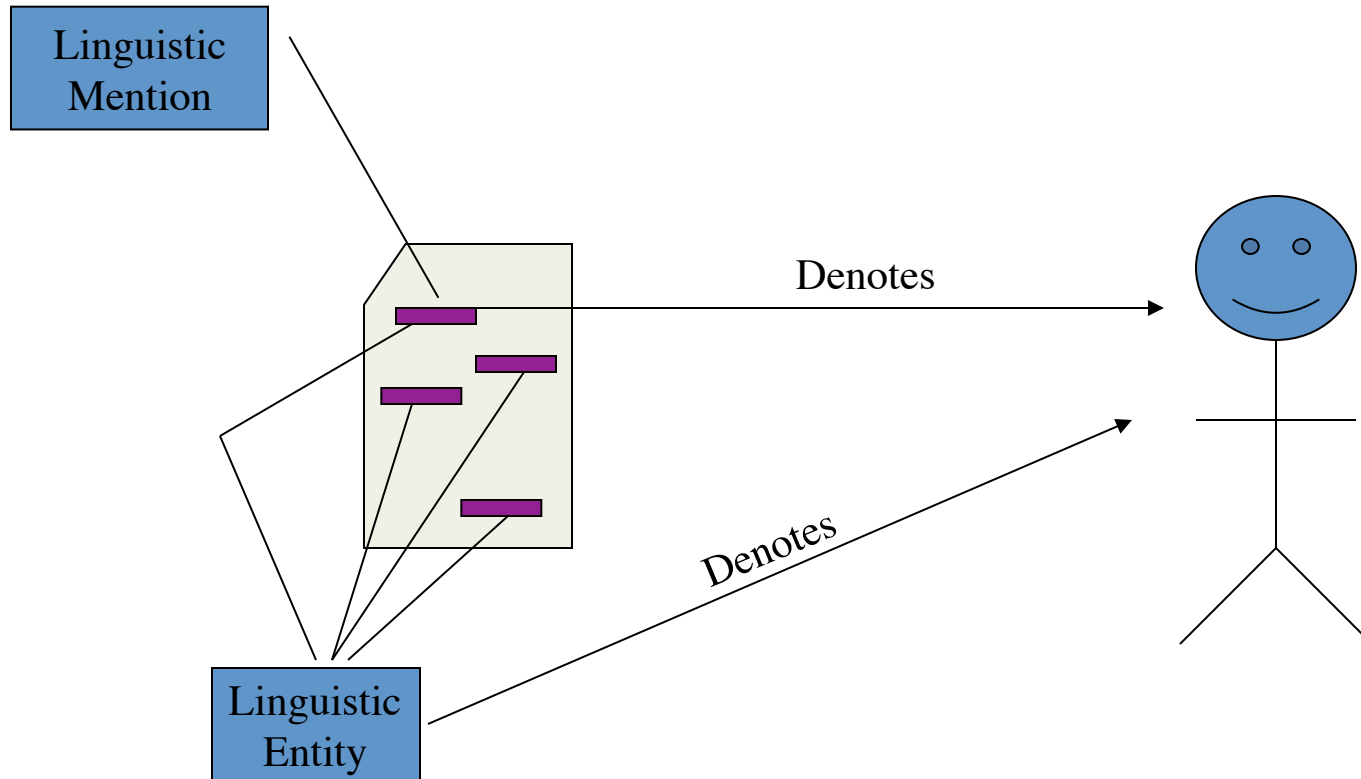
---

- ⊙ Recognition of linguistic mentions
- ⊙ Classification of linguistic mentions into semantic types
- ⊙ Identification of coreference equivalence classes of linguistic entities
- ⊙ Identifying the actual individuals that are mentioned in an article
  - ⊙ Associating linguistic entities with predefined individuals (e.g. a database, or knowledge base)
  - ⊙ Forming equivalence classes of linguistic entities from different documents.

[Appelt, 2003]

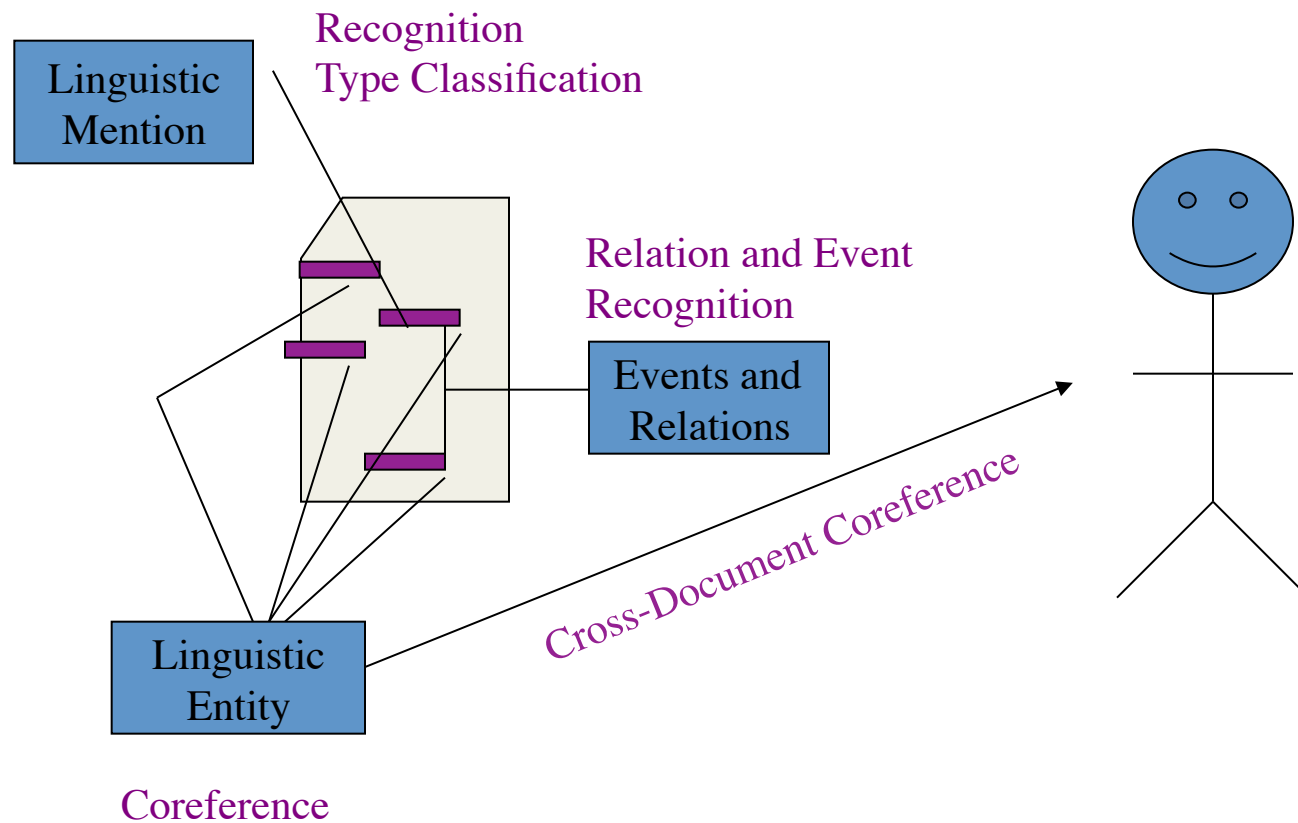


# Language and World Model



[Appelt, 2003]

# NLP Tasks in an Extraction System



[Appelt, 2003]

# Types of Linguistic Mentions

---

## ⊙ Name mentions

- ⊙ The mention uses a proper name to refer to the entity

## ⊙ Nominal mentions

- ⊙ The mention is a noun phrase whose head is a common noun

## ⊙ Pronominal mentions

- ⊙ The mention is a headless noun phrase, or a noun phrase whose head is a pronoun, or a possessive pronoun

## Example of Linguistic Mentions

---

1. Three of the Nobel Prizes for Chemistry during the first decade **were awarded** for pioneering work in organic chemistry.
2. In **1902** **Emil Fischer** (1852-1919), then in Berlin, **was given** the prize for **his** work on sugar and purine syntheses.
3. Another major influence from organic chemistry was the development of the chemical industry, and a chief contributor here was Fischer's teacher, **Adolf von Baeyer** (1835-1917) in Munich, **who was awarded** the prize in **1905**.

# Example

A relation extraction task in the domain *management succession* (MUC-6)

< person\_in, person\_out, position, organisation >

- *person\_in*: the person who obtained the position
- *person\_out*: the person who left the position
- *position*: the job position that the two persons were involved in
- *organisation*: the organisation where the position was located

According to CEO Fred Bell, Acme Inc. hired Peter Smith as COO yesterday, replacing Hans Bloggs.

According to CEO Fred Bell, Acme Inc. hired Peter Smith as COO yesterday, replacing Hans Bloggs.

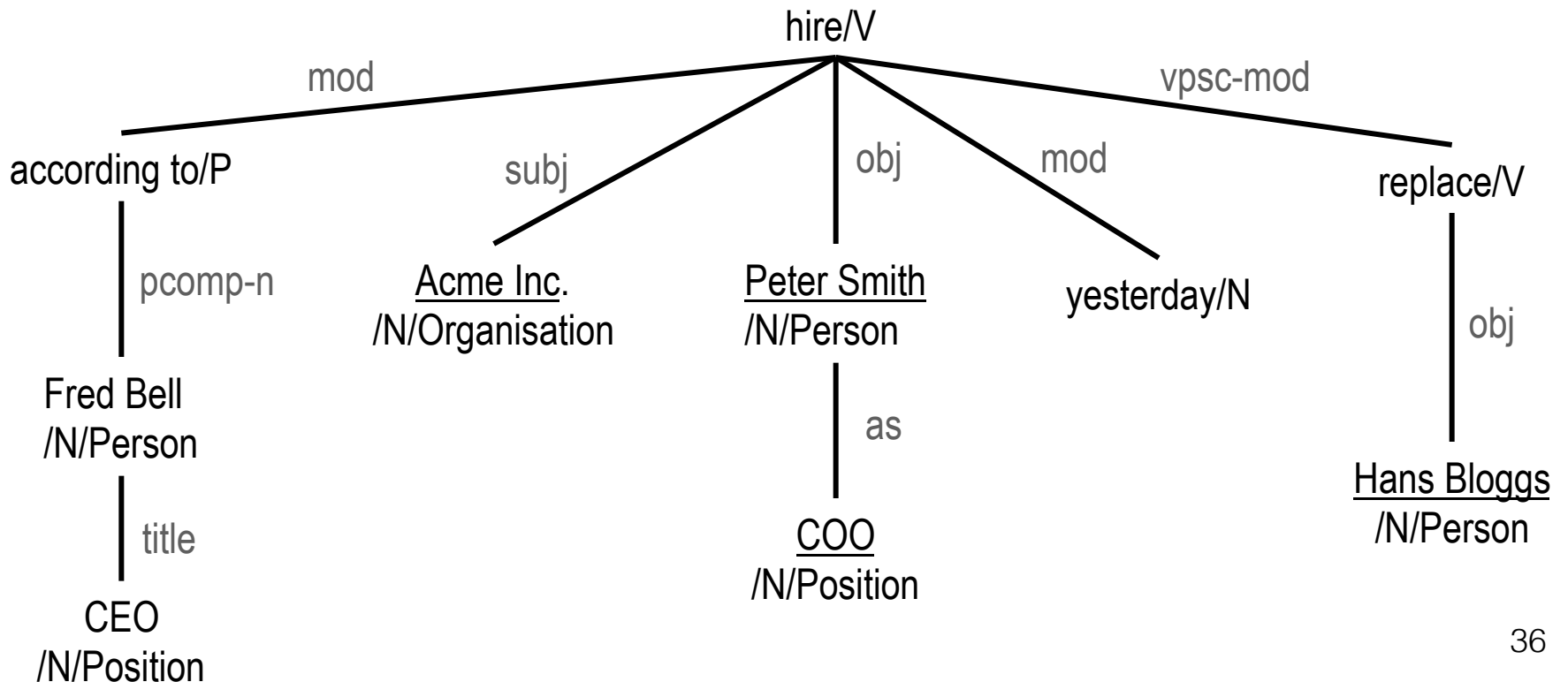
<person\_in, person\_out, position, organisation>

According to CEO Fred Bell, Acme Inc. hired Peter Smith as COO yesterday, replacing Hans Bloggs.

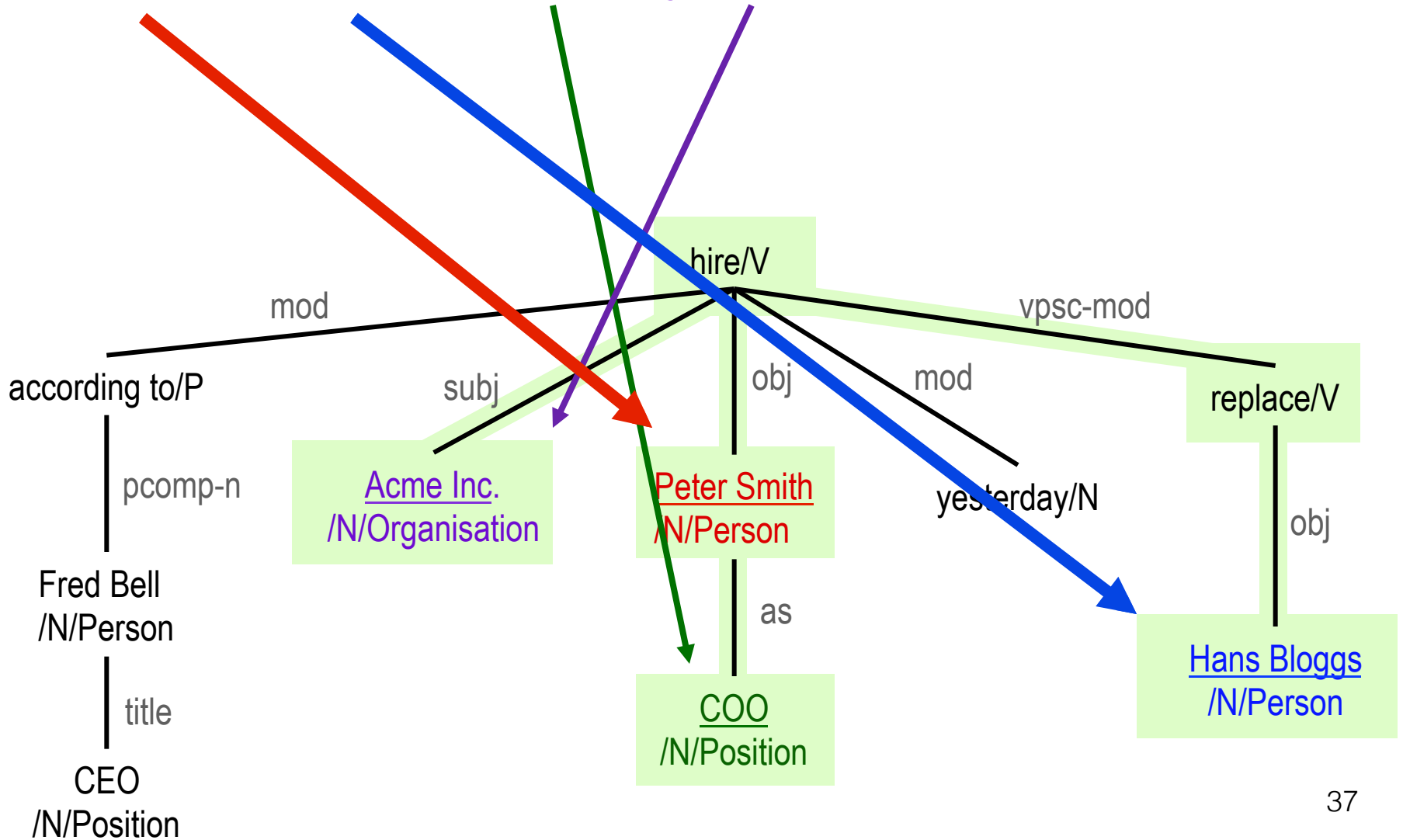
<person\_in, person\_out, position, organisation>

According to CEO Fred Bell, Acme Inc. hired Peter Smith as COO yesterday, replacing Hans Bloggs.

<person\_in, person\_out, position, organisation>



<person\_in, person\_out, position, organisation>



---

# Machine Learning for Relation Extraction



# Motivations of ML

---

- ⊙ Porting to new domains or applications is expensive
- ⊙ Current technology requires IE experts
  - ⊙ Expertise difficult to find on the market
  - ⊙ SME cannot afford IE experts
- ⊙ Machine learning approaches
  - ⊙ Domain portability is relatively straightforward
  - ⊙ System expertise is not required for customization
  - ⊙ “Data driven” rule acquisition ensures full coverage of examples

## Problems

---

- ③ Training data may not exist, and may be very expensive to acquire
- ③ Large volume of training data may be required
- ③ Changes to specifications may require reannotation of large quantities of training data
- ③ Understanding and control of a domain adaptive system is not always easy for non-experts

# Parameters of IE Real-World Tasks

---

- ⊙ Document structure
  - ⊙ Free text
  - ⊙ Semi-structured
  - ⊙ Structured
- ⊙ Linguistic annotation
  - ⊙ Shallow NLP
  - ⊙ Deep NLP
- ⊙ Complexity and specificity of relation
  - ⊙ Unary
  - ⊙ N-ary
- ⊙ Depth of extraction
  - ⊙ Recognition
  - ⊙ Classification
  - ⊙ Semantic role labelling
- ⊙ Degree of automation
  - ⊙ Semi-automatic
  - ⊙ Supervised
  - ⊙ Semi-Supervised
  - ⊙ Minimally-Supervised
  - ⊙ Distant Supervision
  - ⊙ Unsupervised
- ⊙ Human interaction/contribution
- ⊙ Data properties
  - ⊙ Domain relevance
  - ⊙ Redundancy
  - ⊙ Connectivity
- ⊙ Evaluation/validation
  - ⊙ With/without gold standard
  - ⊙ Performance: recall & precision
  - ⊙ Interaction among parameters

---

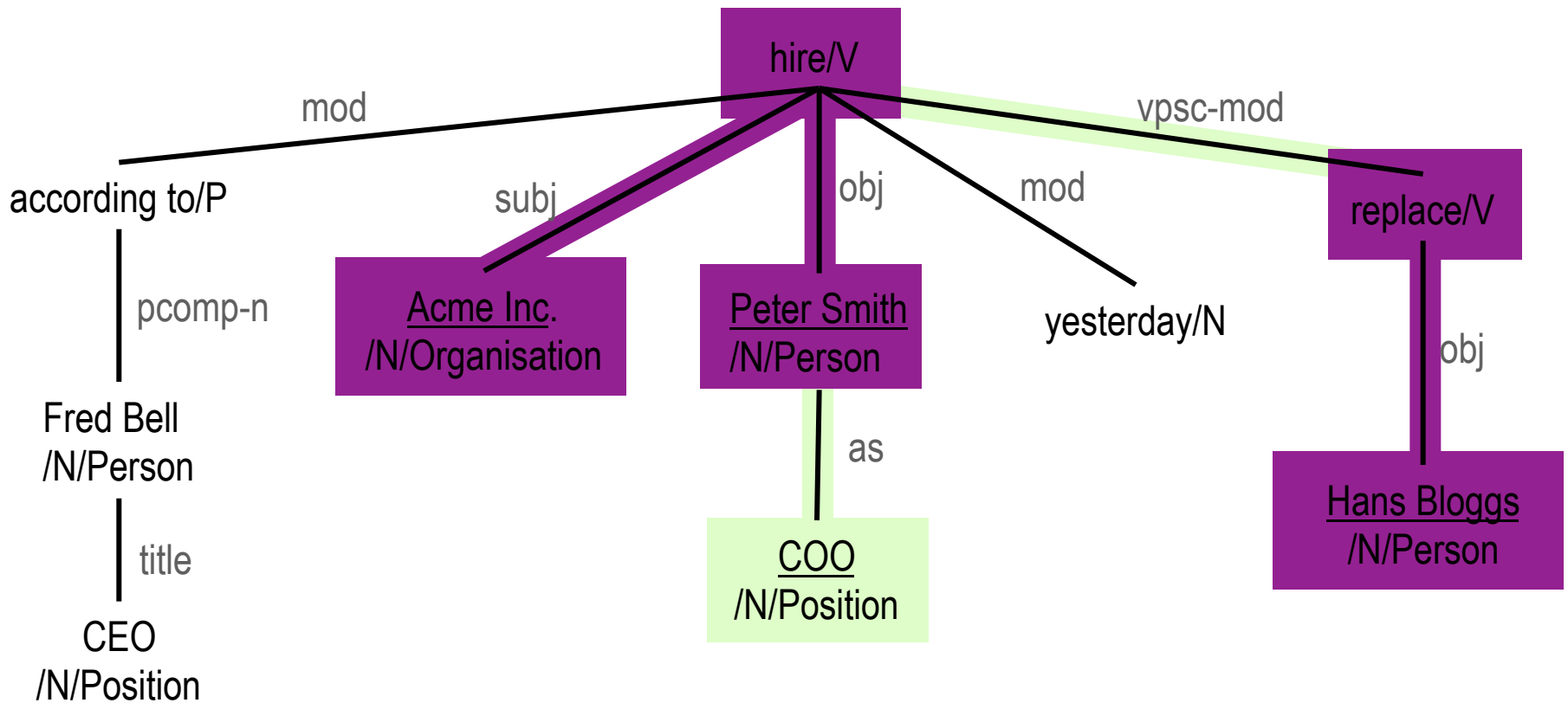
# State of the Art



# Binary Relation Only Approaches

---

- ⊙ Extraction of binary relations only such as
  - ⊙ *author-book*
  - ⊙ *company-location*
- ⊙ Do not employ the existing syntactic and semantic structures among  $n > 2$  arguments and rely on a later component to merge binary relations into complex relations
- ⊙ Approaches
  - ⊙ Ravichandran and Hovy, 2002
  - ⊙ Pantel et al., 2004
  - ⊙ Pasca et al., 2006a; Pasca et al., 2006b



# Surface-oriented Rule Representation

---

- ◎ Shallow linguistic analyses
- ◎ These formalisms are robust and efficient but only handle binary relations.
- ◎ Work best for relations whose arguments usually co-occur in close proximity within a sentence and whose mentions exhibit limited linguistic variation
- ◎ Approaches
  - ◎ Pasca et al., 2006a; Pasca et al., 2006b;
  - ◎ Kozareva et al., 2008; Hovy et al., 2009; Kozareva and Hovy, 2010

# Minimally Supervised (Bootstrapping)

---

- ◎ Based on iterative learning with limited initial knowledge
  - ◎ Start a small number of initial examples of relation instances (or patterns)
  - ◎ Label the free texts during iterations (e.g., Agichtein and Gravano, 2000; Yangarber et al., 2000; Ravichandran and Hovy, 2002; Stevenson and Greenwood, 2005).
- ◎ Often suffer from semantic drift or the propagation of errors occurring during iterations
- ◎ Performance depends on data properties

## Distant Supervision

---

- ⊙ A massive seed-based, one step version of bootstrapping
- ⊙ Rely on a large amount of trustworthy facts
- ⊙ Their performance does not hinge on corpus data properties such as redundancy
- ⊙ Approaches
  - ⊙ (Mintz et al., 2009)
  - ⊙ Others: (Wu and Weld, 2007); (Wu et al., 2008); (Weld et al., 2008); (Hoffmann et al., 2010); (Xu et al., 2011); (Nguyen and Moschitti, 2011)

- ⊙ They do not target given relations.
- ⊙ They are very useful for applications continuously faced with new relation or event types, e.g., online social media monitoring
- ⊙ However, the results of these systems cannot be directly taken for filling knowledge databases, because the semantics of the new relations including the roles of the entities remains unknown.
- ⊙ Example: TextRunner (Banko et al., 2007; Yates et al., 2007)

## Reality in IE Projects

---

- ◎ Our IE users are often not domain experts
- ◎ IE experts have to develop methods and strategies for
  - ◎ Prospecting a domain
  - ◎ Proposing relevant relations
  - ◎ Finding relevant and suitable data