

# Computerlinguistik II

Vorlesung im SoSe 2019  
(M-GSW-10)

**Prof. Dr. Udo Hahn**

Lehrstuhl für Computerlinguistik  
Institut für Germanistische Sprachwissenschaft  
Friedrich-Schiller-Universität Jena

<http://www.julielab.de>

# Two Paradigms for NLP

- Symbolic Specification Paradigm
  - Manual acquisition procedures
  - Lab-internal activities
  - Intuition and (few!) subjectively generated examples drive progress based on individual (competence) judgments
    - “I have a system that parses all of my nine-teen sentences!”

# Symbolic Specification Paradigm

- **Manual rule specification**
  - Source: linguist's intuition
- **Manual lexicon specification**
  - Source: linguist's intuition
- **Each lab has its own (home-grown) set of NLP software**
  - Hampers reusability
  - Limits scientific progress
  - Waste of human and monetary resources (we “burnt” thousands of Ph.D. student all over the world ☹)

# Shortcomings of the “Classical” Linguistic Approach

- Huge amounts of background knowledge req.
  - Lexicons (approx. 100,000 – 150,000 entries)
  - Grammars (>> 15,000 – 20,000 rules)
  - Semantics (>> 15,000 – 20,000 rules)
- As the linguistic and conceptual coverage of classical linguistic systems increases (slowly), it still remains insufficient; systems also reveal ‘spurious’ ambiguity, and, hence, tend to become overly “brittle” and unmaintainable
- More fail-soft behavior is required at the expense of ... ? (e.g., full-depth understanding)

# Two Paradigms for NLP

## • Symbolic Specification Paradigm

- Manual acquisition procedures
- Lab-internal activities
- Intuition and (few!) subjectively generated examples drive progress based on individual (competence) judgments
  - “I have a system that parses all of my nine-teen sentences!”

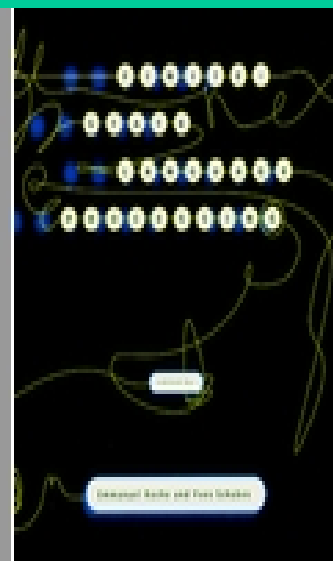
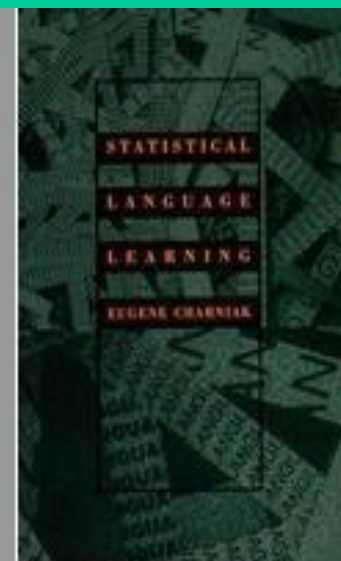
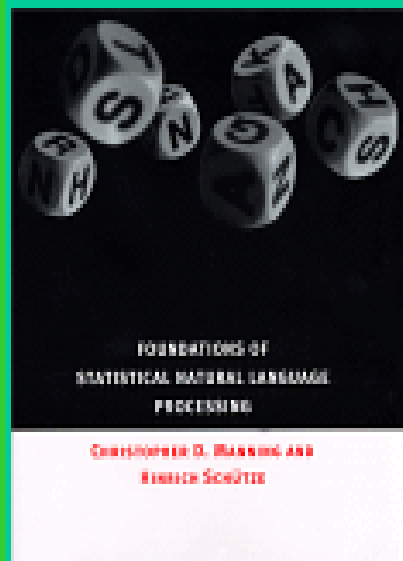
## • Empirical (Learning) Paradigm

- Automatic acquisition procedures
- Community-wide sharing of common knowledge and resources
- Large and ‘representative’ data sets drive progress according to experimental standards
  - “The system was tested on 1,7 million words taken from the WSJ segment of the MUC-7 data set and produced 4.9% parsing errors, thus yielding a statistically significant 1.6% improvement over the best result by parser X on the same data set & a 40.3% improvement over the baseline system!”

# Empirical Paradigm

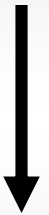
- Large repositories of language data
  - Corpora (plain or annotated, i.e., enriched by meta-data)
- Large, community-wide shared repositories of language processing modules
  - Tokenizers, POS taggers, chunkers, NE recognizers, ...
- Shared repositories of machine learning algos
- Automatic acquisition of linguistic knowledge
  - Applying ML algos to train linguistic processors by using large corpora with valid linguistic metadata (linguist as educated data supplier, „language expert“) rather than manual intuition (linguist as creative rule inventor)
- Shallow analysis rather than deep understanding
- Large, community-wide self-managed, task-oriented competitions, comparative evaluation rounds
- Change of mathematics:
  - Statistics rather than algebra and logics

# Paradigm Shift – We Exchanged our Textbooks...



# POS Tagging

A severe infection ended the pregnancy .



DET ADJ NOUN VERB DET NOUN ST



# Penn Treebank Tag Set

Tag	Description	Examples
.	sentence terminator	. ! ?
DT	determiner	all an many such that the them these this
JJ	adjective, numeral	first oiled separable battery-powered
NN	common noun	cabbage thermostat investment
PRP	personal pronoun	herself him it me one oneself theirs they
IN	preposition	among out within behind into next
VB	verb (base form)	ask assess assign begin break bring
VBD	verb (past tense)	asked assessed assigned began broke
WP	WH-pronoun	that what which who whom

In total,  
45 tags

# Transformation Rules for Tagging [Brill, 1995]

- Initial State: Based on a number of features, guess the most likely POS tag for a given word:
  - die/DET Frau/NOUN ,/COMMA die/DET singt/VFIN
- Learn transformation rules to reduce errors:
  - *Change DET to PREL whenever the preceding word is tagged as COMMA*
- Apply learned transformation rules:
  - die/DET Frau/NOUN,/COMMA die/PREL singt/VFIN

# First 20 Transformation Rules

#	Change Tag		Condition
	From	To	
1	NN	VB	Previous tag is <i>TO</i>
2	VBP	VB	One of the previous three tags is <i>MD</i>
3	NN	VB	One of the previous two tags is <i>MD</i>
4	VB	NN	One of the previous two tags is <i>DT</i>
5	VBD	VBN	One of the previous three tags is <i>VBZ</i>
6	VBN	VBD	Previous tag is <i>PRP</i>
7	VBN	VBD	Previous tag is <i>NNP</i>
8	VBD	VBN	Previous tag is <i>VBD</i>
9	VBP	VB	Previous tag is <i>TO</i>
10	POS	VBZ	Previous tag is <i>PRP</i>
11	VB	VBP	Previous tag is <i>NNS</i>
12	VBD	VBN	One of previous three tags is <i>VBP</i>
13	IN	WDT	One of next two tags is <i>VB</i>
14	VBD	VBN	One of previous two tags is <i>VB</i>
15	VB	VBP	Previous tag is <i>PRP</i>
16	IN	WDT	Next tag is <i>VBZ</i>
17	IN	DT	Next tag is <i>NN</i>
18	JJ	NNP	Next tag is <i>NNP</i>
19	IN	WDT	Next tag is <i>VBD</i>
20	JJR	RBR	Next tag is <i>JJ</i>

Taken from: Brill (1995), Transformation-Based Error-Driven Learning

# Towards Statistical Models of Natural Language Processing ...

# Letter-based Language Models

- **Shannon's Game**
- **Guess the next letter:**
-

# Letter-based Language Models

- **Shannon's Game**
- **Guess the next letter:**
- **W**

# Letter-based Language Models

- **Shannon's Game**
- **Guess the next letter:**
- **Wh**

# Letter-based Language Models

- **Shannon's Game**
- **Guess the next letter:**
- **Wha**



# Letter-based Language Models

- **Shannon's Game**
- **Guess the next letter:**
- **What**

# Letter-based Language Models

- **Shannon's Game**
- **Guess the next letter:**
- **What d**

# Letter-based Language Models

- **Shannon's Game**
- **Guess the next letter:**
- **What do**

# Letter-based Language Models

- **Shannon's Game**
- **Guess the next letter:**
- **What do you think the next letter is?**

# Word-based Language Models

- **Shannon's Game**
- Guess the next letter:
  - What do you think the next letter is?
- Guess the next word:
-

# Word-based Language Models

- **Shannon's Game**
- Guess the next letter:
  - What do you think the next letter is?
- Guess the next word:
  - We

# Word-based Language Models

- **Shannon's Game**
- Guess the next letter:
  - What do you think the next letter is?
- Guess the next word:
  - We are

# Word-based Language Models

- **Shannon's Game**
- Guess the next letter:
  - What do you think the next letter is?
- Guess the next word:
  - We are now



# Word-based Language Models

- **Shannon's Game**
- Guess the next letter:
  - What do you think the next letter is?
- Guess the next word:
  - We are now entering

# Word-based Language Models

- **Shannon's Game**
- Guess the next letter:
  - What do you think the next letter is?
- Guess the next word:
  - We are now entering statistical

# Word-based Language Models

- **Shannon's Game**
- Guess the next letter:
  - What do you think the next letter is?
- Guess the next word:
  - We are now entering statistical territory

# Approximating Natural Language Words

- **zero-order approximation:**  
letter sequences are independent of  
each other and all equally probable:
  - xfoml rxkhrjffjuj zlpwcwkcy  
ffjeyvkcqsghyd

# Approximating Natural Language Words

- **first-order approximation:**  
letters are independent, but occur  
with the frequencies of English text:
  - ocro hli rgwr nmielwis eu ll  
nbnesebya th eei alhenhtppa oobttva  
nah

# Approximating Natural Language Words

- **second-order approximation:**  
the probability that a letter appears depends on the previous letter
  - on ie antsoutinys are t inctore st bes  
deamy achin d ilonasive tucoowe at  
teasonare fuzo tizin andy tobe seace  
ctisbe

# Approximating Natural Language Words

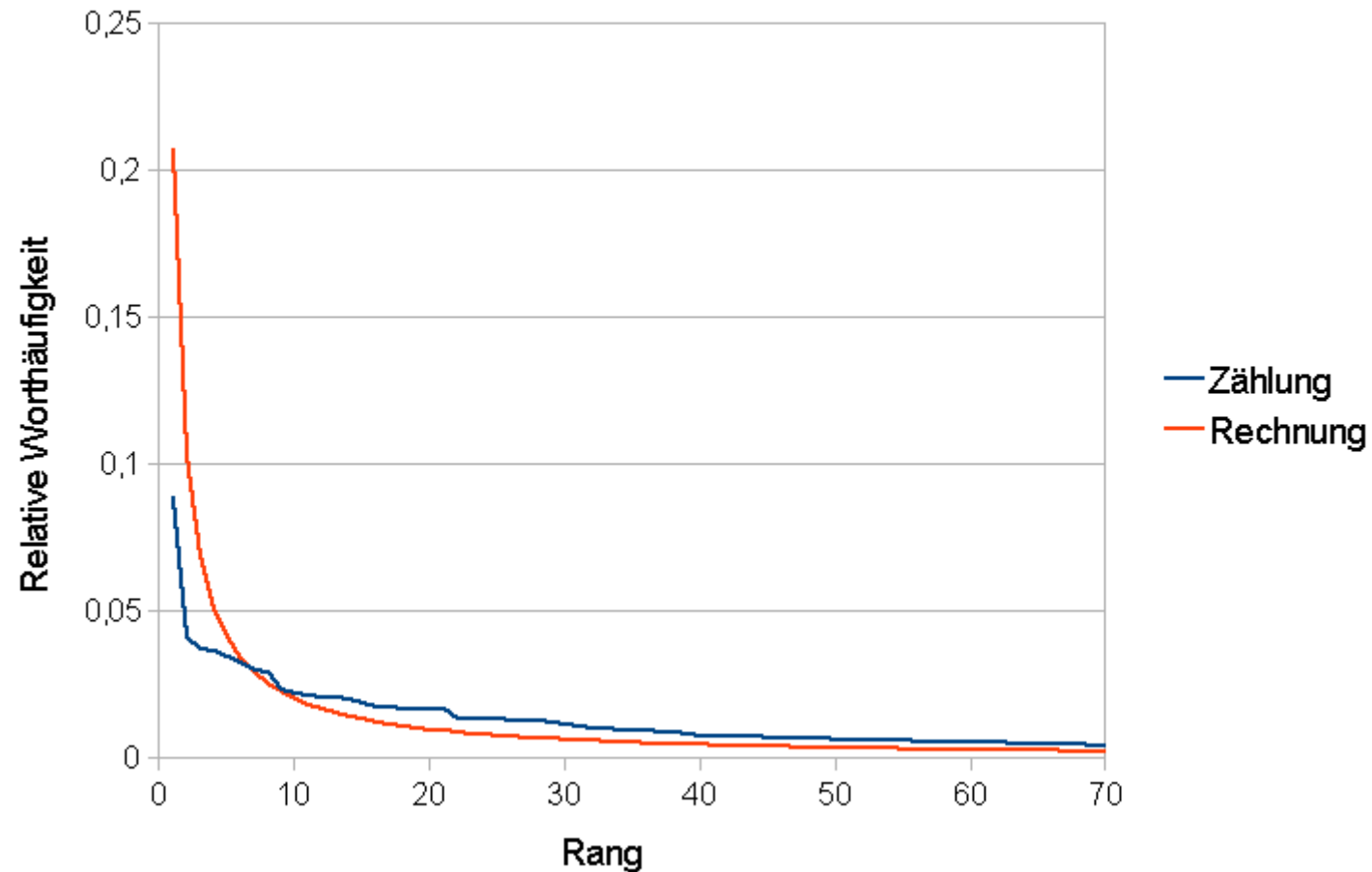
- **third-order approximation:**  
the probability that a certain letter appears depends on the two previous letters
  - in no ist lat whey cratict froure birs  
grocid pondenome of demonstures  
of the reptagin is regoactiona of cre

# Approximating Natural Language Words

- Higher frequency trigrams for different languages:
  - English: THE, ING, ENT, ION
  - German: EIN, ICH, DEN, DER
  - French: ENT, QUE, LES, ION
  - Italian: CHE, ERE, ZIO, DEL
  - Spanish: QUE, EST, ARA, ADO



# Zipfsches Gesetz



Wortverteilung im Vergleich zu einer einfachen Zipf-Verteilung ( $\sim 1/n$ . Wortanzahl: 70;  
Texte aus: <http://www.gutenberg.org/dirs/etext04/8effi10.txt>)

# Terminology

- **Sentence:** unit of written language
- **Utterance:** unit of spoken language
- **Word Form:** the inflected form that appears literally in the corpus
- **Lemma:** lexical forms having the same stem, part of speech, and word sense
- **Types (V):** number of distinct words that might appear in a corpus (vocabulary size)
- **Tokens ( $N_T$ ):** total number of words in a corpus (note:  $V \ll N_T$ )
- **Types seen so far (T):** number of distinct words seen so far in corpus (note:  $T \leq V \ll N_T$ )

# Word-based Language Models

- A model that enables one to compute the probability, or likelihood, of a sentence  $S$ ,  $P(S)$ .
- Simple: Every word follows every other word with equal probability (0-gram)
  - Assume  $|V|$  is the size of the vocabulary  $V$
  - Likelihood of sentence  $S$  of length  $n$  is  $1/|V| \times 1/|V| \dots \times 1/|V|$
  - If English has 100,000 words, the probability of each next word is  $1/100000 = .00001$

# Relative Frequency vs. Conditional Probability

- Smarter: *Relative* Frequency

Probability of each next word is related to word frequency within a corpus (unigram)

- Likelihood of sentence  $S = P(w_1) \times P(w_2) \times \dots \times P(w_n)$
- Assumes probability of each word is independent of probabilities of other words

# Relative Frequency vs. Conditional Probability

- Smarter: *Relative* Frequency

Probability of each next word is related to word frequency within a corpus (unigram)

- Likelihood of sentence  $S = P(w_1) \times P(w_2) \times \dots \times P(w_n)$
- Assumes probability of each word is independent of probabilities of other words

- Even smarter: *Conditional* Probability

Look at probability given previous words (n-gram)

- Likelihood of sentence  $S = P(w_1) \times P(w_2|w_1) \times \dots \times P(w_n|w_{n-1})$
- Assumes probability of each word is dependent on probabilities of previous words

# Generalization of Conditional Probability via Chain Rule

- Conditional Probability for Two Events,  $A_1$  and  $A_2$ 
  - $P(A_1, A_2) = P(A_1) \cdot P(A_2|A_1)$
- **Chain Rule** generalizes to multiple ( $n$ ) events
  - $P(A_1, \dots, A_n) =$   
$$P(A_1) \times P(A_2|A_1) \times P(A_3|A_1, A_2) \times \dots \times P(A_n|A_1 \dots A_{n-1})$$
  - Examples:
    - $P(\text{the dog}) = P(\text{the}) \times P(\text{dog} | \text{the})$
    - $P(\text{the dog bites}) = P(\text{the}) \times P(\text{dog} | \text{the}) \times P(\text{bites} | \text{the dog})$

# Relative Frequencies and Conditional Probabilities

- Relative word frequencies are better than equal probabilities for all words
  - In a corpus with 10K word types, each word would have  $P(w) = 1/10K$
  - Does not match our intuitions that different words are more likely to occur
    - (e.g. “the” vs. “shop” vs. “aardvark”)
- Conditional probability is more useful than individual relative word frequencies
  - **dog** may be relatively rare in a corpus
  - but if we see **barking**,  $P(\text{dog}|\text{barking})$  may be large

# Probability for a Word String

- In general, the probability of a complete string of words  $w_1^n = w_1 \dots w_n$  is

$$\begin{aligned} P(w_1^n) \\ &= P(w_1)P(w_2/w_1)P(w_3/w_1 \ w_2) \dots P(w_n/w_1 \dots w_{n-1}) \\ &= \prod_{k=1}^n P(w_k | w_1^{k-1}) \end{aligned}$$

- But this approach to determining the probability of a word sequence gets to be computationally very expensive and suffers from sparse data



# Markov Assumption (basic idea)

- How do we (efficiently) compute  $P(w_n | w_1^{n-1})$ ?
- Trick (!): Instead of  $P(\text{rabbit} | \text{I saw } \underline{a})$ , we use  $P(\text{rabbit} | \underline{a})$ .
  - This lets us collect statistics in practice via a bigram model:  $P(\text{the barking dog}) = P(\text{the} | \text{<start>}) \times P(\text{barking} | \text{the}) \times P(\text{dog} | \text{barking})$

# Markov Assumption (the very idea)

- Markov models are the class of probabilistic language models that assume that we can predict the probability of some future unit *without looking too far* into the past
  - Specifically, for  $N=2$  (bigram):
    - $P(w_1^n) \approx \prod_{k=1}^n P(w_k | w_{k-1}); w_0 := \text{<start>}$
- Order of a Markov model: length of prior context
  - bigram is first order, trigram is second order, ...

# Statistical HMM-based Tagging

[Brants, 2000]

- *State transition probability*. Likelihood of a tag immediately following n other tags
    - $P_1(\text{Tag}_i \mid \text{Tag}_{i-1} \dots \text{Tag}_{i-n})$
  - *State emission probability*. Likelihood of a word given a tag
    - $P_2(\text{Word}_i \mid \text{Tag}_i)$
- 
- die/DET Frau/NOUN ,/COMMA die/DET or PREL singt/VFIN

# Trigrams for Tagging

- *State transition probabilities (trigrams):*

- $P_1(\text{DET} \mid \text{COMMA NOUN}) = 0.0007$

- $P_1(\text{PREL} \mid \text{COMMA NOUN}) = 0.01$

- *State emission probabilities:*

- $P_2(\text{die} \mid \text{DET}) = 0.7$

- $P_2(\text{die} \mid \text{PREL}) = 0.2$

Taken from  
(POS-  
annotated)  
corpora

- Compute probabilistic evidence for the tag being

- **DET:**  $P_1 \cdot P_2 = 0.0007 \cdot 0.7 = 0.00049$

- **PREL:**  $P_1 \cdot P_2 = 0.01 \cdot 0.2 = 0.002$

• die/DET Frau/NOUN ,/COMMA die/PREL singt/VFIN

# Inside (most) POS Taggers

- Lexicon look-up routines
- Morphological processing (not only deflection!)
- Unknown word handler, if lexicon look-up fails (based on statistical information)
- Ambiguity ranking (priority selection)

# Chunking

Arginine methylation of STAT1 modulates IFN induced transcription

# Chunking

[Arginine methylation] of [STAT1] modulates [IFN induced transcription]

# Shallow Parsing

[Arginine methylation of STAT1]<sub>NP</sub> [modulates]<sub>VP</sub> [IFN induced transcription]<sub>NP</sub>



# Shallow Parsing

[ [Arginine methylation]<sub>NP</sub> [of STAT1]<sub>PP</sub> ]<sub>NP</sub>

[Arginine methylation of STAT1]<sub>NP</sub> [modulates]<sub>VP</sub> [IFN induced transcription]<sub>NP</sub>

# Shallow Parsing

[ [IFN induced]<sub>AP</sub> [transcription]<sub>N</sub> ]<sub>NP</sub>

[ [Arginine methylation]<sub>NP</sub> [of STAT1]<sub>PP</sub> ]<sub>NP</sub>

[Arginine methylation of STAT1]<sub>NP</sub> [modulates]<sub>VP</sub> [IFN induced transcription]<sub>NP</sub>

# Deep Parsing

[ [IFN induced]<sub>AP</sub> [transcription]<sub>N</sub> ]<sub>NP</sub>

[ [ [Arginine]<sub>N</sub> [methylation]<sub>N</sub> ]<sub>NP</sub> [ [of]<sub>P</sub> [STAT1]<sub>N</sub> ]<sub>PP</sub> ]<sub>NP</sub>

[ [Arginine methylation]<sub>NP</sub> [of STAT1]<sub>PP</sub> ]<sub>NP</sub>

[Arginine methylation of STAT1]<sub>NP</sub> [ [modulates]<sub>V</sub> [IFN induced transcription]<sub>NP</sub> ]<sub>VP</sub>

# Deep Parsing

[ [[IFN]<sub>N</sub> [induced]<sub>A</sub>]<sub>AP</sub> [transcription]<sub>N</sub> ]<sub>NP</sub>

[ [IFN induced]<sub>AP</sub> [transcription]<sub>N</sub> ]<sub>NP</sub>

[ [[Arginine]<sub>N</sub> [methylation]<sub>N</sub>]<sub>NP</sub> [[of]<sub>P</sub> [STAT1]<sub>N</sub>]<sub>PP</sub> ]<sub>NP</sub>

[ [Arginine methylation]<sub>NP</sub> [of STAT1]<sub>PP</sub> ]<sub>NP</sub>

[Arginine methylation of STAT1]<sub>NP</sub> [ [modulates]<sub>V</sub> [IFN induced transcription]<sub>NP</sub> ]<sub>VP</sub>

# Chunking Principles

- Goal: divide a sentence into a sequence of chunks (also phrases)
- Chunks are non-overlapping regions of a text
  - *[I] saw [a tall man] in [the park]*
- Chunks are non-exhaustive
  - not all words of a sentence are included in chunks
- Chunks are non-recursive
  - a chunk does not contain other chunks
- Chunks are mostly base NP chunks

[ *[the synthesis]*<sub>NP-base</sub> *of* *[long enhancer transcripts]*<sub>NP-base</sub> ]<sub>NP-complex</sub>

# The Shallow Syntax Pipeline

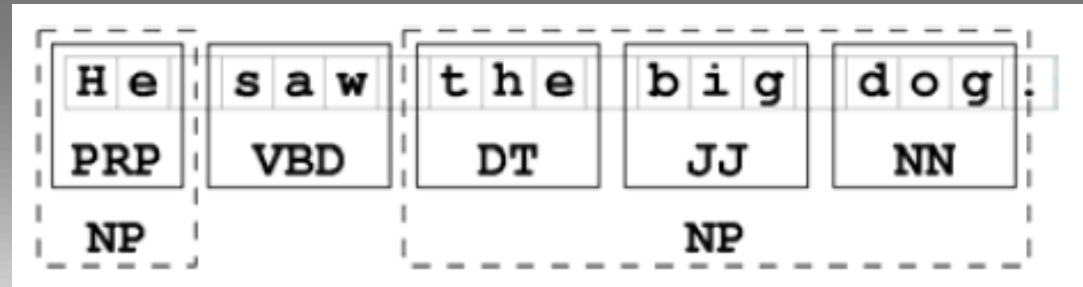
Tagging



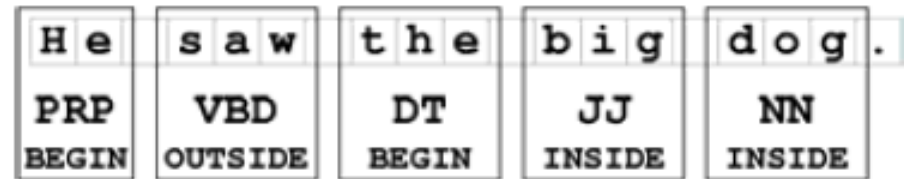
Chunking



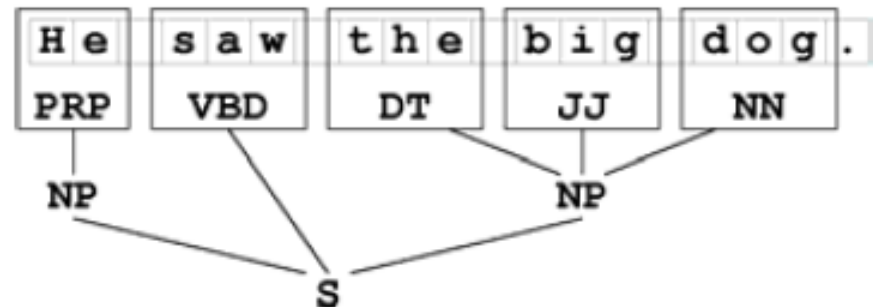
Parsing



BIO (or IOB)



Trees



# BIO Format for Base NPs

a	DT	B
mechanism	NN	I
that	WDT	B
increases	VBZ	O
NF-kappa	NN	B
B/I	NN	I
kappa	NN	I
B	NN	I
dissociation	NN	I
without	IN	O
affecting	VBG	O
the	DT	B
NF-kappa	NN	I
B	NN	I
translocation	NN	I
step	NN	I

# A Simple Chunking Technique

- Simple chunkers usually ignore lexical content
  - Only need to look at part-of-speech tags
- Basic steps in chunking
  - Chunking / Unchunking
  - Chinking
  - Merging / Splitting



# Regular Expression Basics

- “|” OR operator (explicit OR-ing)
  - “[a|e|i|o|u]” matches any occurrence of vowels
- “[abc]” matches any occurrence of either “a”, “b” or “c” (implicit OR-ing)
  - “gr[ae]y” matches “grey” or “gray” (but not “graey”)
- “.” matches arbitrary char
  - “d.g” matches “dag”, “dig”, “dog”, “dkg” ...
- “?” preceding expression/char may or may not occur
  - “colou?r” matches “colour” and “color”
- “+” preceding expression occurs at least one time
  - “(ab)+” matches “ab”, “abab”, “ababab”, ...
- “\*” preceding expression occurs null time or arbitrary often
  - “(ab)\*” matches “\_”, “ab”, “abab”, “ababab”, ...

# Chunking

- Define a regular expression that matches the sequences of tags in a chunk
  - `<DT>? <JJ>* <NN.??>`
- Chunk all matching subsequences
  - *A/DT red/JJ car/NN ran/VBD on/IN the/DT street/NN*
  - *[A/DT red/JJ car/NN] ran/VBD*  
*on/IN [the/DT street/NN]*
- If matching subsequences overlap, the first one gets priority
- **Unchunking** is the opposite of chunking

# Chinking

- A chink is a subsequence of the text that is not a chunk
- Define a regular expression that matches the sequences of tags in a chink
  - ( <VB.??> | <IN> )+
- Chunk anything that is not a matching subsequence
  - *A/DT red/JJ car/NN ran/VBD on/IN the/DT street/NN*
  - [*A/DT red/JJ car/NN*]  

*ran/VBD on/IN [the/DT street/NN]*

  
*chink*

# Merging

- Combine adjacent chunks into a single chunk
- Define a regular expression that matches the sequences of tags on both sides of the point to be merged
  - Merge a chunk ending in “JJ” with a chunk starting with “NN”, i.e. left: <JJ>, right: <NN.>
- Chunk all matching subsequences
  - [A/DT *red*/JJ ] [ *car*/NN] *ran*/VBD  
*on*/IN *the*/DT *street*/NN
  - [A/DT *red*/JJ *car*/NN] *ran*/VBD  
*on*/IN *the*/DT *street*/NN
- **Splitting** is the opposite of merging

# Concluding Remarks

- Chunking – as the weakest form of syntactic structuring – relies on RegExs
- RegExs (formally) belong to the class of regular grammars
- Regular grammars and their (finite-state) automata have linear run-time complexity
- Standard CF grammars and their associated push-down automata have (at best) cubic run-time complexity
- Hence, there is a trade-off between different levels of richness of syntactic structures and gains/losses of run-time behavior

# What are Named Entities?

- Names of persons
  - *Dr. Jonathan Peeko, Professor Johnson*
- Names of companies or organizations
  - *Sony, United Nations, Texas Instruments, General Motors*
- Names of locations
  - *Paris, San Francisco, Rocky Mountains, Yellowstone Park*
- Date and time expressions
  - *Feb 17, 1973; 4.40p.m.; 16.40 Uhr; autumn 2000; last year*
- Addresses
  - *7 Ugly Way, Wolverhampton W40 1Q5*
  - *udo.hahn@uni-jena.de*
- Names of proteins or genes or diseases,
  - *chloramphenicol acetyltransferase, NF-kappa B, SARS*
- Measure expressions
  - *420 kp, 21 l/m<sup>2</sup>, 37%, 900€*

# What are Named Entities?

- Names of persons
  - *Dr. Jonathan Peeko, Professor John*
- Names of companies
  - *Sony, Universal Motors*
- Names of locations
  - *Central Park*
- Dates
  - *1999*
- Addresses
  - *1010 105*
  - *udo.hahn@uni-jena.de*
- Names of proteins or genes or diseases,
  - *chloramphenicol acetyltransferase, NF-kappa B, SARS*
- Measure expressions
  - *420 kp, 21 l/m<sup>2</sup>, 37%, 900€*

**named entities are  
intentionally excluded from  
the lexicon**

# GATE: NER – Examples (1/3)

NYT19980403.0453 NEWS STORY 04/03/1998 21:01:00 CREDIT WARNING BY MOODY'S ON JAPANESE BONDS TOKYO \_ Borrowers in Japan, including even the healthiest corporations, faced a new challenge on Friday as Moody's Investors Service provided a pessimistic outlook on the nation's pristine credit rating. The exchange rate of Japan's currency, the yen, tumbled to a six-and-a-half-year low, and the stock and bond markets fell on the decision by the American-based ratings agency to change its view on Japan \_ whose government debt has been rated triple-A \_ from ``stable'' to ``negative.'' Moody's did not change any existing bond ratings, but the negative outlook may lead to a formal review in 18 months to two years. A lowered rating could raise borrowing costs for all Japanese, from consumers to large corporations, even those with impeccable credit. And such a move could further weaken Japanese banks, which already pay more to borrow because they hold in excess of \$600 billion in bad loans. The step by Moody's was a surprise because even with Japan's economic problems, it is still the world's largest creditor nation and there is little doubt about its ability to repay debts. But the announcement showed that Moody's \_ one of the world's big credit raters, along with Standard & Poor's and Duff & Phelps \_ was beginning to rethink Japan's long-term prospects. In trading here Friday the dollar surged to 135.42 yen, the highest since September 1991, before recovering a little. The benchmark Nikkei index of 225 stocks fell for the third consecutive day \_ to a four-month low of 15,517.78. Bond prices also declined, pushing the yield on the key 10-year Japanese government bond to 1.685 percent, a six-week high. Bond prices and yields move in opposite directions ``The world doesn't trust Japan anymore, even though Japan has lots of money.'' commented Xinyi Lu of Paribas

- ☒ Date
  - ☐ FirstPerson
  - ☐ Identifier
  - ☐ JobTitle
  - ☐ Location
  - ☐ Lookup
  - ☒ Money
  - ☐ Organization
  - ☐ Percent
  - ☐ Person
  - ☐ SpaceToken
  - ☐ Temp
  - ☐ Title
  - ☐ Token
- Original markups



# GATE: NER – Examples (2/3)

NYT19980403.0453 NEWS STORY 04/03/1998 21:01:00 CREDIT WARNING BY  
MOODY'S ON JAPANESE BONDS TOKYO \_ Borrowers in Japan, including even  
the healthiest corporations, faced a new challenge on Friday as  
Moody's Investors Service provided a pessimistic outlook on the  
nation's pristine credit rating. The exchange rate of Japan's  
currency, the yen, tumbled to a six-and-a-half-year low, and the stock  
and bond markets fell on the decision by the American-based ratings  
agency to change its view on Japan \_ whose government debt has been  
rated triple-A \_ from ``stable'' to ``negative.'' Moody's did not  
change any existing bond ratings, but the negative outlook may lead to  
a formal review in 18 months to two years. A lowered rating could  
raise borrowing costs for all Japanese, from consumers to large  
corporations, even those with impeccable credit. And such a move could  
further weaken Japanese banks, which already pay more to borrow  
because they hold in excess of \$600 billion in bad loans. The step by  
Moody's was a surprise because even with Japan's economic problems, it  
is still the world's largest creditor nation and there is little doubt  
about its ability to repay debts. But the announcement showed that  
Moody's \_ one of the world's big credit raters, along with Standard  
&AMP; Poor's and Duff &AMP; Phelps \_ was beginning to rethink Japan's  
long-term prospects. In trading here Friday the dollar surged to  
135.42 yen, the highest since September 1991, before recovering a  
little. The benchmark Nikkei index of 225 stocks fell for the third  
consecutive day \_ to a four-month low of 15,517.78. Bond prices also  
declined, pushing the yield on the key 10-year Japanese government  
bond to 1.685 percent, a six-week high. Bond prices and yields move in  
opposite directions ``The world doesn't trust Japan anymore, even  
though Japan has lots of money.'' commented Xinvi Lu of Paribas

- ☐ Date
  - ☐ FirstPerson
  - ☐ Identifier
  - ☐ JobTitle
  - ☒ Location
  - ☐ Lookup
  - ☐ Money
  - ☒ Organization
  - ☐ Percent
  - ☐ Person
  - ☐ SpaceToken
  - ☐ Temp
  - ☐ Title
  - ☐ Token
- Original markups

# GATE: NER – Examples (3/3)

NYT19980403.0456 NEWS STORY  
04/03/1998 21:02:00 BUOYANT CLINTON TAKES ON GOP SENATORS, BIG TOBACCO  
WASHINGTON \_ Eager to shift the spotlight from Paula Jones back to the  
business of government, President Clinton lambasted the Republican  
Senate budget proposal on Friday and warned tobacco companies to go  
along with a proposed settlement. Tired but buoyant in his first day  
back at the Oval Office after 12 days in Africa, Clinton immediately  
assembled his economic team in the White House Rose Garden this  
morning and signaled an election-year showdown with congressional  
Republicans over the budget for the 1999 fiscal year. While clearly  
emboldened by a federal judge's dismissal on Wednesday of Mrs. Jones'  
sexual misconduct lawsuit, the president vowed not to be distracted by  
such matters, saying, "I am going on with my business." Instead,  
Clinton castigated Senate Republicans for approving a \$1.73 trillion  
spending plan on Thursday night that calls for modest tax cuts and  
excludes virtually all of the president's proposals for new spending.  
And he scolded members of the House for passing a six-year, \$217  
billion transportation bill packed with projects for almost every  
congressional district. "I am very concerned that the budget plan now  
working its way through the Senate will squeeze out critical  
investments in education and children," Clinton said. "I'm also

- ☐ Date
- ☐ FirstPerson
- ☐ Identifier
- ☒ JobTitle
- ☐ Location
- ☐ Lookup
- ☐ Money
- ☐ Organization
- ☐ Percent
- ☒ Person
- ☐ SpaceToken
- ☐ Temp
- ☐ Title
- ☐ Token
- Original markups

# Two Types of NER Methods

## Human Knowledge Engineering (symbolic p.)

- rule based
- developed by experienced language engineers
- based on human intuition
- requires only small amount of plain training data
- development can be very time consuming
- some changes may be hard to accommodate

## (Supervised) Machine Learning Systems (empir.p.)

- use statistics or other machine learning technique
- developers do (almost) not need linguistic expertise
- fully automatic
- requires large amounts of annotated training data
- annotators are cheap (but you get what you pay for!)
- some changes may require re-annotation of the entire training corpus

# Naïve NER Method: List Look-up

- Recognize entities stored in given lists
  - *gazetteers*, e.g., online phone directories, yellow pages)
- Advantages:
  - simple, fast, language independent, easy to retarget (just create lists)
- Disadvantages:
  - impossible to enumerate all names and name variants, collection and maintenance of lists

# NER by Pattern Recognition

- Names often have internal structure - these components can be either stored or guessed, e.g., for "Location" we have RegEx-style constraints such as:

Capitalized Word + {City, Forest, Center, River}

which yields: *Sherwood Forest, Manchester City, Rhine River*

Capitalized Word + {Street, Boulevard, Avenue, Road}

which yields: *Portobello Street, Washington Avenue*

# NER by Expressive Rules

- Context-sensitive rules of the kind:

$$A \rightarrow B \setminus C / D$$

- A is a set of attribute-value expressions and optional score, the attributes refer to elements of the input token feature vector
- B, C, D are sequences of attribute-value pairs and regular expressions; variables are also supported
- B and D are left and right context, respectively, and can be empty (hint: read backwards!)

**Example:** `[syn=NP, sem=ORG] (0.9) →  
          \ [norm="university"], [token="of"],  
          [sem=REGION|COUNTRY|CITY] / ;`

# NER by Machine Learning

- NE task is frequently broken down in two parts:
  - Recognizing the entity boundaries
  - Classifying the entities in the NE categories
- Features are at least as important as the choice of the ML method
  - Simple pattern matching of orthographic features: capitalization, punctuation marks, numerical symbols
  - Windows for lexical features (e.g., “Mr.” for persons)
  - Affix features (“-ase” for proteins, “-ectomy” for medical procedures, etc.)
  - POS info (and chunks)

# Merkmale für die Zuordnung von Named Entities

Feature	Explanation
Lexical items	The token to be labeled
Stemmed lexical items	Stemmed version of the target token
Shape	The orthographic pattern of the target word
Character affixes	Character-level affixes of the target and surrounding words
Part of speech	Part of speech of the word
Syntactic chunk labels	Base-phrase chunk label
Gazetteer or name list	Presence of the word in one or more named entity lists
Predictive token(s)	Presence of predictive words in surrounding text
Bag of words/Bag of N-grams	Words and/or <i>N</i> -grams occurring in the surrounding context

Shape	Example
Lower	cummings
Capitalized	Washington
All caps	IRA
Mixed case	eBay
Capitalized character with period	H.
Ends in digit	A9
Contains hyphen	H-P



# Features for Machine Learning (CoNLL 2003 Shared Task)

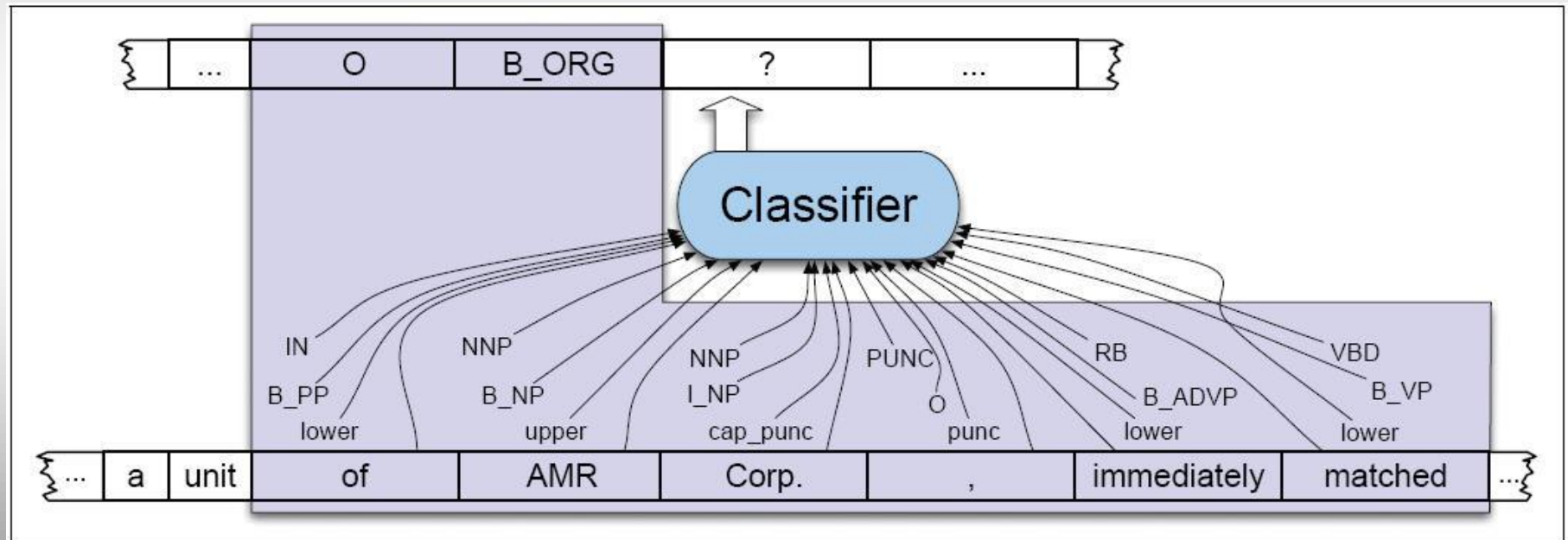
	lex	pos	aff	pre	ort	gaz	chu	pat	cas	tri	bag	quo	doc
Florian	+	+	+	+	+	+	+	-	+	-	-	-	-
Chieu	+	+	+	+	+	+	-	-	-	+	-	+	+
Klein	+	+	+	+	-	-	-	-	-	-	-	-	-
Zhang	+	+	+	+	+	+	+	-	-	+	-	-	-
Carreras (a)	+	+	+	+	+	+	+	+	-	+	+	-	-
Curran	+	+	+	+	+	+	-	+	+	-	-	-	-
Mayfield	+	+	+	+	+	-	+	+	-	-	-	+	-
Carreras (b)	+	+	+	+	+	-	-	+	-	-	-	-	-
McCallum	+	-	-	-	+	+	-	+	-	-	-	-	-
Bender	+	+	-	+	+	+	+	-	-	-	-	-	-
Munro	+	+	+	-	-	-	+	-	+	+	+	-	-
Wu	+	+	+	+	+	+	-	-	-	-	-	-	-
Whitelaw	-	-	+	+	-	-	-	-	+	-	-	-	-
Hendrickx	+	+	+	+	+	+	+	-	-	-	-	-	-
De Meulder	+	+	+	-	+	+	+	-	+	-	-	-	-
Hammerton	+	+	-	-	-	+	+	-	-	-	-	-	-

Table 3: Main features used by the sixteen systems that participated in the CoNLL-2003 shared task sorted by performance on the English test data. Aff: affix information (n-grams); bag: bag of words; cas: global case information; chu: chunk tags; doc: global document information; gaz: gazetteers; lex: lexical features; ort: orthographic information; pat: orthographic patterns (like Aa0); pos: part-of-speech tags; pre: previously predicted NE tags; quo: flag signing that the word is between quotes; tri: trigger words.

# Merkmalskodierung für NEs

Features				Label
American	NNP	B <sub>NP</sub>	cap	B <sub>ORG</sub>
Airlines	NNPS	I <sub>NP</sub>	cap	I <sub>ORG</sub>
,	PUNC	O	punc	O
a	DT	B <sub>NP</sub>	lower	O
unit	NN	I <sub>NP</sub>	lower	O
of	IN	B <sub>PP</sub>	lower	O
AMR	NNP	B <sub>NP</sub>	upper	B <sub>ORG</sub>
Corp.	NNP	I <sub>NP</sub>	cap_punc	I <sub>ORG</sub>
,	PUNC	O	punc	O
immediately	RB	B <sub>ADVP</sub>	lower	O
matched	VBD	B <sub>VP</sub>	lower	O
the	DT	B <sub>NP</sub>	lower	O
move	NN	I <sub>NP</sub>	lower	O
,	PUNC	O	punc	O
spokesman	NN	B <sub>NP</sub>	lower	O
Tim	NNP	I <sub>NP</sub>	cap	B <sub>PER</sub>
Wagner	NNP	I <sub>NP</sub>	cap	I <sub>PER</sub>
said	VBD	B <sub>VP</sub>	lower	O
.	PUNC	O	punc	O

# Named Entity Tagging als Sequence Labeling-Problem



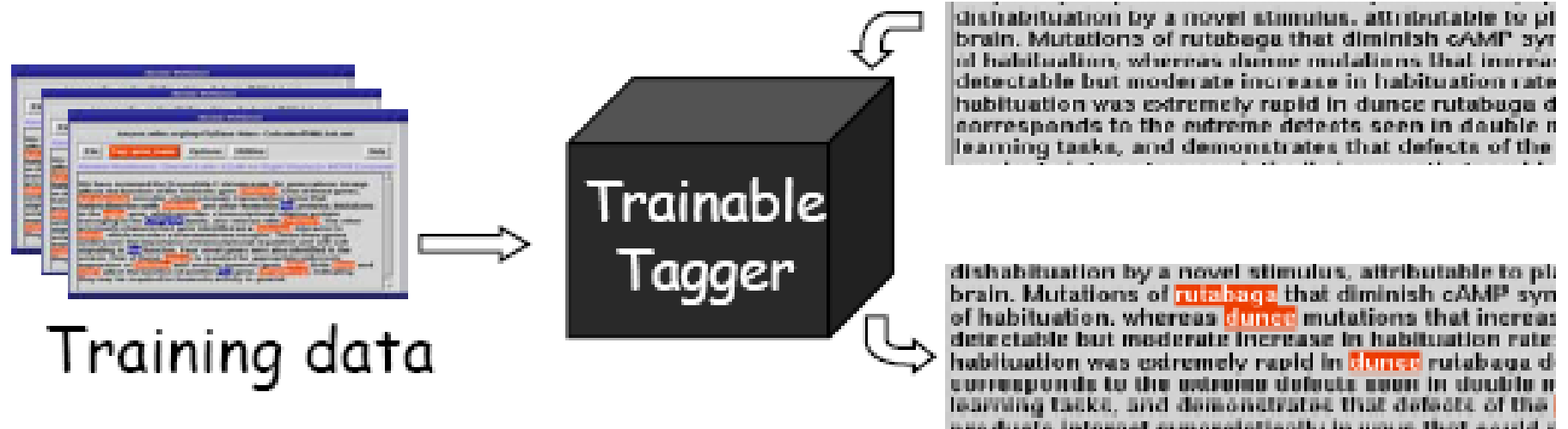
# Systemarchitektur für (überwachtes) Maschinelles Lernen

## Merkmale

= beobachtbare Indikatoren  
(in den Trainingsdaten)

## Algorithmen für Maschinelles Lernen

= Rechenverfahren zur Bestimmung von  
(statistischen) Modellen über die Verteilung von  
Merkmalen (in den Trainingsdaten)



# Algorithmen für (überwachtes) Maschinelles Lernen [Flach 2012, Murphy 2012]

- Einfache Klassifikatoren (Classifier)
  - Naive-Bayes'scher Klassifikator
  - k-Nächster Nachbar (k-nearest neighbor)
  - Entscheidungsbäume (decision trees)
- Hochdimensionale Klassifikatoren (Classifier)
  - Support Vector Machines (SVM)
- (strukturorientierte) Graphische Modelle
  - Hidden-Markov-Modelle
  - Conditional Random Fields (CRF)
  - Bayes'sche Netze
- (Künstliche) neuronale Netze  $\Rightarrow$  Deep Learning
- Genetische Algorithmen

# Machine Learning–General Task

A computer program is said to *learn*

- from experience  $E$  (data in the form of representative examples / instances of the whole input space)
  - with respect to some class of tasks  $T$
  - and performance measure  $P$ ,
  - if its performance at tasks  $T$  as measured by  $P$ , improves with experience  $E$
- 
- Learned hypothesis: model of problem/task  $T$
  - Model quality: accuracy/performance measured by  $P$

# Machine Learning – Two Fundamental Modes

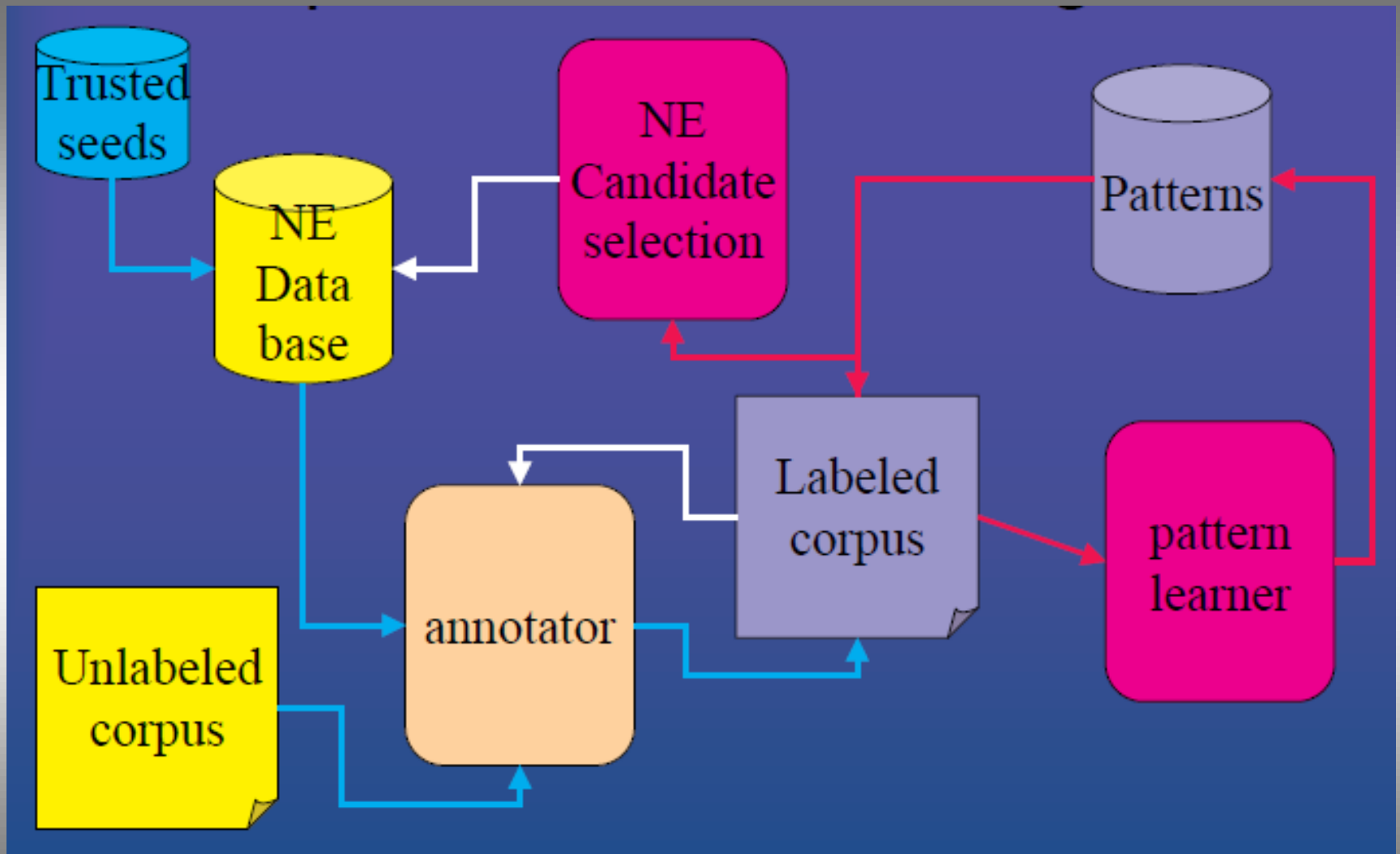
- Supervised learning
  - Given : Training examples (training set T)  
 $\{ (x_1, f(x_1)), (x_2, f(x_2)), \dots (x_n, f(x_n)) \}$   
for some unknown function  $y = f(x)$
  - Find :  $f(x)$
  - Predict  $y' = f(x')$  where  $x'$  is not in the training set but T-wise similar data sets
- Unsupervised learning
  - Given : data (data set D)  
 $\{ x_1, x_2, \dots, x_n \}$   
for some unknown function  $y = f(x)$
  - Find :  $f(x)$
  - Predict  $y = f(x)$  where  $x$  is in the data set or D-wise similar data sets

# Basic Idea for (Almost) Unsupervised NER

- Define manually only a small set of trusted seeds (a bit of ground truth)
- Training then only uses unlabeled data
- Initialize system by labeling the corpus with the seeds
- Extract and generalize patterns from the context of the seeds
- Use the patterns to further label the corpus and to extend the seed set (*bootstrapping*)
- Repeat the process unless no new terms can be identified



# Architecture for (Almost) Unsupervised NER



# Learning Ordered Decision Rules

- The task: to learn a decision list to classify strings as **person**, **location** or **organization**

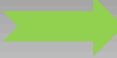
The learned decision list is an *ordered* sequence of if-then rules

... says *Mr. Gates*, founder of *Microsoft* ...

... says *Mr. Gates*, founder of *Microsoft* ...

$R_1$  : if features then **person**  
 $R_2$  : if features then **location**  
 $R_3$  : if features then organization  
...  
 $R_n$  : if features then **person**

# Outline of Unsupervised Co-Training

- Parse an unlabeled document set  syntactic units
- Extract each NP whose head is tagged as Proper Noun (Proper Noun is supertype of NER as subtyping)
- Define a set of relevant features which can be applied to extracted NPs
- Define two separate types of rules on the basis of the feature space
- Determine small initial set of seed rules
- Iteratively extend the rules through co-training

# Two Types of Rules

- **Spelling Rules**
  - Rules which directly specify lexical conditions (e.g., “Mr.”  
⇒PERSON)
- **Contextual Rules**
  - Rules which specify co-occurring lexical or phrasal conditions (e.g., “president” co-occurs with “Mr.”  
⇒PERSON)
- **N.B.: Huge amount of unlabeled data in a corpus gives useful hints!**

# Kinds of Noun Phrases and Spelling-Context Pairs

1. There was an appositive modifier to the NP, whose head is a singular noun (tagged NN).

- ...says [Maury Cooper], [a vice president]...

2. The NP is a complement to a preposition which is the head of a PP. This PP modifies another NP whose head is a singular noun.

- ... fraud related to work on [a federally funded sewage plant] [in [Georgia]].

- ...says Maury Cooper, a vice president...
  - (Maury Cooper, president)

- ... fraud related to work on a federally funded sewage plant in Georgia.
  - (Georgia, plant\_in)

# Features

- Set of spelling features
  - Full-string=x (full-string=Maury Cooper)
  - Contains(x) (contains(Maury))
  - Allcap1 IBM
  - Allcap2 N.Y.
  - Nonalpha=x A.T.&T. (nonalpha=..&.)
- Set of context features
  - Context = x (context = president)
  - Context-type = x appos or prep

# Examples of Features

<u>Sentence</u>	<u>Entities(Spelling/Context)</u>	<u>(Active) Features</u>
But Robert Jordan, a partner at Steptoe & Johnson who took ...	Robert Jordon/partner	Full-string=Robert_Jordan, contains(Robert), contains(Jordan), context=partner, context-type=appos
	Steptoe & Johnson/partner_at	Full-string=Steptoe_&_Johnson, contains(Steptoe), contains(&), contains(Johnson), nonalpha=& , context=partner_at, context-type=prep
By hiring a company like A.T.&T. ...	A.T.&T./company_like	Full-string= A.T.&T., allcap2, nonalpha=..&.. , context=company_like, context-type=prep
Hanson acquired Kidde Incorporated, parent of Kidde Credit, for ...	Kidde Incorporated/parent	Full-string=Kidde_Incorporated, contains(Kidde), contains(Incorporated), context=parent, context-type=appos
	Kidde Credit/parent_of	Full-string=Kidde_Credit, contains(Kidde), contains(Credit), context=parent_of, context-type=prep



# Formal Structure of Rules

## Rules

Two separate types of rules:  
Spelling rules  
Context rules

Feature  $\rightarrow$  NE-type,  $h(\text{Feature}, \text{NE-type})$

$h(x, y)$ : the strength of a rule, defined as

$Count(x, y)$  is the number of times feature  $x$  is seen with label  $y$  in training data,

$$\max_{x, y} \frac{Count(x, y) + \alpha}{Count(x) + k\alpha} \quad \text{where}$$

$$Count(x) = \sum_{y \in Y} Count(x, y)$$

$\alpha$  is a smoothing parameter

$$k = \#NE\text{-types}$$

Is an estimate of the conditional probability of the NE-type given the feature,  $P(y|x)$

The rules ordered according to their strengths  $h$  form a decision list: the sequence of rules are tested in order, and the answer to the **first** satisfied rule is output.



# 7 Seed Rules

## 7 SEED RULES

- Full-string = New York → Location
- Full-string = California → Location
- Full-string = U.S. → Location
- Contains(Mr.) → Person
- Contains(Incorporated) → Organization
- Full-string=Microsoft → Organization
- Full-string=I.B.M. → Organization

Note: only one type of rules used as seed rules, and all NE-types should be covered

# Co-Training Algorithm

1. Set  $N=5$  (max. # of rules of each type induced in each iteration)
2. **Initialize:** Set the **spelling** decision list equal to the set of seed rules. Label the training set using these rules.
3. Use **these** to get contextual rules. ( $x$  = feature,  $y$  = label)
  1. Compute  $h(x,y)$ , and induce at most  $N * K$  rules  $K = \# \text{ NE types}$
  2. all must be above some threshold  $p_{\min}=0.95$
4. Label the training set using the contextual rules.
5. Use these to get  $N*K$  **spelling** rules (same as step 3.)
6. Set **spelling** rules to seed plus the new rules.
7. If  $N < 2500$ , set  $N=N+5$ , and goto step 3.
8. Label the training data with the combined spelling/contextual decision list, then induce a final decision list from the labeled examples where all rules (regardless of strength) are added to the decision list.

# Example

- (IBM, company)
  - ...IBM, the company that makes...
- (General Electric, company)
  - ..General Electric, a leading company in the area,...
- (General Electric, employer )
  - ... joined General Electric, the biggest employer...
- (NYU, employer)
  - NYU, the employer of the famous Ralph Grishman,...

# Power of the Algorithm

- Greedy method
  - At each iteration method increases number of rules
  - While maintaining a high level of agreement between spelling & context rules

For  $n = 2500$ :

1. The two classifiers give both labels on 49.2% of the unlabeled data
  2. And give the *same* label on 99.25% of these cases
- The algorithm maximizes the number of unlabeled examples on which the two decision list agree.

# Evaluation of the Algorithm

- 88,962 (spelling, context) pairs.
  - 971,746 sentences
- 1,000 randomly extracted to be test set.
- Location, person, organization, noise (items outside the other three)
- 186, 289, 402, 123 (- 38 temporal noise).
- Let  $N_c$  be the number of correctly classified examples
  - Noise Accuracy:  $N_c / 962$

# Results

<u>Algorithm</u>	<u>Clean Accuracy</u>	
Baseline	45.8%	
EM	83.1%	
Yarowsky 95	81.3%	
Yarowsky Cautious	91.2%	
DL-CoTrain	91.3%	
CoBoost	91.1%	

# Remarks

- Needs full parsing of unlabeled documents
  - Restricted language independency
  - Need linguistic sophistication for new types of NE
- Slow training
  - In each iteration, full size of training corpus has to be re-labeled

# Resources for NLP

- Empirical (Learning) Paradigm for NLP
- **Types of Resources**
  - Language data (plain, annotated)
  - Systems for acquiring and maintaining language data
  - Computational lexicons and ontologies
  - NLP Core Engines
  - NLP Application Systems
  - Machine Learning Resources
- Methodological Issues of NLP Resources



# Ressourcen für die Sprachverarbeitung

- Referenzkorpora (Nationalkorpora)
  - Standardsprache (Zeitungen, Belletristik)
- Non-Standard-Korpora
  - Informelle Sprache (Chats, Blogs, E-Mails)
  - Fachsprachen (z.B.: klinische Berichte)
- Rohdaten vs. Annotation
  - Linguistische Metadaten
    - Morphologie, Syntax, Semantik, Pragmatik

# Language Data

- **Plain language data**
  - Just text or speech
    - ASCII/UTF-8-compatible, pdf, HTML/SGML
- **Annotated language data**
  - Enriched by linguistic meta-data
    - Linguistic annotation languages (XML)

# Plain Language Data

- **Mixed/Balanced text collections**
  - British National Corpus (BNC)
  - American National Corpus (ANC)
- **Newspaper collections**
  - Wall Street Journal
  - IdS-Korpora (DeReKo\*)
- **The Web**

# British National Corpus (BNC)

- 100M word collection (some 4,050 texts) of 20th century British English
- Written part (90%)
  - Regional and national newspapers
  - Specialist periodicals and journals (various genres)
  - Academic books and popular fiction
  - Letters, memoranda, school and university essays
- Spoken part (10%)
  - Informal conversations (different ages, regions, social classes)
  - Formal business and government meetings
  - Radio shows and phone-ins
- <http://www.natcorp.ox.ac.uk/>

# British National Corpus (BNC)

- Encoding based on ‘Guidelines of the Text Encoding Initiative’ (TEI),
  - using ISO standard 8879 (SGML: Standard Generalized Markup Language)
- Whole collection is POS-tagged
  - using the CLAWS tagger for the C5 tag set (C7 is much more elaborate)
  - Error rate: 1.7%
  - Tagging ambiguity for 4.7% of all tags

# American National Corpus (ANC)

- 15M word collection ( texts) of 20th century American English
- Annotated for structural markup (sections, chapters, etc.) down to the level of paragraph, sentence boundaries, words (tokens) with part of speech annotations and lemma using the Penn tagset, noun and verb chunks, named entities (Person, Location, Organization, Date)
- Written part (80%)
  - Regional and national newspapers
  - Specialist periodicals and journals (various genres)
  - Academic books and popular fiction
  - Governmental docs
- Spoken part (20%)
  - Informal face-to-face conversations (different ages, regions, social classes)
  - Telephone conversations
- <http://www.anc.org/>

# Große deutsche Textkorpora

## (verschriftlichte Sprache)

- **Deutsches Referenzkorpus – DeReKo**
  - Institut für deutsche Sprache (IdS) Mannheim
  - Zeitungen, Belletristik, Handbücher, Parlamentsprotokolle (seit 1956)
  - Umfang: ca. 42 Mrd. Tokens
  - <http://www1.ids-mannheim.de/kl/projekte/korpora/>
- **Digitales Wörterbuch der deutschen Sprache – DWDS**
  - Berlin-Brandenburgische Akademie der Wissenschaften (BBAW)
  - Zeitungen, Belletristik, Gebrauchsliteratur, Wissenschaft (20./21. Jahrhundert)
  - Umfang: 12,8 Mio Dokumente, ca. 5,5 Mrd. Tokens
  - <https://www.dwds.de/>
- **Deutsches Textarchiv – DTA**
  - Historisches Referenzkorpus: 1600-1900
  - 4422 Werke (ca. 900K Seiten), 297 Mio. Tokens
  - Annotiert mit Tokens, Lemmata, POS
  - <http://www.deutschestextarchiv.de/>

# Language Data Repositories

- Linguistic Data Consortium
  - „Catalog“ option
  - „LDC Online“ provides you a guest account

<http://www.ldc.upenn.edu/>



# Language Data Repositories

- Linguist List
  - Open Language Archives Community
  - „Text & Computer Tools“ button
    - Texts and Corpora
  - „Language Resources“ button
    - Texts and Corpora

<http://linguistlist.org/olac>

# Language Data Repositories

- European Language Resources Association (ELRA)
  - „R&D Catalog“ option
  - Spoken LRs
    - Telephone recordings
    - Desktop/mircophone recordings
    - Broadcast resources
    - Speech related resources
  - Written LRs
    - Corpora
    - Mono- and multilingual lexicons
  - (Domain-specific) Terminological resources
  - Multimodal/multimedia LRs

<http://www.elra.info/>

# Language Data Repositories

- Natural Language Software Registry
  - Annotation tools
  - Evaluation tools
  - Language Resources
  - Multimedia
  - Multimodality
  - NLP Development Aid
  - Spoken Language
  - Written Language

<http://registry.dfki.de>

# Annotated Language Data

- **Levels of annotation**
  - **Formal text structure processing**
    - Paragraphs, sentences, tokens
  - **Syntactic mark-up**
    - Parts of speech
    - Shallow syntactic structures: chunks
    - Deep syntactic structures: parses
  - **Semantic mark-up**
    - Named entities
    - Propositions, predicate-argument structures
  - **Discourse mark-up**
    - Referential relations
    - Rhetorical relations

# Annotation Styles

- In-line annotation
  - Mark-ups appear as integral part of the original text
    - This is an `<XMLTag>` `in-line` `<\XMLTag>` annotation
- Stand-off annotation
  - Mark-ups appear distinct from the original text (e.g., in a different window)
    - This is a `stand-off` annotation
      - `<XMLTag StartChar: 11, XMLTag EndChar: 19, XMLTag Type STAND-OFF>`

# General Language Corpora for Syntactic Annotation

- Penn Treebank (U Penn)

- language: English (general language)
- text genre: mostly newspaper articles (*Wall Street Journal*)
- size: 1,200,000 (annotated) tokens
- Syntactic tagging based on set of 45 tags
- Syntactic phrase structures (parse trees) based on Government-Binding grammar
- No named entity annotation
- But propositional annotation: PropBank

<http://www.cis.upenn.edu/~treebank/>

# General Language Corpora for Proposition Annotation

- PropBank (U Penn)

- language: English (general language)
  - text genre: financial newspaper articles (*Wall Street Journal*)
  - size: 300,000 (annotated) tokens
  - proposition format:
    - [ subject - predicate - object ]
  - “semantic” counterpart of Penn Treebank
- <http://www.cis.upenn.edu/~ace/>

# General Language Corpora for Discourse Annotation

- **Penn Discourse TreeBank (PDTB; U Penn)**

- language: English (general language)
- text genre: financial newspaper articles (*Wall Street Journal*)
- size: 1 M tokens (WSJ) and 40k relations
- Annotated with information related to discourse structure and discourse semantics, i.e., temporal, contingency, comparison, and expansion discourse relations (after, when, but, although, if)
- “discourse” counterpart of **Penn Treebank**

<http://www.cis.upenn.edu/~pdtb/>



# General Language Corpora for Discourse Analysis

- **RST Corpus (ISI/USC, USA)**
  - language: English
  - size: 385 documents, i.e., 176,000 tokens;  
21,789 elementary discourse units (EDUs)
  - text genre: newspaper articles (*Wall Street Journal*)
  - Rhetorical Structure Theory (RST)
    - 90 coherence relations

# Penn TreeBank: Sizes and Genres

**Table 4:**  
**Penn Treebank**  
(as of 11/92)

<b>Description</b>	<b>Tagged for Part-of-Speech (Tokens)</b>	<b>Shallow Parsing (Tokens)</b>
Dept. of Energy abstracts	231,404	231,404
Dow Jones Newswire stories	3,065,776	1,061,166
Dept. of Agriculture bulletins	78,555	78,555
Library of America texts	105,652	105,652
MUC-3 messages	111,828	111,828
IBM Manual sentences	89,121	89,121
WBUR radio transcripts	11,589	11,589
ATIS sentences	19,832	19,832
Brown Corpus, retagged	1,172,041	1,172,041
<b>Total:</b>	<b>4,885,798</b>	<b>2,881,188</b>

# Penn TreeBank POS Tag Set

**Table 2:**  
The Penn Treebank POS tagset

1.	CC	Coordinating conjunction	25.	TO	<i>to</i>
2.	CD	Cardinal number	26.	UH	Interjection
3.	DT	Determiner	27.	VB	Verb, base form
4.	EX	Existential <i>there</i>	28.	VBD	Verb, past tense
5.	FW	Foreign word	29.	VBG	Verb, gerund/present participle
6.	IN	Preposition/subord. conjunction	30.	VBN	Verb, past participle
7.	JJ	Adjective	31.	VBP	Verb, non-3rd ps. sing. present
8.	JJR	Adjective, comparative	32.	VBZ	Verb, 3rd ps. sing. present
9.	JJS	Adjective, superlative	33.	WDT	<i>wh</i> -determiner
10.	LS	List item marker	34.	WP	<i>wh</i> -pronoun
11.	MD	Modal	35.	WP\$	Possessive <i>wh</i> -pronoun
12.	NN	Noun, singular or mass	36.	WRB	<i>wh</i> -adverb
13.	NNS	Noun, plural	37.	#	Pound sign
14.	NNP	Proper noun, singular	38.	\$	Dollar sign
15.	NNPS	Proper noun, plural	39.	.	Sentence-final punctuation
16.	PDT	Predeterminer	40.	,	Comma
17.	POS	Possessive ending	41.	:	Colon, semi-colon
18.	PRP	Personal pronoun	42.	(	Left bracket character
19.	PP\$	Possessive pronoun	43.	)	Right bracket character
20.	RB	Adverb	44.	"	Straight double quote
21.	RBR	Adverb, comparative	45.	'	Left open single quote
22.	RBS	Adverb, superlative	46.	"	Left open double quote
23.	RP	Particle	47.	'	Right close single quote
24.	SYM	Symbol (mathematical or scientific)	48.	"	Right close double quote

# PTB POS Annotation Process

- Four annotators: Grad students of linguistics
- Comparison of two annotation styles on a 16,000 word sample:
  - „Tagging“:
    - completely manual annotation
  - „Correcting“:
    - automatical POS tagging and subsequent manual correction
- Inter-annotator disagreement:
  - „Tagging“: 7,2%
  - „Correcting“: 4,1%
- Comparison of accuracy with benchmark version (disagreement):
  - „Tagging“: 5,4%
  - „Correcting“: 4,0%

# Illustration of the „Correcting“ Mode

Battle-tested/NNP\*/JJ Japanese/NNP\*/JJ industrial/JJ managers/NNS here/RB  
always/RB buck/VB\*/VBP up/IN\*/RP nervous/JJ newcomers/NNS with/IN the/DT  
tale/NN of/IN the/DT first/JJ of/IN their/PP\$ countrymen/NNS to/TO visit/VB  
Mexico/NNP ,/, a/DT boatload/NN of/IN samurai/NNS\*/FW warriors/NNS blown/VBN  
ashore/RB 375/CD years/NNS ago/RB ./.

- Training of annotators took 15h
- Annotation speed (after one month of training):  
> 3000 words/h
- Double as fast as „Tagging“ !

# Syntactic Annotation of PTB

- Correction of false automatic parser output as provided by the FIDDITCH parser (Hindle 1989):
  - Outputs only one analysis per sentence
  - No attachments when parser is unsure about attachment decision
  - Alternative solution: decomposition of sentence structure into sets of partial trees
    - partial sentence structure description
  - Good lexicon and grammar coverage
- Task of annotators is mainly to „glue“ (i.e., to attach) partial phrase structure trees
  - Less time-consuming than re-bracketing the entire parser output

# Penn Treebank Phrasal Tag Set

- |     |        |  |
|-----|--------|--|
| 1.  | ADJP   | Adjective phrase   |
| 2.  | ADVP   | Adverb phrase  |
| 3.  | NP     | Noun phrase  |
| 4.  | PP     | Prepositional phrase   |
| 5.  | S      | Simple declarative clause  |
| 6.  | SBAR   | Clause introduced by subordinating conjunction or <i>0</i> (see below) |
| 7.  | SBARQ  | Direct question introduced by <i>wh</i> -word or <i>wh</i> -phrase     |
| 8.  | SINV   | Declarative sentence with subject-aux inversion                        |
| 9.  | SQ     | Subconstituent of SBARQ excluding <i>wh</i> -word or <i>wh</i> -phrase |
| 10. | VP     | Verb phrase  |
| 11. | WHADVP | <i>Wh</i> -adverb phrase   |
| 12. | WHNP   | <i>Wh</i> -noun phrase   |
| 13. | WHPP   | <i>Wh</i> -prepositional phrase  |
| 14. | X      | Constituent of unknown or uncertain category                           |

## Null elements

- |    |     |   |
|----|-----|---|
| 1. | *   | “Understood” subject of infinitive or imperative                        |
| 2. | 0   | Zero variant of <i>that</i> in subordinate clauses                      |
| 3. | T   | Trace—marks position where moved <i>wh</i> -constituent is interpreted  |
| 4. | NIL | Marks position where preposition is interpreted in pied-piping contexts |

# Partially bracketed output from FIDDITCH

```
( (S
  (NP (NBAR (ADJP (ADJ "Battle-tested/JJ")
                  (ADJ "industrial/JJ")))
      (NPL "managers/NNS"))))
  (? (ADV "here/RB"))
  (? (ADV "always/RB"))
  (AUX (TNS *))
  (VP (VPRES "buck/VBP"))
  (? (PP (PREP "up/RP")
        (NP (NBAR (ADJ "nervous/JJ")
                  (NPL "newcomers/NNS")))))
  (? (PP (PREP "with/IN")
        (NP (DART "the/DT")
            (NBAR (N "tale/NN")
                  (PP of/PREP
                    (NP (DART "the/DT")
                        (NBAR (ADJP
                              (ADJ "first/JJ")))))))))
  (? (PP of/PREP
      (NP (PROS "their/PP$")
          (NBAR (NPL "countrymen/NNS"))))
  (? (S (NP (PRO *))
        (AUX to/TNS)
        (VP (V "visit/VB")
            (NP (PNP "Mexico/NNP"))))
  (? (MID ",/,")
  (? (NP (IART "a/DT")
        (NBAR (N "boatload/NN")
              (PP of/PREP
                (NP (NBAR
                    (NPL "warriors/NNS"))))
              (VP (VPPRT "blown/VBN")
                  (? (ADV "ashore/RB"))
                  (NP (NBAR (CARD "375/CD")
                          (NPL "years/NNS")))))
  (? (ADV "ago/RB"))
  (? (FIN "./.")))
```



# Automatic simplification of the output from FIDDITCH

```
( (S
  (NP (ADJP Battle-tested industrial)
      managers)
  (? here)
  (? always)
  (VP buck))
  (? (PP up
      (NP nervous newcomers)))
  (? (PP with
      (NP the tale
        (PP of
          (NP the
            (ADJP first))))))
  (? (PP of
      (NP their countrymen)))
  (? (S (NP *)
      to
      (VP visit
        (NP Mexico))))
  (? ,)
  (? (NP a boatload
      (PP of
        (NP warriors))
      (VP blown
        (? ashore)
        (NP 375 years))))
  (? ago)
  (? .))
```

# After „Correcting“ by the annotators

```
( (S
  (NP Battle-tested industrial managers
    here)
  always
  (VP buck
    up
    (NP nervous newcomers)
    (PP with
      (NP the tale
        (PP of
          (NP (NP the
              (ADJP first
                (PP of
                  (NP their countrymen)))
              (S (NP *)
                to
                (VP visit
                  (NP Mexico))))
            ,
            (NP (NP a boatload
                (PP of
                  (NP (NP warriors)
                    (VP-1 blown
                      ashore
                      (ADVP (NP 375 years)
                        ago))))))
            (VP-1 *pseudo-attach*)))))))))
  .)
```

# TiGer Corpus

- 0,9M word collection (50K sentences) of German language newspaper articles (FR)
- <http://www.ims.uni-stuttgart.de/projekte/TIGER/TIGERCorpus/>
- morphological, POS, parse tree tagging
- Treebank query tool TiGer Search

# TiGer Corpus

**TIGERGraphViewer**

File Graph View Options Help

Icons: [Tree] [Graph] [List] [T] [Folder]

Minister  
NN  
Masc.Nom.Sg  
Minister

heizt  
VVFIN  
3.Sg.Pres.Ind  
heizen

Debatte  
NN  
Fem.Akk.Sg  
Debatte

über  
APPR  
Akk  
über

Sterbehilfe  
NN  
Fem.Akk.Sg  
Sterbehilfe

an  
PTKVZ  
an

Graphs: 200  
Subgraphs: --

Previous 2 Next  
First 1 200 Last

Subgraph: --

**s26743:** Minister heizt Debatte über Sterbehilfe an

Displaying the corpus (200 corpus graphs).

# TiGer Search (NP)

**TIGERGraphViewer**

File Graph View Options Help

Icons: [Tree] [Graph] [View] [Options] [Help]

Graph 1: NP (red) branches to NK (Der, ART, Masc.Nom.Sg, der), NK (sozialdemokratische, ADJA, Pos.Masc.Nom.Sg, sozialdemokratisch), NK (Minister, NN, Masc.Nom.Sg, Minister).

Graph 2: NP (black) branches to NK (dem, ART, Neut.Dat.Sg, das), NK (Partei-Blatt, NN, Neut.Dat.Sg, Partei-Blatt), NK (Doen, NE, \*.Nom.Sg, Doen).

Graph 3: sagte (VVFIN, 3.Sg.Past.Ind, sagen).

Graph 4: NP (black) branches to NK (dem, ART, Neut.Dat.Sg, das), NK (Partei-Blatt, NN, Neut.Dat.Sg, Partei-Blatt), NK (Doen, NE, \*.Nom.Sg, Doen).

Graph 5: , (Comma), 1 (Number), \$ (Dollar sign), C (Currency).

Graph 6: 1 (Number), 1 (Number).

Graphs: 24  
Subgraphs: 24

Navigation: Previous, 1, Next, First, 1, 24, Last

Subgraph: 1 / 1

s26746: Der sozialdemokratische Minister sagte dem Partei-Blatt Doen , 1989 einen Arzt gefragt zu haben , seine unheilbar kranke , im Sterben liegende Mutter von ihrem Leiden zu erlösen .

Displaying matches (24 matching corpus graphs, 24 matching subgraphs).

# STTS Tag Set for German (1/2)

- ADJA attributives Adjektiv [das] große [Haus]
- ADJD adverbiales oder [er fährt] schnell prädikatives Adjektiv [er ist] schnell
- ADV Adverb schon, bald, doch
- APPR Präposition; Zirkumposition links in [der Stadt], ohne [mich]
- APPRART Präposition mit Artikel im [Haus], zur [Sache]
- APPO Postposition [ihm] zufolge, [der Sache] wegen
- APZR Zirkumposition rechts [von jetzt] an
- ART bestimmter oder der, die, das, unbestimmter Artikel ein, eine, ...
- CARD Kardinalzahl zwei [Männer], [im Jahre] 1994 (Ordinalzahlen sind als ADJA getaggt)
- FM Fremdsprachliches Material [Er hat das mit ``] A big fish [`` übersetzt]
- ITJ Interjektion mhm, ach, tja
- KOUJ unterordnende Konjunktion um [zu leben], mit ``zu" und Infinitiv anstatt [zu fragen]
- KOUS unterordnende Konjunktion weil, daß, damit, mit Satz wenn, ob
- KON nebenordnende Konjunktion und, oder, aber
- KOKOM Vergleichskonjunktion als, wie
- NN normales Nomen Tisch, Herr, [das] Reisen
- NE Eigennamen Hans, Hamburg, HSV
- PDS substituierendes Demonstrativ- dieser, jener pronomen
- PDAT attribuierendes Demonstrativ- jener [Mensch] pronomen
- PIS substituierendes Indefinit- keiner, viele, man, niemand pronomen
- PIAT attribuierendes Indefinit- kein [Mensch], pronomen ohne Determiner irgendein [Glas]
- PIDAT attribuierendes Indefinit- [ein] wenig [Wasser], pronomen mit Determiner [die] beiden [Brüder]
- PPER irreflexives Personalpronomen ich, er, ihm, mich, dir
- PPOSS substituierendes Possessiv- meins, deiner pronomen
- PPOSAT attribuierendes Possessivpronomen mein [Buch], deine [Mutter]

# STTS Tag Set for German (2/2)

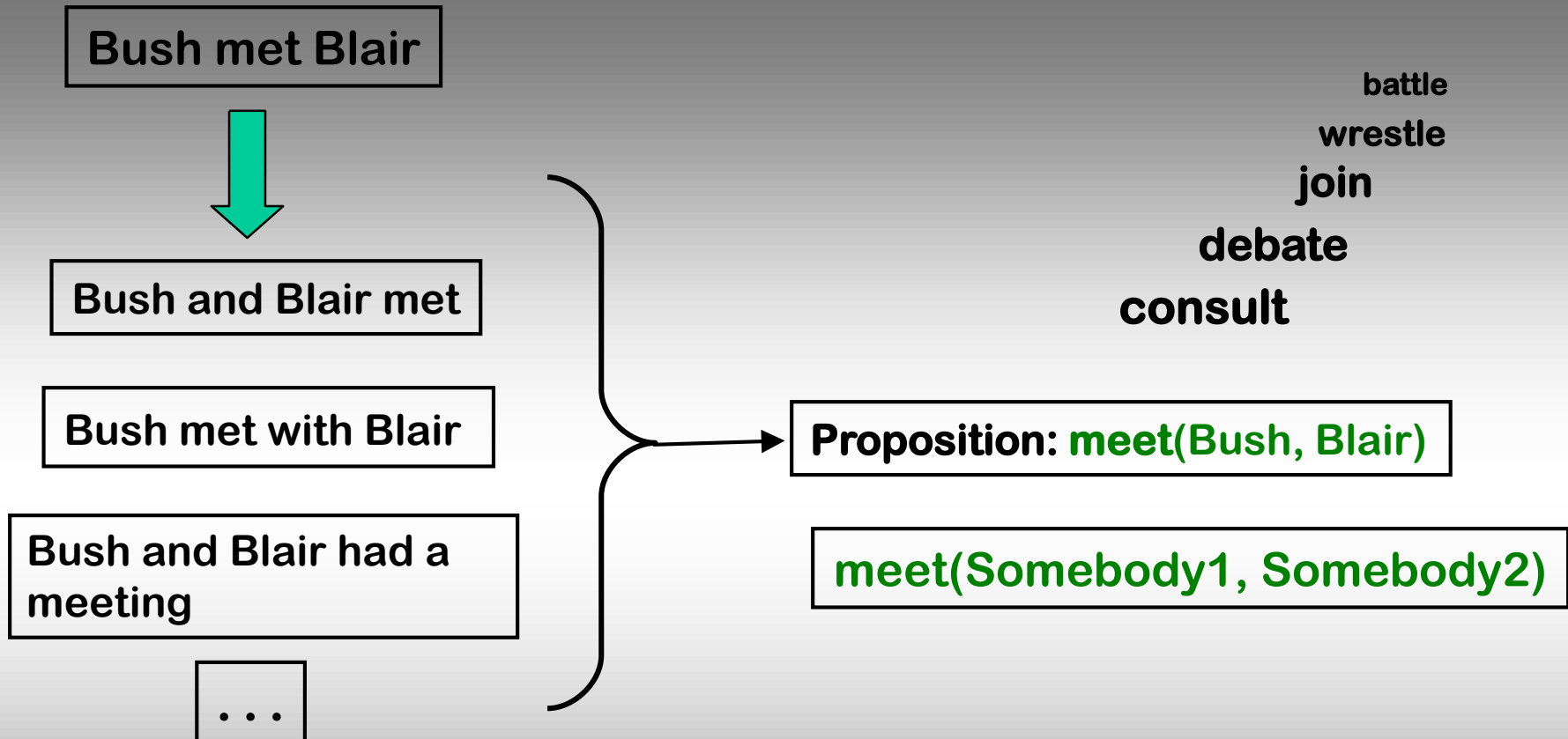
- PRELS substituierendes Relativpronomen [der Hund ,] der
- PRELAT attribuierendes Relativpronomen [der Mann ,] dessen [Hund]
- PRF reflexives Personalpronomen sich, einander, dich, mir
- PWS substituierendes wer, was Interrogativpronomen
- PWAT attribuierendes welche [Farbe], Interrogativpronomen wessen [Hut]
- PWAV adverbiales Interrogativ- warum, wo, wann, oder Relativpronomen worüber, wobei
- PAV Pronominaladverb dafür, dabei, deswegen, trotzdem
- PTKZU ``zu" vor Infinitiv zu [gehen]
- PTKNEG Negationspartikel nicht
- PTKVZ abgetrennter Verbzusatz [er kommt] an, [er fährt] rad
- PTKANT Antwortpartikel ja, nein, danke, bitte
- PTKA Partikel bei Adjektiv am [schönsten], oder Adverb zu [schnell]
- TRUNC Kompositions-Erstglied An- [und Abreise]
- VVFIN finites Verb, voll [du] gehst, [wir] kommen [an]
- VVIMP Imperativ, voll komm [!]
- VVINFIN Infinitiv, voll gehen, ankommen
- VVIZU Infinitiv mit ``zu", voll anzukommen, loszulassen
- VVPP Partizip Perfekt, voll gegangen, angekommen
- VAFIN finites Verb, aux [du] bist, [wir] werden
- VAIMP Imperativ, aux sei [ruhig !]
- VAINFIN Infinitiv, aux werden, sein
- VAPP Partizip Perfekt, aux gewesen
- VMFIN finites Verb, modal dürfen
- VMINFIN Infinitiv, modal wollen
- VMPP Partizip Perfekt, modal gekonnt, [er hat gehen] können
- XY Nichtwort, Sonderzeichen 3:7, H2O, enthaltend D2XW3
- \$ Komma \$ Satzbeendende Interpunktion ? ! : ; \$( sonstige Satzzeichen: satzintern - [ ]()

# Penn Proposition (Prop) Bank (2000 – )

- Predicate/Argument structure (PAS) along syntactic subcategorization frames
  - P:Drink (A: Agent: x)
  - P:Drink (A: Patient: y)
- Focus on verbs (*events*) and their syntactic arguments (*participants*)
  - later phases: nominalizations, adjectives and prepositions
- Linguistic heritage:
  - Verb classes for the English language (Levin 1993)
  - with focus on semantic considerations (semantic or theta roles)
- Large coverage is a major goal



# Example for Propositions (PPB)



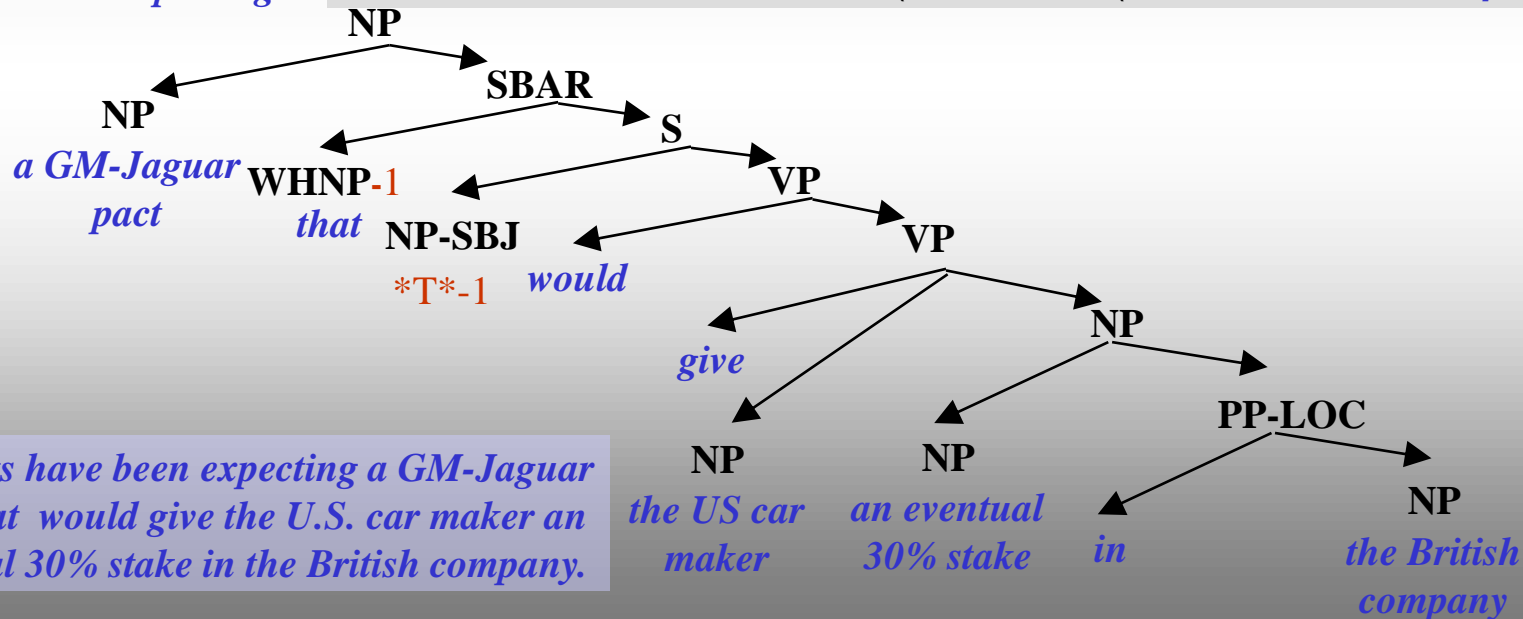
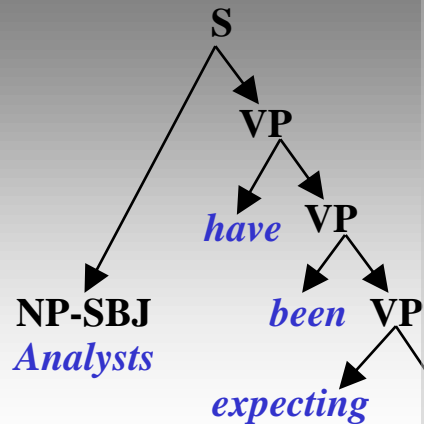
**When Bush met Blair on Thursday**

**they discussed the stabilization of the Iraq.**

**`meet(Bush, Blair)    discuss([Bush, Blair], stabilize(X, Iraq))`**

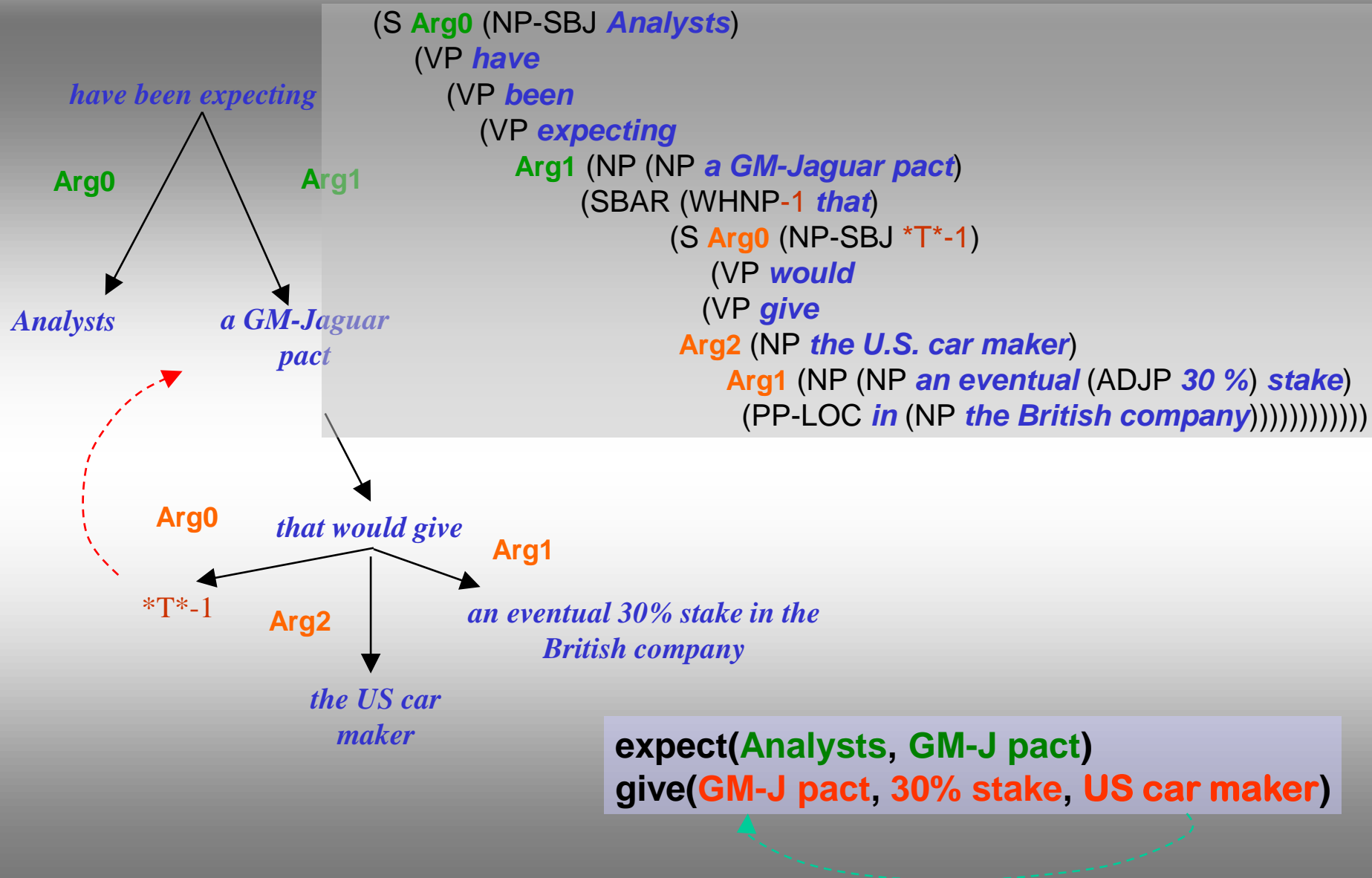
# Penn Treebank Sentence

(S (NP-SBJ *Analysts*)  
 (VP *have*  
 (VP *been*  
 (VP *expecting*  
 (NP (NP *a GM-Jaguar pact*)  
 (SBAR (WHNP-1 *that*)  
 (S (NP-SBJ *\*T\*-1*)  
 (VP *would*  
 (VP *give*  
 (NP *the U.S. car maker*)  
 (NP (NP *an eventual* (ADJP *30 %*) *stake*)  
 (PP-LOC *in* (NP *the British company*))))))))))



*Analysts have been expecting a GM-Jaguar pact that would give the U.S. car maker an eventual 30% stake in the British company.*

# Penn PropBank Sentence



# PPB Annotation Principles

- Search for the most frequently used predicates (verbs) in the PTB
- Survey of the „usage“ of a certain predicate
  - Considering the number of evidences in the corpus
  - Selection of roles which
    - occur frequently
    - are „semantically“ necessary
  - Indexing of roles (arguments) according to the (Arg0 ... Arg5) scheme yields distinct framesets for a verb
    - Arg0: prototypical agent
    - Arg1: prototypical patient or theme
    - Arg2-5: no systematic generalization applies
- Propositional annotation is based on a sentence's PTB parse structure and the availability of the framesets
- Additional annotation of verbs by temporal, aspectual and voice information (ArgMs)

# PPB Annotation Principles: Framesets

Frameset **accept.01** "take willingly"

Arg0: Acceptor

Arg1: Thing accepted

Arg2: Accepted-from

Arg3: Attribute

Ex: [Arg0 He] [ArgM-MOD would] [ArgM-NEG n't] *accept* [Arg1 anything of value] [Arg2 from those he was writing about]. (wsj\_0186)

Frameset **kick.01** "drive or impel with the foot"

Arg0: Kicker

Arg1: Thing kicked

Arg2: Instrument (defaults to foot)

Ex1: [ArgM-DIS But] [Arg0 two big New York banks<sub>i</sub>] seem [Arg0 \*trace\*<sub>i</sub>] to have *kicked* [Arg1 those chances] [ArgM-DIR away], [ArgM-TMP for the moment], [Arg2 with the embarrassing failure of Citicorp and Chase Manhattan Corp. to deliver \$7.2 billion in bank financing for a leveraged buy-out of United Airlines parent UAL Corp]. (wsj\_1619)

Ex2: [Arg0 John<sub>i</sub>] tried [Arg0 \*trace\*<sub>i</sub>] to *kick* [Arg1 the football], but Mary pulled it away at the last moment.

Frameset **edge.01** "move slightly"

Arg0: causer of motion      Arg3: start point

Arg1: thing in motion      Arg4: end point

Arg2: distance moved      Arg5: direction

Ex: [Arg0 Revenue] *edged* [Arg5 up] [Arg2-EXT 3.4%] [Arg4 to \$904 million] [Arg3 from \$874 million] [ArgM-TMP in last year's third quarter]. (wsj\_1210)

- Frames for more than 3,300 verbs exist
- 4,500 framesets exist indicating an average polysemy rate of 1.36
- Classical Zipfian distribution for framesets: ,go' has 20 FSs, ,come', ,get', ,make', ,take', etc. more than a dozen, 2,581 out of 3,342 verbs have only a single one

# PPB Annotation Principles (cont.)

- Extraction of all sentences which contain a given verb
- 1<sup>st</sup> run: automatic tagging

<http://www.cis.upenn.edu/~josephr/TIDES/index.html#lexicon>

- 2<sup>nd</sup> run: “Double blind hand correction”
  - Basically carried out by linguistics students (undergraduates)
  - Tagging tool highlights discrepancies
- 3<sup>rd</sup> run: “Salomonization”
  - Judge’s decision (by project leader?)
  - approximately 5% of the verbs are concerned

# PPB Inter-Annotator Agreement

$$\kappa = \frac{P(A) - P(E)}{1 - P(E)}$$

		$P(A)$	$P(E)$	$\kappa$
including ArgM	role identification	.99	.89	.93
	role classification	.95	.27	.93
	combined decision	.99	.88	.91
excluding ArgM	role identification	.99	.91	.94
	role classification	.98	.41	.96
	combined decision	.99	.91	.93

Subtypes of the ArgM modifier tag

LOC: location	CAU: cause
EXT: extent	TMP: time
DIS: discourse connectives	PNC: purpose
ADV: general-purpose	MNR: manner
NEG: negation marker	DIR: direction
MOD: modal verb	

- **P(A)** probability of interannotator agreement
- **P(E)** agreement expected by chance
- **ArgM** a set of adjunct-like arguments every verb can take in addition to semantic roles from its roleset

# Different Meanings of a Verb

*Mary called John an idiot.*  
(LABEL)

Arg0: Mary

Rel: called

Arg1: John (item being labeled)

Arg3-PRD: an idiot (attribute)

*Mary called John a cab.*  
(SUMMON)<sup>5</sup>

Arg0: Mary

Rel: called

Arg2: John (benefactive)

Arg1: a cab (thing summoned)



# Semantically Related Verbs – Meta Frames

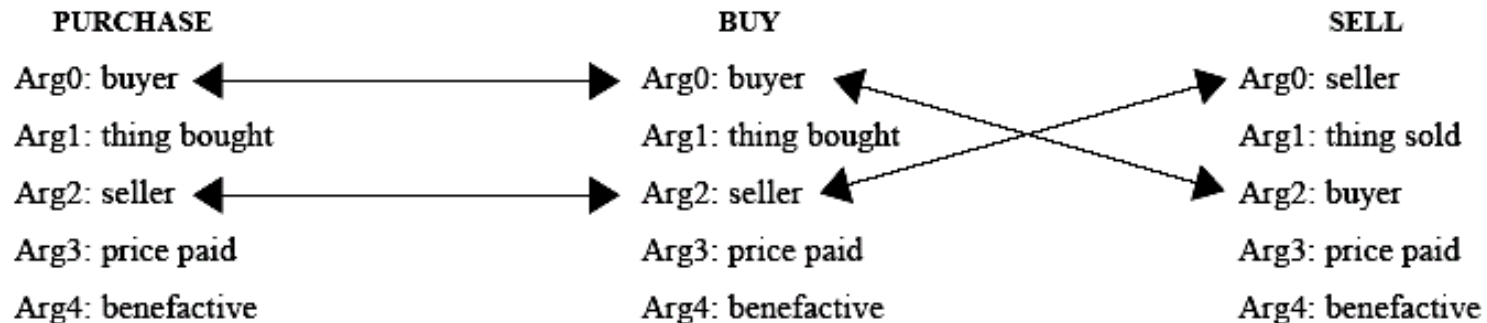


Table 1: comparison of arguments of semantically related verbs

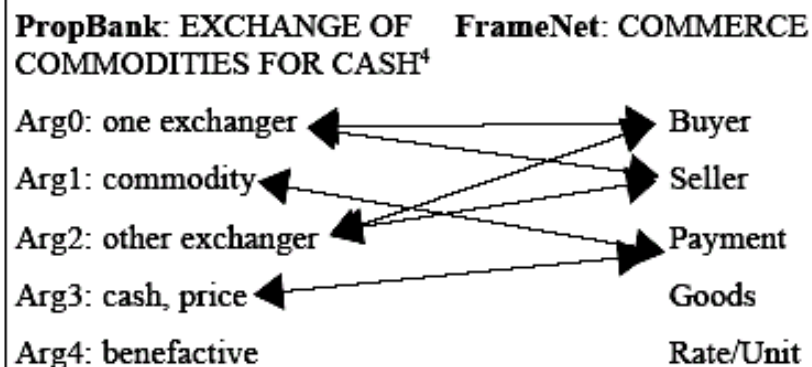


Table 2: comparison on PropBank and Framenet

# PPB Annotation Statistics

- Training time for PropBank annotators: +/- 3 days
  - Less than for syntactic (bracketing) annotators
- Semi-automatic pre-annotation by already existing frames (VerbNet – a generalization of Levin classes)
- Speed statistics
  - 25 verb frames per week
  - 50 (!?) predicates per person and hour
- average inter-annotator agreement: < 80%
  - Still, variance ranges between 60% and 100%
- There exists an arbiter „gold standard“
  - Agreement between annotators and gold standard ranges between 45% and 100%
- The larger the potential number of arguments for a verb, the higher the likelihood of disagreement

# SALSA Corpus

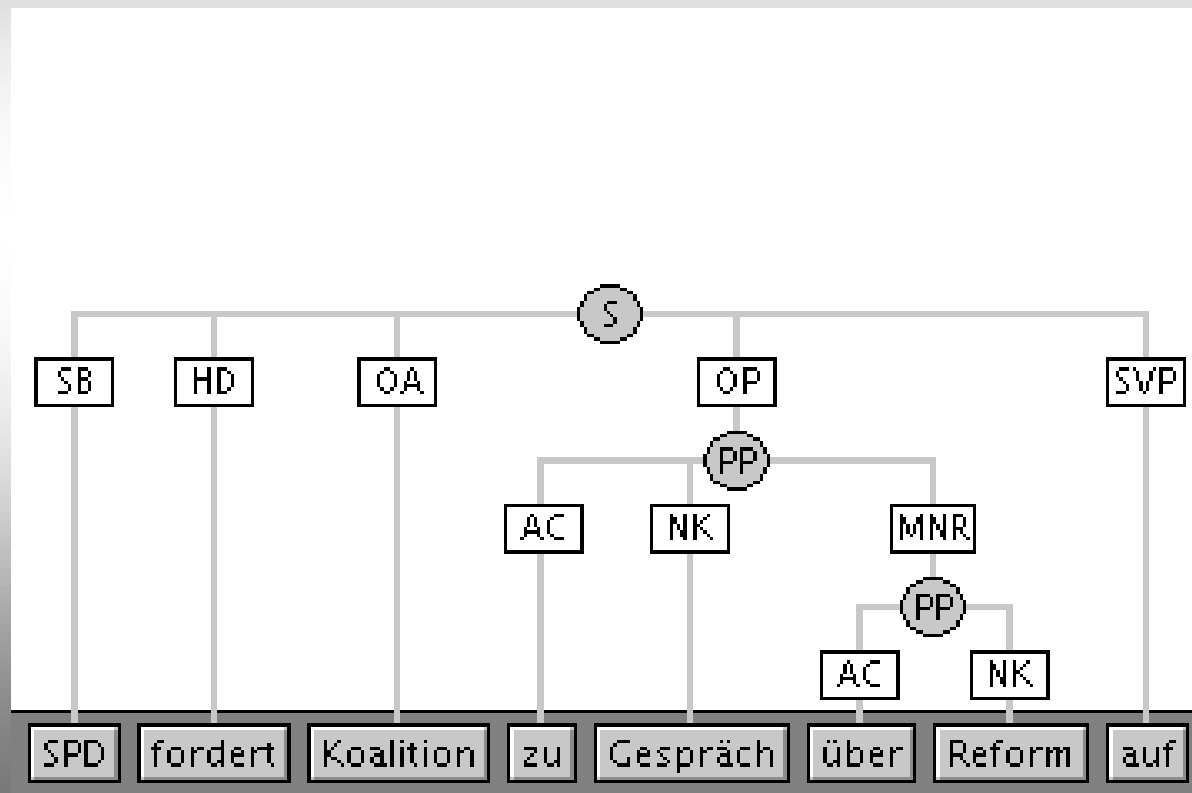
- Saarbrücken **L**exical **S**emantics  
**A**nnotation and Acquisition Project
- Bereitstellung einer großen lexikalisch-  
semantischen Ressource für Prädikat-  
Argument-Struktur im Deutschen
- Verbesserung der semantischen  
Verarbeitung auf der Ebene der Prädikat-  
Argument-Struktur
- <http://www.coli.uni-saarland.de/projects/salsa/>

# SALSA Ziele

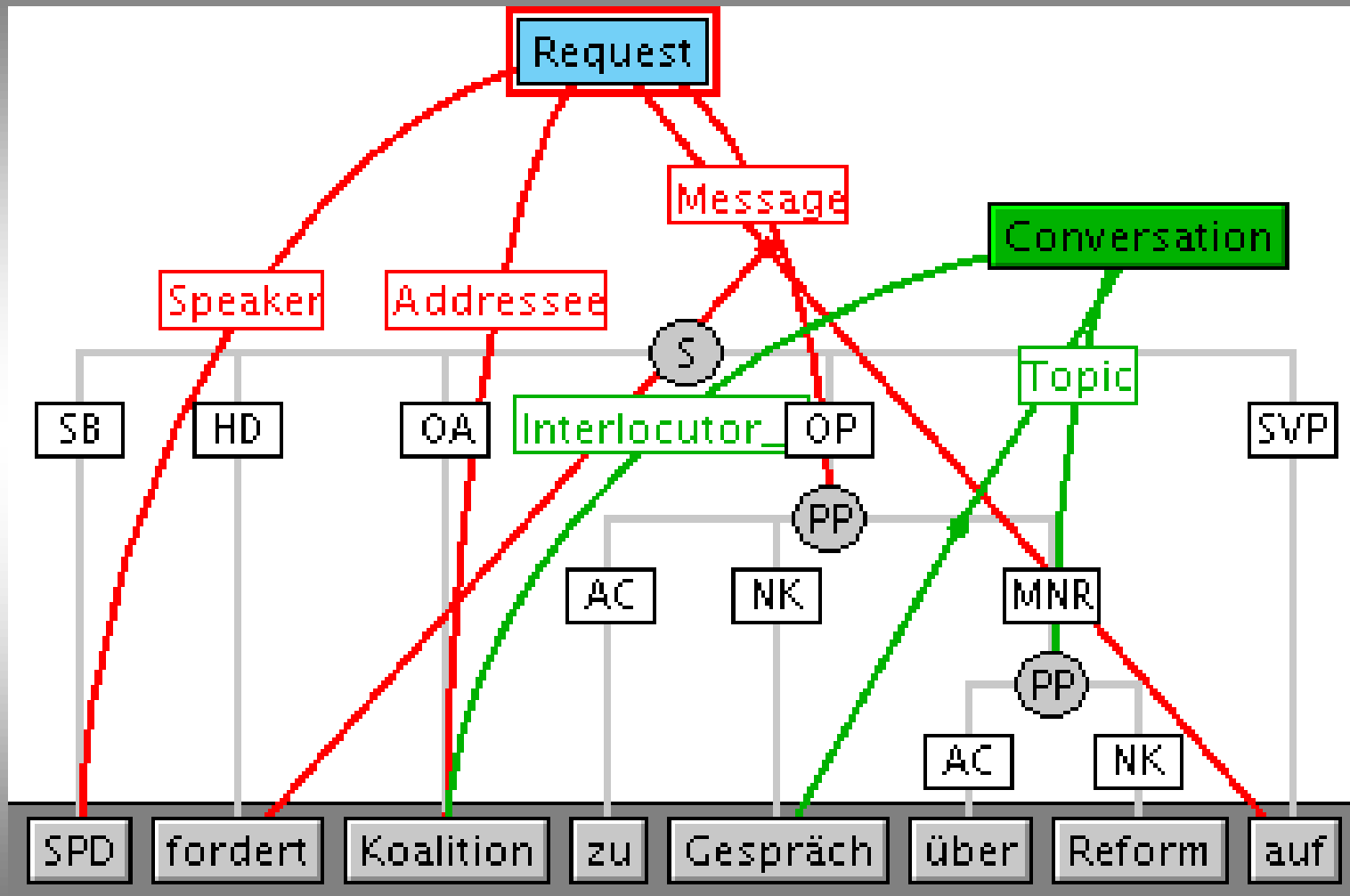
- Bereitstellung einer lexikalisch-semantischen Ressource (Korpus + Lexikon) für das Deutsche mit Informationen über:
  - Wortbedeutungen auf der Ebene Frame-semantischer Klassifikation von Prädikaten
  - Semantische Rollen und syntaktische Realisierungsmuster
- Entwicklung von Verfahren zur
  - Automatischen Akquisition lexikalisch-semantischer Information
  - Auswertung und Anwendung lexikalisch-semantischer Ressourcen

# SALSA Grundlagen

- Berkeley FrameNet-Datenbank
- TIGER-Korpus  
(Saarbrücken/Stuttgart/Potsdam):



# SALSA Annotation auf TiGER Syntax



# Sublanguage Corpora

- GENIA (U Tokyo)

- language: English (biomedical sublanguage)
- text genre: biology articles (*Medline bibliographic database*)
- size: 2,000 annotated abstracts (18,500 sentences, 491,000 tokens)
  - selected from a MeSH term search of “Human”, “Blood Cells” and “Transcription Factors”
- POS tagging based on PTB tag set
- Syntactic phrase structures (beta version); PTB-style treebank (200 abstracts only)
- Named entity annotation based on a subset of substances (peptides, amino acids, DNA), biological locations (organisms, tissues) involved in reactions of proteins (GENIA ontology) — 100,000 bio annotations

# Demo of GENIA

## Example:

*„Preincubation of cells with 1,25-(OH)<sub>2</sub>D<sub>3</sub> augmented IL-1 beta mRNA levels only in U-937 and HL-60 cells.“*



# POS Annotation in GENIA

Preincubation/**NN** of/**IN** cells/**NNS**  
with/**IN** 1,25-(OH)<sub>2</sub>D<sub>3</sub>/**NN**  
augmented/**VBD** IL-1/**NN** beta/**NN**  
mRNA/**NN** levels/**NNS** only/**RB** in/**IN**  
U-937/**NN** and/**CC** HL-60/**NN** cells/**NNS**  
./.

# Syntactic Annotation in GENIA

```
- <S>
- <NP-SBJ>
  <NP>Preincubation/NN</NP>
- <PP>
  of/IN
  <NP>cells/NNS</NP>
</PP>
- <PP>
  with/IN
  <NP>1,25-(OH)2D3/NN</NP>
</PP>
</NP-SBJ>
- <VP>
  augmented/VBD
  <NP>IL-1/NN beta/NN mRNA/NN levels/NNS</NP>
- <PP>
  only/RB in/IN
  - <NP>
    - <NP SYN="COORD">
      <NP>U-937/NN</NP>
      and/CC
      <NP>HL-60/NN</NP>
    </NP>
    cells/NNS
  </NP>
</PP>
</VP>
./
</S>
```

# Named Entity Annotation in GENIA

```
- <sentence>
  Preincubation of cells with
  <cons lex="1,25-(OH)2D3" sem="G#lipid">1,25-(OH)2D3</cons>
  augmented
- <cons lex="IL-1_beta_mRNA_level" sem="G#other_name">
- <cons lex="IL- 1_beta_mRNA" sem="G#RNA_molecule">
  <cons lex="IL-1_beta" sem="G#protein_molecule">IL-1 beta</cons>
  mRNA
  </cons>
  levels
</cons>
only in
<cons lex="U-937" sem="G#cell_line">U-937</cons>
and
<cons lex="HL-60" sem="G#cell_line">HL-60</cons>
cells.
</sentence>
```

# Medical Sublanguage vs. General Language

- Medical language as a sublanguage
  - (ad hoc) abbreviations and acronyms (*o.B.*, *V.a.*, *COPD*)
  - (idiosyncratic) measure units (*mmHg*, *mm Hg*)
  - variable forms of enumeration patterns (*1.,2.,..., a),b)...*)
  - Latin-/Greek-based terminology (*ulcus ventriculi*)
- However: less complexity and variation than general language
- Expect standard general-language-trained off-the-shelf POS taggers to perform 'ok'
- Statistically significant performance gain for biomedical POS taggers when trained on dedicated biomedical corpora (Wermter & Hahn, 2004)

# Infrastructure Requirements

- **Definition of Description Languages for**
  - Tagging/NER: Tag Set (Syntactic, Semantic)
  - Chunking/Parsing: Grammar Format
  - Proposition Analysis: Proposition Format, Ontology (Concept System, Relation Types)
  - Discourse Analysis: Reference and rhetorical relations
- **Manual Creation of Corpora**
  - Training Coders in Applying Description Languages
  - Test of Coder Reliability
- **Benefit:**
  - Solid Foundation for Supervised Learning

# Resources for NLP

- Empirical (Learning) Paradigm for NLP
- **Types of Resources**
  - Language data (plain, annotated)
  - **Systems for rendering language meta data**
  - Computational lexicons and ontologies
  - NLP Core Engines
  - NLP Application Systems
  - Machine Learning Resources
- Methodological Issues of NLP Resources

# Systems for Rendering Language Meta Data

- Software infrastructure which supports the manual annotation processes at all levels
- Easy adaptation to user-defined annotation languages
- Visualization component
  - In-line vs. stand-off
  - Semantics of colors
  - Graphical overlay structures
- Team support mechanisms wrt annotation
  - Comparison of annotator pairs/groups
  - Consensus seeking
  - Built-in quality evaluation schemes (annotator agreement)
- Software engineering standards
  - Version control (of annotation software)
  - Change history (of annotation products)

# Manual Annotation – Workflow

<https://webanno.googlecode.com/svn/tags/latest-stable/docs/user-guide.html>

