

# Discourse Processing

## A Tutorial

**Manfred Stede**  
Applied Computational Linguistics  
EB Cognitive Science  
University of Potsdam  
Germany

- Most of the material is based on

Manfred Stede: Discourse Processing. Morgan  
& Claypool 2011.

## „Discourse Processing“

- Where does it start?
- Treat your text document
  - not as a bag of words,
  - not as a bag of sentences,
  - but as a linear **sequence of connected sentences**, adding **hierarchical structure** where appropriate

© Manfred Stede / NAACL Tutorial 2013

## Discourse-sensitivity in NLP Tasks

- Example: **information extraction** (Wikipedia)

*Angela Dorothea Merkel (born 17 July 1954) is a German politician who has been the Chancellor of Germany since 2005, and the Leader of the Christian Democratic Union (CDU) since 2000. She is the first woman to hold either office.*

© Manfred Stede / NAACL Tutorial 2013

## Discourse-sensitivity in NLP Tasks

- Example: **summarization**

When you know the global structure of your text document, you can make sure to cover relevant portions

(examples will follow...)

© Manfred Stede / NAACL Tutorial 2013

## Discourse-sensitivity in NLP Tasks

- Example: **opinion mining**

*(...) We didn't like the village very much and won't come back because there are only few things to see, and moreover the place is quite dirty. (...)*

© Manfred Stede / NAACL Tutorial 2013

## Goal for today

- Provide an overview of the major problems of discourse processing on text documents, and on the central ideas for tackling them


© Manfred Stede / NAACL Tutorial 2013

## Tutorial overview

- Part 1: „Large“ discourse units
  - Genre-specific structure
  - Topics
- Part 2: Coreference
- Part 3: „Small“ discourse units
  - Local coherence
  - Coherence-relational text structure
- Exploring interconnections

© Manfred Stede / NAACL Tutorial 2013

# Genre: Court decision



**BUNDESGERICHTSHOF**  
IM NAMEN DES VOLKES

**URTEIL**

VI ZR 101/06


Verkündet am:  
27. März 2007  
Holmes,  
Justizangestellte  
als Urkundsbeamtin  
der Geschäftsstelle

in dem Rechtsstreit

Nachschlagewerk: ja  
BGHZ: nein  
BGHR: ja

BGB § 823 Abs. 1, § 1004, SGB § 185, TMG § 10

Ein Unterlassungsanspruch wegen eines in ein Meinungsforum im Internet eingestellten ehrverletzenden Beitrags kann auch dann gegen den Betreiber des Forums gegeben sein, wenn dem Verletzten die Identität des Autors bekannt ist.



**BUNDESGERICHTSHOF**  
IM NAMEN DES VOLKES

**URTEIL**

VI ZR 101/06

Verkündet am:  
27. März 2007  
Holmes,  
Justizangestellte  
als Urkundsbeamtin  
der Geschäftsstelle

in dem Rechtsstreit

Nachschlagewerk: ja  
BGHZ: nein  
BGHR: ja

BGB § 823 Abs. 1, § 1004, SGB § 185, TMG § 10

Ein Unterlassungsanspruch wegen eines in ein Meinungsforum im Internet eingestellten ehrverletzenden Beitrags kann auch dann gegen den Betreiber des Forums gegeben sein, wenn dem Verletzten die Identität des Autors bekannt ist.

© Manfred Stede / NAACL Tutorial 2013

# Genre: Scientific paper

## (Teufel et al. 09)

### Synthesis of pyrazole and pyrimidine Troeger's base-analogues

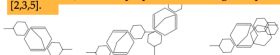
Rodrigo Abenia, Andrea Albernez, Hector Larrabondo, Jairo Quiroga, Braulio Isuasty, Henry Isuasty, Angelina Hormaza, Adolfo Sanchez, and Manuel Nogueras

**1**  
PERKIN

Troeger's-base analogues bearing fused pyrazolic or pyrimidinic rings were prepared in acceptable to good yields through the reaction of 3-alkyl-5-amino-1-arylpiperazines and 6-aminopyrimidin-4(3H)-ones with formaldehyde under mild conditions (i.e. in ethanol at 50°C in the presence of catalytic amounts of acetic acid). Two key intermediates were isolated from the reaction mixtures, which helped us to suggest a sequence of steps for the formation of the Troeger's bases obtained. The structures of the products were assigned by <sup>1</sup>H and <sup>13</sup>C NMR, mass spectra and elemental analysis and confirmed by X-ray diffraction for one of the obtained compounds.

#### Introduction

Although the first Troeger's base 1 was obtained more than a century ago from the reaction of p-toluidine and formaldehyde [1], recently the study of these compounds has gained importance due to their potential applications. They possess a relatively rigid chiral structure which makes them suitable for the development of possible synthetic enzyme and artificial receptor systems [2], chelating and biomimetic systems [3] and transition metal complexes for regio- and stereoselective catalytic reactions [4]. For these reasons, numerous Troeger's-base derivatives have been prepared bearing different types of substituents and structures (i.e. 2-5 Scheme 1), with the purpose of increasing their potential applications [2,3,5].



Scheme 1 The original Troeger's-base 1 and some interesting derivatives and analogues.

However, some of the above methodologies possess tedious work-up procedures or include relatively strong reaction conditions, such as treatment of the starting materials for several hours with an ethanolic solution of conc. hydrochloric acid or TFA solution, with poor to moderate yields, as is the case for analogues 4 and 5 [5].

Considering these potential applications, we now report a simple synthetic method for the preparation of 5,12-dialkyl-3,10-diaryl-1,3,4,8,10,11-hexaazatetracyclo[6.6.1.0.2.6-0.9.13]pentadeca-2(6,4,9)(15,11-tetraenes 8a-e and 4,12-dimethoxy-1,3,5,9,11,13-hexaazatetracyclo[7.7.1.0.2,7.0.10,15-1]heptadeca-2(7,3,10)(15)11-tetraene-6(14)-diones 10a,b based on the reaction of 3-alkyl-5-amino-1-arylpiperazines 6 and 6-aminopyrimidin-4(3H)-ones 9 with formaldehyde in ethanol and catalytic amounts of acetic acid. Compounds 8 and 10 are new Troeger's base analogues bearing heterocyclic rings instead of the usual phenyl rings in their aromatic parts.

#### Results and discussion

In an attempt to prepare the benzotriazolyl derivative 7a, which could be used as an intermediate in the synthesis of new hydroquinolines of interest, [6], a mixture of 5-amino-3-methyl-1-phenylpiperazine 6a, formaldehyde and benzotriazole in 10 ml of ethanol, with catalytic amounts of acetic acid, was heated at 50°C for 5 minutes. A solid precipitated from the solution while it was still hot. However, no consumption of benzotriazole was observed at TLC.

The reaction conditions were modified and the same product was obtained when the reaction was carried out without using benzotriazole, as shown in Scheme 12. On the basis of NMR and mass spectra and X-ray crystallographic analysis we established that the structure is 5,12-dialkyl-3,10-diaryl-1,3,4,8,10,11-hexaazatetracyclo[6.6.1.0.2.6-0.9.13]penta-deca-2(6,4,9)(15,11-tetraene 8, a new pentacyclic Troeger's base analogue.

Co\_Gro Other Aim Gap/Weak Own\_Mthd Own\_Res Own\_Conc

## Teufel et al. 09: Features

- Absolute location
  - pos. of sentence in doc
- Explicit structure
  - pos. of sentence in sct
  - pos. of sentence in para
  - type of headline of sct
- Sentence length
  - >12
- Content features
  - sentence contains words from title or headlines
  - sentence contains tf/idf-prominent terms
- Verb syntax
  - voice of first verb in sentence
  - tense of first verb in sentence
  - Is first verb modified by aux
- Citations
  - is citation present
  - self or other
  - positions
- History
  - most likely previous zone
- Meta-discourse
  - type of formulaic expr (28)
  - type of agent (10)
  - verb class and negation presence (28)

© Manfred Stede / NAACL Tutorial 2013

## Teufel et al. 09: Results

- Chemistry / CompLing
- 110pp guidelines
- 30 papers annotated
- Human Kappa = .71 / .65
- Automatic Kappa = .41 / .39
- Zone classification f-measure: 26%-86%

© Manfred Stede / NAACL Tutorial 2013

## Genre: Newspaper Editorial

### **Should Berlin apply for the 2016 Olympics?**

[2] Hamburg has long understood: [3] Olympic games are worth a lot of gold. [4] Those who draw the Olympics into their city are winners in the world-wide competition for attention. [5] That's why Berlin must not miss the opportunity for the 2016 games. [6] The capital must grab the baton from Leipzig and apply to be the venue. [7] Barcelona has shown that the olympic effect is invaluable. [8] With the 1992 games the city has re-invented itself -- and makes profit up to today: [9] The number of overnight stays has doubled, the economy is still profiting. [10] When Berlin now runs again as applicant, we show the world that we're better now than we were once. [11] After all, today the city offers what a candidate needs: [12] big-city flair, hotel beds, infrastructure. [13] The sports venues planned for 2000, such as Velodrom and Max-Schmeling-Halle, exist, the olympic stadium is in mint condition, the Anschutz arena is nearing completion. [14] Just by re-applying, Berlin would already modernize itself and improve its international profile. [15] Public and private sponsoring money would pour in, millions would follow from IOC. [16] And even if a European city turns out to be the venue for 2012: [17] One has to flex one's muscles in order to win the games, if necessary with the third instead of the second application. [18] Berlin to the starting block: [19] On your mark, ready, go!

### Should Berlin apply for the 2016 Olympics?

[2] Hamburg has long understood: [3] Olympic games are worth a lot of gold. [4] Those who draw the Olympics into their city are winners in the world-wide competition for attention. [5] That's why Berlin must not miss the opportunity for the 2016 games. [6] The capital must grab the baton from Leipzig and apply to be the venue. [7] Barcelona has shown that the olympic effect is invaluable. [8] With the 1992 games the city has re-invented itself -- and makes profit up to today: [9] The number of overnight stays has doubled, the economy is still profit-ing. [10] When Berlin now runs again as applicant, we show the world that we're better now than we were once. [11] After all, today the city offers what a candidate needs: [12] big-city flair, hotel beds, infrastructure. [13] The sports venues planned for 2000, such as Velodrom and Max-Schmeling-Halle, exist, the olympic stadium is in mint condition, the Anschutz arena is nearing completion. [14] Just by re-applying, Berlin would already modernize itself and improve its international profile. [15] Public and private sponsoring money would pour in, millions would follow from IOC. [16] And even if a European city turns out to be the venue for 2012: [17] One has to flex one's muscles in order to win the games, if necessary with the third instead of the second application. [18] Berlin to the starting block: [19] On your mark, ready, go!

## Content zones of *Pro&Contra* Editorial

- (Introduction, exposition of the problem)
- Central thesis of author
- Argument, pro-author
- (Counter-argument, contra-author)
- (Refutation of counter-argument)
- (Background information)
- (Final statement, rhetorical ending)

(optional zone)

linear position fixed



## Genre: Film review

### *The Draughtsman's Contract*

by James Mackenzie

James Mackenzie is currently finishing a Bachelor of Arts at Adelaide University, majoring in philosophy.



*The Draughtsman's Contract* (1982 UK 108 mins)

Source: CAC/NLA Prod Co: BFI/Channel 4 Prod: David Payne Dir, Scr: Peter Greenaway Ph: Curtis Clark Ed: John Wilson Art Dir: Bob Ringwood Mus: Michael Nyman

Cast: Anthony Higgins, Janet Suzman, Anne Louise Lambert, Hugh Fraser, Suzanne Crowley, Neil Cunningham

In 1981, Peter Greenaway spent a warm summer drawing his house on the Welsh Border. So that the house was in the same light every time he sketched it, eight set views were drawn at eight set times of the day. The carefulness of the plan balanced the chaos of its enactment: cows, neighbours and children were equally constant, if pleasurable, interruptions. Anybody who has seen *The Draughtsman's Contract* may find this scenario familiar. It was the basis of the film. But the relevance of any biographical detail stops there.

Like all of Greenaway's work, *The Draughtsman's Contract* addresses issues of representation: between mediums (drawn and photographic representation); between art and nature (the inbreeding between landscape art and ornamental garden design); and the value of classical naturalism in art.

© Manfred Stede / NAACL Tutorial 2013

## Genre structure – Content zones (Bieler et al. 07)

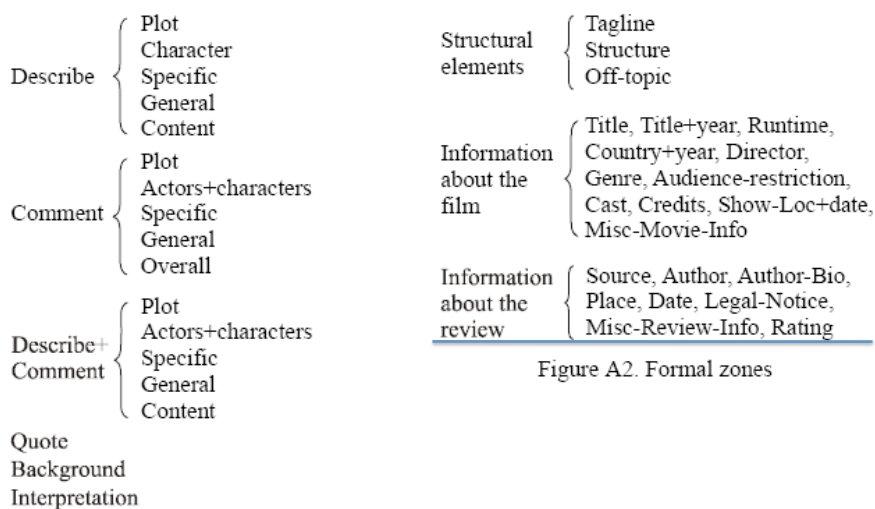


Figure A2. Formal zones

Figure A1. Functional zones

## Describe vs. Comment Paragraphs

Architect Stourley Kracklite (Brian Dennehy) arrives in Rome, where an exhibition of the works of the 18th-century architect Etienne-Louis Boullée is being mounted under Kracklite's supervision. The city - or something - doesn't sit with him; upon arrival, he begins complaining of stomach pains. Cancer? Kracklite is sure of it. Or not: It could be that his wife Louisa (Chloe Webb), with whom he is traveling (and who is pregnant with his child), is poisoning him, a revenge for his self-absorption. She may be further motivated in this by the affair she has taken up with Caspasian Speckler (Lambert Wilson), another architect involved with the exhibition.

Greenaway has an eye for composition, and in *The Belly of an Architect* many formal arrangements stand out for their beauty. Dennehy, always engaging, is slyly illegible in the central role, a stroke of luck maybe for the director, who has shown himself to be disinclined to bother much with actors. But the relentless condescension and self-congratulation with which Greenaway conducts this very private amusement is grotesque. He fosters the worst imaginable relationship with his audience: showing off while condemning those not enlightened enough to cherish his preening. In Kracklite, Greenaway has created a self-obsessed, boorish non-hero on whom to hang his obscurantist ramblings, and his indifference to his audience is so great that he expects us to relish it. Who's the asshole here?

© Manfred Stede / NAACL Tutorial 2013

## Describe vs. Comment: bag-of-terms (Bieler et al. 07)

- Supervised learning (SVM)
- Development set: 100 German reviews from 5 internet sites
- Training/test set: 112 reviews, 60 of which come from three „new“ sites
  - Training set (66%)
  - Test set (33%)
- Features = character 5-grams, weighted with TF/IDF against reference corpus

Zone type	<i>Comment</i>	<i>Describe</i>
Precision	81.6	76.8
Recall	79.7	79.0

© Manfred Stede / NAACL Tutorial 2013

## Classifying zones: English texts (Taboada et al. 09)

- **Manual annotation**

- Three annotators to check agreement

Classes	2-rater kappa	3-rater kappa
Describe/Comment/Describe+Comment/Formal	.82	.73
Describe/Comment/Formal	.92	.84
Describe/Comment/Describe+Comment	.68	.54
Describe/Comment	.84	.69

Table B1. Kappa values for stage annotations

- Then, one annotator annotated 100 texts from RottenTomatoes.com
- 83.000 words in 1.500 paragraphs

Stage	Count
Describe	347
Comment	237
Describe+Comment	237
Background	51
Interpretation	22
Quote	2
Formal	646

Table 1. Stages in 100 text RT corpus

© Manfred Stede / NAACL Tutorial 2013

## Automatic classification: Features

- **Character 5-grams** that appear at least 4 times in the corpus
- **Genre-based**
  1. From Biber (88)
    - 1, 2, 3 person pronouns; demonstrative pronouns
    - Place and time adverbials
    - Intensifiers
    - Modals
  2. Connectives that indicate contrast, comparison, causation, evidence, condition
  3. List of 500 adjectives classified in terms of Appraisal (Martin and White 2005)
    - Appreciation, Judgment or Affect
  4. Text statistics
    - Average length of words and sentences; position of paragraphs in the text

© Manfred Stede / NAACL Tutorial 2013

## Classifier performance

2, 3, and 4 classes

Comment, Describe

Comment, Describe, Formal

Comment, Describe, Comment+Describe, Formal

10-fold cross validation

Classifier	Comment			Describe			Formal			Desc+Comm			Overall Accuracy
	P	R	F	P	R	F	P	R	F	P	R	F	
2-class-5-gram-Bayes	.66	<b>.79</b>	.72	.70	.55	.62	-	-	-	-	-	-	68.0
2-class-5-gram-SVM	.53	.63	.64	.68	.69	<b>.69</b>	-	-	-	-	-	-	66.8
2-class-genre-Bayes	.66	.75	.70	.67	.57	.61	-	-	-	-	-	-	66.2
2-class-genre-SVM	<b>.71</b>	.76	<b>.74</b>	<b>.71</b>	.65	.68	-	-	-	-	-	-	<b>71.1</b>
3-class-5-gram-Bayes	.69	.49	.57	.66	<b>.78</b>	<b>.71</b>	<b>.92</b>	<b>.97</b>	<b>.95</b>	-	-	-	<b>78.1</b>
3-class-5-gram-SVM	.64	.63	.63	<b>.68</b>	.65	.65	.91	<b>.97</b>	.94	-	-	-	77.2
3-class-genre-Bayes	<b>.68</b>	.68	.66	.67	.46	.55	.84	.96	.90	-	-	-	74.0
3-class-genre-SVM	.66	<b>.71</b>	<b>.68</b>	.67	.56	.61	.90	.94	.92	-	-	-	76.8
4-class-5-gram-Bayes	<b>.46</b>	.35	.38	<b>.69</b>	.47	.56	<b>.92</b>	<b>.97</b>	<b>.95</b>	<b>.42</b>	<b>.64</b>	<b>.51</b>	69.0
4-class-5-gram-SVM	.43	<b>.41</b>	.44	.59	.62	.60	.91	<b>.97</b>	.94	.45	.41	.42	<b>69.6</b>
4-class-genre-Bayes	.38	.31	.34	.66	.30	.41	.86	<b>.97</b>	.90	.33	.60	.42	62.3
4-class-genre-SVM	.46	.32	.38	.53	<b>.82</b>	.65	.87	.94	.90	.26	.03	.06	67.4

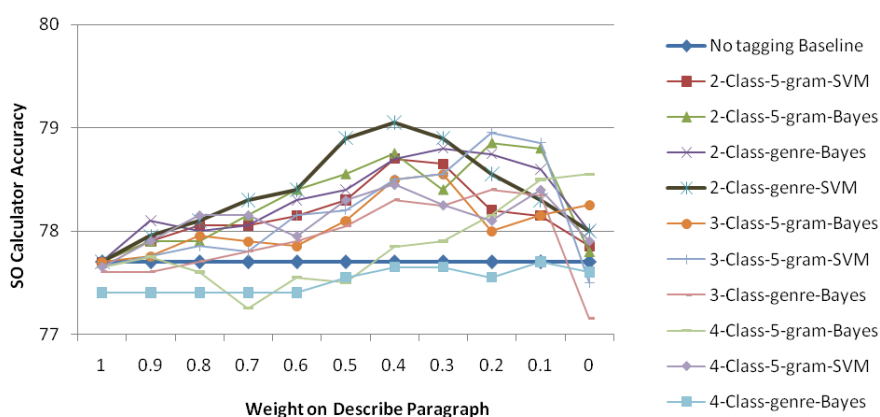
Table 2. Stage identification performance of various categorical classifiers

## Content zones and Sentiment detection

- Idea: Use zone-classifier to improve the performance of text-level sentiment classification by *SO-CAL* (Taboada et al. 11) on movie reviews
  - Disregard **description** altogether (weight of 0)
  - Give it a lower weight than to **comment**
- Evaluation performed on the *Polarity Dataset*, a collection of 2,000 movie reviews (Pang/Lee 04)
  - Run the zone classifiers to label paragraphs
  - Run *SO-CAL* on the texts, with different weights assigned to each type of paragraph

## SO-CAL and weighing paragraphs

- On the Polarity Dataset, 2,000 movie reviews (Pang and Lee 2004)



© Manfred Stede / NAACL Tutorial 2013

## SO-CAL and weighing paragraphs

- Most classifiers improve performance over the 77.7% baseline
- Using manual zone annotations (the original 100 texts used to build the classifiers), performance is boosted by **12%**
- Precision of the classifier is key
  - Need to identify Describe paragraphs accurately
  - In lieu of that, weighting is a better strategy than removing all Describe paragraphs

© Manfred Stede / NAACL Tutorial 2013

## Wrap-up: Definitions

- **Content zone**

A continuous portion of a text document that fulfills a functional role for the text as a whole, contributing to the overall message or purpose, as it is characteristic for the genre of the text.

- **Genre**

A class of texts that fulfill a common function, are being used for a common communicative purpose, and are potentially subject to conventions on various levels of description:

- length
- layout (headlines, pictures, tables, diagrams, enumerations, ...)
- lexical and syntactic features
- internal organization
  - Which zones are obligatory, which are optional
  - Constraints and preferences on zone order

© Manfred Stede / NAACL Tutorial 2013

## Tutorial overview

- Part 1: „Large“ discourse units
  - Genre-specific structure
  - Topics
- Part 2: Coreference
- Part 3: „Small“ discourse units
  - Local coherence
  - Coherence-relational text structure
- Exploring interconnections

© Manfred Stede / NAACL Tutorial 2013

[1.7] A man named Lionel Gaedi went to the Port-au-Prince morgue in search of his brother, Josef, but was unable to find his body among the piles of corpses that had been left there. [1.8]“I don’t see him—it’s a catastrophe,” Gaedi said. [1.9]“God gives, God takes.” [1.10] Chris Rolling, an American missionary and aid worker, tried to extricate a girl named Jacqueline from a collapsed school using nothing more than a hammer. [1.11] He urged her to be calm and pray, and as night fell he promised that he would return with help. [1.12] When he came back the next morning, Jacqueline was dead. [1.13]“The bodies stopped bothering me after a while, but I think what I will always carry with me is the conversation I had with Jacqueline before I left her,” Rolling wrote afterward on his blog. [1.14]“How could I leave someone who was dying, trapped in a building! ...[1.15] She seemed so brave when I left! [1.16] I told her I was going to get help, but I didn’t tell her I would be gone until morning. [1.17] I think this is going to trouble me for a long time.” [1.18] Dozens of readers wrote to comfort Rolling with the view that his story was evidence of divine wisdom and mercy.

Source: *The New Yorker*, 2010 (Copyright Condé Nast)

[1.7] A man named Lionel Gaedi went to the Port-au-Prince morgue in search of his brother, Josef, but was unable to find his body among the piles of corpses that had been left there. [1.8]“I don’t see him—it’s a catastrophe,” Gaedi said. [1.9]“God gives, God takes.” [1.10] Chris Rolling, an American missionary and aid worker, tried to extricate a girl named Jacqueline from a collapsed school using nothing more than a hammer. [1.11] He urged her to be calm and pray, and as night fell he promised that he would return with help. [1.12] When he came back the next morning, Jacqueline was dead. [1.13]“The bodies stopped bothering me after a while, but I think what I will always carry with me is the conversation I had with Jacqueline before I left her,” Rolling wrote afterward on his blog. [1.14]“How could I leave someone who was dying, trapped in a building! ...[1.15] She seemed so brave when I left! [1.16] I told her I was going to get help, but I didn’t tell her I would be gone until morning. [1.17] I think this is going to trouble me for a long time.” [1.18] Dozens of readers wrote to comfort Rolling with the view that his story was evidence of divine wisdom and mercy.

## Definition

- Topic-induced text structure

A sequence of non-overlapping text segments that completely covers the text, i.e., a partitioning. Each unit consists of one or more sentences that address a common topic.

© Manfred Stede / NAACL Tutorial 2013

## Surface Cues

- Paragraph breaks
  - *Prima facie* a good cue, but is to be handled with care
  - Cautious approach: If a text has many paragraph breaks, topic breaks are unlikely to occur *within* paragraphs
- Connectives
  - can indicate continuity: *also, then, so, ...*
  - can indicate discontinuity: *but, still, besides, ...*
- Pronouns
  - often indicate continuity
- Syntax: information structure
  
- => altogether not very reliable (or difficult to compute)

© Manfred Stede / NAACL Tutorial 2013



## Content Words (1): Lexical chains

- Intuition: topic continuity means that words in contiguous sentences are related
- Steps:
  - (a) Compute relations between individual words
  - (b) Build *chains* from those relations
  - (c) Induce boundaries from the topology of chains

© Manfred Stede / NAACL Tutorial 2013

### (a) Establish Lexical Relations

- Which words are to be considered as candidates?
  - Part of speech: usually nouns and NEs
  - Distance:
    - Fixed: a few sentences
    - Flexible: Hirst/St. Onge (98) allow longer distance for stronger relations

© Manfred Stede / NAACL Tutorial 2013

## (a) Establish Lexical Relations

- Measurement
  - WordNet path length (e.g. *WordNet Connect*)  
**Problem:** need to do word sense disambiguation, or not (see Silber/McCoy 02)
  - Distributional similarity

© Manfred Stede / NAACL Tutorial 2013

## (b) Build chains

- Move through text, deciding for each word to
  - start a new chain
  - connect to an existing chain  
**Caveat:** do not compare to the last item only!  
(Morris/Hirst 91: *cow – sheep – wool – scarf – boots – hat – snow*)
- Need weights and thresholds for chain length, chain density, chain distance

© Manfred Stede / NAACL Tutorial 2013

## (c) Induce topic boundaries

- Compute boundary strength for each sentence break
  - Number of chains ending, beginning, crossing
- Accept the top  $n$  boundary candidates
- Some implementations available, e.g. *LCSeg*  
[www.cs.columbia.edu/nlp/tools.cgi](http://www.cs.columbia.edu/nlp/tools.cgi)

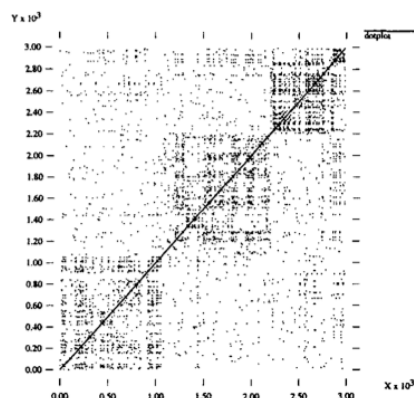
© Manfred Stede / NAACL Tutorial 2013

## Content Words (2): Vocabulary shifts

- Can we find those boundaries without using lexical resources?
- Straightforward idea:
  - Word repetition
  - Introduction of new words

© Manfred Stede / NAACL Tutorial 2013

## Reynar 94: *dot plotting*



- Closed-class words removed
- Forms of *to be*, *to have* removed
- Lemmatization applied

Word appears at pos  $x$  and  $y$   
 $\Rightarrow$  dot at  $(x,x)$   $(y,y)$   $(x,y)$   $(y,x)$

Figure 1: The dotplot of four concatenated *Wall Street Journal* articles.

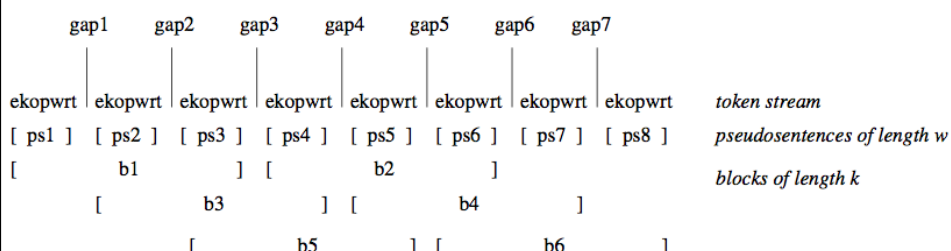
At least for synthetic data, this works at least as good as lexical chaining (e.g., Stokes et al. 04)

© Manfred Stede / NAACL Tutorial 2013

- Skorochood'ko (72):
  - divide the text into sentences,
  - count the overlap in content words between neighboring sentences,
  - postulate topic boundaries on the basis of the overlap count.
- Various implementations, including *Text Tiling* (Hearst 97)
  - work on *expository* text (as opposed to narrative, descriptive, argumentative, instructive)

© Manfred Stede / NAACL Tutorial 2013

## Text tiling (Hearst 1997)



- Recommendation:  $w=20$ ,  $k=6$
- At every pseudosentence boundary, compute similarity between blocks meeting there

© Manfred Stede / NAACL Tutorial 2013

## Text tiling: similarity

- Hearst 97: Similarity via cosine metric

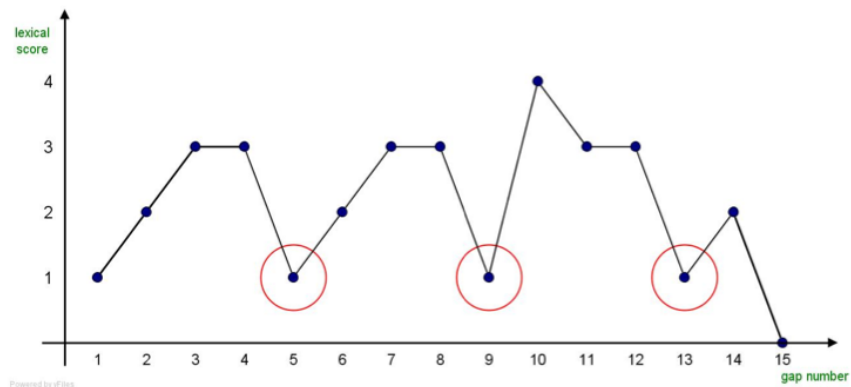
$$\text{score}(\text{gap}_i) = \frac{\sum_t w_{t,b1} w_{t,b2}}{\sqrt{\sum_t w_{t,b1}^2 \sum_t w_{t,b2}^2}}$$

- Choi 00: Need to consider the distribution of words across the whole text – not just locally
- Dias et al. 07: word relevance also depends on *inverse document frequency*, adapted to sentences:

$$\text{tf.isf}(\text{word}) = \frac{\text{stf}(\text{word}, s)}{|s|} * \ln \frac{N_s}{\text{sf}(\text{word})}$$

© Manfred Stede / NAACL Tutorial 2013

## Text tiling: boundaries

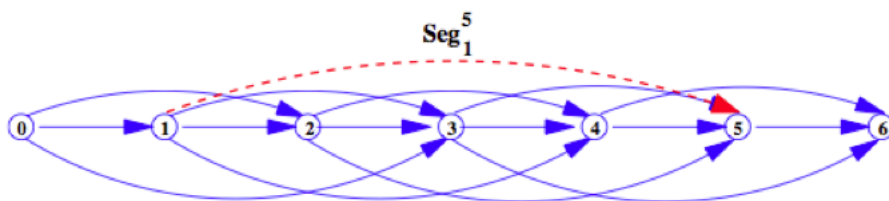


Hearst 97: compute depth score for each gap, then accept the top  $n$  gaps  
 $\text{depth}(\text{gap}_i) = (y_{i-1} - y_i) + (y_{i+1} - y_i)$

Dias et al. 07: compare steepness of the function at valleys

## Other implementations

- Shortest-path problem (Misra et al. 09)



## Other options

- Combine word distribution analysis with surface cues, e.g. Beeferman et al. (99), Galley et al. (03)
- Current line of work: compute *hidden topics*, using latent Dirichlet allocation (LDA), e.g. Eisenstein (09)

© Manfred Stede / NAACL Tutorial 2013

## Tutorial overview

- Part 1: „Large“ discourse units
  - Genre-specific structure
  - Topics
- Part 2: Coreference
- Part 3: „Small“ discourse units
  - Local coherence
  - Coherence-relational text structure
- Exploring interconnections

© Manfred Stede / NAACL Tutorial 2013

## Ref. exprs. / Mentions / Markables

[1.7] A man named Lionel Gaedi went to the Port-au-Prince morgue in search of his brother, Josef, but was unable to find his body among the piles of corpses that had been left there. [1.8]“I don’t see him—it’s a catastrophe,” Gaedi said. [1.9]“God gives, God takes.” [1.10] Chris Rolling, an American missionary and aid worker, tried to extricate a girl named Jacqueline from a collapsed school using nothing more than a hammer. [1.11] He urged her to be calm and pray, and as night fell he promised that he would return with help. [1.12] When he came back the next morning, Jacqueline was dead. [1.13]“The bodies stopped bothering me after a while, but I think what I will always carry with me is the conversation I had with Jacqueline before I left her,” Rolling wrote afterward on his blog. [1.14]“How could I leave someone who was dying, trapped in a building! ...[1.15] She seemed so brave when I left! [1.16] I told her I was going to get help, but I didn’t tell her I would be gone until morning. [1.17] I think this is going to trouble me for a long time.” [1.18] Dozens of readers wrote to comfort Rolling with the view that his story was evidence of divine wisdom and mercy.

© Manfred Stede / NAACL Tutorial 2013

## Referential Chains

[1.7] A man named Lionel Gaedi went to the Port-au-Prince morgue in search of his brother, Josef, but was unable to find his body among the piles of corpses that had been left there. [1.8]“I don’t see him—it’s a catastrophe,” Gaedi said. [1.9]“God gives, God takes.” [1.10] Chris Rolling, an American missionary and aid worker, tried to extricate a girl named Jacqueline from a collapsed school using nothing more than a hammer. [1.11] He urged her to be calm and pray, and as night fell he promised that he would return with help. [1.12] When he came back the next morning, Jacqueline was dead. [1.13]“The bodies stopped bothering me after a while, but I think what I will always carry with me is the conversation I had with Jacqueline before I left her,” Rolling wrote afterward on his blog. [1.14]“How could I leave someone who was dying, trapped in a building! ...[1.15] She seemed so brave when I left! [1.16] I told her I was going to get help, but I didn’t tell her I would be gone until morning. [1.17] I think this is going to trouble me for a long time.” [1.18] Dozens of readers wrote to comfort Rolling with the view that his story was evidence of divine wisdom and mercy.

© Manfred Stede / NAACL Tutorial 2013



## Two tasks

- *Anaphor* = referring expression that cannot be resolved without finding its antecedent in the preceding context
- **Anaphora resolution:**
  - find an antecedent for each anaphor in a text.
  - ‘Anaphora’ is an irreflexive, non-symmetrical relation.
- **Co-/reference resolution:**
  - partition the set of mentions of discourse referents in a text into classes (chains).
  - Since referents are identical, ‘coreference’ is an equivalence relation (reflexive, symmetrical, transitive)

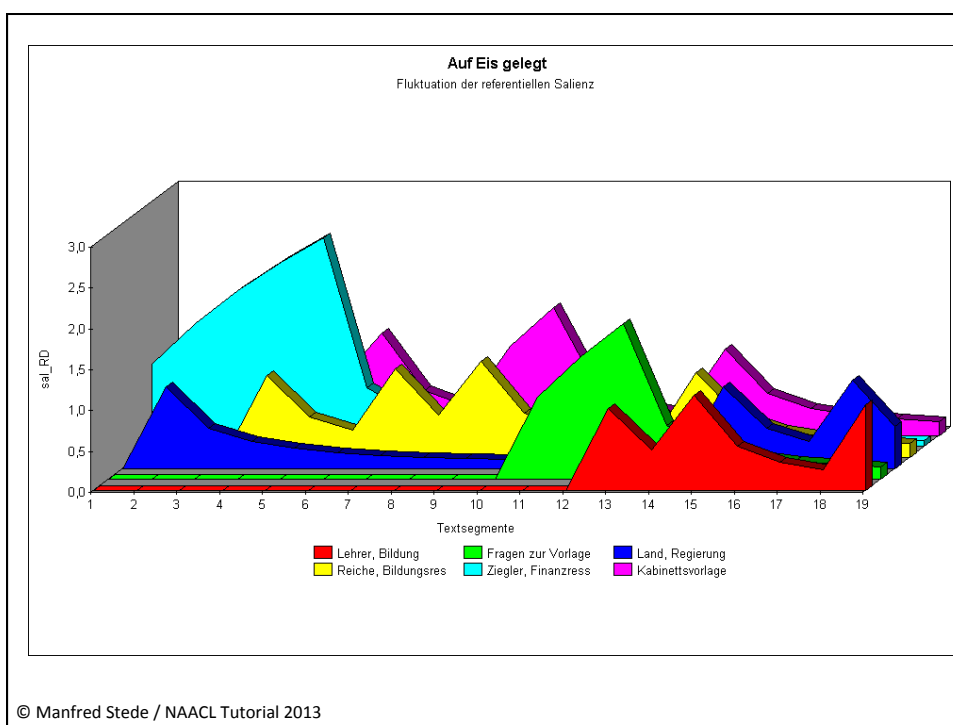
© Manfred Stede / NAACL Tutorial 2013

- Co-referent and anaphoric:  
*Jim forgot his umbrella. He had to return to his house.*
- Co-referent but not anaphoric:  
*Apple, Inc. announced a record surplus for the third quarter of the year. Overall, it has been a very successful year for Apple.*
- Not co-referent but anaphoric:  
*Potsdam University enrolled 1500 new students this year. Therefore, the president is very enthusiastic.*

© Manfred Stede / NAACL Tutorial 2013

## AUF EIS GELEGT

Dagmar Ziegler sitzt in der Schuldenfalle. Auf Grund der dramatischen Kassenlage in Brandenburg hat sie jetzt eine seit mehr als einem Jahr erarbeitete Kabinettsvorlage überraschend auf Eis gelegt und vorgeschlagen, erst 2003 darüber zu entscheiden. Überraschend, weil das Finanz- und das Bildungsressort das Lehrpersonalkonzept gemeinsam entwickelt hatten. Der Rückzieher der Finanzministerin ist aber verständlich. Es dürfte derzeit schwer zu vermitteln sein, weshalb ein Ressort pauschal von künftigen Einsparungen ausgenommen werden soll - auf Kosten der anderen. Reiches Ministerkollegen werden mit Argusaugen darüber wachen, dass das Konzept wasserdicht ist. Tatsächlich gibt es noch etliche offene Fragen. So ist etwa unklar, wer Abfindungen erhalten soll, oder was passiert, wenn zu wenig Lehrer die Angebote des vorzeitigen Ausstiegs nutzen. Dennoch gibt es zu Reiches Personalpapier eigentlich keine Alternative. Das Land hat künftig zu wenig Arbeit für zu viele Pädagogen. Und die Zeit drängt. Der große Einbruch der Schülerzahlen an den weiterführenden Schulen beginnt bereits im Herbst 2003. Die Regierung muss sich entscheiden, und zwar schnell. Entweder sparen um jeden Preis - oder Priorität für die Bildung.



## Hierarchy of Definiteness

(Gundel et al. 93)

In focus	“it”
Activated	“this”/“that”
Familiar	“this N”, “that N”
Identifiable	“the N”
Referential	indefinite “this N”
Identifiable	“a N”

© Manfred Stede / NAACL Tutorial 2013

## Salience factors in the RAP algorithm

(Lappin/Leass 94)

Sentence recency	100
Subject emphasis	80
Existential emphasis	70
Accusative emphasis	50
Indirect object and oblique complement emphasis	40
Non-adverbial emphasis	50
Head noun emphasis	80

© Manfred Stede / NAACL Tutorial 2013

## Sample run: RAP

*Sue found a plastic unicorn in the garden. She handed it to Jill. She liked it very much.*

Step	Referent	Referring Expressions	Value
(1)	Sue	{ <i>Sue</i> }	310
	unicorn	{ <i>a plastic unicorn</i> }	280
	garden	{ <i>the garden</i> }	230
(2)	Sue	{ <i>Sue</i> }	155
	unicorn	{ <i>a plastic unicorn</i> }	140
	garden	{ <i>the garden</i> }	115
(3)	Sue	{ <i>Sue, she</i> }	155+310=465
	unicorn	{ <i>a plastic unicorn, it</i> }	140+280=420
	garden	{ <i>the garden</i> }	115
	Jill	{ <i>Jill</i> }	270
(4)	Sue	{ <i>Sue, she, she</i> }	232.5
	unicorn	{ <i>a plastic unicorn, it, it</i> }	210
	garden	{ <i>the garden</i> }	57.5
	Jill	{ <i>Jill</i> }	135

© Manfred Stede / NAACL Tutorial 2013

## Corpus annotation

- Starting with MUC in the 1990s, several text corpora have been annotated with coreference information and used in competitions
  - MUC: message understanding conference
  - DUC: document understanding conference
  - TREC: text retrieval evaluation conference
  - ...

© Manfred Stede / NAACL Tutorial 2013

## Beware of the task definition

- **Competition scenario**  
Work with standard data sets (MUC, DUC, TREC, ...):  
markables already identified – need to decide on coreference only
- **Real-world scenario**  
Work with authentic raw text, where mentions have to be identified first
- E.g.: Soon et al (01) found that with tagging, chunking, NER, 85% of mentions in the MUC-7 corpus are being detected

© Manfred Stede / NAACL Tutorial 2013

## Mention-pair Models

- Stage 1: Pairwise classification

© Manfred Stede / NAACL Tutorial 2013

## Mention-pair Models: Features

(Soon et al. 01)

- DIST
- SEMCLASS
- NUMBER
- GENDER
- PROPER-NAME
- ALIAS
- ANA-PRONOUN
- DEF-NP
- DEM-NP
- STR-MATCH
- APPOSITIVE
- ANTE-PRONOUN

© Manfred Stede / NAACL Tutorial 2013

## Mention-pair Models

- Stage 1: Pairwise classification
  - NP1 = NP2
  - NP2 = NP3
  - NP1  $\neq$  NP3

© Manfred Stede / NAACL Tutorial 2013

## Mention-pair Models

- Stage 1: Pairwise classification
  - NP1 = NP2
  - NP2 = NP3
  - NP1  $\neq$  NP3
- Stage 2: Clustering
  - Local: consider a small number of pairings
    - Soon et al. 01: „closest first“
  - Global: consider as many pairings as possible

© Manfred Stede / NAACL Tutorial 2013

## Generating training instances

- Given an annotated corpus, what do we learn from?
  - all pairs of mentions => skewed distribution
  - mirror the human disambiguation task:
    - anaphor + antecedent
    - anaphor + all „intervening“ mentions

© Manfred Stede / NAACL Tutorial 2013

## Example

(Soon et al. 01)

**Example 3.6** (Ms. Washington)<sup>73</sup>'s candidacy is being championed by (several powerful lawmakers)<sup>74</sup> including ((her)<sup>76</sup> boss)<sup>75</sup>, (Chairman John Dingell)<sup>77</sup> (D., (Mich.)<sup>78</sup>) of (the House Energy and Commerce Committee)<sup>79</sup>. (She)<sup>80</sup> currently is (a counsel)<sup>81</sup> to (the committee)<sup>82</sup>.

ante	ana	feature vector	class. decision
(several powerful lawmakers) <sup>74</sup>	(her) <sup>76</sup>	0, 1, -, 2, -, -, +, -, -, -, -, -	no
(Ms. Washington) <sup>73</sup>	(her) <sup>76</sup>	0, 1, +, 1, -, -, +, -, -, -, -, -	yes
(the House Energy and C. Committee) <sup>79</sup>	(She) <sup>80</sup>	1, 0, +, 0, -, -, +, -, -, -, -, -	no
(Mich.) <sup>78</sup>	(She) <sup>80</sup>	2, 0, +, 0, -, -, +, -, -, -, -, -	no
(Chairman J.D.) <sup>77</sup>	(She) <sup>80</sup>	3, 1, +, 0, -, -, +, -, -, -, -, -	no
(her) <sup>76</sup>	(She) <sup>80</sup>	3, 1, +, 1, -, -, +, -, -, -, -, +	yes

© Manfred Stede / NAACL Tutorial 2013

## Mention-pair models

Many refinements, different classification techniques, etc., over the years

Fundamental problem, though:

*Jim Miller ... a unicorn ... Laura Smith ... she ...  
her brother ... Miller ... the beast ... he ...*

© Manfred Stede / NAACL Tutorial 2013



## Incremental entity-mention models

(here: Klenner/Tuggener 2011)

- **I**: the chronologically-ordered list of mentions
- **C**: set of coreference sets
- **B**: buffer for non-anaphoric mentions
- $m_i$ : current mention
- **+**: concatenate an item to a list

© Manfred Stede / NAACL Tutorial 2013

## Incremental entity-mention models

Klenner/Tuggener 2011

```

1  for i=1 to length(I)
2    for j=1 to length(C)
3       $r_j :=$  virtual prototype of coreference set  $C_j$ 
4       $Cand := Cand \oplus r_j$  if compatible( $r_j, m_i$ )
5    for k= length(B) to 1
6       $b_k :=$  the k-th licensed buffer element
7       $Cand := Cand \oplus b_k$  if compatible( $b_k, m_i$ )
8    if  $Cand = \{\}$  then  $B := B \oplus m_i$ 
9    if  $Cand \neq \{\}$  then
10      $ante_i :=$  most salient element of  $Cand$ 
11      $C :=$  augment( $C, ante_i, m_i$ )

```

© Manfred Stede / NAACL Tutorial 2013

## One complication: Indirect anaphora / Bridging (Clark 1977)

- **Necessary parts:**  
*I entered the room. The ceiling was high.*
- **Probable parts:**  
*I entered the room. The windows looked out to the bay.*
- **Inducible parts:**  
*I entered the room. The chandeliers sparkled brightly.*
- **Necessary roles:**  
*I went shopping. The time I started was 3pm.*
- **Optional roles:**  
*John was murdered. The knife lay nearby.*
- **Relations like reason, cause, consequence:**  
*An earthquake (...). The suffering people are going through (...)*

© Manfred Stede / NAACL Tutorial 2013

## Some more complications...

- **Non-nominal antecedents**  
*Jim met his mother last week. His father didn't like that.*
- **Reference to sets of objects**  
*Today I saw two roses and later on three tulips. They all were of the same red.*
- **Generic readings**  
*Today I ran into a bunch of squirrels. They are really wonderful animals.*
- **Expletive "it"**  
*It was raining the whole day long.*
- **Cataphora**  
*Before it ran away, the unicorn looked at me sadly.*
- **Ellipsis / one-anaphora**  
*Mike likes black cats, while Paul prefers brown ones.*

© Manfred Stede / NAACL Tutorial 2013

## Tutorial overview

- Part 1: „Large“ discourse units
  - Genre-specific structure
  - Topics
- Part 2: Coreference
- Part 3: „Small“ discourse units
  - Local coherence
  - Coherence-relational text structure
- Exploring interconnections

© Manfred Stede / NAACL Tutorial 2013

- John took a train from Paris to Istanbul. He has family there.
- John took a train from Paris to Istanbul. He likes spinach. (Hobbs 79)
- John took a train from Paris to Istanbul. Turkey has become a popular tourist destination.

© Manfred Stede / NAACL Tutorial 2013

## Coherence relations

- John took a train from Paris to Istanbul. He has family there.
- John took a train from Paris to Istanbul. He likes spinach.
- John took a train from Paris to Istanbul. Turkey has become a popular tourist destination.
- John took a train from Paris to Istanbul, because he has family there.

© Manfred Stede / NAACL Tutorial 2013

## Coherence relations

- John took a train from Paris to Istanbul, but he never arrived.
- Although John took a fast train from Paris to Istanbul, he arrived late.
- ...

© Manfred Stede / NAACL Tutorial 2013

## Coherence relations: levels of description

- If you're thirsty, make sure that you find some water.

SEMANTICS

- If you're thirsty, there's a beer in the fridge.

PRAGMATICS

© Manfred Stede / NAACL Tutorial 2013

## Definition

- **Coherence relation**

A relationship between adjacent units of text, holding on a semantic (proposition) or pragmatic (speech act) level of analysis.

Common groupings of relations are

- causality
- similarity/contrast
- contiguity

© Manfred Stede / NAACL Tutorial 2013

## Issues

- Connectives
- Minimal units of the analysis
- Inventory and definitions of coherence relations
- Corpus: PDTB
- Automatic local coherence analysis
  
- Coherence relations and discourse structure
- Corpora: RST-DT, Discourse Graph Bank
- RST-parsing

© Manfred Stede / NAACL Tutorial 2013

## Connectives

if you use a character reference such as `&#60;`; to insert the `<` character, the formatter will output `&lt;`.

© Manfred Stede / NAACL Tutorial 2013

if you use a character reference such as `&#60;`; to insert the `<` character, the formatter will output `&lt;`;

Coherence relation:

Purpose (you use `&#60` , you insert `<`)

© Manfred Stede / NAACL Tutorial 2013

if you use a character reference such as `&#60;`; to insert the `<` character, the formatter will output `&lt;`;

Coherence relation:

Condition ( you use `&#60` to insert `<` ,  
formatter will output `&lt;` )

© Manfred Stede / NAACL Tutorial 2013

**Because** well-formed XML does not permit raw less-than signs and ampersands, **if** you use a character reference such as `&#60;`; **or** the entity reference `&lt;`; **to** insert the `<` character, the formatter will output `&lt;`; **or** perhaps `&#60;`.

© Manfred Stede / NAACL Tutorial 2013

nonetheless                      notwithstanding that                      all the same  
 though                      despite the fact that                      although  
                     rather                      but                      however                      or else  
 on the other hand                      on the contrary  
                     only                      even so                      in spite of that  
                     in contrast                      else                      despite this  
 yet                      instead                      failing that                      on the other side  
 then again                      unless                      whereas                      admittedly  
                     nevertheless                      apart from that  
                     whilst  
                     by contrast                      alternatively                      otherwise  
 while                      meanwhile                      anyway  
                     still                      even though

(from Knott 96)



## Definition

- **Connective**

A closed-class, non-inflectable word or word group that semantically denotes a two-place relation, where the entities being related can in text be expressed as clauses.

Syntactically, a connective can be

- a subordinating conjunction,
- a coordinating conjunction,
- an adverbial,
- (arguably) a preposition.

© Manfred Stede / NAACL Tutorial 2013

## Connectives: Ambiguity

- Connective or no connective
  - *Since you like ice cream so much, I'll buy one for you.*
  - *Since 1988 I have never had any ice cream.*
- Different coherence relations
  - *Since you like ice cream so much, I'll buy one for you.*
  - *I haven't had any ice cream since you arrived in New York.*

© Manfred Stede / NAACL Tutorial 2013

## Issues

- Connectives
- Minimal units of the analysis
- Inventory and definitions of coherence relations
- Corpus: PDTB
- Automatic local coherence analysis
  
- Coherence relations and discourse structure
- Corpora: RST-DT, Discourse Graph Bank
- RST parsing

© Manfred Stede / NAACL Tutorial 2013

## Minimal units

- Sentences!
  - But:
    - incomplete material, ellipsis  
*Bring me a hammer please. A big one.*
    - complex sentences  
*Bring me a hammer when you pass by the workshop.*
  
- Clauses!
  - But:
    - deal with nominalizations  
*John attended the lecture despite his illness.*
    - deal with embedding  
*John, although he was ill, attended the lecture.*
    - deal with non-/restrictive relative clauses  
*The red car that's parking in front of you belongs to me.*  
*The red car, which was brought here by my Dad, belongs to me.*

© Manfred Stede / NAACL Tutorial 2013

## Minimal units

- **Non-structural definitions**, e.g. Polanyi et al. (04):  
1) „the syntactic constructions that encode a minimum unit of meaning and/or discourse function interpretable relative to a set of contexts.“

**Table 4.1:** Examples of minimal discourse units and their types [Polanyi et al., 2004].

Segment types	Realizations	Examples
Eventualities	clauses	[I heard the dog][that was barking.]
	predication	[California elected Schwarzenegger] [governor.]
	infinitival modif.	[They left] [to get the tickets.]
Interpolations	parentheticals	[The show [(and what a show it was)] lasted 4 hours.]
Fragments	section headings	[4. Discussion]
	list items	[e.g., [hydrogen,] [helium]]
Conj. operators	conjunction	[We arrived] [and] [got seats.]
Discourse operators	“scene setting” modifier	[On Tuesday,] [we will see the sites.]

- 2) units that can be „independently continued“  
(removes many small units from the candidate set)

© Manfred Stede / NAACL Tutorial 2013

## Definition

- **Elementary discourse unit (EDU)**

A span of text, usually a clause, but in general ranging from minimally a (nominalization) NP to maximally a sentence. It denotes a single event or type of event, serving as a complete, distinct unit of information that the surrounding discourse may connect to. An EDU may be structurally embedded in another.

© Manfred Stede / NAACL Tutorial 2013

## Automatic segmentation

- Example: **SLSeg** (Tofiloski et al. 2009)
- Human agreement: kappa 0.85
- Based on syntactic rules, e.g.
  - distinguish sentential complements that should/not be EDUs
  - copy heads of non-restrictive relative clauses

*>The aftermath of the 2008 cyclone in Burma not only betrayed the callous indifference of the ruling junta  
>but demonstrated the vibrancy of civil society there.  
>Haiti's earthquake shows that, whatever the communal spirit of its people at the moment of crisis,  
>the government was not functioning, unable even to bury the dead, much less rescue the living.  
>This vacuum, which had been temporarily filled by the U.N.,  
>This vacuum now poses the threat of chaos.*

## Issues

- Connectives
- Minimal units of the analysis
- **Inventory and definitions of coherence relations**
- Corpus: PDTB
- Automatic local coherence analysis
  
- Coherence relations and discourse structure
- Corpora: RST-DT, Discourse Graph Bank
- RST parsing

## Inventory and definitions of relations

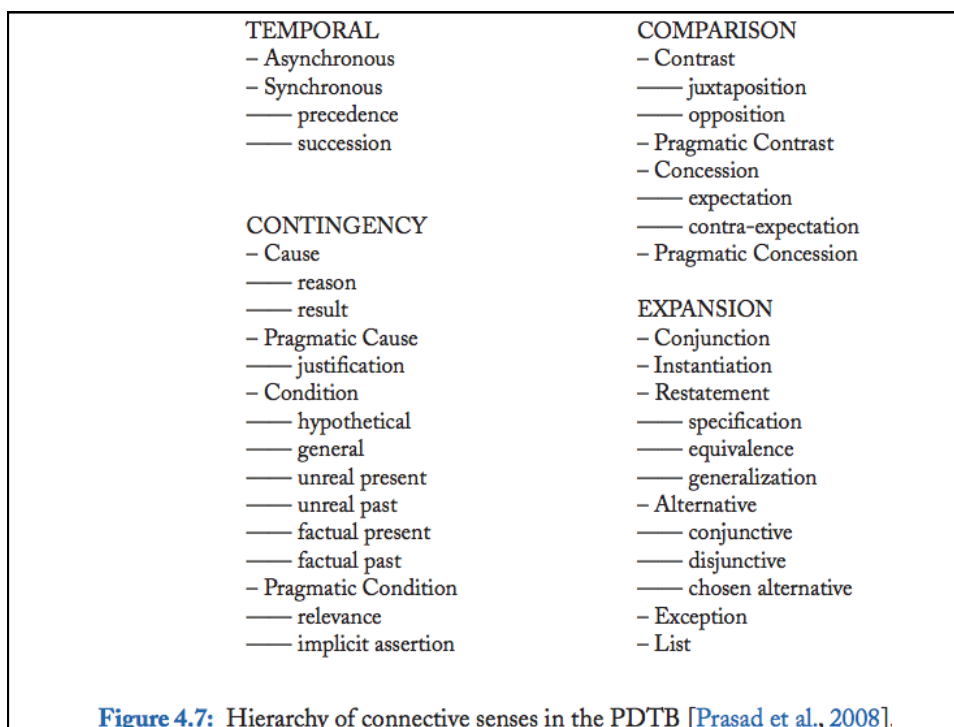
- Many proposals in the (discourse, semantics, pragmatics) literature
- Three influential ones:
  - Segmented Discourse Representation Theory (Asher/Lascarides 03)
  - Rhetorical Structure Theory (Mann/Thompson 88)
  - Connective senses in Penn Discourse Treebank (Prasad et al. 08)

© Manfred Stede / NAACL Tutorial 2013

## Relation definitions

- **Evidence** in RST:
  - *Constraint on Nucleus:*  
Reader might not believe N to a degree satisfactory to Writer
  - *Constraint on Satellite:*  
Reader believes S or will find it credible
  - *Constraints on Nucleus+Satellite:*  
Reader's comprehending S increases her/his belief of N
  - *Intention of Writer:*  
Reader's belief of N is increased

© Manfred Stede / NAACL Tutorial 2013



## Corpus: PDTB (Prasad et al. 08)

- Wall Street Journal portion of Penn Treebank
- 100 connectives (types)
- 18.000 instances
- Implicit connectives: Within paragraphs, check whether a connective *could* be added between segments

**Example 4.26** [Drug makers shouldn't be able to duck liability]<sub>Arg1</sub> [because]<sub>Conn</sub> [people couldn't identify precisely which identical drug was used.]<sub>Arg2</sub>

**Example 4.27** [France's second-largest government-owned insurance company, Assurances Generales de France, has been building its own Navigation Mixte stake]<sub>Arg1</sub> currently thought to be between 8% and 10%. Analysts said [they don't think it is contemplating a takeover]<sub>Arg2</sub>, [however]<sub>Conn</sub>, and its officials couldn't be reached.

- Agreement:
  - identify both arguments of explicit connective: 90.2%
  - identify both arguments of implicit connective: 85.1%
  - identify the relation at top-level (4 groups): 94%

## Automatic local coherence analysis

- E.g., as an extension of opinion mining:

*(...) We didn't like the village very much and won't come back because there are only few things to see, and moreover the place is quite dirty. (...)*

© Manfred Stede / NAACL Tutorial 2013

## Automatic local coherence analysis

- **1** - Disambiguate potential connectives
  - Marcu 00: 1200 of 2100 potential cue phrases *are* cue phrases
  - Pitler/Nenkova 09
    - Of 100 candidates in PDTB, only 11 appear as connectives more than 90% of the time
    - MaxEnt classifier using syntactic features (from manual annotation): f-measure 92.3%
      - POS tag, phrase label of candidate
      - categories of parent and left sibling
      - right sibling: category, presence of VP or trace in subtree

© Manfred Stede / NAACL Tutorial 2013

## Automatic local coherence analysis

- 2 – Sense disambiguation (PDTB-style) for connectives
- Miltsakaki et al. 05
  - *since* (temporal/causal): 89.5% acc.
  - *while* (temporal/opposition/concessive): 71.9%
  - *when* (temporal/conditional): 82.6%
  - Features: tense form of aux. *have* and *be*, tense form of head, presence of modals and temp. expr.
- Pitler/Nenkova 09
  - due to skewed sense distribution, majority class leads to 93.7% acc. already; with syntactic features: 94.2%

© Manfred Stede / NAACL Tutorial 2013

## Automatic local coherence analysis

- 3 – Scope identification
- **Preposition:** Arg2 is the NP following the preposition, Arg1 is the governing clause.  
*Despite [her bad mood.]<sub>Arg2</sub> [Susan decided to go to the party.]<sub>Arg1</sub>*
- **Subordinating conjunction:** Arg2 is the clause following the conjunction, Arg1 is the matrix clause.  
*Because [Susan was in a bad mood.]<sub>Arg2</sub> [she did not go to the party.]<sub>Arg1</sub>*
- **Coordinating conjunction:** Arg2 is the sentence following the conjunction, Arg1 is some text segment preceding it.  
*[Susan was in a bad mood.]<sub>Arg1</sub> but [she decided to go to the party.]<sub>Arg2</sub>*
- **Adverbial:** Arg2 is the clause containing the adverbial; Arg1 is some text segment preceding it.  
*[Susan was in a bad mood.]<sub>Arg1</sub> [She nevertheless decided to go to the party.]<sub>Arg2</sub>*

© Manfred Stede / NAACL Tutorial 2013



## Automatic local coherence analysis

- 3 – Scope identification
- Elwell/Baldrige 08
  - Arg2: acc. 92-94%
  - Arg1: acc. 78-82%
- [Drug makers shouldn't be able [to duck liability]]<sub>Arg1?</sub> **because** [people couldn't identify precisely which identical drug was used.]<sub>Arg2</sub>
- Prasad et al. 08: Arg1 distribution in PDTB
  - 65% in same sentence
  - 30% in immediately-preceding sentence
  - 9% in non-adjacent sentence

© Manfred Stede / NAACL Tutorial 2013

## Automatic local coherence analysis

- 4 – Finding implicit relations
- Proportion of signalled relations
  - roughly 40% (according to several studies)
  - differs considerably between relations!  
E.g., CONCESSION versus BACKGROUND

© Manfred Stede / NAACL Tutorial 2013

## Automatic local coherence analysis

- 4 – Finding implicit relations
- Sporleder/Lascarides 05
  - 5-way classification: CONTRAST, EXPLANATION, RESULT, SUMMARY, CONTINUATION
  - Approach: Remove connective, then classify
  - Features: position, length of arguments, words, parts of speech, WordNet classes, tens/aspect and various syntactic f.s
  - Acc.: 57.6%
  - (Assumption made: contexts of explicit and implicit instances are similar to each other – but are they? See Sporleder/Lascarides 08)

© Manfred Stede / NAACL Tutorial 2013

## A final version of the train example...

- John took a train from Paris to Istanbul, departing from Montparnasse at noon. He has family in Istanbul.

© Manfred Stede / NAACL Tutorial 2013

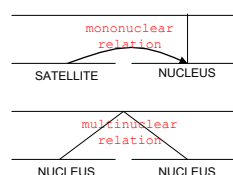
## Issues

- Connectives
- Minimal units of the analysis
- Inventory and definitions of coherence relations
- Corpus: PDTB
- Automatic local coherence analysis
- Coherence relations and discourse structure
- Corpora: RST-DT, Discourse Graph Bank
- RST parsing

© Manfred Stede / NAACL Tutorial 2013

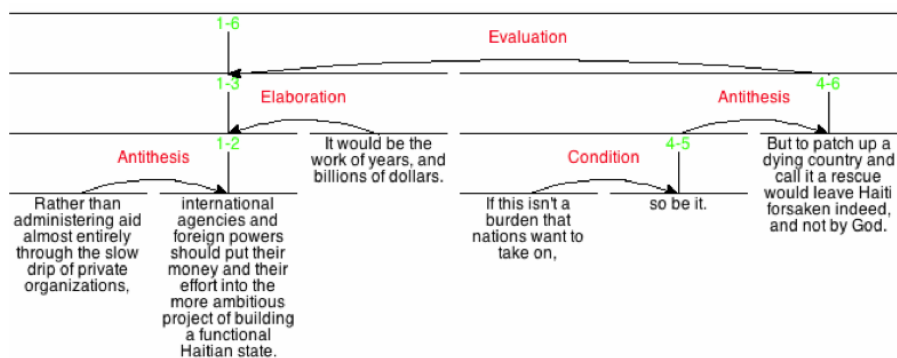
## Rhetorical Structure Theory (RST; Mann/Thompson 88)

- Coherence relations between discourse segments
  - asymmetric (“mononuclear”)
    - one **nucleus**, one **satellite**
  - symmetric (“multinuclear”)
    - multiple **nuclei**
- Resulting structure is a complete tree
  - No gaps
  - No cross-dependencies
- 25 relations  
Cause, Contrast, Elaboration, ...



© Manfred Stede / NAACL Tutorial 2013

## Rhetorical Structure Theory: Example



© Manfred Stede / NAACL Tutorial 2013

## RST Discourse Treebank (Carlson et al. 03)

- 385 Wall Street Journal articles (22k EDUs)
  - 53 mononuclear, 25 multinuclear relations
  - Relations grouped into 16 categories
  - Detailed annotation guidelines; kappa for the 16 categories up to .82
- Attribution
  - Background
  - Cause
  - Comparison
  - Condition
  - Contrast
  - Elaboration
  - Enablement
  - Evaluation
  - Explanation
  - Joint
  - Manner-Means
  - Topic-Comment
  - Summary
  - Temporal
  - TopicChange

© Manfred Stede / NAACL Tutorial 2013

## Structure

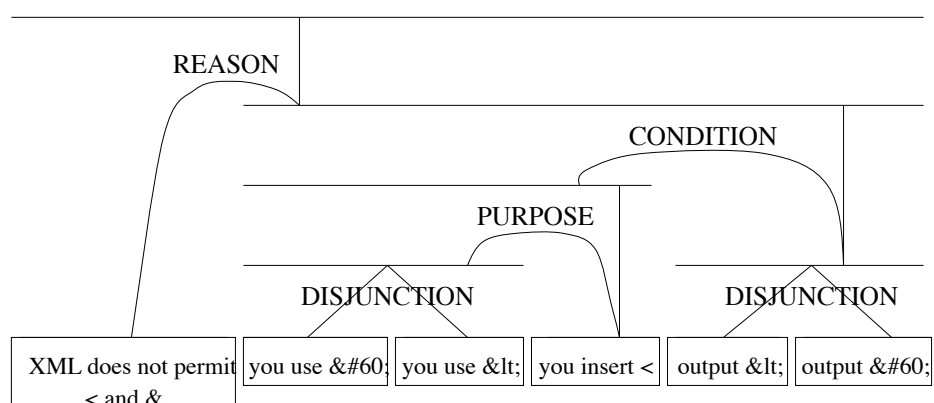
**Because** well-formed XML does not permit raw less-than signs and ampersands, **if** you use a character reference such as `&#60;`; **or** the entity reference `&lt;`; **to** insert the `<` character, the formatter will output `&lt;`; **or** perhaps `&#60;`.

**Because** A, **if** B **or** C **to** D, E **or** F.

(**Because** A, (**if** ((B **or** C) **to** D), (E **or** F))).

© Manfred Stede / NAACL Tutorial 2013

## „Rhetorical structure“



© Manfred Stede / NAACL Tutorial 2013

## Greedy RST parsing (Hernault et al. 10)

```

input:  $L \leftarrow \langle e_1, e_2, \dots, e_n \rangle$  (list of EDUs)
for all  $(l_i, l_{i+1})$  in  $L$  do
   $Scores[i] \leftarrow \text{STRUCT}(l_i, l_{i+1})$ 
end for
while  $|L| > 1$  do
   $i \leftarrow \text{max}(Scores)$ 
   $NewLabel \leftarrow \text{LABEL}(l_i, l_{i+1})$ 
   $NewSubTree \leftarrow \text{CreateTree}(l_i, l_{i+1}, NewLabel)$ 
   $Scores[i - 1] \leftarrow \text{STRUCT}(l_{i-1}, NewSubTree)$ 
   $Scores[i + 2] \leftarrow \text{STRUCT}(NewSubTree, l_{i+2})$ 
   $\text{delete}(Scores[i])$ 
   $\text{delete}(Scores[i + 1])$ 
   $L \leftarrow [l_0, \dots, l_{i-1}, NewSubTree, l_{i+2}, \dots]$ 
end while
return  $l_0$ 

```

© Manfred Stede / NAACL Tutorial 2013

## Greedy RST parsing (Hernault et al. 10)

- SVMs trained on RST-DT
- **STRUCT**: binary classifier scoring whether *any* relation holds between a pair of segments
  - features: sentence and para boundaries, segment size, position in text
- **LABEL**: multi-class classifier assigning relation and nuclearity assignment
  - features similar to those discussed earlier
  - plus relations that have already been assigned

© Manfred Stede / NAACL Tutorial 2013

## Trees or Graphs?

- Webber et al. 99: **structural** versus **anahoric** connectives

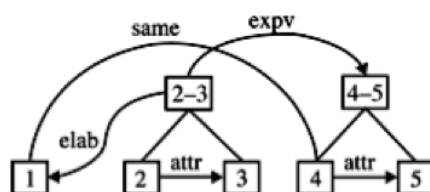
(a) *On the one hand, John loves Barolo.* (b) *So he ordered three cases of the '97.* (c) *On the other hand, because he's broke,* (d) *he then had to cancel the order.*

- See also: Egg/Redeker (10)

© Manfred Stede / NAACL Tutorial 2013

## Discourse Graph Bank (Wolf/Gibson 05)

- 135 texts: AP newswire, Wall Street Journal
- Distinguish directed from undirected relations
- No nuclearity
- Example: (1) *Mr. Baker's assistant for inter-American affairs, Bernard Aronson,* (2) *while maintaining* (3) *that the Sandinistas had also broken the cease-fire,* (4) *acknowledged:* (5) *"It's never very clear who starts what."*



© Manfred Stede / NAACL Tutorial 2013

## EDUs and relations: **where are we?**

- RST parsing: technically interesting (RST-DT), but open issues w.r.t. theoretical status
  - what relations?
  - intentional or informational analysis?
  - what is the right unit for RST tree? paragraph?
- PDTB: much less commitment on overall discourse structure – very useful for local coherence analysis

© Manfred Stede / NAACL Tutorial 2013

## Tutorial overview

- Part 1: „Large“ discourse units
  - Genre-specific structure
  - Topics
- Part 2: Coreference
- Part 3: „Small“ discourse units
  - Local coherence
  - Coherence-relational text structure
- Exploring interconnections

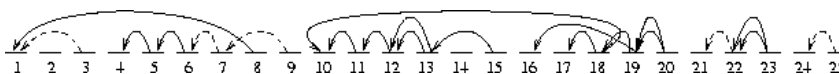
© Manfred Stede / NAACL Tutorial 2013



### Mandatory vaccination against children's diseases?

[1] Today, children don't know anymore what **pox** are. [2] What a joy. [3] When **pox vaccination** was introduced in 1854, [4] quite a few people believed [5] that their head would turn into a cow's head [6] if they got themselves vaccinated. [7] For the vaccine was made from cattle's skin at the times. [8] Nowadays **this dreadful disease** is exterminated. [9] Thanks to a determined, world-wide vaccination campaign. [10] But there still are other diseases: Measles, polio, diphtheria, mumps, rubella, hepatitis B, tuberculosis, pertussis. [11] Millions of children die of these, especially in less developed countries. [12] In Germany, many parents apparently don't take these diseases seriously. [13] Because they don't know them anymore! [14] For it has been achieved with vaccines [15] that these infections hit only rarely today. [16] But those who have experienced [17] how terribly children suffer [18] when they come down with ,just' measles or pertussis, [19] should spare them the agony. [20] As well as the long-term consequences. [21] Only those who have their children vaccinated will contribute to vaccines' becoming superfluous some day. [22] Instead, people rant about side effects [23] that occur very rarely and are known merely from books. [24] Then there is the great argument: This is my child, the government must not prick her. [25] No vaccine can help against such parents.

## Referential Structure



### Mandatory vaccination against children's diseases?

[1] Today, children don't know anymore what pox are. [2] What a joy. [3] When pox vaccination was introduced in 1854, [4] quite a few people believed [5] that their head would turn into a cow's head [6] if they got themselves vaccinated. [7] For the vaccine was made from cattle's skin at the times. [8] Nowadays this dreadful disease is exterminated. [9] Thanks to a determined, world-wide vaccination campaign. [10] But there still are other diseases: Measles, polio, diphteria, mumps, rubella, hepatitis B, tuberculosis, pertussis. [11] Millions of children die of these, especially in less developed countries. [12] In Germany, many parents apparently don't take these diseases seriously. [13] Because they don't know them anymore! [14] For it has been achieved with vaccines [15] that these infections hit only rarely today. [16] But those who have experienced [17] how terribly children suffer [18] when they come down with ,just' measles or pertussis, [19] should spare them the agony. [20] As well as the long-term consequences. [21] Only those who have their children vaccinated will contribute to vaccines' becoming superfluous some day. [22] Instead, people rant about side effects [23] that occur very rarely and are known merely from books. [24] Then there is the great argument: This is my child, the government must not prick her. [25] No vaccine can help against such parents.

## Topic Structure

Pox and pox vaccination									Other diseases and vaccination						Measles and pertussis				Pro Vaccination					
Side effects									Side effects						Side effects				"my child"					
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25

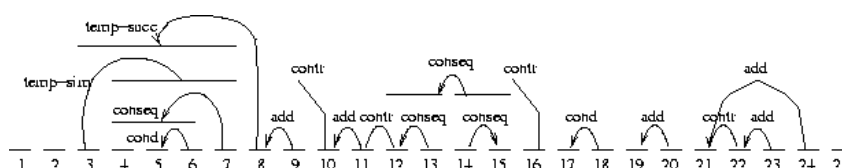
### Mandatory vaccination against children's diseases?

[1] Today, children don't know anymore what pox are. [2] What a joy. [3] **When** pox vaccination was introduced in 1854, [4] quite a few people believed [5] that their head would turn into a cow's head [6] **if** they got themselves vaccinated. [7] **For** the vaccine was made from cattle's skin at the times. [8] Nowadays this dreadful disease is exterminated. [9] **Thanks to** a determined, world-wide vaccination campaign. [10] **But** there still are other diseases: Measles, polio, diphtheria, mumps, rubella, hepatitis B, tuberculosis, pertussis. [11] Millions of children die of these, especially in less developed countries. [12] In Germany, many parents apparently don't take these diseases seriously. [13] **Because** they don't know them anymore! [14] **For** it has been achieved with vaccines [15] that these infections hit only rarely today. [16] **But** those who have experienced [17] how terribly children suffer [18] **when** they come down with ‚just‘ measles or pertussis, [19] should spare them the agony. [20] **As well as** the long-term consequences. [21] Only those who have their children vaccinated will contribute to vaccines' becoming superfluous some day. [22] **Instead**, people rant about side effects [23] that occur very rarely and are known merely from books. [24] **Then** there is the great argument: This is my child, the government must not prick her. [25] No vaccine can help against such parents.

## Conjunctive Relations

- temporal
  - simultaneous, succession
- consequential
  - manner, consequence, condition, purpose, concession
- comparative
  - similarity, contrast, reformulation
- additive
  - addition, alternation

## Conjunctive Relations

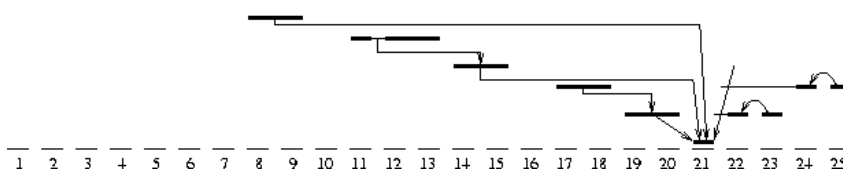


© Manfred Stede / NAACL Tutorial 2013

### Mandatory vaccination against children's diseases?

[1] Today, children don't know anymore what pox are. [2] What a joy. [3] When pox vaccination was introduced in 1854, [4] quite a few people believed [5] that their head would turn into a cow's head [6] if they got themselves vaccinated. [7] For the vaccine was made from cattle's skin at the times. [8] Nowadays this dreadful disease is exterminated. [9] Thanks to a determined, world-wide vaccination campaign. [10] But there still are other diseases: Measles, polio, diphteria, mumps, rubella, hepatitis B, tuberculosis, pertussis. [11] Millions of children die of these, especially in less developed countries. [12] In Germany, many parents apparently don't take these diseases seriously. [13] Because they don't know them anymore! [14] For it has been achieved with vaccines [15] that these infections hit only rarely today. [16] But those who have experienced [17] how terribly children suffer [18] when they come down with ,just' measles or pertussis, [19] should spare them the agony. [20] As well as the long-term consequences. [21] Only those who have their children vaccinated will contribute to vaccines' becoming superfluous some day. [22] Instead, people rant about side effects [23] that occur very rarely and are known merely from books. [24] Then there is the great argument: This is my child, the government must not prick her. [25] No vaccine can help against such parents.

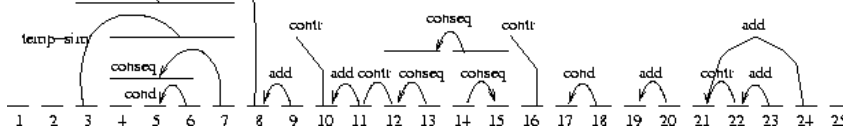
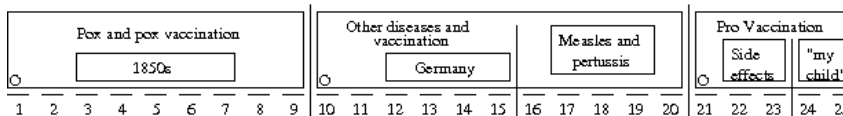
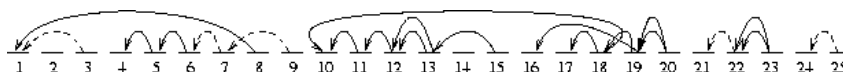
## Argument structure



© Manfred Stede / NAACL Tutorial 2013

(inspired by Freeman 1993)

## Text understanding: Relating levels of analysis



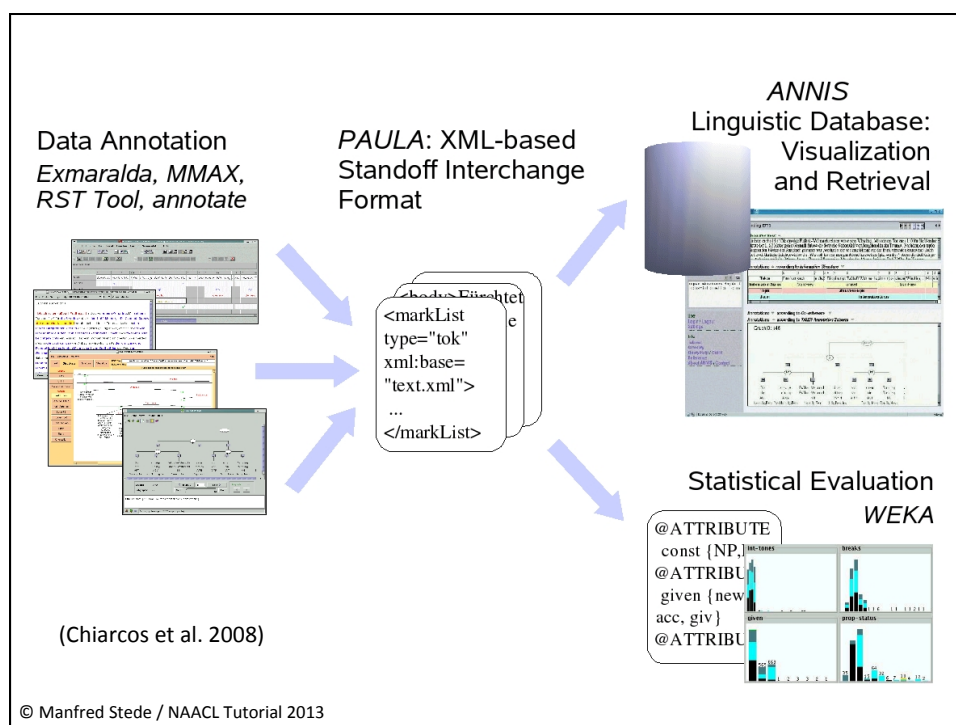
© Manfred Stede / NAACL Tutorial 2013

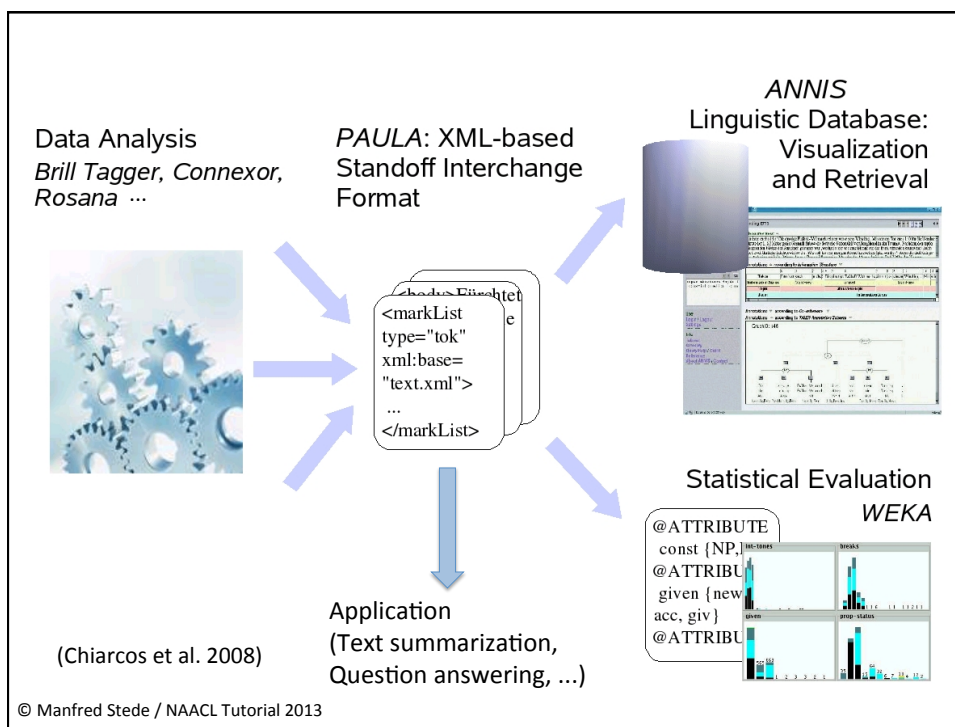
## Potsdam Commentary Corpus (PCC) (Stede 2004)

200 short editorials from two German newspapers

- Syntax
- Coreference
- RST
- Content zones
- Topics
- Information structure
- Conjunctive relations
- Illocutionary status
- Argumentation structure

© Manfred Stede / NAACL Tutorial 2013





Search Result - tok (10, 10)

Page 5 of 20    Token Annotations    Show Citation URL    Displaying Results 81 - 100 of 399

zusammengewürfelt sind sie müssen sich daran gewöhnen , dass sie nun in einer Mannschaft \* Döberitzer Heide \* spielen  
zusammengewürfelt sein sie müssen sich daran gewöhnen , dass sie nun in ein Mannschaft \* Döberitzer Heide \* spielen  
Psp 3.Pl.Pres.Ind -- 3.Nom.Pl.Masc 3.Pl.Pres.Ind 3.Acc.Pl -- Inf -- 3.Nom.Pl.Masc -- Dat.Sg.Fem Dat.Sg.Fem -- Pos.\*.\* Nom.Sg.Fem -- 3.Pl.Pres.Ind --  
VPPP VAFIN \$, PPER VMFIN PRF PROAV VVIN \$, KOUS PPER ADV APPR ART NN \$( ADJA NN \$( VVFIN

- dependencies (arches [SVG capable browsers])
- information structure (grid)
- discourse referents (grid)
- coreference (discourse)
- rhetorical structure (old\_grid)
- tree

zusammengewürfelt sind sie müssen sich daran gewöhnen , dass sie nun in einer Mannschaft \* Döberitzer Heide \* spielen

① sind sie müssen sich daran gewöhnen , dass sie nun in einer Mannschaft \* Döberitzer Heide \* spielen . Und  
3.Pl.Pres.Ind -- 3.Nom.Pl.Masc 3.Pl.Pres.Ind 3.Acc.Pl -- Inf -- 3.Nom.Pl.Masc -- Dat.Sg.Fem Dat.Sg.Fem -- Pos.\*.\* Nom.Sg.Fem -- 3.Pl.Pres.Ind -- und  
VAFIN \$, PPER VMFIN PRF PROAV VVIN \$, KOUS PPER ADV APPR ART NN \$( ADJA NN \$( VVFIN \$, KON

Search Result - tok (10, 10)

Page 5 of 20      Token Annotations      Show Citation URL      Displaying Results 81 - 100 of 399

**zusammengewürfelt** sind , sie müssen sich daran gewöhnen , dass sie nun in einer Mannschaft \* Döberitzer Heide \* spielen

zusammenwürfeln sein sie müssen sich daran gewöhnen , dass sie nun in ein Mannschaft \* Döberitzer Heide \* spielen

3.Pl.Pres.Ind -- 3.Nom.Pl.Masc 3.Pl.Pres.Ind 3.Acc.Pl -- Inf -- 3.Nom.Pl.Masc -- Dat.Sg.Fem Dat.Sg.Fem -- Pos.\*\* Nom.Sg.Fem -- 3.Pl.Pres.Pr

VVPP VAFIN \$, PPER VMFIN PRF PROAV VVINF \$, KOUS PPER ADV APPR ART NN \$( ADJA NN \$( VVFIN

dependencies (arches [SVG capable browsers])

information structure (grid)

discourse referents (grid)

coreference (discourse)

historical structure (old\_grid)

tree

sind , sie müssen sich daran gewöhnen , dass sie nun in einer Mannschaft \* Döberitzer Heide \* spielen . Und  
sein , sie müssen sich daran gewöhnen , dass sie nun in ein Mannschaft \* Döberitzer Heide \* spielen . und  
3.Pl.Pres.Ind -- 3.Nom.Pl.Masc 3.Pl.Pres.Ind 3.Acc.Pl -- Inf -- 3.Nom.Pl.Masc -- Dat.Sg.Fem Dat.Sg.Fem -- Pos.\*\* Nom.Sg.Fem -- 3.Pl.Pres.Ind --  
VAFIN \$, PPER VMFIN PRF PROAV VVINF \$, KOUS PPER ADV APPR ART NN \$( ADJA NN \$( VVFIN \$, KON

dependencies (arches [SVG capable browsers])

information structure (grid)

discourse referents (grid)

coreference (discourse)

historical structure (old\_grid)

tree

dass sie nun in einer Mannschaft \* Döberitzer Heide \* spielen . Und das heißt gemeinsam und nicht gegeneinander . Ermahnung  
dass sie nun in ein Mannschaft \* Döberitzer Heide \* spielen . und der heißen gemeinsam und nicht gegeneinander . Ermahnung  
3.Nom.Pl.Masc -- 3.Nom.Pl.Masc 3.Pl.Pres.Ind 3.Sg.Ne.3.Sg.Pres.Ind -- Pos.\*\* Nom.Sg.Fem -- 3.Pl.Pres.Ind -- Nom.Sg.Ne.3.Sg.Pres.Ind -- Pos --  
KOUS PPER ADV APPR ART NN \$( ADJA NN \$( VVFIN \$, KON PDS VVFIN ADJD KON PTKNEG ADV \$, NN

Search Result - tok (10, 10)

Page 5 of 20      Token Annotations      Show Citation URL      Displaying Results 81 - 100 of 399

**zusammengewürfelt** sind , sie müssen sich daran gewöhnen , dass sie nun in einer Mannschaft \* Döberitzer Heide \* spielen

zusammenwürfeln sein sie müssen sich daran gewöhnen , dass sie nun in ein Mannschaft \* Döberitzer Heide \* spielen

3.Pl.Pres.Ind -- 3.Nom.Pl.Masc 3.Pl.Pres.Ind 3.Acc.Pl -- Inf -- 3.Nom.Pl.Masc -- Dat.Sg.Fem Dat.Sg.Fem -- Pos.\*\* Nom.Sg.Fem -- 3.Pl.Pres.Pr

VVPP VAFIN \$, PPER VMFIN PRF PROAV VVINF \$, KOUS PPER ADV APPR ART NN \$( ADJA NN \$( VVFIN

dependencies (arches [SVG capable browsers])

information structure (grid)

Select Displayed Annotation Levels ▾

Focus_newint	nf-unsc																					
Inf-Stat																						
NP		gfv-active			acc-eggr				nf-unsc		gfv-active		acc-inf									
PP		NP			NP						NP		NP									
Sent	s																					
Topic		ab																				
tok	zusammengewürfelt	sind	,	sie	müssen	sich	daran	gewöhnen	,	dass	sie	nun	in	einer	Mannschaft	*	Döberitzer	Heide	*	spielen	.	Und

discourse referents (grid)

coreference (discourse)

historical structure (old\_grid)

tree

sind , sie müssen sich daran gewöhnen , dass sie nun in einer Mannschaft \* Döberitzer Heide \* spielen . Und  
sein , sie müssen sich daran gewöhnen , dass sie nun in ein Mannschaft \* Döberitzer Heide \* spielen . und  
3.Pl.Pres.Ind -- 3.Nom.Pl.Masc 3.Pl.Pres.Ind 3.Acc.Pl -- Inf -- 3.Nom.Pl.Masc -- Dat.Sg.Fem Dat.Sg.Fem -- Pos.\*\* Nom.Sg.Fem -- 3.Pl.Pres.Ind --  
VAFIN \$, PPER VMFIN PRF PROAV VVINF \$, KOUS PPER ADV APPR ART NN \$( ADJA NN \$( VVFIN \$, KON

dependencies (arches [SVG capable browsers])

information structure (grid)

discourse referents (grid)

coreference (discourse)

historical structure (old\_grid)

tree

dass sie nun in einer Mannschaft \* Döberitzer Heide \* spielen . Und das heißt gemeinsam und nicht gegeneinander . Ermahnung  
dass sie nun in ein Mannschaft \* Döberitzer Heide \* spielen . und der heißen gemeinsam und nicht gegeneinander . Ermahnung  
3.Nom.Pl.Masc -- 3.Nom.Pl.Masc 3.Pl.Pres.Ind 3.Sg.Ne.3.Sg.Pres.Ind -- Pos.\*\* Nom.Sg.Fem -- 3.Pl.Pres.Ind -- Nom.Sg.Ne.3.Sg.Pres.Ind -- Pos --  
KOUS PPER ADV APPR ART NN \$( ADJA NN \$( VVFIN \$, KON PDS VVFIN ADJD KON PTKNEG ADV \$, NN



Search Result - tok (10, 10) Displaying Results 81 - 100 of 399

Page 5 of 20 Token Annotations Show Citation URL

zusammengewürfelt sind , sie müssen sich daran gewöhnen , dass sie nun in einer Mannschaft \* Döberitzer Heide \* spielen  
 zusammenwürfeln sein , sie müssen sich daran gewöhnen , dass sie nun in ein Mannschaft \* Döberitzer Heide \* spielen  
 3.Pl.Pres.Ind -- 3.Nom.Pl.Masc 3.Pl.Pres.Ind 3.Acc.Pl -- Inf -- -- 3.Nom.Pl.Masc -- -- Dat.Sg.Fem Dat.Sg.Fem -- Pos.\*\*.\* Nom.Sg.Fem -- 3.Pl.Pres.Ind  
 VVPP VAFIN \$, PPER VMFIN PRF PROAV VVINF \$, KOUS PPER ADV APPR ART NN \$( ADJA NN \$( VVFIN

dependencies (arches [SVG capable browsers])  
 information structure (grid)  
 discourse referents (grid)  
 coreference (discourse)

Stollpass Wunder gibt es immer wieder ! Erst spielen die Dallgower Gemeindevertreter so statisch und verzagt wie die deutsche Abwehrreihe der Fußballkicker . Und dann kommt aus der Tiefe solch ein fulminanter Stollpass , von dem man hofft , dass die Seeburger oder Groß-Glienicker Mitspieler ihn aufnehmen können . Ein Befreiungsschlag ist es allerdings nicht , weil es vorerst keine Gefahr fürs Dallgower Tor gab . Die Seeburger und einige Groß-Glienicker haben den Ball erst zurückgespielt und dann um so drängender wieder gefordert . Nun sollen sie zeigen , wie sie die Chance verwerten . Eine Diskussion , wo künftig die Trainerkabine stehen soll , wäre in der jetzigen Spielsituation verheerend . Und eine Parallele zu den deutschen Grotten-Kickern gibt es immer noch . Auch wenn die Spieler aus den verschiedenen Vereinen zusammengewürfelt sind , sie müssen sich daran gewöhnen , dass sie nun in einer Mannschaft \* Döberitzer Heide \* spielen . Und das heißt gemeinsam und nicht gegeneinander . Ermahnungen von der Seitenlinie , miteinander fair umzugehen und sich nicht beim kleinsten Schuss gegenseitig zu zerfleischen , sind normalerweise überflüssig . Vorerst allerdings hilfreich .

rhetorical structure (old\_grid)  
 tree

sind , sie müssen sich daran gewöhnen , dass sie nun in einer Mannschaft \* Döberitzer Heide \* spielen . Und  
 sein , sie müssen sich daran gewöhnen , dass sie nun in ein Mannschaft \* Döberitzer Heide \* spielen . und  
 3.Pl.Pres.Ind -- 3.Nom.Pl.Masc 3.Pl.Pres.Ind 3.Acc.Pl -- Inf -- -- 3.Nom.Pl.Masc -- -- Dat.Sg.Fem Dat.Sg.Fem -- Pos.\*\*.\* Nom.Sg.Fem -- 3.Pl.Pres.Ind -- --  
 VAFIN \$, PPER VMFIN PRF PROAV VVINF \$, KOUS PPER ADV APPR ART NN \$( ADJA NN \$( VVFIN \$, KON

dependencies (arches [SVG capable browsers])  
 information structure (grid)  
 discourse referents (grid)  
 coreference (discourse)  
 rhetorical structure (old\_grid)  
 tree

dass sie nun in einer Mannschaft \* Döberitzer Heide \* spielen . Und das heißt gemeinsam und nicht gegeneinander . Ermahnung  
 dass sie nun in ein Mannschaft \* Döberitzer Heide \* spielen . und der heißen gemeinsam und nicht gegeneinander . Ermahnung  
 -- 3.Nom.Pl.Masc -- -- Dat.Sg.Fem Dat.Sg.Fem -- Pos.\*\*.\* Nom.Sg.Fem -- 3.Pl.Pres.Ind -- -- Nom.Sg.Neut 3.Sg.Pres.Ind -- -- Pos -- --  
 KOUS PPER ADV APPR ART NN \$( ADJA NN \$( VVFIN \$, KON PDS VVFIN ADJD KON PTKNEG ADV \$, NN

dependencies (arches [SVG capable browsers])

## ANNIS2 (Zeldes et al. 09)

- ExtJS (JavaScript) web interface, PostgreSQL backend
- Types of annotations that can be merged:
  - spans with labels
  - DAGs with labelled edges
  - pointing relations between terminals and non-terminals
- Query language
  - exact match, regexp match on primary data and annotations
  - relations between elements
    - overlapping/contained/adjacent spans
    - hierarchical dominance (direct/indirect, common ancestors, left-/rightmost child, edges with specific labels, etc.)
- Search across different layers (using namespaces)
- Imported „standard“ corpora: TIGER, TüBa/D-Z, Ontonotes, ...
- <http://www.sfb632.uni-potsdam.de/~d1/annis/>

## Example: multi-level querying

Pronominal anaphors with subject antecedent in an RST-satellite, please. (ProCon10 corpus)

mmax:node &	#1: anaphor
tiger:pos=/P.*/ & #1 == #2 &	#2: anaphor is pronoun
mmax:node & #1 ->anaph_antec #3 &	#3: antecedent (coref)
tiger:node & #4 == #3 &	#4: antecedent (syn)
tiger:cat=„S“ & #5 >[func=„SB“]#4 &	#5: antec. sentence with #4 as subject
rst:cat=„segment“ & #6 _i_ #4 &	#6: antec. edu incl. #4
rst:cat & #7 >[func=/.*SAT/] #6	#7: parent of antec. edu is satellite

© Manfred Stede / NAACL Tutorial 2013

## Example: multi-level querying

Pronominal anaphors with subject antecedent in an RST-satellite, please. (ProCon10 corpus)

mmax:node &	#1: anaphor
tiger:pos=/P.*/ & #1 == #2 &	#2: anaphor is pronoun
mmax:node & #1 ->anaph_antec #3 &	#3: antecedent (coref)
tiger:node & #4 == #3 &	#4: antecedent (syn)
tiger:cat=„S“ & #5 >[func=„SB“]#4 &	#5: antec. sentence with #4 as subject
rst:cat=„segment“ & #6 _i_ #4 &	#6: antec. edu incl. #4
rst:cat & #7 >[func=/.*SAT/] #6	#7: parent of antec. edu is satellite

Expectation: pronouns prefer **subject antecedents** and **nucleus antecedents**

 Das Bild kann nicht angezeigt werden. Dieser Computer verfügt möglicherweise über zu wenig Arbeitsspeicher, um das Bild zu öffnen, oder das Bild ist beschädigt. Starten Sie den Computer neu, und öffnen Sie dann erneut die Datei. Wenn weiterhin das rote x angezeigt wird, müssen Sie das Bild möglicherweise löschen und dann erneut einfügen.

(Chiarcos, subm.)

## Accessibility in discourse: Distance-based approaches (Krasavina/Chiarcos 05)

- Referential Distance  
the number of clauses between anaphor and antecedent has an effect on the form of the anaphor (cf. Givón 1983)
- But  
*It is **not simple distance** that triggers the use of one anaphoric device over the other. Rather, it is **the rhetorical organisation** of that distance that determines whether a pronoun or a full NP is appropriate. (Fox 1987)*  
⇒ „Rhetorical Distance“

© Manfred Stede / NAACL Tutorial 2013

Thank you for your attention!



## References

- Nicholas Asher and Alex Lascarides. *Logics of Conversation*. Cambridge University Press, Cambridge, 2003
- Doug Beeferman, Adam Berger, and John Lafferty. Statistical models for text segmentation. *Machine Learning*, 34 (1-3):177–210, 1999
- Douglas Biber. A typology of English texts. *Linguistics*, 27:3–43, 1989.
- H. Bieler, S. Dipper, M. Stede. *Identifying formal and functional zones in film reviews*. In: Proc. SIGdial WS, Antwerpen, 2007
- Lynn Carlson, Daniel Marcu, and Mary Ellen Okurowski. Building a discourse-tagged corpus in the framework of Rhetorical Structure Theory. In Jan van Kuppevelt and Ronnie Smith, editors, *Current Directions in Discourse and Dialogue*. Kluwer, Dordrecht, 2003
- C. Chiarcos, S. Dipper, M. Götze, U. Leser, A. Lüdeling, J. Ritz, M. Stede. *A flexible framework for integrating annotations from different tools and tagsets*. In: *Traitement Automatique des Langues*, 49:217-246, 2008
- Freddy Choi. Advances in domain independent linear text segmentation. In Proc. of the 1st Meeting of the North American Chapter of the Association for Computational Linguistics (NAACL), pages 26–33, Seattle/WA, 2000
- Herbert Clark. Inferences in comprehension. In D. Laberge and S. J. Samuels, editors, *Basic Processes in Reading: Perception and Comprehension*, pages 243–263. Lawrence Erlbaum, Hillsdale/NJ, 1977
- Gaél Dias, Elsa Alves, and José Gabriel Pereira Lopes. Topic segmentation algorithms for text summarization and passage retrieval: An exhaustive evaluation. In Proc. of the 22nd national conference on Artificial intelligence (AAAI), pages 1334–1339, 2007
- Markus Egg and Gisela Redeker. How complex is discourse structure? In Proc. of the 7th Conference on International Language Resources and Evaluation (LREC), Malta, 2010
- Jacob Eisenstein. Hierarchical text segmentation from multi-scale lexical cohesion. In *Proceedings of ACL/HLT 2009*, pages 353–361, Boulder, Colorado, 2009.
- Robert Elwell and Jason Baldridge. Discourse connective argument identification with connective specific rankers. In Proc. of the IEEE Conference on Semantic Computing (ICSC), Santa Clara/CA, 2008
- Michel Galley, Kathleen R. McKeown, Eric Fosler-Lussier, and Hongyan Jing. Discourse segmentation of multi-party conversation. In Proc. of ACL 2003, pages 562–569, Sapporo, Japan, 2003
- Gundel, Jeanette K., Nancy Hedberg, Ron Zacharski. 1993. Cognitive status and the form of referring expressions in discourse. *Language* 69:274-307

- Marti A. Hearst. TextTiling: Segmenting text into multi-paragraph subtopic passages. *Computational Linguistics*, 23 (1):33–64, 1997
- Hugo Hernaut, Hmut Prendinger, David duVerle, and Mitsuru Ishizuka. Hilda: A discourse parser using support vector machine classification. *Dialogue and Discourse*, 1(3):1–33, 2010.
- Graeme Hirst and David St-Onge. Lexical chains as representations of context for the detection and correction of malapropisms. In Christiane Fellbaum, editor, *WordNet: An electronic lexical database*. MIT Press, 1998
- Jerry Hobbs. Coherence and coreference. *Cognitive Science*, 3:67–90, 1979.
- Manfred Klenner and Don Tuggener. An incremental entity-mention model for coreference resolution with restrictive antecedent accessibility. In Proc. of the Conference on Recent Advances in Natural Language Processing (RANLP), pages 178–185. Hissar, Bulgaria, 2011
- Alistair Knott and Robert Dale. Using linguistic phenomena to motivate a set of coherence relations. *Discourse Processes*, 18(1):35–62, 1994
- Shalom Lappin and Herbert J. Leass. An algorithm for pronominal anaphora resolution. *Computational Linguistics*, 20(4):535–561, 1994
- William Mann and Sandra Thompson. Rhetorical structure theory: Towards a functional theory of text organization. *TEXT*, 8:243–281, 1988
- Daniel Marcu. The rhetorical parsing of unrestricted texts: A surface-based approach. *Computational Linguistics*, 26(3):395–448, 2000
- James R. Martin. *English text: system and structure*. John Benjamins, Philadelphia/Amsterdam, 1992.
- Eleni Miltsakaki, Nikhil Dinesh, Rashmi Prasad, Aravind Joshi, and Bonnie Webber. Experiments on sense annotations and sense disambiguation of discourse connectives. In Proc. of the 4th Workshop on Treebanks and Linguistic Theories (TLT), Barcelona, 2005
- Hemant Misra, Francois Yvon, Joemon Jose, and Olivier Cappe. Text segmentation via topic modeling: an analytical study. In Proc. of the 18th ACM Conference on Information and Knowledge Management, Hong Kong, 2009
- Jane Morris and Graeme Hirst. Lexical cohesion computed by thesaural relations as an indicator of the structure of text. *Computational Linguistics*, 17(1):21–48, 1991
- Pang, Bo and Lillian Lee. 2004. A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts. In *Proceedings of ACL 2004*, Barcelona

- Emily Pitler and Ani Nenkova. Using syntax to disambiguate explicit discourse connectives in text. In Proc. of the ACL-IJCNLP 2009 Conference Short Papers. Suntec, Singapore, 2009
- Livia Polanyi, Chris Culy, Martin van den Berg, Gian Lorenzo Thione, and David Ahn. A rule based approach to discourse parsing. In Proc. of the SIGDIAL '04 Workshop, Cambridge/MA, 2004
- R. Prasad, N. Dinesh, A. Lee, E. Miltsakaki, L. Robaldo, A. Joshi, and B. Webber. The Penn Discourse Treebank 2.0. In Proc. of the 6<sup>th</sup> LREC, Marrakech, Morocco, 2008
- Jeffrey C. Reynar. An automatic method of finding topic boundaries. In Proc. of the 32nd Annual Meeting of the Association for Computational Linguistics, pages 331–333, Las Cruces/NM, June 1994
- H. Gregory Silber and Kathleen F. McCoy. Efficiently computed lexical chains as an intermediate representation for automatic text summarization. *Computational Linguistics*, 28(4):488–496, 2002
- Wee Meng Soon, Hwee Tou Ng, and Daniel Chung Yong Lim. A machine learning approach to coreference resolution of noun phrases. *Computational Linguistics*, 27(4):521–544, 2001
- Caroline Sporleder and Alex Lascarides. Exploiting linguistic cues to classify rhetorical relations. In Proc. of the Conference on Recent Advances in Natural Language Processing (RANLP), 2005.
- Caroline Sporleder and Alex Lascarides. Using automatically labelled examples to classify rhetorical relations: An assessment. *Natural Language Engineering*, 14(3):369–416, 2008
- Manfred Stede. The Potsdam Commentary Corpus. In Proc. of the ACL Workshop on Discourse Annotation, pages 96–102, Barcelona, 2004
- Nicola Stokes, Joe Carthy, and Alan F. Smeaton. Select: a lexical cohesion based news story segmentation system. *AI Communication*, 17:3–12, 2004
- Maite Taboada, Julian Brooke, and Manfred Stede. Genre-based paragraph classification for sentiment analysis. In Proc. of the SIGDIAL 2009 Conference, page 62–70, London, UK, September 2009
- M. Taboada, J. Brooke, M. Tofiloski, K. Voll, M. Stede. *Lexicon-based methods for sentiment analysis* In: *Computational Linguistics*, 37(2):267–307, 2011
- S. Teufel, A. Siddharthan, C. Batchelor. 2009. Towards Discipline-Independent Argumentative Zoning: Evidence from Chemistry and Computational Linguistics. In: *Proceedings of EMNLP-09, Singapore*
- Milan Tofiloski, Julian Brooke, and Maite Taboada. A syntactic and lexical-based discourse segmenter. In Proc. of ACL/IJCNLP (short papers), Suntec, Singapore, 2009
- Bonnie Webber, Alistair Knott, and Aravind Joshi. Multiple discourse connectives in a lexicalized grammar for discourse. In 3rd International Workshop on Computational Semantics (IWCS), NL- Tilburg, 1999
- Florian Wolf and Edward Gibson. Representing discourse coherence: a corpus-based study. *Computational Linguistics*, 31(2):249–287, 2005
- Zeldes, Amir, Ritz, Julia, Lüdeling, Anke & Chiarcos, Christian (2009), "ANNIS: A Search Tool for Multi-Layer Annotated Corpora". In: *Proceedings of Corpus Linguistics 2009*, July 20-23, Liverpool, UK