# Computational Discourse

## Chapter 21

Lecture #15

November 2012

1

# Discourse

- Consists of collocated, structured, coherent groups of sentences
  - What makes something a discourse as opposed to a set of unrelated sentences?
  - How can text be structured (related)?

- \* Monologue: a speaker (writer) and hearer (reader) with communication flow in one direction only
- Dialogue: each participant takes turn being the speaker and the hearer (so 2-way participation)
  - Human-human dialogue
  - Human-computer dialogue (conversational agent)

2

# Discourse Phenomina: Coreference Resolution

- The Tin Woodman went to the Emerald City to see the Wizard of Oz and ask for a heart. After he asked for it, the Woodman waited for the Wizard's response.

3

# Discourse Phenomina: Coreference Resolution

- The Tin Woodman went to the Emerald City to see the Wizard of Oz and ask for a heart. After he asked for it, the Woodman waited for the Wizard's response.

- What do we need to resolve?
- Why is it important?
  - Information extraction, summarization, conversational agents

4

# Coherence Relations: Coreference

- First Union Corp is continuing to wrestle with severe problems. According to industry insiders at Pain Webber, their president, John R. Georgius, is planning to announce his retirement tomorrow.

5

# Coherence Relations: Coreference

- First Union Corp is continuing to wrestle with severe problems. According to industry insiders at Pain Webber, their president, John R. Georgius, is planning to announce his retirement tomorrow.

- First Union Corp is continuing to wrestle with severe problems. According to industry insiders at Pain Webber, their president, John R. Georgius, believes Pain Webber can be instrumental in solving most of First Union's problems.

6

## Coherence Relations (Discourse Structure)

- First Union Corp is continuing to wrestle with severe problems. According to industry insiders at Pain Webber, their president, John R. Georgius, is planning to announce his retirement tomorrow.

- Reasonable summary:
  First Union President John R. Georgius is planning to announce his retirement tomorrow.

What you need to know: coherence relations between text segment – the first sentence is providing background for the more important 2nd sentence.

7

## Coherence (relation based)

- John hid Bill's car keys. He was drunk.

- ?? John hid Bill's car keys. He likes spinach.

- Coherence Relations – relations such as EXPLANATION or CAUSE that exists between two coherent sentences. Connections between utterances.

8

## More Coherence (entity based)

a) John went to his favorite music store to buy a piano.
b) He had frequented the store for many years.
c) He was excited that he could finally buy a piano.
d) He arrived just as the store was closing for the day.

e) John went to his favorite music store to buy a piano.
f) It was a store John had frequented for many years.
g) He was excited that he could finally buy a piano.
h) It was closing just as John arrived.

9

## More Coherence (entity based)

a) John went to his favorite music store to buy a piano.
b) He had frequented the store for many years.
c) He was excited that he could finally buy a piano.
d) He arrived just as the store was closing for the day.

e) John went to his favorite music store to buy a piano.
f) It was a store John had frequented for many years.
g) He was excited that he could finally buy a piano.
h) It was closing just as John arrived.

10

## Discourse Segmentation

- We want to separate a document into a linear sequence of subtopics

- Unsupervised Discourse Segmentation: Marti Hearst's TextTiling (done in early 90's)

11

## Consider a 23 paragraph article broken into segments (subtopics):

- 1-2 Intro to Magellan space probe
- 3-4 Intro to Venus
- 5-7 Lack of craters
- 8-11 Evidence of volcanic action
- 12-15 River Styx
- 16-18 Crustal spreading
- 19-21 Recent volcanism
- 22-23 Future of Magellan

Wants to do this in an unsupervised fashion – how?
Text Cohesion

12

## Text Cohesion

Halliday and Hasan (1976): "The use of certain linguistic devices to link or tie together textual units"

- Lexical cohesion: Indicated by relations between words in the two units (identical word, synonym, hypernym)
  - Before winter I built a chimney, and shingled the sides of my house..
  - I thus have a tight shingled and plastered house.

- Non-lexical cohesion like anaphora
  - Peel, core and slice the pears and the apples.
  - Add the fruit to the skillet.

13

## Intuition to a Cohesion-based approach to segmentation

- Sentences or paragraphs in a subtopic are cohesive with each other, but not with paragraphs in a neighboring subtopic.
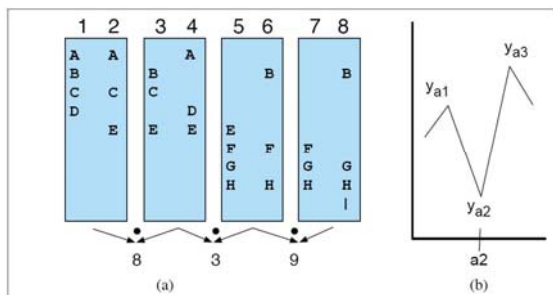
14

## From Hearst 1997



15

## TextTiling (Hearst, 1997)

1. Tokenization – convert words to lower case, remove stop words, stem words, group into pseudo-sentences
2. Lexical Score Determination – check scores between each pair of sentences = average similarity of the words in the pseudo-sentences before the gap to the pseudo-sentences after the gap

$$\mathrm{sim}_{\mathrm{cosine}}(\vec{b}, \vec{a}) = \frac{\vec{b} \cdot \vec{a}}{|\vec{b}||\vec{a}|} = \frac{\sum_{i=1}^{N} b_i \times a_i}{\sqrt{\sum_{i=1}^{N} b_i^2}\sqrt{\sum_{i=1}^{N} a_i^2}}$$

3. Boundary Identification – assign a cut-off distance to identify a new segment.

16

**Figure 21.1**

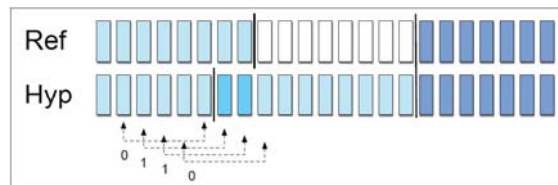## Supervised Discourse Segmentation

- To be used when it is relatively easy to acquire boundary-labeled training data
  - News stories from TV broadcasts
  - Paragraph segmentation

- Lots of different classifiers have been used
  - Feature set; generally a superset of those used for unsupervised segmentation
  - + discourse markers and cue words

- Discourse Markers generally domain specific

18

## Supervised Discourse Segmentation

- Supervised machine learning

  - Label segment boundaries in training and test set

  - Extract features in training

  - Learn a classifier

  - In testing, apply features to predict boundaries

- Evaluation – usual measures of precision, recall, and F-measure don't work – need to be sensitive to near-misses.

19

---

**Figure 21.2**

---

## What makes a text coherent?

- Appropriate use of coherence relations between subparts of the discourse --rhetorical structure

- Appropriate sequencing of subparts of the discourse --discourse/topic structure

- Appropriate use of referring expressions

21

---

## Coherence Relations

- Possible connections between utterances in a discourse. Such as in Hobbs 1997.

- Result: Infer that the state or event asserted by $S_0$ causes or could cause the state or event asserted in $S_1$.
  - The Tin Woodman was caught in the rain. His joints rusted.

- Explanation: Infer that the state or event asserted by $S_1$ causes or could cause the state or event asserted by $S_0$.
  - John hid Bill's car keys. He was drunk.

22

---

## Coherence Relations

- Parallel: Infer $p(a_1, a_2, \ldots)$ from the assertion of $S_0$ and $p(b_1, b_2, \ldots)$ from the assertion of $S_1$, where $a_i$ and $b_i$ are similar, for all i.
  - The scarecrow wanted some brains. The Tin Woodman wanted a heart.

- Elaboration: Infer the same proposition P from the assertions of $S_0$ and $S_1$.
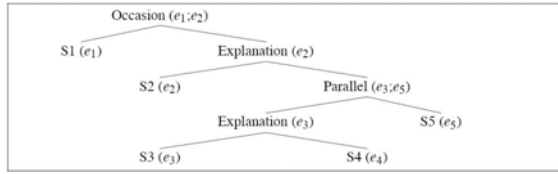  - Dorothy was from Kansas. She lived in the midst of the great Kansas prairies.

23

---

## Coherence Relations

- Occasion: A change of state can be inferred from the assertion of $S_0$, whose final state can be inferred from $S_1$, or a change of state can be inferred from the assertion of $S_1$, whose initial state can be inferred from $S_0$.
  - Dorothy picked up the oil-can. She oiled the Tin Woodman's joints.

24

---

4

## Hierarchical structures

(S1) John went to the bank to deposit his paycheck.

(S2) He then took a train to Bill's car dealership.

(S3) He needed to buy a car.

(S4) The company he works for now isn't near any public transportation.

(S5) He also wanted to talk to Bill about their softball league.



---

## Rhetorical Structure Theory

- See old slides

- See old slides on referring and discourse models

---

## 5 Types of Referring Expressions

- **Indefinite Noun Phrases**: Introduces into discourse context entities that are new to the hearer.
  - A man, some walnuts, this new computer

- **Definite Noun Phrases**: refers to an entity that is identifiable to the hearer (e.g., been mentioned previously or well known, in set of beliefs about the world).
  - "a big dog…. the dog…", the sun

---

## 5 Types of Referring Expressions

- **Pronouns**: another form of definite reference, generally stronger constraints on use than standard definite reference.
  - He, she, him, it, they…
- **Demonstratives**: demonstrative pronouns (this, that) can be alone or as determiners.
- **Names**: Common method of referring including people, organizations, and locations.

---

## Features for Filtering Potential Referents

- **Number Agreement**: pronoun and referent must agree in number (single, plural)
- **Person Agreement**: 1st, 2nd, 3rd
- **Gender Agreement**: male, female, nonpersonal (it)
- **Binding Theory Constraints**: constraints by syntactic relationships between a referential expression and a possible antecedent noun phrase in the same sentence.
  - John bought himself a new Ford.
  - John bought him a new Ford.
  - He said that he bought John a new Ford.

---

## Preferences in Pronoun Interpretation

- Recency
- Grammatical Role
- Repeated Mention
- Parallelism
- Verb Semantics
- Selectional Restrictions