Discourse: Reference

Ling571 Deep Processing Techniques for NLP March 2, 2011

What is a Discourse?

- Discourse is:
 - Extended span of text
 - Spoken or Written
 - One or more participants
 - Language in Use
 - Goals of participants
 - Processes to produce and interpret

Why Discourse?

- Understanding depends on context
 - Referring expressions: it, that, the screen
 - Word sense: plant
 - Intention: Do you have the time?
- Applications: Discourse in NLP
 - Question-Answering
 - Information Retrieval
 - Summarization
 - Spoken Dialogue
- Automatic Essay Grading

U: Where is A Bug's Life playing in Summit?
S: A Bug's Life is playing at the Summit theater.
U: When is it playing there?
S: It's playing at 2pm, 5pm, and 8pm.
U: I'd like 1 adult and 2 children for the first show. How much would that cost?

• Knowledge sources:

U: Where is A Bug's Life playing in Summit?
S: A Bug's Life is playing at the Summit theater.
U: When is it playing there?
S: It's playing at 2pm, 5pm, and 8pm.
U: I'd like 1 adult and 2 children for the first show. How much would that cost?

- Knowledge sources:
 - Domain knowledge

U: Where is A Bug's Life playing in Summit?
S: A Bug's Life is playing at the Summit theater.
U: When is it playing there?
S: It's playing at 2pm, 5pm, and 8pm.
U: I'd like 1 adult and 2 children for the first show. How much would that cost?

- Knowledge sources:
 - Domain knowledge
 - Discourse knowledge

U: Where is A Bug's Life playing in Summit?
S: A Bug's Life is playing at the Summit theater.
U: When is it playing there?
S: It's playing at 2pm, 5pm, and 8pm.
U: I'd like 1 adult and 2 children for the first show. How much would that cost?

- Knowledge sources:
 - Domain knowledge
 - Discourse knowledge
 - World knowledge

Coherence

- First Union Corp. is continuing to wrestle with severe problems. According to industry insiders at PW, their president, John R. Georgius, is planning to announce his retirement tomorrow.
- Summary:
- First Union President John R. Georgius is planning to announce his retirement tomorrow.
- Inter-sentence coherence relations:
 - Second sentence: main concept (nucleus)
 - First sentence: subsidiary, background

Different Parameters of Discourse

- Number of participants
 - Multiple participants -> Dialogue
- Modality
 - Spoken vs Written
- Goals
 - Transactional (message passing) vs Interactional (relations, attitudes)
 - Cooperative task-oriented rational interaction

Spoken vs Written Discourse

Speech

- Paralinguistic effects
 - Intonation, gaze, gesture
- Transitory
- Real-time, on-line
- Less "structured"
 - Fragments
 - Simple, Active, Declarative
 - Topic-Comment
 - Non-verbal referents
 - Disfluencies
 - Self-repairs
 - False Starts
 - Pauses

- Written text
 - No paralinguistic effects
 - "Permanent"
 - Off-line. Edited, Crafted
 - More "structured"
 - Full sentences
 - Complex sentences
 - Subject-Predicate
 - Complex modification
 - More structural markers
 - No disfluencies

Spoken vs Written: Representation

- Spoken "text" "same" if:
 - Recorded (Audio/Video Tape)
 - Transcribed faithfully
 - Always some interpretation
 - Text (normalized) transcription
 - Map paralinguistic features
 - e.g. pause = -,+,++
 - Notate accenting, pitch

- Written text "same" if:
 - Same words
 - Same order
 - Same punctuation (headings)
 - Same lineation

Agenda

- Coherence: Holding discourse together
 - Coherence types and relations
- Reference resolution
 - Referring expressions
 - Information status and structure
 - Features and Preferences for resolution
 - Knowledge-rich, deep analysis approaches
 - Lappin&Leass,
 - Hobbs

Coherence Relations

- John hid Bill's car keys. He was drunk.
- ?? John hid Bill's car keys. He likes spinach.
- Why odd?
 - No obvious relation between sentences
 - Readers often try to construct relations
- How are first two related?
 - Explanation/cause
- Utterances should have meaningful connection
 - Establish through coherence relations

Entity-based Coherence

- John went to his favorite music store to buy a piano.
- He had frequented the store for many years.
- He was excited that he could finally buy a piano.
- VS
 - John went to his favorite music store to buy a piano.
 - It was a store John had frequented for many years.
 - He was excited that he could finally buy a piano.
 - It was closing just as John arrived.
- Which is better? Why?
 - 'about' one entity vs two, focuses on it for coherence

- Match referring expressions to referents
- Syntactic & semantic constraints
- Syntactic & semantic preferences

• Reference resolution algorithms

- U: Where is A Bug's Life playing in Summit?
 S: A Bug's Life is playing at the Summit theater.
 U: When is it playing there?
 S: It's playing at 2pm, 5pm, and 8pm.
 U: I'd like 1 adult and 2 children for the first show. How much would that cost?
- Knowledge sources:
 - Domain knowledge
 - Discourse knowledge
 - World knowledge

Reference Resolution: Global Focus/ Task

- (From Grosz "Typescripts of Task-oriented Dialogues")
- E: Assemble the air compressor.
- •
- ... 30 minutes later...
- E: Plug it in / See if it works

- (From Grosz)
- E: Bolt the pump to the base plate
- A: What do I use?
-
- A: What is a ratchet wrench?
- E: Show me the table. The ratchet wrench is [...]. Show it to me.
- A: It is bolted. What do I do now?

Relation Recognition: Intention

- A: You seem very quiet today; is there a problem?
- A: Would you be interested in going to dinner tonight?

• B: I have a headache.

• B: I have a headache.

Answer

• Reject

Reference

 Queen Elizabeth set about transforming her husband, King George VI, into <u>a viable monarch</u>. Logue, a renowned speech therapist, was summoned to help the King overcome his <u>speech</u> <u>impediment</u>...

Referring expression: (refexp) Linguistic form that picks out entity in some model That entity is the "referent" When introduces entity, "evokes" it Set up later reference, "antecedent" 2 refexps with same referent "co-refer"

Reference (terminology)

- Queen Elizabeth set about transforming her husband, King George VI, into <u>a viable monarch</u>. Logue, a renowned speech therapist, was summoned to help the King overcome his <u>speech</u> <u>impediment</u>...
- Anaphor:
 - Abbreviated linguistic form interpreted in context
 - Her, his, the King
 - Refers to previously introduced item ("accesses")
 - Referring expression is then anaphoric

Referring Expressions

- Many alternatives:
 - Queen Elizabeth, she, her, the Queen, etc
 - Possible correct forms depend on discourse context
 - E.g. she, her presume prior mention, or presence in world
- Interpretation (and generation) requires:
 - Discourse Model with representations of:
 - Entities referred to in the discourse
 - Relationships of these entities
 - Need way to construct, update model
 - Need way to map refexp to hearer's beliefs

Reference and Model





 Queen Elizabeth set about transforming her husband, King George VI, into <u>a viable monarch</u>. Logue, a renowned speech therapist, was summoned to help the King overcome his <u>speech</u> <u>impediment</u>...

Coreference resolution:

Find all expressions referring to same entity, 'corefer' Colors indicate coreferent sets Pronominal anaphora resolution: Find antecedent for given pronoun

Referring Expressions

- Indefinite noun phrases (NPs): e.g. "a cat"
 - Introduces new item to discourse context
- Definite NPs: e.g. "the cat"
 - Refers to item identifiable by hearer in context
 - By verbal, pointing, or environment availability; implicit
- Pronouns: e.g. "he", "she", "it"
 - Refers to item, must be "salient"
- Demonstratives: e.g. "this", "that"
 - Refers to item, sense of distance (literal/figurative)
- Names: e.g. "Miss Woodhouse","IBM"
 New or old entities

Information Status

- Some expressions (e.g. indef NPs) introduce new info
- Others refer to old referents (e.g. pronouns)
- Theories link form of refexp to given/new status

I ne givenr	iess nierarch	y:			
			uniquely		type
in focus $>$	activated >	familiar >	identifiable >	referential >	identifiable
{it}	$\left\{\begin{array}{c} that\\ this\\ this\\ this \ N\end{array}\right\}$	{that N}	{the N}	{indef. <i>this</i> N}	$\{a \mathbf{N}\}$

- Accessibility:
 - More salient elements easier to call up, can be shorter Correlates with length: more accessible, shorter refexp

Complicating Factors

• Inferrables:

- Refexp refers to inferentially related entity
 - I bought a car today, but the door had a dent, and the engine was noisy.
 - E.g. car -> door, engine
- Generics:
 - I want to buy a Mac. They are very stylish.
 - General group evoked by instance.
- Non-referential cases:
 - It's raining.

Syntactic Constraints for Reference Resolution

- Some fairly rigid rules constrain possible referents
- Agreement:
 - Number: Singular/Plural
 - Person: 1st: I,we; 2nd: you; 3rd: he, she, it, they
 - Gender: he vs she vs it

Syntactic & Semantic Constraints

- Binding constraints:
 - Reflexive (x-self): corefers with subject of clause
 - Pronoun/Def. NP: can't corefer with subject of clause

"Selectional restrictions":

- "animate": The cows eat grass.
- "human": The author wrote the book.
- More general: drive: John drives a car....

Syntactic & Semantic Preferences

- Recency: Closer entities are more salient
 - The doctor found an old map in the chest. Jim found an even older map on the shelf. It described an island.
- Grammatical role: Saliency hierarchy of roles
 - e.g. Subj > Object > I. Obj. > Oblique > AdvP
 - Billy Bones went to the bar with Jim Hawkins. He called for a glass of rum. [he = Billy]
 - Jim Hawkins went to the bar with Billy Bones. He called for a glass of rum. [he = Jim]

Syntactic & Semantic Preferences

- Repeated reference: Pronouns more salient
 - Once focused, likely to continue to be focused
 - Billy Bones had been thinking of a glass of rum. He hobbled over to the bar. Jim Hawkins went with him. He called for a glass of rum. [he=Billy]
- Parallelism: Prefer entity in same role
 - Silver went with Jim to the bar. Billy Bones went with him to the inn. [him = Jim]
 - Overrides grammatical role
- Verb roles: "implicit causality", thematic role match,...
 - John telephoned Bill. He lost the laptop.
 - John criticized Bill. He lost the laptop.

Reference Resolution Approaches

Common features

- "Discourse Model"
 - Referents evoked in discourse, available for reference
 - Structure indicating relative salience
- Syntactic & Semantic Constraints
- Syntactic & Semantic Preferences

• Differences:

 Which constraints/preferences? How combine? Rank?

A Resolution Algorithm (Lappin & Leass)

- Discourse model update:
 - Evoked entities:
 - Equivalence classes: Coreferent referring expressions
 - Salience value update:
 - Weighted sum of salience values:
 - Based on syntactic preferences
- Pronoun resolution:
 - Exclude referents that violate syntactic constraints
 - Select referent with highest salience value

Salience Factors (Lappin & Leass 1994)

Weights empirically derived from corpus

- Recency: 100
- Subject: 80
- Existential: 70
- Object: 50
- Indirect Object/Oblique: 40
- Non-adverb PP: 50
- Head noun: 80
- Parallelism: 35, Cataphora: -175
- Divide by 50% for each sentence distance

- John saw a beautiful Acura Integra in the dealership.
- He showed it to Bob.
- He bought it.

• John saw a beautiful Acura Integra in the dealership.

Referent	Phrases	Value
John	{John}	310
Integra	{a beautiful Acura Integra}	280
Dealership	{the dealership}	230

• He showed it to Bob.

Referent	Phrases	Value
John	{John, he1}	465
Integra	{a beautiful Acura Integra}	140
Dealership	{the dealership}	115

Referent	Phrases	Value
John	{John, he1}	465
Integra	{a beautiful Acura Integra}	420
Dealership	{the dealership}	115

• He showed it to Bob.

Referent	Phrases	Value
John	{John, he1}	465
Integra	{a beautiful Acura Integra}	140
Bob	{Bob}	270
Dealership	{the dealership}	115

Referent	Phrases	Value
John	{John, he1}	232.5
Integra	{a beautiful Acura Integra}	210
Bob	{Bob}	135
Dealership	{the dealership}	57.5

• He bought it.

Referent	Phrases	Value
John	{John, he1}	542.5
Integra	{a beautiful Acura Integra}	490
Bob	{Bob}	135
Dealership	{the dealership}	57.5

Hobbs' Resolution Algorithm

- Requires:
 - Syntactic parser
 - Gender and number checker
- Input:
 - Pronoun
 - Parse of current and previous sentences
- Captures:
 - Preferences: Recency, grammatical role
 - Constraints: binding theory, gender, person, number

Hobbs Algorithm

Intuition:

- Start with target pronoun
- Climb parse tree to S root
- For each NP or S
 - Do breadth-first, left-to-right search of children
 - Restricted to left of target
 - For each NP, check agreement with target
- Repeat on earlier sentences until matching NP found

Hobbs Algorithm Detail

- Begin at NP immediately dominating pronoun
- Climb tree to NP or S: X=node, p = path
- Traverse branches below X, and left of p
 - Breadth-first, Left-to-Right
 - If find NP, propose as antecedent
 - If separated from X by NP or S
- Loop: If X highest S in sentence, try previous sentences.
- If X not highest S, climb to next NP or S: X = node
- If X is NP, and p not through X's nominal, propose X
- Traverse branches below X, left of p: BF,LR
 - Propose any NP
- If X is S, traverse branches of X, right of p: BF, LR
 - Do not traverse NP or S; Propose any NP
 - Go to Loop

Hobbs Example



Lyn's mom is a gardener. Craige likes her.

Another Hobbs Example



Hobbs Algorithm

- Results: 88% accuracy ; 90+% intrasential
 - On perfect, manually parsed sentences
- Useful baseline for evaluating pronominal anaphora
- Issues:
 - Parsing:
 - Not all languages have parsers
 - Parsers are not always accurate
 - Constraints/Preferences:
 - Captures: Binding theory, grammatical role, recency
 - But not: parallelism, repetition, verb semantics, selection

Reference Resolution: Agreements

- Knowledge-based
 - Deep analysis: full parsing, semantic analysis
 - Enforce syntactic/semantic constraints
 - Preferences:
 - Recency
 - Grammatical Role Parallelism (ex. Hobbs)
 - Role ranking
 - Frequency of mention
- Local reference resolution
- Little/No world knowledge
- Similar levels of effectiveness

Questions

- 80% on (clean) text. What about...
 - Conversational speech?
 - Ill-formed, disfluent
 - Dialogue?
 - Multiple speakers introduce referents
 - Multimodal communication?
 - How else can entities be evoked?
 - Are all equally salient?

More Questions

- 80% on (clean) (English) text: What about..
 - Other languages?
 - Salience hierarchies the same
 - Other factors
 - Syntactic constraints?
 - E.g. reflexives in Chinese, Korean,...
 - Zero anaphora?
 - How do you resolve a pronoun if you can't find it?

Reference Resolution Algorithms

- Many other alternative strategies:
 - Linguistically informed, saliency hierarchy
 - Centering Theory
 - Machine learning approaches:
 - Supervised: Maxent
 - Unsupervised: Clustering
 - Heuristic, high precision:
 - Cogniac

Reference Resolution: Extensions

Cross-document co-reference

- (Baldwin & Bagga 1998)
- Break "the document boundary"
- Question: "John Smith" in A = "John Smith" in B?
- Approach:
 - Integrate:
 - Within-document co-reference
 - with
 - Vector Space Model similarity

Cross-document Coreference

- Run within-document co-reference (CAMP)
 Produce chains of all terms used to refer to entity
- Extract all sentences with reference to entity
 Pseudo per-entity summary for each document
- Use Vector Space Model (VSM) distance to compute similarity between summaries

Cross-document Coreference

• Experiments:

• 197 NYT articles referring to "John Smith"

- 35 different people, 24: 1 article each
- With CAMP: Precision 92%; Recall 78%
- Without CAMP: Precision 90%; Recall 76%
- Pure Named Entity: Precision 23%; Recall 100%

Conclusions

- Co-reference establishes coherence
- Reference resolution depends on coherence
- Variety of approaches:
 - Syntactic constraints, Recency, Frequency, Role
- Similar effectiveness different requirements
- Co-reference can enable summarization within and across documents (and languages!)