

# Computational Discourse

Yejin Choi  
University of Washington

Many slides from Aravind Joshi, Rashmi Prasad, Bonnie Weber,  
Regina Barzilay, Edward Loper, Luke Zettlemoyer

# Plan

1. Textual Coherence
2. Rhetorical Structure Theory (RST)
3. Penn Discourse Tree Bank (PDTB)
4. Coreference Resolution / Entity-grid Model

# Textual Coherence

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

# Textual 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.
  - He arrived just as the store was closing for the day.
- 
- 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.

# Why Model Coherence



*"How much wood could a woodchuck chuck if a woodchuck would chuck wood."*

*It depends on whether you are talking about African or European woodchucks.*

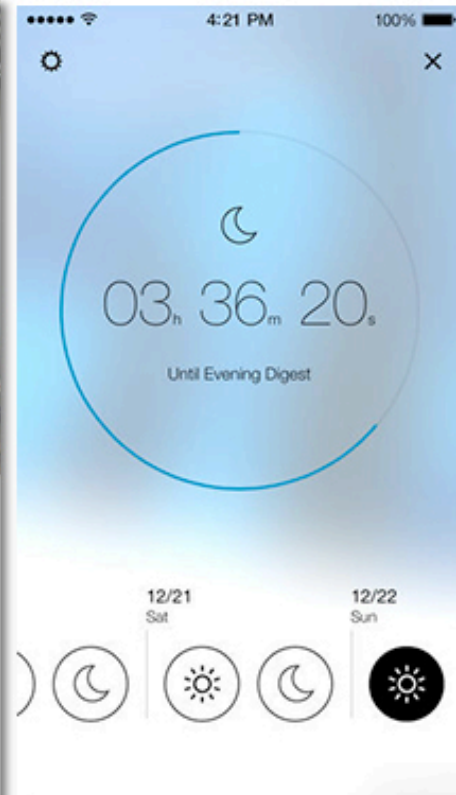
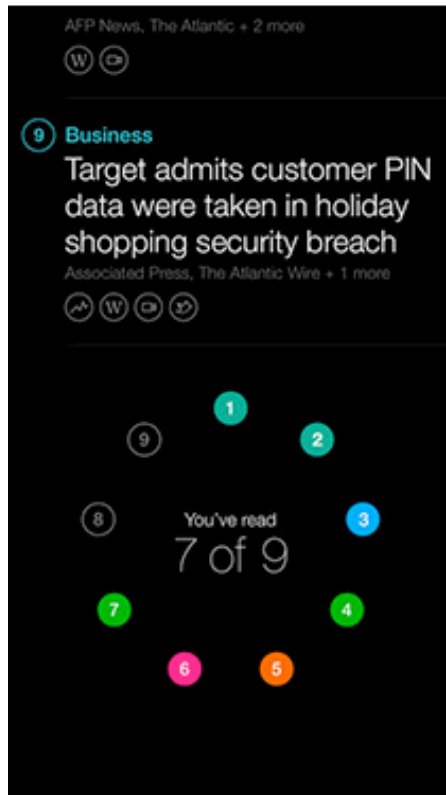
*"European woodchucks"*

*I found 8 European restaurants fairly close to you.*

# Long-term Coherent Conversation



# News aggregation and summary app



# Journalism: Robot or Human?

Despite an expected dip in profit, analysts are generally optimistic about **Steelcase** as it prepares to reports its third-quarter earnings on Monday, December 22, 2014. The consensus earnings per share estimate is 26 cents per share.

The consensus estimate remains unchanged over the past month, but it has decreased from three months ago when it was 27 cents. Analysts are expecting earnings of 85 cents per share for the fiscal year. Revenue is projected to be 5% above the year-earlier total of \$784.8 million at \$826.1 million for the quarter. For the year, revenue is projected to come in at \$3.11 billion.

The company has seen revenue grow for three quarters straight. The less than a percent revenue increase brought the figure up to \$786.7 million in the most recent quarter. Looking back further, revenue increased 8% in the first quarter from the year earlier and 8% in the fourth quarter.

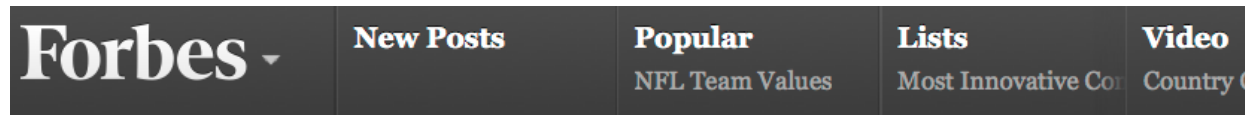
The majority of analysts (100%) rate Steelcase as a buy. This compares favorably to the analyst ratings of three similar companies, which average 57% buys. Both analysts rate Steelcase as a buy.

[Forbes.com; Dec 19, 2014]



# Writer-bots for earthquake & financial reports

While far from op-ed, some of the formulaic news articles are now written by computers.



2 FREE issues of Forbes



Forbes Partner

## Narrative Science

+ [Follow](#) (83)

**NEWS**

[Social](#)


[Archive](#)

Post 19 hours ago | 364 views

## Oracle Earnings Projected to Increase

Analysts expect higher profit for **Oracle** when the company reports its first quarter results on Thursday, September 18, 2014. The consensus estimate is calling for profit of 60 cents a share, reflecting a rise from 56 cents per share a year ago.

For the fiscal year, analysts are expecting earnings of \$3.01 per share. [read »](#)

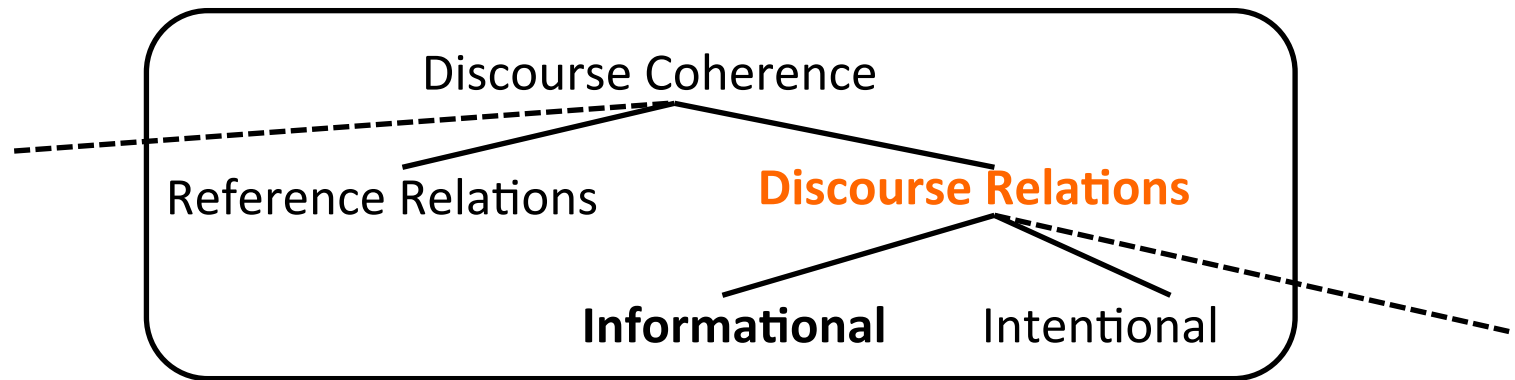
 **Narrative Science**, Partner

Post 19 hours ago | 246 views

## Rite Aid Profit Expected to Slip

# What is "discourse"?





Discourse is a coherent structured group of textual units



# Discourse “Relations”

- John hid Bill’s car keys. He was drunk.  
→ “Explanation” relation
- John hid Bill’s car keys. He likes spinach.  
→ ??? relation

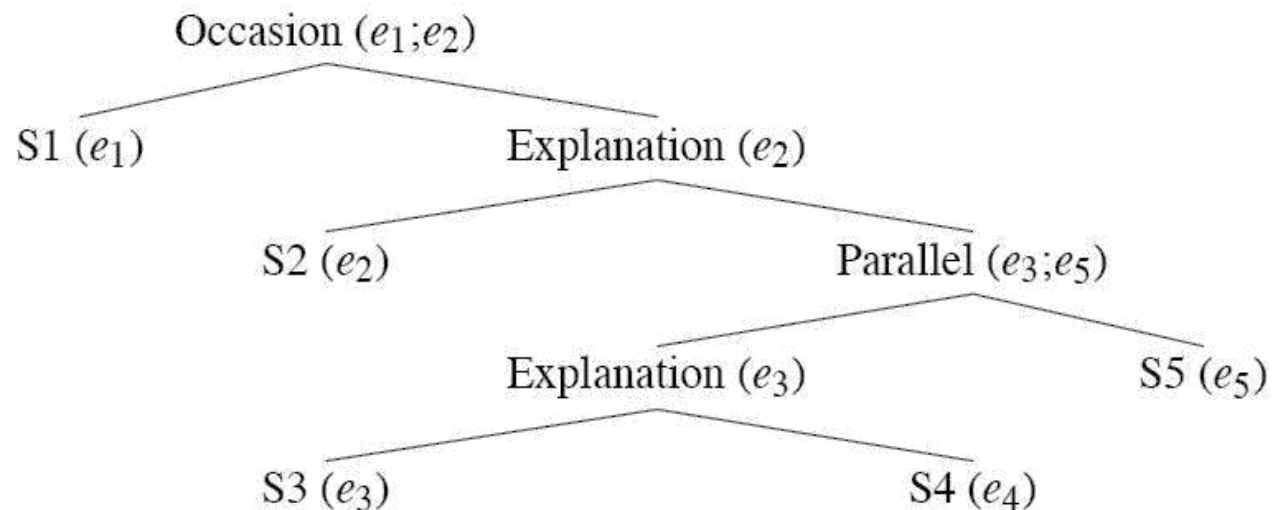
# Discourse “Relations”

- Dorothy was from Kansas. She lived on the Kansas prairies. 
- The tin woodman was caught in the rain. His joints rusted. 
- The scarecrow wanted some brains. The tin woodsman wanted a heart. 
- Dorothy picked up the oil-can. She oiled the Tin Woodman's joints. 

- **Result**
  - “*as a result ...*”
- **Occasion**
  - “*and then ...*”
- **Elaboration**
  - “*more specifically ...*”
- **Parallel**

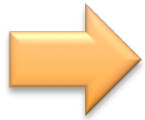
# Discourse Parsing: Tree of Relations

- Explanation
  - Elaboration
  - Result
  - Parallel
  - Occasion
- John went to the bank to deposit the paycheck. (e1)
  - He then took a train to Bill's car dealership. (e2)
  - He needed to buy a car. (e3)
  - The company he works for now isn't near any public transportation. (e4)
  - He also wanted to talk to Bill about their softball league. (e5)



# Plan

1. Textual Coherence



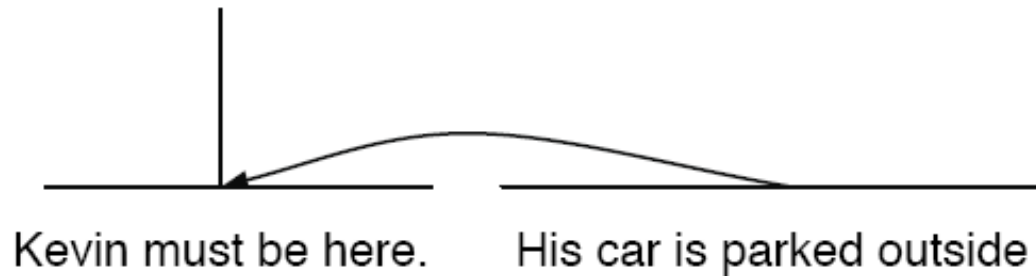
2. **Theory**: Rhetorical Structure Theory (RST)

3. **Corpus**: Penn Discourse Tree Bank (PDTB)

4. Coreference Resolution

# Rhetorical structure theory (RST)

Mann and Thompson, 1987



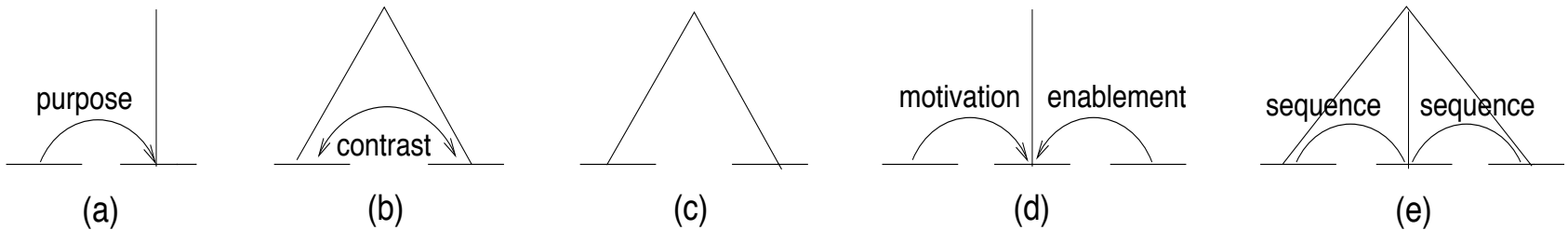
- **N**ucleus – the central unit, interpretable independently.
- **S**atellite – interpretation depends on N
- RST relation --- a set of constraints on the nucleus and satellite, w.r.t. the goals/beliefs/effects of the writer (W) and the reader (R)

<b>Relation Name:</b>	Evidence
<b>Constraints on N:</b>	R might not believe N to a degree satisfactory to W
<b>Constraints on S:</b>	R believes S or will find it credible
<b>Constraints on N+S:</b>	R's comprehending S increases R's belief of N
<b>Effects:</b>	R's belief of N is increased

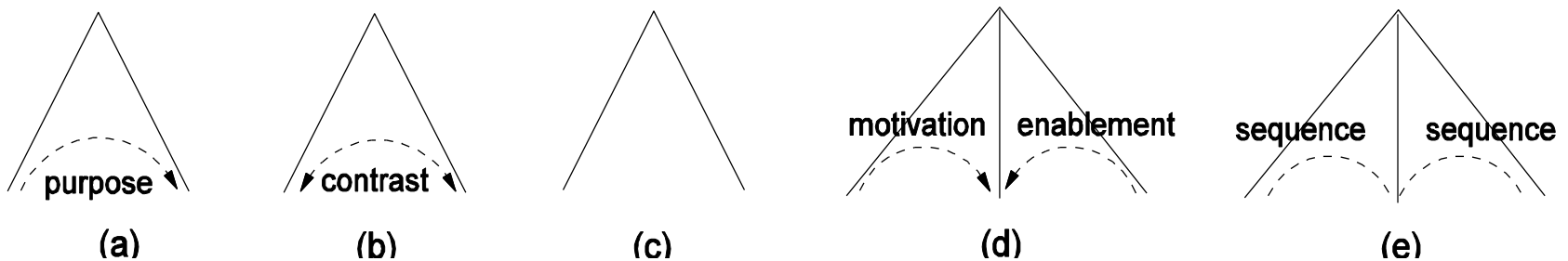
# Types of Schemas in RST

RST **schemas** := context-free rules for discourse structure

- whether or not the schema has binary, ternary, or arbitrary branching.
- whether or not the RHS has a head (called a *nucleus*);
- what rhetorical relation, if any, hold between right-hand side (RHS) sisters;



RST schema types in RST annotation

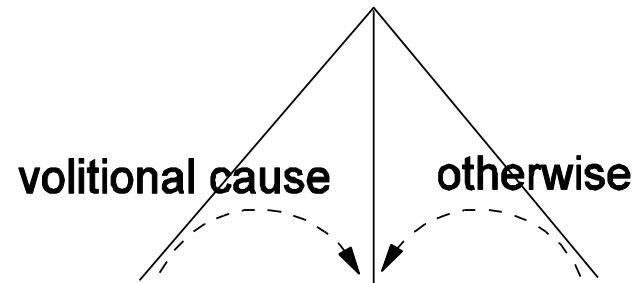
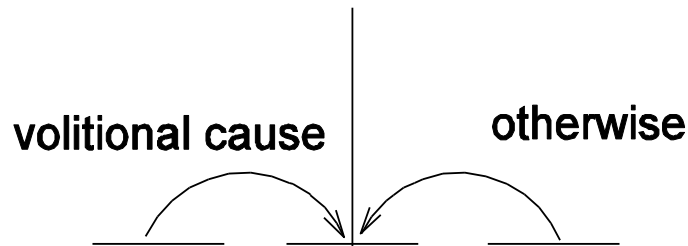


RST schema types in standard tree notation



## RST Example

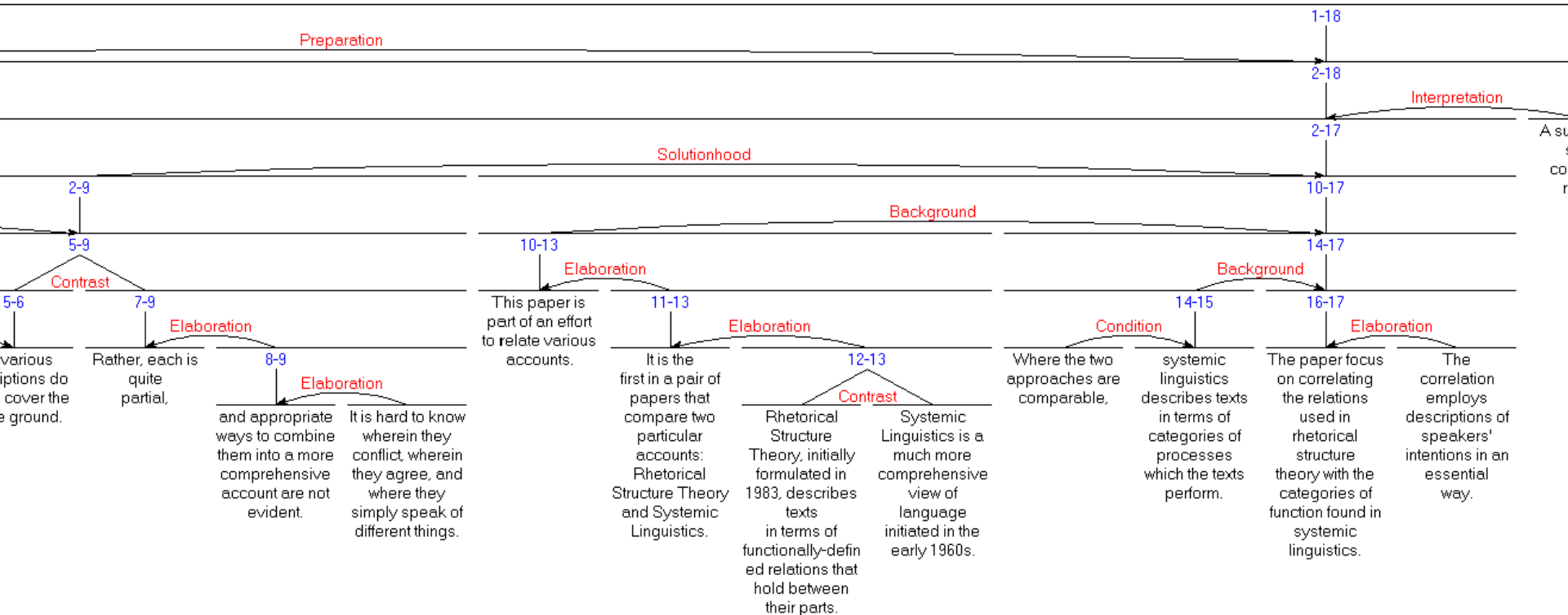
- (1) George Bush supports big business.
- (2) He's sure to veto House Bill 1711.
- (3) Otherwise, big business won't support him.



### Discourse structure as a **tree**:

- Leaf := an *elementary discourse unit* (a *continuous text span*)
- non-terminal := a contiguous, non-overlapping text span
- root := a complete, non-overlapping cover of the text

# RST Example



# From Theory to TreeBank

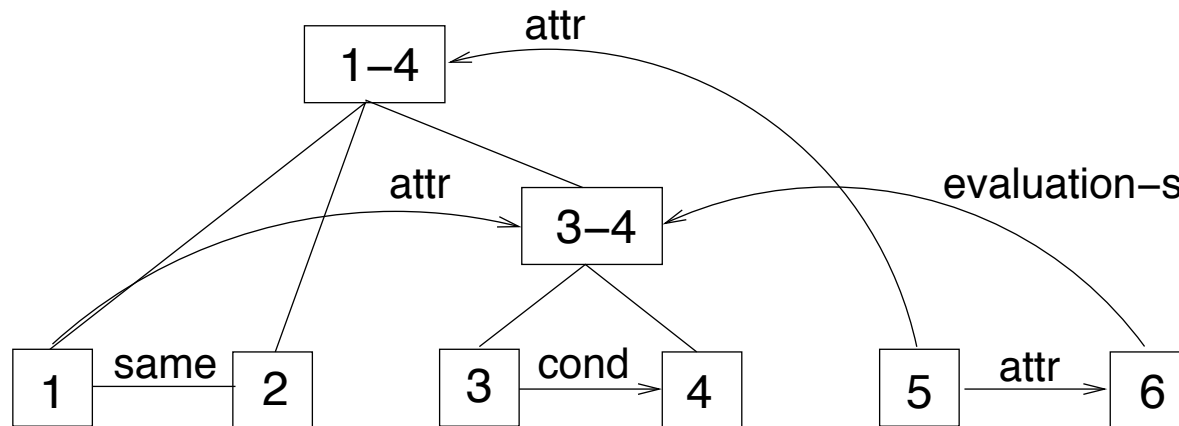
- Rhetorical Structure Theory: Mann and Thompson (1987)
- RST TreeBank: Carlson et al., (2001) defines 78 different RST relations, grouped into 16 classes.

<b>Relation Name:</b>	Evidence
<b>Constraints on N:</b>	R might not believe N to a degree satisfactory to W
<b>Constraints on S:</b>	R believes S or will find it credible
<b>Constraints on N+S:</b>	R's comprehending S increases R's belief of N
<b>Effects:</b>	R's belief of N is increased

## Discourse GraphBank [Wolf & Gibson 2005]

## Graph instead of a tree:

- (1) The administration should now state  
(2) that  
(3) if the February election is voided by the Sandinistas  
(4) they should call for military aid,  
(5) said former Assistant Secretary of State Elliot Abrams.  
(6) In these circumstances, I think they'd win.



# Penn Discourse Treebank (PDTB)

Discourse relations defined over “abstract objects”

Abstract Objects:

events, states, propositions (Asher, 1993)

Example of discourse relations:

Cause, temporal, contrast, condition

A discourse relation holds between  
*two and only two* AO arguments:



# Explicit Connectives

Explicit connectives are the lexical items that trigger discourse relations.

- Subordinating conjunctions (e.g., *when, because, although*, etc.)
    - The federal government suspended sales of U.S. savings bonds because Congress hasn't lifted the ceiling on government debt.
  - Coordinating conjunctions (e.g., *and, or, so, nor*, etc.)
    - The subject will be written into the prime-time shows, and viewers will be given a 900 number to call.
  - Discourse adverbials (e.g., *then, however, as a result*, etc.)
    - In the past, the socialist policies of the government strictly limited the profits businessmen could make. As a result, industry operated out of highly inefficient industrial units.
- **Arg2**: the argument with which connective is syntactically associated
  - **Arg1**: the other argument

# Explicit Connectives

Explicit connectives are the lexical items that trigger discourse relations.

- Subordinating conjunctions (e.g., *when, because, although*, etc.)
    - *The federal government suspended sales of U.S. savings bonds* because Congress hasn't lifted the ceiling on government debt.
  - Coordinating conjunctions (e.g., *and, or, so, nor*, etc.)
    - The subject will be written into the prime-time shows, and viewers will be given a 900 number to call.
  - Discourse adverbials (e.g., *then, however, as a result*, etc.)
    - In the past, the socialist policies of the government strictly limited the profits businessmen could make. As a result, industry operated out of highly inefficient industrial units.
- **Arg2**: the argument with which connective is syntactically associated
  - **Arg1**: the other argument

# Explicit Connectives

Explicit connectives are the lexical items that trigger discourse relations.

- Subordinating conjunctions (e.g., *when, because, although*, etc.)
    - *The federal government suspended sales of U.S. savings bonds* because Congress hasn't lifted the ceiling on government debt.
  - Coordinating conjunctions (e.g., *and, or, so, nor*, etc.)
    - *The subject will be written into the prime-time shows,* and viewers will be given a 900 number to call.
  - Discourse adverbials (e.g., *then, however, as a result*, etc.)
    - In the past, the socialist policies of the government strictly limited the profits businessmen could make. As a result, industry operated out of highly inefficient industrial units.
- **Arg2**: the argument with which connective is syntactically associated
  - **Arg1**: the other argument



# Explicit Connectives

Explicit connectives are the lexical items that trigger discourse relations.

- Subordinating conjunctions (e.g., *when, because, although*, etc.)
  - *The federal government suspended sales of U.S. savings bonds because Congress hasn't lifted the ceiling on government debt.*
- Coordinating conjunctions (e.g., *and, or, so, nor*, etc.)
  - *The subject will be written into the prime-time shows, and viewers will be given a 900 number to call.*
- Discourse adverbials (e.g., *then, however, as a result*, etc.)
  - *In the past, the socialist policies of the government strictly limited the profits businessmen could make. As a result, industry operated out of highly inefficient industrial units.*
- **Arg2**: the argument with which connective is syntactically associated
- **Arg1**: the other argument

# Argument Labels and Order

- **Arg2** is the argument with which connective is syntactically associated.
- **Arg1** is the other argument.
- Most oil companies, when they set exploration and production budgets for this year, forecast revenue of \$15 for each barrel of crude produced.
- The chief culprits, he says, are big companies and business groups that buy huge amounts of land "not for their corporate use, but for resale at huge profit." ... The Ministry of Finance, as a result, has proposed a series of measures that would restrict business investment in real estate even more tightly than restrictions aimed at individuals.

# Argument Labels and Order

- **Arg2** is the argument with which connective is syntactically associated.
- **Arg1** is the other argument.
- *Most oil companies, when they set exploration and production budgets for this year, forecast revenue of \$15 for each barrel of crude produced.*
- The chief culprits, he says, are big companies and business groups that buy huge amounts of land "not for their corporate use, but for resale at huge profit." ... The Ministry of Finance, as a result, has proposed a series of measures that would restrict business investment in real estate even more tightly than restrictions aimed at individuals.

# Argument Labels and Order

- **Arg2** is the argument with which connective is syntactically associated.
- **Arg1** is the other argument.
- *Most oil companies*, when they set exploration and production budgets for this year, *forecast revenue of \$15 for each barrel of crude produced.*
- *The chief culprits*, he says, *are big companies and business groups that buy huge amounts of land "not for their corporate use, but for resale at huge profit."* ... The Ministry of Finance, as a result, has proposed a series of measures that would restrict business investment in real estate even more tightly than restrictions aimed at individuals.

Relative location of Arg1?

# Finding Arg 1: Preliminary Experiment

- where to we find Arg 1 the most often?

CONN	Same	Previous	Multiple Previous	Distant
nevertheless				
otherwise				
as a result				
therefore				
instead				

## Finding Arg 1: Preliminary Experiment

- where to we find Arg 1 the most often?
- which connective has highest % of "distant" arg-1 ?

CONN	Same	Previous	Multiple Previous	Distant
nevertheless	9.7%	54.8%		
otherwise	11.1%	<b>77.8%</b>		
as a result	<b>4.8%</b>	69.8%		
therefore	<b>55%</b>	<b>35%</b>		
instead	22.7%	63.9%		

## Finding Arg 1: Preliminary Experiment

- where to we find Arg 1 the most often?
- which connective has highest % of "distant" arg-1 ?

CONN	Same	Previous	Multiple Previous	Distant
nevertheless	9.7%	54.8%	9.7%	<b>25.8%</b>
otherwise	11.1%	<b>77.8%</b>	5.6%	5.6%
as a result	<b>4.8%</b>	69.8%	7.9%	19%
therefore	<b>55%</b>	<b>35%</b>	5%	<b>5%</b>
instead	22.7%	63.9%	2.1%	11.3%

# Hierarchy of PDTB Discourse Relations

## CONTINGENCY

- Cause
  - Reason
  - Result
- Condition
  - Hypothetical
  - ...
- ...

## COMPARISON

- Contrast
  - Juxtaposition
  - Opposition
- Concession
  - Expectation
  - Contra-expectation
- ...

## TEMPORAL

- Asynchronous
- Synchronous
  - Precedence
  - Succession

Operating revenue rose 69% to A\$8.48 billion from A\$5.01 billion.

**But** the net interest bill jumped 85% to A\$686.7 million from A\$371.1 million.

The Texas oilman has acquired a 26.2% stake valued at more than \$1.2 billion in an automotive lighting company, Koito Manufacturing Co.

**But** he has failed to gain any influence at the company.

## EXPANSION

- Conjunction
- Instantiation
- Restatement
  - Specification
  - Equivalence
  - Generalization
- ...
- Exception
- List



# Annotation Overview (PDTB 1.0): Explicit Connectives

- All WSJ sections (25 sections; 2304 texts)
- 100 distinct types
  - Subordinating conjunctions – 31 types
  - Coordinating conjunctions – 7 types
  - Discourse Adverbials – 62 types

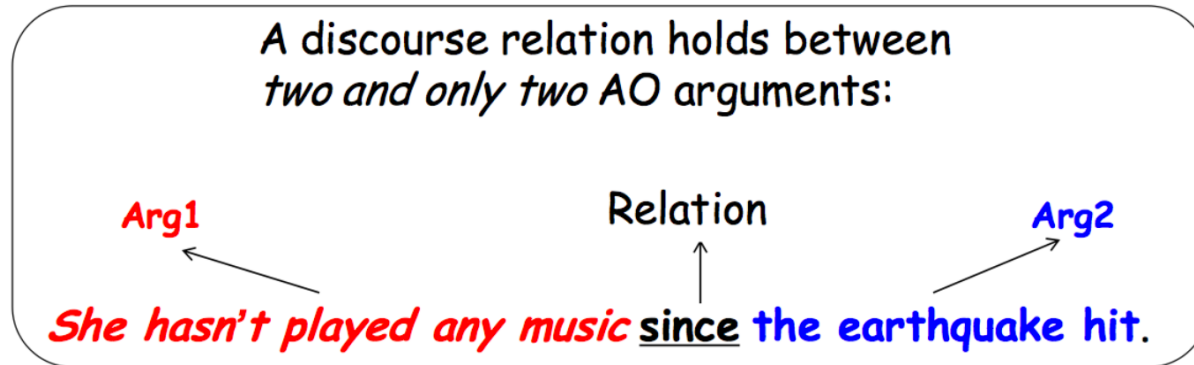
Some additional types will be annotated for PDTB-2.0.

- 18505 distinct tokens

# Natural Language Generation: Sentence Planning

Discourse analysis can help enhancing NLG. How?

- the relative **linear order** of component semantic units
- whether or not to explicitly realize discourse relations (**occurrence**), and if so, how to realize them (**lexical selection and placement**)



## NLG: Preliminary Experiment 2

**Question:** Given a subordinating conjunction and its arguments, in what relative order should the arguments be realized? Arg1-Arg2? Arg2-Arg1?

☞ Different patterns for different connectives

- When almost equally distributed:  
54% (Arg1-Arg2) and 46% (Arg2-Arg1)
- Although and (even) though have opposite patterns:  
Although: 37% (Arg1-Arg2) and 63% (Arg2-Arg1)  
(Even) though: 72% (Arg1-Arg2) and 28% (Arg2-Arg1)

## NLG: Preliminary Experiment 2

**Question:** What constrains the lexical choice of a connective for a given discourse relation? (Prasad et al., 2005)

- Testing a prediction for lexical choice rule for CAUSAL because and since (Elhadad and McKeown, 1990):
  - **Assumption:** New information tends to be placed at the end and *given* information at the beginning.
  - **Claim:** Because presents *new* information, and since presents *given* information
  - **Lexical choice rule:** Use because when subordinate clause is postposed (Arg1-Arg2); use since when subordinate clause is preposed (Arg2-Arg1)
- ☞ Because does tend to appear with Arg1-Arg2 order (90%), but CAUSAL since is equally distributed as Arg1-Arg2 and Arg2-Arg1.

# Sense Disambiguation of Connectives

Some discourse connectives are **polysemous**, e.g.,

- While: comparative, oppositive, concessive
- Since: temporal, causal, temporal/causal
- When: temporal/causal, conditional

Sense disambiguation is required for many applications:

- **Discourse parsing**: identification of arguments
- **NLG**: relative order of arguments
- **MT**: choice of connective in target language

# Sense Disambiguation: Preliminary Experiment

- Features (from raw text and PTB):
  - Form of auxiliary *have* - *Has*, *Have*, *Had* or *Not Found*.
  - Form of auxiliary *be* – *Present* (am, is, are), *Past* (was, were), *Been*, or *Not Found*.
  - Form of the head - *Present* (part-of-speech VBP or VBZ), *Past* (VBD), *Past Participial* (VBN), *Present Participial* (VBG).
  - Presence of a modal - *Found* or *Not Found*.
  - Relative position of Arg1 and Arg2: preposed, postposed
  - If the same verb was used in both arguments
  - If the adverb “not” was present in the head verb phrase of a single argument

- MaxEnt classifier (McCallum, 2002)
- Baseline: most frequent sense (**CAUSAL**)
- 10-fold cross-validation

Experiment	Accuracy	Baseline
(T,C,T/C)	75.5%	53.6%
({T,T/C}, C)	90.1%	53.6%
(T,{C,T/C})	74.2%	65.6%
(T,C)	89.5%	60.9%

T=temporal, C=causal, T/C=temporal/causal

👉 **15-20% improvement over baseline across the board, with state of the art.**

# Robot or Human?

Despite an expected dip in profit, analysts are generally optimistic about **Steelcase** as it prepares to reports its third-quarter earnings on Monday, December 22, 2014. The consensus earnings per share estimate is 26 cents per share.

The consensus estimate remains unchanged over the past month, but it has decreased from three months ago when it was 27 cents. Analysts are expecting earnings of 85 cents per share for the fiscal year. Revenue is projected to be 5% above the year-earlier total of \$784.8 million at \$826.1 million for the quarter. For the year, revenue is projected to come in at \$3.11 billion.

The company has seen revenue grow for three quarters straight. The less than a percent revenue increase brought the figure up to \$786.7 million in the most recent quarter. Looking back further, revenue increased 8% in the first quarter from the year earlier and 8% in the fourth quarter.

The majority of analysts (100%) rate Steelcase as a buy. This compares favorably to the analyst ratings of three similar companies, which average 57% buys. Both analysts rate Steelcase as a buy.

# Plan

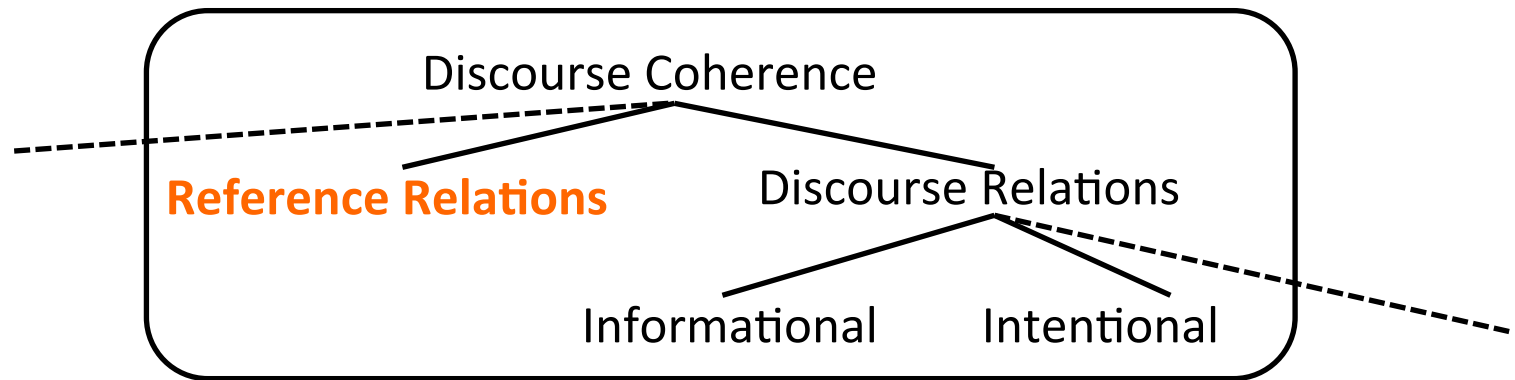
1. Textual Coherence
2. Rhetorical Structure Theory (RST)
3. Penn Discourse Tree Bank (PDTB)
4. Coreference Resolution





# Discourse Coherence

Discourse is a coherent structured group of textual units



# The Problem: Find and Cluster Mentions

Victoria Chen, Chief Financial Officer of Megabucks banking corp since 2004, saw her pay jump 20%, to \$1.3 million, as the 37 year old also became the Denver-based financial services company's president. It has been ten years since she came to Megabucks from rival Lotsabucks.



Mention Detection

[Victoria Chen], [Chief Financial Officer of [Megabucks banking corp] since 2004], saw [[her] pay] jump 20%, to \$1.3 million, as [the 37 year old] also became the [[Denver-based financial services company] 's president]. It has been ten years since [she] came to [Megabucks] from rival [Lotsabucks].

# The Problem: Find and Cluster Mentions

[Victoria Chen], [Chief Financial Officer of [Megabucks banking corp] since 2004], saw [[her] pay] jump 20%, to \$1.3 million, as [the 37 year old] also became the [[Denver-based financial services company] 's president]. It has been ten years since [she] came to [Megabucks] from rival [Lotsabucks].



Co-reference chains:

- 1 {Victoria Chen, Chief Financial Officer...since 2004, her, the 37-year-old, the Denver-based financial services company's president}
- 2 {Megabucks Banking Corp, Denver-based financial services company, Megabucks}
- 3 {her pay}
- 4 {rival Lotsabucks}

# Types of Coreference (I)

- Types of coreferent phrase:
  - Referential (“semantically definite”) NPs

The author of the book walked in.  
His name was John Smith.  
Mr. Smith said...
  - Anaphors

Mr. Smith walked in.  
He talked about his car.
  - Descriptive NPs

The stock price fell from \$4.02 to \$3.85.

# Types of Coreference (II)

- Types of antecedent:
  - Non-generic referring NPs
    - Mr. Smith likes his car.
  - Generic referring NPs
    - People like their cars.
  - Non-referring NPs
    - No one talked about their car.
  - Clauses
    - Driving fast isn't safe, but it's fun.

# Coreference as Clustering

The coreference problem can be solved by assigning all NPs in the text to equivalence classes, i.e., by clustering.  
[Cardie and Wagstaff, 1999]

We need:

- a *representation* of NPs (as a set of features)
- a *distance metric*
- a clustering *algorithm*.

# Representing Mentions

Each NP is represented as a set of features:

- **head noun**: last word of the NP;
- **position** in the document;
- **pronoun type**: nominative, accusative, possessive, ambiguous;
- **article**: indefinite, definite, none;
- **appositive**: based on heuristics (commas, etc.)
- **number**: plural, singular;
- **proper name**: based on heuristics (capitalization, etc.);
- **semantic class**: based on Wordnet;
- **gender**: masculine, feminine, either, neuter;
- **animacy**: based on semantic class.

# Example Mentions

Words, Head Noun (in <b>bold</b> )	Position	Pronoun Type	Article	Appositive	Number	Proper Name	Semantic Class	Gender	Animacy
John <b>Simon</b> Chief Financial <b>Officer</b> Prime <b>Corp.</b> <b>1986</b> <b>his</b> <b>pay</b> <b>20%</b> <b>\$1.3 million</b> the <b>37-year-old</b> the financial-services <b>company</b> <b>president</b>	1	NONE	NONE	NO	SING	YES	HUMAN	MASC	ANIM
	2	NONE	NONE	NO	SING	NO	HUMAN	EITHER	ANIM
	3	NONE	NONE	NO	SING	NO	COMPANY	NEUTER	INANIM
	4	NONE	NONE	NO	PLURAL	NO	NUMBER	NEUTER	INANIM
	5	POSS	NONE	NO	SING	NO	HUMAN	MASC	ANIM
	6	NONE	NONE	NO	SING	NO	PAYMENT	NEUTER	INANIM
	7	NONE	NONE	NO	PLURAL	NO	PERCENT	NEUTER	INANIM
	8	NONE	NONE	NO	PLURAL	NO	MONEY	NEUTER	INANIM
	9	NONE	DEF	NO	SING	NO	HUMAN	EITHER	ANIM
	10	NONE	DEF	NO	SING	NO	COMPANY	NEUTER	INANIM
	11	NONE	NONE	NO	SING	NO	HUMAN	EITHER	ANIM



# Clustering

Distance Metric

$$\text{dist}(NP_1, NP_2) = \sum_{f \in F} w_f \cdot \text{incompatibility}_f(NP_1, NP_2)$$

Feature $f$	Weight	Incompatibility function
Words	10.0	(# of mismatching words <sup>a</sup> ) / (# of words in the longer NP)
Head Noun	1.0	1 if the head nouns differ; else 0
Position	5.0	(difference in position) / (maximum difference in document)
Pronoun	$r$	1 if $NP_i$ is a pronoun and $NP_j$ is not; else 0
Article	$r$	1 if $NP_j$ is indefinite and not appositive; else 0
Words-Substring	$-\infty$	1 if $NP_i$ subsumes (entirely includes as a substring) $NP_j$ ;
Appositive	$-\infty$	1 if $NP_j$ is appositive and $NP_i$ is its immediate predecessor; else 0
Number	$\infty$	1 if they do not match in number; else 0
Proper Name	$\infty$	1 if both are proper names, but mismatch on every word; else 0
Semantic Class	$\infty$	1 if they do not match in class; else 0
Gender	$\infty$	1 if they do not match in gender (allows EITHER to match MASC or FEM); else 0
Animacy	$\infty$	1 if they do not match in animacy; else 0

compatible classes, compute transitive closure

# Pairwise Model: Features matter! [Bengston & Roth, 2008]

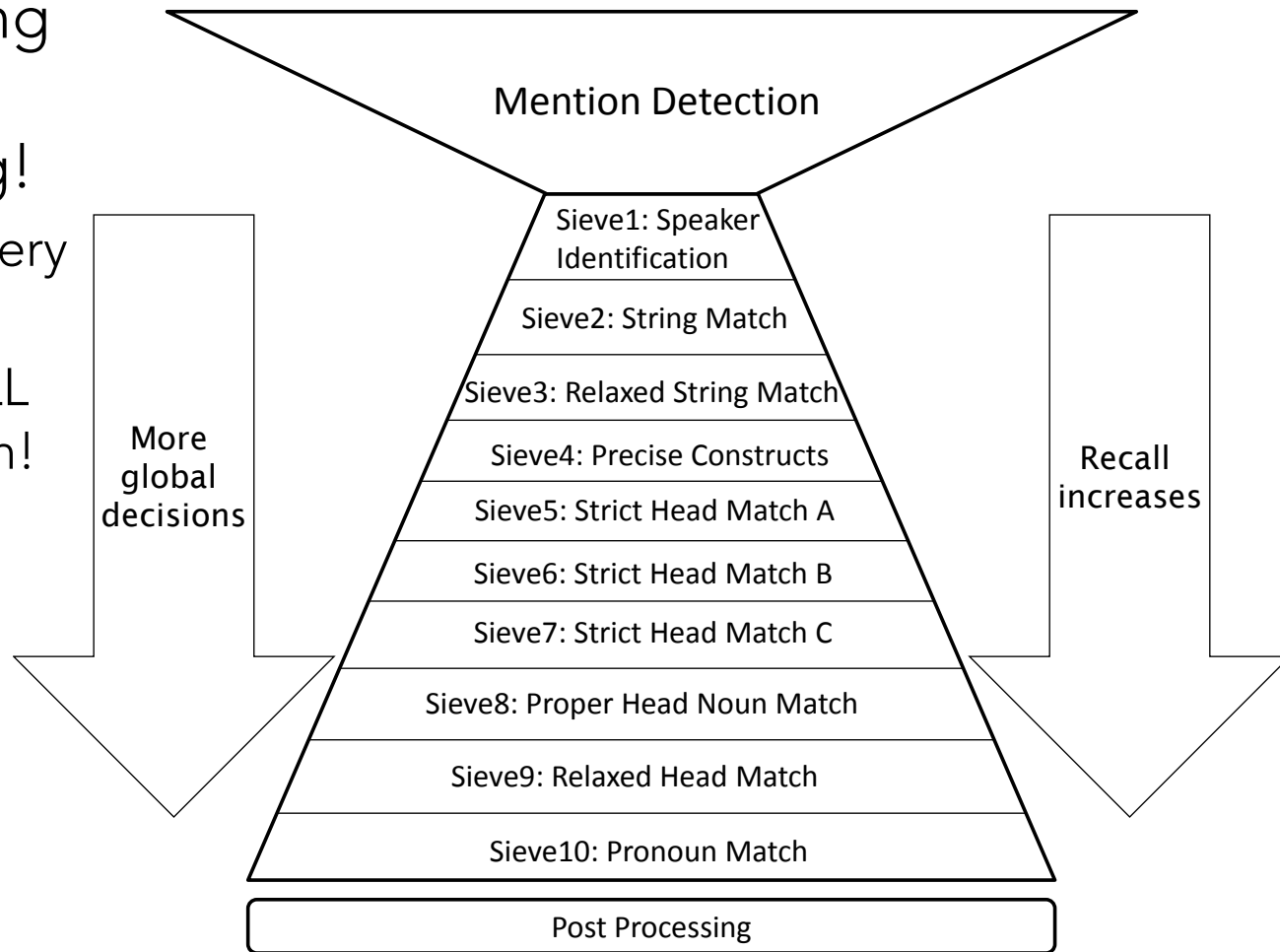
Category	Feature	Source
Mention Types	Mention Type Pair	Annotation and tokens
String Relations	Head Match	Tokens
	Extent Match	Tokens
	Substring	Tokens
	Modifiers Match	Tokens
	Alias	Tokens and lists
Semantic	Gender Match	WordNet and lists
	Number Match	WordNet and lists
	Synonyms	WordNet
	Antonyms	WordNet
	Hypernyms	WordNet
	Both Speak	Context
Relative Location	Apposition	Positions and context
	Relative Pronoun	Positions and tokens
	Distances	Positions
Learned	Anaphoricity	Learned
	Name Modifiers Predicted Match	Learned
Aligned Modifiers	Aligned Modifiers Relation	WordNet and lists
Memorization	Last Words	Tokens
Predicted Entity Types	Entity Types Match	Annotation and tokens
	Entity Type Pair	WordNet and tokens

# Two Recent Supervised Learners

- Linear Model
  - [Bengston & Roth 2008]
  - Pairwise classification
  - Careful experimental setup with tons of features!
  - 80.8 B<sup>3</sup> F1
- FOL-based approach
  - [Culotta et al. 2007]
  - Includes global constraints on clusters
  - 79.3 B<sup>3</sup> F1

# Multi-pass Sieve

- Basically, a ranking model with no machine learning!
  - 10 sieves, each very simple
  - Winner of CONLL 2011 competition!



# A Carefully Constructed Example

Input:	John is a musician. He played a new song. A girl was listening to the song. "It is my favorite," John said to her.
Mention Detection:	[John] <sub>1</sub> <sup>1</sup> is [a musician] <sub>2</sub> <sup>2</sup> . [He] <sub>3</sub> <sup>3</sup> played [a new song] <sub>4</sub> <sup>4</sup> . [A girl] <sub>5</sub> <sup>5</sup> was listening to [the song] <sub>6</sub> <sup>6</sup> . "[It] <sub>7</sub> <sup>7</sup> is [[my] <sub>9</sub> <sup>9</sup> favorite] <sub>8</sub> <sup>8</sup> ," [John] <sub>10</sub> <sup>10</sup> said to [her] <sub>11</sub> <sup>11</sup> .
Speaker Sieve:	[John] <sub>1</sub> <sup>1</sup> is [a musician] <sub>2</sub> <sup>2</sup> . [He] <sub>3</sub> <sup>3</sup> played [a new song] <sub>4</sub> <sup>4</sup> . [A girl] <sub>5</sub> <sup>5</sup> was listening to [the song] <sub>6</sub> <sup>6</sup> . "[It] <sub>7</sub> <sup>7</sup> is [[my] <sub>9</sub> <sup>9</sup> favorite] <sub>8</sub> <sup>8</sup> ," [ <b>John</b> ] <sub>10</sub> <sup>9</sup> said to [her] <sub>11</sub> <sup>11</sup> .
String Match:	<b>[John]</b> <sub>1</sub> <sup>1</sup> is [a musician] <sub>2</sub> <sup>2</sup> . [He] <sub>3</sub> <sup>3</sup> played [a new song] <sub>4</sub> <sup>4</sup> . [A girl] <sub>5</sub> <sup>5</sup> was listening to [the song] <sub>6</sub> <sup>6</sup> . "[It] <sub>7</sub> <sup>7</sup> is [[my] <sub>9</sub> <sup>1</sup> favorite] <sub>8</sub> <sup>8</sup> ," [ <b>John</b> ] <sub>10</sub> <sup>1</sup> said to [her] <sub>11</sub> <sup>11</sup> .
Relaxed String Match:	[John] <sub>1</sub> <sup>1</sup> is [a musician] <sub>2</sub> <sup>2</sup> . [He] <sub>3</sub> <sup>3</sup> played [a new song] <sub>4</sub> <sup>4</sup> . [A girl] <sub>5</sub> <sup>5</sup> was listening to [the song] <sub>6</sub> <sup>6</sup> . "[It] <sub>7</sub> <sup>7</sup> is [[my] <sub>9</sub> <sup>1</sup> favorite] <sub>8</sub> <sup>8</sup> ," [John] <sub>10</sub> <sup>1</sup> said to [her] <sub>11</sub> <sup>11</sup> .
Precise Constructs:	<b>[John]</b> <sub>1</sub> <sup>1</sup> is <b>[a musician]</b> <sub>2</sub> <sup>1</sup> . [He] <sub>3</sub> <sup>3</sup> played [a new song] <sub>4</sub> <sup>4</sup> . [A girl] <sub>5</sub> <sup>5</sup> was listening to [the song] <sub>6</sub> <sup>6</sup> . " <b>[It]</b> <sub>7</sub> <sup>7</sup> is <b>[[my]<sub>9</sub><sup>1</sup> favorite]</b> <sub>8</sub> <sup>7</sup> ," [John] <sub>10</sub> <sup>1</sup> said to [her] <sub>11</sub> <sup>11</sup> .
Strict Head Match A:	[John] <sub>1</sub> <sup>1</sup> is [a musician] <sub>2</sub> <sup>1</sup> . [He] <sub>3</sub> <sup>3</sup> played <b>[a new song]</b> <sub>4</sub> <sup>4</sup> . [A girl] <sub>5</sub> <sup>5</sup> was listening to <b>[the song]</b> <sub>6</sub> <sup>4</sup> . "[It] <sub>7</sub> <sup>7</sup> is [[my] <sub>9</sub> <sup>1</sup> favorite] <sub>8</sub> <sup>7</sup> ," [John] <sub>10</sub> <sup>1</sup> said to [her] <sub>11</sub> <sup>11</sup> .
Strict Head Match B,C:	[John] <sub>1</sub> <sup>1</sup> is [a musician] <sub>2</sub> <sup>1</sup> . [He] <sub>3</sub> <sup>3</sup> played [a new song] <sub>4</sub> <sup>4</sup> . [A girl] <sub>5</sub> <sup>5</sup> was listening to [the song] <sub>6</sub> <sup>4</sup> . "[It] <sub>7</sub> <sup>7</sup> is [[my] <sub>9</sub> <sup>1</sup> favorite] <sub>8</sub> <sup>7</sup> ," [John] <sub>10</sub> <sup>1</sup> said to [her] <sub>11</sub> <sup>11</sup> .
Proper Head Noun Match:	[John] <sub>1</sub> <sup>1</sup> is [a musician] <sub>2</sub> <sup>1</sup> . [He] <sub>3</sub> <sup>3</sup> played [a new song] <sub>4</sub> <sup>4</sup> . [A girl] <sub>5</sub> <sup>5</sup> was listening to [the song] <sub>6</sub> <sup>4</sup> . "[It] <sub>7</sub> <sup>7</sup> is [[my] <sub>9</sub> <sup>1</sup> favorite] <sub>8</sub> <sup>7</sup> ," [John] <sub>10</sub> <sup>1</sup> said to [her] <sub>11</sub> <sup>11</sup> .
Relaxed Head Match:	[John] <sub>1</sub> <sup>1</sup> is [a musician] <sub>2</sub> <sup>1</sup> . [He] <sub>3</sub> <sup>3</sup> played [a new song] <sub>4</sub> <sup>4</sup> . [A girl] <sub>5</sub> <sup>5</sup> was listening to [the song] <sub>6</sub> <sup>4</sup> . "[It] <sub>7</sub> <sup>7</sup> is [[my] <sub>9</sub> <sup>1</sup> favorite] <sub>8</sub> <sup>7</sup> ," [John] <sub>10</sub> <sup>1</sup> said to [her] <sub>11</sub> <sup>11</sup> .
Pronoun Match:	<b>[John]</b> <sub>1</sub> <sup>1</sup> is [a musician] <sub>2</sub> <sup>1</sup> . <b>[He]</b> <sub>3</sub> <sup>1</sup> played [a new song] <sub>4</sub> <sup>4</sup> . <b>[A girl]</b> <sub>5</sub> <sup>5</sup> was listening to <b>[the song]</b> <sub>6</sub> <sup>4</sup> . " <b>[It]</b> <sub>7</sub> <sup>4</sup> is [[my] <sub>9</sub> <sup>1</sup> favorite] <sub>8</sub> <sup>4</sup> ," [John] <sub>10</sub> <sup>1</sup> said to <b>[her]</b> <sub>11</sub> <sup>5</sup> .
Post Processing:	[John] <sub>1</sub> <sup>1</sup> is <b>a musician</b> . [He] <sub>3</sub> <sup>1</sup> played [a new song] <sub>4</sub> <sup>4</sup> . [A girl] <sub>5</sub> <sup>5</sup> was listening to [the song] <sub>6</sub> <sup>4</sup> . "[It] <sub>7</sub> <sup>4</sup> is <b>[my]</b> <sub>9</sub> <sup>1</sup> <b>favorite</b> ," [John] <sub>10</sub> <sup>1</sup> said to [her] <sub>11</sub> <sup>5</sup> .
Final Output:	[John] <sub>1</sub> <sup>1</sup> is a musician. [He] <sub>3</sub> <sup>1</sup> played [a new song] <sub>4</sub> <sup>4</sup> . [A girl] <sub>5</sub> <sup>5</sup> was listening to [the song] <sub>6</sub> <sup>4</sup> . "[It] <sub>7</sub> <sup>4</sup> is [my] <sub>9</sub> <sup>1</sup> favorite," [John] <sub>10</sub> <sup>1</sup> said to [her] <sub>11</sub> <sup>5</sup> .

Table 1

# The Most Useful Sieves

- **2: Exact string match** -- e.g., *[the Shahab 3 ground-ground missile]* and *[the Shahab 3 ground-ground missile]*. Precision is over 90% B3 **[+16 F1]**
- **5: Entity head match** – The mention head word matches *any* head word of mentions in the antecedent entity. Also, looks at modifiers, e.g. to separate *Harvard University* and *Yale University*. **[+3 F1]**
- **10: Pronominal Coreference Resolution** – observe constraints on number, gender, person, animacy, and NER types. Link to closest, with a maximum distance. **[+10 F1]**
- Most others get between 0-2 points improvement, but are cumulative

# Some Results

System	MUC			B <sup>3</sup>		
	R	P	F1	R	P	F1
<b>ACE2004-Culotta-Test</b>						
This paper	70.2	82.7	75.9	74.5	88.7	81.0
Haghighi and Klein (2009)	77.7	74.8	79.6	78.5	79.6	79.0
Culotta et al. (2007)	–	–	–	73.2	86.7	79.3
Bengston and Roth (2008)	69.9	82.7	75.8	74.5	88.3	80.8
<b>ACE2004-nwire</b>						
This paper	75.1	84.6	79.6	74.1	87.3	80.2
Haghighi and Klein (2009)	75.9	77.0	76.5	74.5	79.4	76.9
Poon and Domingos (2008)	70.5	71.3	70.9	–	–	–
Finkel and Manning (2008)	58.5	78.7	67.1	65.2	86.8	74.5
<b>MUC6-Test</b>						
This paper	69.1	90.6	78.4	63.1	90.6	74.4
Haghighi and Klein (2009)	77.3	87.2	81.9	67.3	84.7	75.0
Poon and Domingos (2008)	75.8	83.0	79.2	–	–	–
Finkel and Manning (2008)	55.1	89.7	68.3	49.7	90.9	64.3

**Table 5**

Comparison of our system with the other reported results on the ACE and MUC corpora. All these systems use gold mention boundaries.

[Lee et al, 2013]

# Back to ... Textual 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.
  - He arrived just as the store was closing for the day.
- 
- 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.
- ➔ Same content, different realization through different syntactic choices



# Centering Theory

(Grosz et al.,1983)

- 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.
- He arrived just as the store was closing for the day.

❑ Focus is the most salient entity in a discourse segment

❑ Constraints on linguistic realization of focus

- Focus is more likely to be realized as subject or object
- Focus is more likely to be referred to with anaphoric expression

❑ Constraints on the entity distribution in a coherent text

- Transition between adjacent sentences is characterized in terms of focus switch

# Entity-grid Model

1. [Former Chilean dictator Augusto Pinochet]**S**, was arrested in [London]**X** on [October 14th]**X** 1998.
2. [Pinochet]**S**, 82, was recovering from [surgery]**X**.
3. [The arrest]**S** was in [response]**X** to [an extradition warrant]**X** served by [a Spanish judge]**S**.
4. [Pinochet]**S** was charged with murdering [thousands]**O**, including many [Spaniards]**O**.
5. [He]**S** is awaiting [a hearing]**O**, [his fate]**X** in [the balance]**X**.
6. [American scholars]**S** applauded the [arrest]**O**.

Notation: **S**=subjects, **O**=object, **X**=other

# Entity-grid Model

1. [Former Chilean dictator Augusto Pinochet]**S**, was arrested in [London]**X** on [October 14]**X** 1998.
2. [Pinochet]**S**, 82, was recovering from [surgery]**X**.
3. [The arrest]**S** was in [response]**X** to [an extradition warrant]**X** served by [a Spanish judge]**S**.
4. [Pinochet]**S** was charged with murdering [thousands]**O**, including many [Spaniards]**O**.
5. [He]**S** is awaiting [a hearing]**O**, [his fate]**X** in [the balance]**X**.
6. [American scholars]**S** applauded the [arrest]**O**.

# Entity-grid Model

[illegible]

# Comparing Grids

S	S	S	X	X	-	-	-	-	-	-	-	-	-	-
-	-	S	-	-	X	-	-	-	-	-	-	-	-	-
-	-	-	-	-	-	S	X	X	O	-	-	-	-	-
-	-	S	-	-	-	-	-	-	-	O	O	-	-	-
-	-	S	-	-	-	-	-	-	-	-	-	O	X	X
-	-	-	-	-	-	O	-	-	-	-	-	-	-	-

S	S	X	X	X	-	-	-	-	-	-	-	-	-	X
-	-	X	-	-	X	-	-	-	-	-	-	-	-	X
-	-	X	-	-	-	-	X	X	O	-	-	-	-	X
-	-	X	-	-	-	-	-	-	-	O	O	-	-	X
-	-	X	-	-	-	-	-	-	-	-	-	O	X	X
-	-	X	-	-	-	O	-	-	-	-	-	-	-	X

# Quantifying Textual Coherence

- Text is encoded as **a distribution over entity transition types**
- Entity transition type —  $\{s, o, x, -\}^n$

	s	o	x	-	s	o	x	-	s	o	x	-	s	o	x	-
	s	s	s	s	o	o	o	o	x	x	x	x	-	-	-	-
$d_{i1}$	0	0	0	.03	0	0	0	.02	.07	0	0	.12	.02	.02	.05	.25
$d_{i2}$	.02	0	0	.03	0	0	0	.06	0	0	0	.05	.03	.07	.07	.29

**How to select relevant transition types?:**

- Use all the unigrams, bigrams, ... over  $\{s, o, x, -\}$
- Do feature selection

# Evaluation / Applications

Goal: recover the most coherent sentence ordering

Basic set-up:

- Input: a pair of a source document and a permutation of its sentences
- Task: find a source document via coherence ranking

Data: Training 4000 pairs, Testing 4000 pairs (Natural disasters and Transportation Safety Reports)

# Conclusion

- Computational modeling of discourse coherence
- **Theories:** Rhetorical structure theory / Centering theory
- **Corpus:** Penn Discourse Tree Bank
- **Applications:**
  - better summarization
  - automatic ESL grading
  - better QA (sharp et al., NAACL 2010)
  - better machine translation (this workshop!)
- Coreference Relations
- Entity-grid Models

