

#### COS 484: Natural Language Processing

## **Coreference Resolution**

Fall 2019

(Slides adapted from Christopher Manning)

#### Announcements

- No lecture on **Dec 12th**.
- All the project meetings scheduled on **Dec 10th** 
  - Submit checkpoint reports by Dec 9th (optional but recommended!)
  - Sign up project meetings!

## Overview

- What is coreference resolution?
- Four types of coreference models
- End-to-end coreference resolution
- Evaluation and results

## Discourse

- Text is more than the sum of its individual sentences/utterances.
- Discourse processing: NLP beyond the sentence/utterance boundary.
  - Monologue
  - Dialogue (multi-party)
- Topics:
  - Topic segmentation
  - Coreference resolution
  - Coherence
  - Discourse relations
  - ...

• Identify all mentions that refer to the same real world entity

• Identify all mentions that refer to the same real world entity

• Identify all mentions that refer to the same real world entity



• Identify all mentions that refer to the same real world entity



• Identify all mentions that refer to the same real world entity

InputBarack Obama nominated Hillary RodhamClinton as his secretary of state on Monday. Hechose her because she had foreign affairsexperience as a former First Lady.

#### Output

- {Barack Obama, his, He}
- {Hillary Rodham Clinton, secretary of state, her, she, First Lady}

# Try coref systems yourself!

#### • https://corenlp.run

- Text to annotate -Barack Obama nominated Hillary Rodham Clinton as his secretary of state on Monday. He chose her because she had foreign affairs experience as a former First Lady. - Annotations - Language – Submit English coreference × Ŧ Coreference: Mention Mentio Mention 1 Barack Obama nominated Hillary Rodham Clinton as his secretary of state on Monday . -----coref-------coref--- Mention He chose her because she had foreign affairs experience as a former First Lady .

# Try coref systems yourself!

#### <u>https://demo.allennlp.org/coreference-resolution/</u>

#### Document

Paul Allen was born on January 21, 1953, in Seattle, Washington, to Kenneth Sam Allen and Edna Faye Allen. Allen attended Lakeside School, a private school in Seattle, where he befriended Bill Gates, two years younger, with whom he shared an enthusiasm for computers. Paul and Bill used a teletype terminal at their high school, Lakeside, to develop their programming skills on several time-sharing computer systems.

Run >
<ul> <li>Paul Allen was born on January 21, 1953, in 1 Seattle, Washington, to Kenneth Sam Allen and Edna Faye Allen.</li> <li>Allen attended Lakeside School, a private school in 1 Seattle, where <ul> <li>he befriended</li> <li>Bill Gates</li> <li>two years younger, with whom <ul> <li>he shared an enthusiasm for computers.</li> <li>Paul and <ul> <li>Bill</li> <li>Bill</li> <li>Bill</li> <li>Bill</li> <li>Bill</li> <li>Bill</li> <li>Bill</li> <li>Bill</li> </ul> </li> </ul></li></ul></li></ul>
Document
Barack Obama nominated Hillary Rodham Clinton as his secretary of state on Monday. He chose her because she had foreign affairs experience as a former First Lady.
Run >
Barack Obama nominated • Hillary Rodham Clinton as • his secretary of state on Monday. • He chose • her because • she had foreign affairs experience as a former First Lady.

# Applications

- Full text understanding
- Information extraction, question answering, summarization..

#### Barack Obama

From Wikipedia, the free encyclopedia

Obama was born in Honolulu, Hawaii. After graduating from Columbia University in 1983, he worked as a community organizer in Chicago. In 1988, **he** enrolled in Harvard Law School, where he was the first black president of the *Harvard Law Review*. After graduating, he became a civil rights attorney and an academic, teaching constitutional law at the University of Chicago Law School from 1992 to 2004.

- he = Barack Obama
- (Barack Obama, schools\_attended, Harvard University)

# Applications

• Languages have different features for gender, number, dropped pronouns, etc.

☆A Text     ■ Documents			
SPANISH - DETECTED ENGLISH SPANISH FRENCH	~ +	- ENGLISH SPANISH ARABIC V	
"Me incanta el conocimiento", dice.	×	"I love knowledge," he says.	\$
Showing translation for "Me <i>encanta</i> el conocimiento", dice. Translate instead "Me incanta el conocimiento", dice.			
	35/5000 🧪	•	

- "I love knowledge," he said.
- "I love knowledge," she said.



#### **Coreference vs Entity Linking**



en.wikipedia.org/wiki/Freddie\_Mac



PERSON

ORGANIZATION PERSON

<u>Donald Layton<sub>1</sub></u> took the helm of <u>Freddie Mac<sub>2</sub></u> after <u>his<sub>1</sub></u> ...

(Image credit: Greg Durrett)

# Winograd Schema problems

• Some cases of coreference require world knowledge or common sense reasoning to solve.

They city council denied the demonstrators a permit because they feared violence.

They city council denied the demonstrators a permit because they advocated violence.

The trophy didn't fit into the suitcase because it was too large.

The trophy didn't fit into the suitcase because it was too small.



Winograd 1972

# Winograd Schema problems

The Winograd Schema Challenge: recently proposed as an alternative to the Turing test

- (Levesque 2013): "On our best behavior"
- If you've fully solved coreference, arguably you've solved AI



http://commonsensereasoning.org/winograd.html

- **Coreference** is when two mentions refer to the same entity in the world
  - Barack Obama traveled to.. Obama...
- A related linguistic concept is **anaphora**: when a term (anaphor) refers to another term (antecedent)
  - The interpretation of the anaphor is in some way determined by the interpretation of the antecedent
  - Barack Obama said he would sign the bill.
     antecedent anaphor

• Coreference with named entities



• Anaphora



• Not all anaphoric relations are coreferential



Anaphora vs cataphora

- Usually the antecedent comes before the anaphor (e.g., a pronoun) but not always
- Cataphora: mentioned before the their referents

Even before she saw it, Dorothy had been thinking about the Emerald City every day.

More linguistic background in J & M: Chapter 22.1

#### Coreference resolution in two steps

• Mention detection (easy)

"[I] voted for [Nader] because [he] was most aligned with [[my] values]," [she] said

Mentions can be nested!

• Mention clustering (hard)

"[I] voted for [Nader] because [he] was most aligned with [[my] values]," [she] said

## Mention detection

- Mention: span of text referring to some entity
- Three kinds of mentions:
  - #1: Pronouns
    - I, your, it, she, him
  - #2: Named entities
    - "Barack Obama" "Princeton University"
  - #3: Noun phrases
    - "a dog", "the big fluffy cat stuck in the tree"

Part-of-speech tagger

named entity recognizer

constituency parer

### Mention detection: not so simple

- Making all pronouns, named entities, and noun phrases as mentions over-generates mentions
- Are these mentions?
  - It has been ten years
  - Every student
  - No student
  - 100 miles

## Mention detection: not so simple

- We can train a classifier to filter out spurious mentions
- Much more common: keep all mentions as "candidate mentions"
  - After your coreference system is done, discard all singleton mentions (ones that have not been marked as coreference with anything else)
- Recent work: it is possibly to do mention detection and clustering end-to-end instead of in two steps
  - We will cover this later!

#### Four kinds of coreference models

- Rule-based (pronominal anaphora resolution)
- Machine learning models:
  - Mention pair
  - Mention ranking 🛩
  - Clustering-based

"I voted for Nader because he was most aligned with my values," she said.

Current state-of-the-art



# Hobbs algorithm (1978)



- An algorithm that searches parse trees for antecedents of a pronoun
  - Starting at the NP node immediately dominating the pronoun
  - Search previous trees, in order of recency, left-to-right, breadth-first
  - Looking for the first match of the correct gender and number (male vs female, singular vs plural)
- 72.7% accuracy based on the 100 examples he analyzed

# Hobbs algorithm (1978)



The castle in Camelot remained the residence of the king until he moved it to London.



• Train a binary classifier that assigns every pair of mentions a probability of being coreferent:  $p(m_i, m_j)$ 

"I voted for Nader because he was most aligned with my values," she said.



"I voted for Nader because he was most aligned with my values," she said.



Positive examples: want  $p(m_i, m_j)$  to be near 1

"I voted for Nader because he was most aligned with my values," she said.



Negative examples: want  $p(m_i, m_j)$  to be near 0



#### **Testing time:**

- Run your own mention detector
- Pick some threshold (e.g., 0.5) and add coreference links
- Take the transitive closure to get the clustering



**Disadvantage:** suppose we have a long document with the following mentions: Ralph Nader... he... his.. him.. <several paragraphs>... voted for Nader because he..



Many mentions only have one clear antecedent but we are asking the model to predict all of them.

**Solution**: instead train the model to predict only one antecedent for each mention (more linguistically plausible)

# Mention-ranking models

- Assign each mention its highest scoring candidate antecedent according to the model
- Add a dummy NA mention to decline linking the current mention to anything ("singleton" or "first" mention)



## Mention-ranking models

• Training time: only clustering information is observed (no annotation of "antecedent"), so we optimize the marginal log-likelihood of all the correct antecedents.



• Testing time: same as mention-pair but we only pick one antecedent for each mention

- Coreference is a clustering task, so let's use a clustering algorithm directly.
- Start with each mention in its own singleton cluster
- Merge a pair of clusters at each step and use a model to score which cluster merges are good.

Google recently ... the company announced Google Plus ... the product features ...



Mention-pair decision is difficult



#### Cluster-pair decision is easier



- Current candidate cluster merges depend on previous ones it already made
  - We can't use supervised learning.
  - Researchers used reinforcement learning to train the model
  - Reward function is the change in a coreference evaluation metric
- Some nice attempts but NOT the dominating approach

## Features

- Non-neural coreference models:
  - Person/Number/Gender agreement
    - Jack gave Mary a gift. She was excited.
  - Semantic compatibility
    - ... the mining conglomerate ... the company ...
  - Certain syntactic constraints
    - John bought him a new car. [him can not be John]
  - More recently mentioned entities preferred for referenced
    - John went to a movie. Jack went as well. He was not busy.
  - Grammatical Role: Prefer entities in the subject position
    - John went to a movie with Jack. He was not busy.
  - Parallelism:
    - John went with Jack to a movie. Joe went with him to a bar.

### Features

• Neural coreference models:



- A mention-ranking model
- Joint mention detection and clustering so you don't need an additional mention detector (parser/part-of-speech tagger)

$$J = \sum_{i=2}^{N} - \log \left( \sum_{j=1}^{i-1} \mathbbm{1}(y_{ij} = 1)p(m_j, m_i) \right)$$
  
Iterate over all the mentions  
in the document Usual trick of taking negative  
log to go from likelihood to loss

We consider all the possible spans + {NA}  

$$M = \frac{T(T+1)}{2} + 1$$

$$p(m_j, m_i) = \frac{\exp(s(m_j, m_i))}{\sum_{j' < i} \exp(m_{j'}, m_i)}$$
T: number of words

#### (Lee et al, 2017): End-to-end Neural Coreference Resolution



Let's compute a vector representation  $\mathbf{g}_i \in \mathbb{R}^d$  for each span *i* 

 $\phi(i, j)$ : manual features such speaker/gender information



#### (Lee et al, 2017): End-to-end Neural Coreference Resolution

Attention scores

 $\alpha_t = \boldsymbol{w}_{lpha} \cdot \text{FFNN}_{lpha}(\boldsymbol{x}_t^*)$ 

dot product of weight vector and transformed hidden state

span



just a softmax over attention scores for the span

**Final representation** 

$$\hat{\boldsymbol{x}}_i = \sum_{t= ext{start}(i)}^{ ext{end}(i)} a_{i,t} \cdot \boldsymbol{x}_t$$

Attention-weighted sum of word embeddings



(Lee et al, 2017): End-to-end Neural Coreference Resolution

#### **Testing time:**

- $O(T^2)$  spans of text in a document
- $O(T^4)$  possible pairs of spans too expensive!
- Use  $s_m(i)$  for aggressive pruning

## Coreference evaluation

- You need to get both "mentions" and "clusters" correctly.
- Standard practice: we use 3 types of metrics
  - B<sup>3</sup>: mention-based
  - MUC: link-based (pair of mentions)
  - CEAF: entity-based
  - .. and finally take the average of these 3 F1 scores

## **B-cubed evaluation metric**

- For each mention in the reference chain, compute a precision and a recall (e.g., *#* of mentions in the same reference chain with the current mention)
- The final precision/recall is an average of all the mentions



## System performance

- Evaluation on English Ontonotes (CoNLL-2012 Shared Task)
- #Train: 2,802 / #Dev: 343 / #Test: 348 documents

	MUC			$\mathbf{B}^3$			$\text{CEAF}_{\phi_A}$			
	Р	R	F1	Р	R	F1	Р	R	F1	Avg. F1
Lee et al. (2017) (single model)	78.4	73.4	75.8	68.6	61.8	65.0	62.7	59.0	60.8	67.2
Clark and Manning (2016a)	79.2	70.4	74.6	69.9	58.0	63.4	63.5	55.5	59.2	65.7
Clark and Manning (2016b)	79.9	69.3	74.2	71.0	56.5	63.0	63.8	54.3	58.7	65.3
Wiseman et al. (2016)	77.5	69.8	73.4	66.8	57.0	61.5	62.1	53.9	57.7	64.2
Wiseman et al. (2015)	76.2	69.3	72.6	66.2	55.8	60.5	59.4	54.9	57.1	63.4
Fernandes et al. (2014)	75.9	65.8	70.5	77.7	65.8	71.2	43.2	55.0	48.4	63.4
Clark and Manning (2015)	76.1	69.4	72.6	65.6	56.0	60.4	59.4	53.0	56.0	63.0
Martschat and Strube (2015)	76.7	68.1	72.2	66.1	54.2	59.6	59.5	52.3	55.7	62.5
Durrett and Klein (2014)	72.6	69.9	71.2	61.2	56.4	58.7	56.2	54.2	55.2	61.7
Björkelund and Kuhn (2014)	74.3	67.5	70.7	62.7	55.0	58.6	59.4	52.3	55.6	61.6
Durrett and Klein (2013)	72.9	65.9	69.2	63.6	52.5	57.5	54.3	54.4	54.3	60.3

#### "Eash Victories and Uphill Battles in Coreference Resolution"

#### System performance

Avg. F1 83 79.6 77 77.1 73.0 72 71.0 67.2 66 60 (Lee et al, 2017) SpanBERT (Lee et al, 2018) 00 (Joshi & Chen, 2019)