### Graph-Cut-Based Anaphoricity Determination for Coreference Resolution

Vincent Ng
Human Language Technology Research Institute
University of Texas at Dallas

Identify the noun phrases (NPs) that refer to the same entity

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

Identify the noun phrases (NPs) that refer to the same entity

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

Identify the noun phrases (NPs) that refer to the same entity

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

Identify the noun phrases (NPs) that refer to the same entity

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

Identify the noun phrases (NPs) that refer to the same entity

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

### Step 1: Classification

- train a coreference model, P<sub>c</sub>, to determine the probability that two NPs are coreferent
- two NPs are classified as coreferent iff probability ≥ 0.5

#### Step 1: Classification

- train a coreference model, P<sub>c</sub>, to determine the probability that two NPs are coreferent
- two NPs are classified as coreferent iff probability ≥ 0.5

#### Step 2: Antecedent selection

find an antecedent for each NP<sub>i</sub>

#### Step 1: Classification

- train a coreference model, P<sub>c</sub>, to determine the probability that two NPs are coreferent
- two NPs are classified as coreferent iff probability ≥ 0.5

#### Step 2: Antecedent selection

- find an antecedent for each NP<sub>i</sub>
  - choose the closest preceding noun phrase that is classified as coreferent with NP<sub>i</sub>

#### Step 1: Classification

- train a coreference model, P<sub>c</sub>, to determine the probability that two NPs are coreferent
- two NPs are classified as coreferent iff probability ≥ 0.5

#### Step 2: Antecedent selection

- find an antecedent for each NP<sub>i</sub>
  - choose the closest preceding noun phrase that is classified as coreferent with NP<sub>i</sub>

coref

[Queen Elizabeth] set about transforming [her] [husband], ...

#### Step 1: Classification

- train a coreference model, P<sub>c</sub>, to determine the probability that two NPs are coreferent
- two NPs are classified as coreferent iff probability ≥ 0.5

#### Step 2: Antecedent selection

- find an antecedent for each NP<sub>j</sub>
  - choose the closest preceding noun phrase that is classified as coreferent with NP<sub>i</sub>

coref

[Queen Elizabeth] set about transforming [her] [husband], ...

- Step 1: Classification
  - train a coreference model two NPs are coreferent

find an antecedent for each NPj that has an antecedent

at

- two NPs are classified as coreferent of probability ≥ 0.5
- Step 2: Antecedent selection
  - find an antecedent for each NP<sub>i</sub>
    - choose the closest preceding noun phrase that is classified as coreferent with NP<sub>i</sub>

coref

[Queen Elizabeth] set about transforming [her] [husband], ...

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

Queen Elizabeth set about transforming her husband,

King George VI, into a viable monarch.

A renowned speech therapist, was summoned to help

Queen Elizabeth set about transforming her husband,

King George VI, into a viable monarch.

A renowned speech therapist, was summoned to help

Queen Elizabeth set about transforming her husband,

King George VI, into a viable monarch.

A renowned speech therapist was summoned to help

Queen Elizabeth set about transforming her busband,

King George VI, into a viable monarch.

A renowned speech therapist was summoned to help

Queen Elizabeth set about transforming her busband,

King George VI, into a viable monarch.

A renowned speech therapist was summoned to help

Any NP that is part of a coref chain but is not the head of the chain has an antecedent. It's an anaphoric NP.

Queen Elizabeth set about transforming her husband,

King George VI, into a viable monarch.

A renowned speech therapist was summoned to help

## **Anaphoricity Determination**

- determines whether an NP is anaphoric or not
- helps improve the precision of a coreference system

### Goal

Improve learning-based coreference systems using automatically acquired anaphoricity information, by proposing a new approach to anaphoricity determination

### Plan for the Talk

- Existing methods for computing and using anaphoricity info
- Our graph-cut-based approach to anaphoricity determination
- Evaluation

# Methods for Computing and Using Anaphoricity Information

- Five existing methods
  - Ng & Cardie (2002)
  - Ng (2004)
  - Luo (2007)
  - Denis & Baldridge (2007)
  - Kleener (2007), Finkel & Manning (2008)

## What do the methods have in common?

### What do the methods have in common?

- Training an anaphoricity model (P<sub>A</sub>)
  - determines the probability that an NP is anaphoric
  - classifies an NP as anaphoric iff probability ≥ 0.5

### What do the methods have in common?

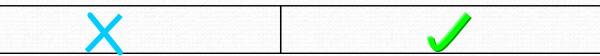
- Training an anaphoricity model (P<sub>A</sub>)
  - determines the probability that an NP is anaphoric
  - classifies an NP as anaphoric iff probability ≥ 0.5
- Training data creation
  - texts annotated with coreference information
  - one instance for each NP
    - positive if the NP is part of a coref chain but not head of chain
    - negative otherwise

- They differ in terms of
  - whether they improve the output of P<sub>A</sub>

- They differ in terms of
  - whether they improve the output of P<sub>A</sub>
    - if so, how?

- They differ in terms of
  - whether they improve the output of P<sub>A</sub>
    - if so, how?
  - how anaphoricity info is used by the coreference system
    - as hard constraints or as soft constraints?

Improve P<sub>A</sub>'s output? Used as hard constraint?

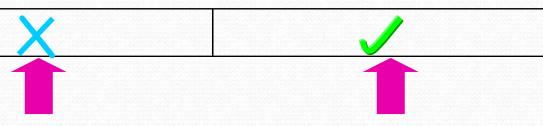


Improve P<sub>A</sub>'s output? Used as hard constraint?



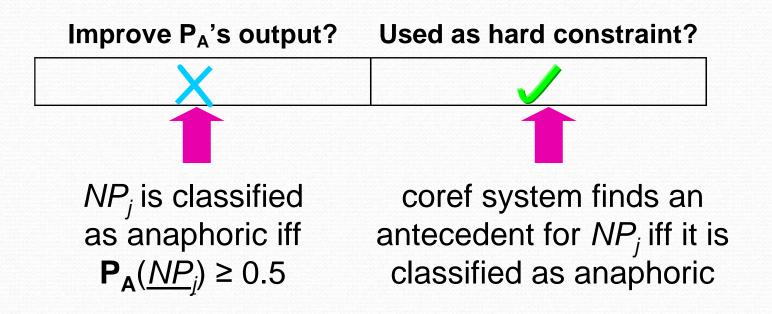
 $NP_j$  is classified as anaphoric iff  $\mathbf{P_A}(\underline{NP_i}) \ge 0.5$ 

Improve  $P_A$ 's output? Used as hard constraint?

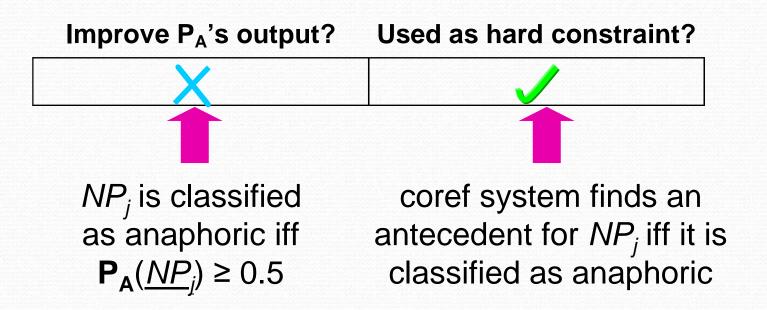


 $NP_j$  is classified as anaphoric iff  $\mathbf{P_A}(\underline{NP_j}) \ge 0.5$ 

coref system finds an antecedent for  $NP_j$  iff it is classified as anaphoric

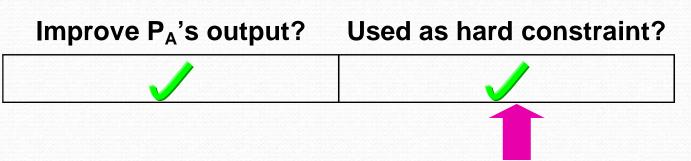


- Problem
  - many anaphoric NPs are misclassified (as non-anaphoric)



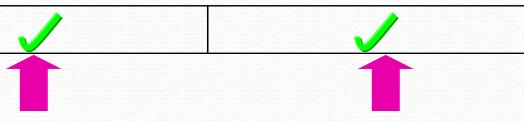
- Problem
  - many anaphoric NPs are misclassified (as non-anaphoric)
    - P<sub>A</sub> is overly conservative in classifying an NP as anaphoric

Improve P<sub>A</sub>'s output? Used as hard constraint?



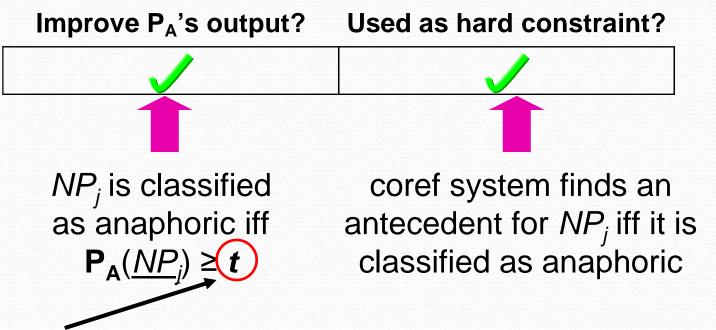
coref system finds an antecedent for  $NP_j$  iff it is classified as anaphoric



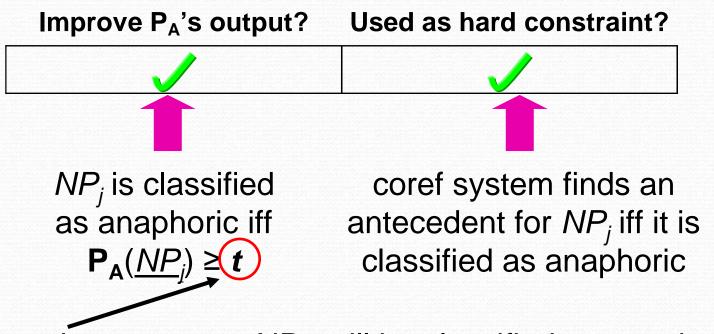


 $NP_j$  is classified as anaphoric iff  $\mathbf{P}_{\mathbf{A}}(\underline{NP_j}) \geq \mathbf{t}$ 

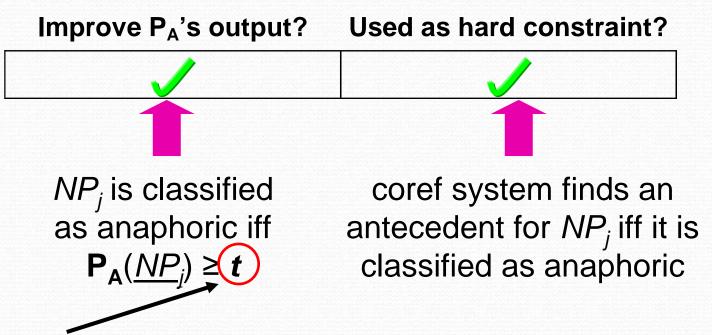
coref system finds an antecedent for  $NP_j$  iff it is classified as anaphoric



- decreasing t \( \mathbb{t} \) more NPs will be classified as anaphoric
- increasing t ± fewer NPs will be classified as anaphoric



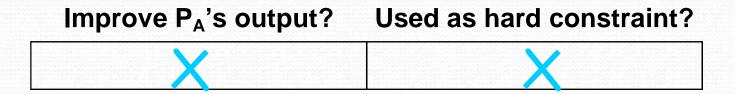
- decreasing t \( \mathbb{t} \) more NPs will be classified as anaphoric
- increasing t ½ fewer NPs will be classified as anaphoric
- t is the "conservativeness" parameter



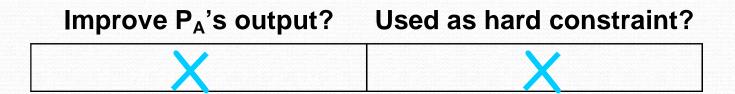
- decreasing t \( \mathbb{L} \) more NPs will be classified as anaphoric
- increasing t ½ fewer NPs will be classified as anaphoric
- select t to use held-out data to maximize coreference performance (i.e., F-measure)

Improve P<sub>A</sub>'s output? Used as hard constraint?

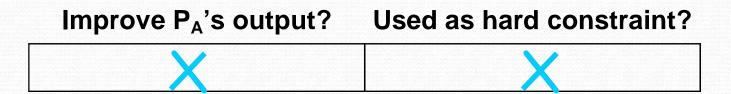




- Goal
  - score an NP partition
    - by multiplying the probabilities provided by  $P_A$  and  $P_C$



- Goal
  - score an NP partition
    - by multiplying the probabilities provided by P<sub>A</sub> and P<sub>C</sub>
  - find highest-scored NP partition



- Goal
  - score an NP partition
    - by multiplying the probabilities provided by P<sub>A</sub> and P<sub>C</sub>
  - find highest-scored NP partition
    - by performing a beam search through the Bell tree

# Denis & Baldridge (2007)

Improve P<sub>A</sub>'s output? Used as hard constraint?

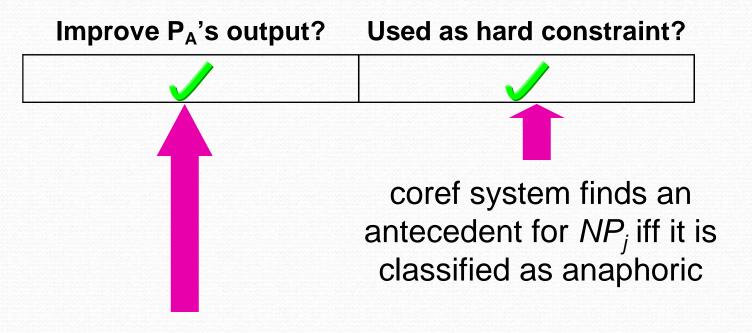
## Denis & Baldridge (2007)

Improve P<sub>A</sub>'s output?

Used as hard constraint?

coref system finds an antecedent for  $NP_j$  iff it is classified as anaphoric

## Denis & Baldridge (2007)



 use Integer Linear Programming (ILP) to perform joint inference for anaphoricity determination and coreference

- 3 NPs: 1, 2, 3
- $P_c(1, 2) = 0.6$ ,  $P_c(1, 3) = 0.2$ ,  $P_c(2, 3) = 0.9$ •  $P_A(1) = 0.1$ ,  $P_A(2) = 0.9$ ,  $P_A(3) = 0.2$

- 3 NPs: 1, 2, 3
- $P_c(1, 2) = 0.6$ ,  $P_c(1, 3) = 0.2$ ,  $P_c(2, 3) = 0.9$ •  $P_A(1) = 0.1$ ,  $P_A(2) = 0.9$ ,  $P_A(3) = 0.2$

NP 3 is anaphoric!!!

• 3 NPs: 1, 2, 3

•  $P_c(1, 2) = 0.6$ ,  $P_c(1, 3) = 0.2$ ,  $P_c(2, 3) = 0.9$ •  $P_A(1) = 0.1$ ,  $P_A(2) = 0.9$ ,  $P_A(3) = 0.2$ 

NP 3 is **not** anaphoric!

NP 3 is anaphoric!!!

- 3 NPs: 1, 2, 3
- $P_c(1, 2) = 0.6$ ,  $P_c(1, 3) = 0.2$ ,  $P_c(2, 3) = 0.9$ •  $P_A(1) = 0.1$ ,  $P_A(2) = 0.9$ ,  $P_A(3) = 0.2$

NP 3 is **not** anaphoric!

NP 3 is anaphoric!!!

• P<sub>A</sub> and P<sub>C</sub>'s outputs don't seem to be consistent. Why???

- 3 NPs: 1, 2, 3
- $P_c(1, 2) = 0.6$ ,  $P_c(1, 3) = 0.2$ ,  $P_c(2, 3) = 0.9$

 $P_A(1) = 0.1, P_A(2) = 0.9, P_A(3) = 0.2$ 

NP 3 is **not** anaphoric!

NP 3 is anaphoric!!!

- $P_A$  and  $P_C$ 's outputs don't seem to be consistent. Why???
  - Because they are trained independently of each other

- 3 NPs: 1, 2, 3
- $P_c(1, 2) = 0.6$ ,  $P_c(1, 3) = 0.2$ ,  $P_c(2, 3) = 0.9$

 $P_A(1) = 0.1, P_A(2) = 0.9, P_A(3) = 0.2$ 

NP 3 is **not** anaphoric!

NP 3 is anaphoric!!!

- P<sub>A</sub> and P<sub>C</sub>'s outputs don't seem to be consistent. Why???
  - Because they are trained independently of each other
- Certain hard constraints need to be enforced
  - If P<sub>c</sub> determines that NP<sub>j</sub> is not coreferent with any NP, then P<sub>A</sub> should determine that NP<sub>j</sub> is non-anaphoric

• ...

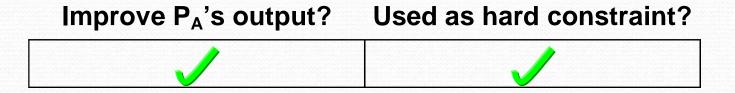
## **ILP for Anaphoricity and Coreference**

- Goal
  - jointly determine anaphoricity and coreference decisions such that all the desired constraints are satisfied
- improve anaphoricity decisions with automatically computed coreference information and manually-specified constraints

# Kleener (2007), Finkel & Manning (2008)

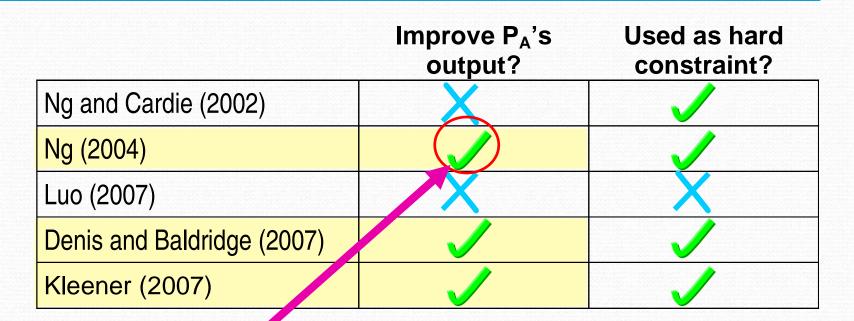
Improve P<sub>A</sub>'s output? Used as hard constraint?

# Kleener (2007), Finkel & Manning (2008)

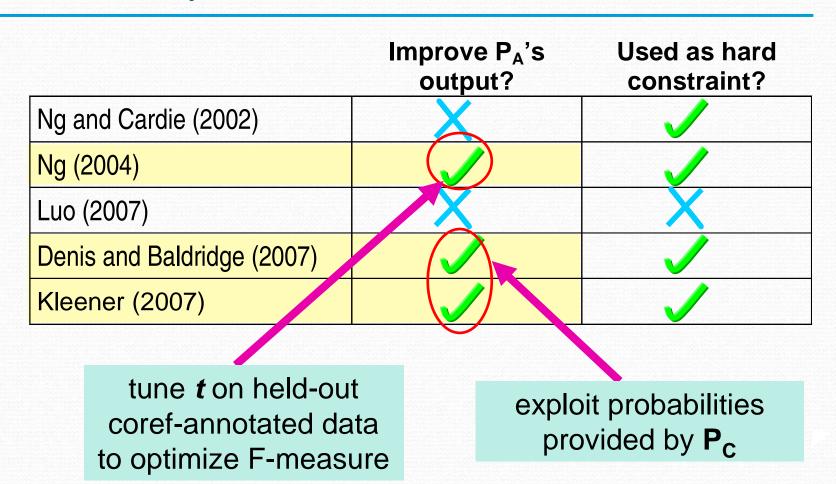


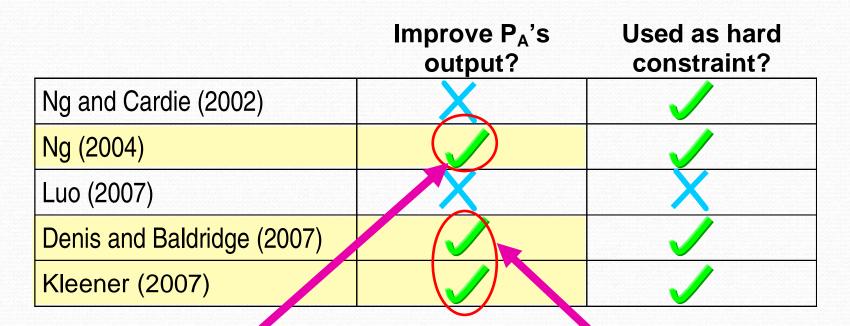
- Also employ ILP, but additionally impose the transitivity constraint on the coreference decisions
  - A, B are coref and B,C are coref
     A,C are coreferent

	Improve P <sub>A</sub> 's output?	Used as hard constraint?
Ng and Cardie (2002)		
Ng (2004)		
Luo (2007)	X	X
Denis and Baldridge (2007)		
Kleener (2007)		



tune *t* on held-out coref-annotated data to optimize F-measure

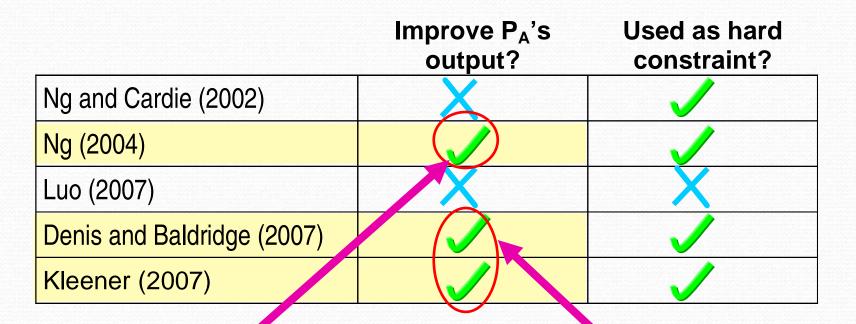




tune *t* on held-out coref-annotated data to optimize F-measure

exploit probabilities provided by **P**<sub>C</sub>

do **not** optimize F-measure



tune *t* on held-out coref-annotated data to optimize F-measure

does not exploit Pc

exploit probabilities provided by **P**<sub>C</sub>

do **not** optimize F-measure

## Summary of the Five

Can we have a method that optimizes F-measure and exploits P<sub>c</sub>?

 Ng and Cardie (2002)
 Constraint?

 Ng (2004)
 Image: Constraint in the constraint in

tune *t* on held-out coref-annotated data to optimize F-measure

does not exploit Pc

exploit probabilities provided by **P**<sub>C</sub>

do **not** optimize F-measure

### **Cut-Based Anaphoricity Determination**

- Motivated by our desire to have a method that can
  - optimize the desired coreference evaluation metric
  - exploit probabilities provided by P<sub>c</sub>

Want to partition a set of objects, {x<sub>1</sub>, x<sub>2</sub>, ..., x<sub>n</sub>}, into two sets, S and T

- Want to partition a set of objects, {x<sub>1</sub>, x<sub>2</sub>, ..., x<sub>n</sub>}, into two sets, S and T
- Given two types of scores:
  - Membership scores:  $mem_S(x_i)$ ,  $mem_T(x_i)$ 
    - captures the affinity of  $x_i$  to S and T, respectively
    - large  $mem_S(x_i)$   $x_i$  is likely to be in S
  - Similarity scores: sim(x<sub>i</sub>, x<sub>i</sub>)
    - captures the similarity between  $x_i$  and  $x_j$

- Want to partition a set of objects, {x<sub>1</sub>, x<sub>2</sub>, ..., x<sub>n</sub>}, into two sets, S and T
- Given two types of scores:
  - Membership scores:  $mem_S(x_i)$ ,  $mem_T(x_i)$ 
    - captures the affinity of  $x_i$  to S and T, respectively
    - large  $mem_S(x_i)$   $x_i$  is likely to be in S
  - Similarity scores: sim(x<sub>i</sub>, x<sub>i</sub>)
    - captures the similarity between  $x_i$  and  $x_j$

- Want to partition a set of objects, {x<sub>1</sub>, x<sub>2</sub>, ..., x<sub>n</sub>}, into two sets, S and T
- Given two types of scores:
  - Membership scores:  $mem_S(x_i)$ ,  $mem_T(x_i)$ 
    - captures the affinity of  $x_i$  to S and T, respectively
    - large  $mem_S(x_i)$   $x_i$  is likely to be in S
  - Similarity scores: sim(x<sub>i</sub>, x<sub>i</sub>)
    - captures the similarity between  $x_i$  and  $x_j$
- Goal:

$$\operatorname{Maximize} \sum_{x_i \in S, x_j \in T} sim(x_i, x_j) + \sum_{x \in S} mem_S(x) + \sum_{x \in T} mem_T(x)$$

- Want to partition a set of objects, {x<sub>1</sub>, x<sub>2</sub>, ..., x<sub>n</sub>}, into two sets, S and T
- Given two types of scores:
  - Membership scores:  $mem_S(x_i)$ ,  $mem_T(x_i)$ 
    - captures the affinity of x<sub>i</sub> to S and T, respectively
    - Put similar objects is likely to be in S
  - S into the same set  $(x_i, x_i)$ 
    - captures the similarity between  $x_i$  and  $x_j$
- Goal: Maximize  $\sum_{x_i \in S, x_i \in T} sim(x_i, x_j) + \sum_{x \in S} mem_S(x) + \sum_{x \in T} mem_T(x)$

- Want to partition a set of objects, {x<sub>1</sub>, x<sub>2</sub>, ..., x<sub>n</sub>}, into two sets, S and T
- Given two types of scores:
  - Membership scores:  $mem_S(x_i)$ ,  $mem_T(x_i)$ 
    - captures the affinity of x<sub>i</sub> to S and T, respectively
    - Put similar objects is likely to be Put an object to the set
  - S into the same set  $(x_i, x_j)$  where its membership score is high
    - captures the similarity between  $x_i$  and  $x_j$
- Goal:

$$\operatorname{Maximize} \sum_{x_i \in S, x_j \in T} sim(x_i, x_j) + \sum_{x \in S} mem_S(x) + \sum_{x \in T} mem_T(x)$$

# Solving this Problem Using MinCut





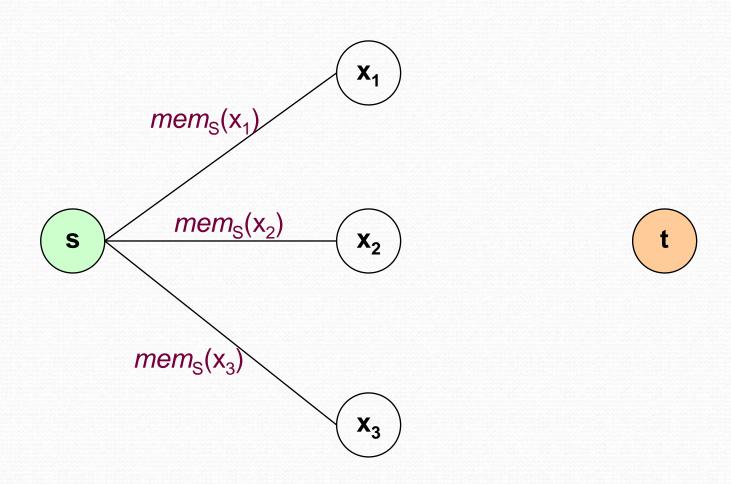


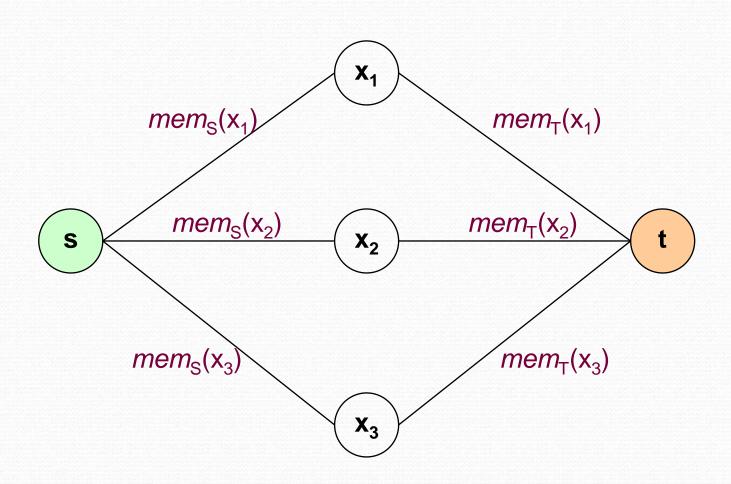


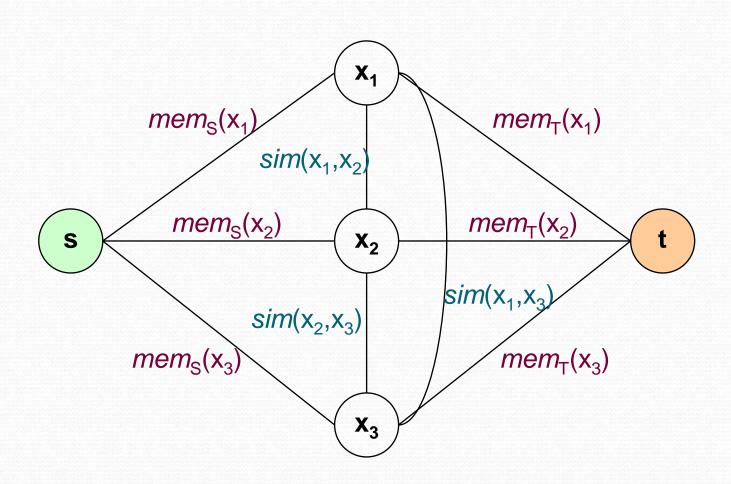


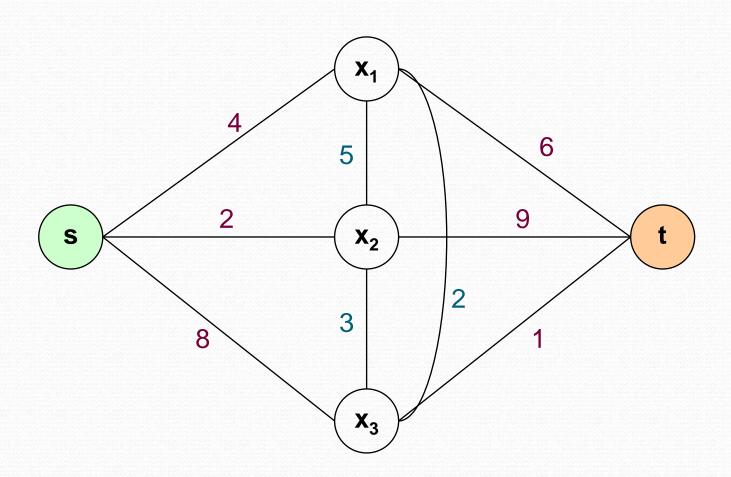


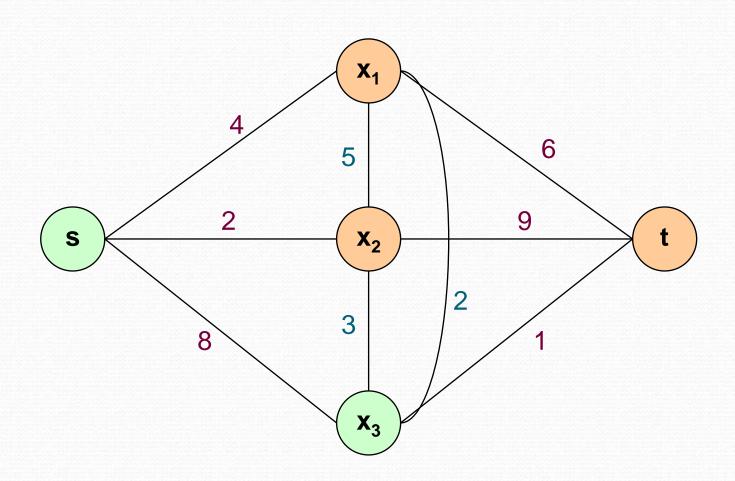


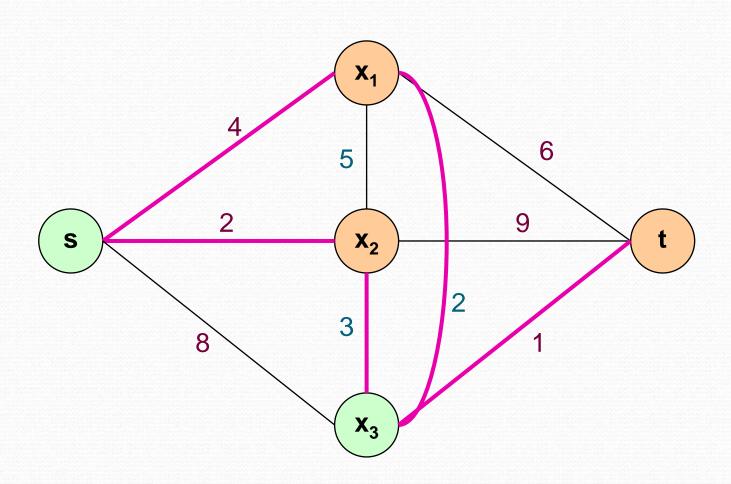


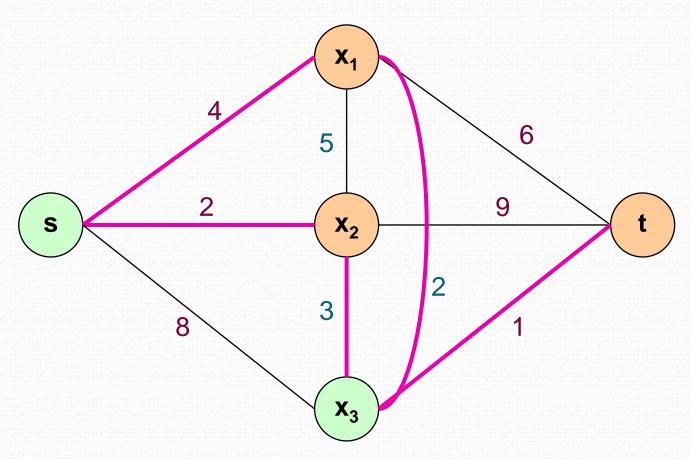






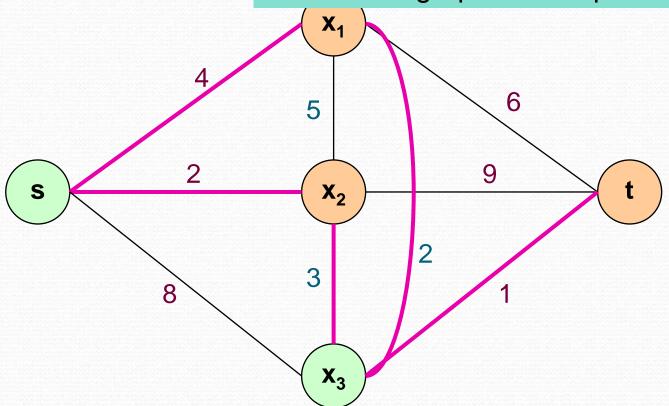




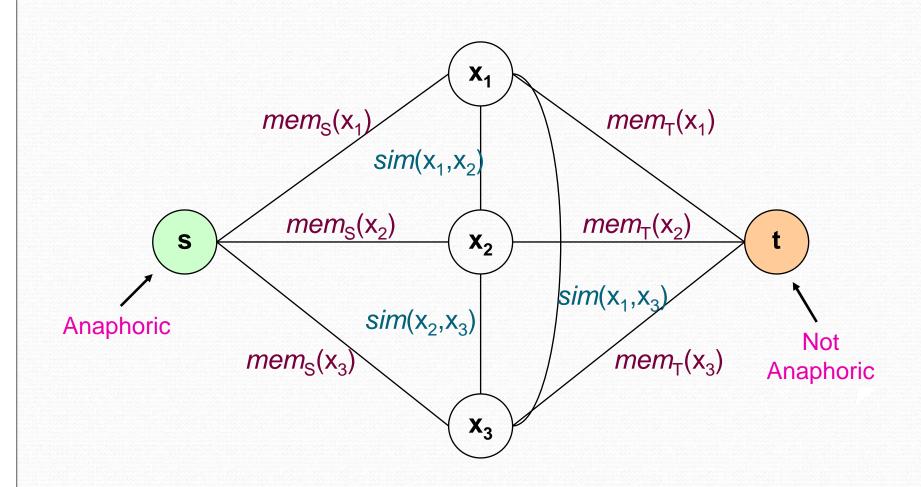


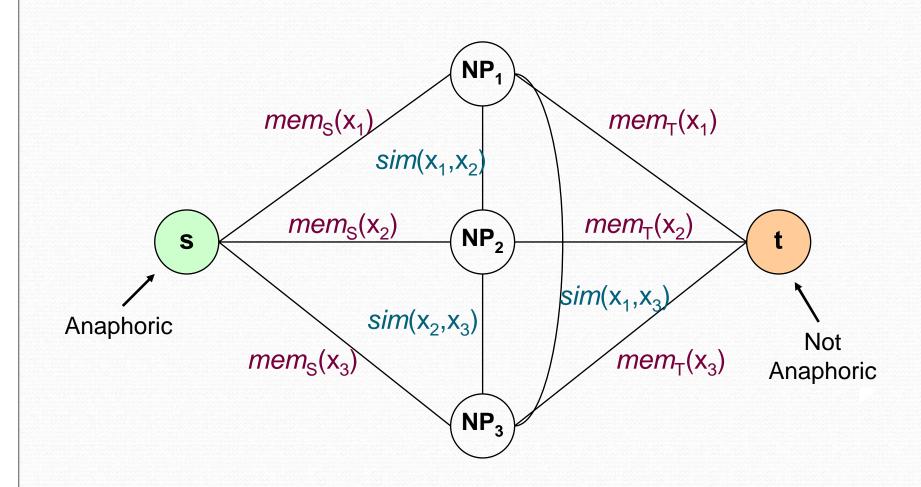
Efficient algorithms for finding the mincut exist

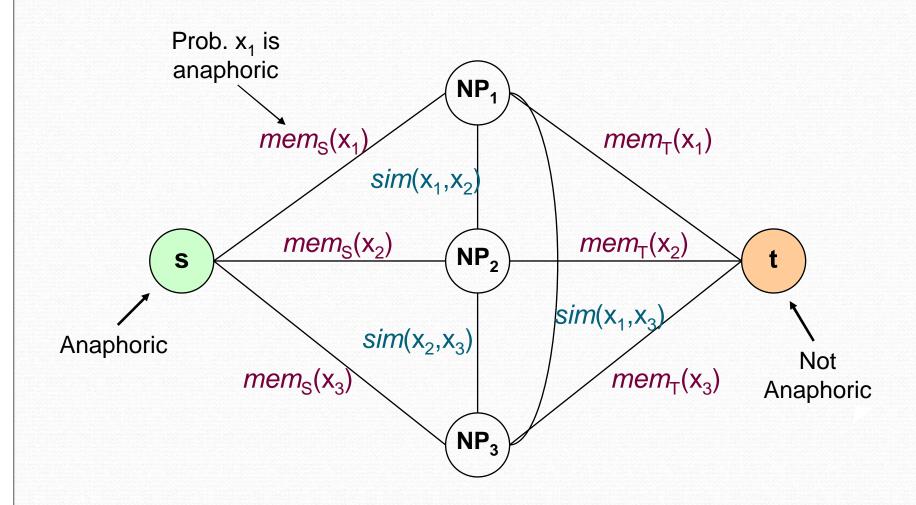
How to recast anaphoricity determination as a graph mincut problem?

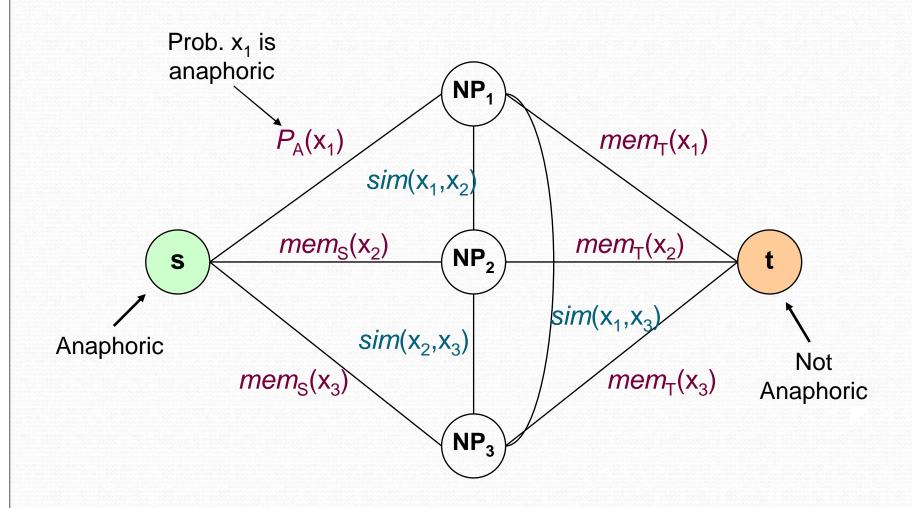


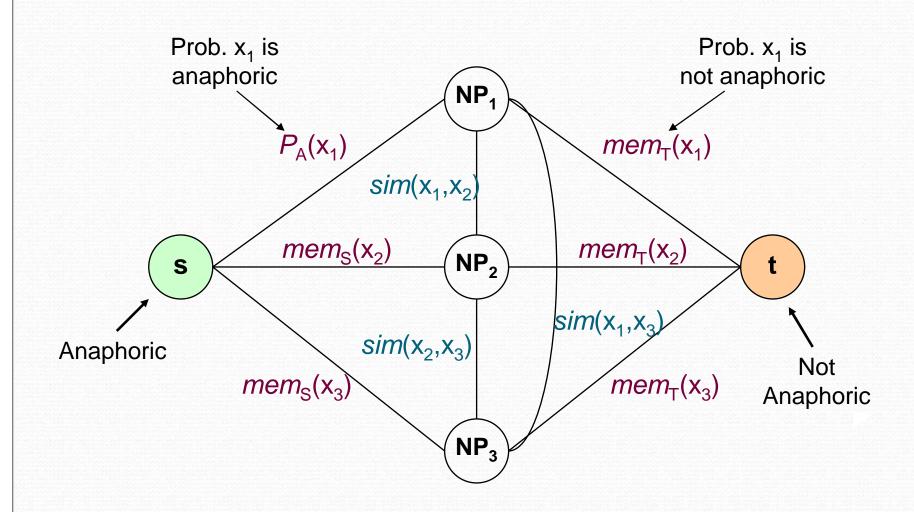
Efficient algorithms for finding the mincut exist

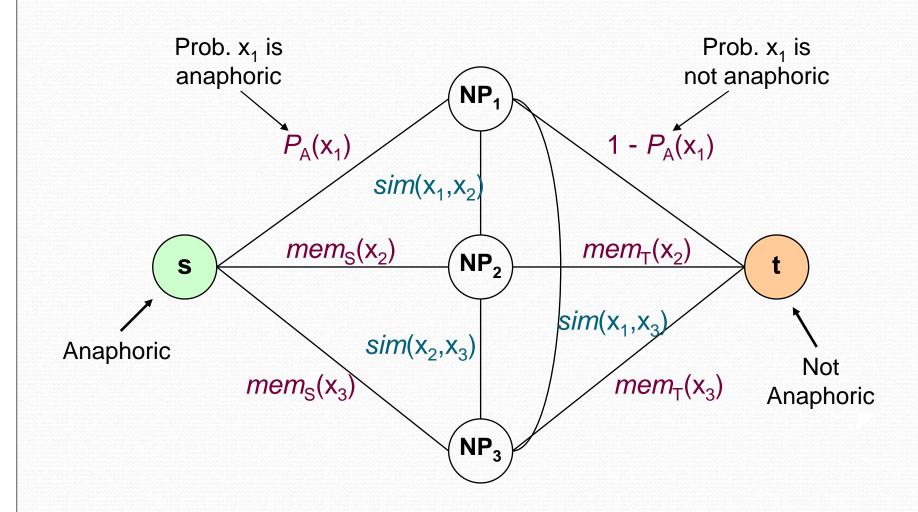


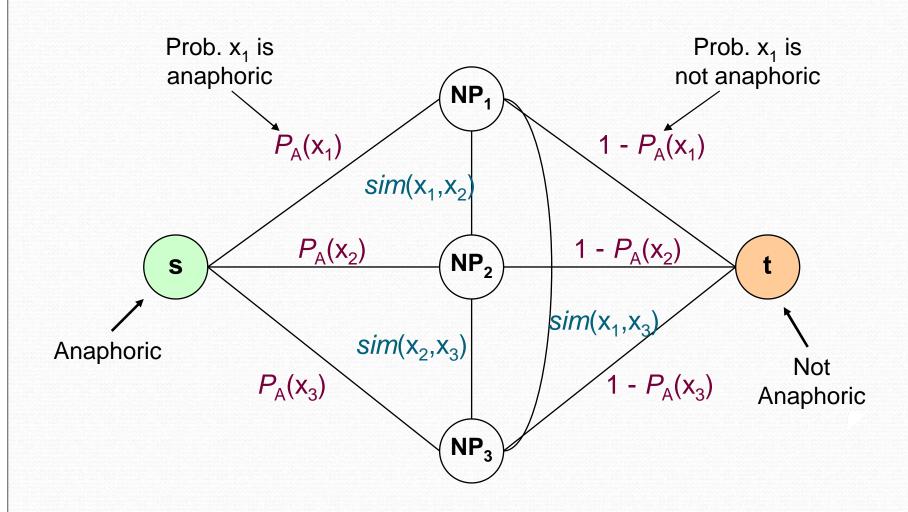


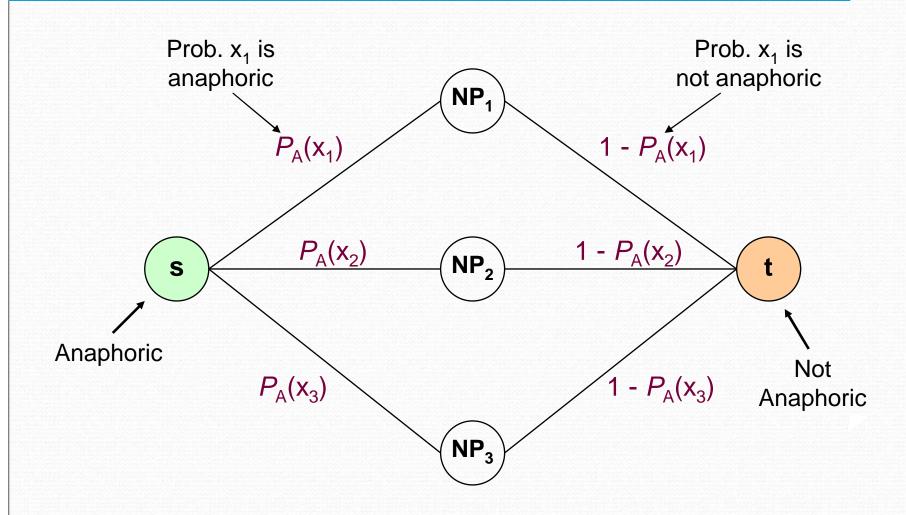




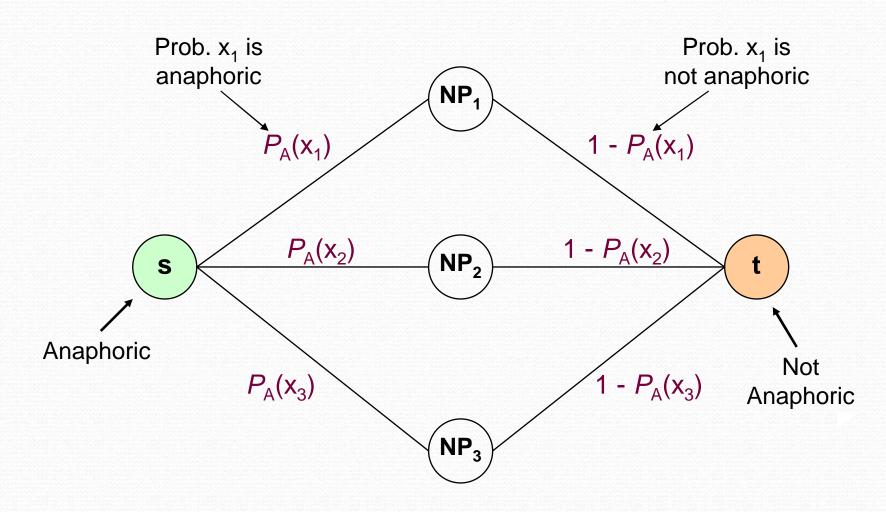




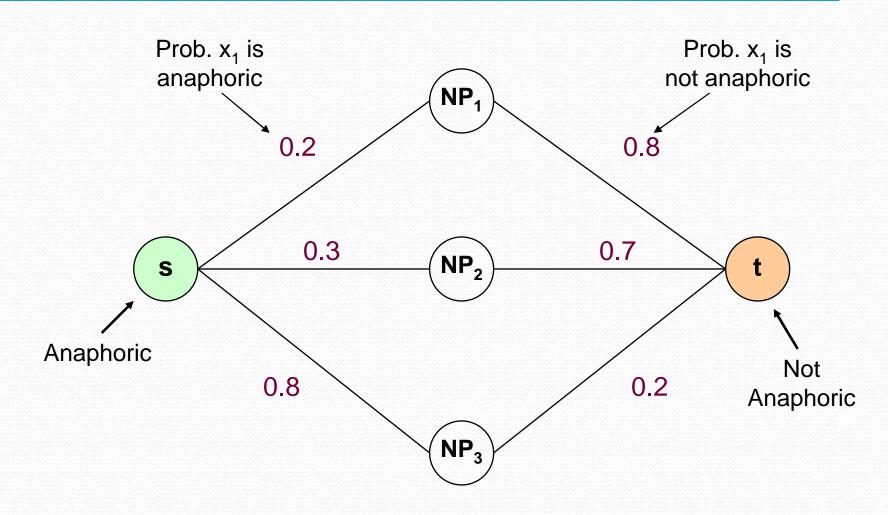




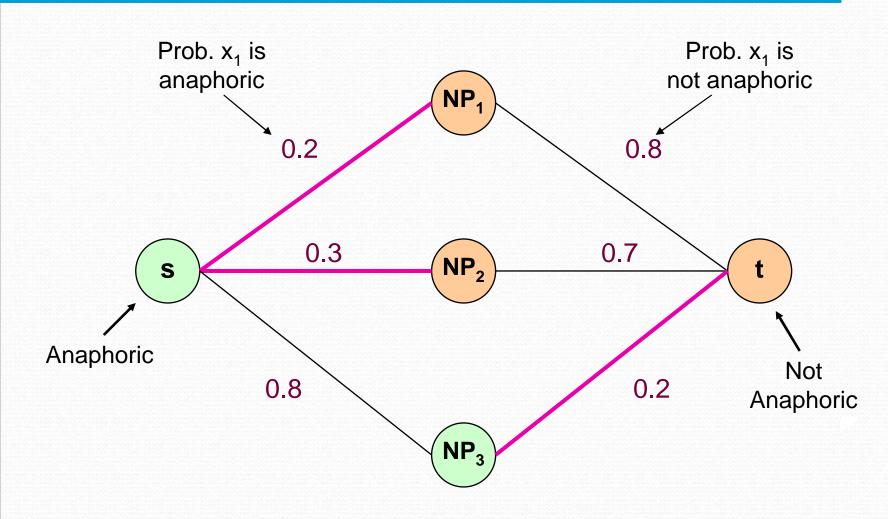
What if we run the mincut finding algorithm on this graph?

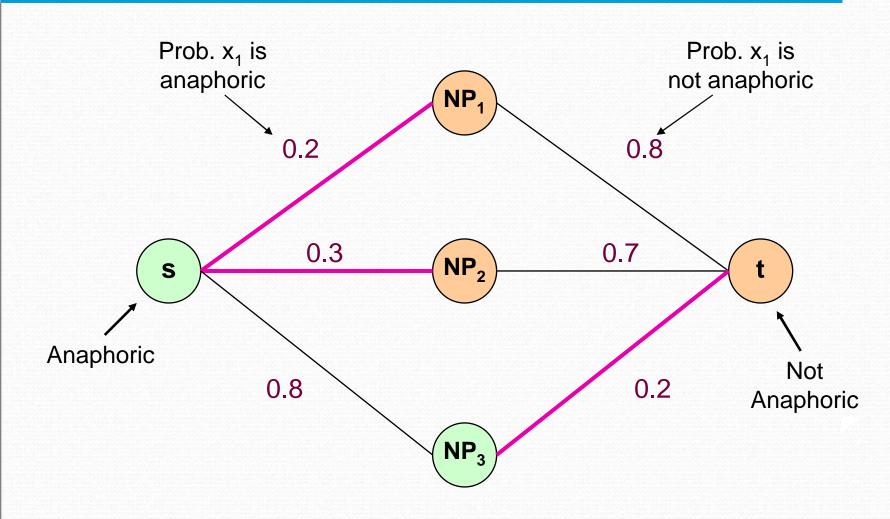


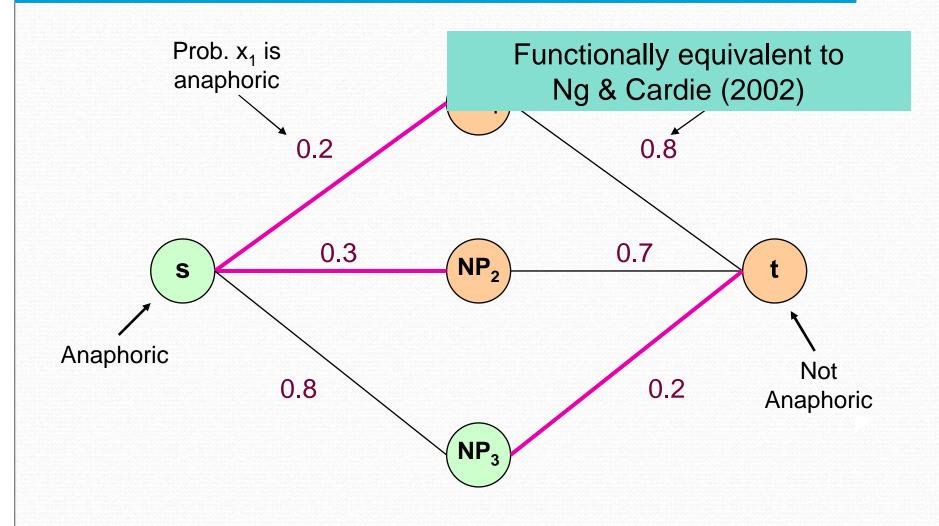
## What if we run the mincut finding algorithm on this graph?

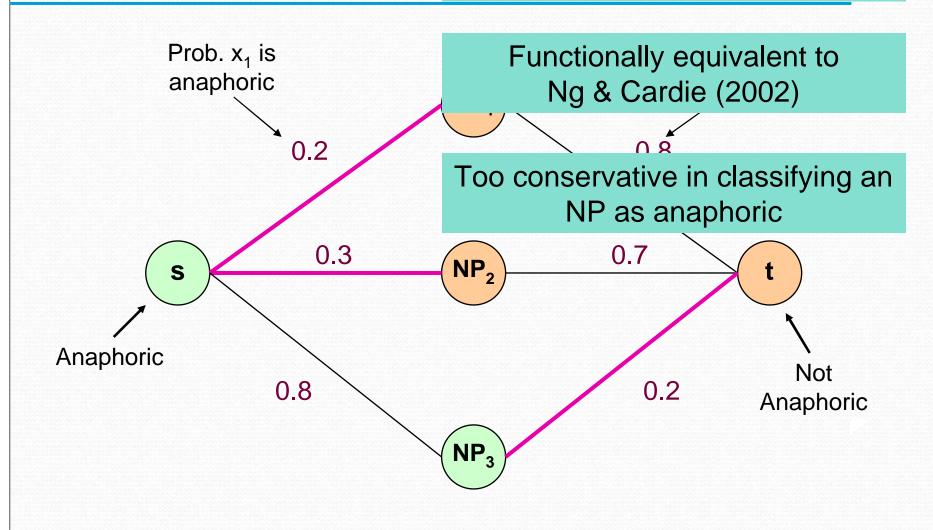


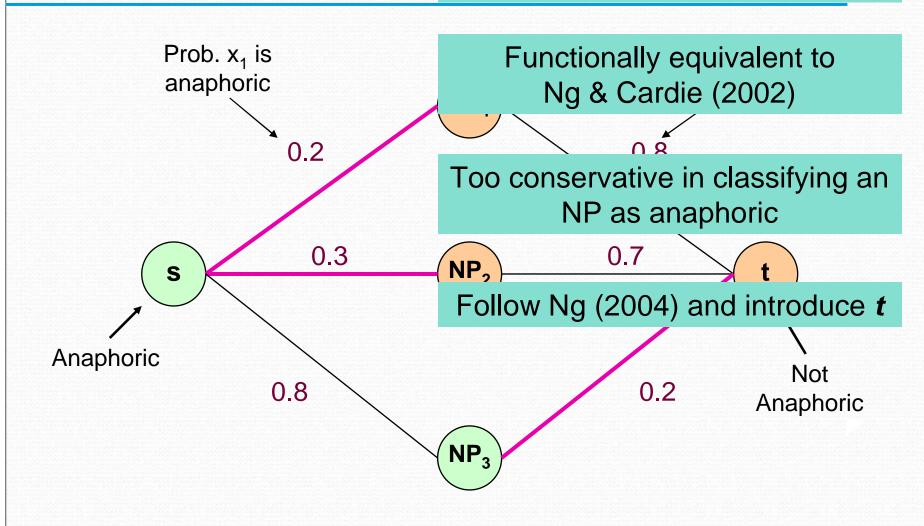
## What if we run the mincut finding algorithm on this graph?











NP<sub>i</sub> is classified as anaphoric iff P<sub>A</sub>(NP<sub>i</sub>) ≥ t

- NP<sub>i</sub> is classified as anaphoric iff P<sub>A</sub>(NP<sub>i</sub>) ≥ t
  - *t* is the "conservativeness" parameter
    - tuned on held-out data to maximize coreference F-measure

- NP<sub>i</sub> is classified as anaphoric iff P<sub>A</sub>(NP<sub>i</sub>) ≥ t
  - *t* is the "conservativeness" parameter
    - tuned on held-out data to maximize coreference F-measure
- Goal: modify the edge weights s.t. NP<sub>i</sub> will be assigned to the Anaphoric class iff P<sub>A</sub>(NP<sub>i</sub>) ≥ t

- NP<sub>i</sub> is classified as anaphoric iff P<sub>A</sub>(NP<sub>i</sub>) ≥ t
  - *t* is the "conservativeness" parameter
    - tuned on held-out data to maximize coreference F-measure
- Goal: modify the edge weights s.t. NP<sub>i</sub> will be assigned to the Anaphoric class iff P<sub>A</sub>(NP<sub>i</sub>) ≥ t
- How? Do a linear transformation

- NP<sub>i</sub> is classified as anaphoric iff P<sub>A</sub>(NP<sub>i</sub>) ≥ t
  - *t* is the "conservativeness" parameter
    - tuned on held-out data to maximize coreference F-measure
- Goal: modify the edge weights s.t. NP<sub>i</sub> will be assigned to the Anaphoric class iff P<sub>A</sub>(NP<sub>i</sub>) ≥ t
- How? Do a linear transformation
- But ... we are not happy with just mimicking Ng (2004)

• What is a good candidate similarity function?

- What is a good candidate similarity function?
- Two observations:
  - If  $NP_i$  and  $NP_j$  are likely to be coreferent according to  $\mathbf{P_C}$ , then  $NP_i$  and  $NP_j$  are likely to be both anaphoric (except if  $NP_i$  is the head of a coreference chain)

- What is a good candidate similarity function?
- Two observations:
  - If NP<sub>i</sub> and NP<sub>j</sub> are likely to be coreferent according to P<sub>c</sub>, then NP<sub>i</sub> and NP<sub>j</sub> are likely to be both anaphoric (except if NP<sub>i</sub> is the head of a coreference chain)
    - want mincut finder to assign the NPs to the same set

- What is a good candidate similarity function?
- Two observations:
  - If NP<sub>i</sub> and NP<sub>j</sub> are likely to be coreferent according to P<sub>c</sub>, then NP<sub>i</sub> and NP<sub>j</sub> are likely to be both anaphoric (except if NP<sub>i</sub> is the head of a coreference chain)
    - want mincut finder to assign the NPs to the same set
  - If  $NP_i$  and  $NP_j$  are unlikely to be coreferent according to  $P_c$ , it's hard to claim anything regarding their anaphoricity

- What is a good candidate similarity function?
- Two observations:
  - If NP<sub>i</sub> and NP<sub>j</sub> are likely to be coreferent according to P<sub>c</sub>, then NP<sub>i</sub> and NP<sub>j</sub> are likely to be both anaphoric (except if NP<sub>i</sub> is the head of a coreference chain)
    - want mincut finder to assign the NPs to the same set
  - If  $NP_i$  and  $NP_j$  are unlikely to be coreferent according to  $P_c$ , it's hard to claim anything regarding their anaphoricity
- Use  $\mathbf{P_c}(NP_i, NP_j)$  as the  $sim(NP_i, NP_j)$ , but only if  $\mathbf{P_c}(NP_i, NP_j) > t_2$

- What is a good candidate similarity function?
- Two observations:
  - If NP<sub>i</sub> and NP<sub>j</sub> are likely to be coreferent according to P<sub>c</sub>, then NP<sub>i</sub> and NP<sub>j</sub> are likely to be both anaphoric (except if NP<sub>i</sub> is the head of a coreference chain)
    - want mincut finder to assign the NPs to the same set
  - If  $NP_i$  and  $NP_j$  are unlikely to be coreferent according to  $\mathbf{P_c}$ , it's hard to claim anything regarding their anaphoricity
- Use  $P_{\mathbf{c}}(NP_i, NP_j)$  as the  $sim(NP_i, NP_j)$ , but only if  $P_{\mathbf{c}}(NP_i, NP_i) > t_2$  [tuned jointly with t on held-out data]

- What is a good candidate similarity function?
- Two observations:
  - If  $NP_i$  and  $NP_j$  are likely to be coreferent according to  $\mathbf{P_c}$ , then  $NP_i$  and  $NP_j$  are likely to be both anaphoric (except if  $NP_i$  is the head of a coreference chain)
    - want mincut finder to assign the NPs to the same set
  - If  $NP_i$  and  $NP_j$  are unlikely to be coreferent according to  $\mathbf{P_c}$ , it's hard to claim anything regarding their anaphoricity
- Use  $P_{\mathbf{C}}(NP_i, NP_j)$  as the  $sim(NP_i, NP_j)$ , but only if  $P_{\mathbf{C}}(NP_i, NP_j) > t_2$  [tuned jointly with t on held-out data]
- Otherwise, set sim(NP<sub>i</sub>, NP<sub>i</sub>) to 0

#### **Cut-Based Anaphoricity Determination**

- An anaphoricity determination method that can
  - maximize the desired coreference metric (by tuning t and  $t_2$ )
  - exploit probabilities provided by P<sub>c</sub> (via the sim function)

#### Plan for the Talk

- Existing methods for computing and using anaphoricity info
- Our graph-cut-based approach to anaphoricity determination
- Evaluation

#### **Evaluation: Goal**

 compare our cut-based method for anaphoricity determination with existing methods w.r.t. their effectiveness in improving a learning-based coreference system

#### **Experimental Setup**

- Coreference system [Ng, 2007]
  - implements the standard machine learning framework
  - 34 features per instance
- Features for anaphoricity determination [Ng & Cardie, 2002]
  - 37 features per instance
- Learning algorithm
  - Maximum entropy for training P<sub>A</sub> and P<sub>C</sub>

### Experimental Setup (Cont')

- The ACE coreference corpus
  - 3 data sets (Broadcast News, Newspaper, Newswire)
  - each data set comprises a training set and a test set
- NPs extracted automatically
- Scoring programs
  - MUC (Vilain et al., 1995) and CEAF (Luo, 2005)
    - recall, precision, F-measure

	Broa	<b>Broadcast News</b>			Newspaper			Newswire		
	R	P	F	R	P	F	R	P	F	
No anaphoricity	63.2	49.2	55.3	64.5	54.3	59.0	67.3	56.1	61.2	

	Broa	<b>Broadcast News</b>			Newspaper			Newswire			
	R	P	F	R	P	F	R	P	F		
No anaphoricity	63.2	49.2	55.3	64.5	54.3	59.0	67.3	56.1	61.2		

	Broa	<b>Broadcast News</b>			Newspaper			Newswire		
	R	P	F	R	P	F	R	P	F	
No anaphoricity	63.2	49.2	55.3	64.5	54.3	59.0	67.3	56.1	61.2	

	Broa	Broadcast News			ewspap	er	Newswire			
	R	P	F	R	P	F	R	P	F	
No anaphoricity	63.2	49.2	55.3	64.5	54.3	59.0	67.3	56.1	61.2	

### CEAF Results: Ng & Cardie (2002) Baseline

	Broa	dcast	News	Ne	Newspaper			Newswire		
	R	P	F	R	P	F	R	P	F	
No anaphoricity	63.2	49.2	55.3	64.5	54.3	59.0	67.3	56.1	61.2	
Ng & Cardie (2002)	55.9	53.3	54.5	60.7	56.3	58.3	60.6	58.2	59.4	

#### CEAF Results: Ng & Cardie (2002) Baseline

	Broa	<b>Broadcast News</b>			Newspaper			Newswire		
	R	P	F	R	P	F	R	P	F	
No anaphoricity	63.2	49.2	55.3	64.5	54.3	59.0	67.3	56.1	61.2	
Ng & Cardie (2002)	55.9	53.3	54.5	60.7	56.3	58.3	60.6	58.2	59.4	

- F-measure drops slightly in all cases
  - large drops in recall accompanied by smaller gains in precision
  - many anaphoric NPs were misclassified

### CEAF Results: Ng (2004) Baseline

	Broa	<b>Broadcast News</b>			ewspap	er	Newswire		
	R	P	F	R	P	F	R	P	F
No anaphoricity	63.2	49.2	55.3	64.5	54.3	59.0	67.3	56.1	61.2
Ng & Cardie (2002)	55.9	53.3	54.5	60.7	56.3	58.3	60.6	58.2	59.4
Ng (2004)	62.5	49.9	55.5	63.5	57.0	61.0	65.6	56.3	60.6

### CEAF Results: Ng (2004) Baseline

	Broa	<b>Broadcast News</b>			Newspaper			Newswire		
	R	P	F	R	P	F	R	P	F	
No anaphoricity	63.2	49.2	55.3	64.5	54.3	59.0	67.3	56.1	61.2	
Ng & Cardie (2002)	55.9	53.3	54.5	60.7	56.3	58.3	60.6	58.2	59.4	
Ng (2004)	62.5	49.9	55.5	63.5	57.0	61.0	65.6	56.3	60.6	

- requires tuning of t
  - reserve 1/3 of the training data for parameter tuning
  - train P<sub>A</sub> and P<sub>C</sub> on remaining 2/3 of the training data

### CEAF Results: Ng (2004) Baseline

	Broa	<b>Broadcast News</b>			ewspap	er	Newswire		
	R	P	F	R	P	F	R	P	F
No anaphoricity	63.2	49.2	55.3	64.5	54.3	59.0	67.3	56.1	61.2
Ng & Cardie (2002)	55.9	53.3	54.5	60.7	56.3	58.3	60.6	58.2	59.4
Ng (2004)	62.5	49.9	55.5	63.5	57.0	61.0	65.6	56.3	60.6

- requires tuning of t
  - reserve 1/3 of the training data for parameter tuning
  - train P<sub>A</sub> and P<sub>C</sub> on remaining 2/3 of the training data
- mixed results in comparison to "No Anaphoricity" baseline
  - F-measure gains by 0.2% for BNEWS, 1.1% for NPAPER, but drops by 0.6% for NWIRE

### CEAF Results: Luo (2007) Baseline

	Broa	<b>Broadcast News</b>			Newspaper			Newswire		
	R	P	F	R	P	F	R	P	F	
No anaphoricity	63.2	49.2	55.3	64.5	54.3	59.0	67.3	56.1	61.2	
Ng & Cardie (2002)	55.9	53.3	54.5	60.7	56.3	58.3	60.6	58.2	59.4	
Ng (2004)	62.5	49.9	55.5	63.5	57.0	61.0	65.6	56.3	60.6	
Luo (2007)	62.7	51.1	56.3	64.6	55.4	59.6	67.0	56.8	61.5	

### CEAF Results: Luo (2007) Baseline

	Broa	<b>Broadcast News</b>			Newspaper			Newswire		
	R	P	F	R	P	F	R	P	F	
No anaphoricity	63.2	49.2	55.3	64.5	54.3	59.0	67.3	56.1	61.2	
Ng & Cardie (2002)	55.9	53.3	54.5	60.7	56.3	58.3	60.6	58.2	59.4	
Ng (2004)	62.5	49.9	55.5	63.5	57.0	61.0	65.6	56.3	60.6	
Luo (2007)	62.7	51.1	56.3	64.6	55.4	59.6	67.0	56.8	61.5	

- in comparison to "No Anaphoricity" baseline
  - F-measure improves insignificantly (by 0.3-1.0%)

### Results: Denis & Baldridge (2007) Baseline

	<b>Broadcast News</b>			Ne	ewspap	er	Newswire		
	R	P	F	R	P	F	R	P	F
No anaphoricity	63.2	49.2	55.3	64.5	54.3	59.0	67.3	56.1	61.2
Ng & Cardie (2002)	55.9	53.3	54.5	60.7	56.3	58.3	60.6	58.2	59.4
Ng (2004)	62.5	49.9	55.5	63.5	57.0	61.0	65.6	56.3	60.6
Luo (2007)	62.7	51.1	56.3	64.6	55.4	59.6	67.0	56.8	61.5
Denis & Baldridge (2007)	63.8	51.4	56.9	62.6	53.6	57.8	67.0	56.8	61.5

#### Results: Denis & Baldridge (2007) Baseline

	<b>Broadcast News</b>			Newspaper			Newswire		
	R	P	F	R	P	F	R	P	F
No anaphoricity	63.2	49.2	55.3	64.5	54.3	59.0	67.3	56.1	61.2
Ng & Cardie (2002)	55.9	53.3	54.5	60.7	56.3	58.3	60.6	58.2	59.4
Ng (2004)	62.5	49.9	55.5	63.5	57.0	61.0	65.6	56.3	60.6
Luo (2007)	62.7	51.1	56.3	64.6	55.4	59.6	67.0	56.8	61.5
Denis & Baldridge (2007)	63.8	51.4	56.9	62.6	53.6	57.8	67.0	56.8	61.5

- mixed results in comparison to "No Anaphoricity" baseline
  - F-measure rises significantly for BNEWS, drop insignificantly for NPAPER, and rises insignificantly for NWIRE

### CEAF Results: Kleener (2007) Baseline

	<b>Broadcast News</b>			Newspaper			Newswire		
	R	P	F	R	P	F	R	P	F
No anaphoricity	63.2	49.2	55.3	64.5	54.3	59.0	67.3	56.1	61.2
Ng & Cardie (2002)	55.9	53.3	54.5	60.7	56.3	58.3	60.6	58.2	59.4
Ng (2004)	62.5	49.9	55.5	63.5	57.0	61.0	65.6	56.3	60.6
Luo (2007)	62.7	51.1	56.3	64.6	55.4	59.6	67.0	56.8	61.5
Denis & Baldridge (2007)	63.8	51.4	56.9	62.6	53.6	57.8	67.0	56.8	61.5
Kleener (2007)	63.2	51.3	56.7	62.6	53.6	57.8	66.7	56.7	61.3

### CEAF Results: Kleener (2007) Baseline

	<b>Broadcast News</b>			Newspaper			Newswire		
	R	P	F	R	P	F	R	P	F
No anaphoricity	63.2	49.2	55.3	64.5	54.3	59.0	67.3	56.1	61.2
Ng & Cardie (2002)	55.9	53.3	54.5	60.7	56.3	58.3	60.6	58.2	59.4
Ng (2004)	62.5	49.9	55.5	63.5	57.0	61.0	65.6	56.3	60.6
Luo (2007)	62.7	51.1	56.3	64.6	55.4	59.6	67.0	56.8	61.5
Denis & Baldridge (2007)	63.8	51.4	56.9	62.6	53.6	57.8	67.0	56.8	61.5
Kleener (2007)	63.2	51.3	56.7	62.6	53.6	57.8	66.7	56.7	61.3

- in comparison to Denis & Baldridge baseline,
  - F-measure never improves, recall slightly deteriorates
  - transitivity constraints tend to produce smaller clusters
  - enforcing transitivity does not improve coreference performance

	<b>Broadcast News</b>			Newspaper			Newswire		
	R	P	F	R	P	F	R	P	F
No anaphoricity	63.2	49.2	55.3	64.5	54.3	59.0	67.3	56.1	61.2
Ng & Cardie (2002)	55.9	53.3	54.5	60.7	56.3	58.3	60.6	58.2	59.4
Ng (2004)	62.5	49.9	55.5	63.5	57.0	61.0	65.6	56.3	60.6
Luo (2007)	62.7	51.1	56.3	64.6	55.4	59.6	67.0	56.8	61.5
Denis & Baldridge (2007)	63.8	51.4	56.9	62.6	53.6	57.8	67.0	56.8	61.5
Kleener (2007)	63.2	51.3	56.7	62.6	53.6	57.8	66.7	56.7	61.3
Graph Minimum Cut	61.4	57.6	59.4	64.1	59.4	61.7	65.7	61.9	63.8

• 1/3 of training data for joint tuning of t and  $t_2$ ; 2/3 for training  $P_A$  and  $P_C$ 

	<b>Broadcast News</b>			Newspaper			Newswire		
	R	P	F	R	P	F	R	P	F
No anaphoricity	63.2	49.2	55.3	64.5	54.3	59.0	67.3	56.1	61.2
Ng & Cardie (2002)	55.9	53.3	54.5	60.7	56.3	58.3	60.6	58.2	59.4
Ng (2004)	62.5	49.9	55.5	63.5	57.0	61.0	65.6	56.3	60.6
Luo (2007)	62.7	51.1	56.3	64.6	55.4	59.6	67.0	56.8	61.5
Denis & Baldridge (2007)	63.8	51.4	56.9	62.6	53.6	57.8	67.0	56.8	61.5
Kleener (2007)	63.2	51.3	56.7	62.6	53.6	57.8	66.7	56.7	61.3
Graph Minimum Cut	61.4	57.6	59.4	64.1	59.4	61.7	65.7	61.9	63.8

- 1/3 of training data for joint tuning of t and  $t_2$ ; 2/3 for training  $P_A$  and  $P_C$
- significant improvement over "No Anaphoricity" baseline
  - large gains in precision and smaller drops in recall

	<b>Broadcast News</b>			Newspaper			Newswire		
	R	P	F	R	P	F	R	P	F
No anaphoricity	63.2	49.2	55.3	64.5	54.3	59.0	67.3	56.1	61.2
Ng & Cardie (2002)	55.9	53.3	54.5	60.7	56.3	58.3	60.6	58.2	59.4
Ng (2004)	62.5	49.9	55.5	63.5	57.0	61.0	65.6	56.3	60.6
Luo (2007)	62.7	51.1	56.3	64.6	55.4	59.6	67.0	56.8	61.5
Denis & Baldridge (2007)	63.8	51.4	56.9	62.6	53.6	57.8	67.0	56.8	61.5
Kleener (2007)	63.2	51.3	56.7	62.6	53.6	57.8	66.7	56.7	61.3
Graph Minimum Cut	61.4	57.6	59.4	64.1	59.4	61.7	65.7	61.9	63.8

- 1/3 of training data for joint tuning of t and  $t_2$ ; 2/3 for training  $P_A$  and  $P_C$
- significant improvement over "No Anaphoricity" baseline
  - large gains in precision and smaller drops in recall
- significant improvement over best baseline (D&B)

	<b>Broadcast News</b>			Newspaper			Newswire		
	R	P	F	R	P	F	R	P	F
No anaphoricity	63.2	49.2	55.3	64.5	54.3	59.0	67.3	56.1	61.2
Ng & Cardie (2002)	55.9	53.3	54.5	60.7	56.3	58.3	60.6	58.2	59.4
Ng (2004)	62.5	49.9	55.5	63.5	57.0	61.0	65.6	56.3	60.6
Luo (2007)	62.7	51.1	56.3	64.6	55.4	59.6	67.0	56.8	61.5
Denis & Baldridge (2007)	63.8	51.4	56.9	62.6	53.6	57.8	67.0	56.8	61.5
Kleener (2007)	63.2	51.3	56.7	62.6	53.6	57.8	66.7	56.7	61.3
Graph Minimum Cut	61.4	57.6	59.4	64.1	59.4	61.7	65.7	61.9	63.8

- 1/3 of training data for joint tuning of t and  $t_2$ ; 2/3 for training  $P_A$  and  $P_C$
- significant improvement over "No Anaphoricity" baseline
  - large gains in precision and smaller drops in recall
- significant improvement over best baseline (D&B)
- best F-measure score achieved for each dataset

#### Summary

- Proposed a graph-cut-based approach to anaphoricity determination that
  - directly optimizes the desired coreference evaluation metric
  - exploits the probabilities provided by the coreference model
  - achieves the best results on all three ACE datasets according to both the MUC scorer and the CEAF scorer
  - provides a flexible mechanism for co-ordinating anaphoricity and coreference decisions