

# Review Slides

Introduction to Natural Language Processing  
Computer Science 585—Fall 2009  
University of Massachusetts Amherst

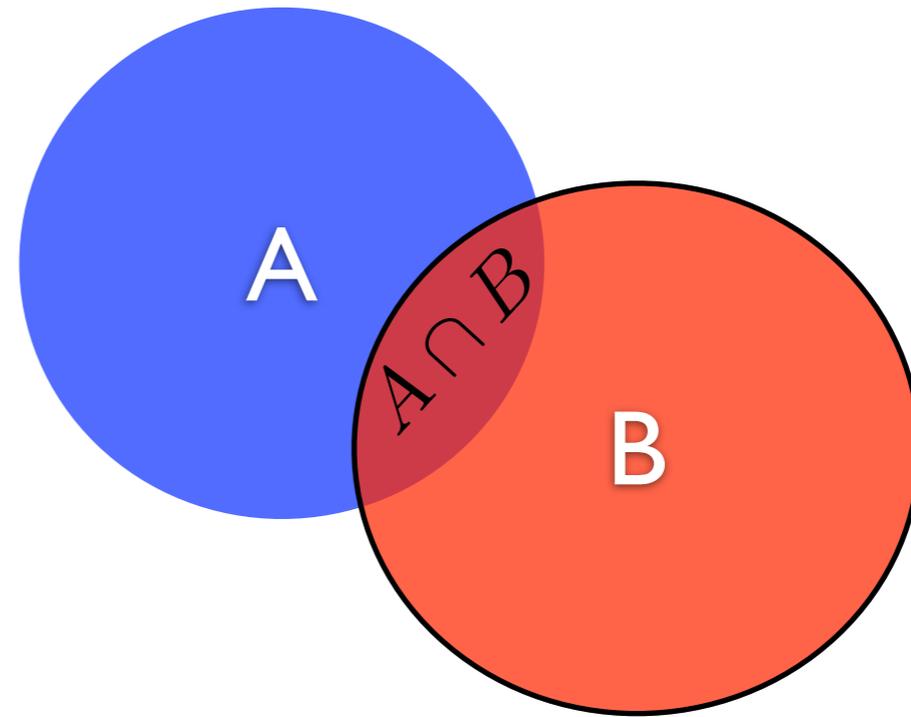
David Smith

# Final Exam

- Wednesday, Dec. 16, 10:30, CS 142
- At least 2/3 from course's second half
- Focus on modeling techniques, such as:
  - Log-linear models
  - Sequence labeling, e.g. for information extraction
  - Formal semantics, simple  $\lambda$ -expressions
  - Word clustering
  - Simple machine translation algorithms: IBM Model-1, ITG

# Conditional Probability

$$P(A | B) = \frac{P(A, B)}{P(B)}$$



$$P(A, B) = P(B)P(A | B) = P(A)P(B | A)$$

$$P(A_1, A_2, \dots, A_n) = P(A_1)P(A_2 | A_1)P(A_3 | A_1, A_2) \dots P(A_n | A_1, \dots, A_{n-1})$$

*Chain rule*

# Independence

$$P(A, B) = P(A)P(B)$$

$\Leftrightarrow$

$$P(A | B) = P(A) \quad \wedge \quad P(B | A) = P(B)$$

In coding terms, knowing  $B$  doesn't help in decoding  $A$ , and vice versa.

# Another View of Markov Models

$$p(w_1, w_2, \dots, w_n) = p(w_1)p(w_2 | w_1)p(w_3 | w_1, w_2) \\ p(w_4 | w_1, w_2, w_3) \cdots p(w_n | p_1, \dots, p_{n-1})$$

# Another View of Markov Models

$$p(w_1, w_2, \dots, w_n) = p(w_1)p(w_2 | w_1)p(w_3 | w_1, w_2) \\ p(w_4 | w_1, w_2, w_3) \cdots p(w_n | w_1, \dots, w_{n-1})$$

Markov independence assumption

$$p(w_i | w_1, \dots, w_{i-1}) \approx p(w_i | w_{i-1})$$

# Another View of Markov Models

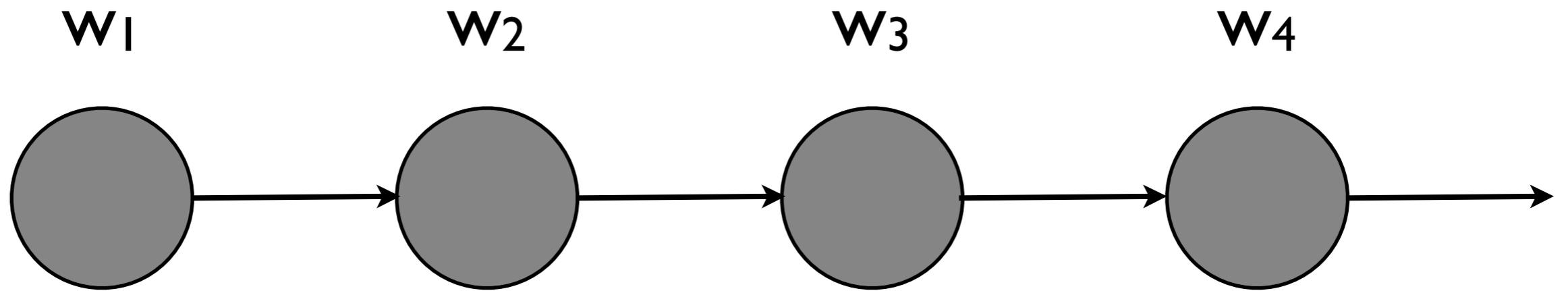
$$p(w_1, w_2, \dots, w_n) = p(w_1)p(w_2 | w_1)p(w_3 | w_1, w_2) \\ p(w_4 | w_1, w_2, w_3) \cdots p(w_n | p_1, \dots, p_{n-1})$$

Markov independence assumption

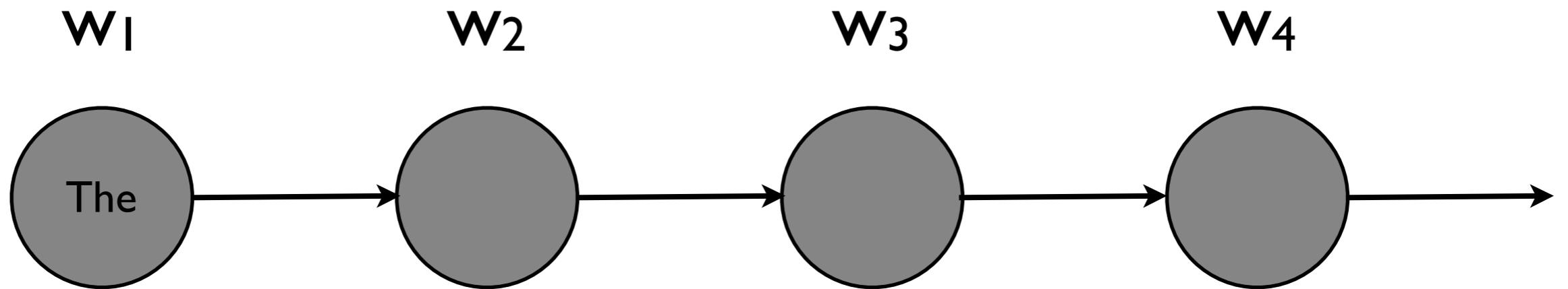
$$p(w_i | w_1, \dots, w_{i-1}) \approx p(w_i | w_{i-1})$$

$$p(w_1, w_2, \dots, w_n) \approx p(w_1)p(w_2 | w_1)p(w_3 | w_2) \\ p(w_4 | w_3) \cdots p(w_n | p_{n-1})$$

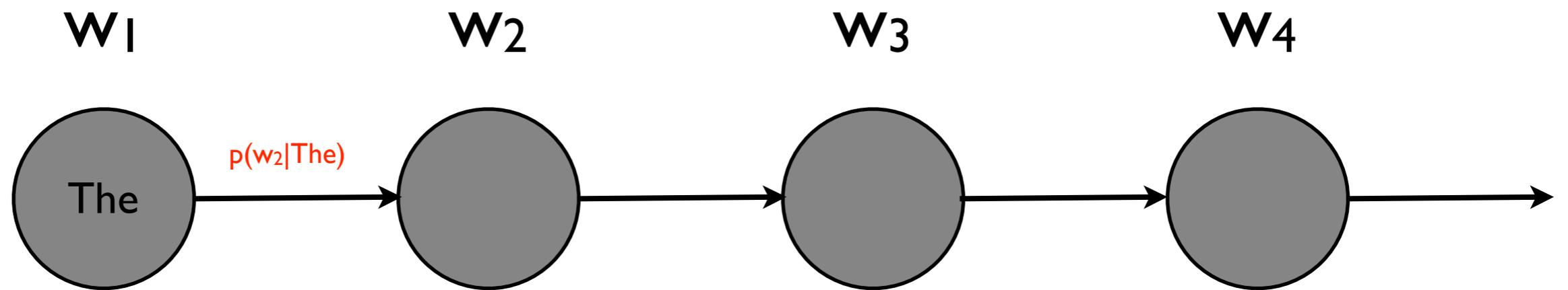
# Yet Another View



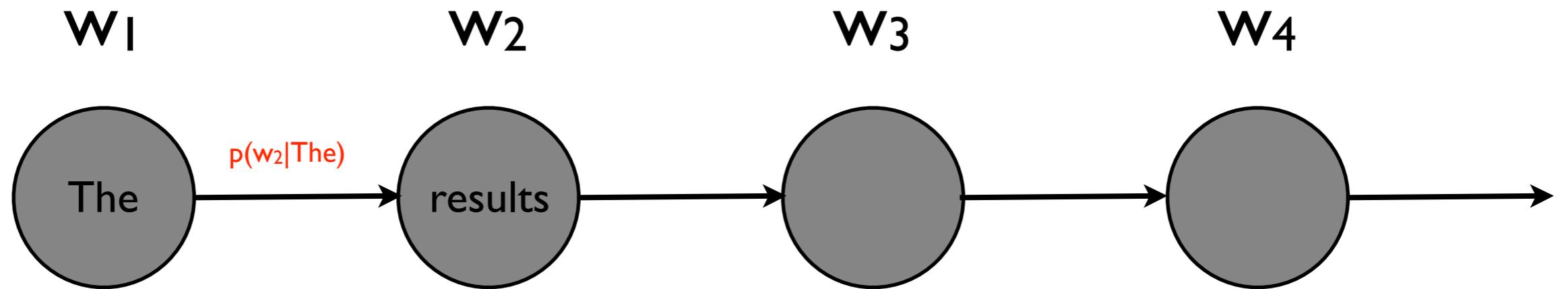
# Yet Another View



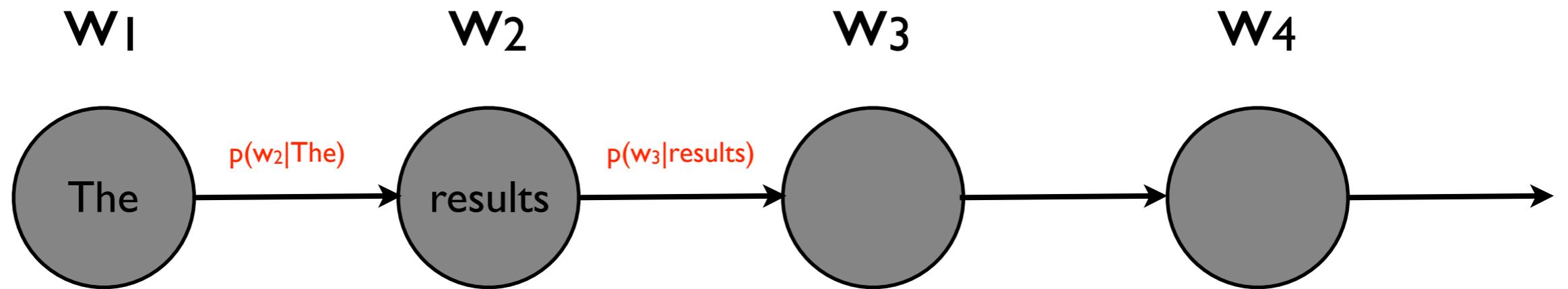
# Yet Another View



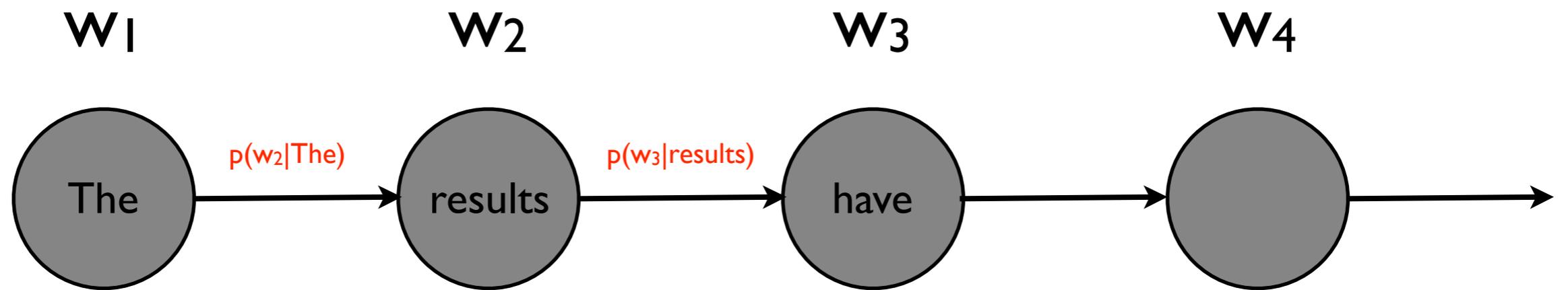
# Yet Another View



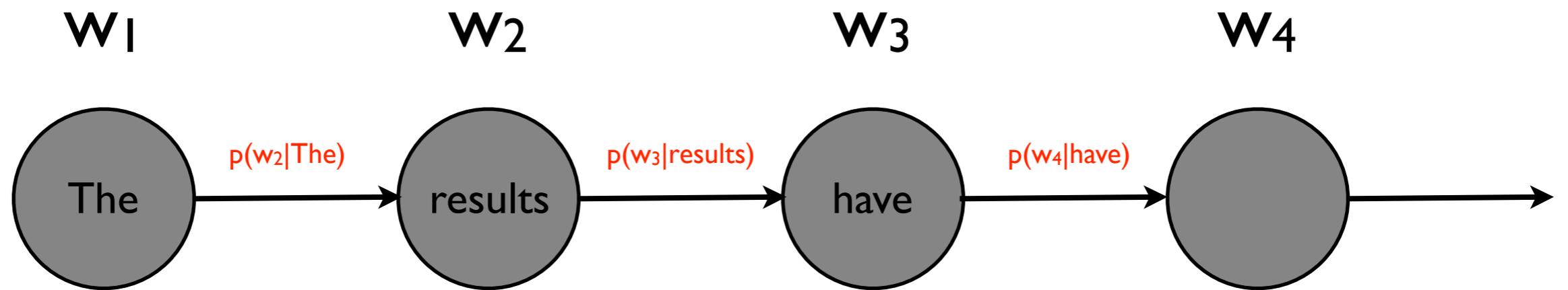
# Yet Another View



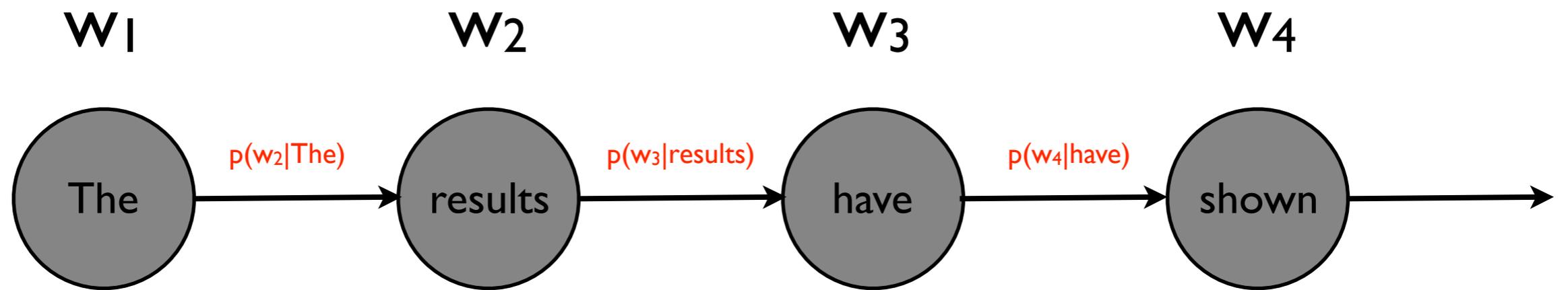
# Yet Another View



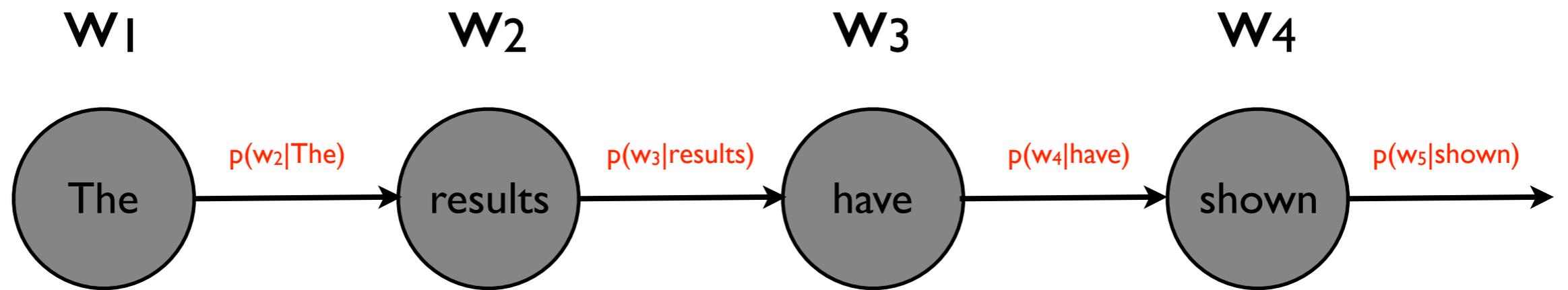
# Yet Another View



# Yet Another View

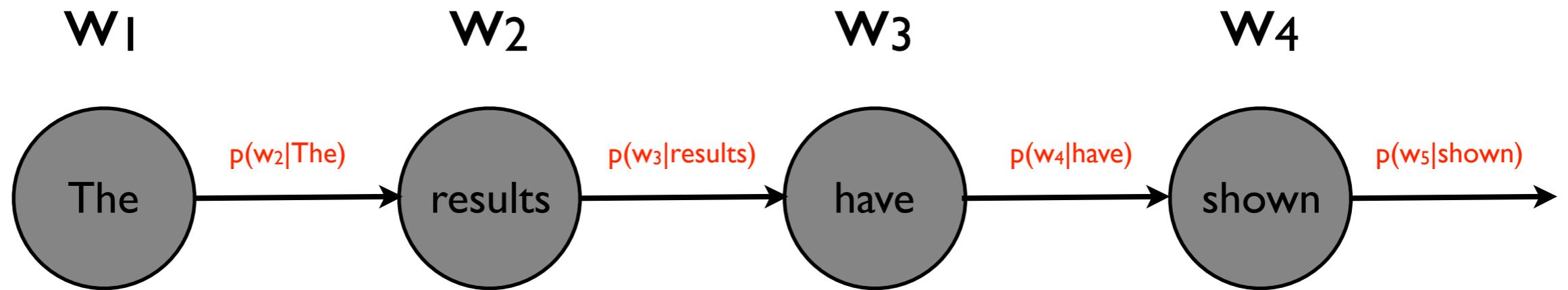


# Yet Another View



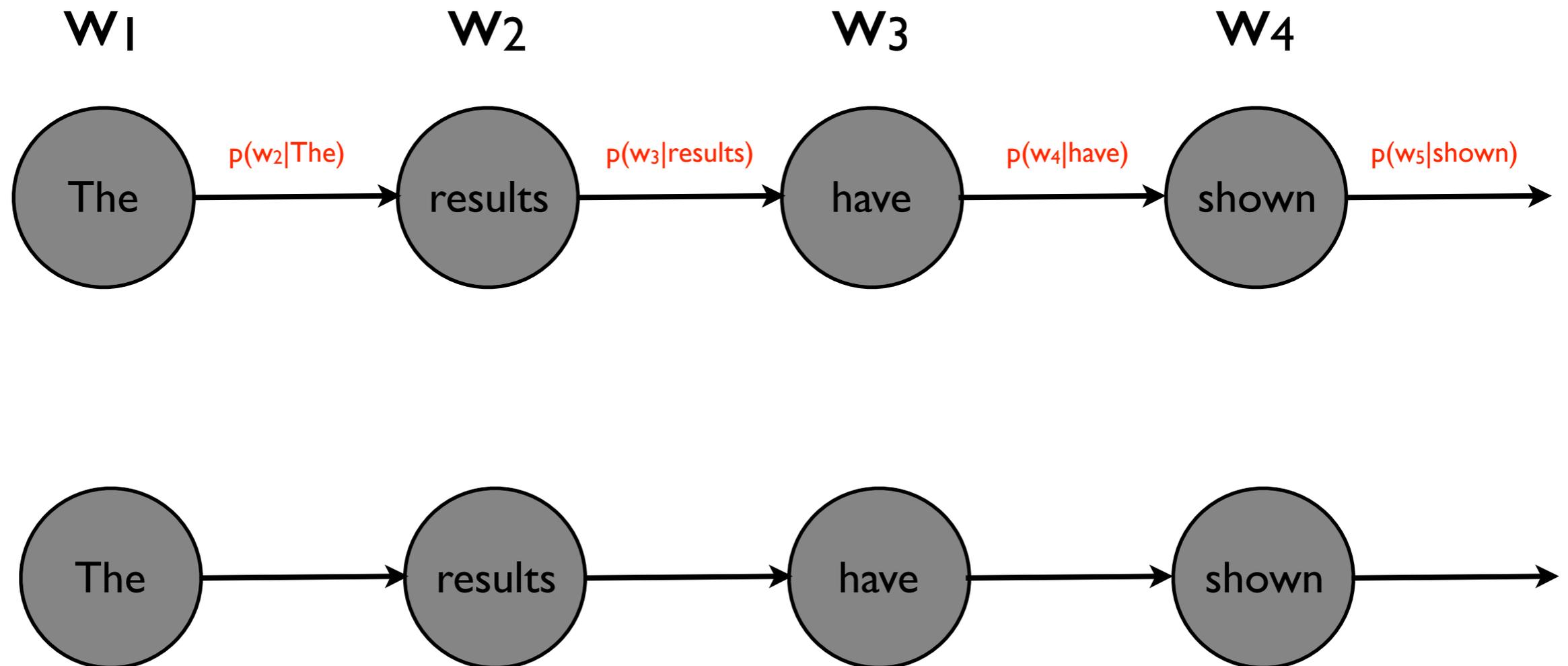
# Yet Another View

Directed graphical models: *lack of edge* means conditional independence



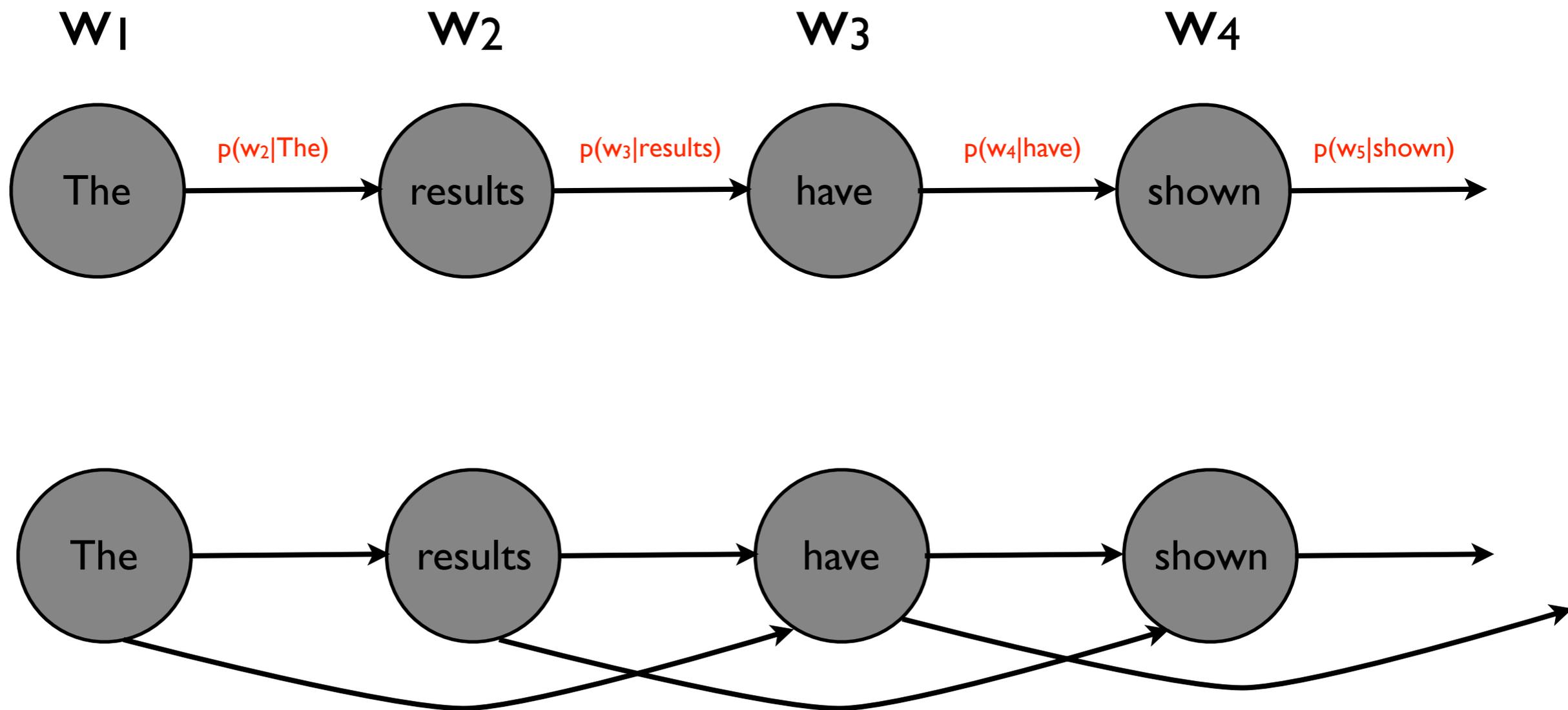
# Yet Another View

Directed graphical models: *lack of edge means conditional independence*



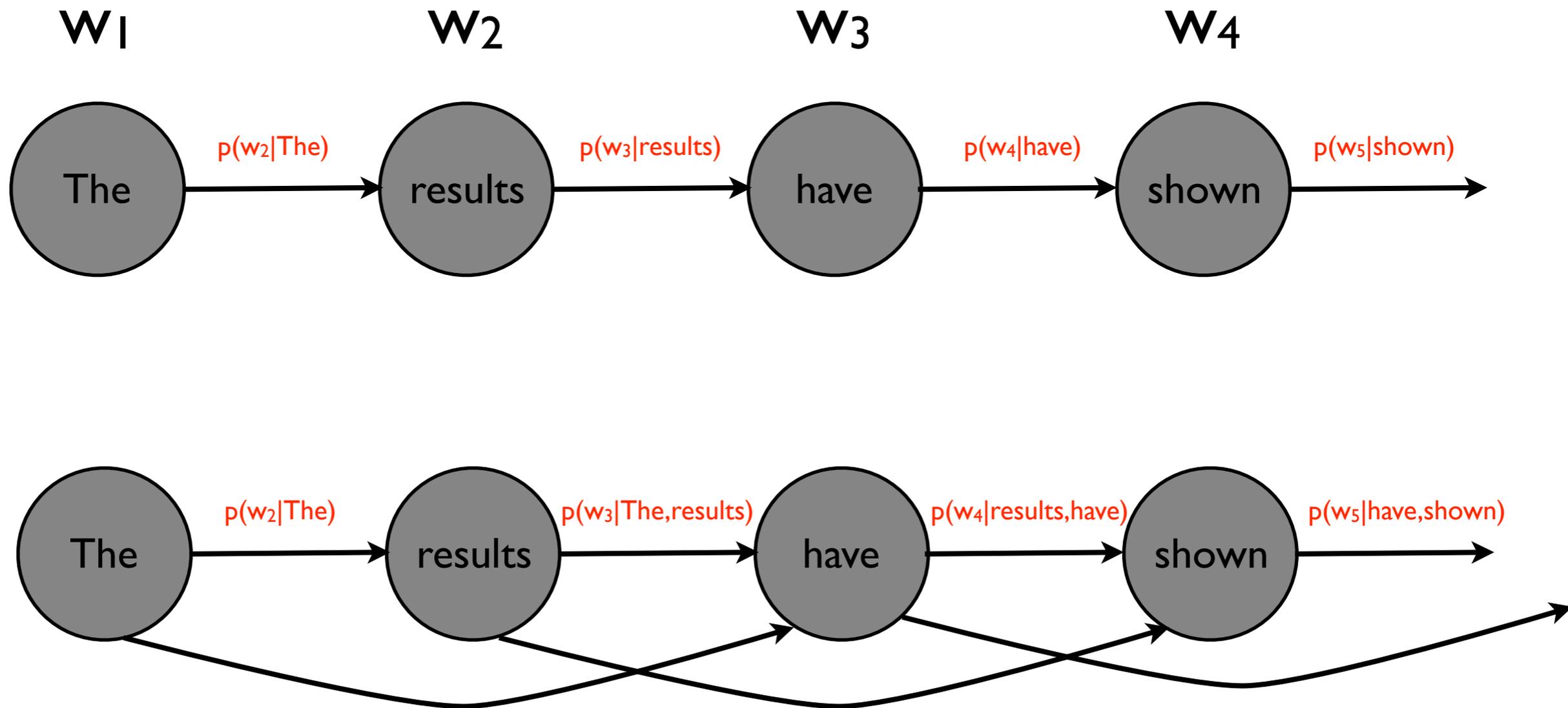
# Yet Another View

Directed graphical models: *lack of edge means conditional independence*

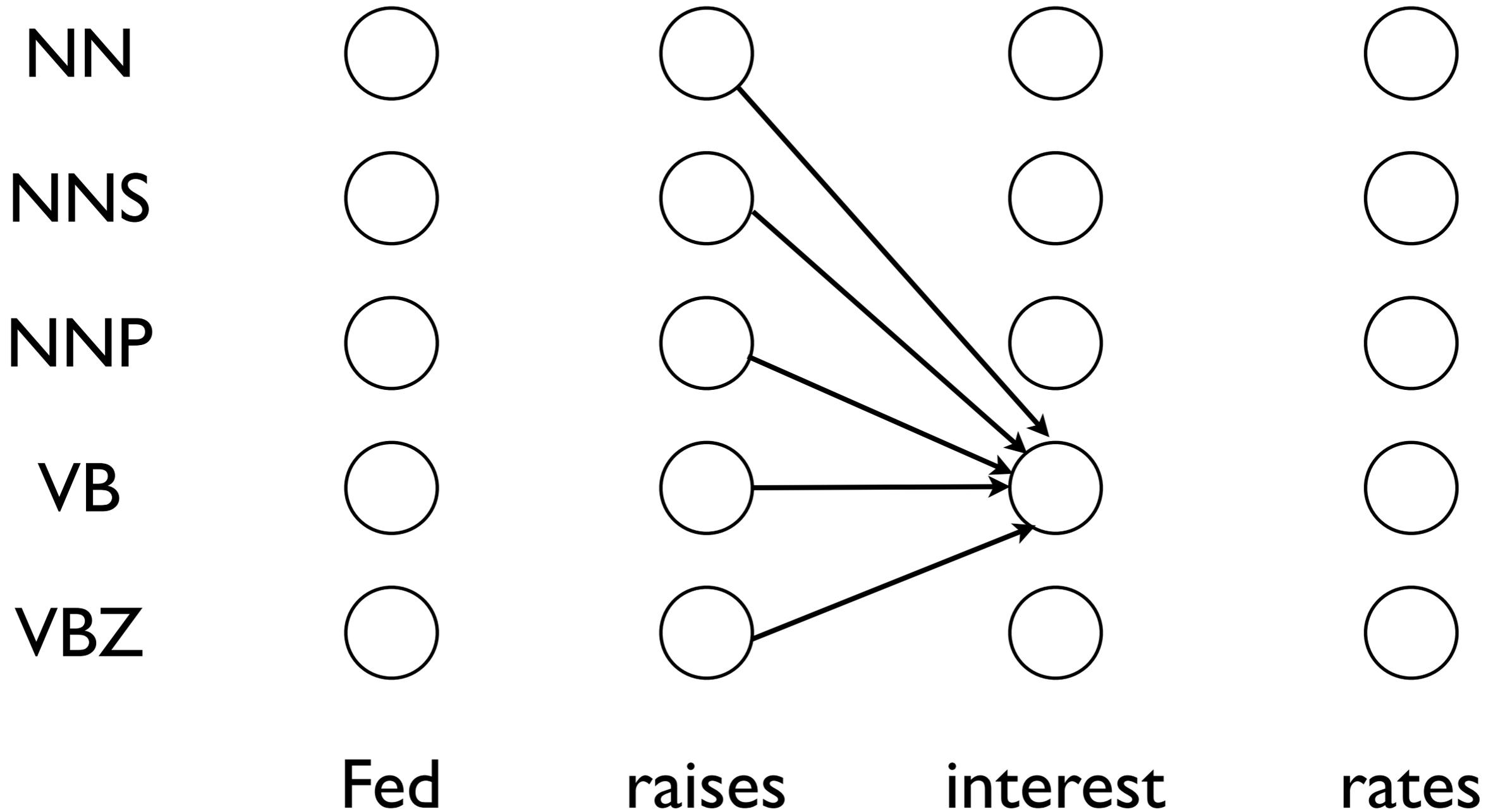


# Yet Another View

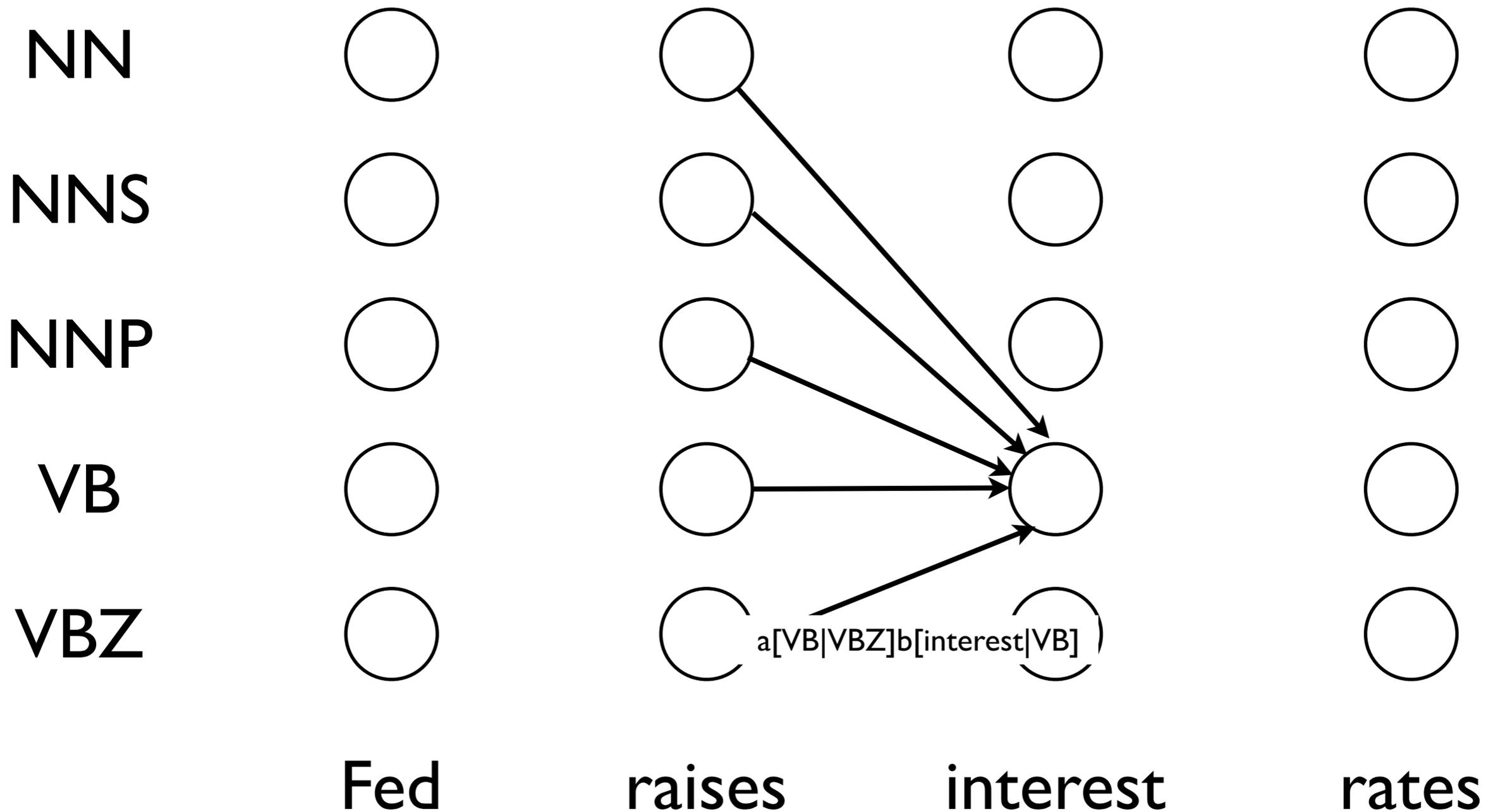
Directed graphical models: *lack of edge means conditional independence*



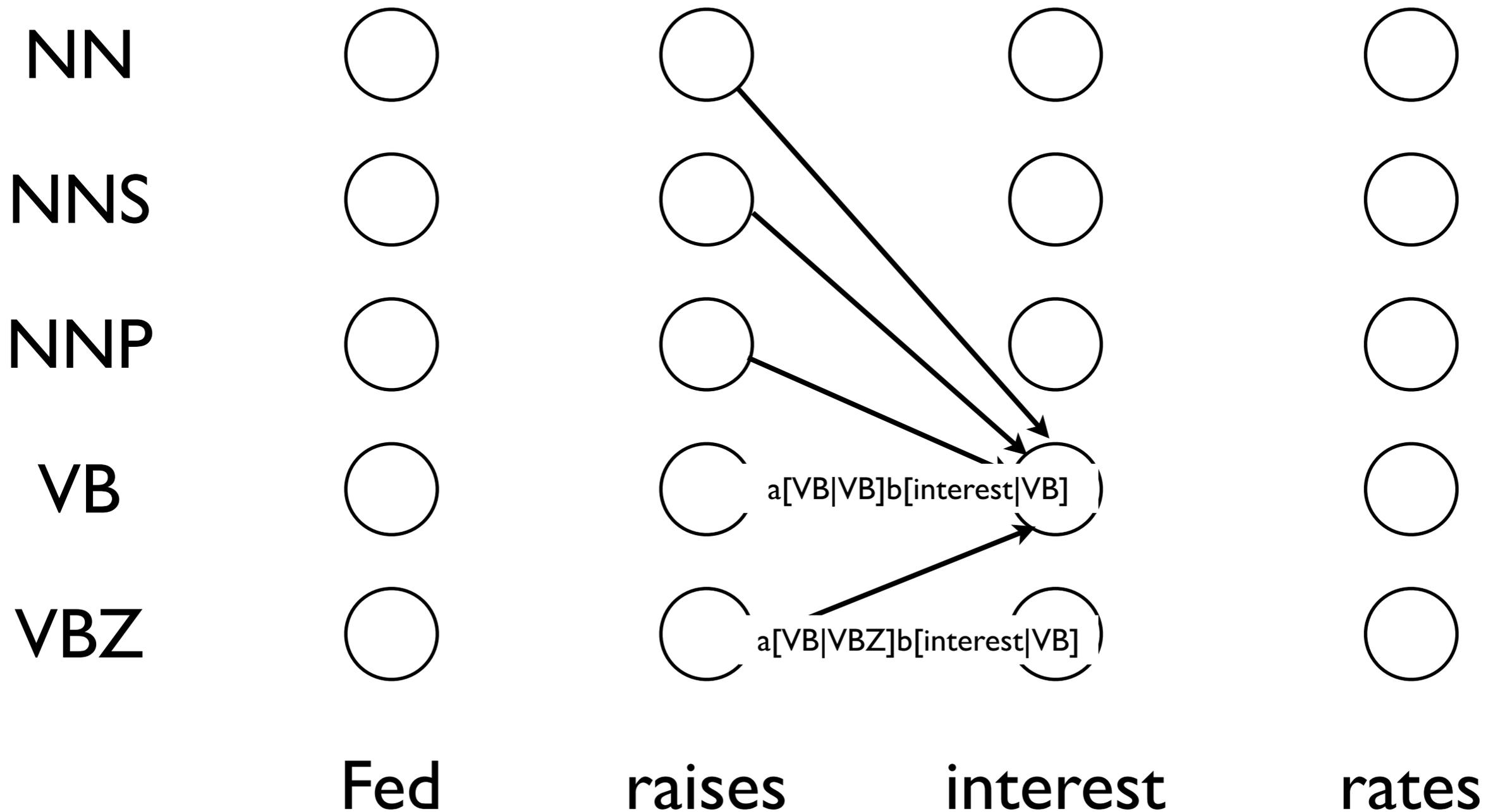
# Forward Algorithm (LM)



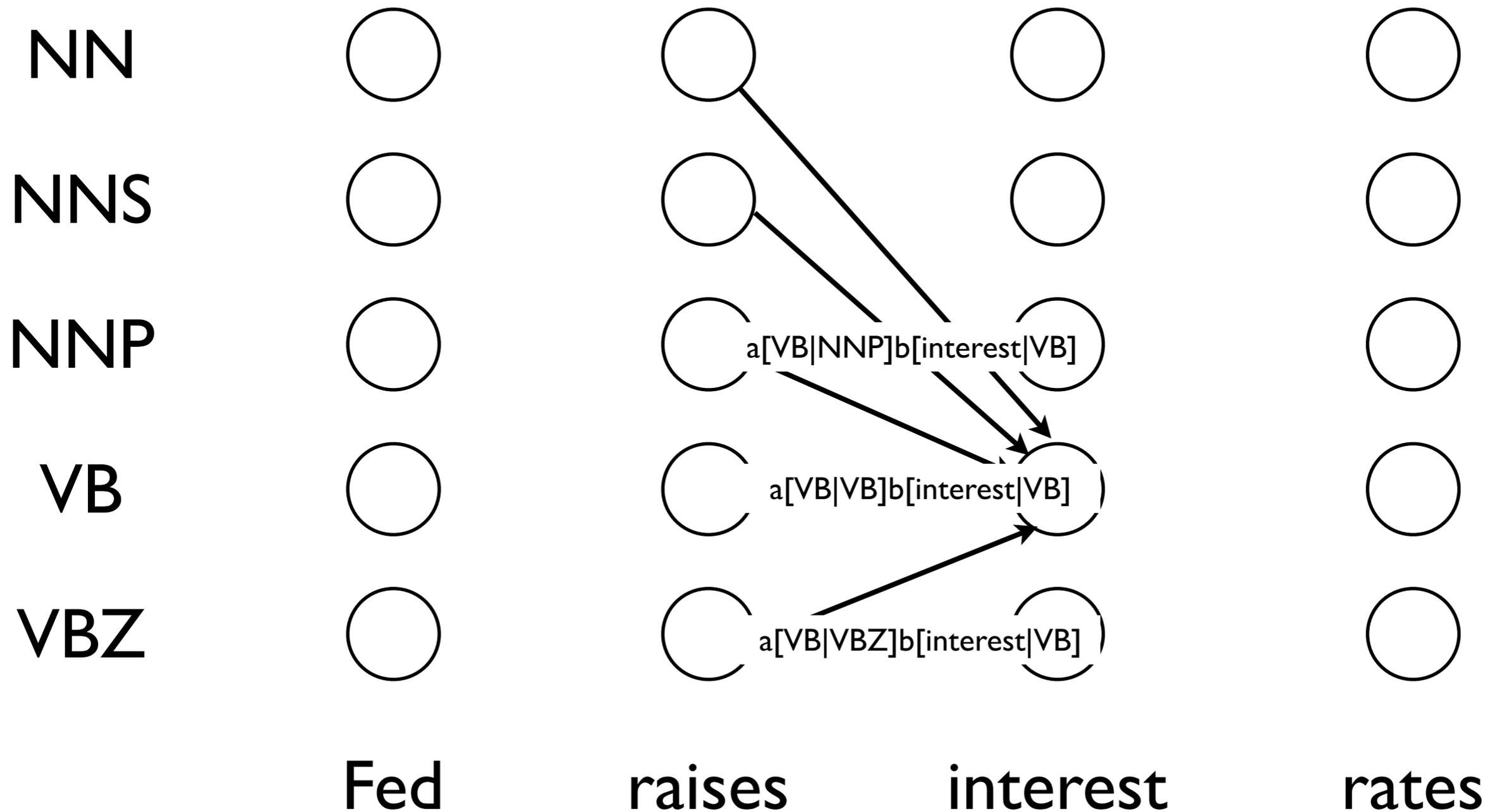
# Forward Algorithm (LM)



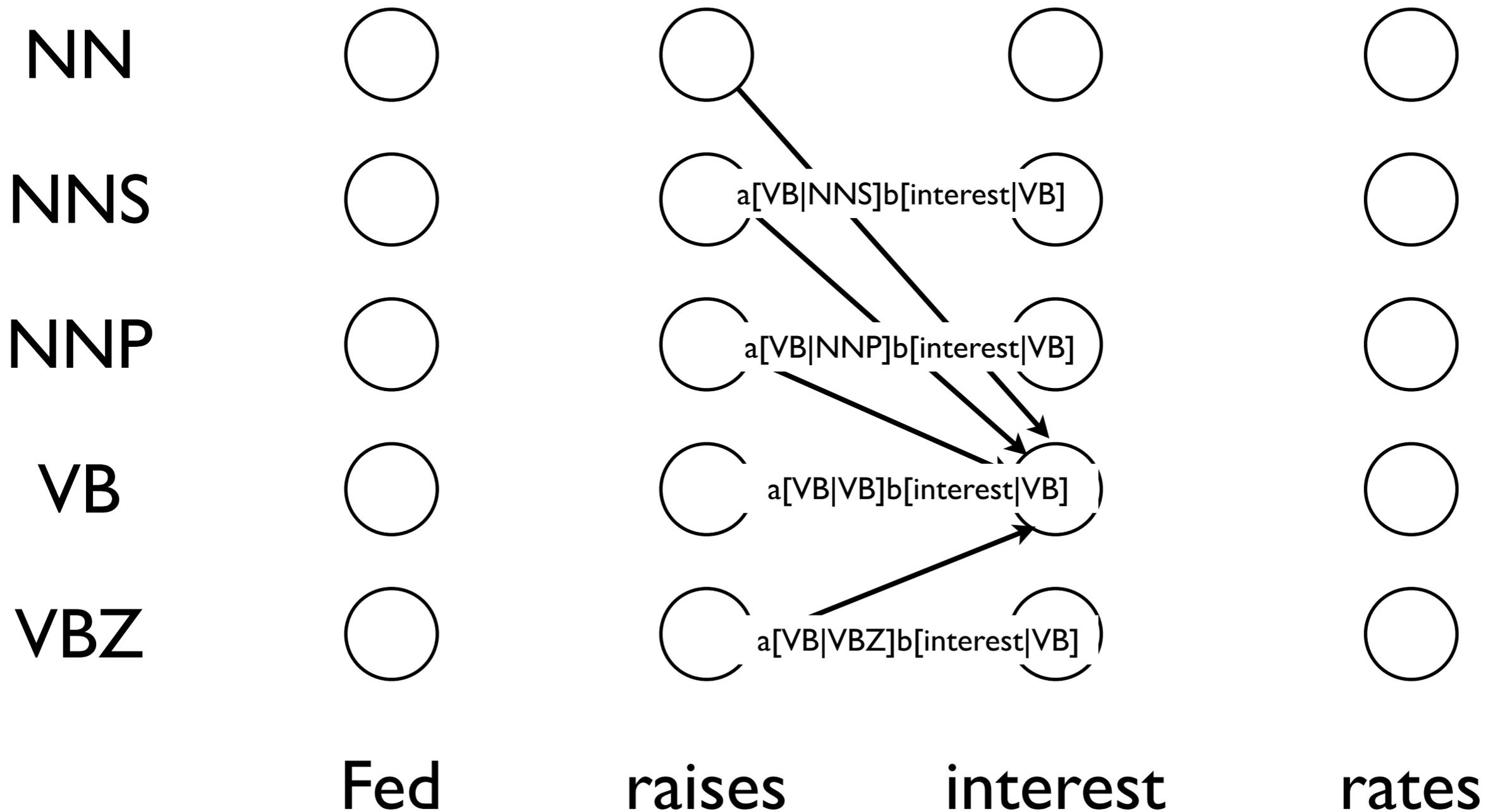
# Forward Algorithm (LM)



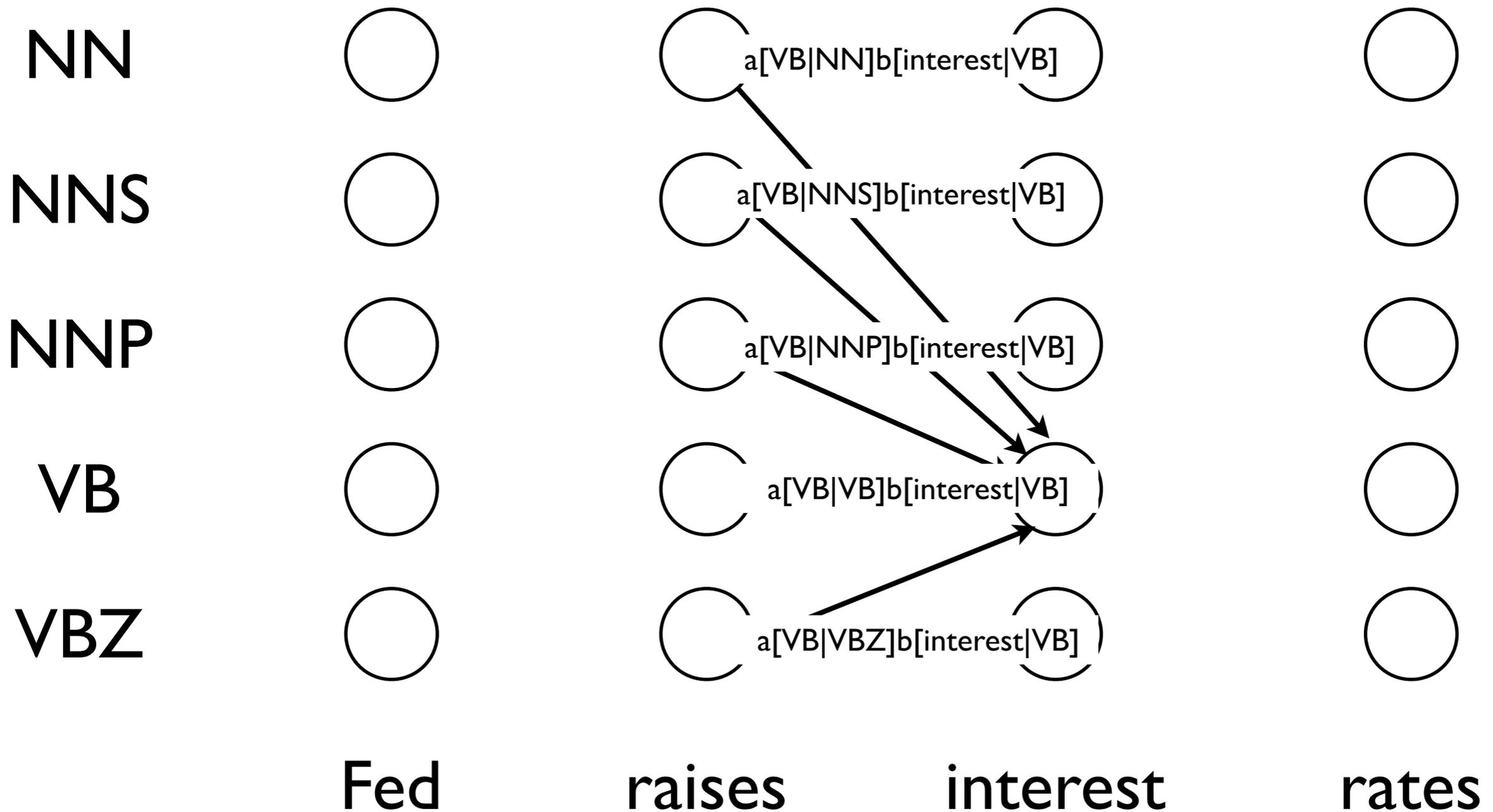
# Forward Algorithm (LM)



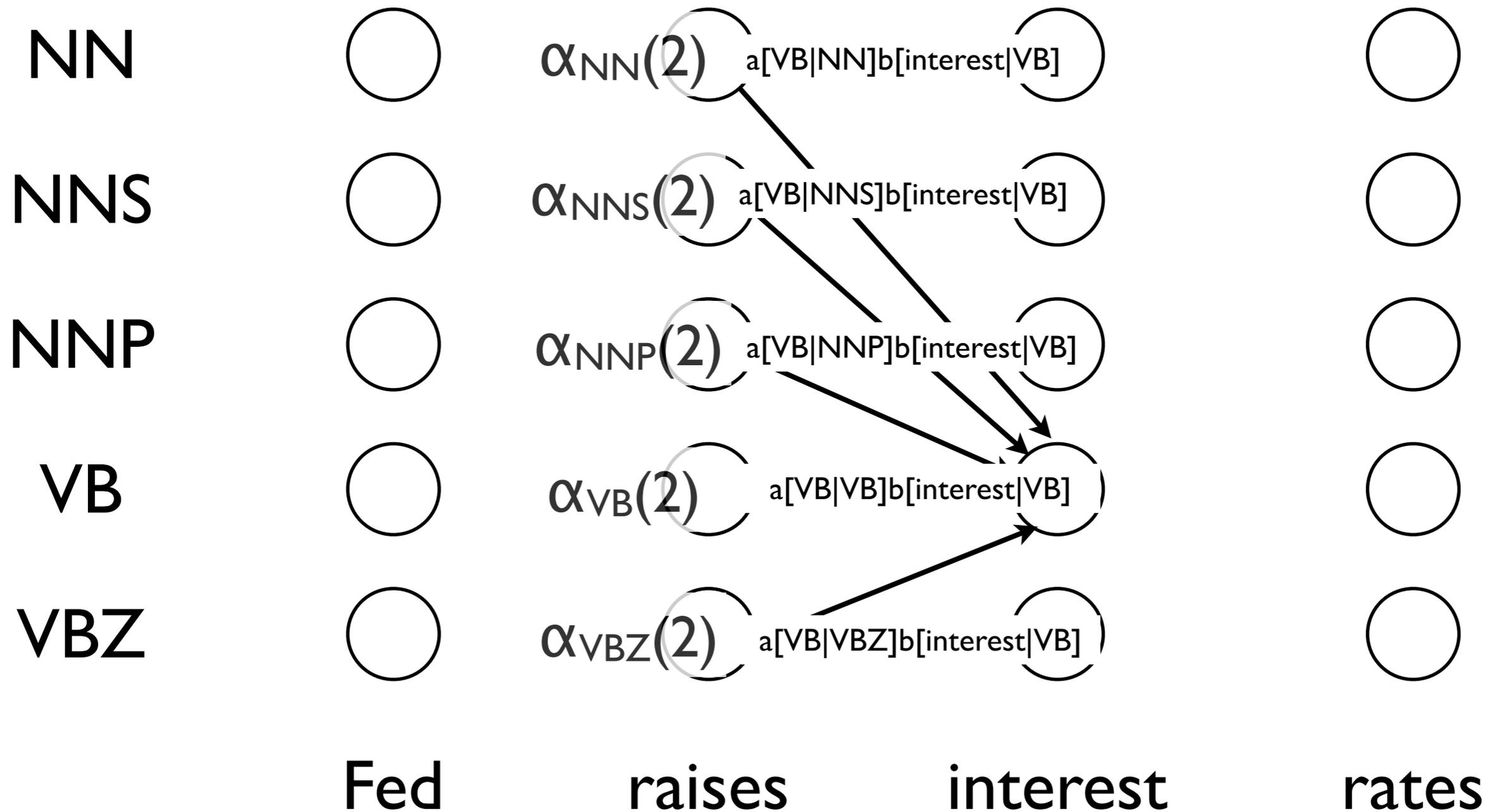
# Forward Algorithm (LM)



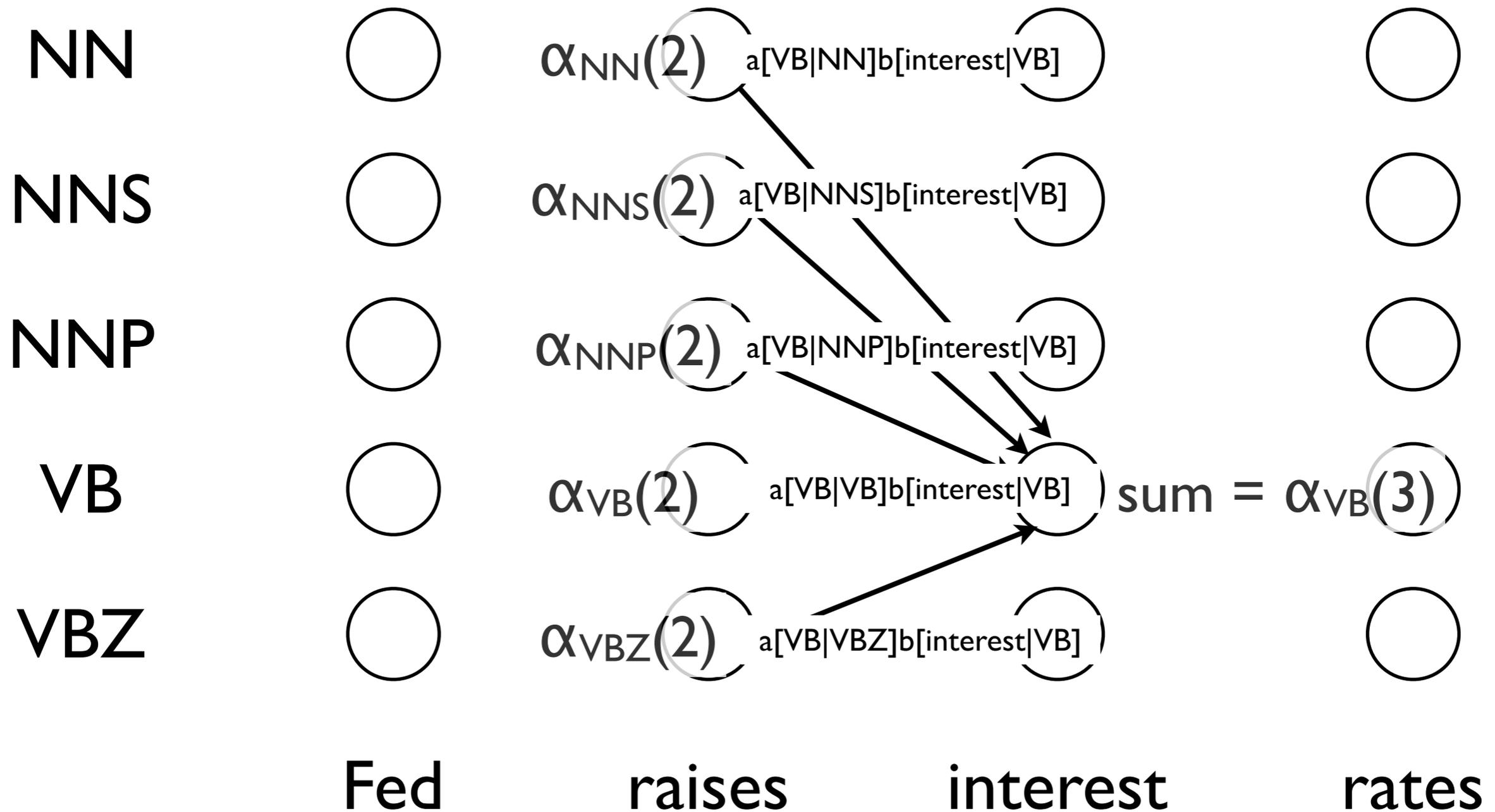
# Forward Algorithm (LM)



# Forward Algorithm (LM)



# Forward Algorithm (LM)



# Setting up a Classifier

# Setting up a Classifier

- What we want:

$$p(\text{😊} \mid w_1, w_2, \dots, w_n) > p(\text{😞} \mid w_1, w_2, \dots, w_n) ?$$

# Setting up a Classifier

- What we want:

$$p(\text{😊} \mid w_1, w_2, \dots, w_n) > p(\text{😞} \mid w_1, w_2, \dots, w_n) ?$$

- What we know how to build:

# Setting up a Classifier

- What we want:

$$p(\text{😊} \mid w_1, w_2, \dots, w_n) > p(\text{😞} \mid w_1, w_2, \dots, w_n) ?$$

- What we know how to build:
  - A language model for each class

# Setting up a Classifier

- What we want:

$$p(\text{😊} \mid w_1, w_2, \dots, w_n) > p(\text{😞} \mid w_1, w_2, \dots, w_n) ?$$

- What we know how to build:
  - A language model for each class
    - $p(w_1, w_2, \dots, w_n \mid \text{😊})$

# Setting up a Classifier

- What we want:

$$p(\text{😊} \mid w_1, w_2, \dots, w_n) > p(\text{😞} \mid w_1, w_2, \dots, w_n) ?$$

- What we know how to build:
  - A language model for each class
    - $p(w_1, w_2, \dots, w_n \mid \text{😊})$
    - $p(w_1, w_2, \dots, w_n \mid \text{😞})$

# Bayes' Theorem

By the definition of conditional probability:

$$P(A, B) = P(B)P(A | B) = P(A)P(B | A)$$

we can show:

$$P(A | B) = \frac{P(B | A)P(A)}{P(B)}$$

Seemingly trivial result from 1763;  
interesting consequences...

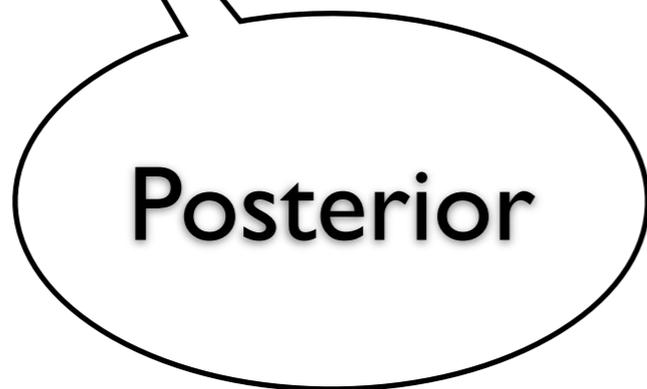


REV. T. BAYES

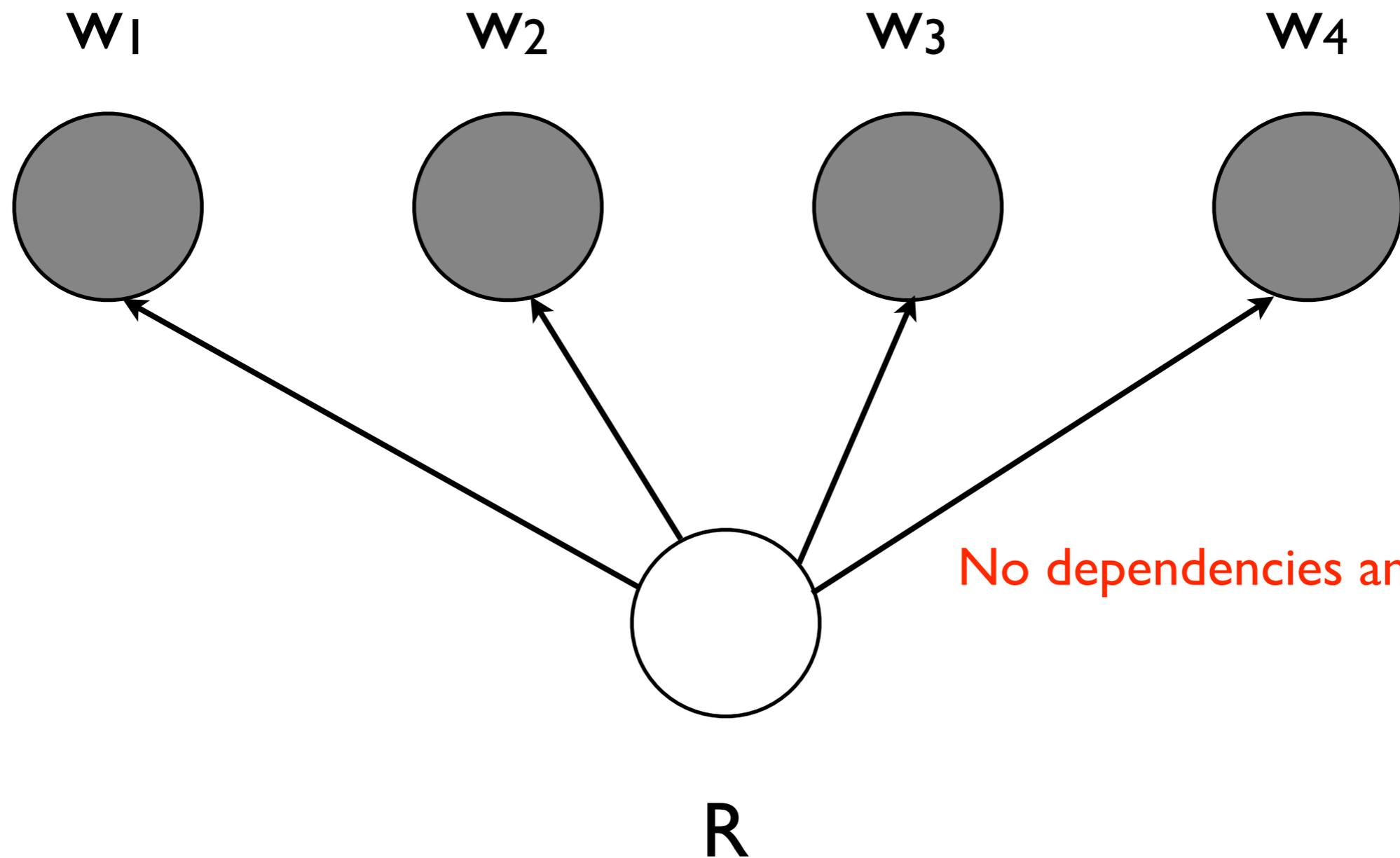
# A “Bayesian” Classifier

$$p(R | w_1, w_2, \dots, w_n) = \frac{p(R)p(w_1, w_2, \dots, w_n | R)}{p(w_1, w_2, \dots, w_n)}$$

$$\max_{R \in \{\overset{\circ}{\smile}, \overset{\circ}{\frown}\}} p(R | w_1, w_2, \dots, w_n) = \max_{R \in \{\overset{\circ}{\smile}, \overset{\circ}{\frown}\}} p(R)p(w_1, w_2, \dots, w_n | R)$$



# *Naive Bayes Classifier*



# NB on Movie Reviews

- Train models for positive, negative
- For each review, find higher posterior
- Which word probability ratios are highest?

```
>>> classifier.show_most_informative_features(5)
```

```
classifier.show_most_informative_features(5)
```

```
Most Informative Features
```

contains(outstanding) = True	pos : neg	=	14.1 : 1.0
contains(mulan) = True	pos : neg	=	8.3 : 1.0
contains(seagal) = True	neg : pos	=	7.8 : 1.0
contains(wonderfully) = True	pos : neg	=	6.6 : 1.0
contains(damon) = True	pos : neg	=	6.1 : 1.0

# What's Wrong With NB?

- What happens for word dependencies are strong?
- What happens when some words occur only once?
- What happens when the classifier sees a new word?

# Generative vs. Conditional

- What is the most likely label for a given input?
- How likely is a given label for a given input?
- What is the most likely input value?
- How likely is a given input value?
- How likely is a given input value with a given label?
- What is the most likely label for an input that might have one of two values (but we don't know which)?

# Generative vs. Conditional

- What is the most likely label for a given input?
- How likely is a given label for a given input?
- What is the most likely input value?
- How likely is a given input value?
- How likely is a given input value with a given label?
- What is the most likely label for an input that might have one of two values (but we don't know which)?

# Sequence Labeling

- Inputs:  $x = (x_1, \dots, x_n)$
- Labels:  $y = (y_1, \dots, y_n)$
- Typical goal: Given  $x$ , predict  $y$
  
- Example sequence labeling tasks
  - Part-of-speech tagging
  - Named-entity-recognition (NER)
    - Label people, places, organizations

# NER Example:

## Red Sox and Their Fans Let Loose



Elise Amendola/Associated Press

Fans of the slugger David Ortiz in Boston's Copley Square.

By [PETE THAMEL](#)

Published: October 31, 2007

[BOSTON](#), Oct. 30 — [Jonathan Papelbon](#) turned Boston's World Series victory parade into a full-scale dance party Tuesday as the [Red Sox](#) put an exclamation point on the 2007 season.

 E-MAIL

 PRINT

 REPRINTS

 SAVE

# First Solution:

## Maximum Entropy Classifier

- Conditional model  $p(y|x)$ .
  - Do not waste effort modeling  $p(x)$ , since  $x$  is given at test time anyway.
  - Allows more complicated input features, since we do not need to model dependencies between them.
- Feature functions  $f(x,y)$ :
  - $f_1(x,y) = \{ \text{word is Boston} \ \& \ y=\text{Location} \}$
  - $f_2(x,y) = \{ \text{first letter capitalized} \ \& \ y=\text{Name} \}$
  - $f_3(x,y) = \{ x \text{ is an HTML link} \ \& \ y=\text{Location} \}$

# First Solution: MaxEnt Classifier

- How should we choose a classifier?
- Principle of maximum entropy
  - We want a classifier that:
    - Matches feature constraints from training data.
    - Predictions maximize entropy.
- There is a unique, exponential family distribution that meets these criteria.

# First Solution: MaxEnt Classifier

- Problem with using a maximum entropy classifier for sequence labeling:
- It makes decisions at each position independently!

## Second Solution: HMM

$$P(\mathbf{y}, \mathbf{x}) = \prod_t P(y_t | y_{t-1}) P(x | y_t)$$

- Defines a generative process.
- Can be viewed as a weighted finite state machine.

## Second Solution: HMM

- HMM problems: (ON BOARD)
  - Probability of an input sequence.
  - Most likely label sequence given an input sequence.
  - Learning with known label sequences.
  - Learning with unknown label sequences?

## Second Solution: HMM

- How can represent we multiple features in an HMM?
  - Treat them as conditionally independent given the class label?
    - The example features we talked about are not independent.
  - Try to model a more complex generative process of the input features?
    - We may lose tractability (i.e. lose a dynamic programming for exact inference).

## Second Solution: HMM

- Let's use a conditional model instead.

# Third Solution: MEMM

- Use a series of maximum entropy classifiers that know the previous label.
- Define a Viterbi algorithm for inference.

$$P(\mathbf{y} \mid \mathbf{x}) = \prod_t P(y_t \mid y_{t-1}, \mathbf{x})$$

## Third Solution: MEMM

- Combines the advantages of maximum entropy and HMM!
- But there is a problem...

# Problem with MEMMs: Label Bias

- In some state space configurations, MEMMs essentially completely ignore the inputs.
- Example (ON BOARD).
- This is not a problem for HMMs, because the input sequence is generated by the model.

# Fourth Solution: Conditional Random Field

- Conditionally-trained, undirected graphical model.
- For a standard linear-chain structure:

$$P(\mathbf{y} \mid \mathbf{x}) = \frac{1}{Z} \prod_t \Psi_t(y_t, y_{t-1}, \mathbf{x})$$

Bigram model:  
potentials consider  
pairs of labels

$$\Psi_t(y_t, y_{t-1}, \mathbf{x}) = \exp \left[ \sum_k \lambda_k f_k(y_t, y_{t-1}, \mathbf{x}) \right]$$

Dot-product of  
weights and features

$$Z = \sum_{\mathbf{y}'} \prod_t \Psi(y'_t, y'_{t-1}, \mathbf{x})$$

Normalize over all  
possible outputs  
using forward alg.

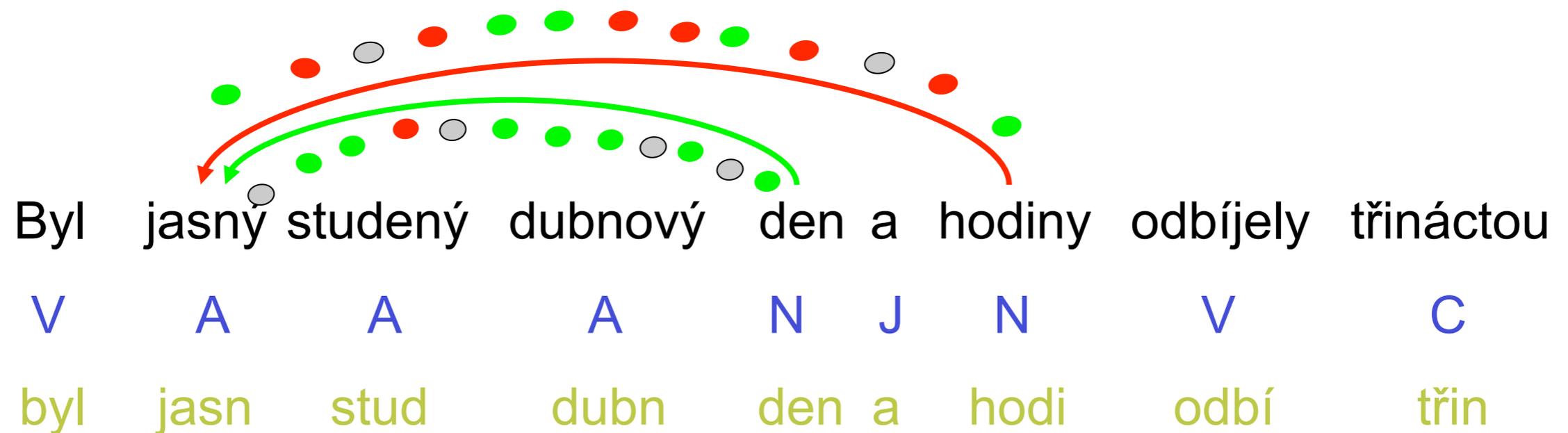
## Fourth Solution: CRF

- Have the advantages of MEMMs, but avoid the label bias problem.
- CRFs are globally normalized, whereas MEMMs are locally normalized.
- Widely used and applied. CRFs give state-the-art results in many domains.

# Example Applications

- CRFs have been applied to:
  - Part-of-speech tagging
  - Named-entity-recognition
  - Table extraction
  - Gene prediction
  - Chinese word segmentation
  - Extracting information from research papers.
  - Many more...

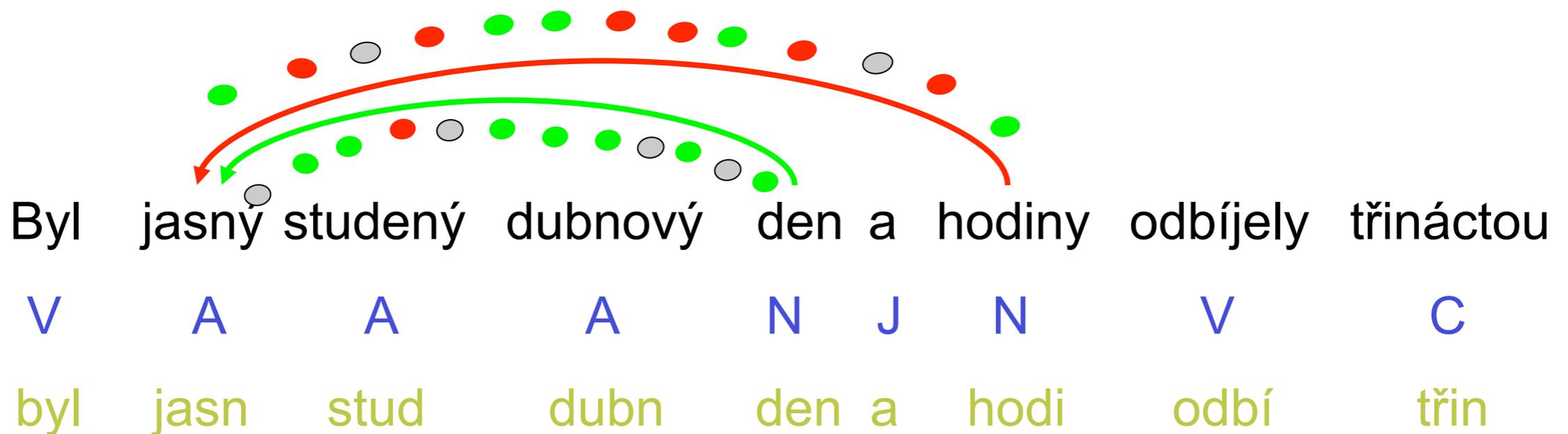
# Edge-Factored Parsers (McDonald et al. 2005)



“It was a bright cold day in April and the clocks were striking thirteen”

# Edge-Factored Parsers (McDonald et al. 2005)

- Which edge is better?



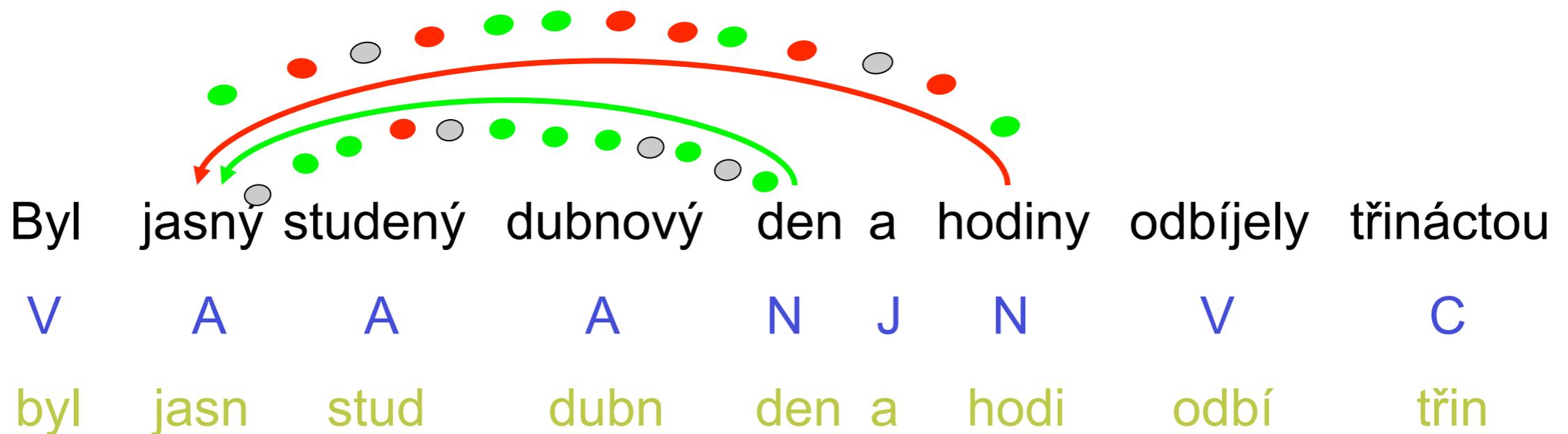
“It was a bright cold day in April and the clocks were striking thirteen”

# Edge-Factored Parsers (McDonald et al. 2005)

- Which edge is better?



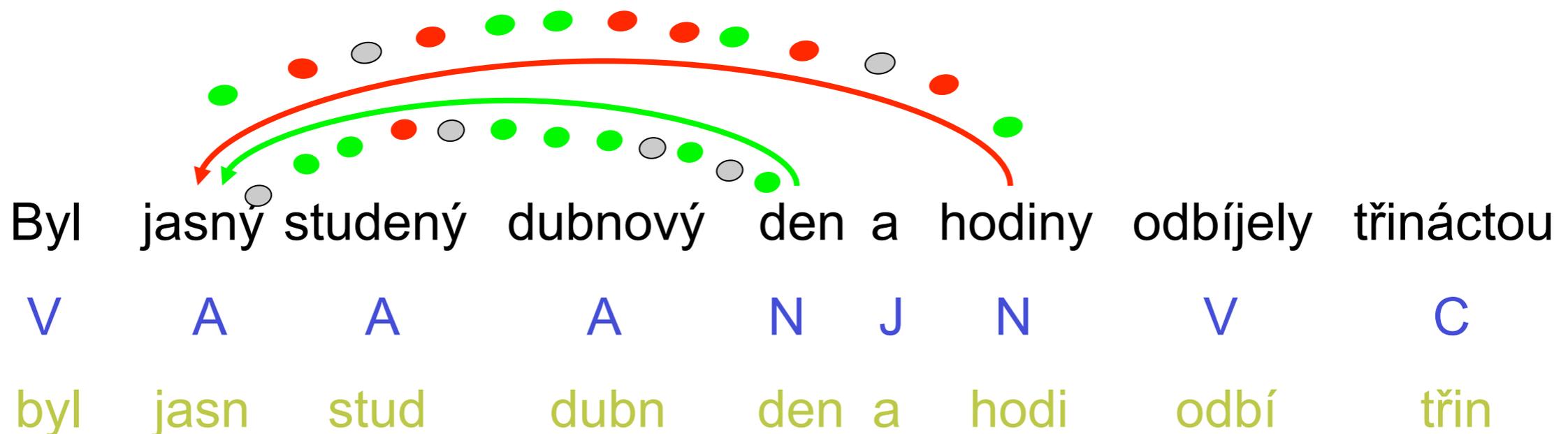
our current weight vector



“It was a bright cold day in April and the clocks were striking thirteen”

# Edge-Factored Parsers (McDonald et al. 2005)

- Which edge is better?
  - Score of an edge  $e = \theta \cdot \text{features}(e)$
- our current weight vector

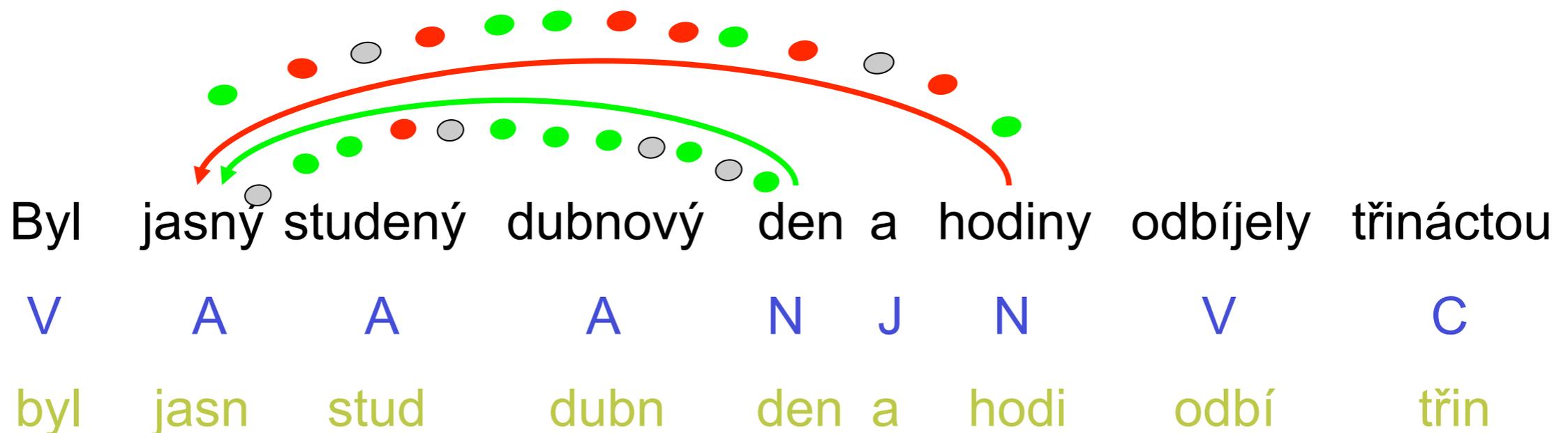


“It was a bright cold day in April and the clocks were striking thirteen”

# Edge-Factored Parsers (McDonald et al. 2005)

- Which edge is better?
- Score of an edge  $e = \theta \cdot \text{features}(e)$
- Standard algos  $\rightarrow$  valid parse with max total score

our current weight vector



“It was a bright cold day in April and the clocks were striking thirteen”

# Recipe for Conditional Training of $p(y | x)$

1. Gather constraints/features from training data

$$\alpha_{iy} = \tilde{E}[f_{iy}] = \sum_{x_j, y_j \in D} f_{iy}(x_j, y_j)$$

2. Initialize all parameters to zero

3. Classify training data with current parameters; calculate expectations

$$E_{\Theta}[f_{iy}] = \sum_{x_j \in D} \sum_{y'} p_{\Theta}(y' | x_j) f_{iy}(x_j, y')$$

4. Gradient is  $\tilde{E}[f_{iy}] - E_{\Theta}[f_{iy}]$

5. Take a step in the direction of the gradient

6. Repeat from 3 until convergence

# Recipe for Conditional Training of $p(y | x)$

1. Gather constraints/features from training data

$$\alpha_{iy} = \tilde{E}[f_{iy}] = \sum_{x_j, y_j \in D} f_{iy}(x_j, y_j)$$

2. Initialize all parameters to zero

3. Classify training data with current parameters; calculate expectations

$$E_{\Theta}[f_{iy}] = \sum_{x_j \in D} \sum_{y'} p_{\Theta}(y' | x_j) f_{iy}(x_j, y')$$

4. Gradient is  $\tilde{E}[f_{iy}] - E_{\Theta}[f_{iy}]$

5. Take a step in the direction of the gradient

6. Repeat from 3 until convergence

Where have we seen expected counts before?

# Recipe for Conditional Training of $p(y | x)$

1. Gather constraints/features from training data

$$\alpha_{iy} = \tilde{E}[f_{iy}] = \sum_{x_j, y_j \in D} f_{iy}(x_j, y_j)$$

2. Initialize all parameters to zero

3. Classify training data with current parameters; calculate expectations

$$E_{\Theta}[f_{iy}] = \sum_{x_j \in D} \sum_{y'} p_{\Theta}(y' | x_j) f_{iy}(x_j, y')$$

4. Gradient is  $\tilde{E}[f_{iy}] - E_{\Theta}[f_{iy}]$

5. Take a step in the direction of the gradient

**EM!**

6. Repeat from 3 until convergence

Where have we seen expected counts before?

# Gradient-Based Training

- $\lambda := \lambda + \text{rate} * \text{Gradient}(F)$
- After all training examples? (batch)
- After every example? (on-line)
- Use second derivative for faster learning?
- A big field: numerical optimization

# Overfitting

- If we have too many features, we can choose weights to model the training data perfectly
- If we have a feature that only appears in spam training, not ham training, it will get weight  $\infty$  to maximize  $p(\text{spam} \mid \text{feature})$  at 1.
- These behaviors
  - Overfit the training data
  - Will probably do poorly on test data

# Solutions to Overfitting

- Throw out rare features.
  - Require every feature to occur  $> 4$  times, and  $> 0$  times with legit, and  $> 0$  times with spam.
- Only keep, e.g., 1000 features.
  - Add one at a time, always greedily picking the one that most improves performance on held-out data.
- Smooth the observed feature counts.
- Smooth the weights by using a prior.
  - $\max p(\lambda|\text{data}) = \max p(\lambda, \text{data}) = p(\lambda)p(\text{data}|\lambda)$
  - decree  $p(\lambda)$  to be high when most weights close to 0

# Smoothing with Priors

- What if we had a prior expectation that parameter values wouldn't be very large?
- We could then balance evidence suggesting large (or infinite) parameters against our prior expectation.
- The evidence would never totally defeat the prior, and parameters would be smoothed (and kept finite)
- We can do this explicitly by changing the optimization objective to maximum posterior likelihood:

$$\log P(y, \lambda \mid x) = \log P(\lambda) + \log P(y \mid x, \lambda)$$

Posterior

Prior

Likelihood

# **(First Order) Logic**

## **Some Preliminaries**

# **(First Order) Logic**

## **Some Preliminaries**

Three major kinds of objects

# **(First Order) Logic**

## **Some Preliminaries**

Three major kinds of objects

1. Booleans

- Roughly, the semantic values of sentences

# (First Order) Logic

## Some Preliminaries

### Three major kinds of objects

#### 1. Booleans

- Roughly, the semantic values of sentences

#### 2. Entities

- Values of NPs, e.g., objects like this slide
- Maybe also other types of entities, like times

# (First Order) Logic

## Some Preliminaries

### Three major kinds of objects

#### 1. Booleans

- Roughly, the semantic values of sentences

#### 2. Entities

- Values of NPs, e.g., objects like this slide
- Maybe also other types of entities, like times

#### 3. Functions of various types

- Functions from booleans to booleans (and, or, not)
- A function from entity to boolean is called a “predicate” – e.g., `frog(x)`, `green(x)`
- Functions might return other functions!

# (First Order) Logic

## Some Preliminaries

### Three major kinds of objects

#### 1. Booleans

- Roughly, the semantic values of sentences

#### 2. Entities

- Values of NPs, e.g., objects like this slide
- Maybe also other types of entities, like times

#### 3. Functions of various types

- Functions from booleans to booleans (and, or, not)
- A function from entity to boolean is called a “predicate” – e.g., `frog(x)`, `green(x)`
- Functions might return other functions!
- Function might take other functions as arguments!

# First-order Representations

# First-order Representations

- Gilly swallowed a goldfish
  - First attempt: `swallowed(Gilly, goldfish)`

# First-order Representations

- Gilly swallowed a goldfish
  - First attempt: `swallowed(Gilly, goldfish)`
- Better:  $\exists g$  `goldfish(g) AND swallowed(Gilly, g)`

# First-order Representations

- Gilly swallowed a goldfish
  - First attempt: `swallowed(Gilly, goldfish)`
- Better:  $\exists g$  `goldfish(g) AND swallowed(Gilly, g)`
- Or using one of our quantifier predicates:
  - `exists( $\lambda g$  goldfish(g),  $\lambda g$  swallowed(Gilly,g))`
  - Equivalently: `exists(goldfish, swallowed(Gilly))`
    - “In the set of goldfish there exists one swallowed by Gilly”

# First-order Representations

- `Gilly` swallowed a `goldfish`
  - First attempt: `swallowed(Gilly, goldfish)`
- Better:  $\exists g$  `goldfish(g)` AND `swallowed(Gilly, g)`
- Or using one of our quantifier predicates:
  - `exists( $\lambda g$  goldfish(g),  $\lambda g$  swallowed(Gilly,g))`
  - Equivalently: `exists(goldfish, swallowed(Gilly))`
    - “In the set of goldfish there exists one swallowed by Gilly”
- Here `goldfish` is a predicate on entities
  - This is the same semantic type as `red`
  - But note: `goldfish` is noun and `red` is adjective

# Compositional Semantics

# Compositional Semantics

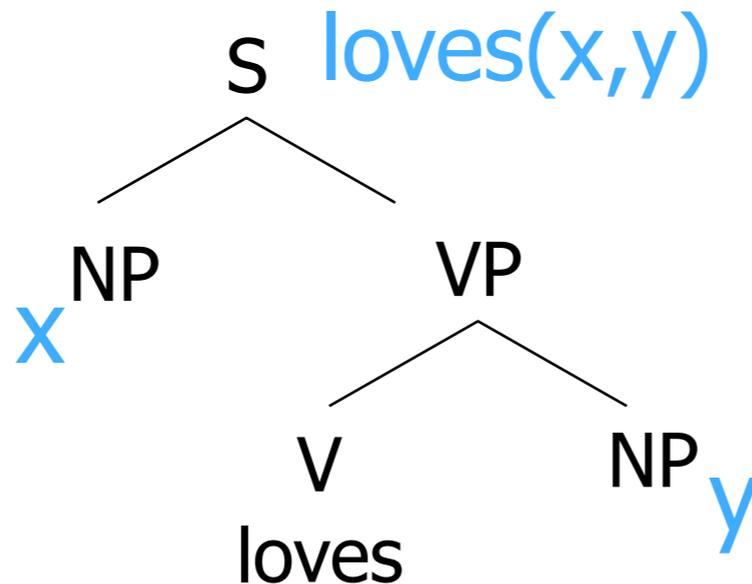
- Add a “sem” feature to each context-free rule
  - $S \rightarrow NP \text{ loves } NP$
  - $S[\text{sem}=\text{loves}(x,y)] \rightarrow NP[\text{sem}=x] \text{ loves } NP[\text{sem}=y]$
  - Meaning of S depends on meaning of NPs

# Compositional Semantics

- Add a “sem” feature to each context-free rule
  - $S \rightarrow NP \text{ loves } NP$
  - $S[\text{sem}=\text{loves}(x,y)] \rightarrow NP[\text{sem}=x] \text{ loves } NP[\text{sem}=y]$
  - Meaning of S depends on meaning of NPs
- TAG version:

# Compositional Semantics

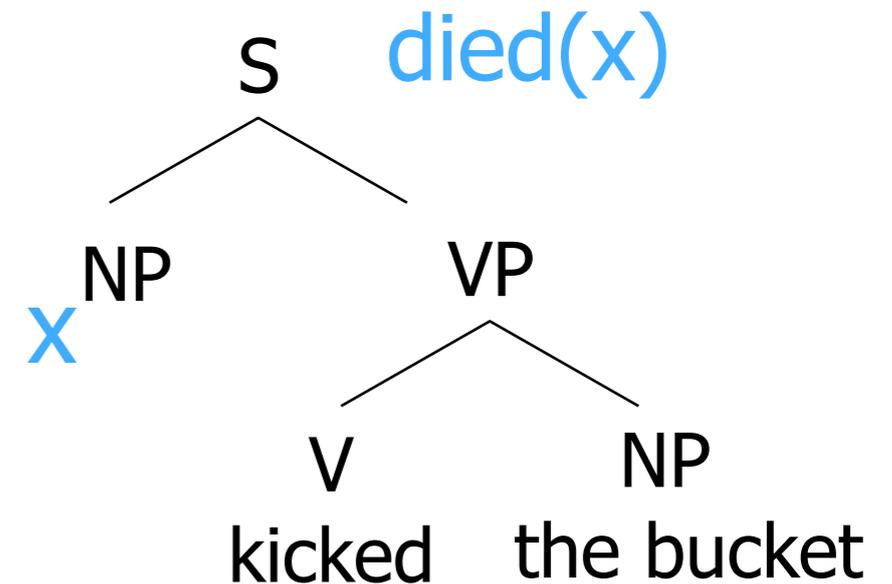
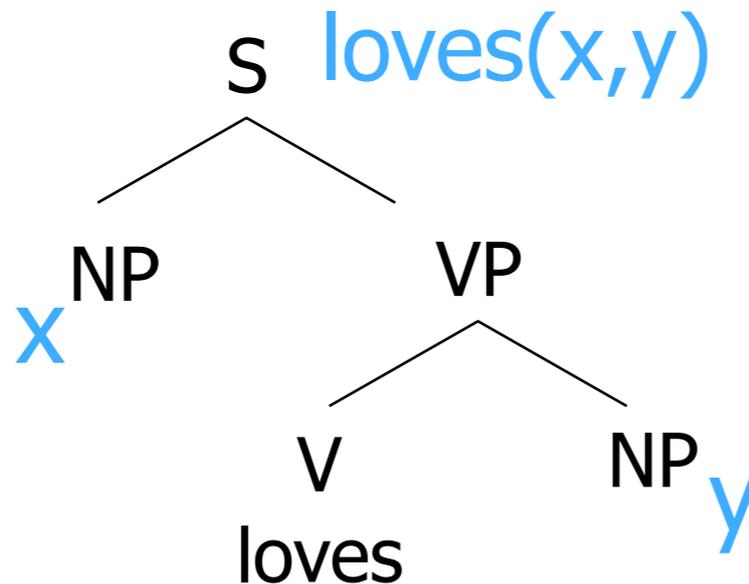
- Add a “sem” feature to each context-free rule
  - $S \rightarrow NP \text{ loves } NP$
  - $S[\text{sem}=\text{loves}(x,y)] \rightarrow NP[\text{sem}=x] \text{ loves } NP[\text{sem}=y]$
  - Meaning of S depends on meaning of NPs
- TAG version:



# Compositional Semantics

- Add a “sem” feature to each context-free rule
  - $S \rightarrow NP \text{ loves } NP$
  - $S[\text{sem}=\text{loves}(x,y)] \rightarrow NP[\text{sem}=x] \text{ loves } NP[\text{sem}=y]$
  - Meaning of S depends on meaning of NPs

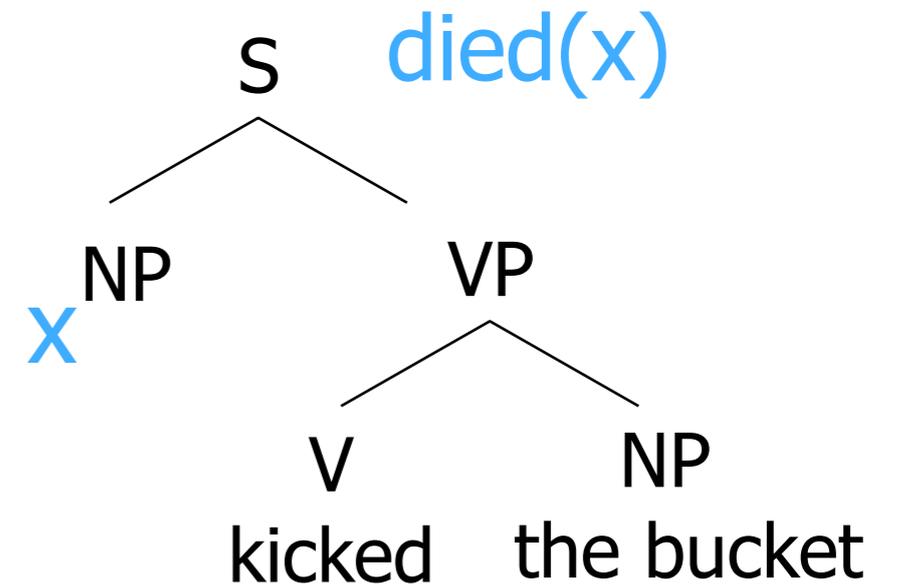
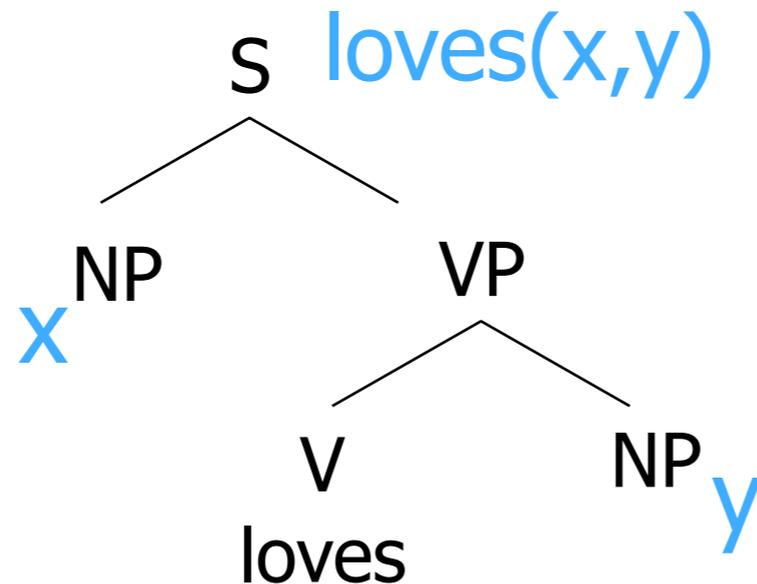
- TAG version:



# Compositional Semantics

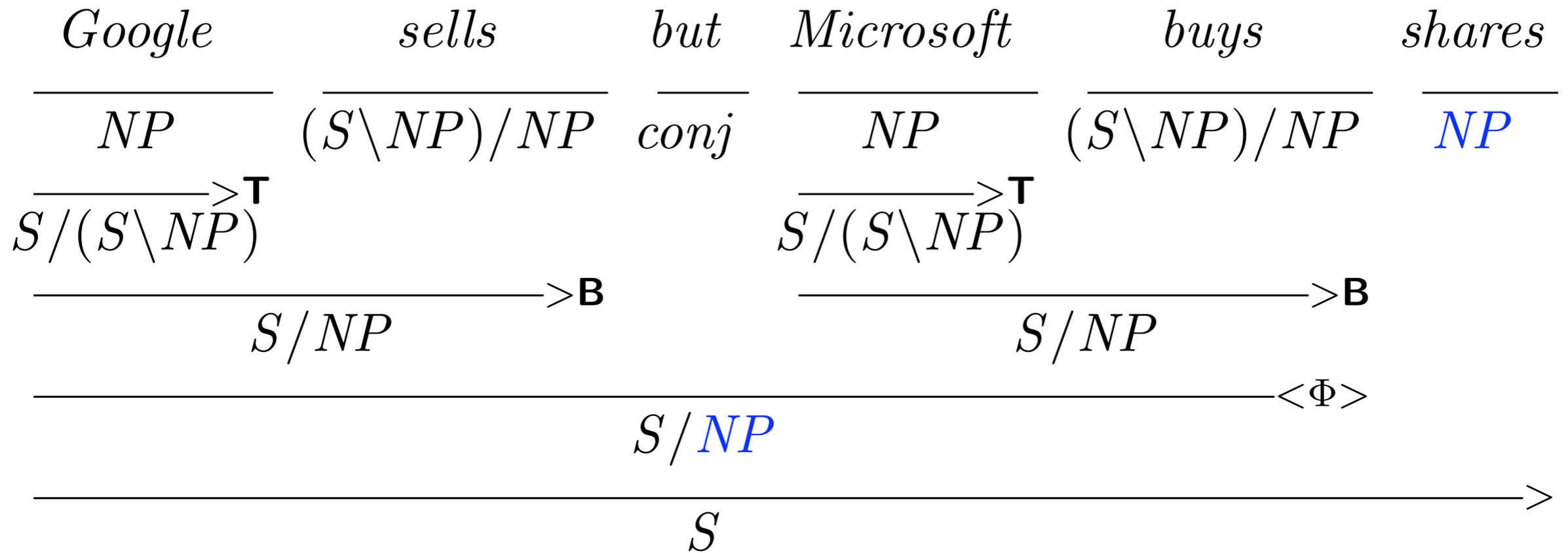
- Add a “sem” feature to each context-free rule
  - $S \rightarrow NP \text{ loves } NP$
  - $S[\text{sem}=\text{loves}(x,y)] \rightarrow NP[\text{sem}=x] \text{ loves } NP[\text{sem}=y]$
  - Meaning of S depends on meaning of NPs

- TAG version:



- Template filling:  $S[\text{sem}=\text{showflights}(x,y)] \rightarrow$   
I want a flight from  $NP[\text{sem}=x]$  to  $NP[\text{sem}=y]$

# “Non-constituents” in CCG – Right Node Raising



# CCG Semantics

- Categories encode argument sequences
- Parallel syntactic combinator operations and lambda calculus semantic operations

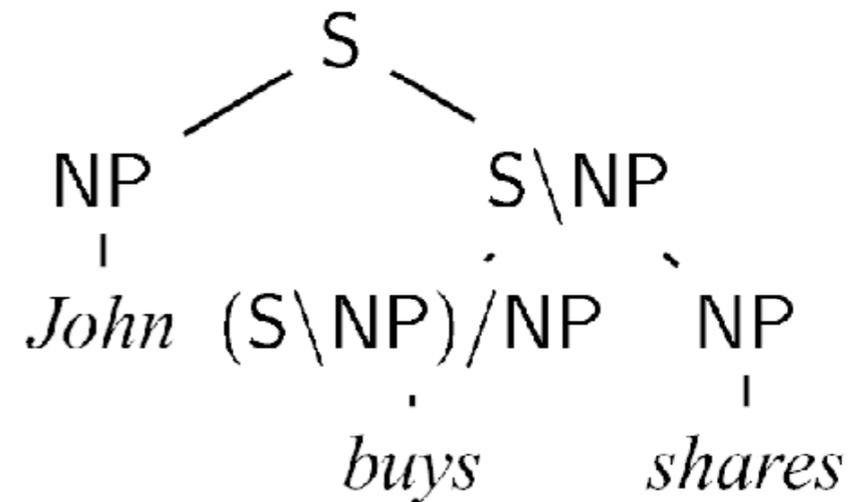
*John*  $\vdash$  NP : *john'*

*shares*  $\vdash$  NP : *shares'*

*buys*  $\vdash$  (S\NP)/NP :  $\lambda x.\lambda y.buys'xy$

*sleeps*  $\vdash$  S\NP :  $\lambda x.sleeps'x$

*well*  $\vdash$  (S\NP)\(S\NP) :  $\lambda f.\lambda x.well'(fx)$



# Words as Vectors

- Represent each word **type**  $w$  by a point in  $k$ -dimensional space
  - e.g.,  $k$  is size of vocabulary
  - the 17<sup>th</sup> coordinate of  $w$  represents **strength** of  $w$ 's association with vocabulary word 17

# Words as Vectors

- Represent each word **type**  $w$  by a point in  $k$ -dimensional space
  - e.g.,  $k$  is size of vocabulary
  - the 17<sup>th</sup> coordinate of  $w$  represents **strength** of  $w$ 's association with vocabulary word 17

= party

# Words as Vectors

- Represent each word **type**  $w$  by a point in  $k$ -dimensional space
  - e.g.,  $k$  is size of vocabulary
  - the 17<sup>th</sup> coordinate of  $w$  represents **strength** of  $w$ 's association with vocabulary word 17

(0, 0, 3, 1, 0, 7, . . . , 1, 0)

# Words as Vectors

- Represent each word **type**  $w$  by a point in  $k$ -dimensional space
  - e.g.,  $k$  is size of vocabulary
  - the 17<sup>th</sup> coordinate of  $w$  represents **strength** of  $w$ 's association with vocabulary word 17

*= party*

aardvark  
abacus  
abandoned  
abbot  
abduct  
above

(0, 0, 3, 1, 0, 7, ...

zygote  
zymurgy

1, 0)

# Words as Vectors

- Represent each word **type**  $w$  by a point in  $k$ -dimensional space
  - e.g.,  $k$  is size of vocabulary
  - the 17<sup>th</sup> coordinate of  $w$  represents **strength** of  $w$ 's association with vocabulary word 17

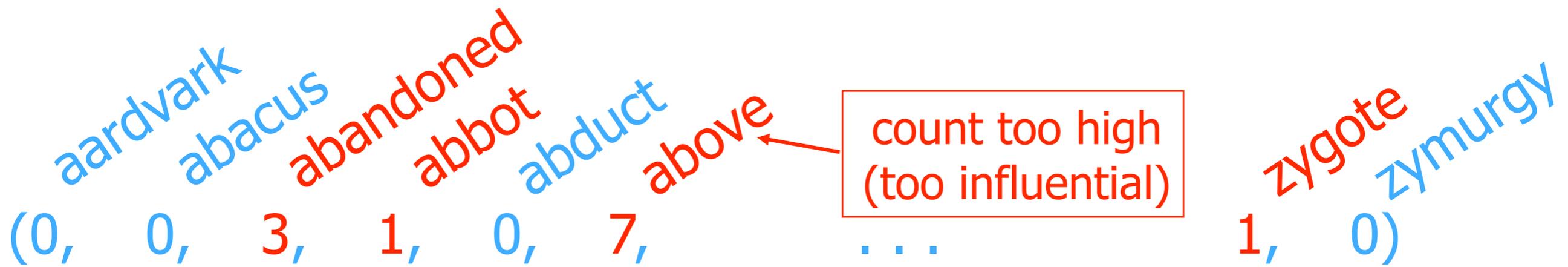
aardvark (0, 0, 3, 1, 0, 7, ... )  
abacus  
abandoned  
abbot  
abduct  
above  
zygote  
zymurgy

From  
corpus:

Arlen Specter **abandoned** the Republican party.  
There were lots of **abbots** and nuns dancing at that party.  
The party **above** the art gallery was, **above** all, a laboratory  
for synthesizing **zygotes** and beer.

# Words as Vectors

- Represent each word **type**  $w$  by a point in  $k$ -dimensional space
  - e.g.,  $k$  is size of vocabulary
  - the 17<sup>th</sup> coordinate of  $w$  represents **strength** of  $w$ 's association with vocabulary word 17

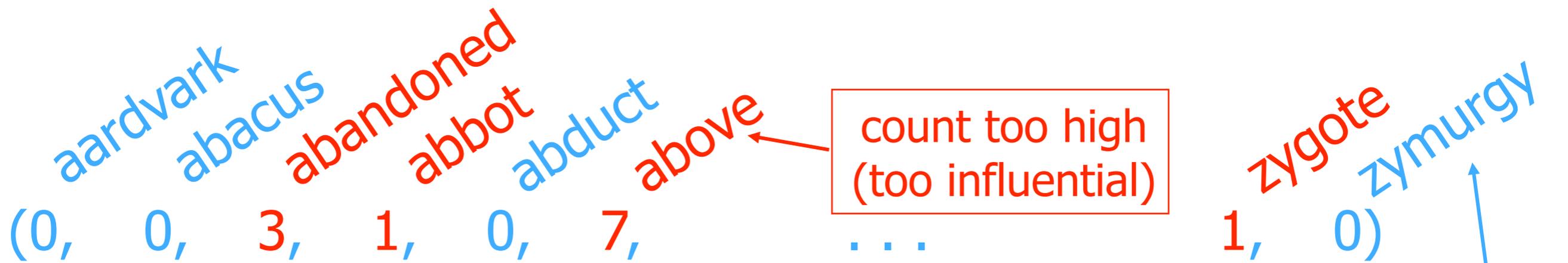


From  
corpus:

Arlen Specter **abandoned** the Republican **party**.  
There were lots of **abbots** and nuns dancing at that **party**.  
The **party** **above** the art gallery was, **above** all, a laboratory  
for synthesizing **zygotes** and beer.

# Words as Vectors

- Represent each word **type**  $w$  by a point in  $k$ -dimensional space
  - e.g.,  $k$  is size of vocabulary
  - the 17<sup>th</sup> coordinate of  $w$  represents **strength** of  $w$ 's association with vocabulary word 17



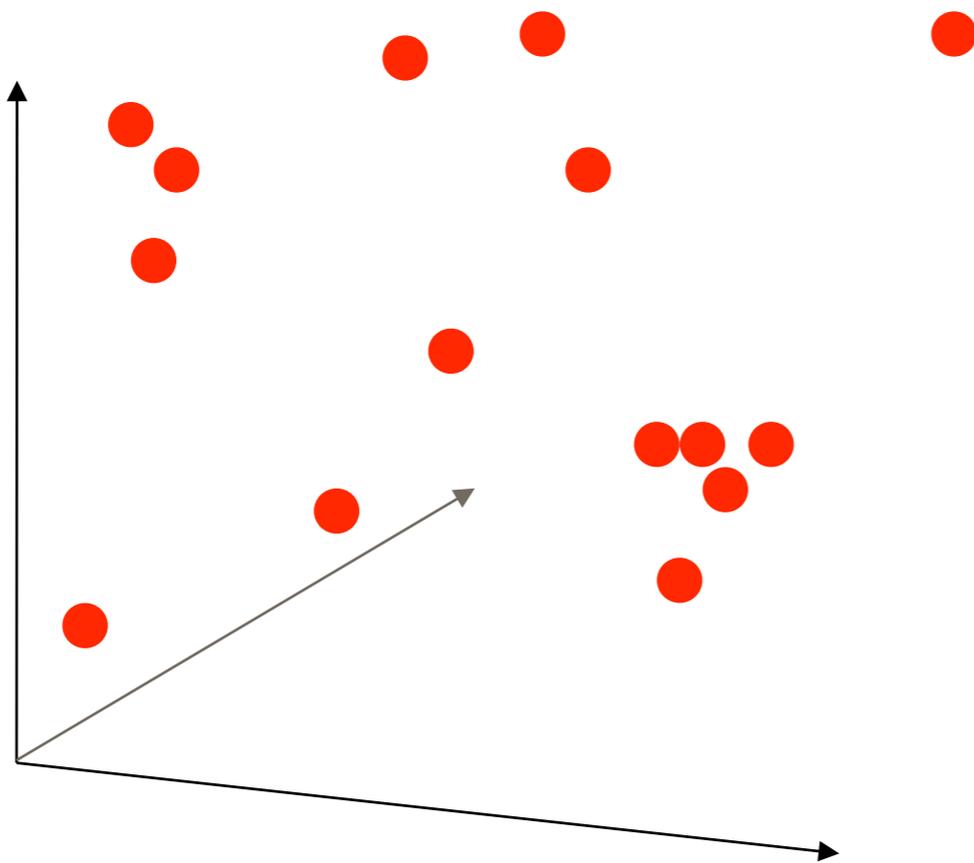
From  
corpus:

Arlen Specter **abandoned** the Republican **party**.  
There were lots of **abbots** and nuns dancing at that **party**.  
The **party** **above** the art gallery was, **above** all, a laboratory  
for synthesizing **zygotes** and beer.

# Learning Classes by Clustering

- Plot all word types in k-dimensional space
- Look for **clusters** of close-together types

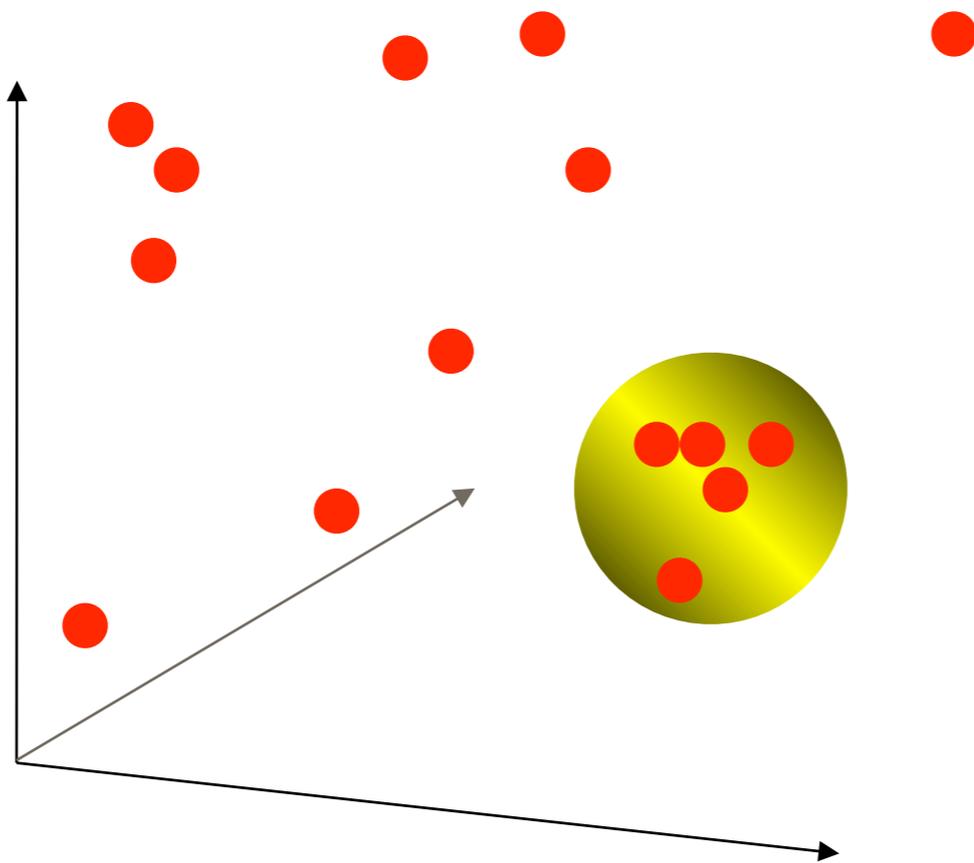
Plot in k dimensions (here  $k=3$ )



# Learning Classes by Clustering

- Plot all word types in k-dimensional space
- Look for **clusters** of close-together types

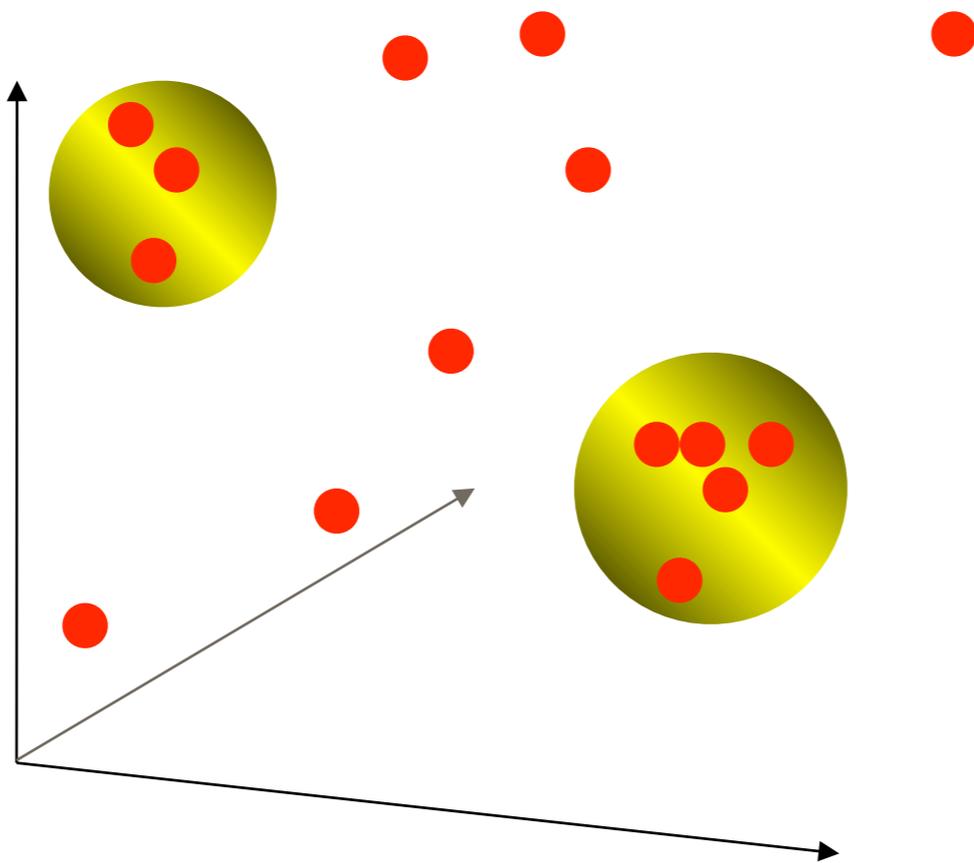
Plot in k dimensions (here k=3)



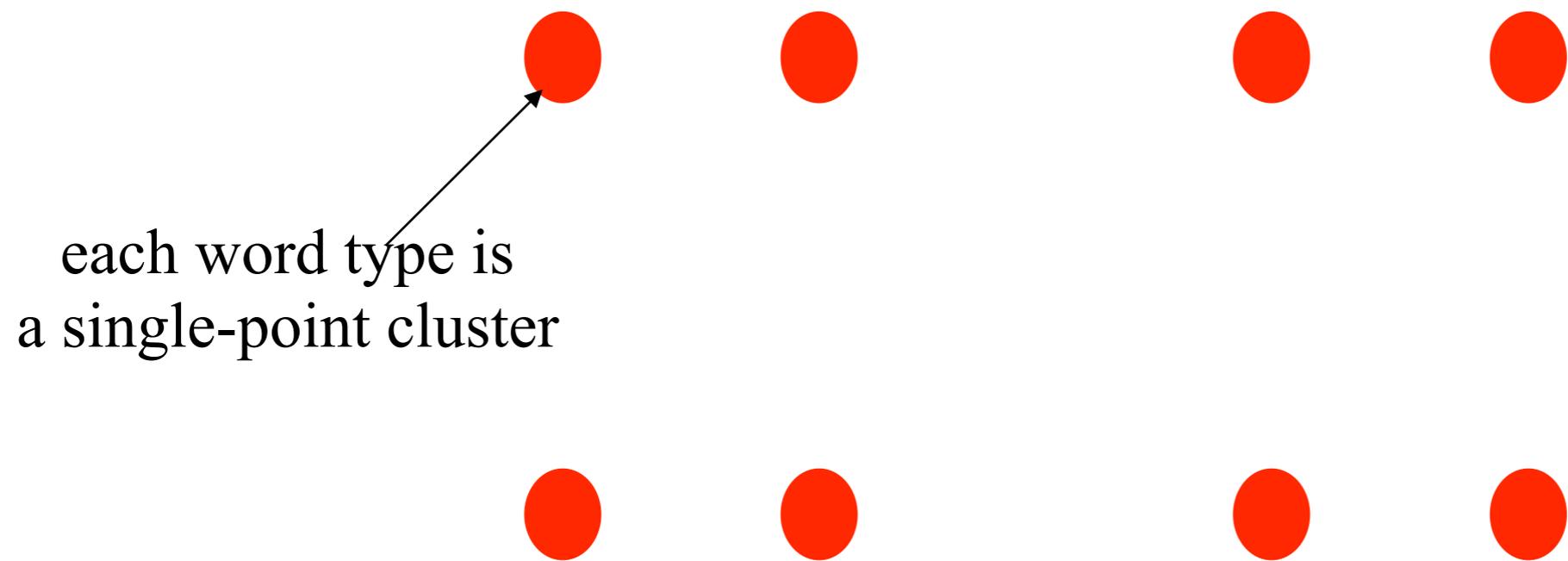
# Learning Classes by Clustering

- Plot all word types in k-dimensional space
- Look for **clusters** of close-together types

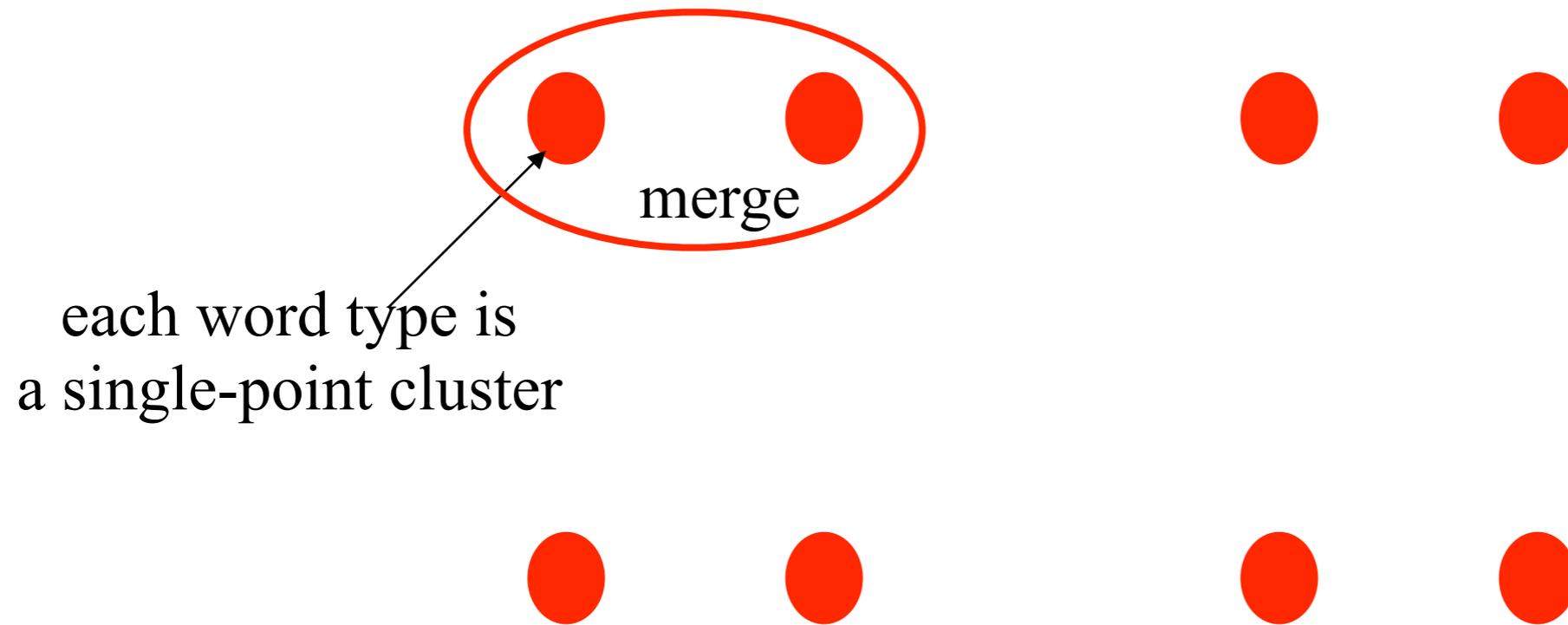
Plot in k dimensions (here  $k=3$ )



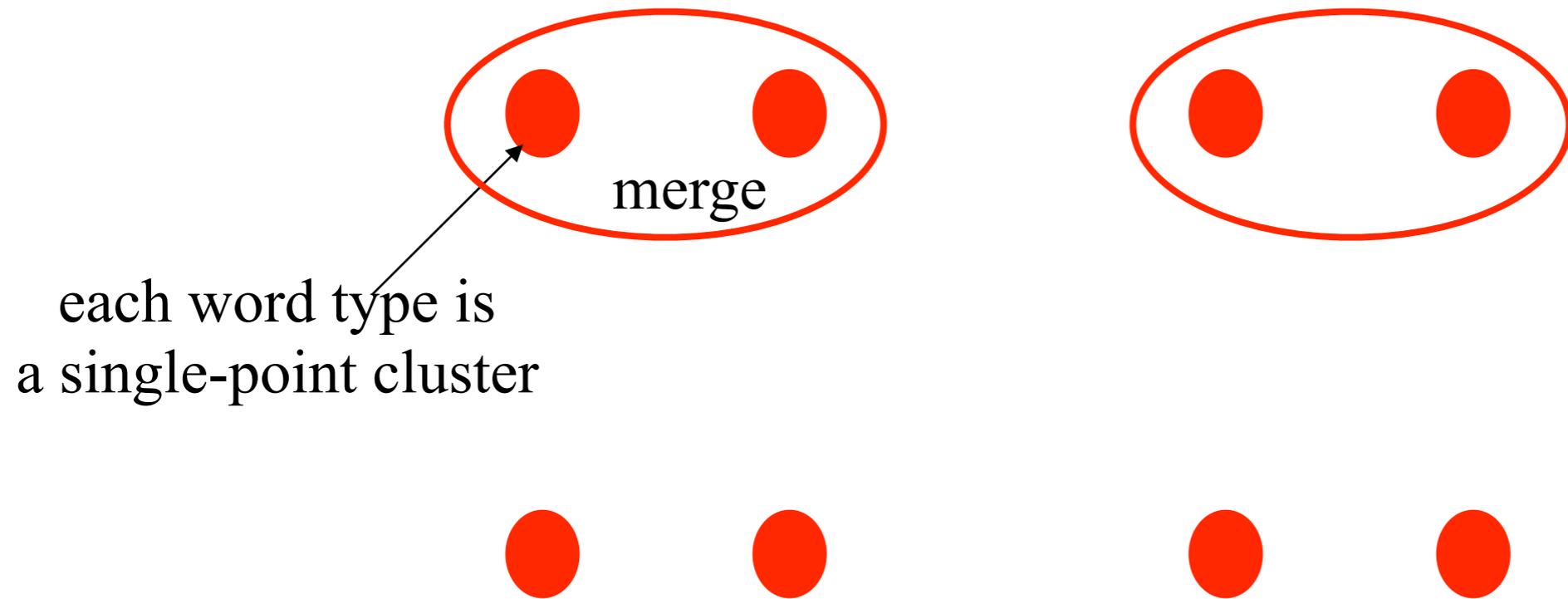
# Bottom-Up Clustering – Single-Link



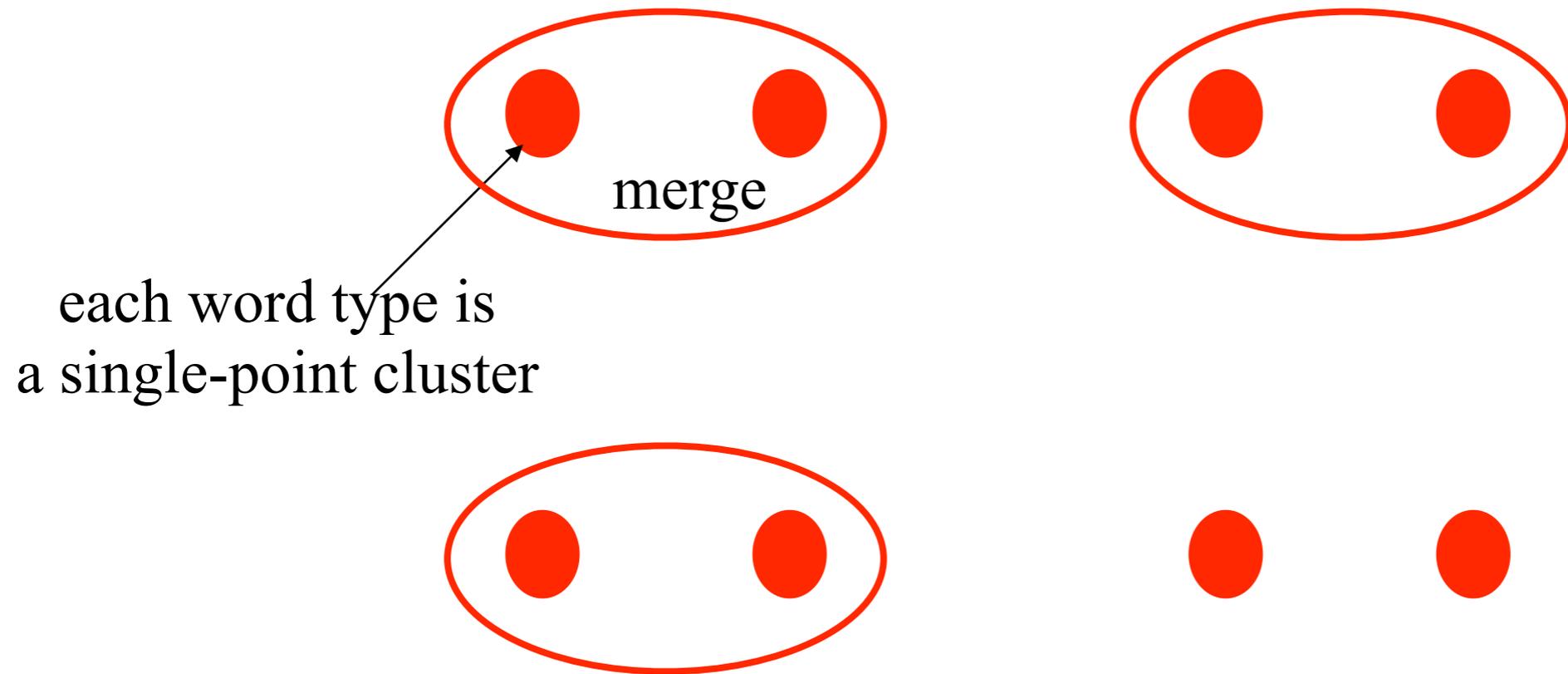
# Bottom-Up Clustering – Single-Link



# Bottom-Up Clustering – Single-Link

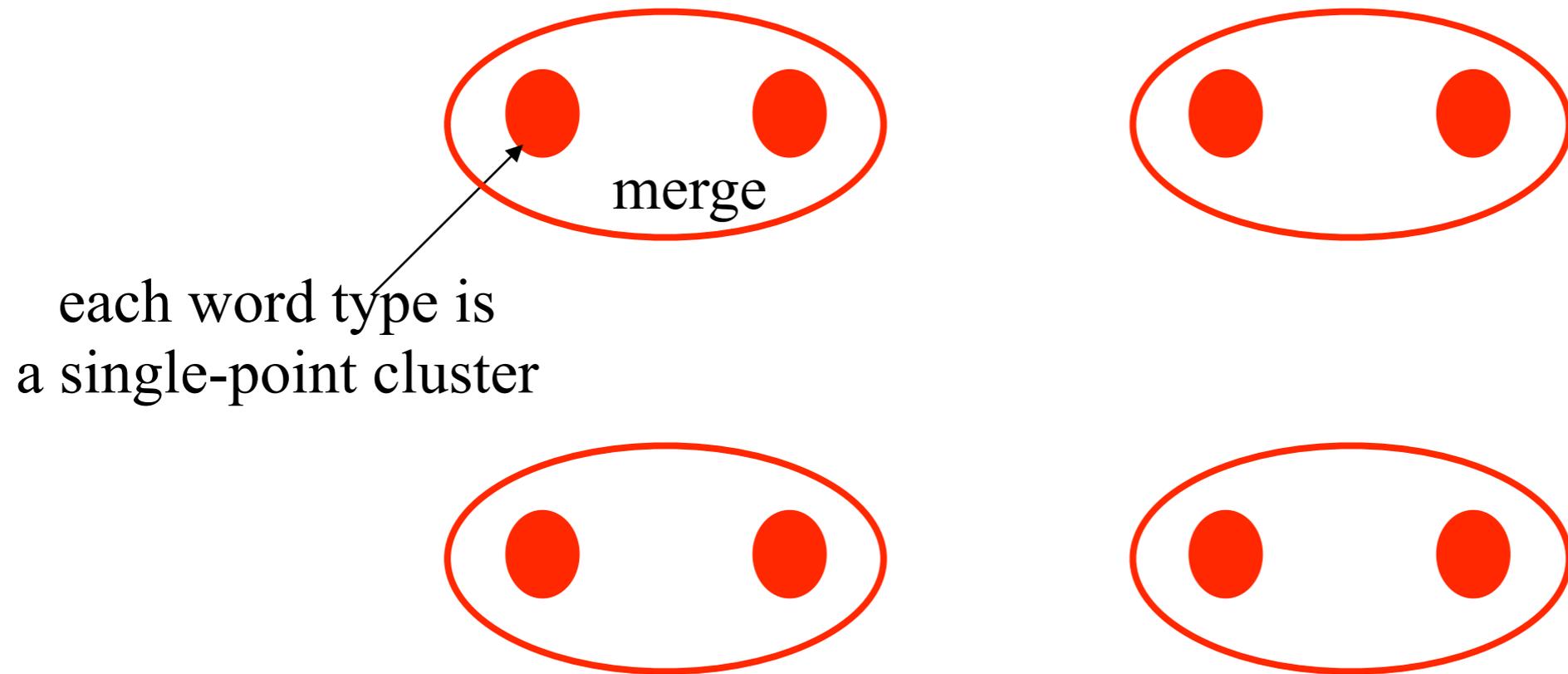


# Bottom-Up Clustering – Single-Link

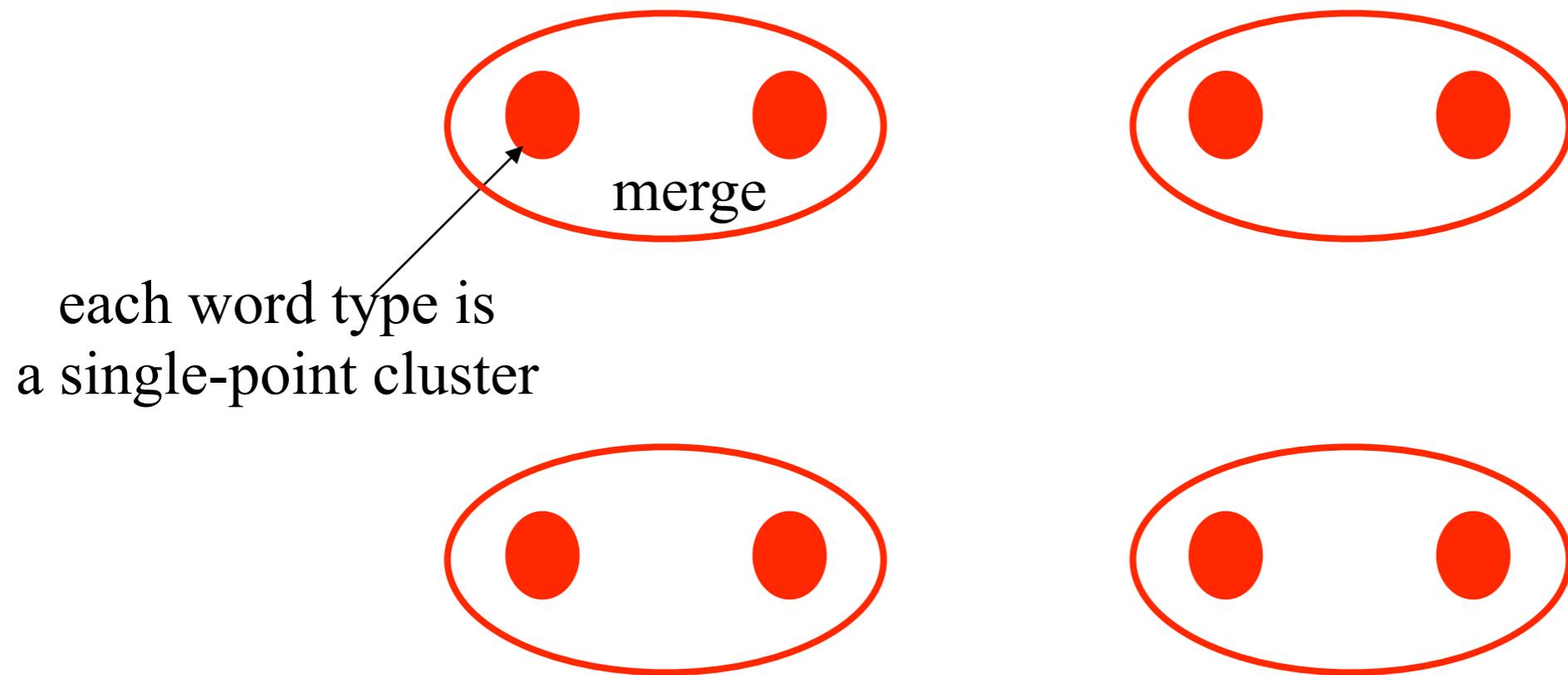


each word type is  
a single-point cluster

# Bottom-Up Clustering – Single-Link



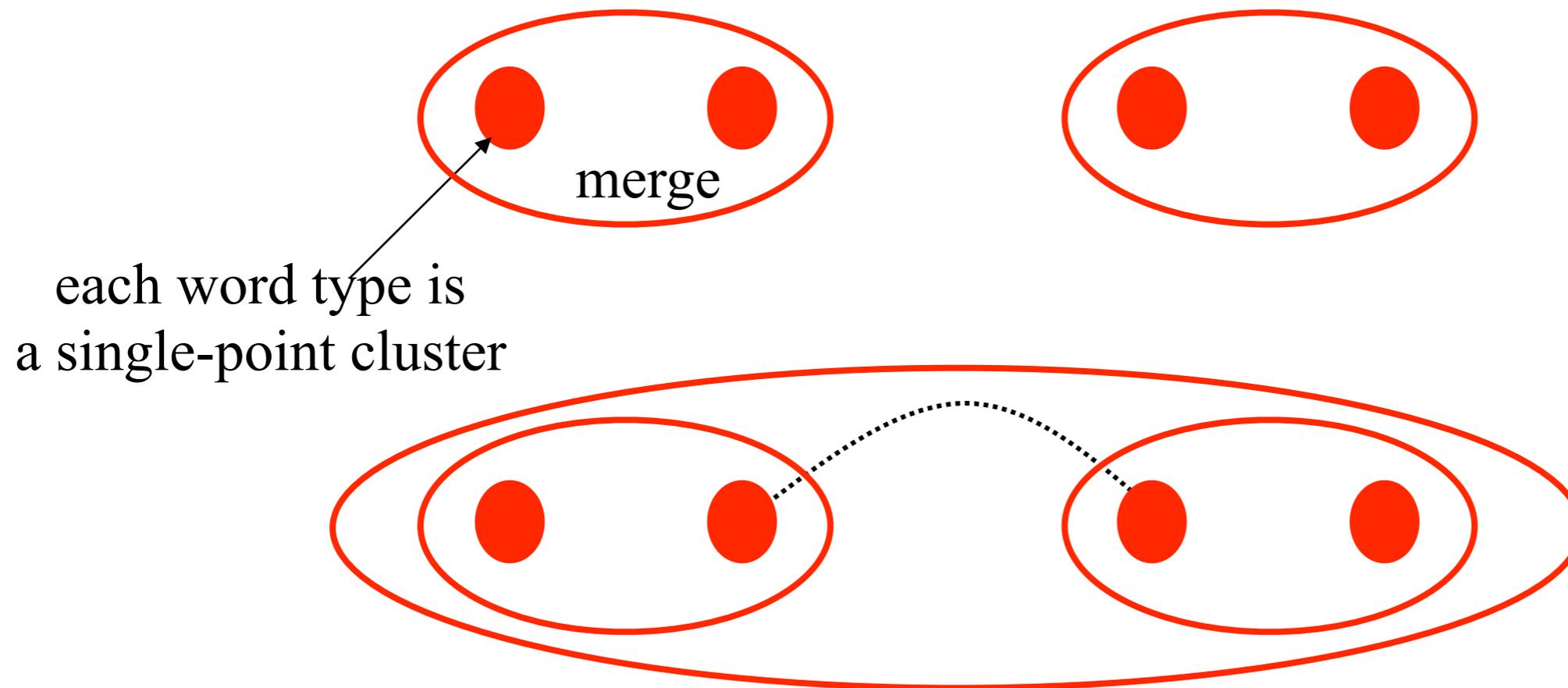
# Bottom-Up Clustering – Single-Link



Again, merge closest pair of clusters:

**Single-link:** clusters are close if **any** of their points are  
 $\text{dist}(A,B) = \min \text{dist}(a,b)$  for  $a \in A, b \in B$

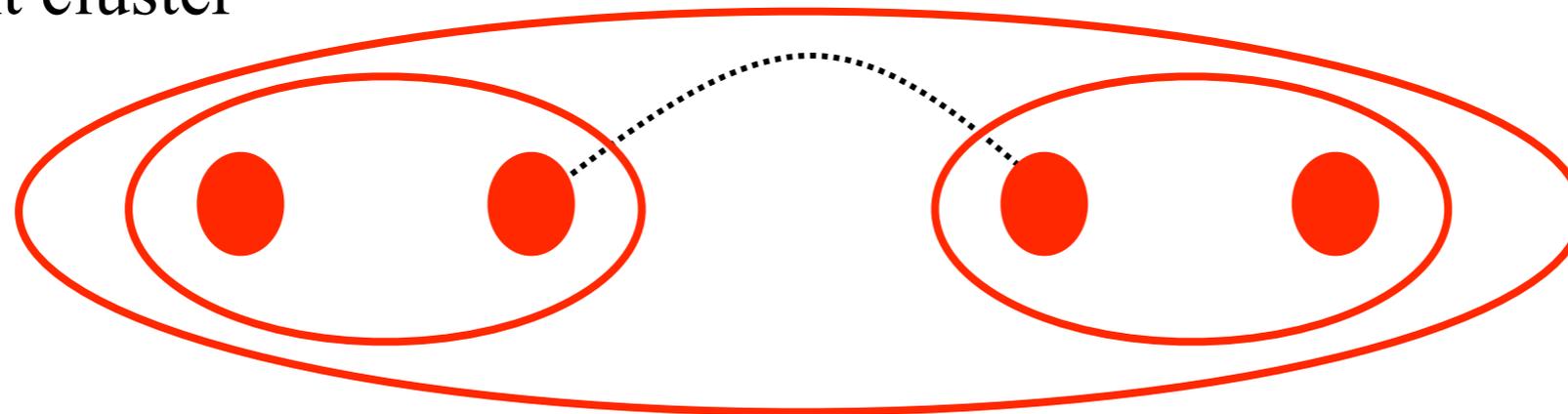
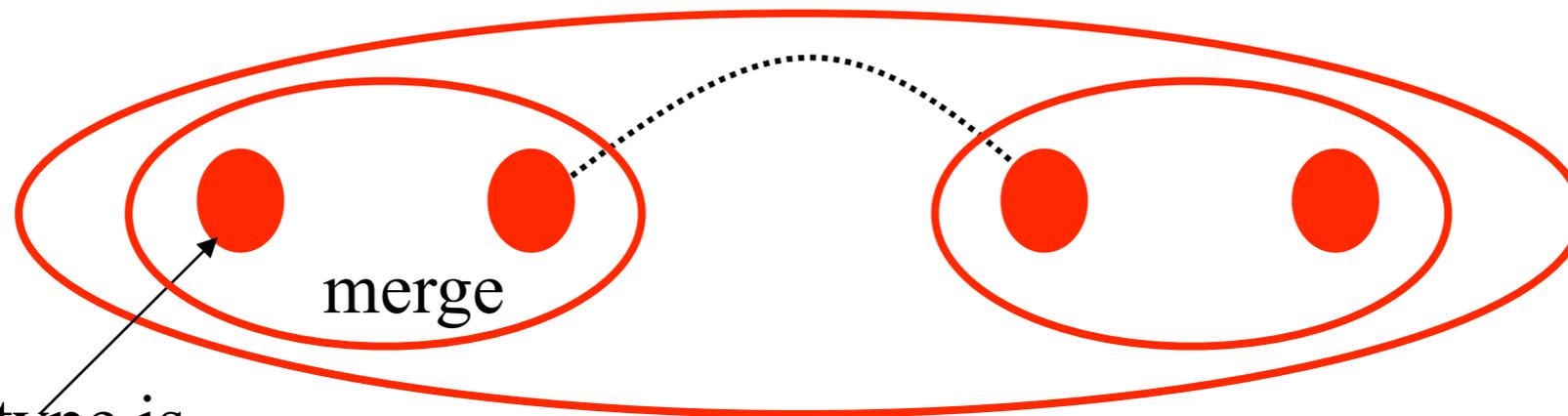
# Bottom-Up Clustering – Single-Link



Again, merge closest pair of clusters:

**Single-link:** clusters are close if **any** of their points are  
 $\text{dist}(A,B) = \min \text{dist}(a,b)$  for  $a \in A, b \in B$

# Bottom-Up Clustering – Single-Link



Again, merge closest pair of clusters:

**Single-link:** clusters are close if **any** of their points are  
 $\text{dist}(A,B) = \min \text{dist}(a,b)$  for  $a \in A, b \in B$

# Bottom-Up Clustering

- Start with one cluster per point
- Repeatedly merge 2 closest clusters
  - **Single-link:**  $\text{dist}(A,B) = \min \text{dist}(a,b)$  for  $a \in A, b \in B$
  - **Complete-link:**  $\text{dist}(A,B) = \max \text{dist}(a,b)$  for  $a \in A, b \in B$ 
    - too slow to update cluster distances after each merge; but  $\exists$  alternatives!

# Bottom-Up Clustering

- Start with one cluster per point
- Repeatedly merge 2 closest clusters
  - **Single-link:**  $\text{dist}(A,B) = \min \text{dist}(a,b)$  for  $a \in A, b \in B$
  - **Complete-link:**  $\text{dist}(A,B) = \max \text{dist}(a,b)$  for  $a \in A, b \in B$ 
    - too slow to update cluster distances after each merge; but  $\exists$  alternatives!
  - **Average-link:**  $\text{dist}(A,B) = \text{mean dist}(a,b)$  for  $a \in A, b \in B$
  - **Centroid-link:**  $\text{dist}(A,B) = \text{dist}(\text{mean}(A), \text{mean}(B))$

# Bottom-Up Clustering

- Start with one cluster per point
- Repeatedly merge 2 closest clusters
  - **Single-link:**  $\text{dist}(A,B) = \min \text{dist}(a,b)$  for  $a \in A, b \in B$
  - **Complete-link:**  $\text{dist}(A,B) = \max \text{dist}(a,b)$  for  $a \in A, b \in B$ 
    - too slow to update cluster distances after each merge; but  $\exists$  alternatives!
  - **Average-link:**  $\text{dist}(A,B) = \text{mean dist}(a,b)$  for  $a \in A, b \in B$
  - **Centroid-link:**  $\text{dist}(A,B) = \text{dist}(\text{mean}(A), \text{mean}(B))$
- Stop when clusters are “big enough”
  - e.g., provide adequate support for backoff (on a development corpus)

# Bottom-Up Clustering

- Start with one cluster per point
- Repeatedly merge 2 closest clusters
  - **Single-link:**  $\text{dist}(A,B) = \min \text{dist}(a,b)$  for  $a \in A, b \in B$
  - **Complete-link:**  $\text{dist}(A,B) = \max \text{dist}(a,b)$  for  $a \in A, b \in B$ 
    - too slow to update cluster distances after each merge; but  $\exists$  alternatives!
  - **Average-link:**  $\text{dist}(A,B) = \text{mean dist}(a,b)$  for  $a \in A, b \in B$
  - **Centroid-link:**  $\text{dist}(A,B) = \text{dist}(\text{mean}(A), \text{mean}(B))$
- Stop when clusters are “big enough”
  - e.g., provide adequate support for backoff (on a development corpus)
- Some flexibility in defining  $\text{dist}(a,b)$ 
  - Might not be Euclidean distance; e.g., use vector angle

# EM Clustering (for $k$ clusters)

# EM Clustering (for k clusters)

- EM algorithm
  - Viterbi version – called “k-means clustering”
  - Full EM version – called “Gaussian mixtures”

# EM Clustering (for k clusters)

- EM algorithm
  - Viterbi version – called “k-means clustering”
  - Full EM version – called “Gaussian mixtures”

# EM Clustering (for k clusters)

- EM algorithm
  - Viterbi version – called “k-means clustering”
  - Full EM version – called “Gaussian mixtures”
- **Expectation step:** Use current parameters (and observations) to reconstruct hidden structure

# EM Clustering (for k clusters)

- EM algorithm
  - Viterbi version – called “k-means clustering”
  - Full EM version – called “Gaussian mixtures”
- **Expectation step:** Use current parameters (and observations) to reconstruct hidden structure
- **Maximization step:** Use that hidden structure (and observations) to reestimate parameters

# EM Clustering (for k clusters)

- EM algorithm
  - Viterbi version – called “k-means clustering”
  - Full EM version – called “Gaussian mixtures”
- **Expectation step:** Use current parameters (and observations) to reconstruct hidden structure
- **Maximization step:** Use that hidden structure (and observations) to reestimate parameters

# EM Clustering (for k clusters)

- EM algorithm
  - Viterbi version – called “k-means clustering”
  - Full EM version – called “Gaussian mixtures”
- **Expectation step:** Use current parameters (and observations) to reconstruct hidden structure
- **Maximization step:** Use that hidden structure (and observations) to reestimate parameters
- **Parameters:** k points representing cluster centers

# EM Clustering (for k clusters)

- EM algorithm
  - Viterbi version – called “k-means clustering”
  - Full EM version – called “Gaussian mixtures”
- **Expectation step:** Use current parameters (and observations) to reconstruct hidden structure
- **Maximization step:** Use that hidden structure (and observations) to reestimate parameters
- **Parameters:** k points representing cluster centers
- **Hidden structure:** for each data point (word type), which center generated it?

# Lexical translation

- How to translate a word → look up in dictionary

**Haus** — *house, building, home, household, shell.*

- *Multiple translations*

- some more frequent than others
- for instance: *house*, and *building* most common
- special cases: *Haus* of a *snail* is its *shell*

- Note: During all the lectures, we will translate from a foreign language into English

## Collect statistics

- Look at a *parallel corpus* (German text along with English translation)

Translation of <i>Haus</i>	Count
<i>house</i>	8,000
<i>building</i>	1,600
<i>home</i>	200
<i>household</i>	150
<i>shell</i>	50

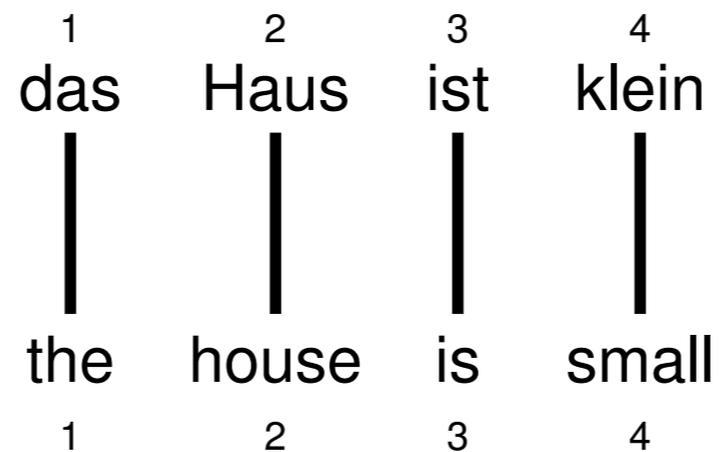
## Estimate translation probabilities

- *Maximum likelihood estimation*

$$p_f(e) = \begin{cases} 0.8 & \text{if } e = \textit{house}, \\ 0.16 & \text{if } e = \textit{building}, \\ 0.02 & \text{if } e = \textit{home}, \\ 0.015 & \text{if } e = \textit{household}, \\ 0.005 & \text{if } e = \textit{shell}. \end{cases}$$

## Alignment

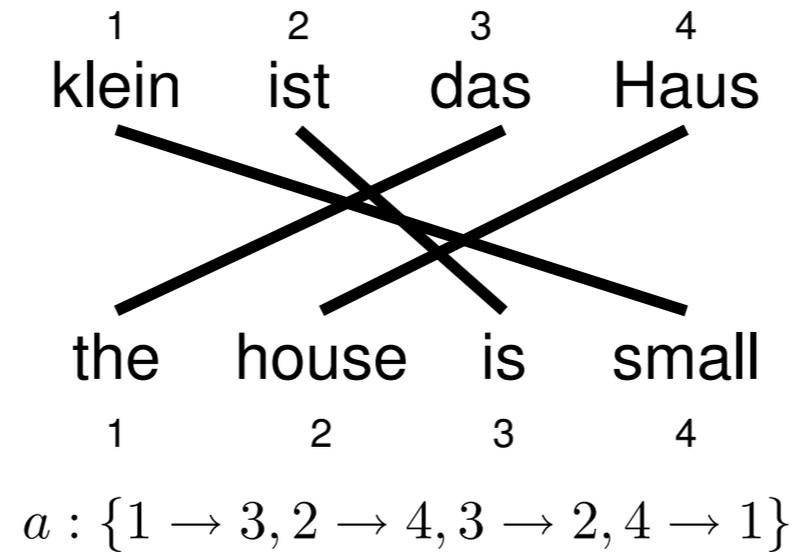
- In a parallel text (or when we translate), we **align** words in one language with the words in the other



- Word *positions* are numbered 1–4

## Reordering

- Words may be **reordered** during translation



## IBM Model 1

- *Generative model*: break up translation process into smaller steps
  - **IBM Model 1** only uses *lexical translation*
- Translation probability
  - for a foreign sentence  $\mathbf{f} = (f_1, \dots, f_{l_f})$  of length  $l_f$
  - to an English sentence  $\mathbf{e} = (e_1, \dots, e_{l_e})$  of length  $l_e$
  - with an alignment of each English word  $e_j$  to a foreign word  $f_i$  according to the alignment function  $a : j \rightarrow i$

$$p(\mathbf{e}, a | \mathbf{f}) = \frac{\epsilon}{(l_f + 1)^{l_e}} \prod_{j=1}^{l_e} t(e_j | f_{a(j)})$$

- parameter  $\epsilon$  is a *normalization constant*

## Example

<i>das</i>		<i>Haus</i>		<i>ist</i>		<i>klein</i>	
<i>e</i>	$t(e f)$	<i>e</i>	$t(e f)$	<i>e</i>	$t(e f)$	<i>e</i>	$t(e f)$
<i>the</i>	0.7	<i>house</i>	0.8	<i>is</i>	0.8	<i>small</i>	0.4
<i>that</i>	0.15	<i>building</i>	0.16	<i>'s</i>	0.16	<i>little</i>	0.4
<i>which</i>	0.075	<i>home</i>	0.02	<i>exists</i>	0.02	<i>short</i>	0.1
<i>who</i>	0.05	<i>household</i>	0.015	<i>has</i>	0.015	<i>minor</i>	0.06
<i>this</i>	0.025	<i>shell</i>	0.005	<i>are</i>	0.005	<i>petty</i>	0.04

$$\begin{aligned} p(e, a|f) &= \frac{\epsilon}{4^3} \times t(\text{the}|\text{das}) \times t(\text{house}|\text{Haus}) \times t(\text{is}|\text{ist}) \times t(\text{small}|\text{klein}) \\ &= \frac{\epsilon}{4^3} \times 0.7 \times 0.8 \times 0.8 \times 0.4 \\ &= 0.0028\epsilon \end{aligned}$$

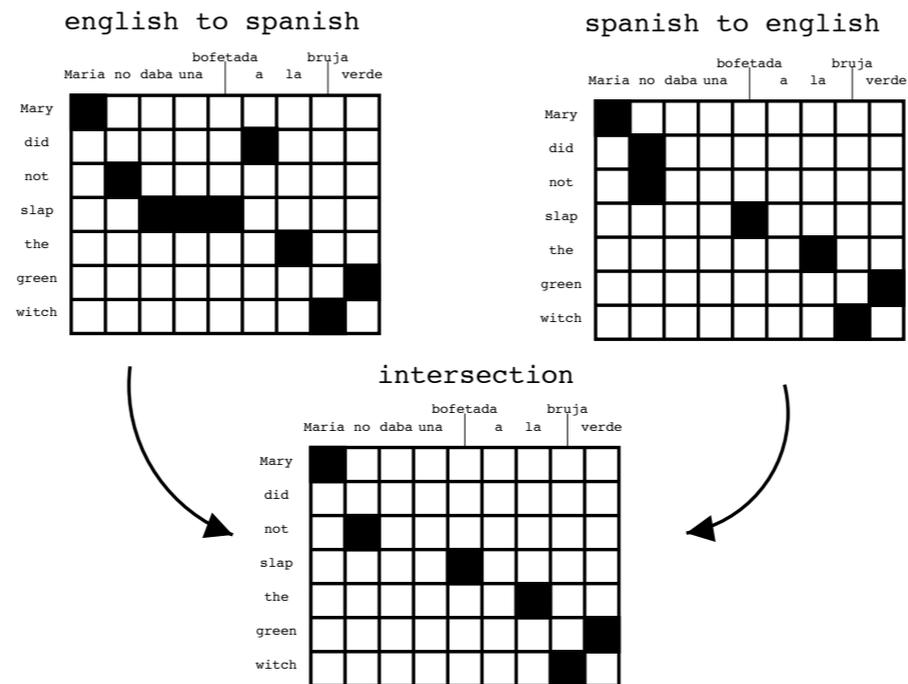
## Learning lexical translation models

- We would like to *estimate* the lexical translation probabilities  $t(e|f)$  from a parallel corpus
- ... but we do not have the alignments
- **Chicken and egg problem**
  - if we had the *alignments*,
    - we could estimate the *parameters* of our generative model
  - if we had the *parameters*,
    - we could estimate the *alignments*

## EM algorithm

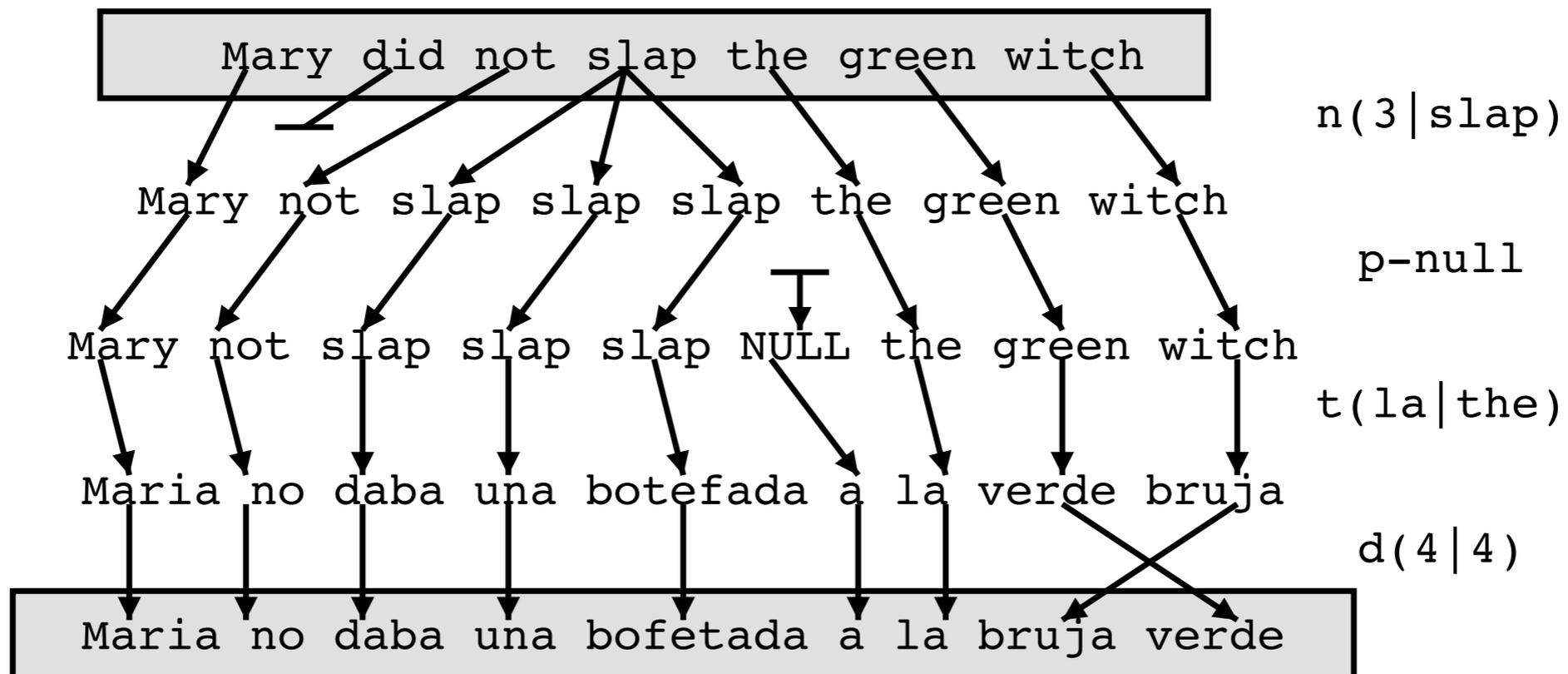
- **Incomplete data**
  - if we had *complete data*, would could estimate *model*
  - if we had *model*, we could fill in the *gaps in the data*
- **Expectation Maximization (EM)** in a nutshell
  - initialize model parameters (e.g. uniform)
  - assign probabilities to the missing data
  - estimate model parameters from completed data
  - iterate

# Symmetrizing word alignments



- *Intersection* of GIZA++ bidirectional alignments

# IBM Model 4



# Phrase Models

I									
did									
not									
unfortunately									
receive									
an									
answer									
to									
this									
question									
	Auf	diese	Frage	habe	ich	leider	keine	Antwort	bekom men

Some good phrase pairs.

# Phrase Models

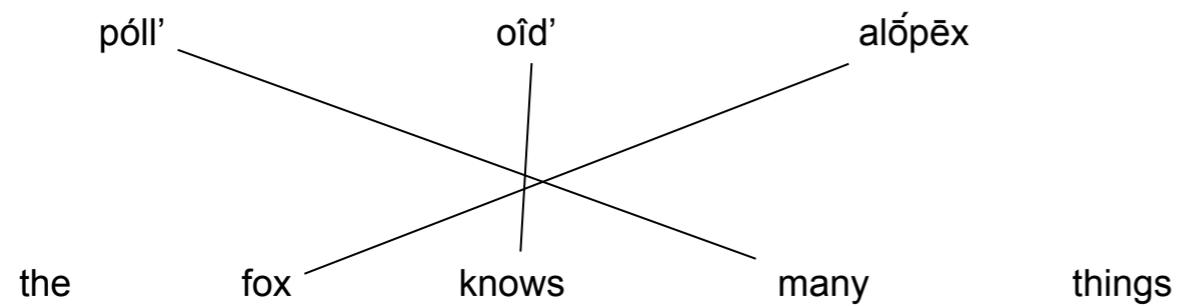
I								
did								
not								
unfortunately								
receive								
an								
answer								
to								
this								
question								
	Auf	diese	Frage	habe	ich	leider	keine	Antwort bekommen

Some bad phrase pairs.

# Synchronous Grammars

- Just like monolingual grammars except...
  - Each rule involves pairs (tuples) of nonterminals
  - Tuples of elementary trees for TAG, etc.
- First proposed for source-source translation in compilers
- Can be constituency, dependency, lexicalized, etc.
- Parsing speedups for monolingual grammar don't necessarily work
  - E.g., no split-head trick for lexicalized parsing
- Binarization less straightforward

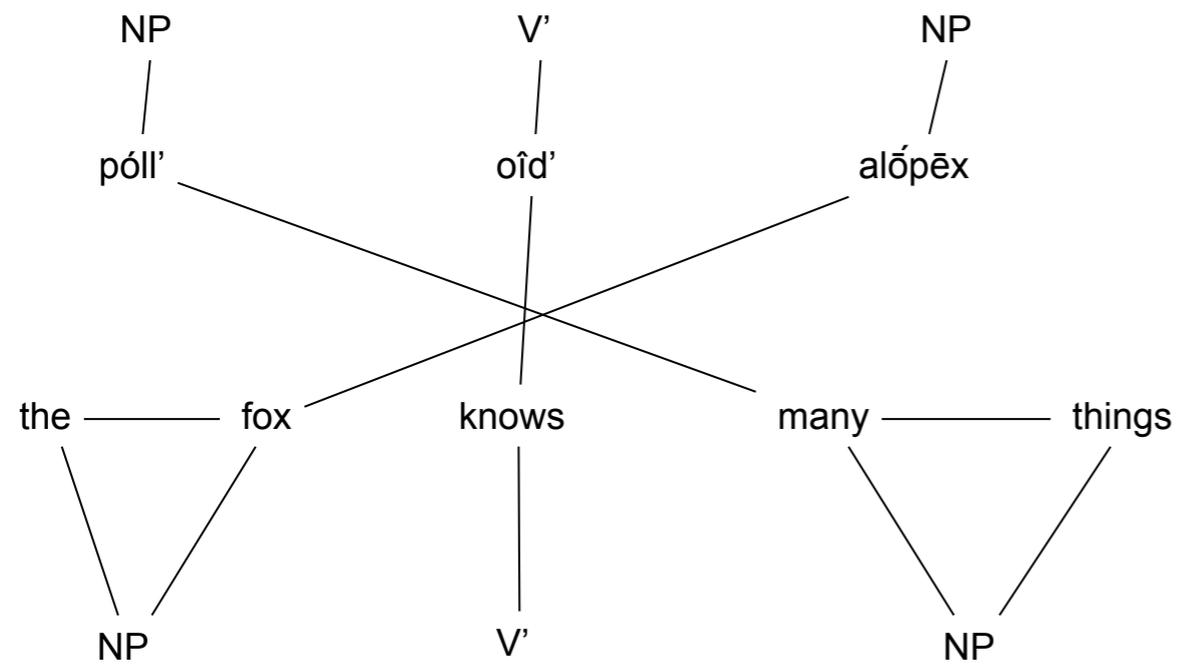
# Bilingual Parsing



	póll'	oîd'	alópēx
the			
fox			NN/NN
knows		VB/VB	
many	JJ/JJ		
things			

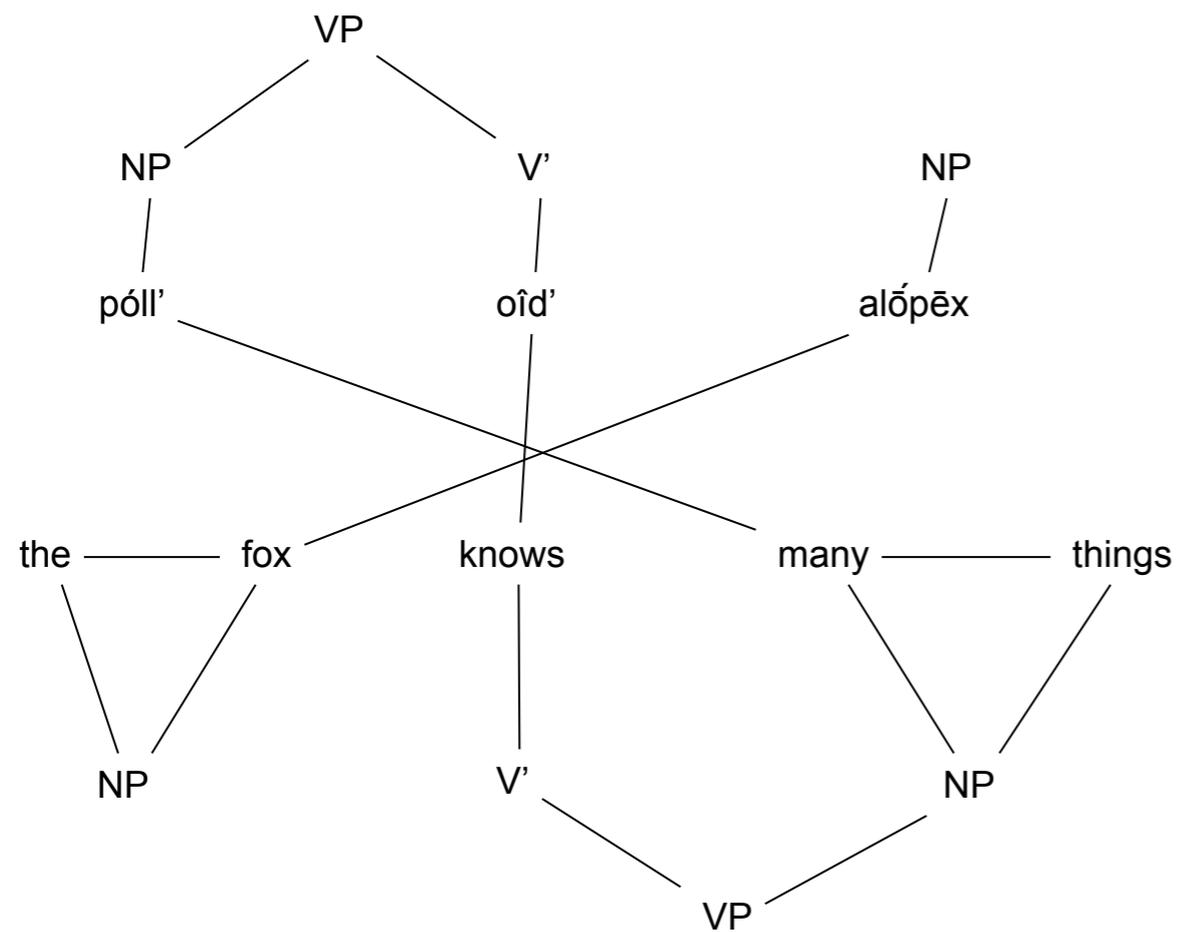
*A variant of CKY chart parsing.*

# Bilingual Parsing



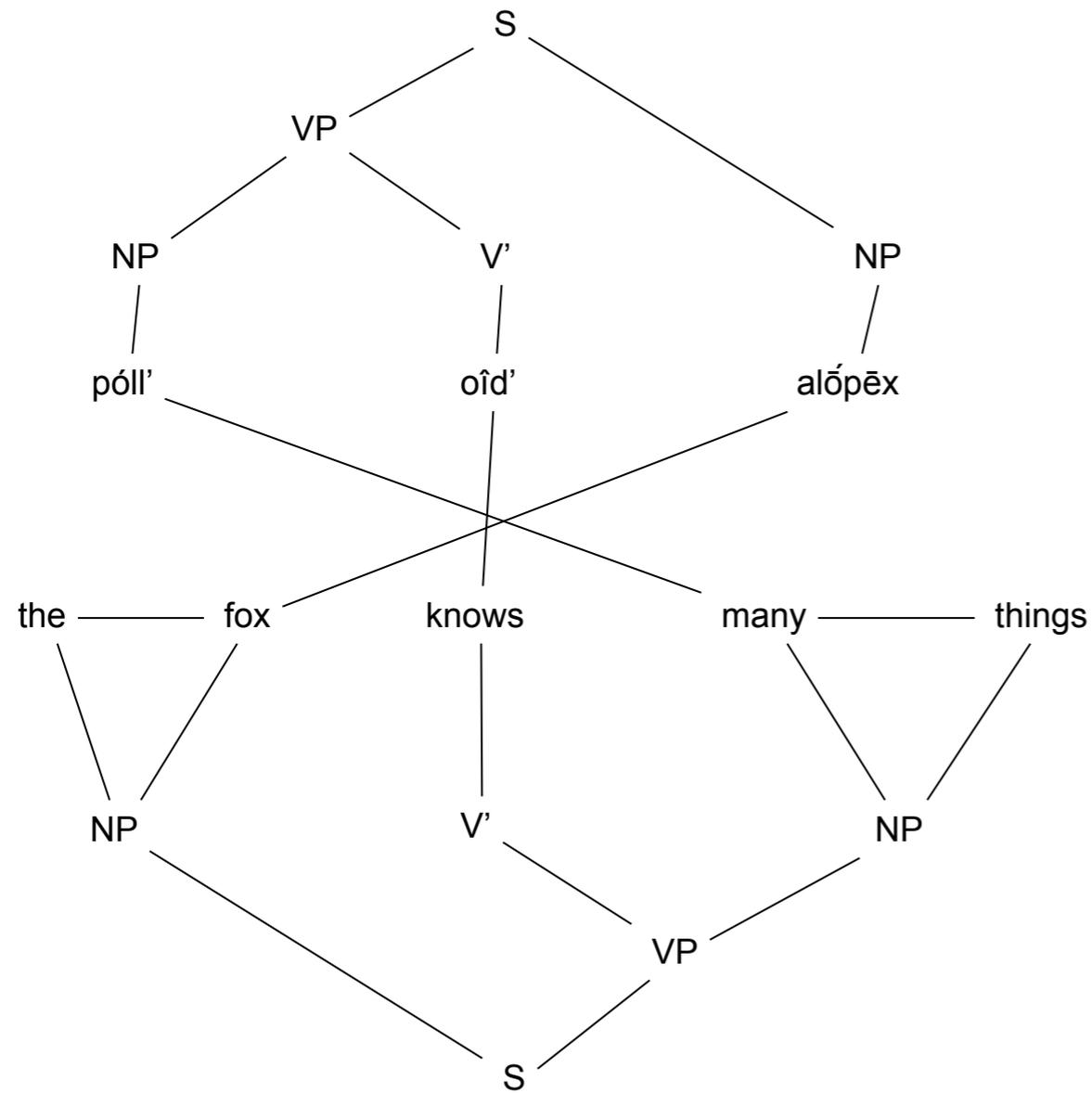
	pól'	oîd'	alópēx
the			NP/NP
fox			
knows		VP/VP	
many	NP/NP		
things			

# Bilingual Parsing



	póll'	oîd'	alópēx
the			NP/NP
fox			
knows	VP/VP		
many			
things			

# Bilingual Parsing



	póll'	oîd'	alópēx
the	S/S		
fox			
knows			
many			
things			