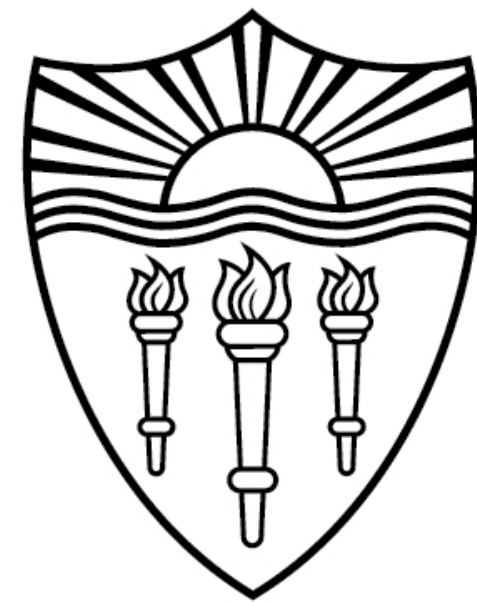CSCI 544: Applied Natural Language Processing

# Sequence Labeling-I

Xuezhe Ma (Max)

Xuezhe Ma (Max)

USC University of Southern California

# Logistical Points

- **For Team Projects**
  - If you have not joined a team, please write you name on the <u>spreadsheet</u> below all teams.
  - We will assign you randomly

# Overview

- **The Sequence Labeling Problem**
  - General Structured Prediction Tasks
  - Part-of-speech Tagging: A case study
  - Generative Models vs. Discriminative Models
  - Maximum Likelihood Estimation (MLE)
- **Hidden Markov Model (HMM)**
  - Basic definitions
  - Parameter estimation
  - The Viterbi algorithm
- **Log-Linear Models**
  - Maximum Entropy Markov Models (MEMMs)
  - Conditional Random Fields (CRFs)

# Overview

- **The Sequence Labeling Problem**
  - General Structured Prediction Tasks
  - Part-of-speech Tagging: A case study
  - Generative Models vs. Discriminative Models
  - Maximum Likelihood Estimation (MLE)
- Hidden Markov Model (HMM)
  - Basic definitions
  - Parameter estimation
  - The Viterbi algorithm
- Log-Linear Models
  - Maximum Entropy Markov Models (MEMMs)
  - Conditional Random Fields (CRFs)

# What is Structured Prediction

- **Unstructured Prediction**
  - Output Y consists of a single component



$X$ $\longrightarrow$ {Cat, Dog, Finch, Owl}

$Y$

# Structured Prediction

- Y consists of **multiple components** $Y = \{y_1, y_2, \ldots, y_n\}$



Max is teaching NLP

$\longrightarrow$

**subject**

**aux**   **object**

Max   is   teaching   NLP

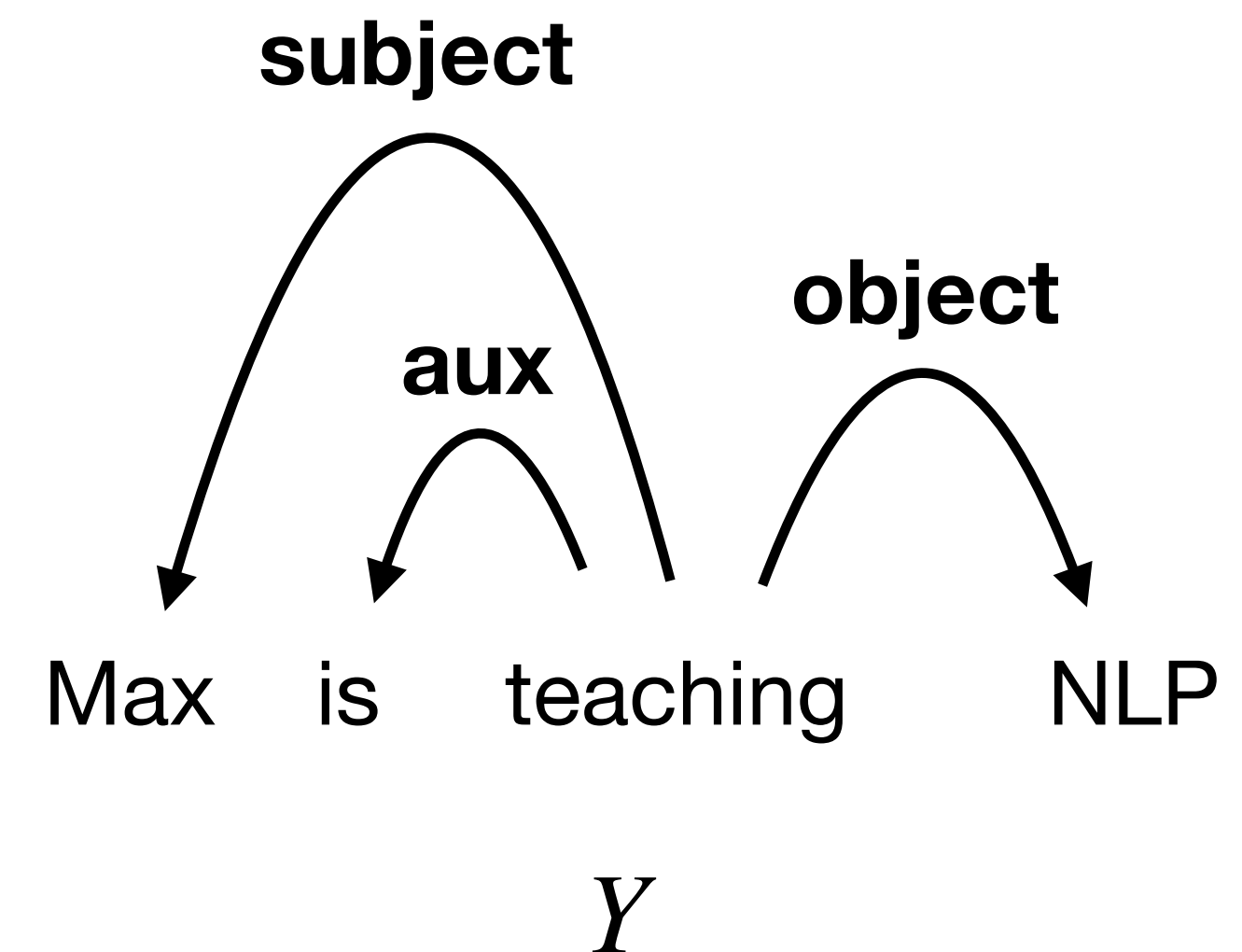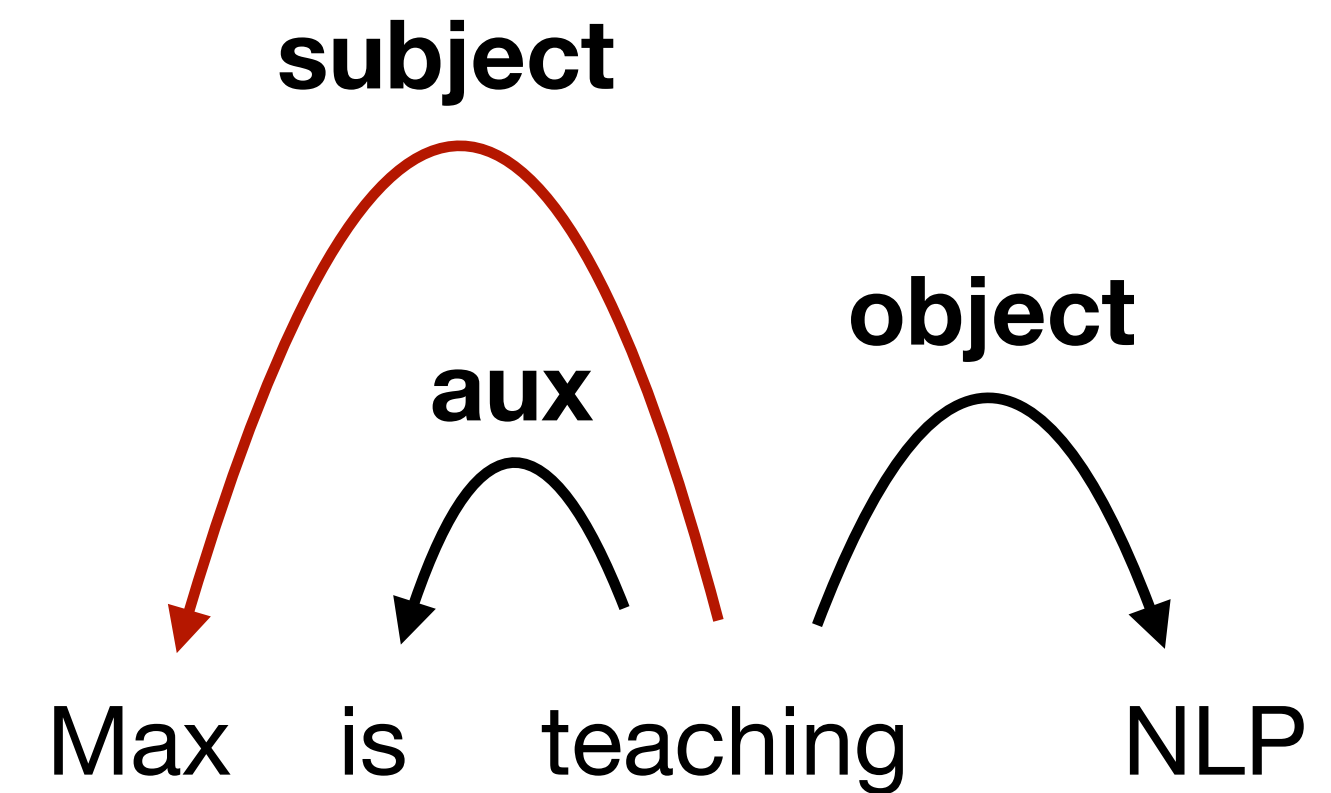$X$                              $Y$

# Structured Prediction

- Y consists of **multiple components** $Y = \{y_1, y_2, \ldots, y_n\}$



Max is teaching NLP

$X$

$$Y_1 = < \textbf{teaching} \rightarrow \textbf{Max} >$$

# Structured Prediction

- Y consists of **multiple components** $Y = \{y_1, y_2, \ldots, y_n\}$

麦克斯 在 南加大 工作

$X$

$\longrightarrow$

**Max is working at USC**

$Y_1 \quad Y_2 \qquad Y_3 \qquad Y_4 \quad Y_5$

# Structured Prediction

- Y consists of **multiple components** $Y = \{y_1, y_2, \ldots, y_n\}$
- **(Strong) correlations** between output components

麦克斯 在 南加大 工作

$X$

$\longrightarrow$

**Max is working at USC**

$Y_1 \quad Y_2 \qquad Y_3 \qquad Y_4 \quad Y_5$

# Structured Prediction

- Y consists of **multiple components** $Y = \{y_1, y_2, \ldots, y_n\}$
- **(Strong) correlations between output components**
- **Exponential output space**

  - Decoding: $y^* = \mathrm{argmax}_{y \in \mathcal{Y}} \, p(y \,|\, x)$
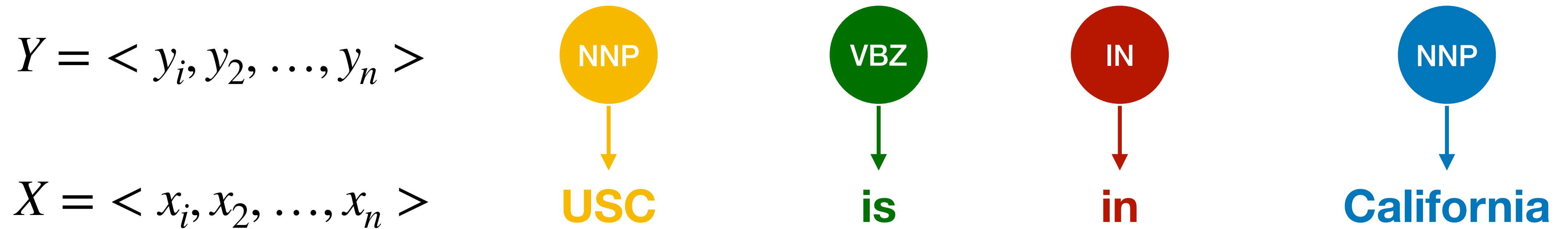
| 麦克斯 在 南加大 工作 | → | **Max is working at USC** |

$X$

$Y_1 \quad Y_2 \qquad Y_3 \qquad Y_4 \quad Y_5$

# What is Sequence Labeling?

**A type of structured prediction tasks**

$$Y = <y_i, y_2, \ldots, y_n>$$

$$X = <x_i, x_2, \ldots, x_n>$$

| NNP | VBZ | IN | NNP |
|-----|-----|-----|-----|
| USC | is | in | California |

Assigning each token of $X$, e.g. $x_i$ a corresponding label $y_i$

# Why Sequence Labeling?

## Part-of-Speech Tagging



PRP VBD DT NN IN DT NN
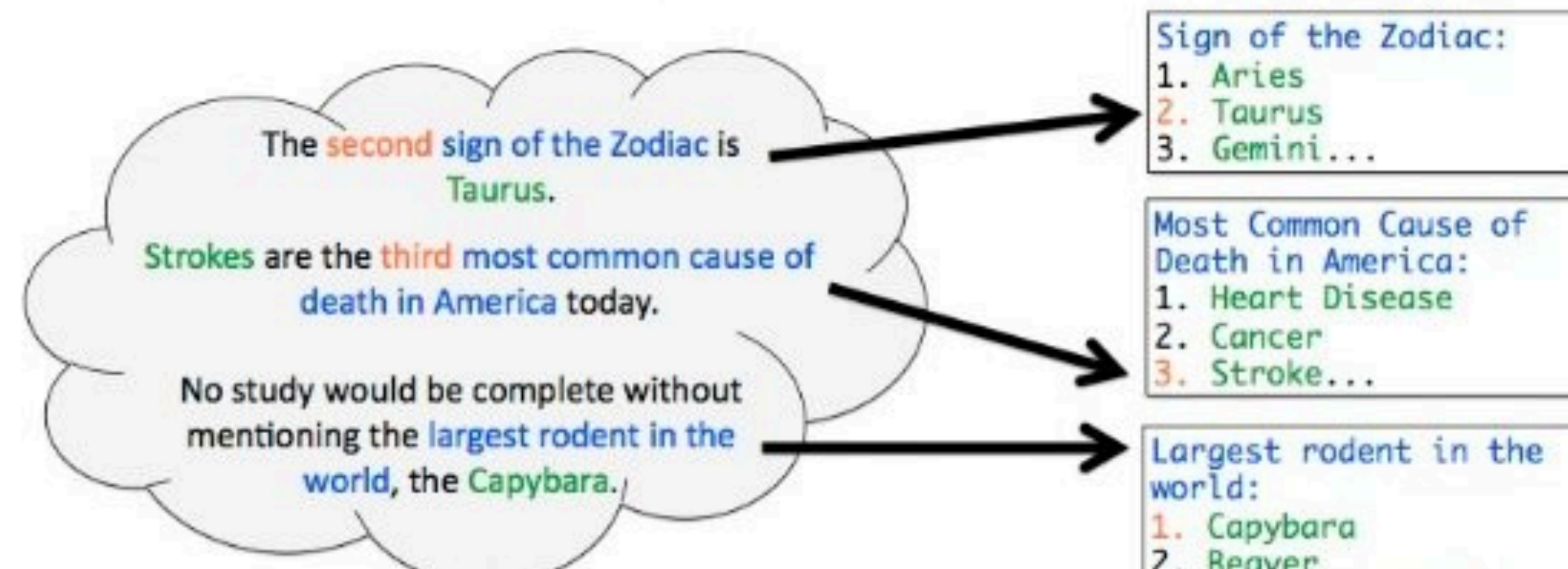I saw a girl with a telescope

## Named Entity Recognition



PERSON — Arsène Wenger was named TITLE manager of ORGANIZATION Arsenal in DATE 1996 .

## Information Extraction



**Unstructured Web Text** ➡ **Structured Sequences**

The second sign of the Zodiac is Taurus.

Strokes are the third most common cause of death in America today.

No study would be complete without mentioning the largest rodent in the world, the Capybara.

Sign of the Zodiac:
1. Aries
2. Taurus
3. Gemini...

Most Common Cause of Death in America:
1. Heart Disease
2. Cancer
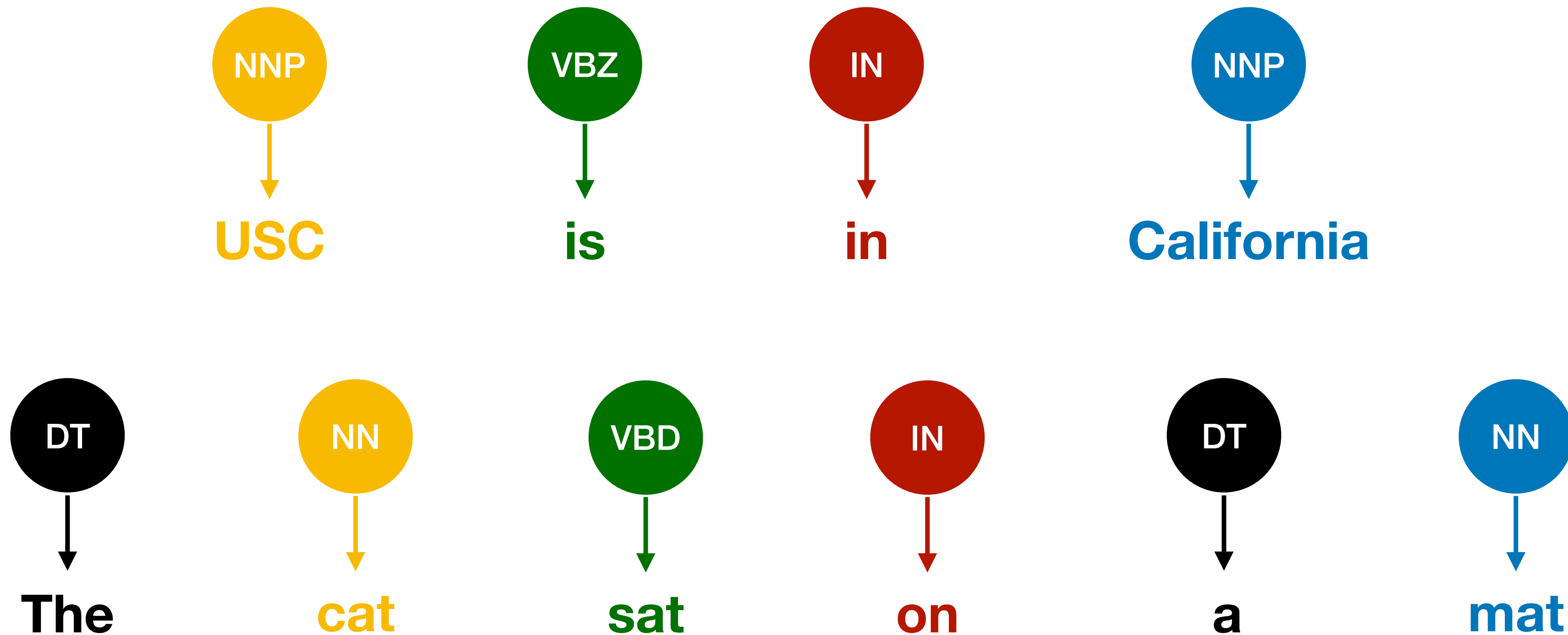3. Stroke...

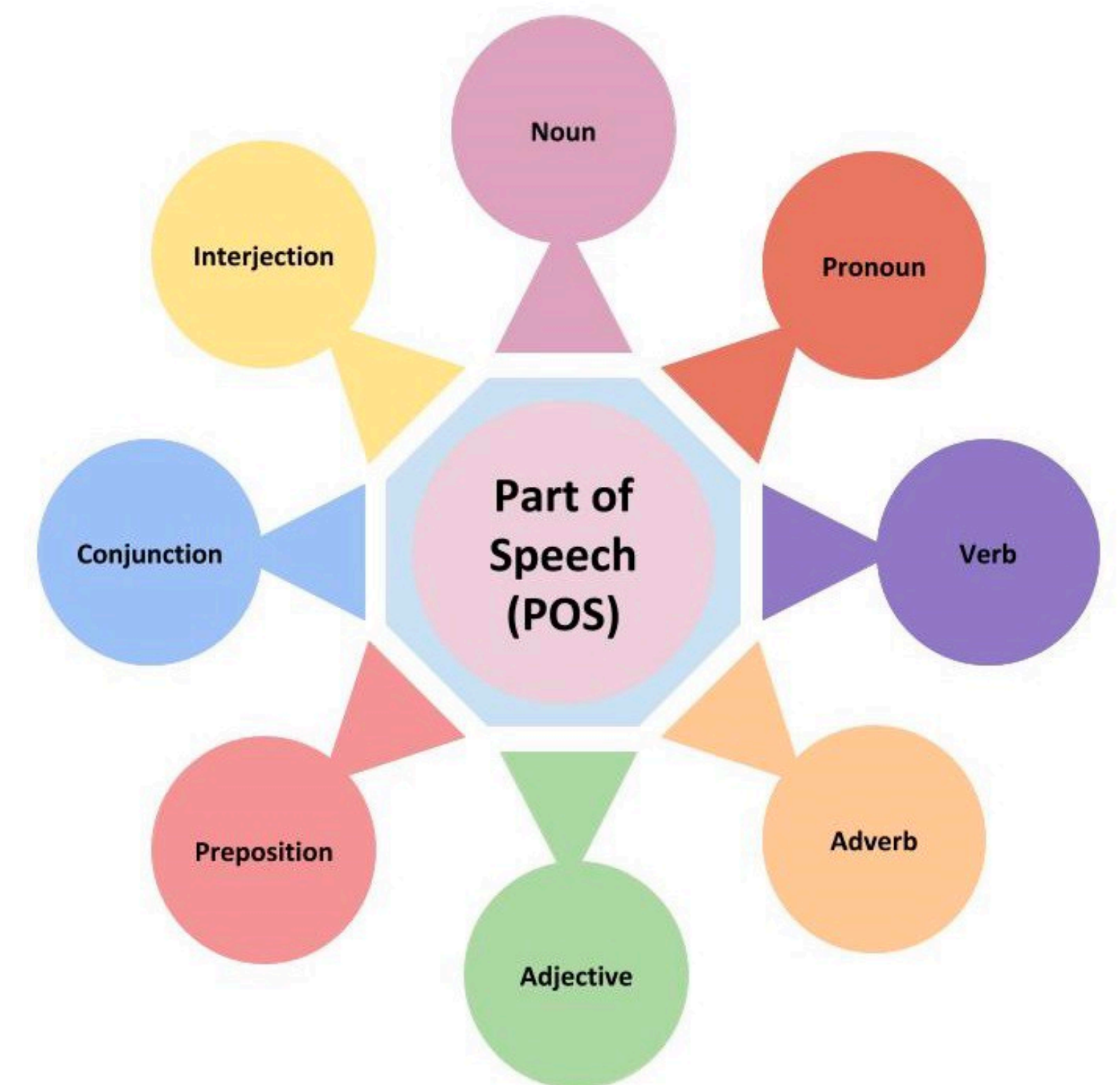Largest rodent in the world:
1. Capybara
2. Beaver

# What are Part-of-Speech (POS) Tags

- **Word classes or syntactic categories**
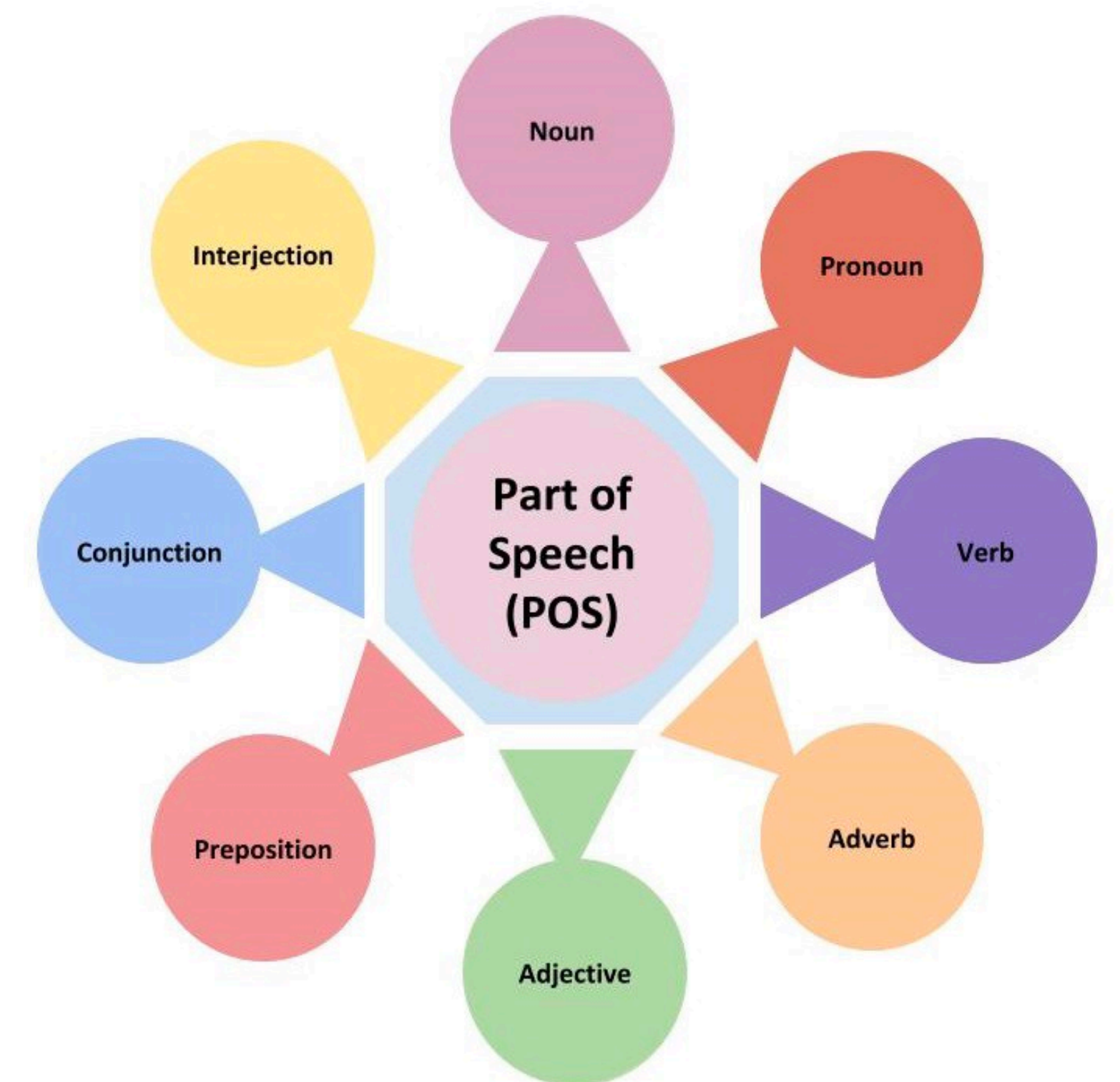- **Reveal useful information about the syntactic role of a word** (and its neighbors!)

NNP → **USC**

VBZ → **is**

IN → **in**

NNP → **California**

DT → **The**

NN → **cat**

VBD → **sat**

IN → **on**

DT → **a**

NN → **mat**

# Part of Speech

- **Different words have different syntactic functions**

- **Can be roughly divided into two classes**

  – Closed class: fixed membership, function words

    - e.g. prepositions (in, on, of), determiners (a, the)

  – Open class: New words get added frequently

    - e.g. nouns (Twitter, Facebook), verbs (google), adjectives and adverbs.

# Part of Speech

- How many part of speech tags do you think English has?
  - A. < 10
  - B. 10 - 30
  - C. 30 - 50
  - D. > 50

# Penn Tree Bank Tagset

| Tag | Description | Example | Tag | Description | Example | Tag | Description | Example |
|---|---|---|---|---|---|---|---|---|
| CC | coordinating conjunction | *and, but, or* | PDT | predeterminer | *all, both* | VBP | verb non-3sg present | *eat* |
| CD | cardinal number | *one, two* | POS | possessive ending | *'s* | VBZ | verb 3sg pres | *eats* |
| DT | determiner | *a, the* | PRP | personal pronoun | *I, you, he* | WDT | wh-determ. | *which, that* |
| EX | existential 'there' | *there* | PRP$ | possess. pronoun | *your, one's* | WP | wh-pronoun | *what, who* |
| FW | foreign word | *mea culpa* | RB | adverb | *quickly* | WP$ | wh-possess. | *whose* |
| IN | preposition/ subordin-conj | *of, in, by* | RBR | comparative adverb | *faster* | WRB | wh-adverb | *how, where* |
| JJ | adjective | *yellow* | RBS | superlatv. adverb | *fastest* | $ | dollar sign | *$* |
| JJR | comparative adj | *bigger* | RP | particle | *up, off* | # | pound sign | *#* |
| JJS | superlative adj | *wildest* | SYM | symbol | *+,%, &* | " | left quote | *' or "* |
| LS | list item marker | *1, 2, One* | TO | "to" | *to* | " | right quote | *' or "* |
| MD | modal | *can, should* | UH | interjection | *ah, oops* | ( | left paren | *[, (, {, <* |
| NN | sing or mass noun | *llama* | VB | verb base form | *eat* | ) | right paren | *], ), }, >* |
| NNS | noun, plural | *llamas* | VBD | verb past tense | *ate* | , | comma | *,* |
| NNP | proper noun, sing. | *IBM* | VBG | verb gerund | *eating* | . | sent-end punc | *. ! ?* |
| NNPS | proper noun, plu. | *Carolinas* | VBN | verb past part. | *eaten* | : | sent-mid punc | *: ; ... – -* |

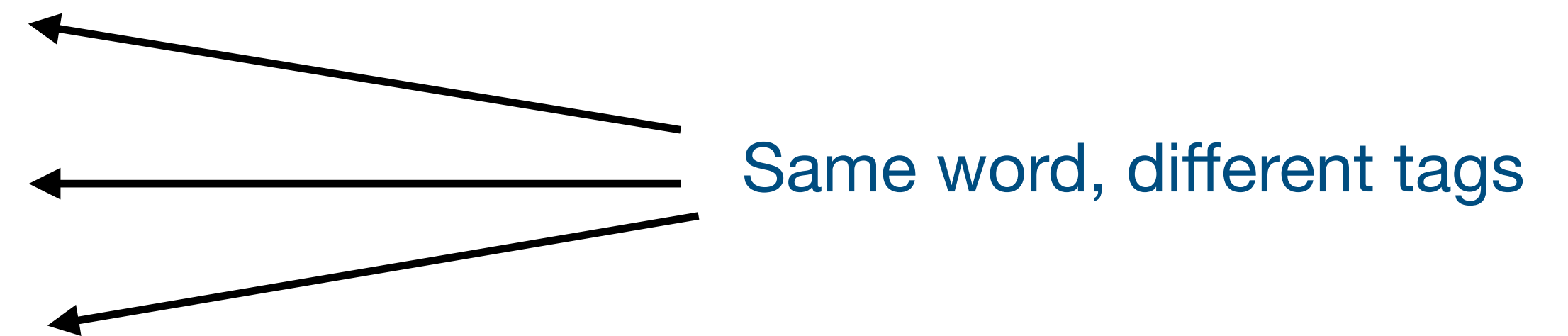## 45 tags!
## (Marcus et al., 1993)

# The Task of Part of Speech Tagging

- Tag each word with its part of speech
- Disambiguation task: each word might have different senses/functions

  – The/DT  back/ADJ  door/NN

  – On/IN   my/PRP$   back/NN          Same word, different tags

  – Win/VB  the/DT    voters/NNS  back/RP

| Types: | | WSJ | Brown |
|---|---|---|---|
| Unambiguous | (1 tag) | 44,432 (**86%**) | 45,799 (**85%**) |
| Ambiguous | (2+ tags) | 7,025 (**14%**) | 8,050 (**15%**) |
| Tokens: | | | |
| Unambiguous | (1 tag) | 577,421 (**45%**) | 384,349 (**33%**) |
| Ambiguous | (2+ tags) | 711,780 (**55%**) | 786,646 (**67%**) |

**Figure 8.2**   Tag ambiguity for word types in Brown and WSJ, using Treebank-3 (45-tag) tagging. Punctuation were treated as words, and words were kept in their original case.

# A Simple Baseline

- **Many words might be easy to disambiguate**
- **Most Frequent Class: Assign each token (word) to the class it occurred most in the training data**. (e.g. student/NN)
  - Entirely discarding contextual information
- **How accurate do you think this baseline would be at tagging words?**
  - A. < 50%
  - B. 50% - 75%
  - C. 75% - 90%
  - D. > 90%

Accurately tags **92.34%** of word tokens on Wall Street Journal (WSJ)

# POS Tagging Not Solved!

- **State of the art:** ∼ **97%**

- **Sentence level accuracies**

  – Average length of English sentence ∼ 14 words

  – $0.92^{14} = 31\%$ vs. $0.97^{14} = 65\%$

- **Highly relying on domain information**

  – Training data and testing data mush be from the same domain

  – < 70% on data from social media

# Some Observations

- **The function (or POS) of a word depends on its context**
  - The/DT  back/ADJ  door/NN

  - On/IN    my/PRP$    back/NN

  - Win/VB  the/DT      voters/NNS    back/RP

- **Certain POS combinations are extremely unlikely**

  - *<JJ, DT>* ("good the") or *<DT, IN>* ("the in")

- **Better to make predictions on entire sentences instead of individual words**

**Sequence Labeling Models!**

# Overview

- **The Sequence Labeling Problem**
  - General Structured Prediction Tasks
  - Part-of-speech Tagging: A case study
  - Generative Models vs. Discriminative Models
  - Maximum Likelihood Estimation (MLE)
- **Hidden Markov Model (HMM)**
  - Basic definitions
  - Parameter estimation
  - The Viterbi algorithm
- **Log-Linear Models**
  - Maximum Entropy Markov Models (MEMMs)
  - Conditional Random Fields (CRFs)

# Hidden Markov Models

# Markov Sequences

▶ Consider a sequence of random variables $X_1, X_2, \ldots, X_m$ where $m$ is the length of the sequence

▶ Each variable $X_i$ can take any value in $\{1, 2, \ldots, k\}$

▶ How do we model the joint distribution

$$P(X_1 = x_1, X_2 = x_2, \ldots, X_m = x_m)$$

?

23

# The Markov Assumption

$$P(X_1 = x_1, X_2 = x_2, \ldots, X_m = x_m)$$

$$= P(X_1 = x_1) \prod_{j=2}^{m} P(X_j = x_j | X_1 = x_1, \ldots, X_{j-1} = x_{j-1})$$

$$= P(X_1 = x_1) \prod_{j=2}^{m} P(X_j = x_j | X_{j-1} = x_{j-1})$$

Markov assumption

▶ The first equality is exact (by the chain rule).

▶ The second equality follows from *the Markov assumption*: for all $j = 2 \ldots m$,

$$P(X_j = x_j | X_1 = x_1, \ldots, X_{j-1} = x_{j-1}) = P(X_j = x_j | X_{j-1} = x_{j-1})$$

# Markov Sequences

- A Generative Model for Sequences



Pick $x_1$ at random from the distribution $P(X_1)$

Pick $x_2$ at random from the distribution $P(X_2 \mid X_1 = x_1)$

Pick $x_t$ at random from the distribution $P(X_t \mid X_{t-1} = x_{t-1})$

# Modeling Pairs of Sequences

- In Sequence Labeling, we need to model pairs of sequences

$$S = S_i, S_2, \ldots, S_n$$

| NNP | VBZ | IN | NNP |

$$X = X_i, X_2, \ldots, X_n$$

**USC**          **is**          **in**          **California**

Hidden Markov Models (HMMs) allow us to *jointly* reason over $X$ and $S$

# Hidden Markov Models

▶ We have two sequences of random variables:
$X_1, X_2, \ldots, X_m$ and $S_1, S_2, \ldots, S_m$

▶ Intuitively, each $X_i$ corresponds to an "observation" and each $S_i$ corresponds to an underlying "state" that generated the observation. Assume that each $S_i$ is in $\{1, 2, \ldots k\}$, and each $X_i$ is in $\{1, 2, \ldots o\}$

▶ How do we model the joint distribution

$$P(X_1 = x_1, \ldots, X_m = x_m, S_1 = s_1, \ldots, S_m = s_m)$$

?

# The HMM Assumptions



1. **Markov Assumption on $S$**

$$P(S_j = s_j \mid S_{j-1} = s_{j-1}, \ldots, S_1 = s_1) = P(S_j = s_j \mid S_{j-1} = s_{j-1})$$
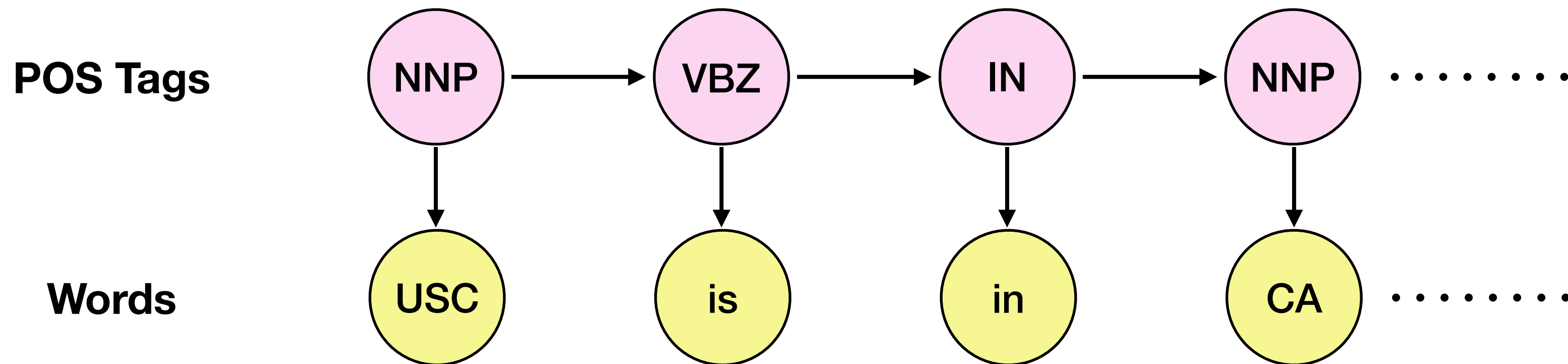
Transition Probabilities

2. **Conditional Independence on $X$ given $S$**

$$P(X_1 = x_j, \ldots, X_m = x_m \mid S_1 = s_1, \ldots, S_m = s_m) = \prod_{j=1}^{m} P(X_j = x_j \mid S_j = s_j)$$

Emission Probabilities

# The HMM Assumptions



**POS Tags**: NNP → VBZ → IN → NNP ·········

**Words**: USC, is, in, CA ·········

1. **Markov Assumption on $S$**

$$P(S_3 = \text{IN} \mid S_2 = \text{VBZ}, S_1 = \text{NNP}) = P(S_3 = \text{IN} \mid S_2 = \text{VBZ})$$

2. **Conditional Independence on $X$ given $S$**

$$P(\text{USC is in CA} \mid \text{NNP VBZ IN NNP}) = P(\text{USC} \mid \text{NNP})P(\text{is} \mid \text{VBZ})P(\text{in} \mid \text{IN})P(\text{CA} \mid \text{NNP})$$

Which assumption do you think is stronger?

Second assumption is stronger..why? Becoz its more
powerful than the markov assumption

29

# Joint Distribution of Sequence Pairs in HMMs

$$P(X_1 = x_j, \ldots, X_m = x_m, S_1 = s_1, \ldots, S_m = s_m)$$

$$= P(X_1 = x_j, \ldots, X_m = x_m \mid S_1 = s_1, \ldots, S_m = s_m)$$

Output Independence

$$\times P(S_1 = s_1, \ldots, S_m = s_m)$$

Markov Assumption

$$= \prod_{j=1}^{m} P(X_j = x_j \mid S_j = s_j)$$

How to model $P(X_j = x_j \mid S_j = s_j)$ and $P(S_j = s_j \mid S_{j-1} = s_{j-1})$?

$$\times P(S_1 = s_1) \prod_{j=1}^{m} P(S_j = s_j \mid S_{j-1} = s_{j-1})$$

30

# Homogeneous HMMs

- In a *homogeneous* HMM, we make an additional assumption:

$$P(S_j = s_j \mid S_{j-1} = s_{j-1}) = t(s_j \mid s_{j-1})$$

$$P(X_j = x_j \mid S_j = s_j) = e(x_j \mid s_j)$$

- Idea behind this assumption: the transition and emission probabilities do not depend on the position in the Markov chain (do not depend on the index $j$)
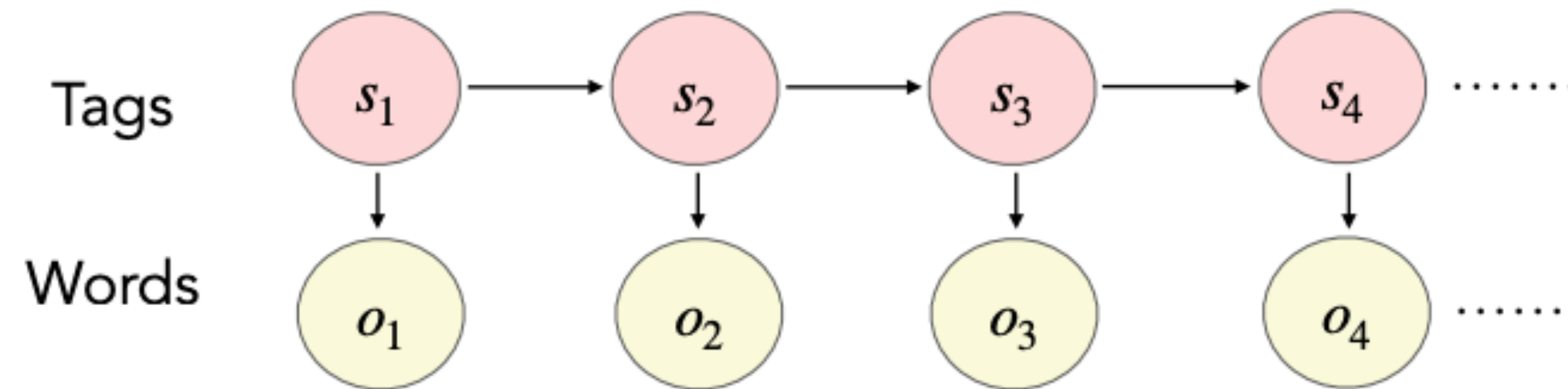
# The Model Form for Homogeneous HMMs

▶ The model takes the following form:

$$p(x_1 \ldots x_m, s_1 \ldots s_m; \underline{\theta}) = t(s_1) \prod_{j=2}^{m} t(s_j|s_{j-1}) \prod_{j=1}^{m} e(x_j|s_j)$$

▶ Parameters in the model:

1. Initial state parameters $t(s)$ for $s \in \{1, 2, \ldots, k\}$

2. Transition parameters $t(s'|s)$ for $s, s' \in \{1, 2, \ldots, k\}$

3. Emission parameters $e(x|s)$ for $s \in \{1, 2, \ldots, k\}$ and $x \in \{1, 2, \ldots, o\}$

# Example: Sequence Probability



Tags

$s_1 \rightarrow s_2 \rightarrow s_3 \rightarrow s_4$ .......

Words

$o_1 \quad o_2 \quad o_3 \quad o_4$ .......

Dummy start state

$s_{t+1}$

| $s_t$ | | DT | NN |
|---|---|---|---|
| | $\emptyset$ | 0.8 | 0.2 |
| | DT | 0.2 | 0.8 |
| | NN | 0.3 | 0.7 |

$o_t$

| | | the | cat |
|---|---|---|---|
| | DT | 0.9 | 0.1 |
| | NN | 0.5 | 0.5 |

*What is the joint probability*
*$P$(the cat, DT NN)?*

*A)  $(0.8 * 0.8) * (0.9 * 0.5)$*
*B)  $(0.2 * 0.8) * (0.9 * 0.5)$*
*C)  $(0.3 * 0.7) * (0.5 * 0.5)$*

# HMMs are Generative Models



1. Pick $s_1$ at random from the distribution $t(s)$. Pick $x_1$ from the distribution $e(x|s_1)$

2. For $j = 2 \ldots m$:

   ▸ Choose $s_j$ at random from the distribution $t(s|s_{j-1})$

   ▸ Choose $x_j$ at random from the distribution $e(x|s_j)$

# HMMs are Generative Models



1. Pick $s_1$ at random from the distribution $t(s)$. Pick $x_1$ from the distribution $e(x|s_1)$

2. For $j = 2 \ldots m$:

   ▶ Choose $s_j$ at random from the distribution $t(s|s_{j-1})$

   ▶ Choose $x_j$ at random from the distribution $e(x|s_j)$

# Learning a Hidden Markov Model

# Parameter Estimation

- **Assuming we have fully observed data** $\{X_i, S_i\}_{i=1}^{N}$, **e.g. WSJ**

**Training set:**

1 Pierre/NNP Vinken/NNP ,/, 61/CD years/NNS old/JJ ,/
join/VB the/DT board/NN as/IN a/DT nonexecutive/JJ di
Nov./NNP 29/CD ./.
2 Mr./NNP Vinken/NNP is/VBZ chairman/NN of/IN Elsev
N.V./NNP ,/, the/DT Dutch/NNP publishing/VBG group/
3 Rudolph/NNP Agnew/NNP ,/, 55/CD years/NNS old/JJ
chairman/NN of/IN Consolidated/NNP Gold/NNP Fields/N
,/, was/VBD named/VBN a/DT nonexecutive/JJ director/
this/DT British/JJ industrial/JJ conglomerate/NN ./.
...
38,219 It/PRP is/VBZ also/RB pulling/VBG 20/CD peopl
of/IN Puerto/NNP Rico/NNP ,/, who/WP were/VBD help
Huricane/NNP Hugo/NNP victims/NNS ,/, and/CC sendin
them/PRP to/TO San/NNP Francisco/NNP instead/RB ./

**Maximum Likelihood Estimate:**

$$\max_{t(\cdot|\cdot), e(\cdot|\cdot)} \prod_{i=1}^{N} P(X_i, S_i)$$

$$t(s'|s) = \frac{\text{count}(s \to s')}{\text{count}(s)}$$
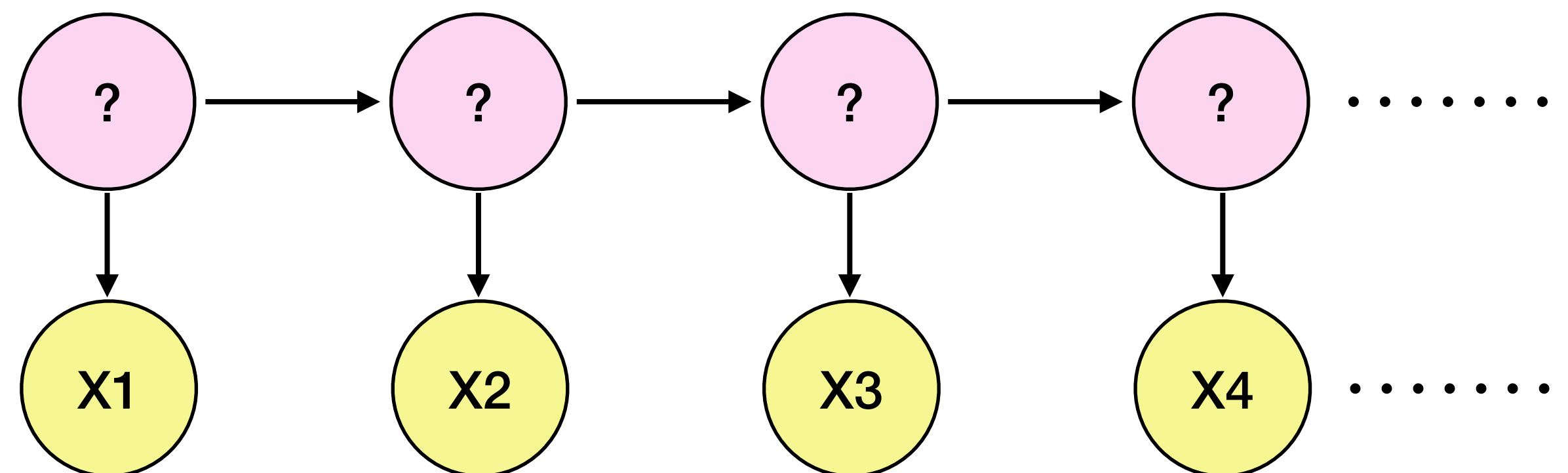
$$e(x|s) = \frac{\text{count}(s \to x)}{\text{count}(s)}$$

# Learning Example

1. the/DT cat/NN sat/VBD on/IN the/DT mat/NN

2. Princeton/NNP is/VBZ in/IN New/NNP Jersey/NNP

3. the/DT old/NN man/VB the/DT boats/NNS

$$t(\textbf{NN}\,|\,\textbf{DT}) = \frac{3}{4}$$

$$e(\textbf{cat}\,|\,\textbf{NN}) = \frac{1}{3}$$

**Maximum Likehood Estimate:**

$$\max_{t(\cdot\,|\,\cdot),e(\cdot\,|\,\cdot)} \prod_{i=1}^{N} P(X_i, S_i)$$

$$t(s'\,|\,s) = \frac{\text{count}(s \rightarrow s')}{\text{count}(s)}$$

$$e(x\,|\,s) = \frac{\text{count}(s \rightarrow x)}{\text{count}(s)}$$

# Decoding with HMMs
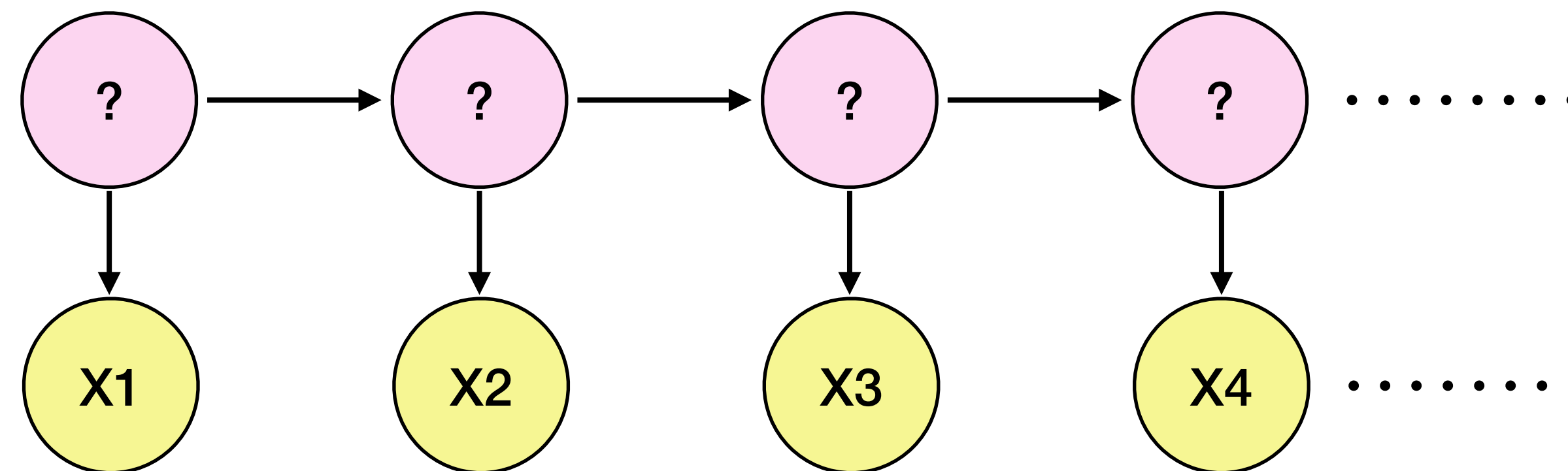
# Decoding with HMMs

▶ Goal: for a given input sequence $x_1, \ldots, x_m$, find

$$\arg\max_{s_1,\ldots,s_m} p(x_1 \ldots x_m, s_1 \ldots s_m; \underline{\theta})$$

▶ This is the most likely state sequence $s_1 \ldots s_m$ for the given input sequence $x_1 \ldots x_m$
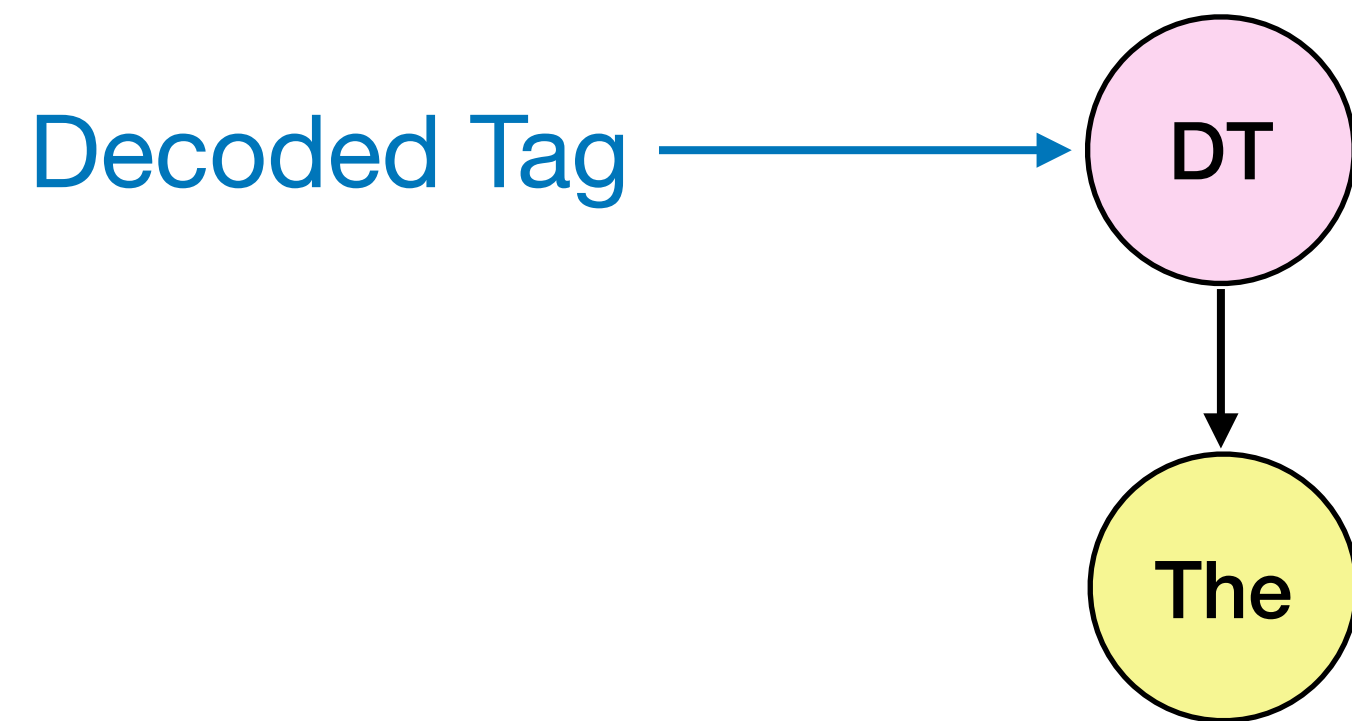
# Decoding with HMMs



$$S^* = \arg\max_{s_1,\ldots,s_m} p(x_1, \ldots, x_m, s_1, \ldots, s_m) = t(s_1)\prod_{j=2}^{m} t(s_j \,|\, s_{j-1})\prod_{j=1}^{m} e(x_j \,|\, s_j)$$

How can we maximize this over all state sequences?

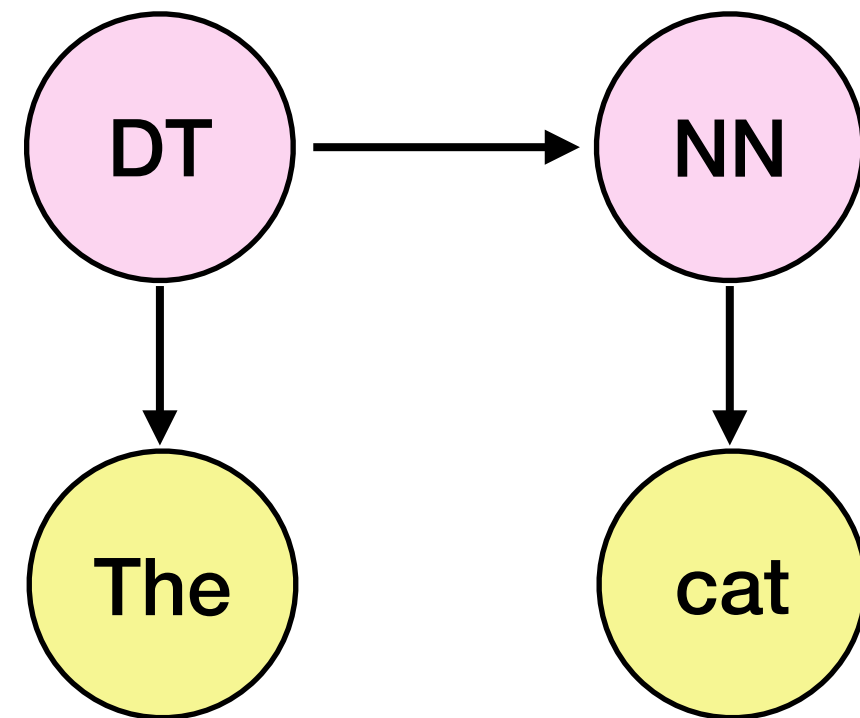# Greedy Decoding

Decoded Tag → ( DT )

( DT ) → ( The )

Decode/reveal one state at a time

$$s_1^* = \arg\max_{s_1} t(s_1)e(x_1 \mid s_1)$$

$$S^* = \arg\max_{s_1,\ldots,s_m} p(x_1, \ldots, x_m, s_1, \ldots, s_m) = t(s_1)\prod_{j=2}^{m} t(s_j \mid s_{j-1})\prod_{j=1}^{m} e(x_j \mid s_j)$$
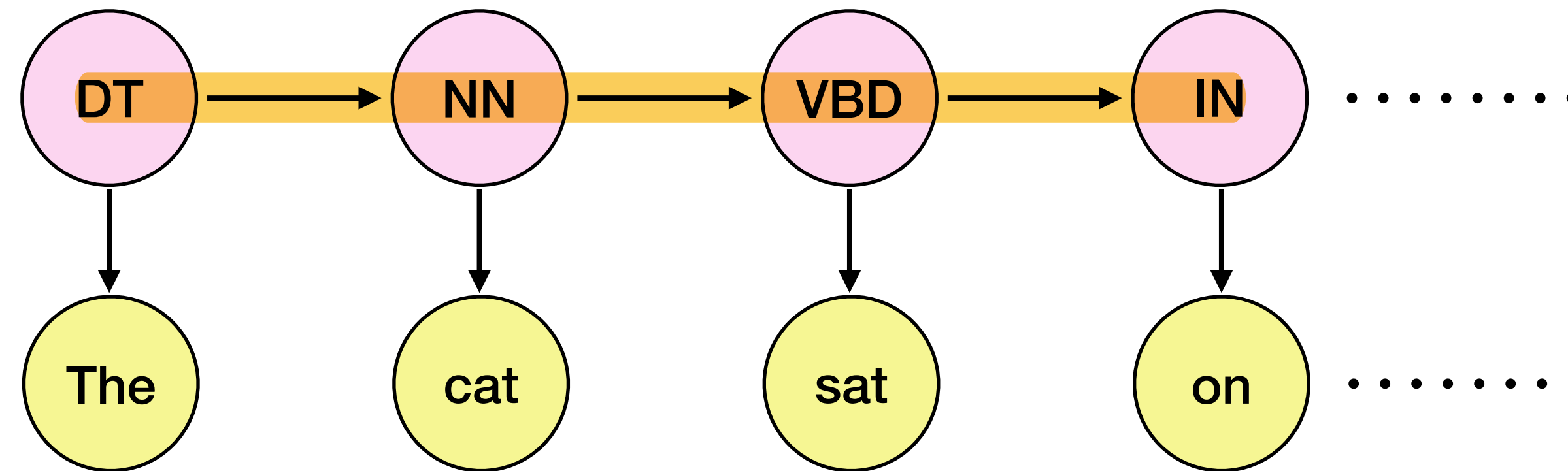
# Greedy Decoding



Decode/reveal one state at a time

$$s_2^* = \arg\max_{s_2} t(s_2 \mid s_1^*)e(x_2 \mid s_2)$$

$$S^* = \arg\max_{s_1,\ldots,s_m} p(x_1, \ldots, x_m, s_1, \ldots, s_m) = t(s_1)\prod_{j=2}^{m} t(s_j \mid s_{j-1})\prod_{j=1}^{m} e(x_j \mid s_j)$$

# Greedy Decoding



$$s_j^* = \arg\max_{s_j} t(s_j \mid s_{j-1}^*)e(x_j \mid s_j), \quad \forall j$$

TIME COMPLEXITY OF THE
ALGO: O(KN) where k is the
number of the number of the
tags

- Local decisions
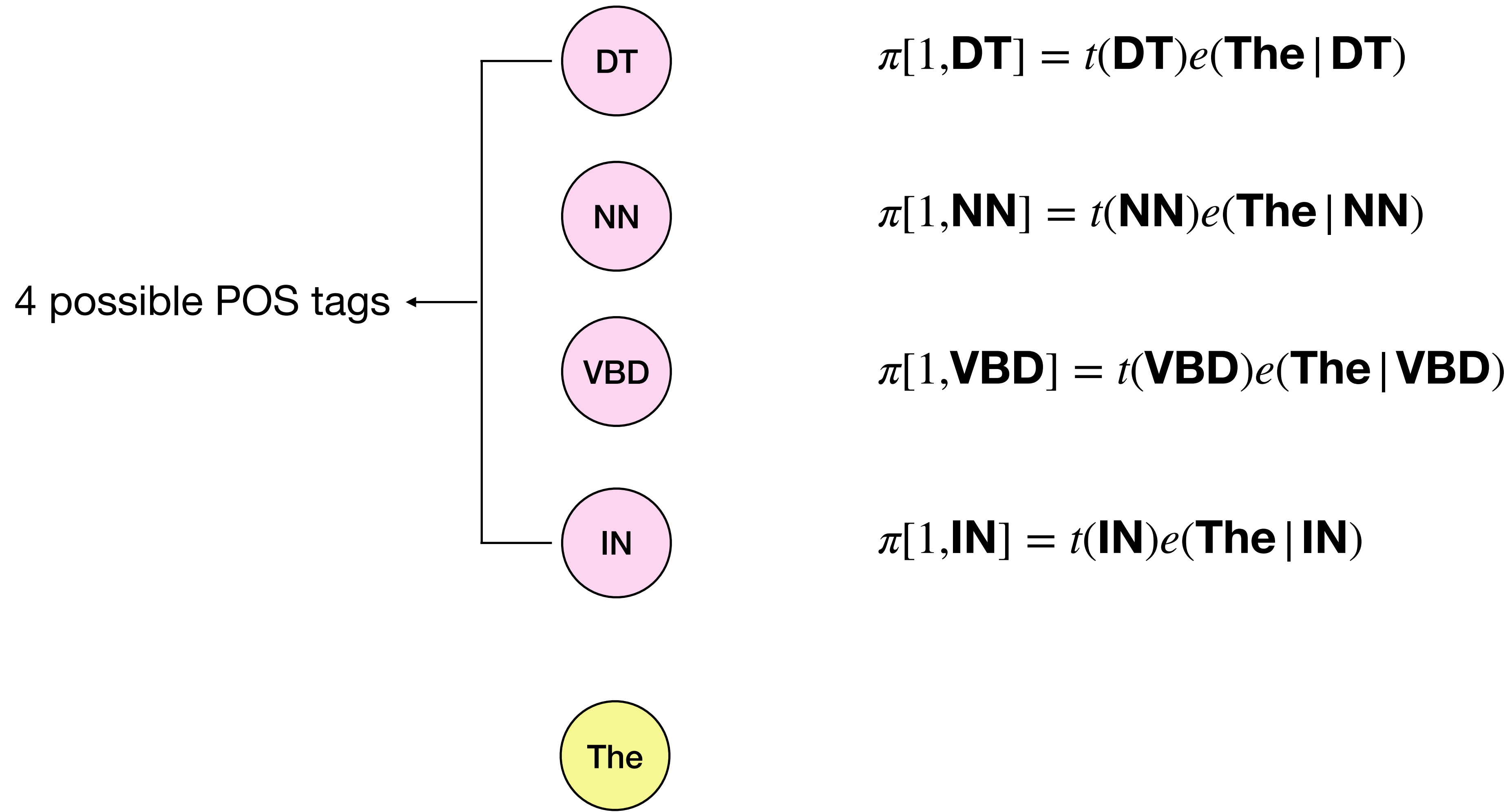- Not guaranteed to produce the overall optimal sequence

# Viterbi Decoding

▶ The *Viterbi algorithm* is a dynamic programming algorithm. Basic data structure:
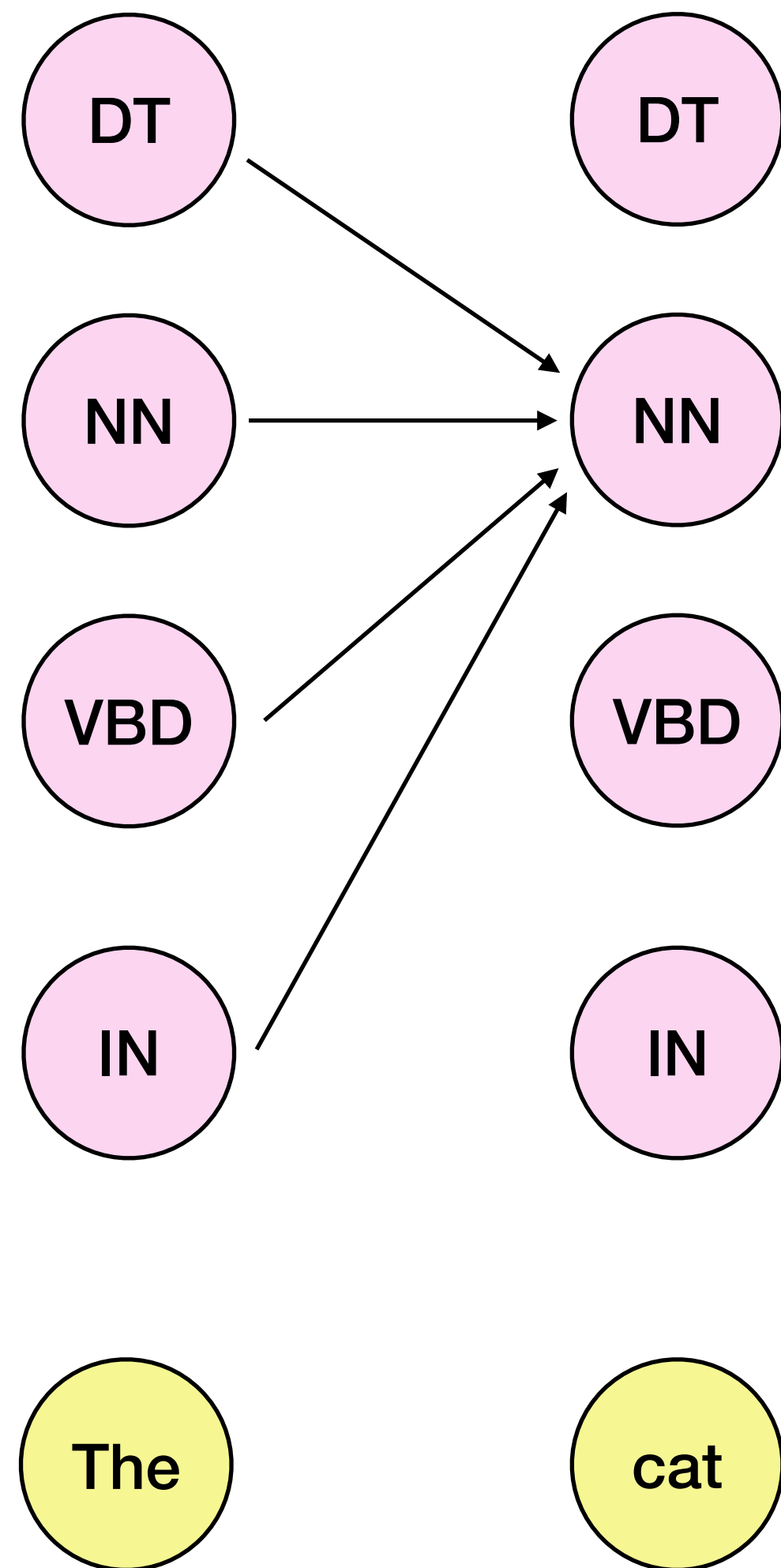
$$\pi[j, s]$$

will be a table entry that stores the maximum probability for any state sequence ending in state $s$ at position $j$. More formally: $\pi[1, s] = t(s)e(x_1|s)$, and for $j > 1$,

$$\pi[j, s] = \max_{s_1 \ldots s_{j-1}} \left[ t(s_1)e(x_1|s_1) \left( \prod_{k=2}^{j-1} t(s_k|s_{k-1})e(x_k|s_k) \right) \boxed{t(s|s_{j-1})} e(x_j|s) \right]$$

# Viterbi Decoding

$\pi[1, \textbf{DT}] = t(\textbf{DT})e(\textbf{The}|\textbf{DT})$

$\pi[1, \textbf{NN}] = t(\textbf{NN})e(\textbf{The}|\textbf{NN})$

4 possible POS tags

$\pi[1, \textbf{VBD}] = t(\textbf{VBD})e(\textbf{The}|\textbf{VBD})$

$\pi[1, \textbf{IN}] = t(\textbf{IN})e(\textbf{The}|\textbf{IN})$
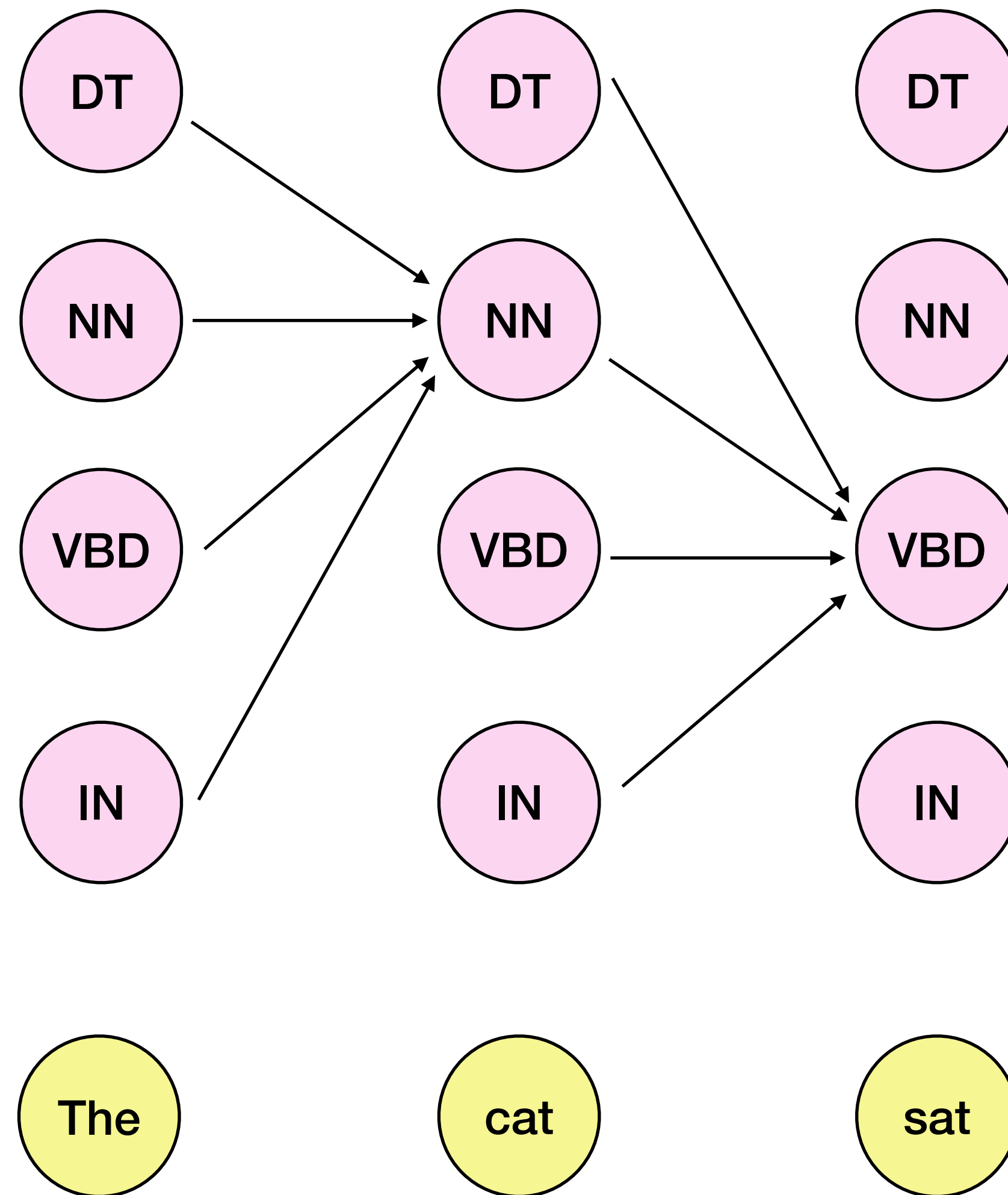
DT

NN

VBD

IN

The

# Viterbi Decoding



$$\pi[2,\mathbf{NN}] = \max_{s_1} \boxed{\pi[1, \ s_1]} \ t(\mathbf{NN}|s_1)e(\mathbf{cat}|\mathbf{NN})$$

# Viterbi Decoding



$$\pi[2,\mathbf{NN}] = \max_{s_1} \pi[1, \; s_1] \; t(\mathbf{NN}|s_1)e(\mathbf{cat}|\mathbf{NN})$$

$$\pi[3,\mathbf{VBD}] = \max_{s_2} \pi[2, \; s_2] \; t(\mathbf{VBD}|s_2)e(\mathbf{sat}|\mathbf{VBD})$$

# Viterbi Decoding

▶ Initialization: for $s = 1 \ldots k$

$$\pi[1, s] = t(s)e(x_1|s)$$

▶ For $j = 2 \ldots m$, $s = 1 \ldots k$:

$$\pi[j, s] = \max_{s' \in \{1 \ldots k\}} [\pi[j-1, s'] \times t(s|s') \times e(x_j|s)]$$

▶ We then have

$$\max_{s_1 \ldots s_m} p(x_1 \ldots x_m, s_1 \ldots s_m; \underline{\theta}) = \max_s \pi[m, s]$$

▶ The algorithm runs in $O(mk^2)$ time

# Pros and Cons

- **HMMs are simple to train**
  - Just need to compile counts from the training corpus

- **Performs relatively well**
  - > 96% on POS tagging (92.3% of most frequent class)
  - > 90% on Named Entity Recognition

- **Main difficulty is modeling $e(word|tag)$**
  - Words are very complex
  - Unknown words

# Reading Materials

- **Notes from Michael Collins:**
  - Sequence Labeling and HMMs
  - EM Algorithm