



CMPT 825: Natural Language Processing

# Sequence Models

Spring 2020

Adapted from slides from Danqi Chen and Karthik Narasimhan  
(Princeton COS 484)

# Overview

- Hidden markov models (HMM)
- Viterbi algorithm
- Maximum entropy markov models (MEMM)

# Sequence Tagging

Input: sequence of words; Output: sequence of labels

Input	British	left	waffles	on	Falkland	Islands
-------	---------	------	---------	----	----------	---------

Output1	N	N	V	P	N	N
---------	---	---	---	---	---	---

Output2	N	V	N	P	N	N
---------	---	---	---	---	---	---

:

N Noun, e.g. islands

V Verb, e.g. leave, left

P Preposition, e.g. on

# What are POS tags

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

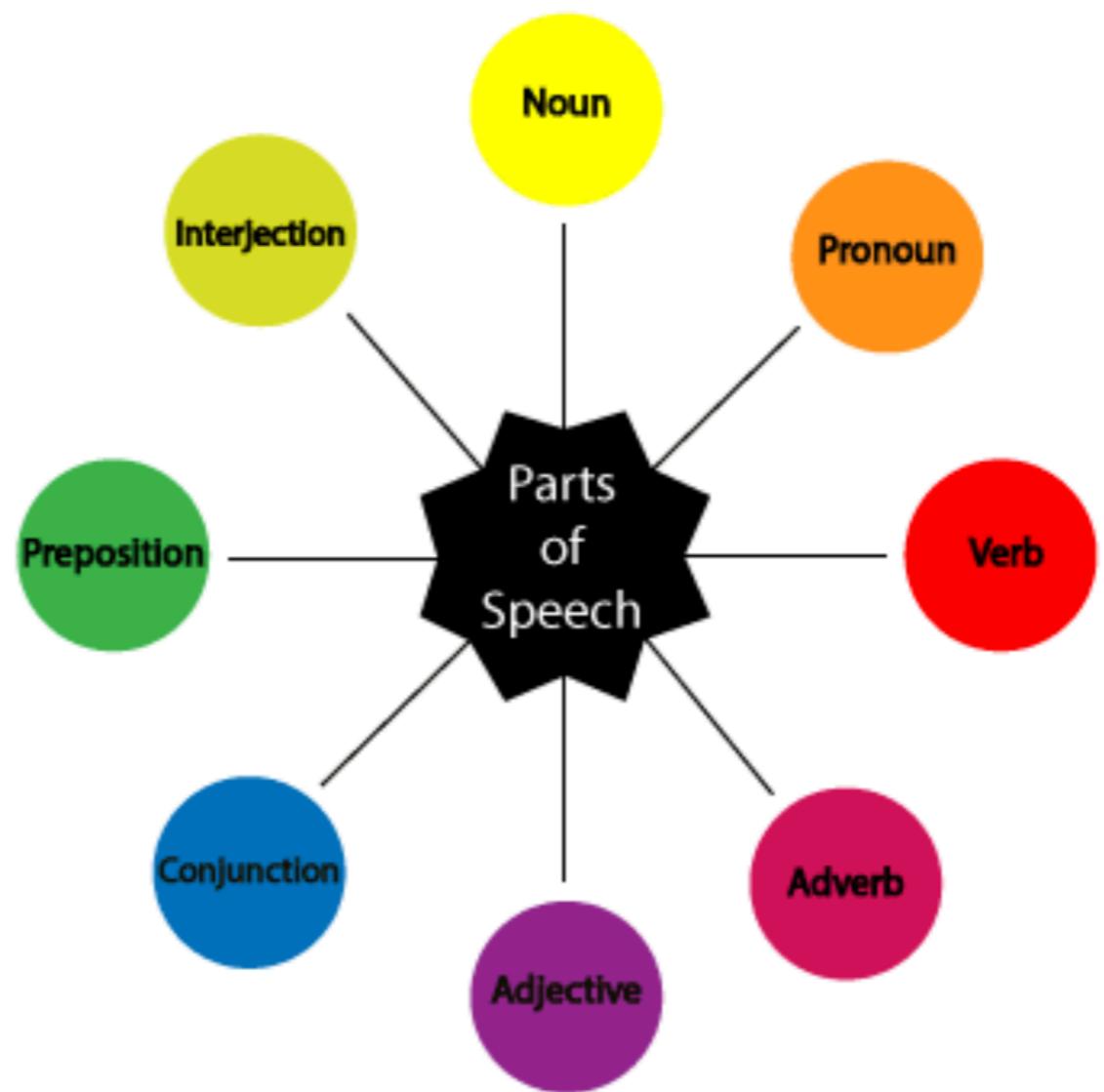
The/**DT** cat/**NN** sat/**VBD** on/**IN** the/**DT** mat/**NN**

British/**NNP** left/**NN** waffles/**NNS** on/**IN** Falkland/**NNP** Islands/**NNP**

The/**DT** old/**NN** man/**VB** the/**DT** boat/**NN**

# Parts of Speech

- Different words have different functions
- Closed class: fixed membership,**function words**
  - e.g. prepositions (*in, on, of*), determiners (*the, a*)
- Open class: New words get added frequently
  - e.g. nouns (Twitter, Facebook), verbs (google), adjectives, adverbs



# Penn Tree Bank tagset

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

[45 tags]

**Figure 8.1** Penn Treebank part-of-speech tags (including punctuation).

(Marcus et al., 1993)

Other corpora: Brown, WSJ, Switchboard

# Part of Speech Tagging

- Disambiguation task: each word might have different senses/functions
  - The/DT man/NN bought/VBD a/DT boat/NN
  - The/DT old/NN man/VB the/DT boat/NN

Types:	WSJ	Brown
<b>Unambiguous</b> (1 tag)	44,432 (86%)	45,799 (85%)
<b>Ambiguous</b> (2+ tags)	7,025 (14%)	8,050 (15%)
Tokens:		
<b>Unambiguous</b> (1 tag)	577,421 (45%)	384,349 (33%)
<b>Ambiguous</b> (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.

# Part of Speech Tagging

- Disambiguation task: each word might have different senses/functions
  - The/DT man/NN bought/VBD a/DT boat/NN
  - The/DT old/NN man/VB the/DT boat/NN

earnings growth took a **back/JJ** seat  
a small building in the **back/NN**  
a clear majority of senators **back/VBP** the bill  
Dave began to **back/VB** toward the door  
enable the country to buy **back/RP** about debt  
I was twenty-one **back/RB** then

Some words have  
many functions!

# 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 set. (e.g. man/NN)
- Accurately tags **92.34%** of word tokens on Wall Street Journal (WSJ)!
- State of the art ~ 97%
- Average English sentence ~ 14 words
  - Sentence level accuracies:  $0.92^{14} = \textbf{31\%}$  vs  $0.97^{14} = \textbf{65\%}$
- POS tagging not solved yet!

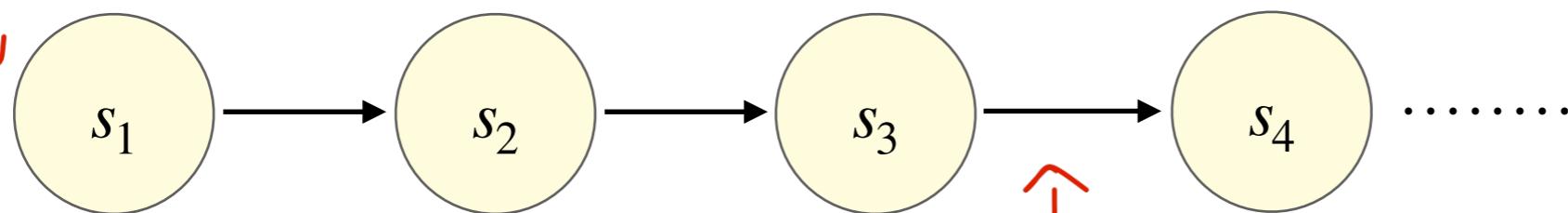
# Hidden Markov Models

# Some observations

- The function (or POS) of a word depends on its context
  - The/DT old/NN man/VB the/DT boat/NN
  - The/DT old/JJ man/NN bought/VBD the/DT boat/NN
- Certain POS combinations are extremely unlikely
  - $\langle JJ, DT \rangle$  or  $\langle DT, IN \rangle$
- Better to make decisions on entire sequences instead of individual words (**Sequence modeling!**)

# Markov chains

$\pi(s_1)$ : Initial distribution

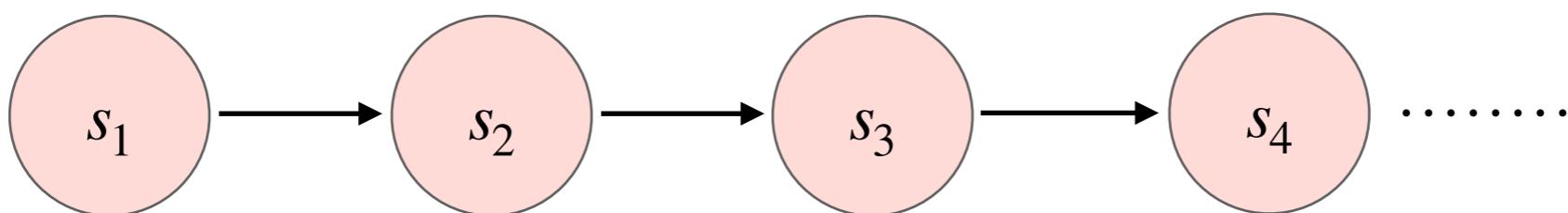


$p(s_t | s_{t-1})$ : Transition probability

- Model probabilities of sequences of variables
- Each state can take one of K values ( $\{1, 2, \dots, K\}$  for simplicity)
- Markov assumption:  $P(s_t | s_{<t}) \approx P(s_t | s_{t-1})$

Where have we seen this before?

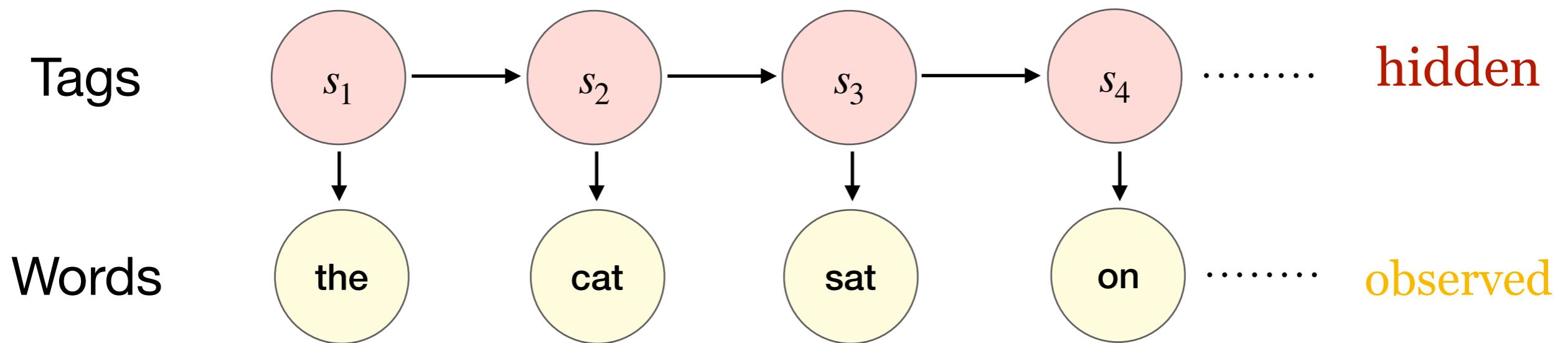
# Markov chains



The/**??** cat/**??** sat/**??** on/**??** the/**??** mat/**??**

- We don't observe POS tags at test time

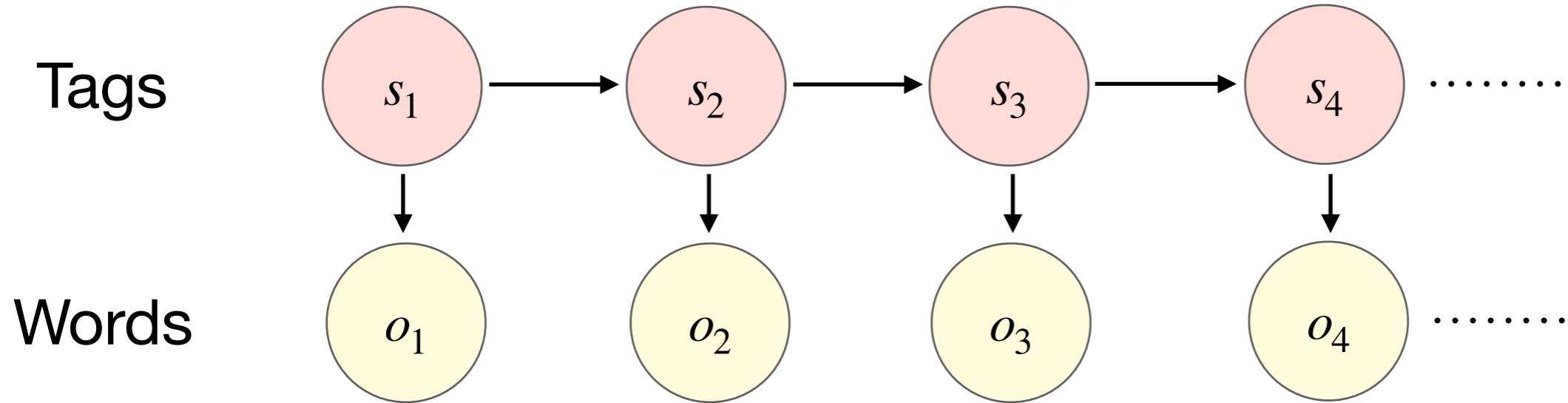
# Hidden Markov Model (HMM)



The/?? cat/?? sat/?? on/?? the/?? mat/??

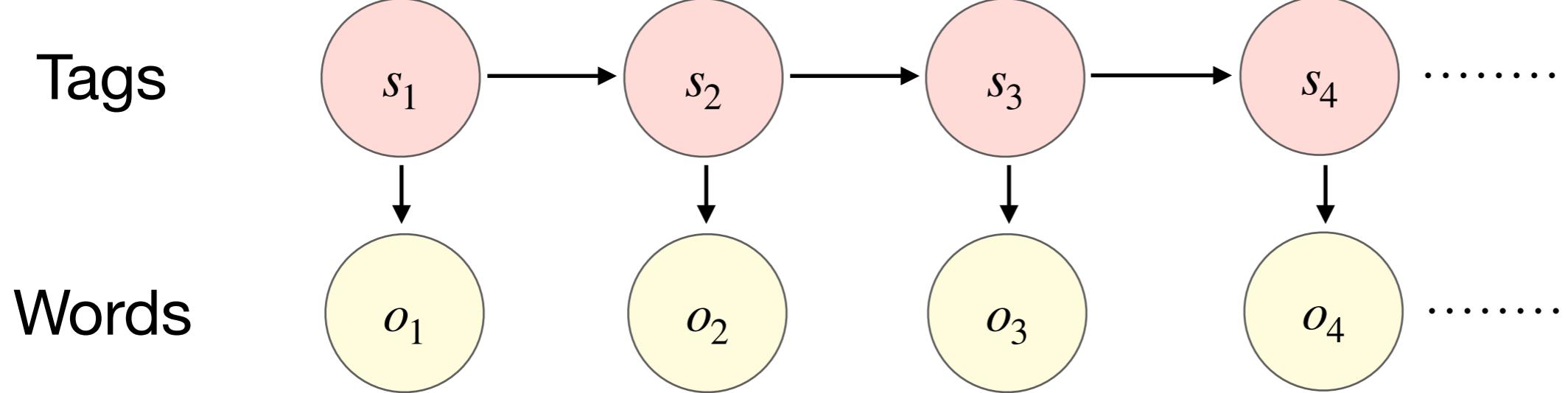
- We don't observe POS tags at test time
- But we do observe the words!
- HMM allows us to *jointly reason* over both **hidden** and **observed** events.

# Components of an HMM



1. Set of states  $S = \{1, 2, \dots, K\}$  and observations  $O$
2. Initial state probability distribution  $\pi(s_1)$
3. Transition probabilities  $P(s_{t+1} | s_t)$
4. Emission probabilities  $P(o_t | s_t)$

# Assumptions



1. Markov assumption:

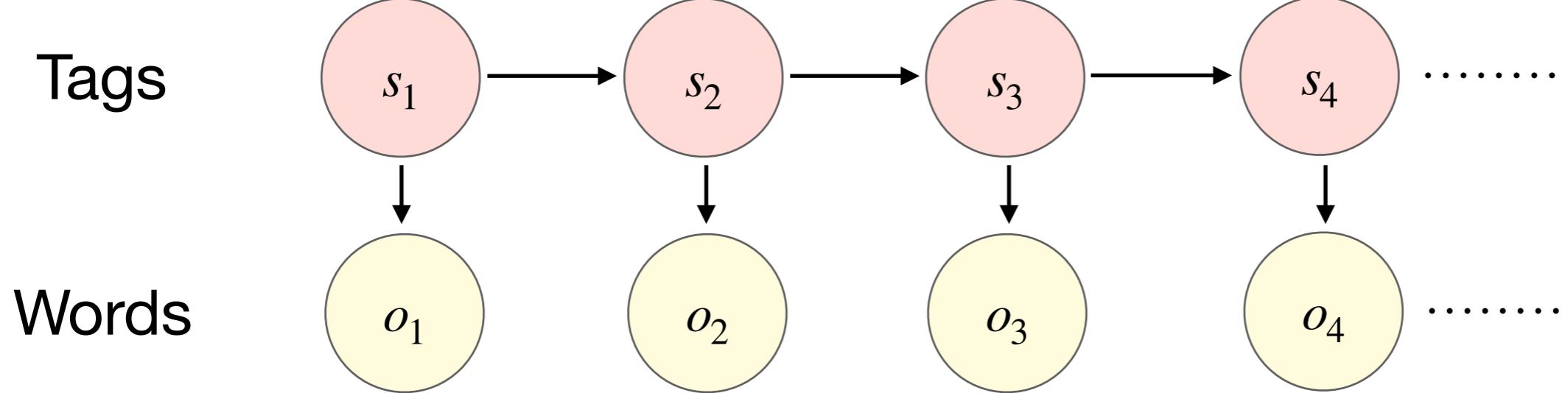
$$P(s_{t+1} | s_1, \dots, s_t) = P(s_{t+1} | s_t)$$

2. Output independence:

$$P(o_t | s_1, \dots, s_t) = P(o_t | s_t)$$

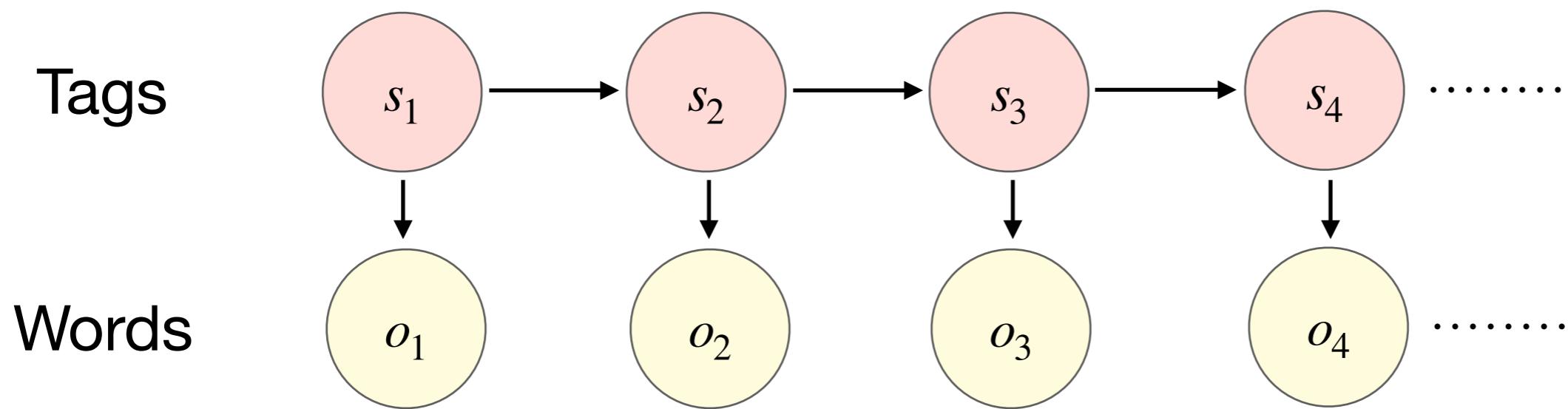
Which is a stronger assumption?

# Sequence likelihood



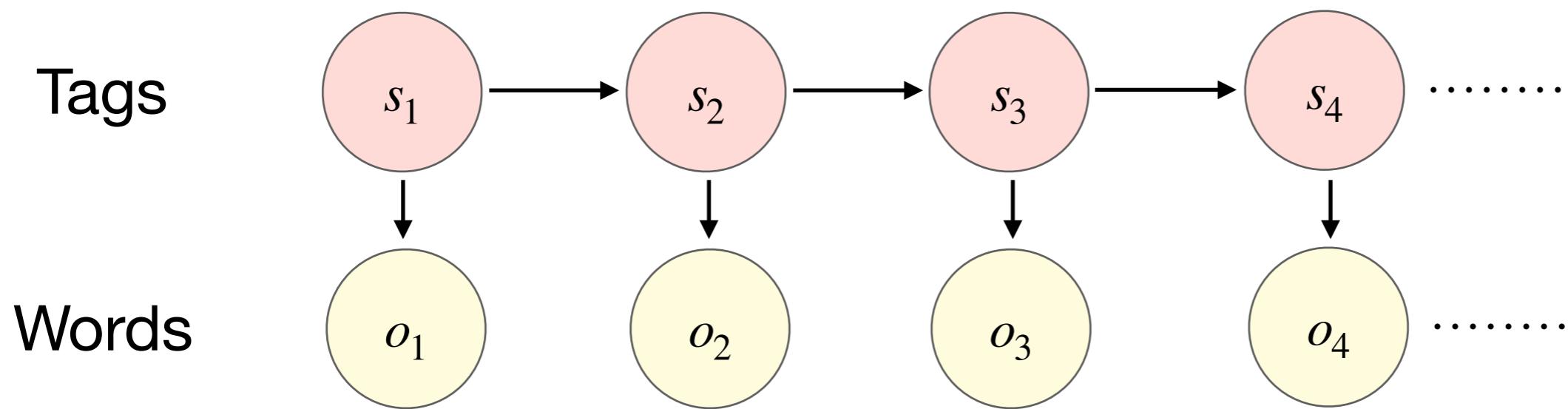
$$P(s, o) = P(s_1, s_2, \dots, s_n, o_1, o_2, \dots, o_n)$$

# Sequence likelihood



$$\begin{aligned} P(s, o) &= P(s_1, s_2 \dots s_n, o_1, o_2 \dots o_n) \\ &= \pi(s_1) P(o_1 | s_1) \overbrace{\prod_{i=2}^n P(s_i, o_i | s_{i-1})} \end{aligned}$$

# Sequence likelihood



$$\begin{aligned} P(s, o) &= P(s_1, s_2 \dots s_n, o_1, o_2 \dots o_n) \\ &= \pi(s_1) P(o_1 | s_1) \overbrace{\prod_{i=2}^n P(s_i, o_i | s_{i-1})}^{\text{Product of hidden states and observed words}} \\ &= \pi(s_1) P(o_1 | s_1) \overbrace{\prod_{i=2}^n P(s_i | s_{i-1}) P(o_i | s_i)}^{\text{Product of hidden state transitions and observed words}} \end{aligned}$$

# Learning

## Training set:

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

# Learning

## Training set:

- 1 Pierre/NNP Vinken/NNP ,/, 61/CD year  
join/VB the/DT board/NN as/IN a/DT no  
Nov./NNP 29/CD ./.
  - 2 Mr./NNP Vinken/NNP is/VBZ chairman  
N.V./NNP ,/, the/DT Dutch/NNP publish
  - 3 Rudolph/NNP Agnew/NNP ,/, 55/CD ye  
chairman/NN of/IN Consolidated/NNP Go  
,/, was/VBD named/VBN a/DT nonexecut  
this/DT British/JJ industrial/JJ conglomer
- ...
- 38,219 It/PRP is/VBZ also/RB pulling/VB  
of/IN Puerto/NNP Rico/NNP ,/, who/WP  
Hurricane/NNP Hugo/NNP victims/NNS ,/  
them/PRP to/TO San/NNP Francisco/NN

Maximum likelihood  
estimate:

$$P(s_i | s_j) = \frac{C(s_j, s_i)}{C(s_j)}$$

$$P(o | s) = \frac{C(s, o)}{C(s)}$$

# Example: POS tagging

the/?? cat/?? sat/?? on/?? the/?? mat/??

$$\pi(DT) = 0.8$$

$s_{t+1}$

$o_t$

	DT	NN	IN	VBD
DT	0.5	0.8	0.05	0.1
NN	0.05	0.2	0.15	0.6
IN	0.5	0.2	0.05	0.25
VBD	0.3	0.3	0.3	0.1

	the	cat	sat	on	mat
DT	0.5	0	0	0	0
NN	0.01	0.2	0.01	0.01	0.2
IN	0	0	0	0.4	0
VBD	0	0.01	0.1	0.01	0.01

# Example: POS tagging

the/?? cat/?? sat/?? on/?? the/?? mat/??

$$\pi(DT) = 0.8$$

$s_{t+1}$

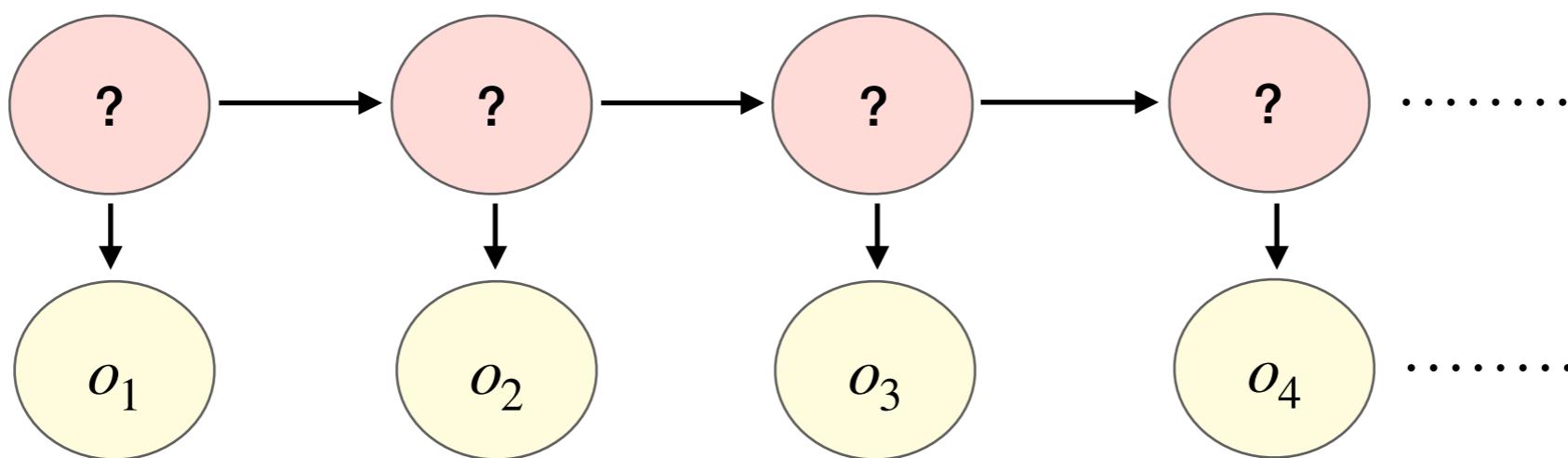
$o_t$

	DT	NN	IN	VBD
DT	0.5	0.8	0.05	0.1
NN	0.05	0.2	0.15	0.6
IN	0.5	0.2	0.05	0.25
VBD	0.3	0.3	0.3	0.1

	the	cat	sat	on	mat
DT	0.5	0	0	0	0
NN	0.01	0.2	0.01	0.01	0.2
IN	0	0	0	0.4	0
VBD	0	0.01	0.1	0.01	0.01

$$P(\text{the/DT, cat/NN, sat/VBD, on/IN, the/DT, mat/NN}) \\ = 1.84 * 10^{-5}$$

# Decoding with HMMs

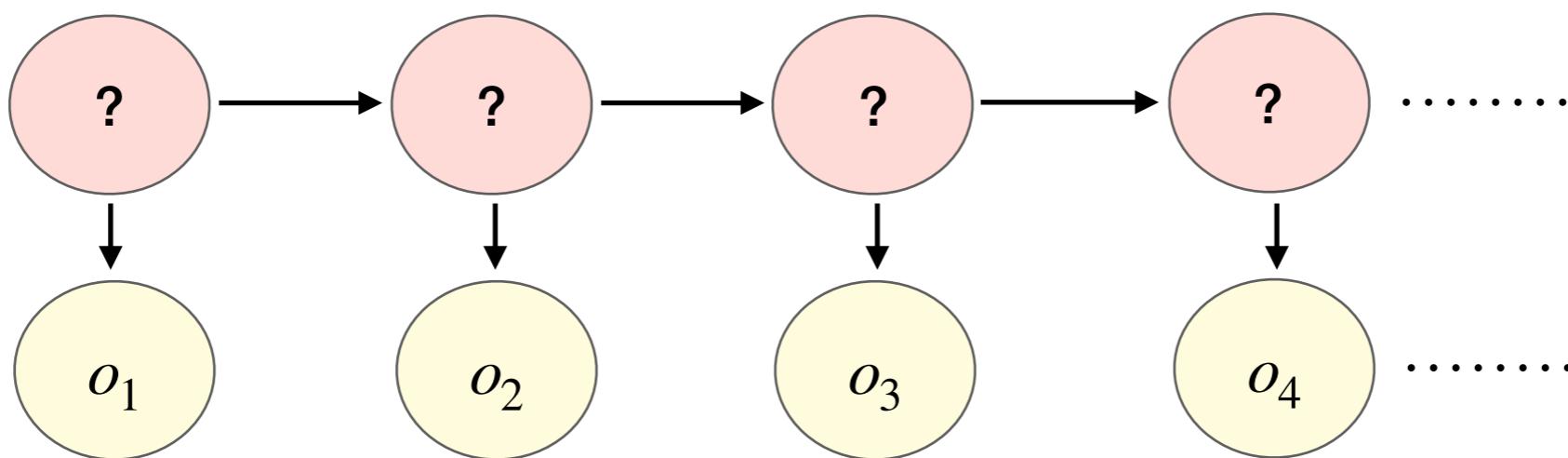


- **Task:** Find the most probable sequence of states  $\langle s_1, s_2, \dots, s_n \rangle$  given the observations  $\langle o_1, o_2, \dots, o_n \rangle$

$$\hat{s} = \underset{s}{\operatorname{argmax}} P(s|o) = \underset{s}{\operatorname{argmax}} \frac{P(s) P(o|s)}{P(o)}$$

[Bayes]

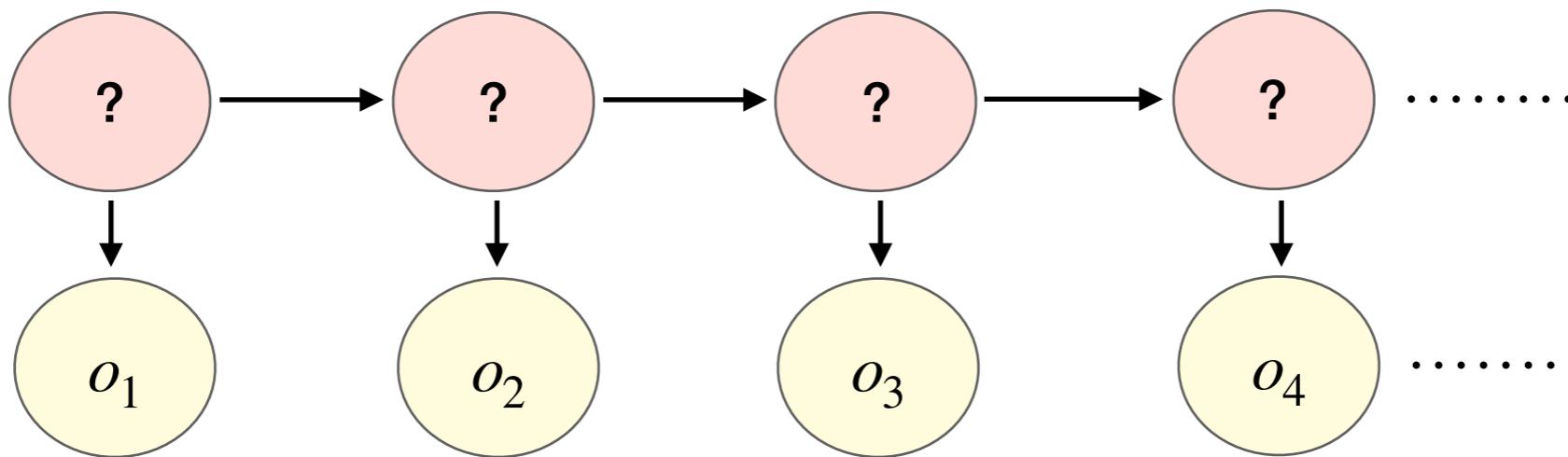
# Decoding with HMMs



- **Task:** Find the most probable sequence of states  $\langle s_1, s_2, \dots, s_n \rangle$  given the observations  $\langle o_1, o_2, \dots, o_n \rangle$

$$\hat{s} = \underset{s}{\operatorname{argmax}} P(s|o) = \underset{s}{\operatorname{argmax}} \frac{P(s) P(o|s)}{P(o)} \quad [\text{Bayes}]$$
$$= \underset{s}{\operatorname{argmax}} P(s) P(o|s)$$

# Decoding with HMMs



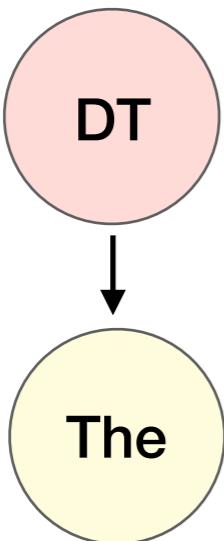
- **Task:** Find the most probable sequence of states  $\langle s_1, s_2, \dots, s_n \rangle$  given the observations  $\langle o_1, o_2, \dots, o_n \rangle$

$$\hat{s} = \arg \max_s P(s) P(o|s)$$

$$= \arg \max_s \prod_{i=1}^n P(s_i | s_{i-1}) \underbrace{P(o_i | s_i)}_{\text{Emission}}$$

$\underbrace{\hspace{10em}}$  Transition

# Greedy decoding



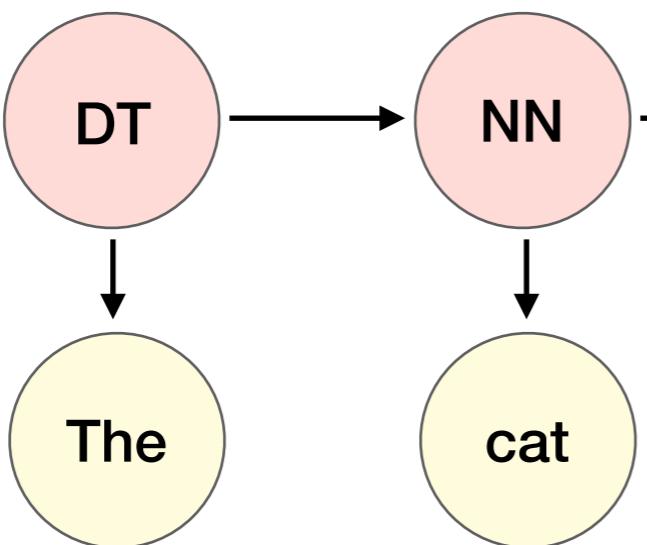
$$\underset{S}{\operatorname{argmax}} \pi(s_i = s) P(\text{The} | s) \\ = \text{'DT'}$$

$$\hat{s} = \underset{S}{\operatorname{argmax}} p(s) P(o|s)$$

$$= \underset{S}{\operatorname{argmax}} \prod_{i=1}^n P(s_i | s_{i-1}) \underbrace{P(o_i | s_i)}_{\text{Emission}}$$

Transition

# Greedy decoding



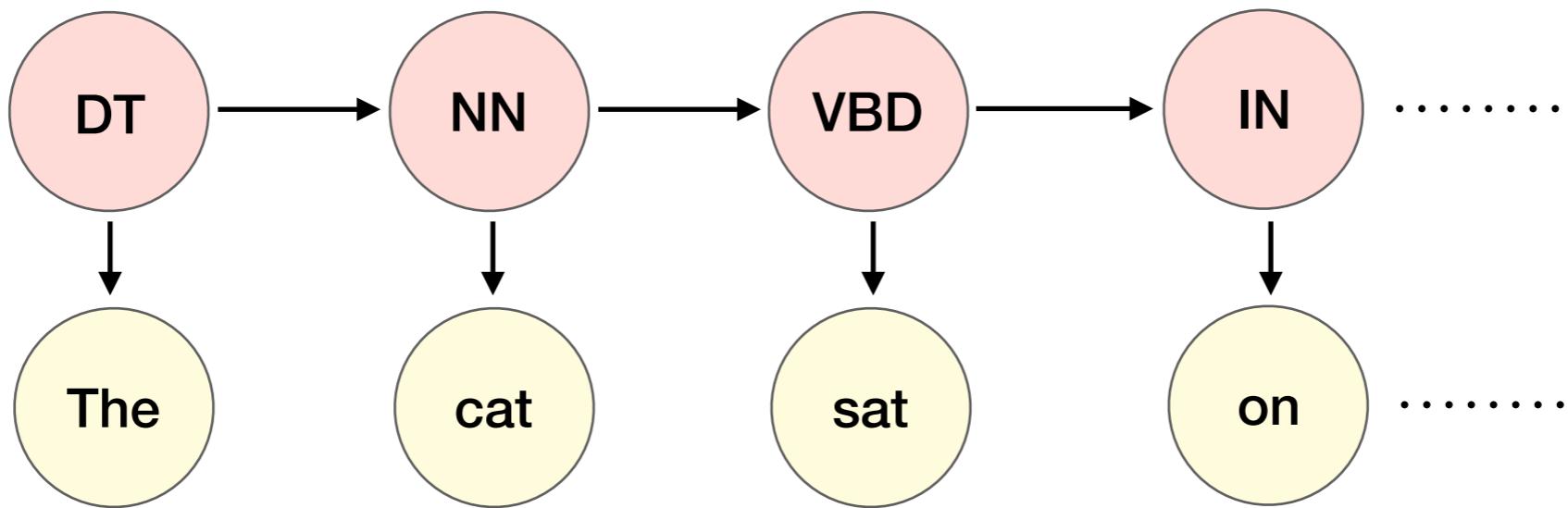
$$\underset{s}{\operatorname{argmax}} \ P(s_2=s \mid DT) P(cat \mid s) \\ = 'NN'$$

$$\hat{s} = \underset{s}{\operatorname{argmax}} \ P(s) P(o \mid s)$$

$$= \underset{s}{\operatorname{argmax}} \ \prod_{i=1}^n P(s_i \mid s_{i-1}) \underbrace{P(o_i \mid s_i)}_{\text{Emission}}$$

Transition

# Greedy decoding



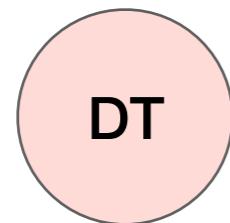
$$\forall t, \hat{s}_{t+1} = \operatorname{argmax}_s P(s | \hat{s}_t) P(o_{t+1} | s)$$

- Not guaranteed to be optimal!
- Local decisions

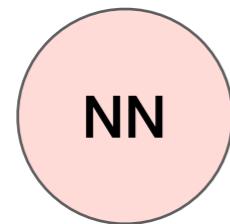
# Viterbi decoding

- Use dynamic programming!
- Probability lattice,  $M[T, K]$ 
  - $T$  : Number of time steps
  - $K$  : Number of states
- $M[i, j]$  : Most probable sequence of states ending with state **j** at time **i**

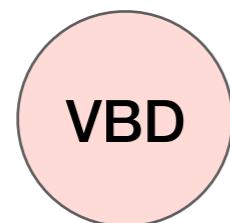
# Viterbi decoding



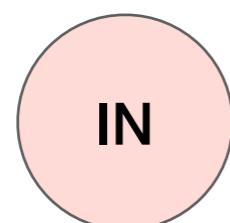
$$M[1,DT] = \pi(DT) P(\mathbf{the} | DT)$$



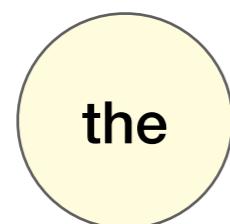
$$M[1,NN] = \pi(NN) P(\mathbf{the} | NN)$$



$$M[1,VBD] = \pi(VBD) P(\mathbf{the} | VBD)$$

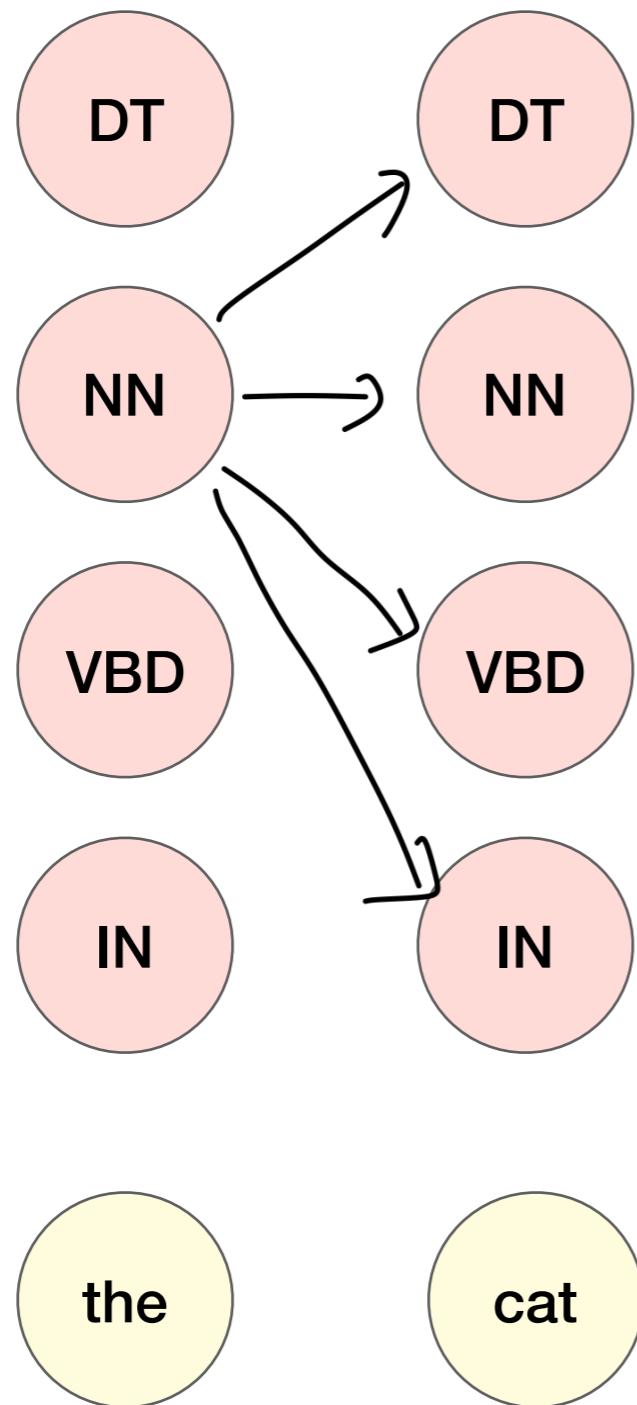


$$M[1,IN] = \pi(IN) P(\mathbf{the} | IN)$$



*Forward*

# Viterbi decoding



$$M[2,DT] = \max_k M[1,k] P(DT|k) P(\mathbf{cat}|DT)$$

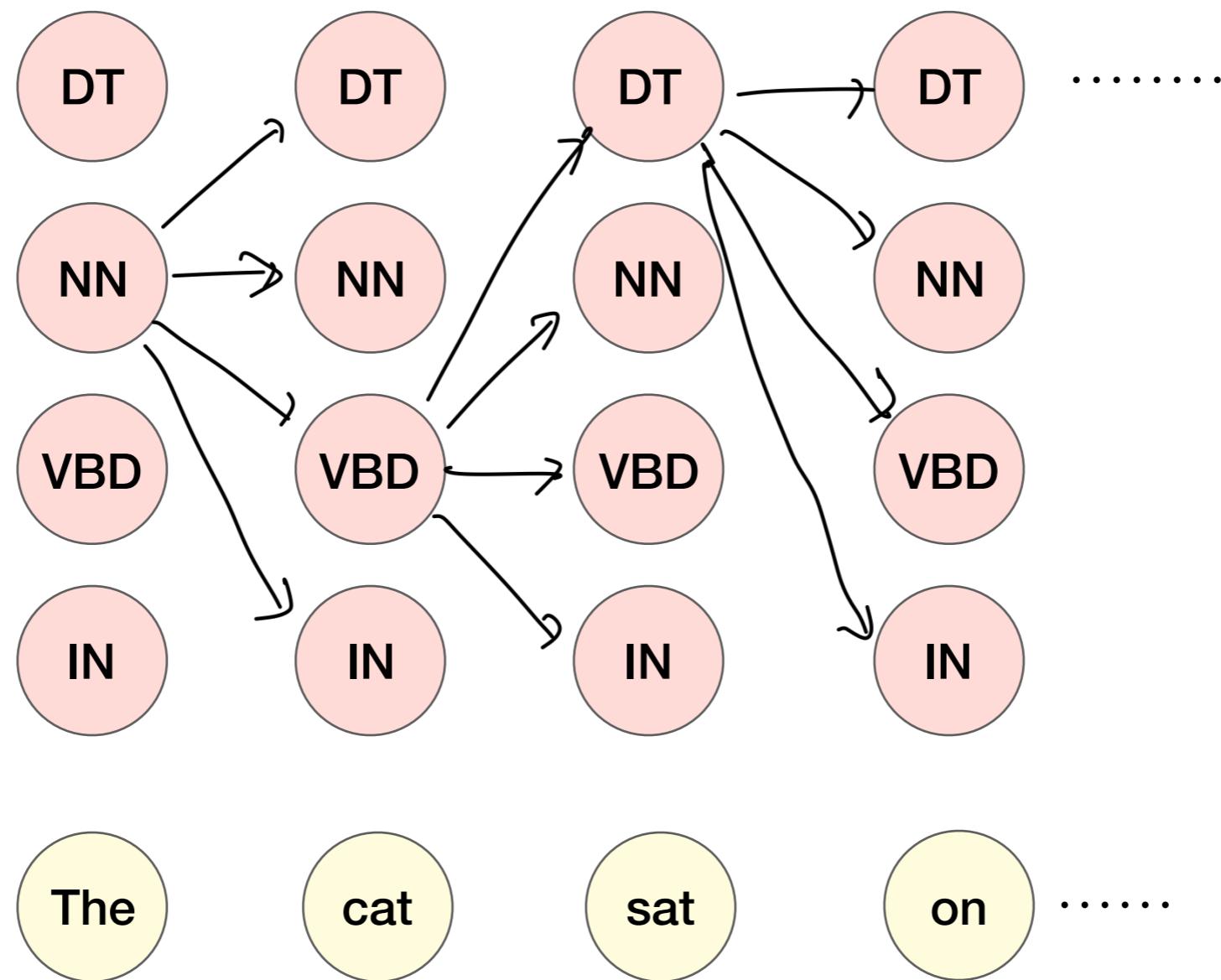
$$M[2,NN] = \max_k M[1,k] P(NN|k) P(\mathbf{cat}|NN)$$

$$M[2,VBD] = \max_k M[1,k] P(VBD|k) P(\mathbf{cat}|VBD)$$

$$M[2,IN] = \max_k M[1,k] P(IN|k) P(\mathbf{cat}|IN)$$

*Forward*

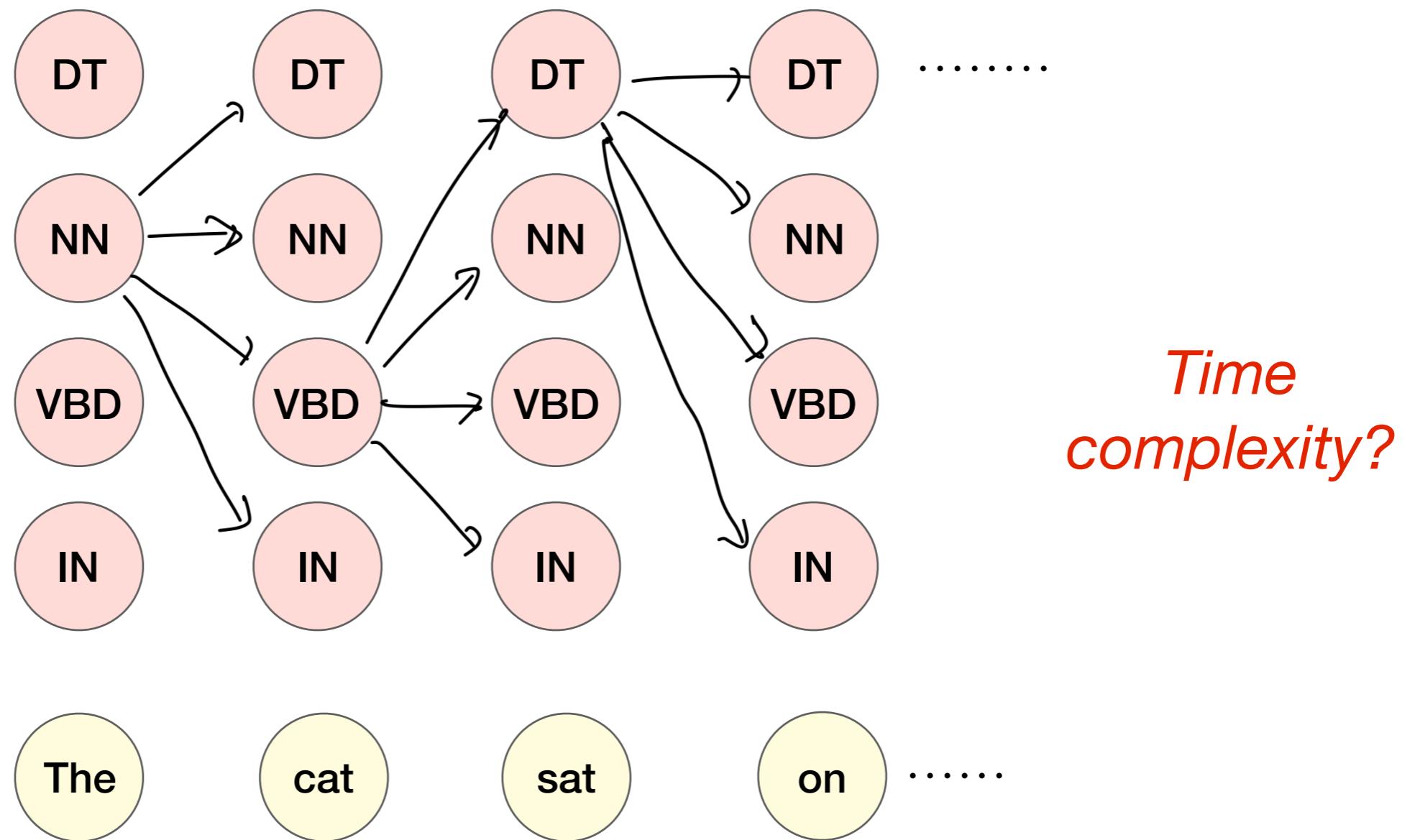
# Viterbi decoding



$$M[i, j] = \max_k M[i - 1, k] P(s_j | s_k) P(o_i | s_j) \quad 1 \leq k \leq K \quad 1 \leq i \leq n$$

*Backward:* Pick  $\max_k M[n, k]$  and backtrack

# Viterbi decoding

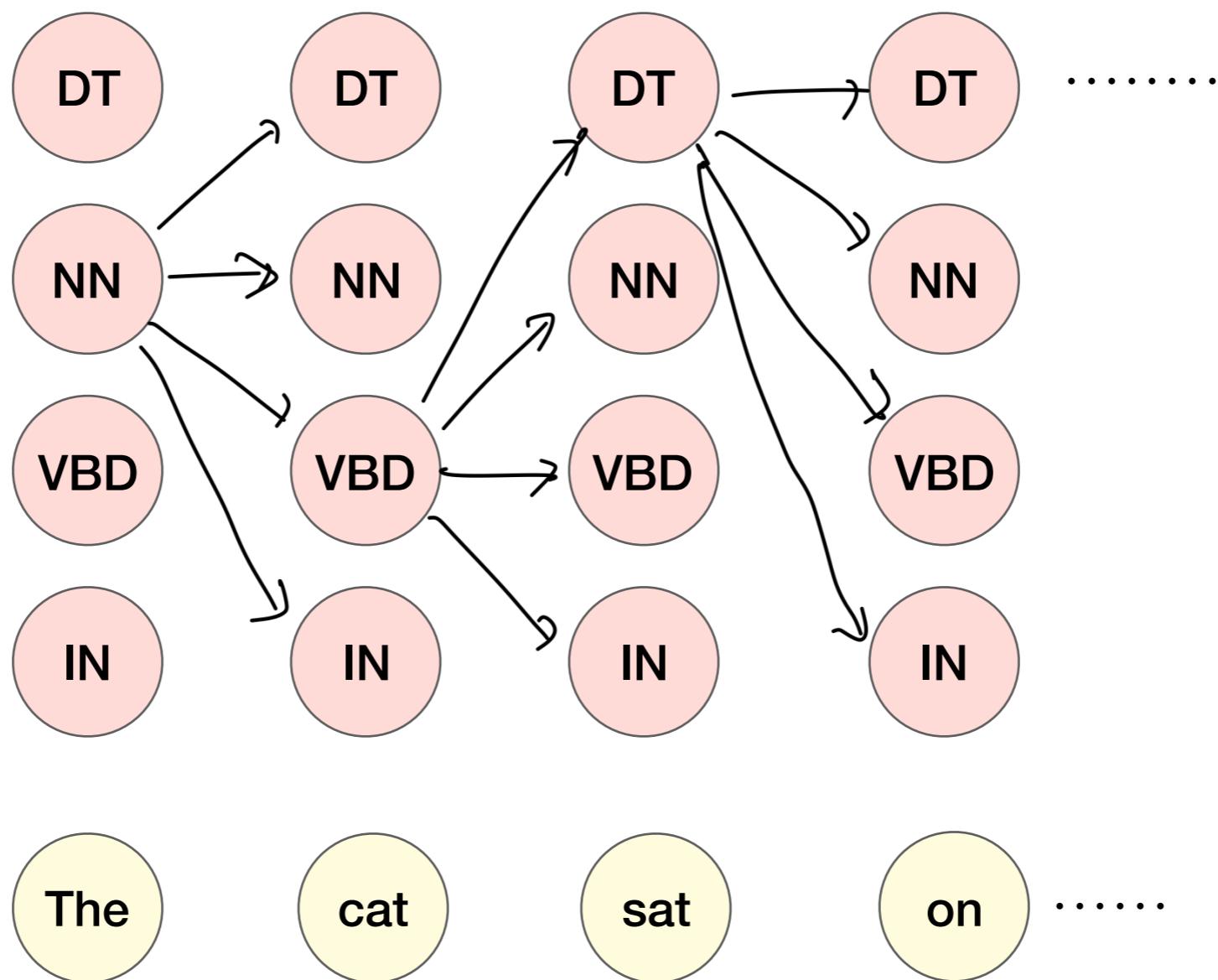


$$M[i, j] = \max_k M[i - 1, k] P(s_j | s_k) P(o_i | s_j) \quad 1 \leq k \leq K \quad 1 \leq i \leq n$$

*Backward:* Pick  $\max_k M[n, k]$  and backtrack

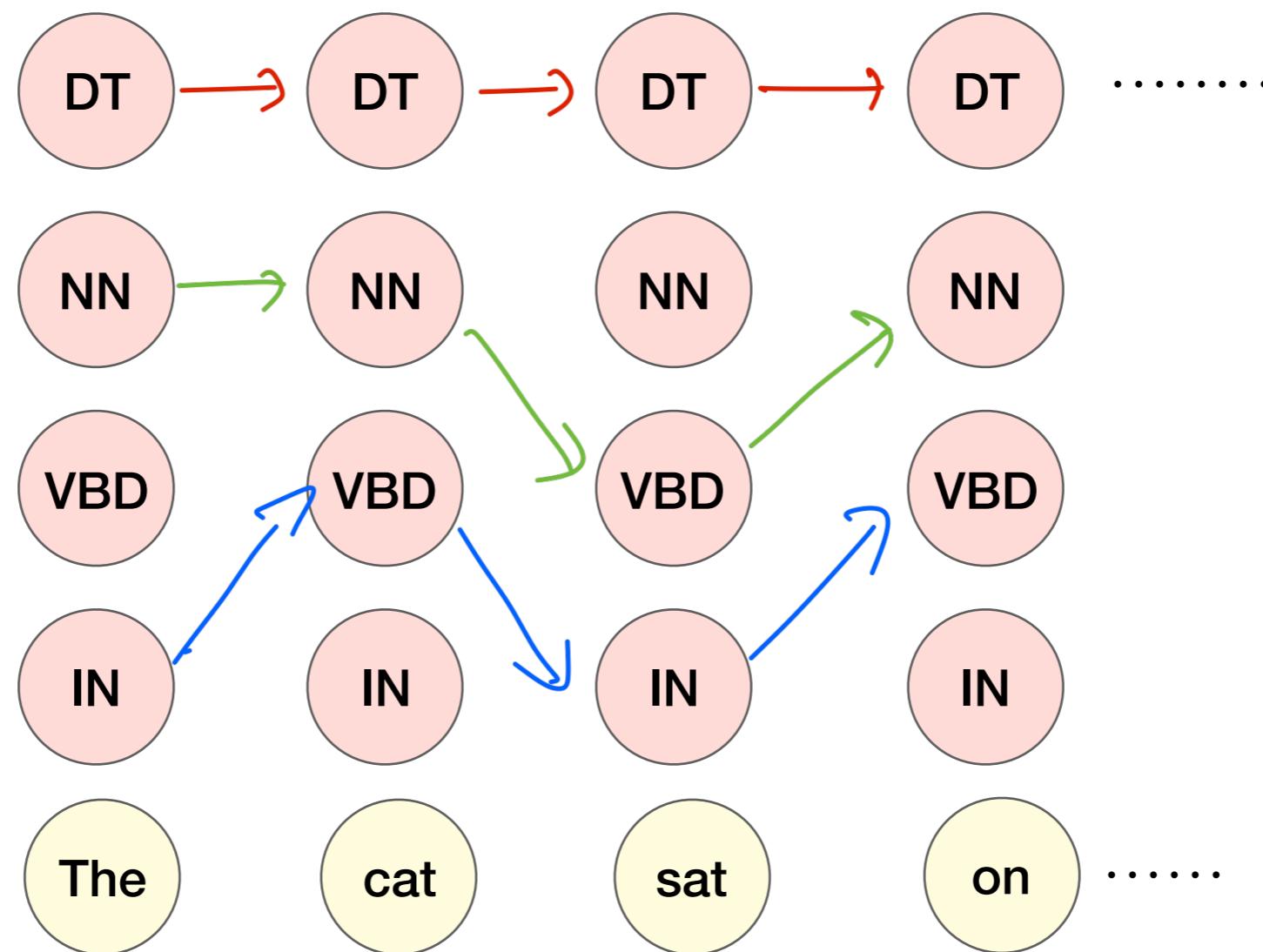
# Beam Search

- If K (number of states) is too large, Viterbi is too expensive!



# Beam Search

- If K (number of states) is too large, Viterbi is too expensive!



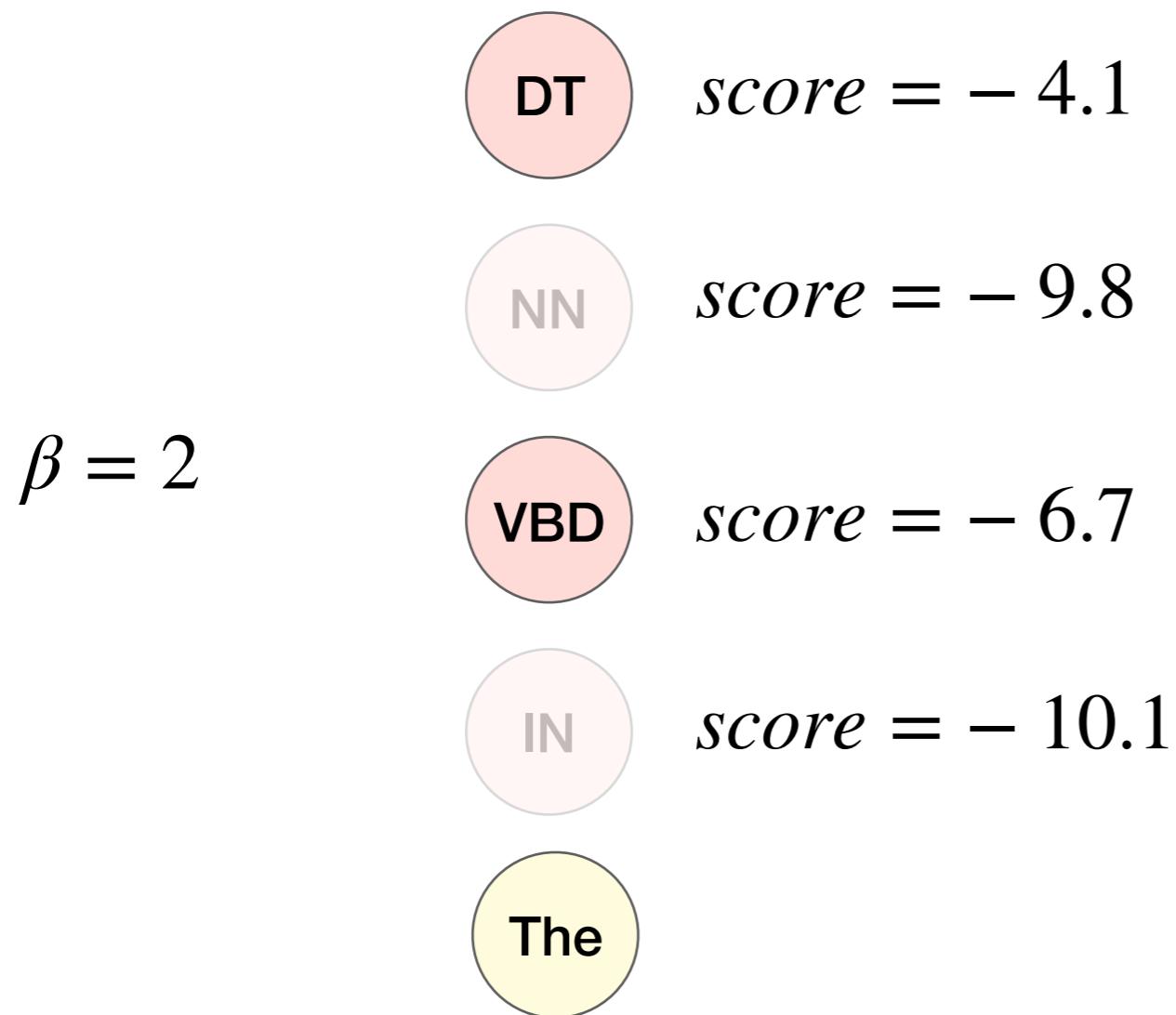
*Many paths have very low likelihood!*

# Beam Search

- If  $K$  (number of states) is too large, Viterbi is too expensive!
- Keep a fixed number of hypotheses at each point
  - Beam width,  $\beta$

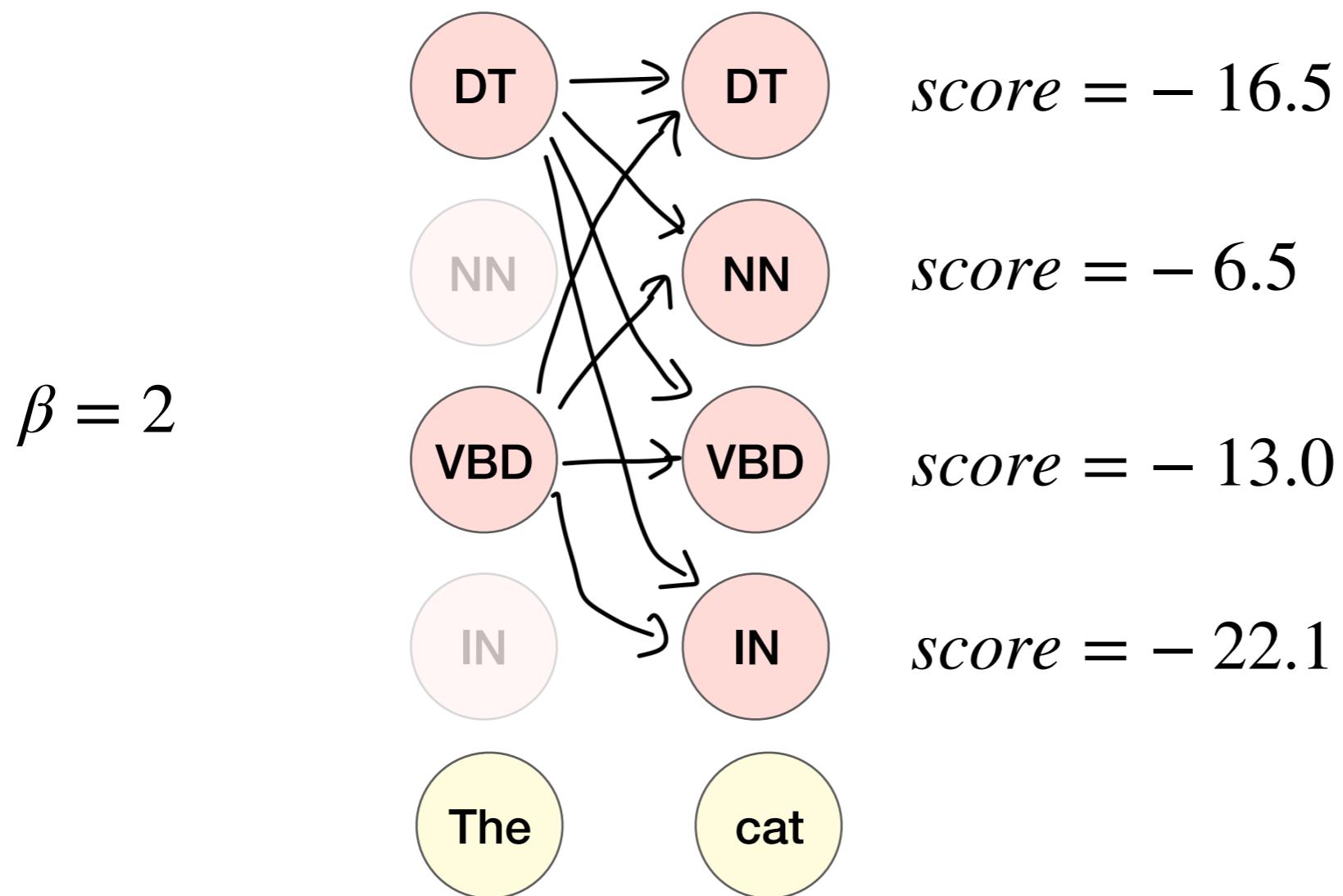
# Beam Search

- Keep a fixed number of hypotheses at each point



# Beam Search

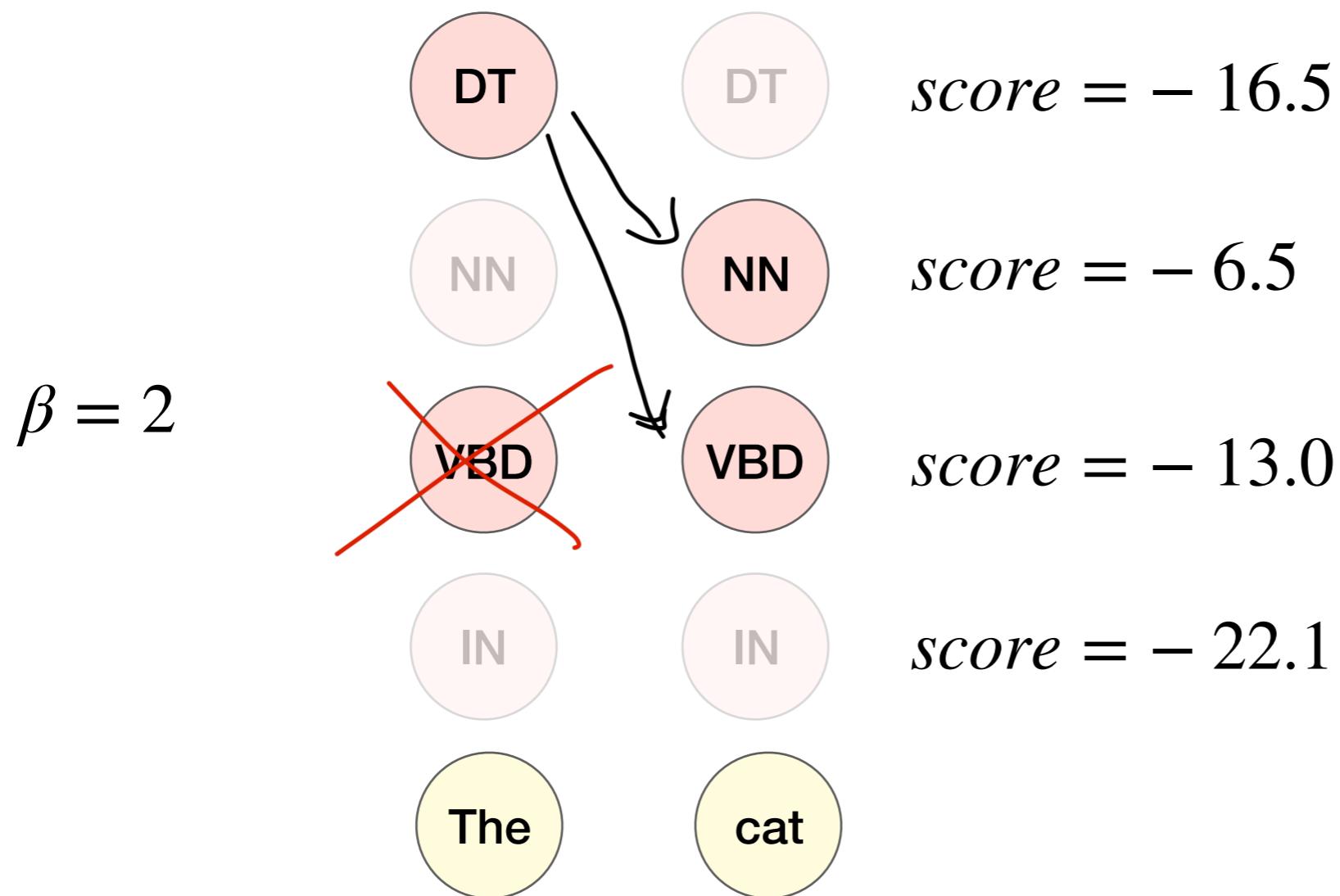
- Keep a fixed number of hypotheses at each point



Step 1: Expand all partial sequences in current beam

# Beam Search

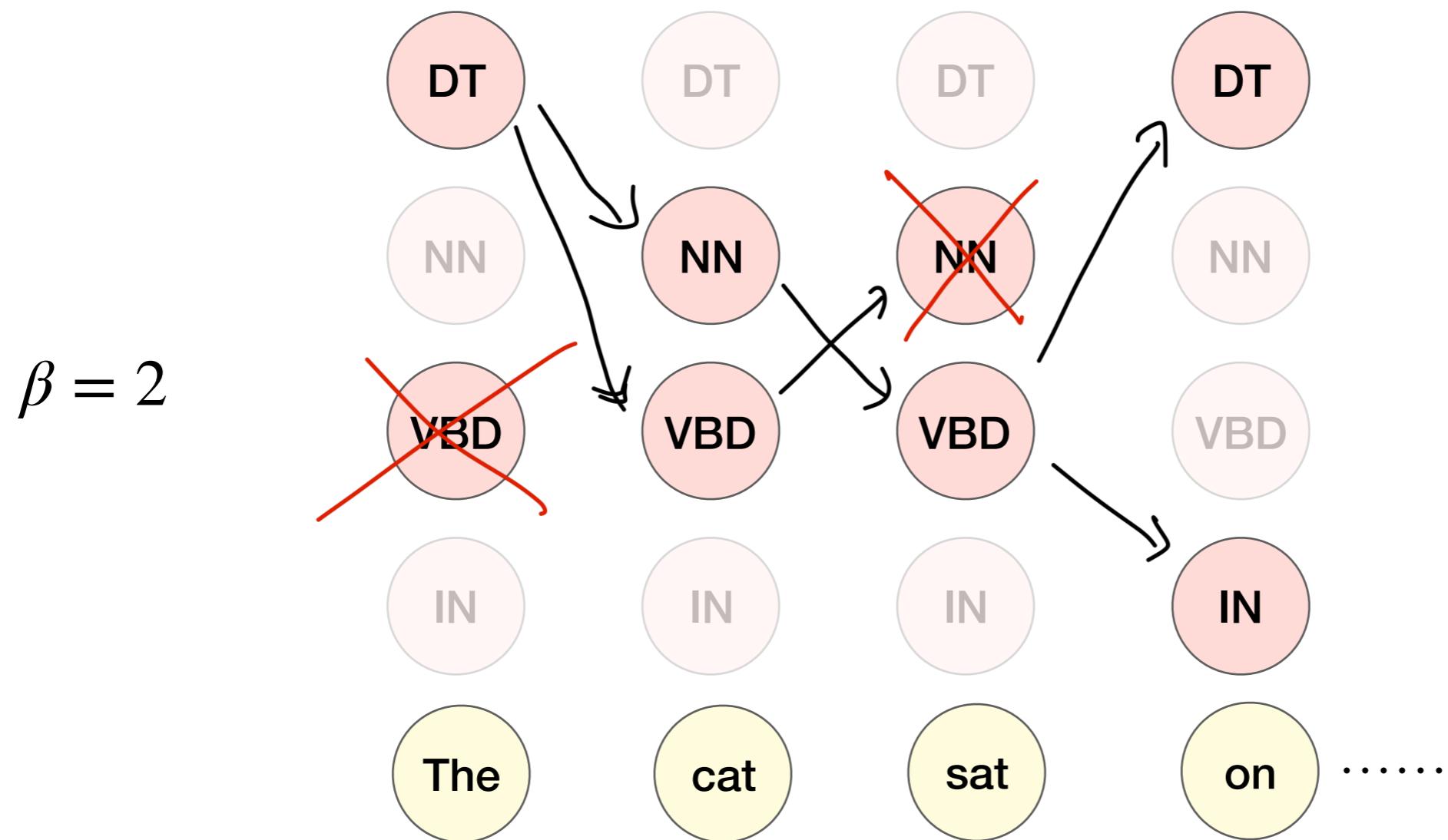
- Keep a fixed number of hypotheses at each point



**Step 2:** Prune set back to top  $\beta$  sequences

# Beam Search

- Keep a fixed number of hypotheses at each point



Pick  $\max_k M[n, k]$  from within beam and backtrack

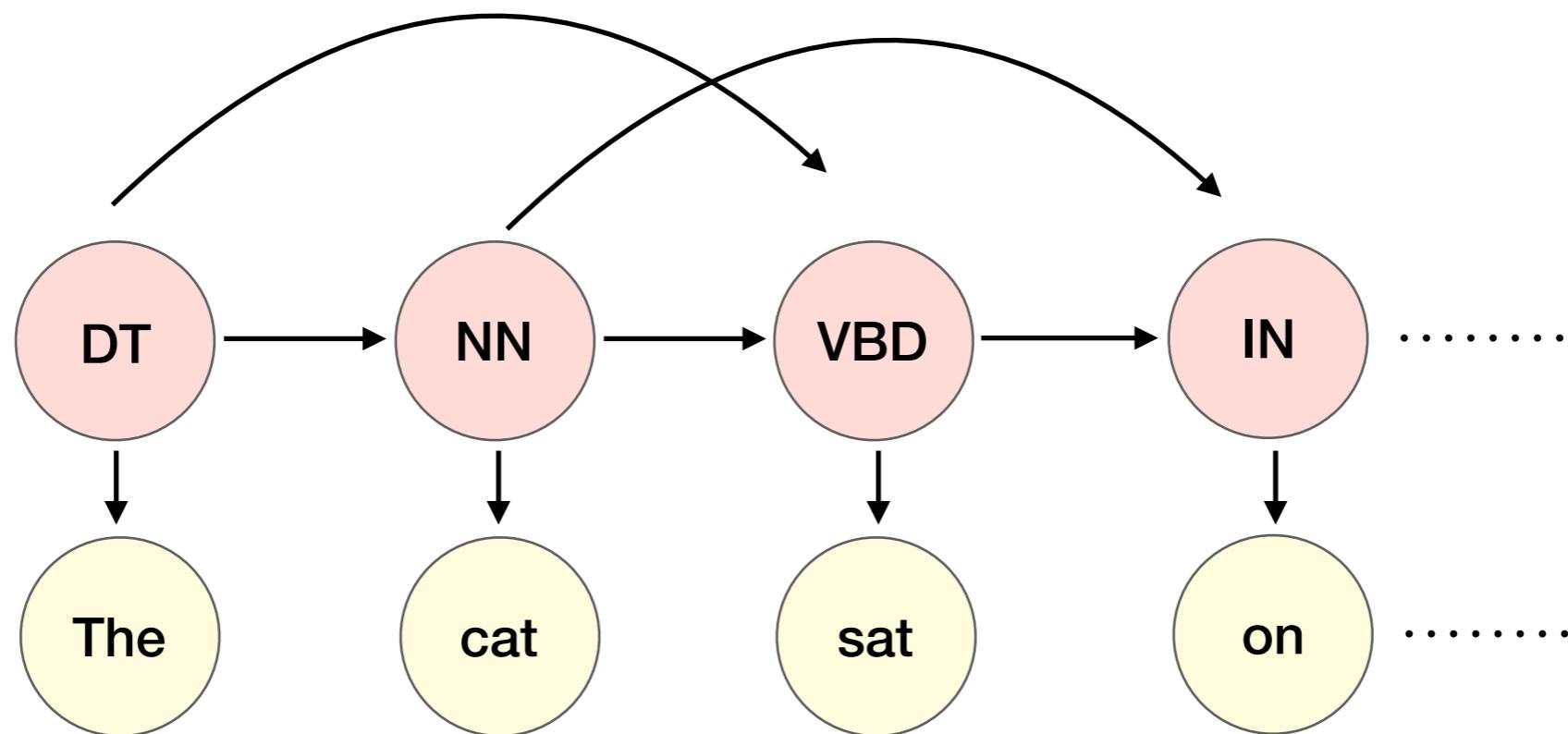
# Beam Search

- If  $K$  (number of states) is too large, Viterbi is too expensive!
- Keep a fixed number of hypotheses at each point
  - Beam width,  $\beta$
- Trade-off computation for (some) accuracy

*Time complexity?*

# Beyond bigrams

- Real-world HMM taggers have more relaxed assumptions
- Trigram HMM:  $P(s_{t+1} | s_1, s_2, \dots, s_t) \approx P(s_{t+1} | s_{t-1}, s_t)$



Pros?

Cons?

# Maximum Entropy Markov Models

# Generative vs Discriminative

- HMM is a *generative* model
- Can we model  $P(s_1, \dots, s_n | o_1, \dots, o_n)$  directly?

Generative

Naive Bayes:

$$P(c)P(d | c)$$

HMM:

$$P(s_1, \dots, s_n)P(o_1, \dots, o_n | s_1, \dots, s_n)$$

Discriminative

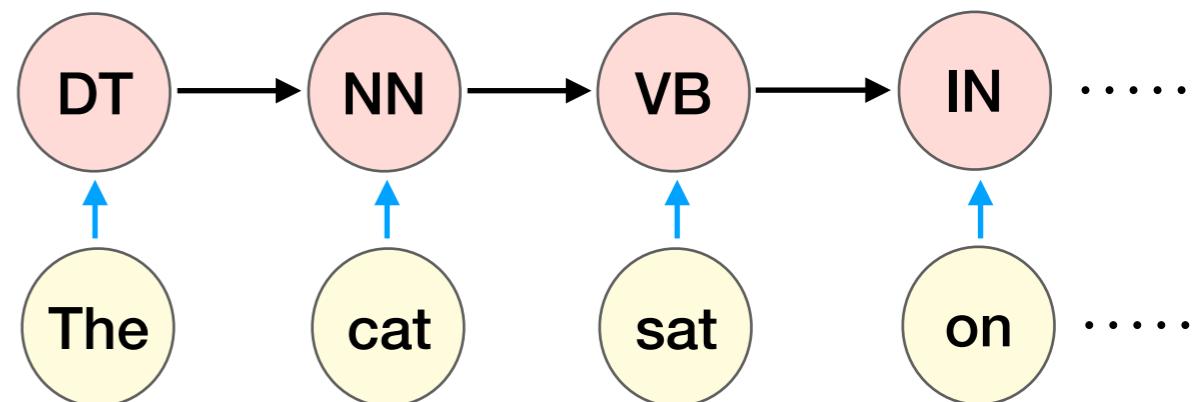
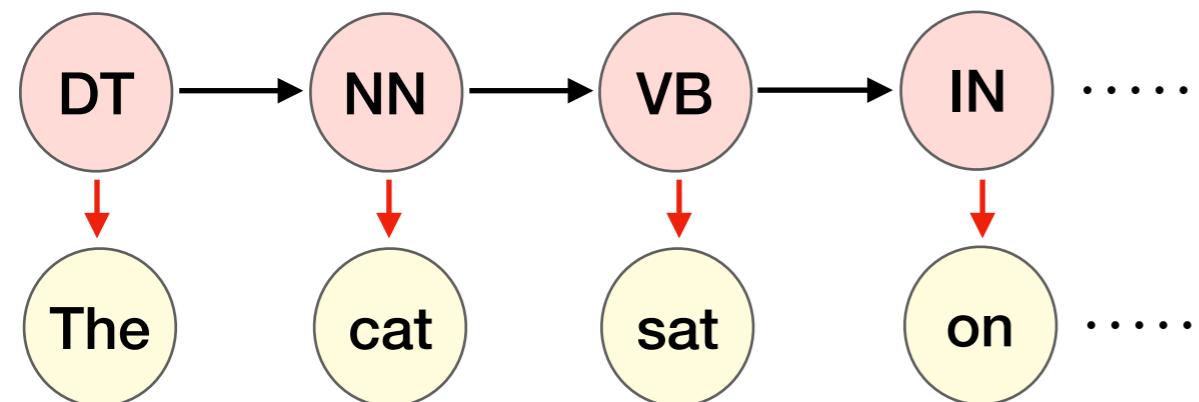
Logistic Regression:

$$P(c | d)$$

MEMM:

$$P(s_1, \dots, s_n | o_1, \dots, o_n)$$

# MEMM



- Compute the posterior directly:

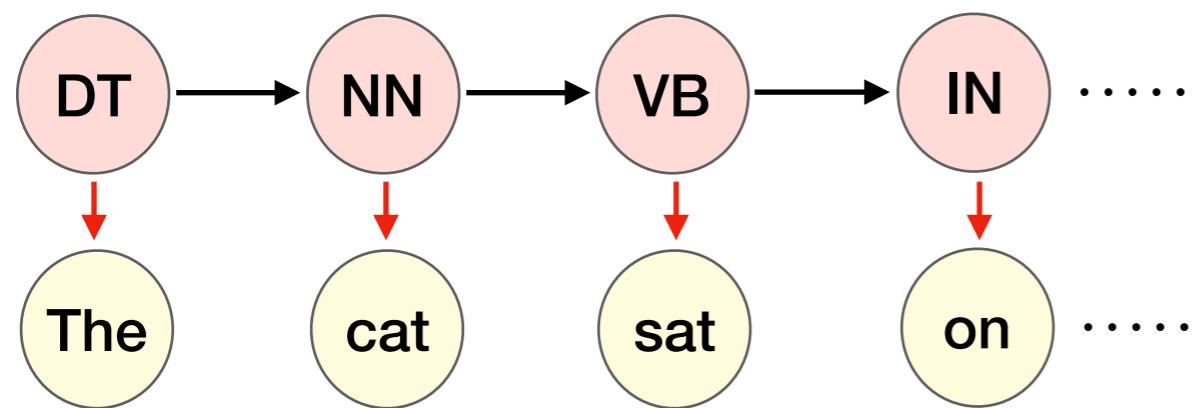
$$\hat{S} = \arg \max_S P(S | O) = \arg \max_S \prod_i P(s_i | o_i, s_{i-1})$$

*Features*

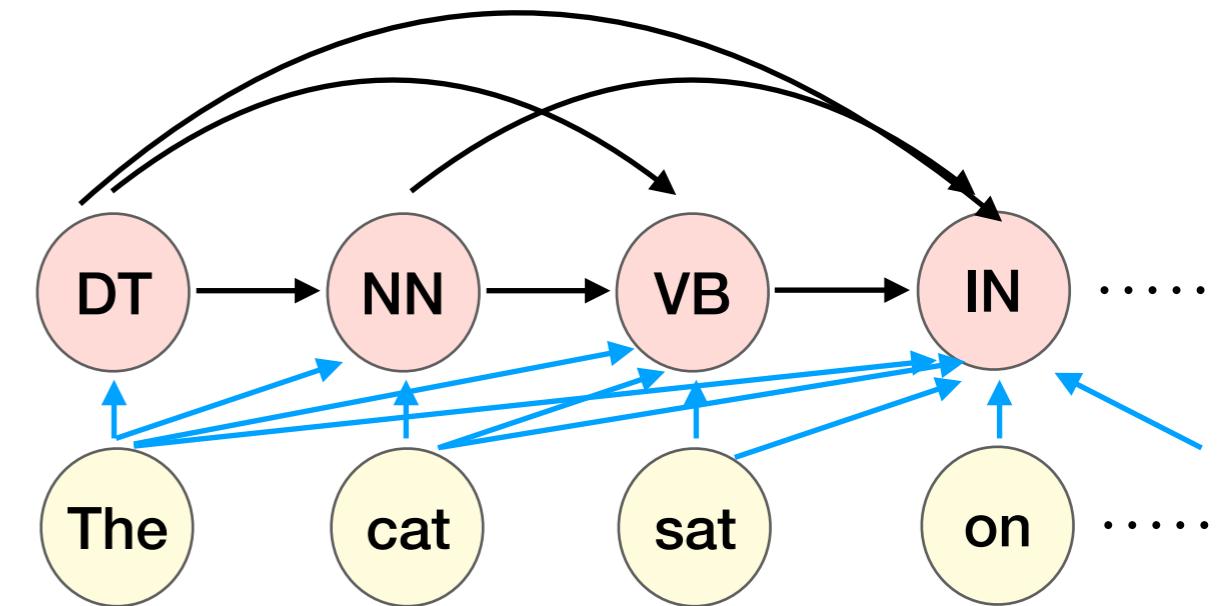
- Use features:  $P(s_i | o_i, s_{i-1}) \propto \exp(w \cdot f(s_i, o_i, s_{i-1}))$

*weights*

# MEMM



HMM



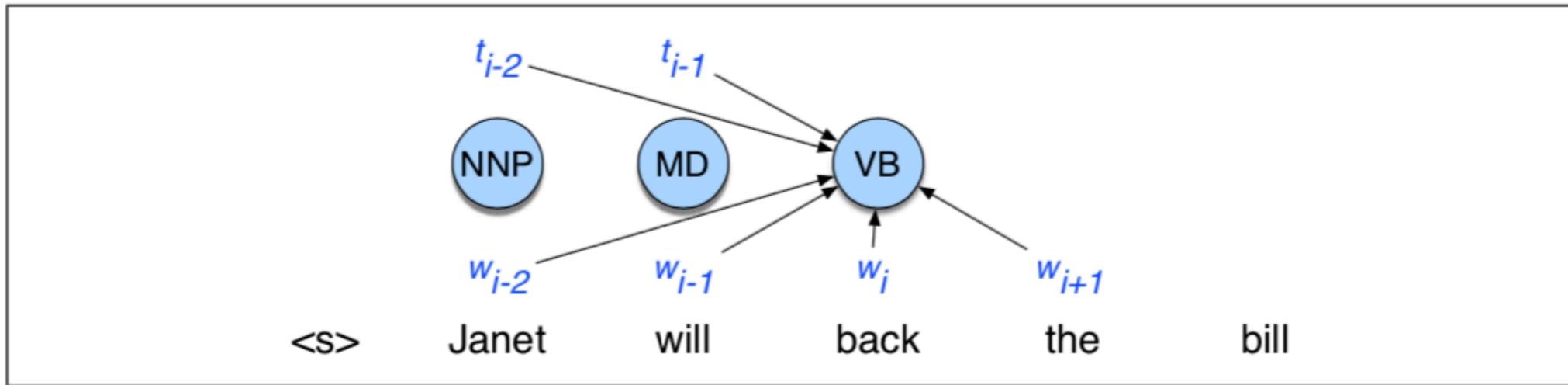
MEMM

- In general, we can use all observations and all previous states:

$$\hat{S} = \arg \max_S P(S | O) = \arg \max_S \prod_i P(s_i | o_n, o_{i-1}, \dots, o_1, s_{i-1}, \dots, s_1)$$

$$P(s_i | s_{i-1}, \dots, s_1, O) \propto \exp(w \cdot f(s_i, s_{i-1}, \dots, s_1, O))$$

# Features in an MEMM



**Figure 8.13** An MEMM for part-of-speech tagging showing the ability to condition on more features.

$\langle t_i, w_{i-2} \rangle, \langle t_i, w_{i-1} \rangle, \langle t_i, w_i \rangle, \langle t_i, w_{i+1} \rangle, \langle t_i, w_{i+2} \rangle$   
 $\langle t_i, t_{i-1} \rangle, \langle t_i, t_{i-2}, t_{i-1} \rangle,$   
 $\langle t_i, t_{i-1}, w_i \rangle, \langle t_i, w_{i-1}, w_i \rangle \langle t_i, w_i, w_{i+1} \rangle,$

Feature templates

$t_i = \text{VB and } w_{i-2} = \text{Janet}$   
 $t_i = \text{VB and } w_{i-1} = \text{will}$   
 $t_i = \text{VB and } w_i = \text{back}$   
 $t_i = \text{VB and } w_{i+1} = \text{the}$   
 $t_i = \text{VB and } w_{i+2} = \text{bill}$   
 $t_i = \text{VB and } t_{i-1} = \text{MD}$   
 $t_i = \text{VB and } t_{i-1} = \text{MD and } t_{i-2} = \text{NNP}$   
 $t_i = \text{VB and } w_i = \text{back and } w_{i+1} = \text{the}$

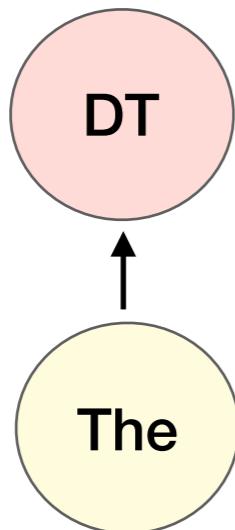
Features

# MEMMs: Decoding

$$\hat{S} = \arg \max_S P(S | O) = \arg \max_S \prod_i P(s_i | o_i, s_{i-1})$$

(assume features only on previous time step and current obs)

- Greedy decoding:

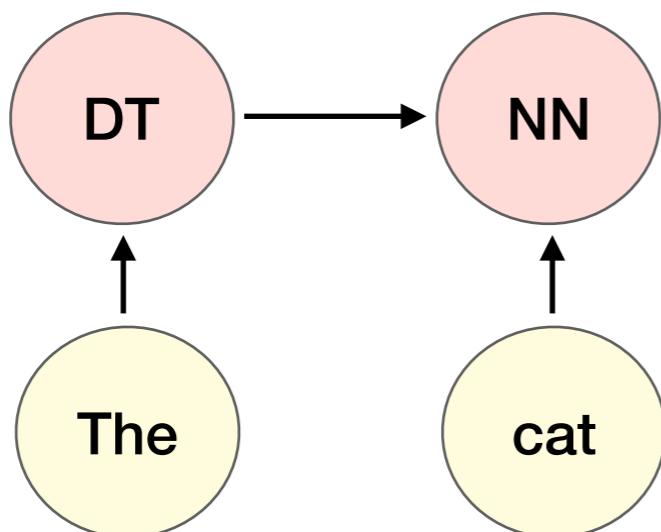


$$\begin{aligned}\hat{s}_i &= \arg \max_S P(s_i | \text{The}) \\ &= \text{DT}\end{aligned}$$

# MEMMs: Decoding

$$\hat{S} = \arg \max_S P(S | O) = \arg \max_S \prod_i P(s_i | o_i, s_{i-1})$$

- Greedy decoding:

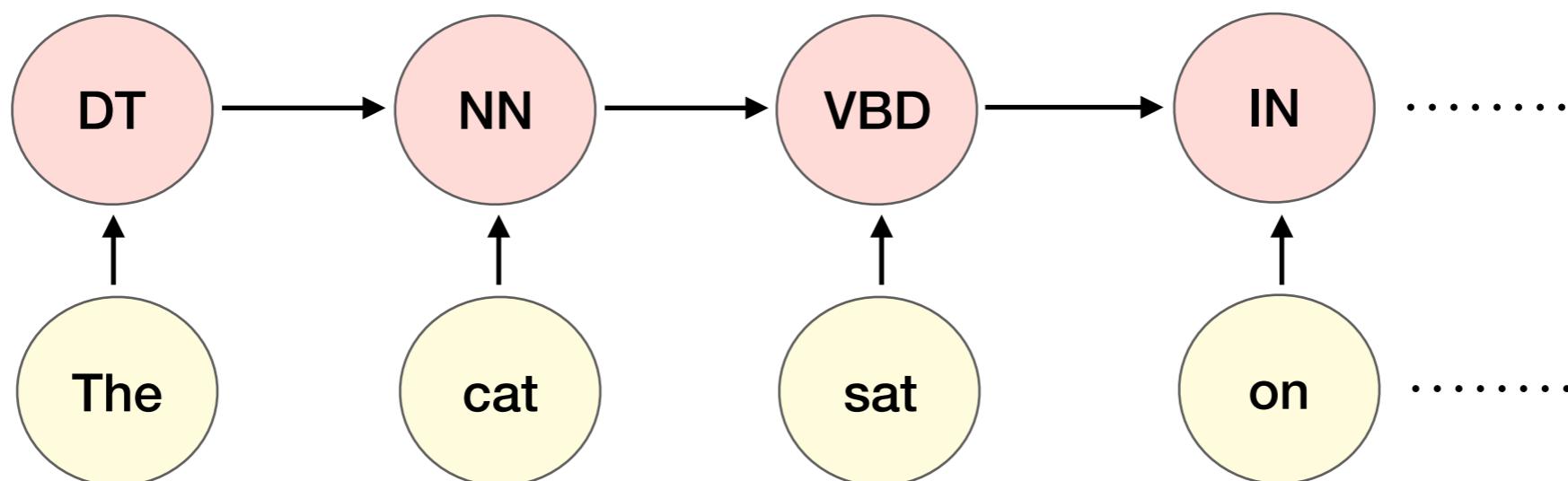


$$\begin{aligned}\hat{s}_2 &= \underset{S}{\operatorname{argmax}} P(s | \text{cat}, \text{DT}) \\ &= \text{NN}\end{aligned}$$

# MEMMs: Decoding

$$\hat{S} = \arg \max_S P(S | O) = \arg \max_S \prod_i P(s_i | o_i, s_{i-1})$$

- Greedy decoding:



$$\forall t, \hat{s}_{t+1} = \arg \max_S P(S | o_{t+1}, \hat{s}_t)$$

# MEMMs: Decoding

$$\hat{S} = \arg \max_S P(S | O) = \arg \max_S \prod_i P(s_i | o_i, s_{i-1})$$

- Greedy decoding
- Viterbi decoding:

$$M[i, j] = \max_k M[i - 1, k] P(s_j | o_i, s_k) \quad 1 \leq k \leq K \quad 1 \leq i \leq n$$

DP Lattice      # states      # timesteps

# MEMM: Learning

- Gradient descent: similar to logistic regression!

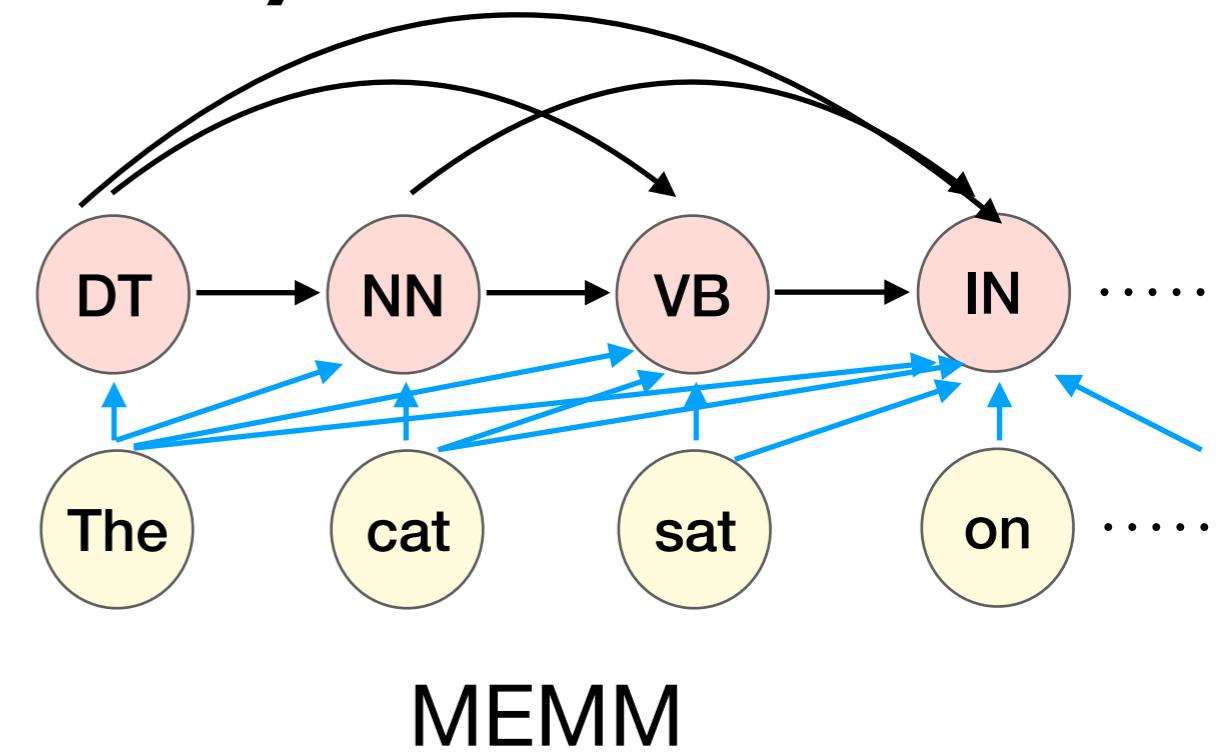
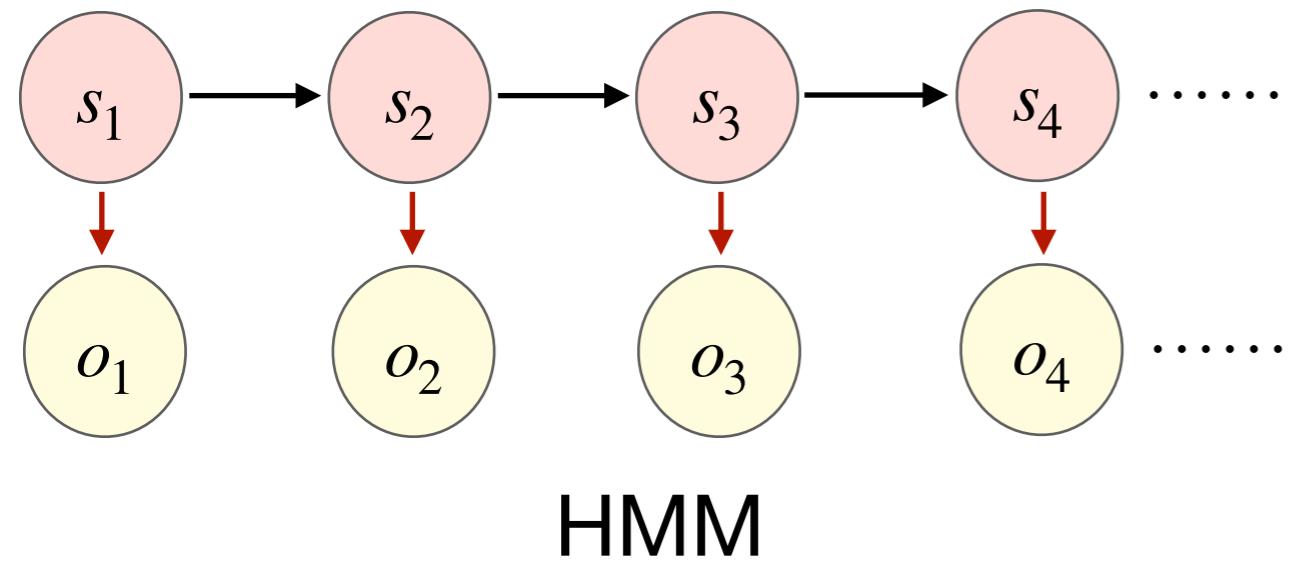
$$P(s_i | s_1, \dots, s_{i-1}, O) \propto \exp(w \cdot f(s_1, \dots, s_i, O))$$

- Given: pairs of  $(S, O)$  where each  $S = \langle s_1, s_2, \dots, s_n \rangle$

Loss for one sequence,  $L = - \sum_i \log P(s_i | s_1, \dots, s_{i-1}, O)$

- Compute gradients with respect to weights  $w$  and update

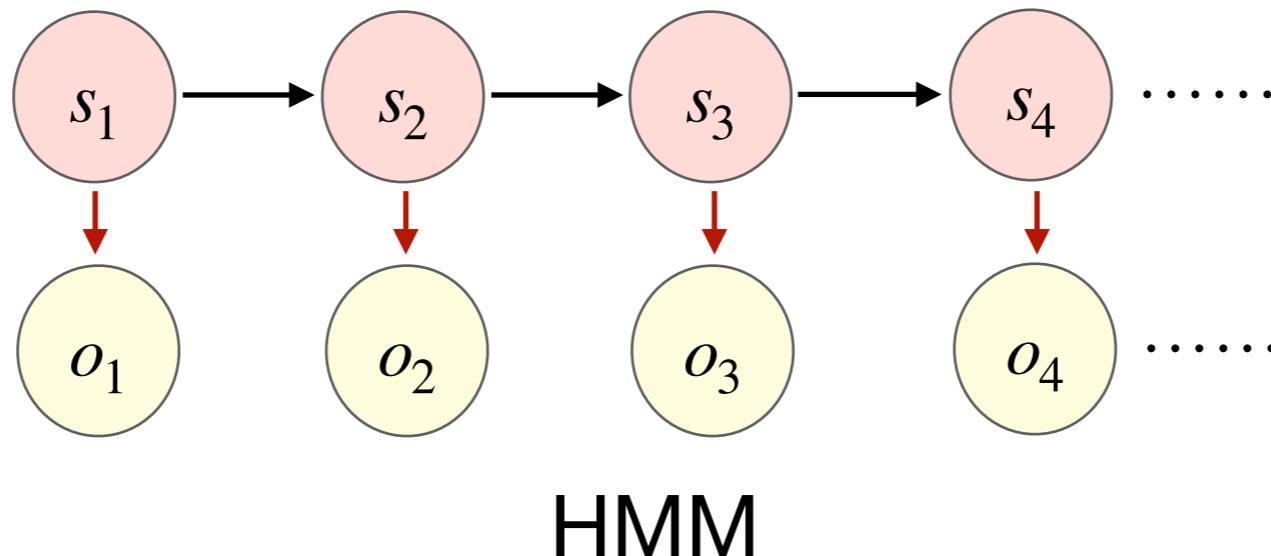
# Bidirectionality



Both HMM and MEMM assume left-to-right processing

*Why can this be undesirable?*

# Bidirectionality



The/? old/? man/? the/? boat/?

$$\begin{array}{cccccc} P(JJ | DT) & \boxed{P(\mathbf{old} | JJ)} & P(NN | JJ) & \boxed{P(\mathbf{man} | NN)} & P(DT | NN) \\ P(NN | DT) & \boxed{P(\mathbf{old} | NN)} & P(VB | NN) & \boxed{P(\mathbf{man} | VB)} & P(DT | VB) \end{array}$$

Observation bias

# Stanford Parser

Please enter a sentence to be parsed:

The old man the boat

Language: English

Sample Sentence

Parse

## Your query

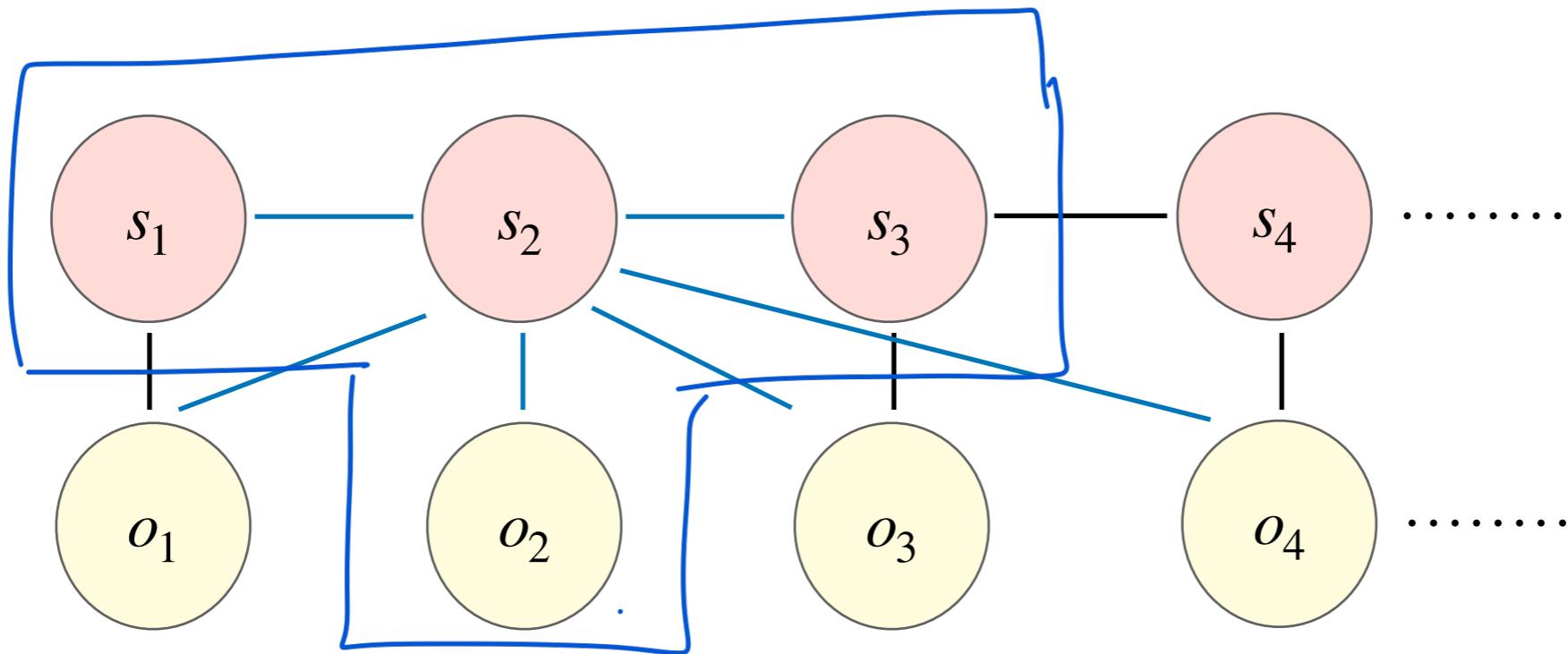
*The old man the boat*

## Tagging

The/DT old/JJ man/NN the/DT boat/NN

Observation bias

# Conditional Random Field (advanced)



- Compute log-linear functions over cliques
- Lesser independence assumptions
- Ex:  $P(s_t \mid \text{everything else}) \propto \exp(w \cdot f(s_{t-1}, s_t, s_{t+1}, O))$

