



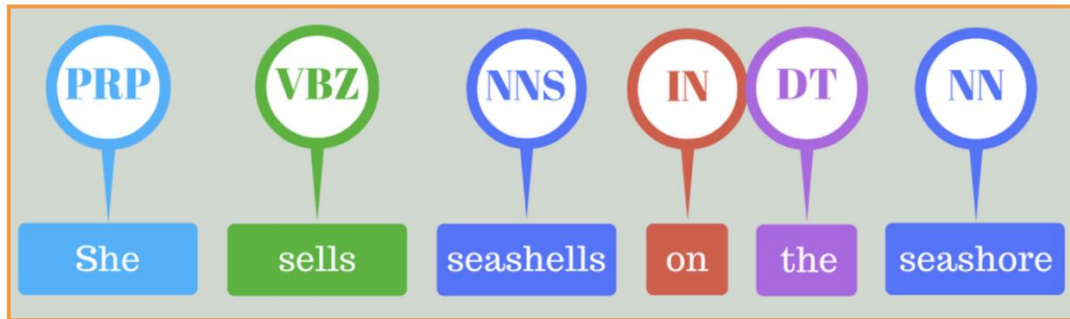
**CSE 4392 SPECIAL TOPICS**  
**NATURAL LANGUAGE PROCESSING**

# Sequence Models

1

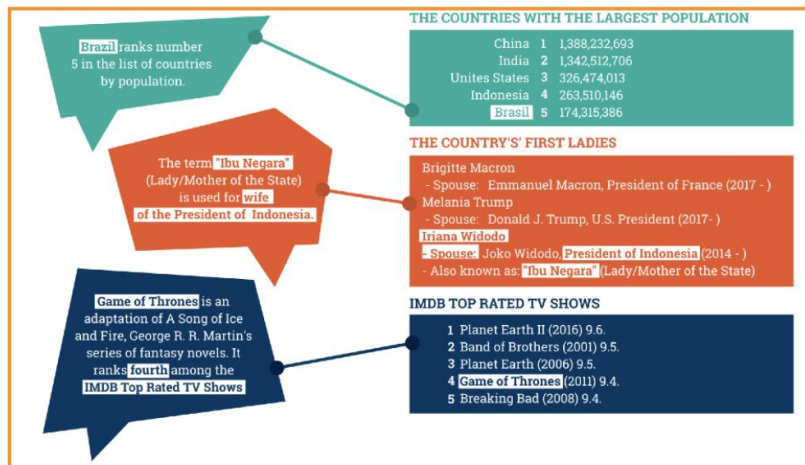
2025 Spring

# WHY MODEL SEQUENCES?



Part-of-speech tagging

Name Entity Recognition



Information extraction

# OVERVIEW

- Hidden Markov Models (HMM)
- Viterbi algorithm
- Conditional Random Field (CRF)

# WHAT ARE POS TAGS?

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

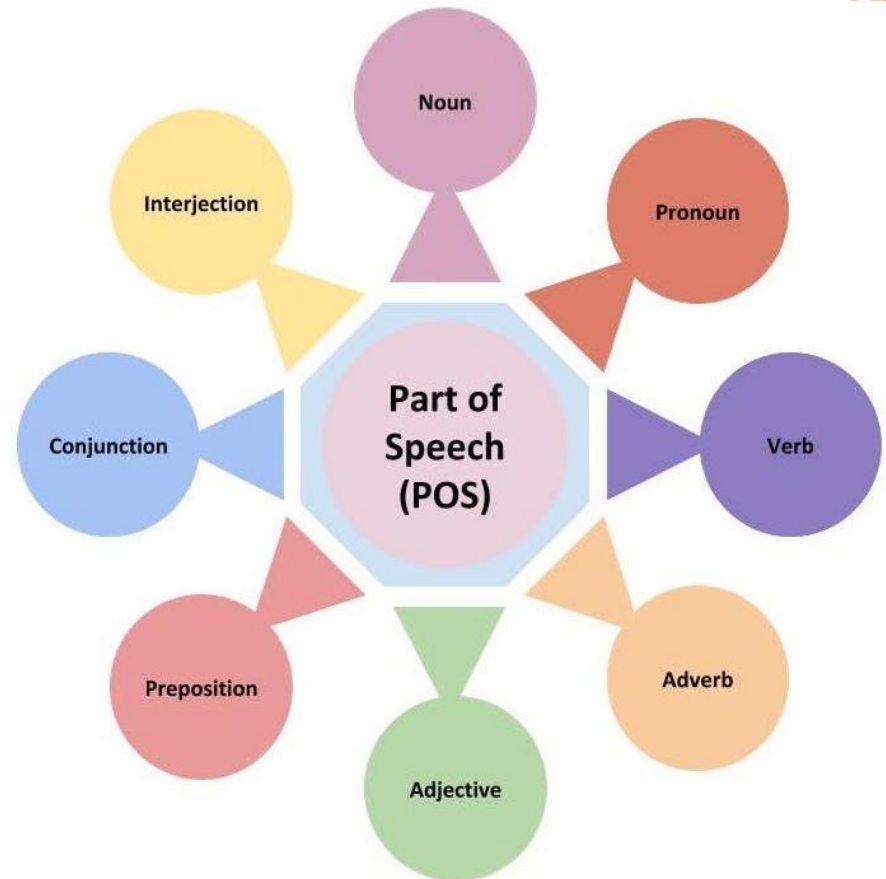
The/**DT** cat/**NN** sat/**VBD** on/**IN** the/**DT** mat/**NN**

Fort/**NNP** Worth/**NNP** is/**VBZ** in/**IN** Texas/**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 TAG SET

## 45 Tags

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	<i>'s</i>	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	<i>\$</i>
JJR	comparative adj	<i>bigger</i>	RP	particle	<i>up, off</i>	#	pound sign	<i>#</i>
JJS	superlative adj	<i>wildest</i>	SYM	symbol	<i>+, %, &amp;</i>	“	left quote	<i>' or “</i>
LS	list item marker	<i>1, 2, One</i>	TO	“to”	<i>to</i>	”	right quote	<i>' or ”</i>
MD	modal	<i>can, should</i>	UH	interjection	<i>ah, oops</i>	(	left paren	<i>[, (, {, &lt;</i>
NN	sing or mass noun	<i>llama</i>	VB	verb base form	<i>eat</i>	)	right paren	<i>], ), }, &gt;</i>
NNS	noun, plural	<i>llamas</i>	VBD	verb past tense	<i>ate</i>	,	comma	<i>,</i>
NNP	proper noun, sing.	<i>IBM</i>	VBG	verb gerund	<i>eating</i>	.	sent-end punc	<i>. ! ?</i>
NNPS	proper noun, plu.	<i>Carolinas</i>	VBN	verb past part.	<i>eaten</i>	:	sent-mid punc	<i>: ; ... - -</i>

(Marcus et al., 1993)

Other corpora: Brown, WSJ, Switchboard

# PART OF SPEECH TAGGING

- A disambiguation task: each word may have different senses/functions
  - The/DT **man/NN** bought/VBD a/DT boat/NN
  - The/DT old/NN **man/VB** the/DT boat/NN
- Some words have MANY functions:

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

# A SIMPLE BASELINE

- Most words are easy to disambiguate
- **Most frequency class:** assign each word (token) its most frequently used class in the training set. (e.g., man/NN)
- Accuracy: 92.34% on the Wall Street Journal (WSJ) dataset!
- State of the art: ~ 97%
- Average English sentence: ~ 14 words
  - Sentence level accuracy:  $0.92^{14} = 31\%$  vs  $0.97^{14} = 65\%$
- POS tagging not solved yet!



# HIDDEN MARKOV MODELS

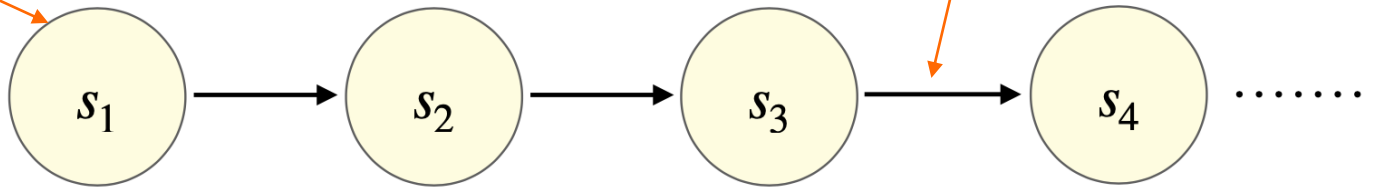
# SOME OBSERVATIONS

- 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 prob. dist.

$P(s_t | s_{t-1})$ : transitional prob.



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

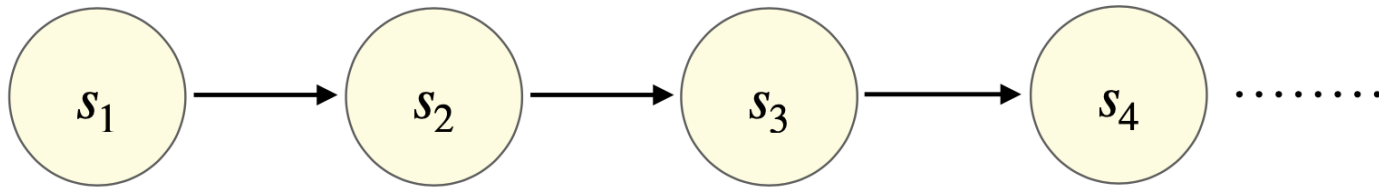
$$P(s_t | s_{<t}) \approx P(s_t | s_{t-1})$$

## QUIZ: MARKOV ASSUMPTION

$$P(s_t | s_{<t}) \approx P(s_t | s_{t-1})$$

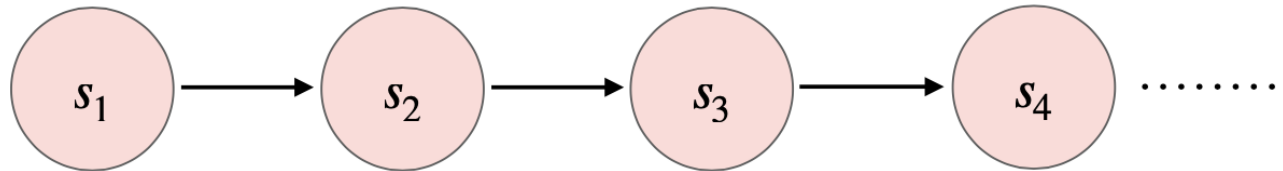
- Where have we seen this before?
  - a) Logistic regression
  - b) Linear regression
  - c) Large language model
  - d) N-gram language model

# MARKOV CHAINS



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

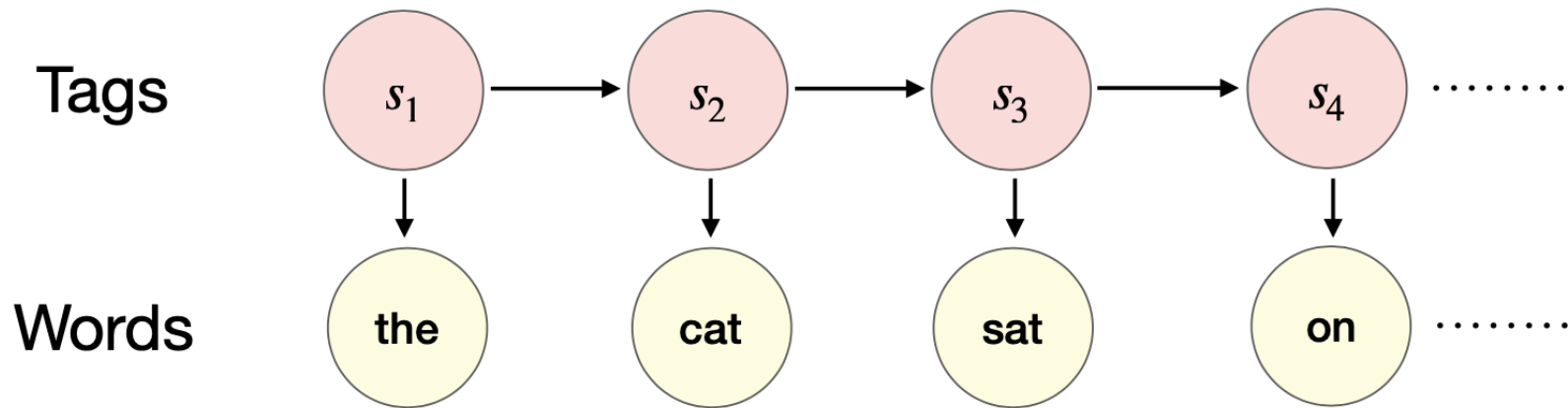
# MARKOV CHAINS



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

- We don't know the tags in the corpus.

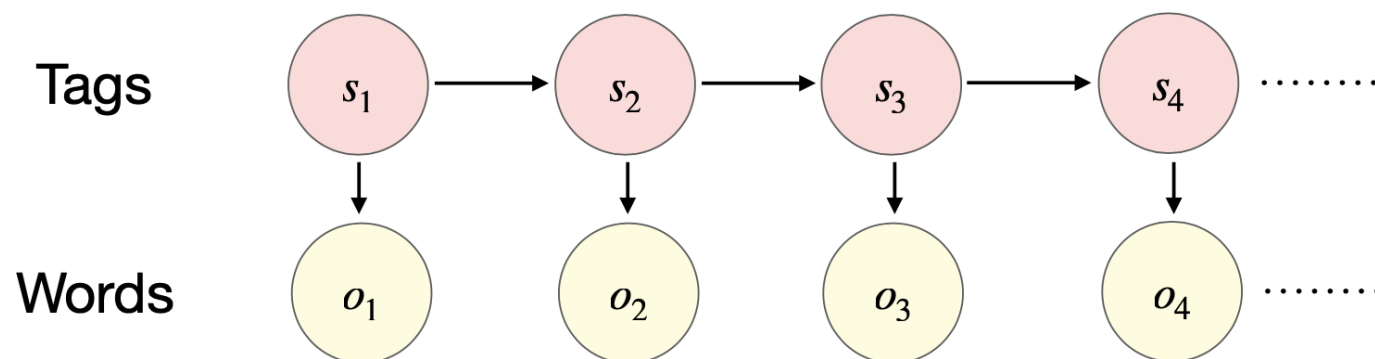
# MARKOV CHAINS



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

- We don't know the tags in the corpus.
- But we do observe the words!
- HMM allows us to jointly reason over both *hidden* and *observed* events.

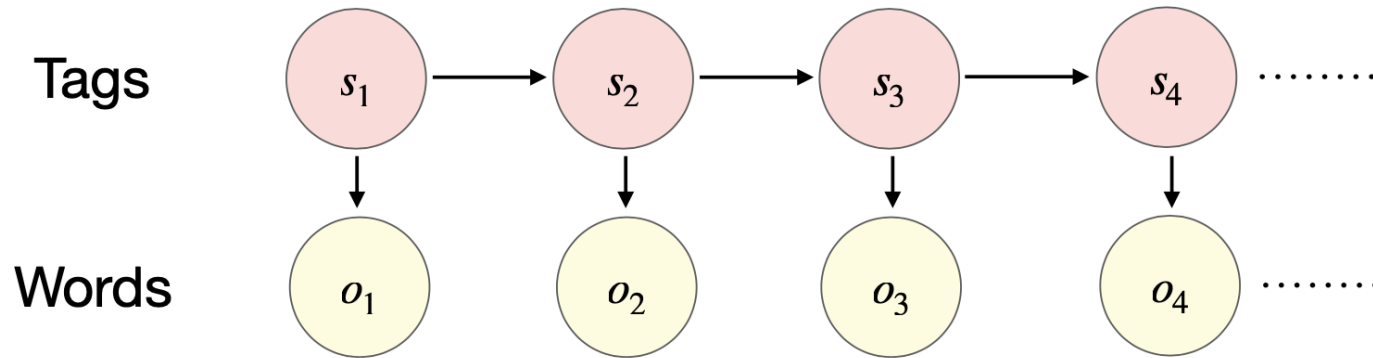
# COMPONENTS OF AN HMM



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



# ASSUMPTIONS



1. Markov assumption:

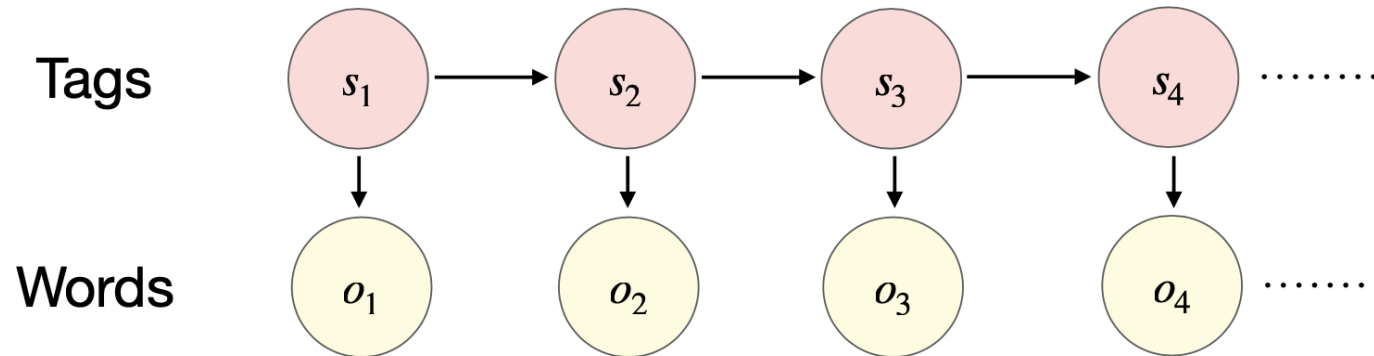
$$P(s_{t+1} \mid s_1, \dots, s_t) = P(s_{t+1} \mid s_t)$$

2. Output independence assumption:

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

**Quiz:** Which one of the two assumptions is stronger, and why?

# SEQUENCE LIKELIHOOD



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

# LEARNING

## Training Set:

1 Pierre/**NNP** Vinken/**NNP** ,/, 61/**CD** years/**NNS** old/**JJ** and/**CC** chairman/**NN** of/**IN** Elsevier/**NNP** N.V./**NNP** ,/, the/**DT** Dutch/**NNP** publishing/**VBG** group/**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** ./.  
...

Maximum likelihood estimate:

Transition prob:  $P(s_i | s_j) = \frac{c(s_i, s_j)}{c(s_j)}$

Emission Prob:  $P(o | s) = \frac{c(s, o)}{c(s)}$

# EXAMPLE: POS TAGGING

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

$$\pi(DT) = 0.8$$

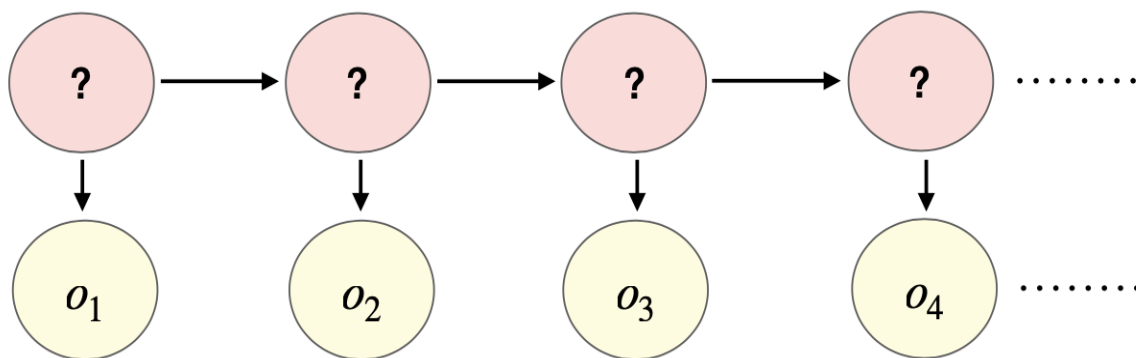
 $s_{t+1}$  $o_t$ 

$s_t$		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

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

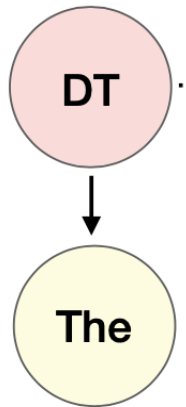
# DECODING WITH HMMs



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

$$\begin{aligned}\hat{S} &= \underset{S}{\operatorname{argmax}} P(S|O) = \underset{S}{\operatorname{argmax}} \frac{P(S)P(O|S)}{P(O)} \quad \text{constant} \\ &= \underset{S}{\operatorname{argmax}} P(S)P(O|S) \\ &= \underset{S}{\operatorname{argmax}} \prod_{i=1}^n P(s_i|s_{i-1})P(o_i|s_i) \\ &\quad \text{transition} \quad \text{emission}\end{aligned}$$

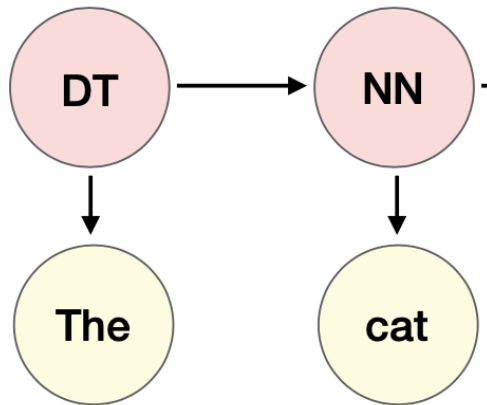
# GREEDY DECODING



$$\underset{s}{\operatorname{argmax}} \Pi(s_1 = s)P(\textit{The} | s) = \textit{'DT'}$$

$$\hat{S} = \underset{s}{\operatorname{argmax}} \prod_{i=1}^n P(s_i | s_{i-1}) P(o_i | s_i)$$

# GREEDY DECODING

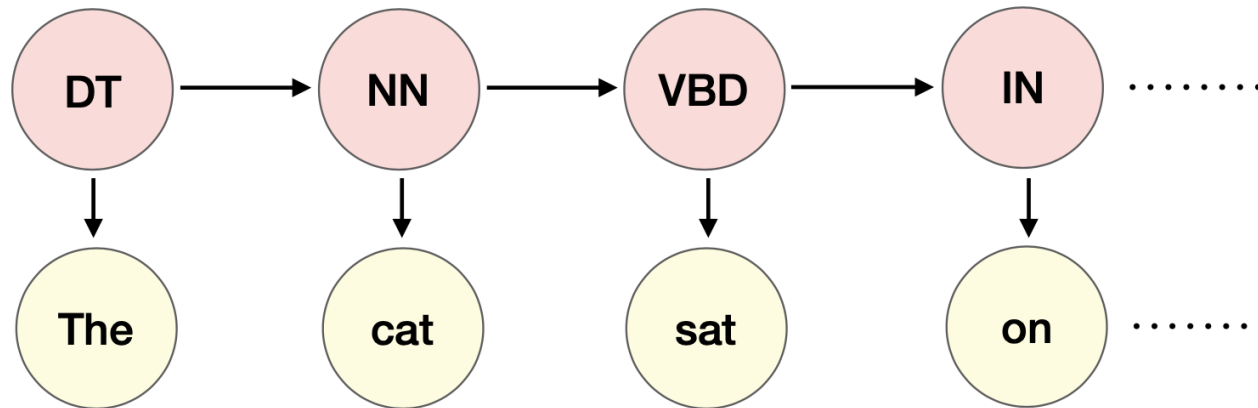


$$\underset{s}{\operatorname{argmax}} \Pi(s_1 = s)P(\textit{The} | s) = \textit{'DT'}$$

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

$$\hat{S} = \underset{s}{\operatorname{argmax}} \prod_{i=1}^n P(s_i | s_{i-1})P(o_i | s_i)$$

# GREEDY DECODING



$$\underset{s}{\operatorname{argmax}} \Pi(s_1 = s)P(\textit{The} | s) = \textit{'DT'}$$

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

$$\forall i, \hat{s}_{i+1} = \underset{s}{\operatorname{argmax}} P(s | \hat{s}_i)P(o_{i+1} | s)$$

Not guaranteed to be optimal: local decision only!



# 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

DT

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

NN

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

VBD

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

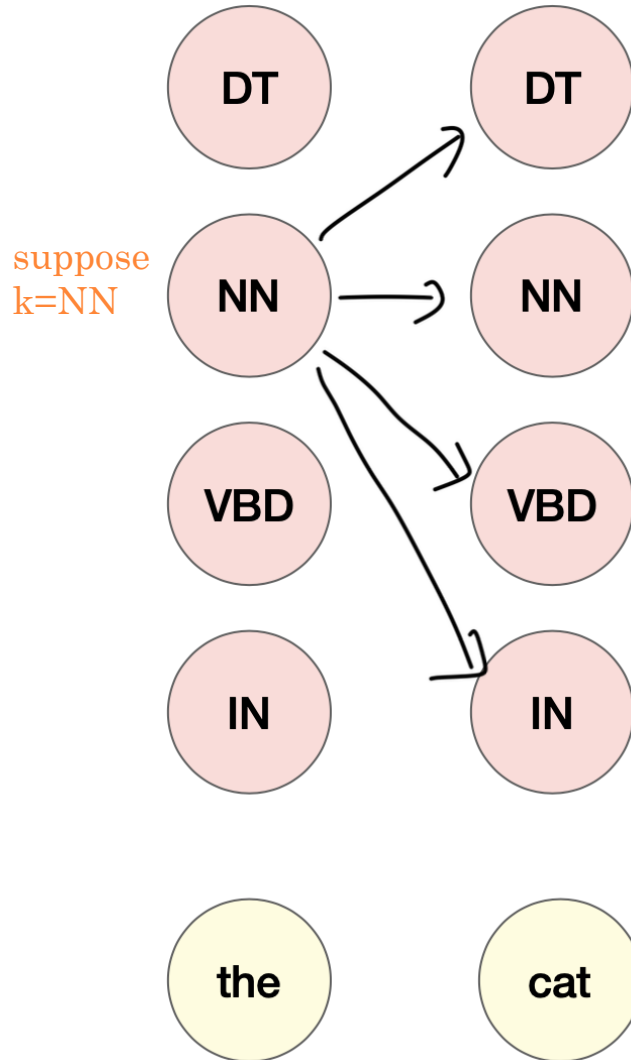
IN

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

the

Forward →

# VITERBI DECODING



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

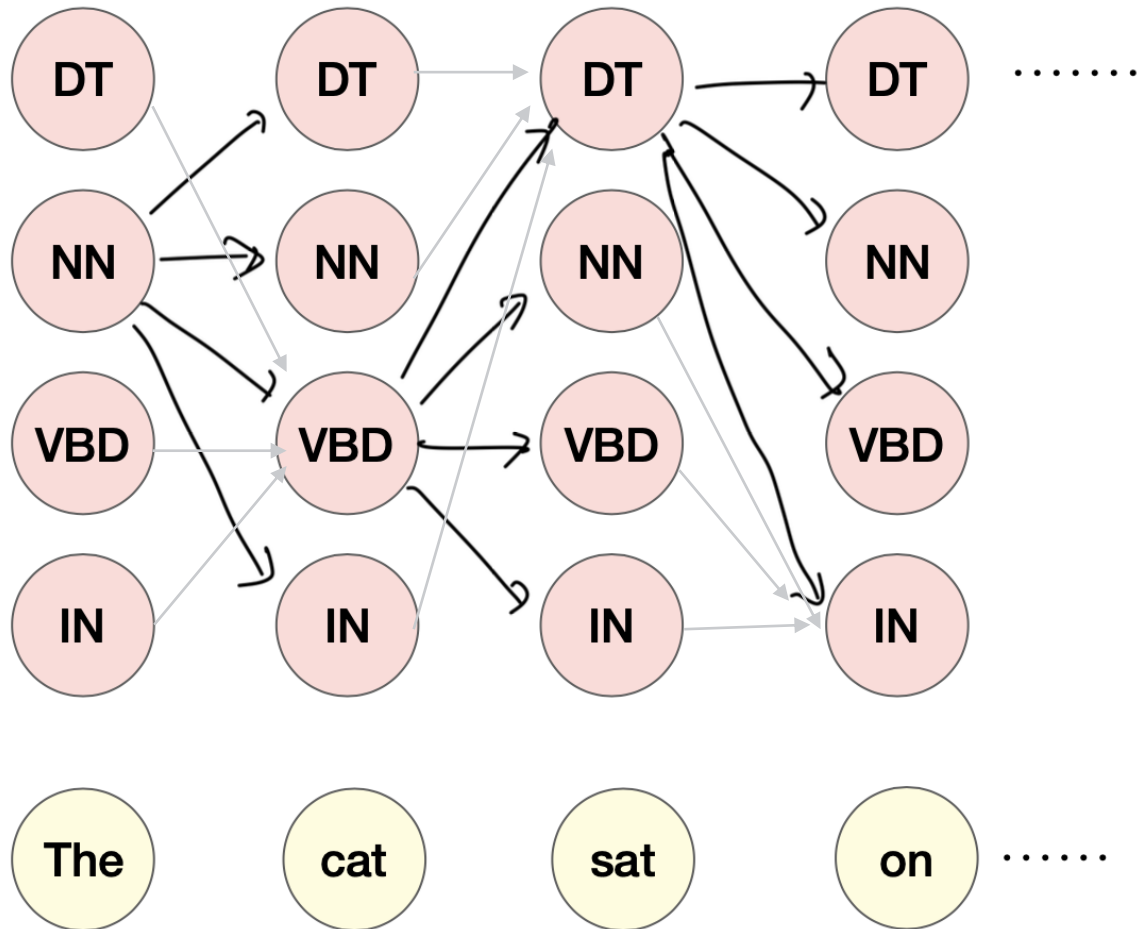
$$M[2,NN] = \max_k M[1,k] P(NN|k) P(\mathbf{cat}|NN)$$

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

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

Forward →

# VITERBI DECODING



This is a recursive process!

Viterbi Algorithm needs to backtrack.

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

# QUIZ: VITERBI ALGORITHM

Assume

$T$  : Number of time steps (sequence length)

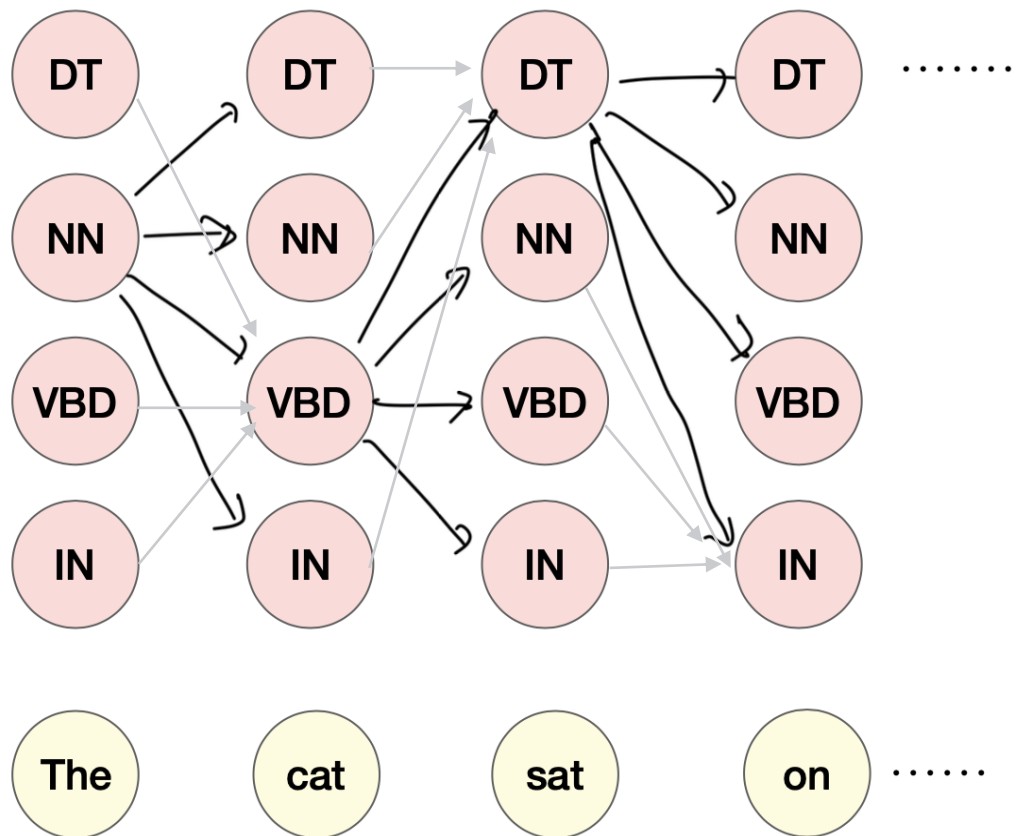
$K$  : Number of states

What is the time complexity of the Viterbi algorithm (in Big O)?

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

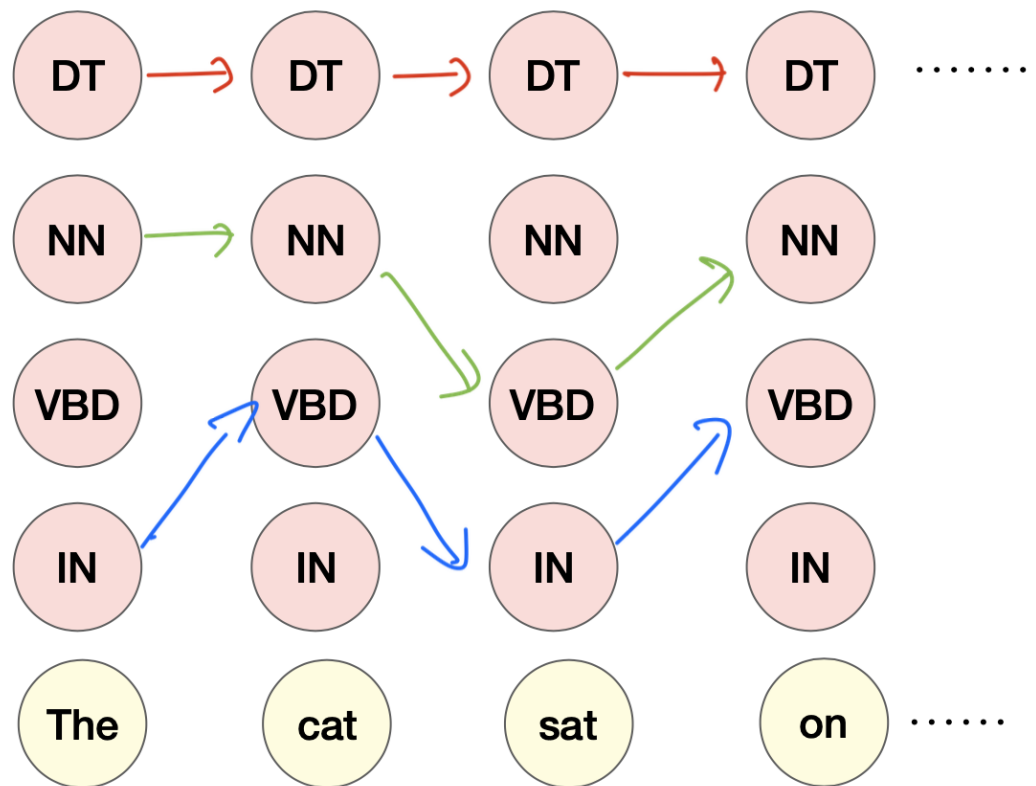
# BEAM SEARCH

- When  $K$  (the number of states) is large, Viterbi algorithm is very expensive!



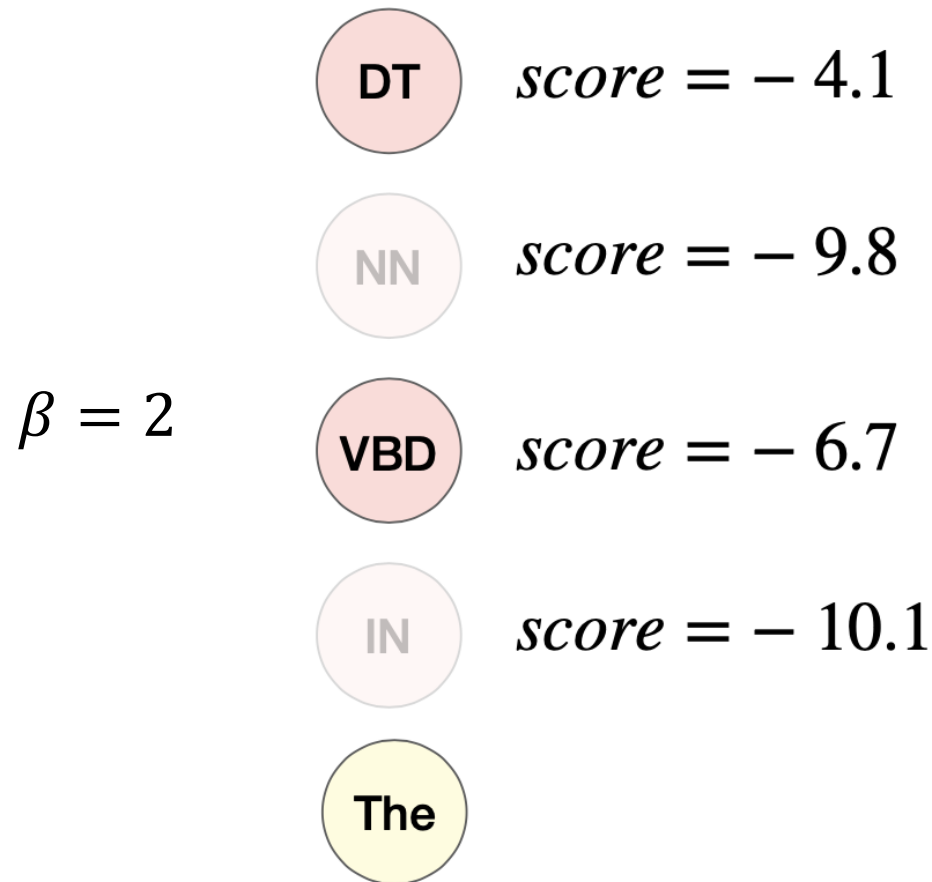
# BEAM SEARCH

- But many paths have very low likelihood!



# BEAM SEARCH

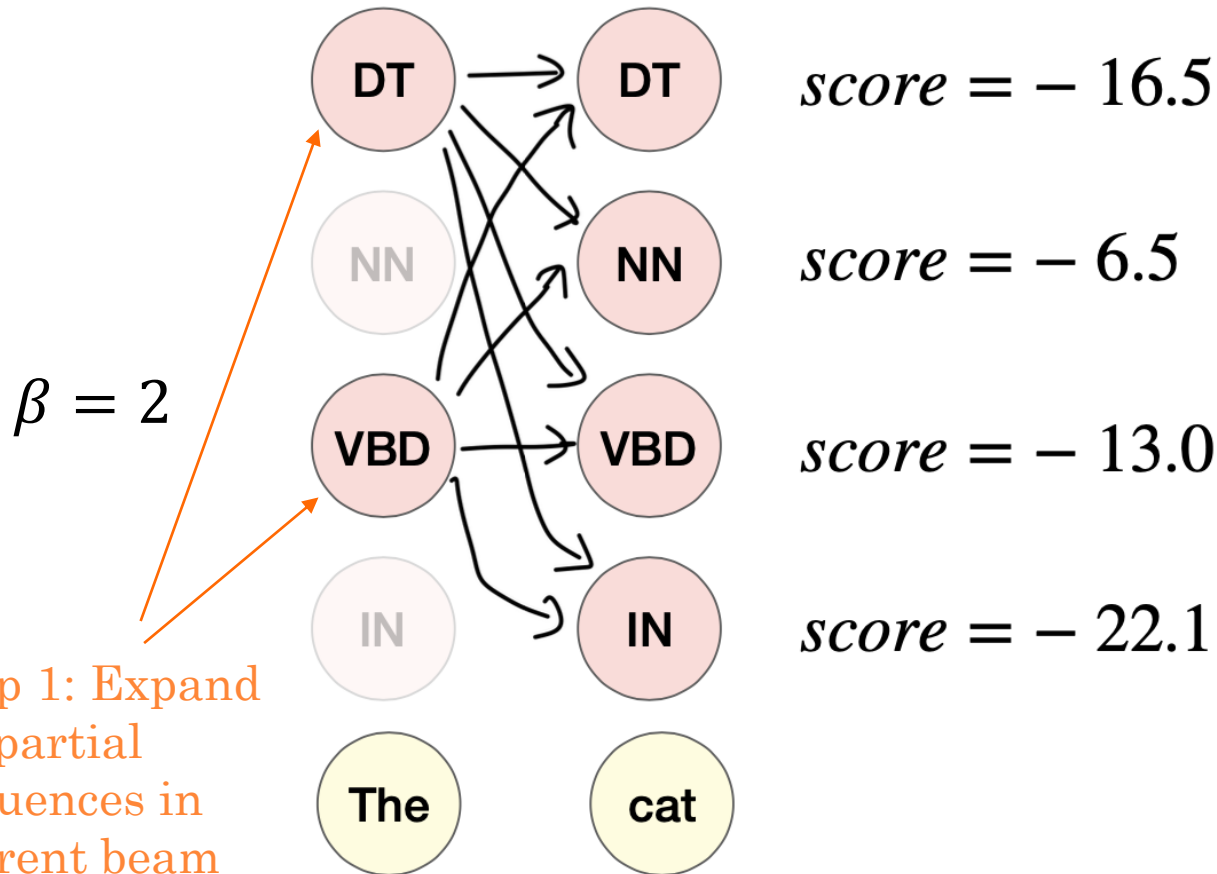
- Keep a fix number  $\beta$  of hypotheses at each stage:





# BEAM SEARCH

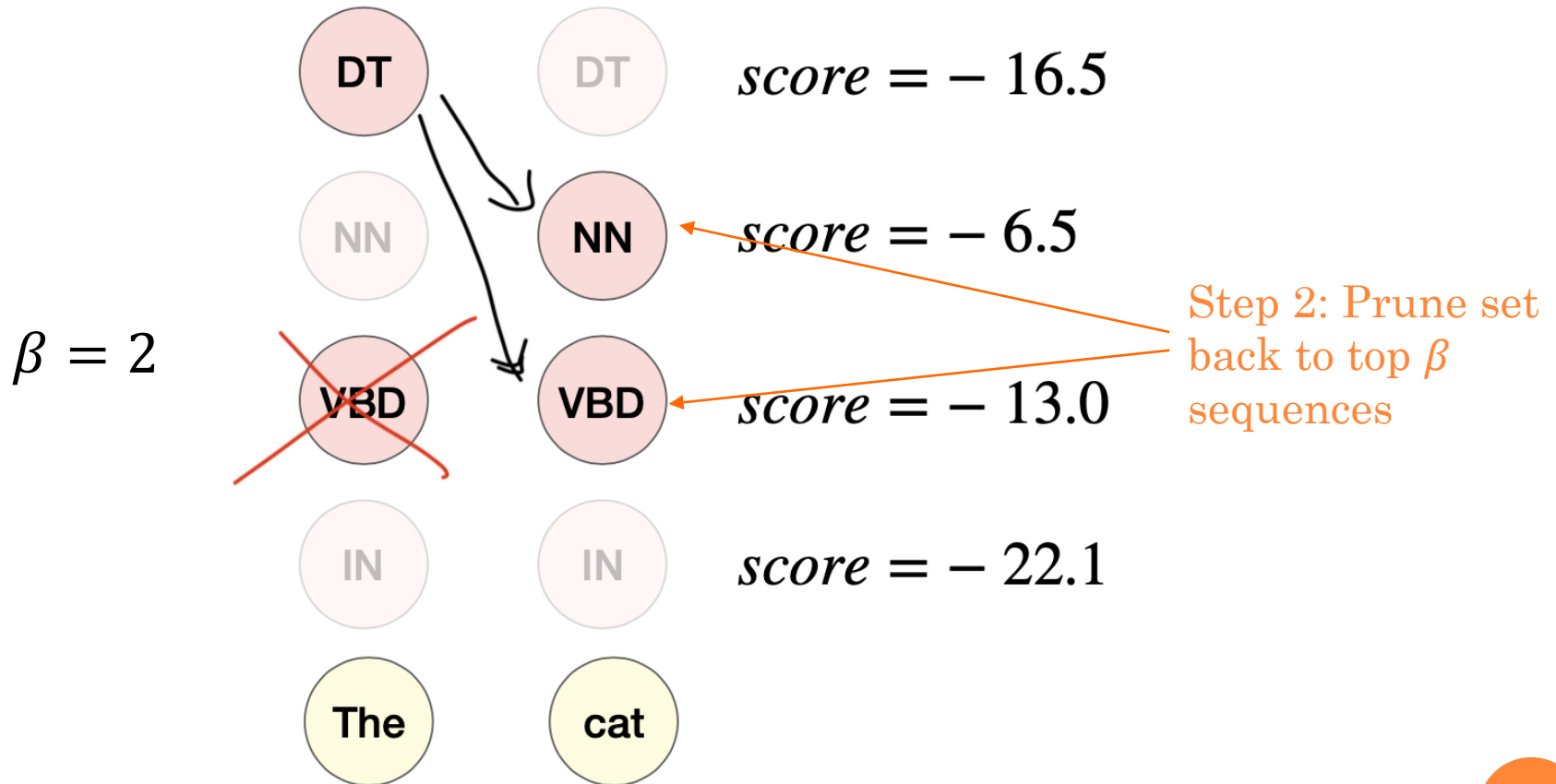
- Keep a fix number  $\beta$  of hypotheses at each stage:



Step 1: Expand  
all partial  
sequences in  
current beam

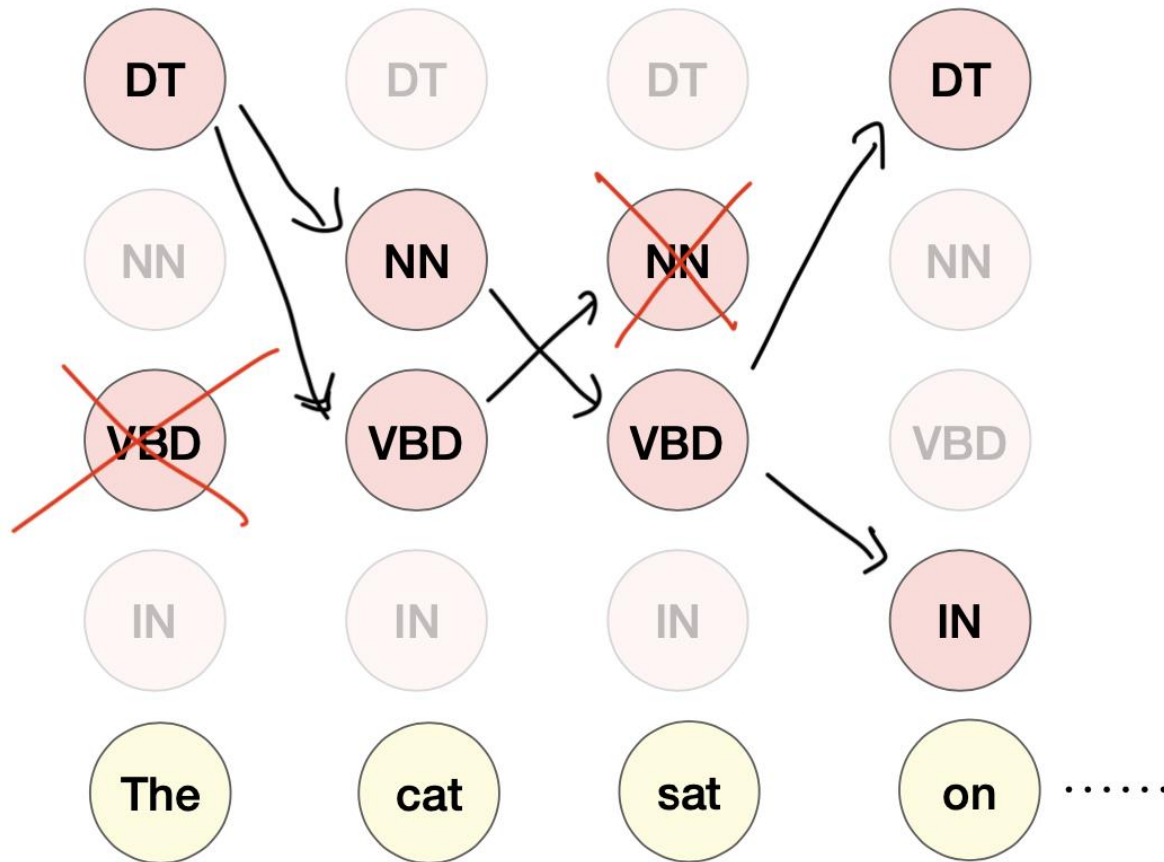
# BEAM SEARCH

- Keep a fix number  $\beta$  of hypotheses at each stage:



# BEAM SEARCH

- Keep a fix number  $\beta$  of hypotheses at each stage:



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