

# Natural Language Processing

Lecture 4: Sequence Labeling with Hidden Markov  
Models. Part-of-Speech Tagging.

11/6/2020

COMS W4705  
Yassine Benajiba

# Garden-Path Sentences

- *The horse raced past the barn.*
- *The horse raced past the barn fell.*
- *The old dog the footsteps of the young.*

# Garden-Path Sentences

- Why does this happen?

past tense verb  
VBD  
*The horse **raced** past the barn **fell*** ???

- **raced** can be a past tense verb or a a past participle (indicating passive voice).
- The verb interpretation is more likely before *fell* is read.

# Garden-Path Sentences

- Why does this happen?

past participle  
VBN  
[The horse **raced** past the barn] fell  
NP VBD

- **raced** can be a past tense verb or a a past participle (indicating passive voice).
- Once *fell* is read, the verb interpretation is impossible.

# Garden-Path Sentences

- Why does this happen?

adjective  
JJ NN  
*[The old dog]* *[the footsteps of the young]*  
NP NP

- **dog** can be a noun or a verb (plural, present tense)

# Garden-Path Sentences

- Why does this happen?

          NN    VB  
*[The old] dog [the footsteps of the young]*  
          NP                                  NP

- **dog** can be a noun or a verb (plural, present tense)

# Parts-of-Speech

- Classes of words that behave alike:
  - Appear in similar contexts.
  - Perform a similar grammatical function in the sentence.
  - Undergo similar morphological transformations.
- ~9 traditional parts-of-speech:
  - noun, pronoun, determiner, adjective, verb, adverb, preposition, conjunction, interjection

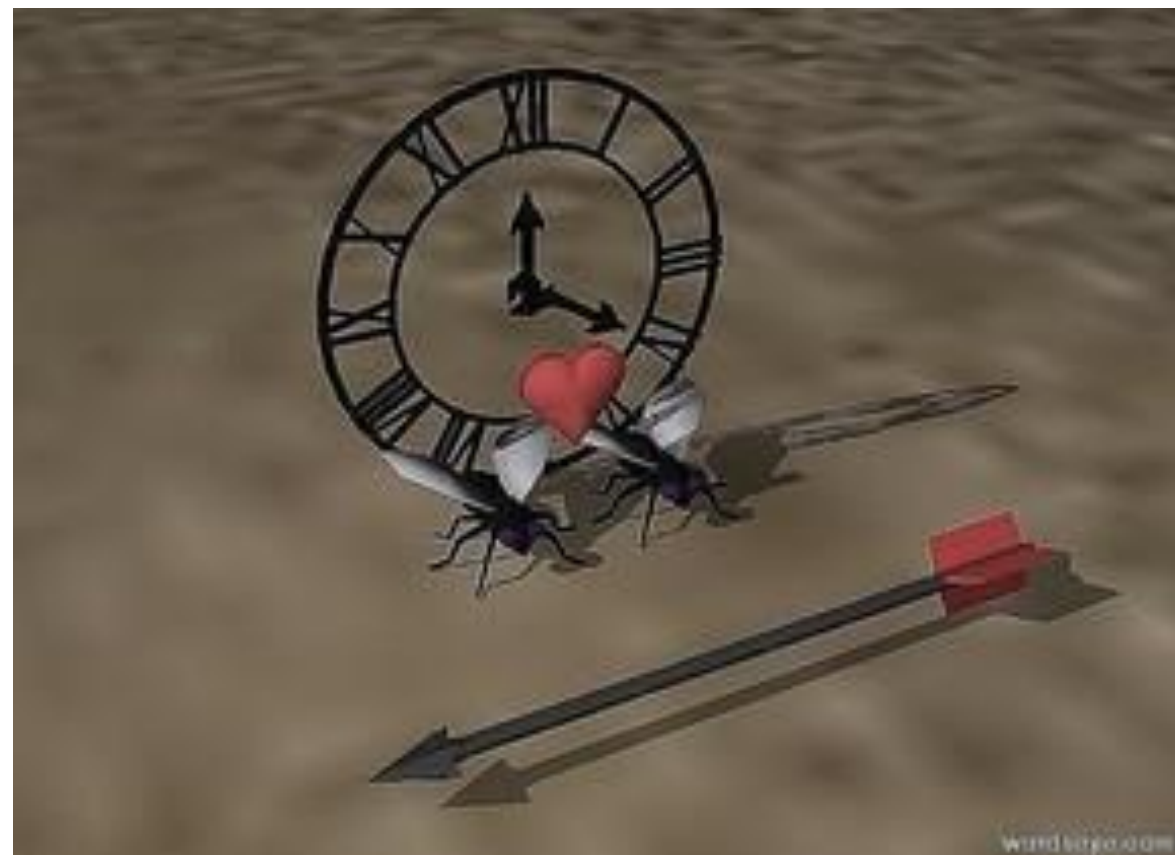
# Syntactic Ambiguities and Parts-of-Speech

- | N / V?      | N / V?       | V / Preposition |                  |
|-------------|--------------|-----------------|------------------|
| <i>Time</i> | <i>flies</i> | <i>like</i>     | <i>an arrow.</i> |



# Syntactic Ambiguities and Parts-of-Speech

- |               |                |             |                  |
|---------------|----------------|-------------|------------------|
| N             | N              | V           |                  |
| [ <i>Time</i> | <i>flies</i> ] | <i>like</i> | <i>an arrow.</i> |
| NP            |                |             |                  |



# Why do we need P.O.S.?

- Interacts with most levels of linguistic representation.
- Speech processing:
  - ***object, object***
  - ***content, content***
- Syntactic parsing
- ...
- P.O.S. tag-set should contain morphological and maybe syntactic information.

# Penn Treebank Tagset

CC	Coordinating conjunction	PRP\$	Possessive pronoun
CD	Cardinal number	RB	Adverb
DT	Determiner	RBR	Adverb, comparative
EX	Existential <i>there</i>	RBS	Adverb, superlative
FW	Foreign word	RP	Particle
IN	Preposition or subordinating conjunction	SYM	Symbol
JJ	Adjective	TO	<i>to</i>
JJR	Adjective, comparative	UH	Interjection
JJS	Adjective, superlative	VB	Verb, base form
LS	List item marker	VBD	Verb, past tense
MD	Modal	VBG	Verb, gerund or present participle
NN	Noun, singular or mass	VBN	Verb, past participle
NNS	Noun, plural	VBP	Verb, non-3rd person singular present
NNP	Proper noun, singular	VBZ	Verb, 3rd person singular present
NNPS	Proper noun, plural	WP	Wh-pronoun
PDT	Predeterminer	WP\$	Possessive wh-pronoun
POS	Possessive ending	WRB	Wh-adverb
PRP	Personal pronoun	plus punctuation symbols	

# P.O.S. Tagsets

- Tagset is language specific.
- Some language capture more morphological information which should be reflected in the tag set.
- “Universal Part Of Speech Tags?”
  - Petrov et al. 2011: Mapping of 25 language specific tag-sets to a common set of 12 universal tags

# Part-of-Speech Tagging

- Goal: Assign a part-of-speech label to each word in a sentence.

<i><b>DT</b></i>	<i><b>NN</b></i>	<i><b>VBD</b></i>	<i><b>DT</b></i>	<i><b>NNS</b></i>	<i><b>IN</b></i>	<i><b>DT</b></i>	<i><b>NN</b></i>	<i><b>.</b></i>
<i>the</i>	<i>koala</i>	<i>put</i>	<i>the</i>	<i>keys</i>	<i>on</i>	<i>the</i>	<i>table</i>	<i>.</i>

- This is an example of a **sequence labeling** task.
- Think of this as a translation task from a sequence of words  $(w_1, w_2, \dots, w_n) \in V^*$ , to a sequence of tags  $(t_1, t_2, \dots, t_n) \in T^*$ .

# Determining Part-of-Speech

- *A blue seat / A child seat:* noun or adj?
- Syntactic tests:
  - *A very **blue** seat*
  - *\*A very **child** seat*
  - *This seat is **blue***
  - *\*This seat is **child***
- Morphological Tests
  - *bluer*
  - *\*childer*

# Determining Part-of-Speech

- Preposition or Particle?

- *He threw **out** the garbage.*

- He threw the garbage **out**.*

**out** is a particle

- *He threw the garbage **out** the door.*

- \*He threw the garbage the door **out***

**out** is a preposition

# Part-of-Speech Tagging

- Goal: Translate from a sequence of words  $(w_1, w_2, \dots, w_n) \in V^*$ , to a sequence of tags  $(t_1, t_2, \dots, t_n) \in T^*$ .
- NLP is full of translation problems from one structure to another. Basic solution:
  - For each translation step:
    1. Construct search space of possible translations.
    2. Find best paths through this space (decoding) according to some performance measure.



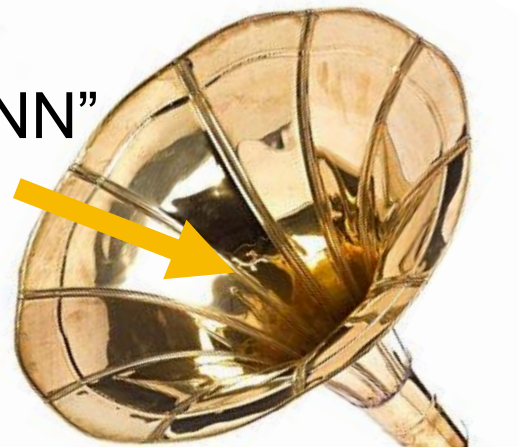
# Bayesian Inference for Sequence Labeling

- Recall Bayesian Inference (Generative Models): Given some observation, infer the value of some *hidden variable*. (see Naive Bayes')
- We can apply this approach to sequence labeling:
  - Assume each word  $w_i$  in the observed sequence  $(w_1, w_2, \dots, w_n) \in V^*$  was generated by some hidden variable  $t_i$ .
  - Infer the most likely sequence of hidden variables given the sequence of observed words.

# Noisy Channel Model

“NN VBZ IN DT NN”

$P(\text{tags})$



$P(\text{words} | \text{tags})$

“time flies like an arrow”

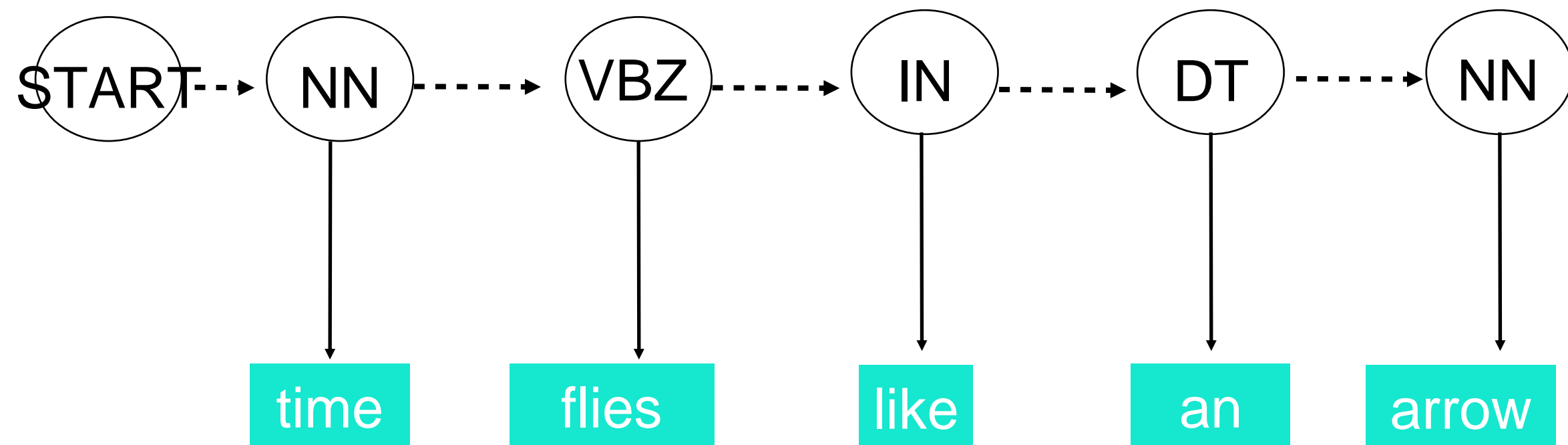
- Goal: figure out what the original input to the the channel was. Use Bayes' rule:

$$\arg \max_{\text{tags}} P(\text{tags} | \text{words}) = \arg \max_{\text{tags}} \frac{P(\text{tags}) \cdot P(\text{word} | \text{tags})}{P(\text{words})}$$

- This model is used widely (speech recognition, MT)

# Hidden Markov Models (HMMs)

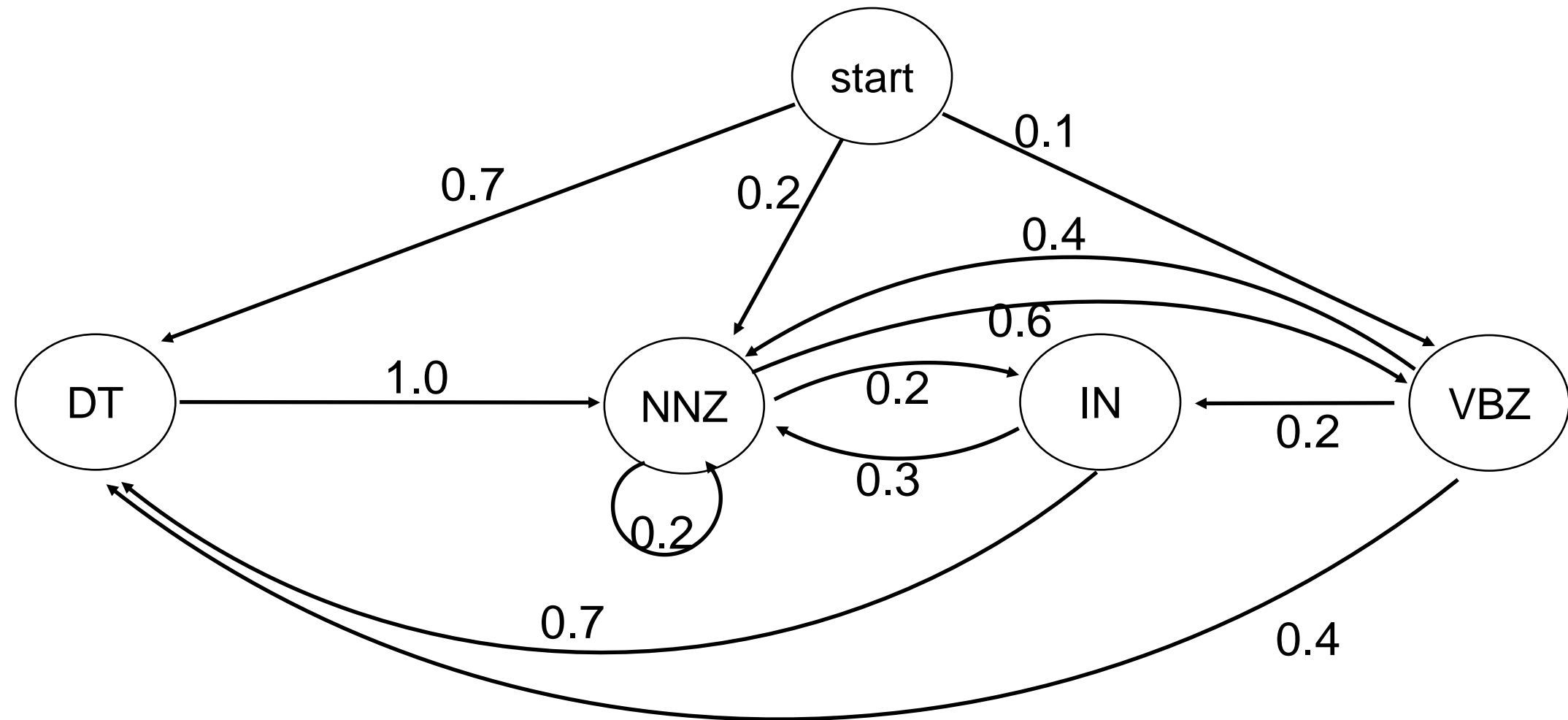
- Generative (Bayesian) probability model.  
Observations: sequences of words.  
Hidden states: sequence of part-of-speech labels.



- Hidden sequence is generated by an n-gram language model (typically a bi-gram model)

$$t_0 = START \quad P(t_1, t_2, \dots, t_n) = \prod_{i=1}^n P(t_i | t_{i-1})$$

# Markov Chains



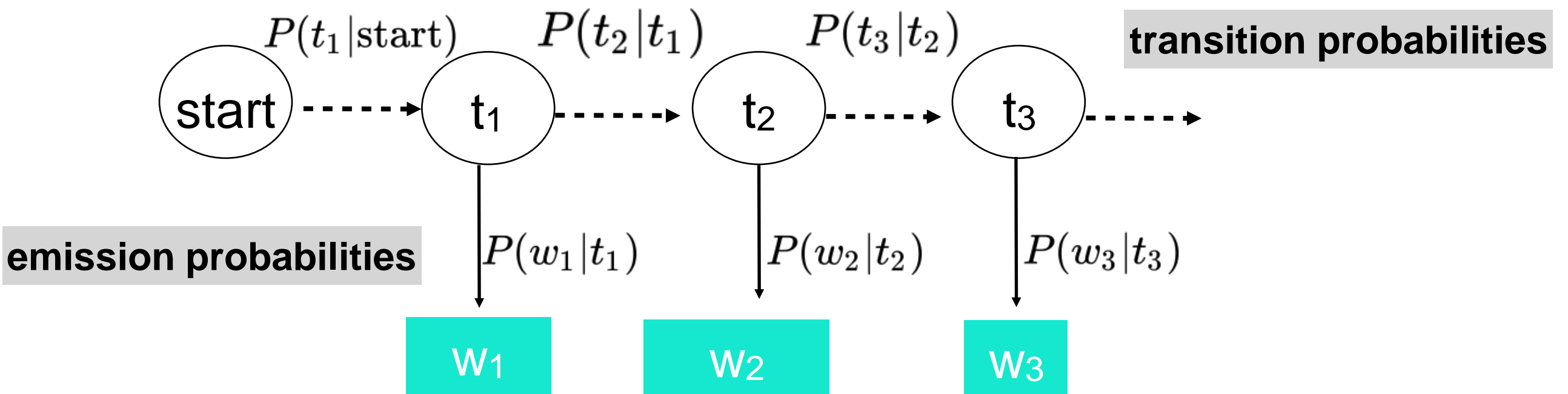
- A **Markov chain** is a sequence of random variables  $X_1, X_2, \dots$
- The domain of these variables is a set of states.
- **Markov assumption:** Next state depends only on current state.

$$P(X_{n+1} | X_1, X_2, \dots, X_n) = P(X_{n+1} | X_n)$$

- This is a special case of a weighted finite state automaton (WFSA).

# Hidden Markov Models (HMMs)

- There are two types of probabilities:  
Transition probabilities and Emission Probabilities.



$$P(t_1, t_2, \dots, t_n, w_1, w_2, \dots, w_n) =$$

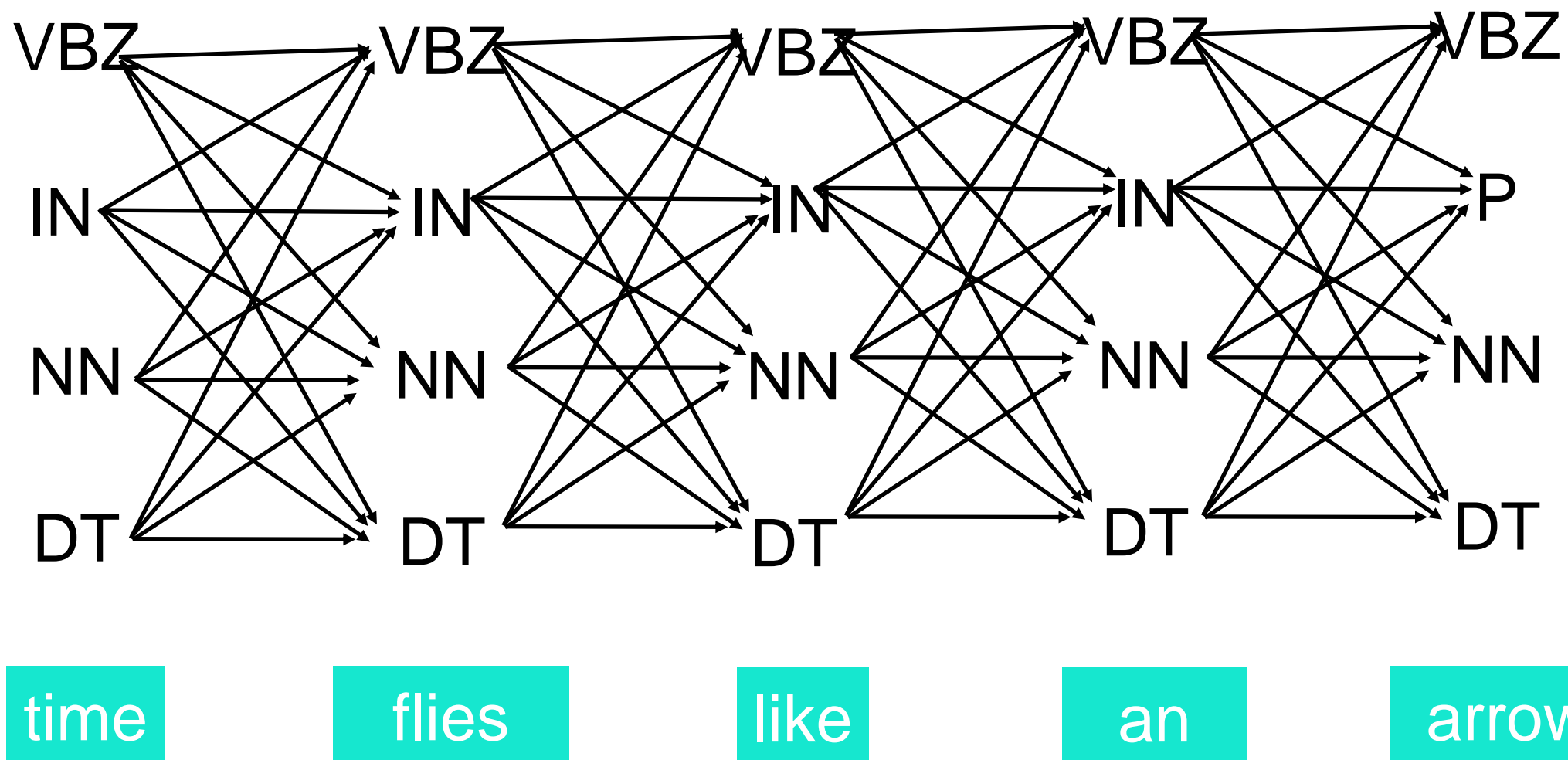
$$P(t_1 | \text{start}) P(w_1 | t_1) P(t_2 | t_1) P(w_2 | t_2) \cdots P(t_n | t_{n-1}) P(w_n | t_n) \cdot P(\text{End} | t_n)$$

$$= \left( \prod_{i=1}^n P(t_i | t_{i-1}) P(w_i | t_i) \right) \cdot P(t_{n+1} | t_n)$$

# Important Tasks on HMMs

- **Decoding:** Given a sequence of words, find the *most likely* probability sequence.  
(Bayesian inference using *Viterbi algorithm*).
- **Evaluation:** Given a sequence of words, find the *total probability for this word sequence* given an HMM.  
Note that we can view the HMM as another type of language model. (Forward algorithm)
- **Training:** Estimate emission and transition probabilities from training data. (MLE, Forward-Backward a.k.a Baum-Welch algorithm)

# Decoding HMMs



Goal: Find the path with the highest total probability (given the words)

$$\arg \max_{t_1, \dots, t_n} \prod_{i=1}^n P(t_i | t_{i-1}) P(w_i | t_i)$$

There are  $d^n$  paths for  $n$  words and  $d$  tags.

# Viterbi Algorithm

- **Input:** Sequence of observed words  $w_1, \dots, w_n$
- Create a table  $\pi$ , such that each entry  $\pi[k, t]$  contains the score of the highest-probability sequence ending in tag  $t$  at time  $k$ .
- initialize  $\pi[0, \text{start}] = 1.0$  and  $\pi[0, t] = 0.0$  for all tags  $t \in T$ .
- for  $k = 1$  to  $n$ :
  - for  $t \in T$ :
    - $\pi[k, t] \leftarrow \max_s \pi[k - 1, s] \cdot P(t|s) \cdot P(w_k|t)$ 
      - transition probability
      - emission probability
- **return**  $\max_s \pi[n, s]$



# Emission Probabilities

- $P(\text{time} \mid \text{VB}) = 0.2$   
 $P(\text{flies} \mid \text{VB}) = 0.3$   
 $P(\text{like} \mid \text{VB}) = 0.5$
- $P(\text{time} \mid \text{NN}) = 0.3$   
 $P(\text{flies} \mid \text{NN}) = 0.2$   
 $P(\text{arrow} \mid \text{NN}) = 0.5$
- $P(\text{like} \mid \text{P}) = 1.0$
- $P(\text{an} \mid \text{DT}) = 1.0$

# Viterbi Algorithm

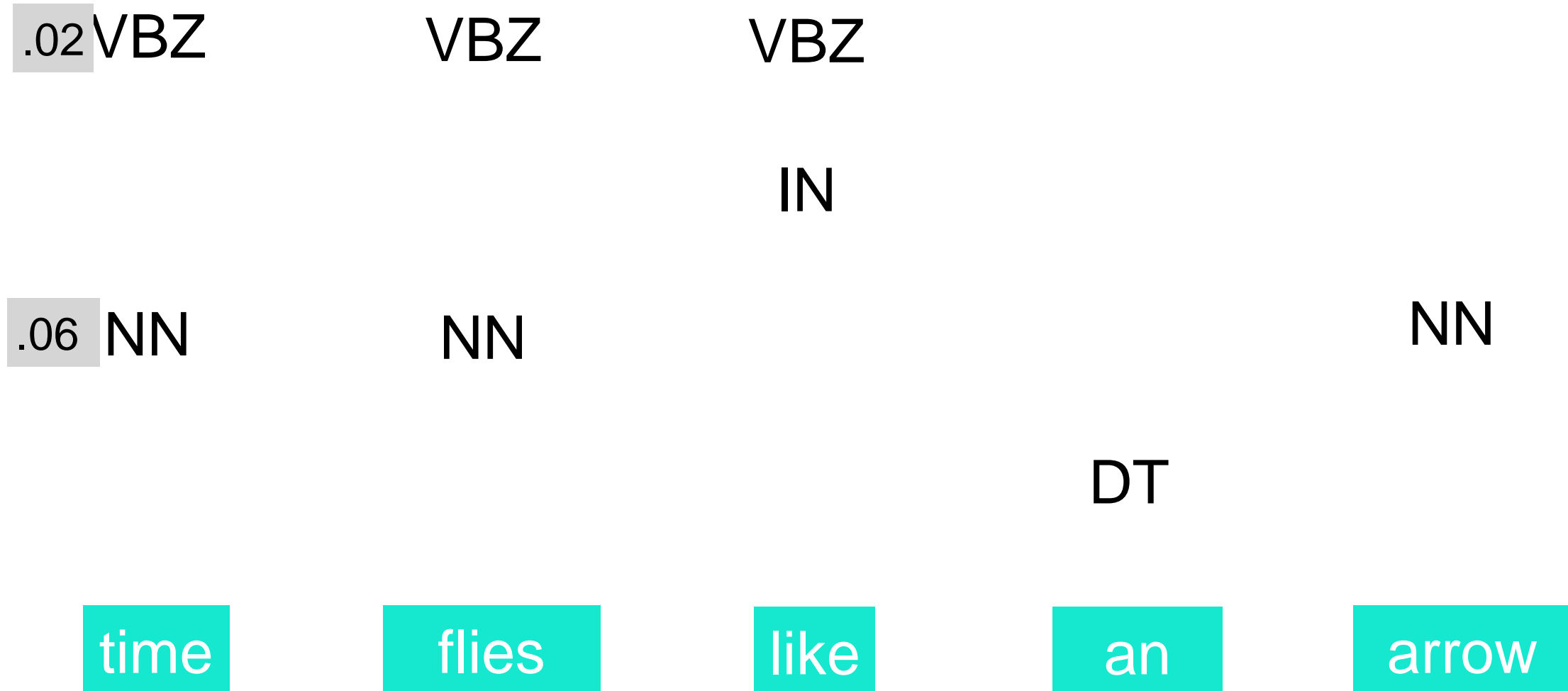
- $P(\text{time} \mid \text{VB}) = 0.2$   
 $P(\text{flies} \mid \text{VB}) = 0.3$   
 $P(\text{like} \mid \text{VB}) = 0.5$
- $P(\text{time} \mid \text{NN}) = 0.3$   
 $P(\text{flies} \mid \text{NN}) = 0.2$   
 $P(\text{arrow} \mid \text{NN}) = 0.5$
- $P(\text{like} \mid \text{P}) = 1.0$
- $P(\text{an} \mid \text{DT}) = 1.0$

$.1 \times .2 = .02$	VBZ	VBZ	VBZ	VBZ	VBZ
0	IN	IN	IN	IN	IN
$.2 \times .3 = .06$	NN	NN	NN	NN	NN
0	DT	DT	DT	DT	DT
time	flies	like	an	arrow	

- Idea: Because of the Markov assumption, we only need the probabilities for  $X_n$  to compute the probabilities for  $X_{n+1}$ . This suggests a **dynamic programming** algorithm.

# Viterbi Algorithm

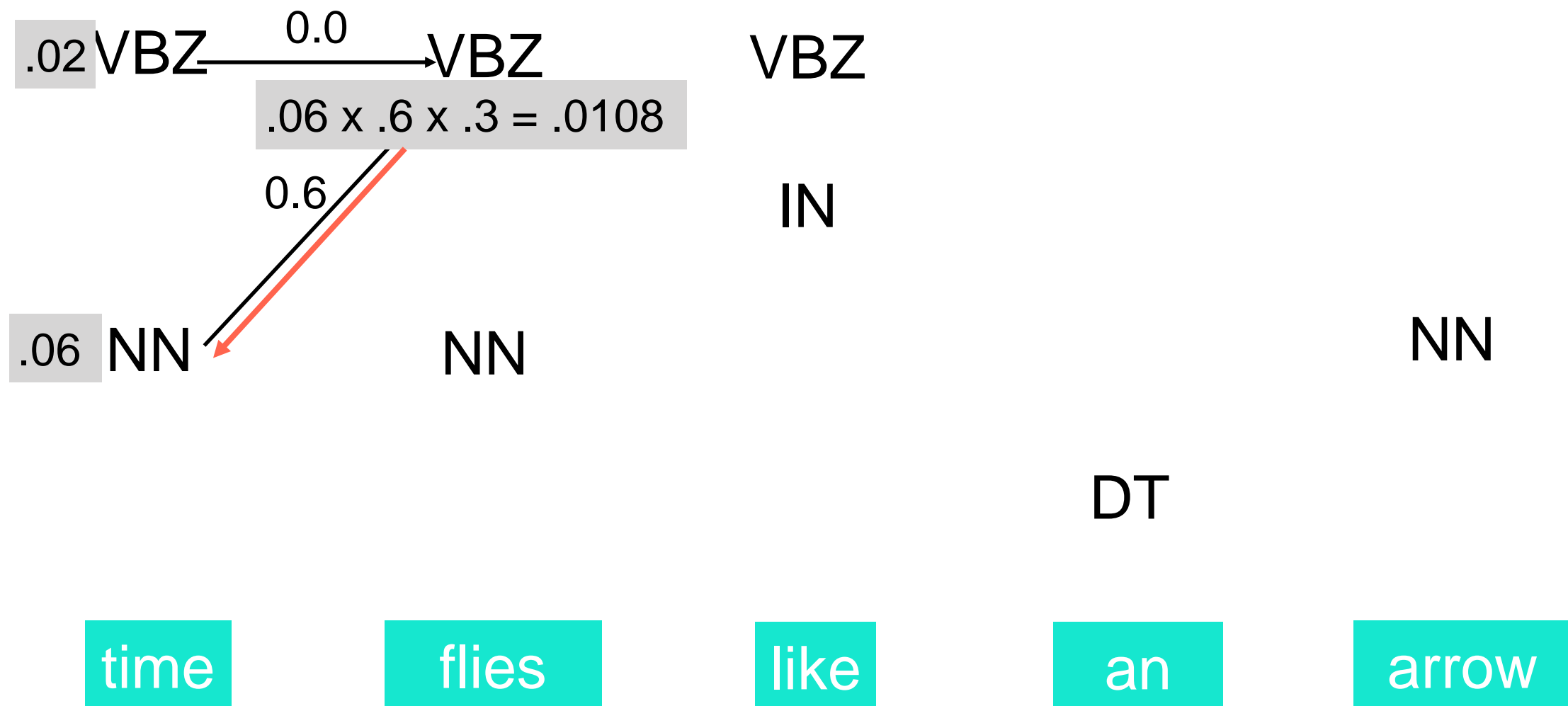
- $P(\text{time} \mid \text{VB}) = 0.2$   
 $P(\text{flies} \mid \text{VB}) = 0.3$   
 $P(\text{like} \mid \text{VB}) = 0.5$
- $P(\text{time} \mid \text{NN}) = 0.3$   
 $P(\text{flies} \mid \text{NN}) = 0.2$   
 $P(\text{arrow} \mid \text{NN}) = 0.5$
- $P(\text{like} \mid \text{P}) = 1.0$
- $P(\text{an} \mid \text{DT}) = 1.0$



- Idea: Because of the Markov assumption, we only need the probabilities for  $X_n$  to compute the probabilities for  $X_{n+1}$ . This suggests a **dynamic programming** algorithm.

# Viterbi Algorithm

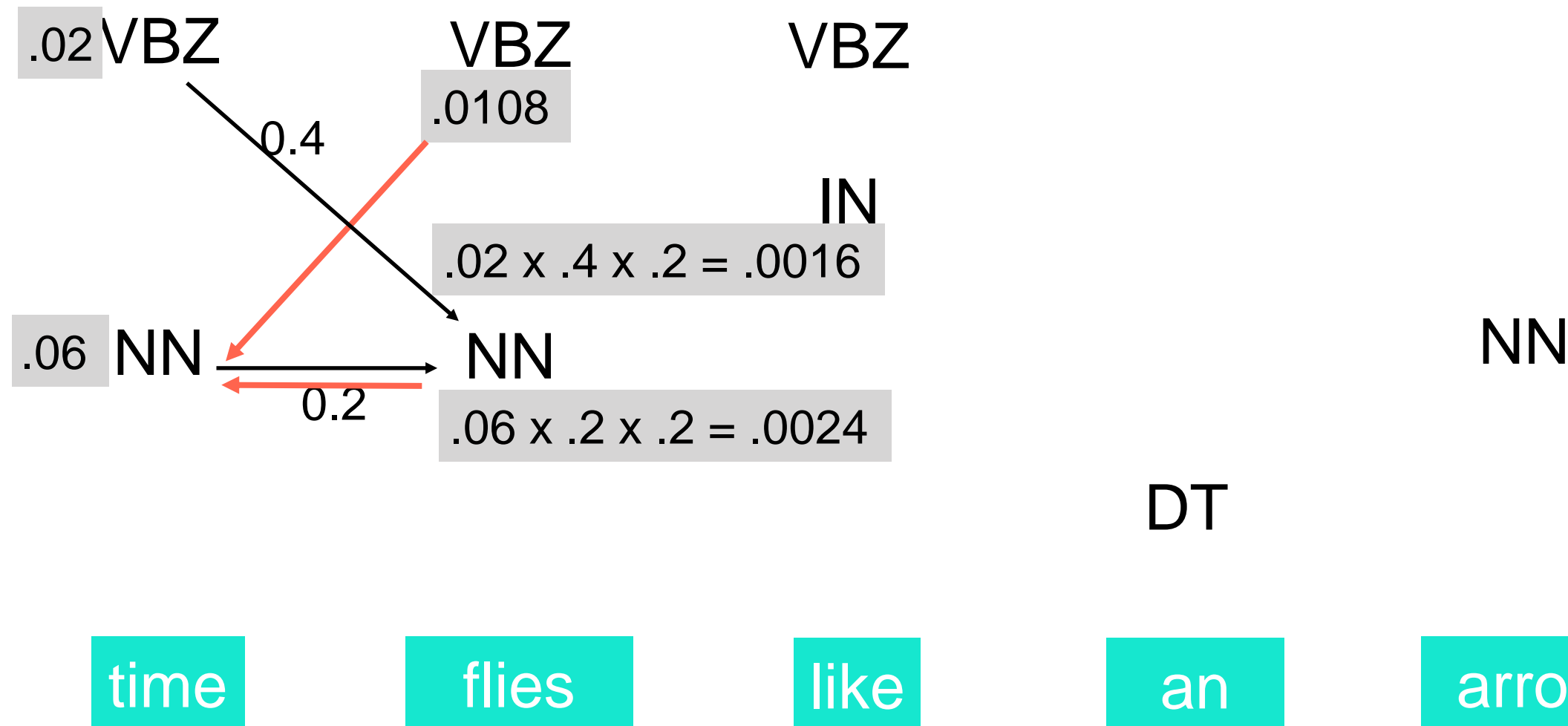
- $P(\text{time} \mid \text{VB}) = 0.2$   
 $P(\text{flies} \mid \text{VB}) = 0.3$   
 $P(\text{like} \mid \text{VB}) = 0.5$
- $P(\text{time} \mid \text{NN}) = 0.3$   
 $P(\text{flies} \mid \text{NN}) = 0.2$   
 $P(\text{arrow} \mid \text{NN}) = 0.5$
- $P(\text{like} \mid \text{P}) = 1.0$
- $P(\text{an} \mid \text{DT}) = 1.0$



- Idea: Because of the Markov assumption, we only need the probabilities for  $X_n$  to compute the probabilities for  $X_{n+1}$ . This suggests a **dynamic programming** algorithm.

# Viterbi Algorithm

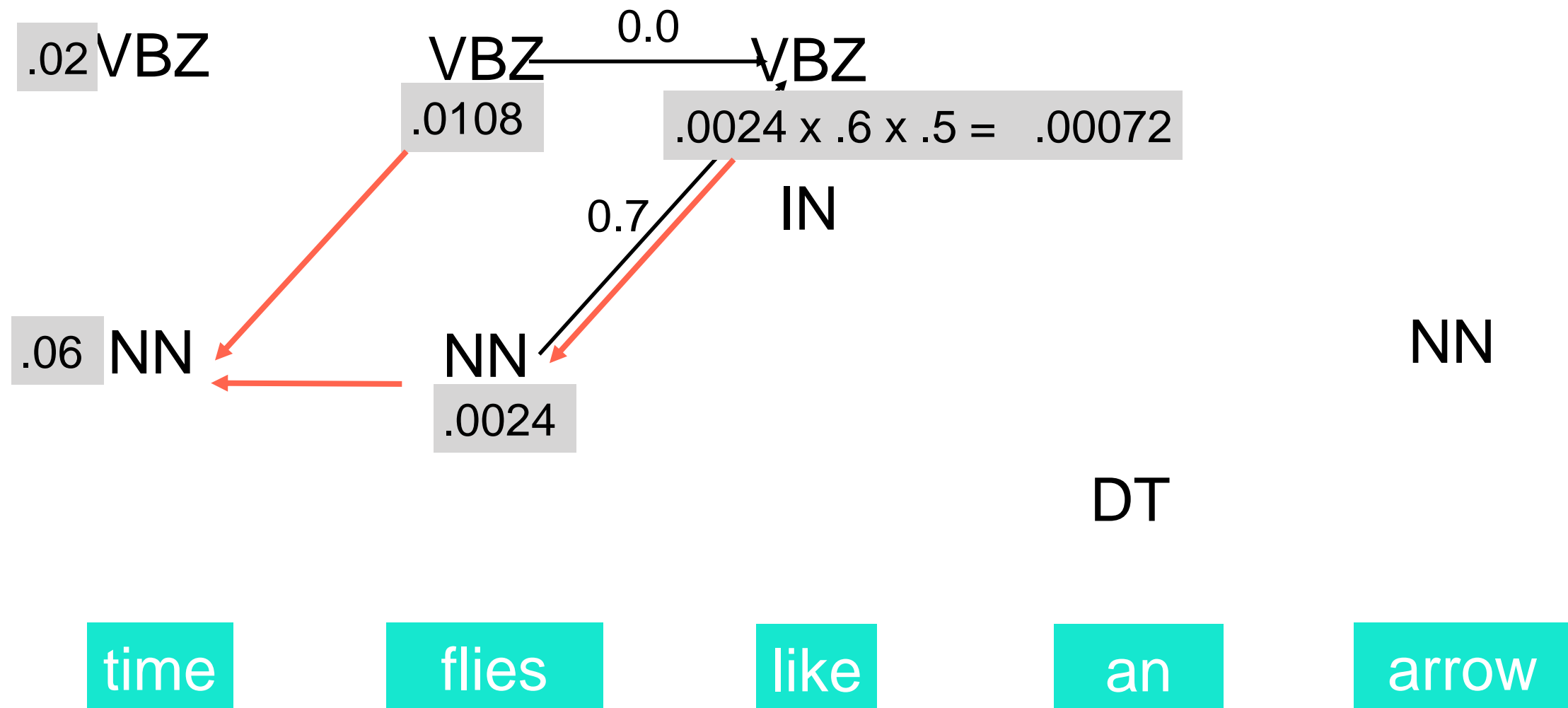
- $P(\text{time} \mid \text{VB}) = 0.2$   
 $P(\text{flies} \mid \text{VB}) = 0.3$   
 $P(\text{like} \mid \text{VB}) = 0.5$
- $P(\text{time} \mid \text{NN}) = 0.3$   
 $P(\text{flies} \mid \text{NN}) = 0.2$   
 $P(\text{arrow} \mid \text{NN}) = 0.5$
- $P(\text{like} \mid \text{P}) = 1.0$
- $P(\text{an} \mid \text{DT}) = 1.0$



- Idea: Because of the Markov assumption, we only need the probabilities for  $X_n$  to compute the probabilities for  $X_{n+1}$ . This suggests a **dynamic programming** algorithm.

# Viterbi Algorithm

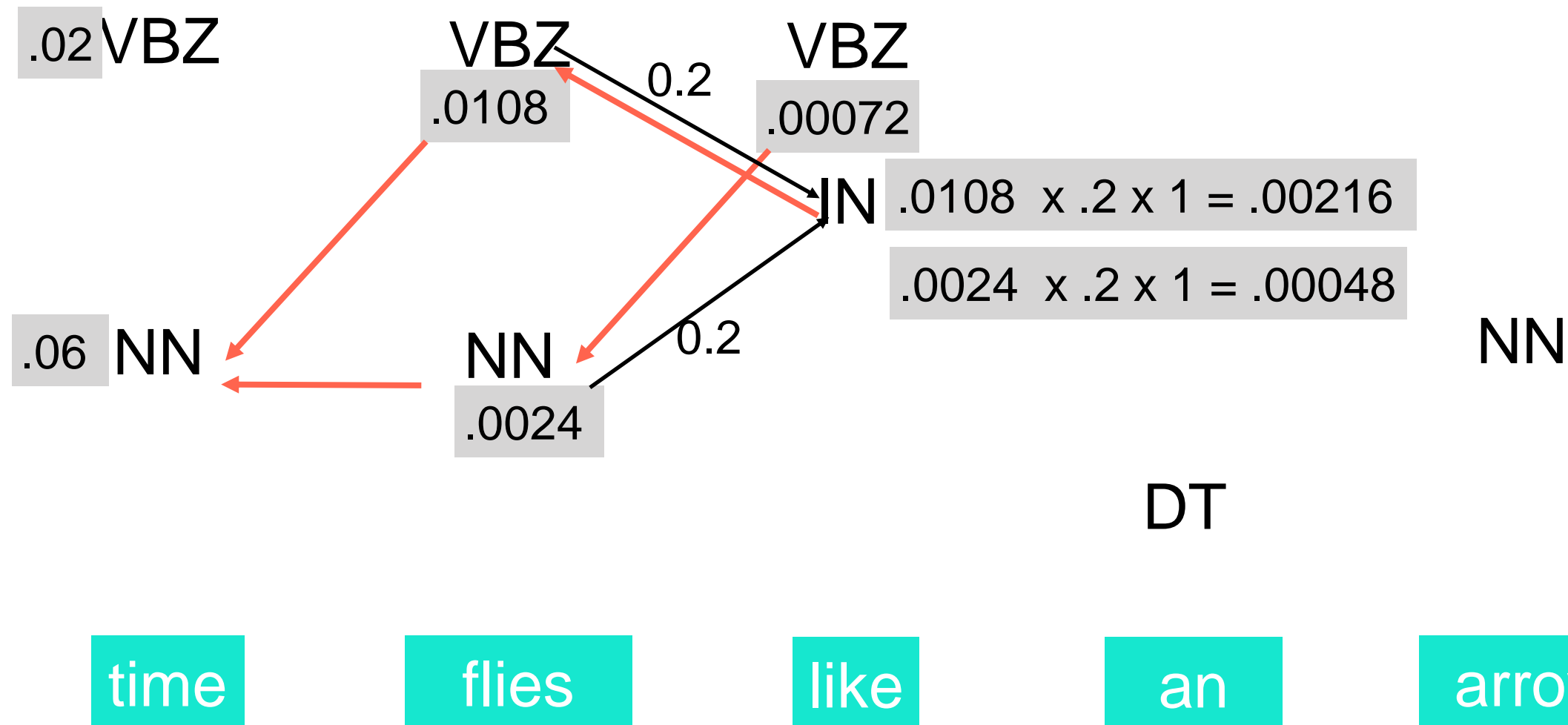
- $P(\text{time} \mid \text{VB}) = 0.2$   
 $P(\text{flies} \mid \text{VB}) = 0.3$   
 $P(\text{like} \mid \text{VB}) = 0.5$
- $P(\text{time} \mid \text{NN}) = 0.3$   
 $P(\text{flies} \mid \text{NN}) = 0.2$   
 $P(\text{arrow} \mid \text{NN}) = 0.5$
- $P(\text{like} \mid \text{P}) = 1.0$
- $P(\text{an} \mid \text{DT}) = 1.0$



- Idea: Because of the Markov assumption, we only need the probabilities for  $X_n$  to compute the probabilities for  $X_{n+1}$ . This suggests a **dynamic programming** algorithm.

# Viterbi Algorithm

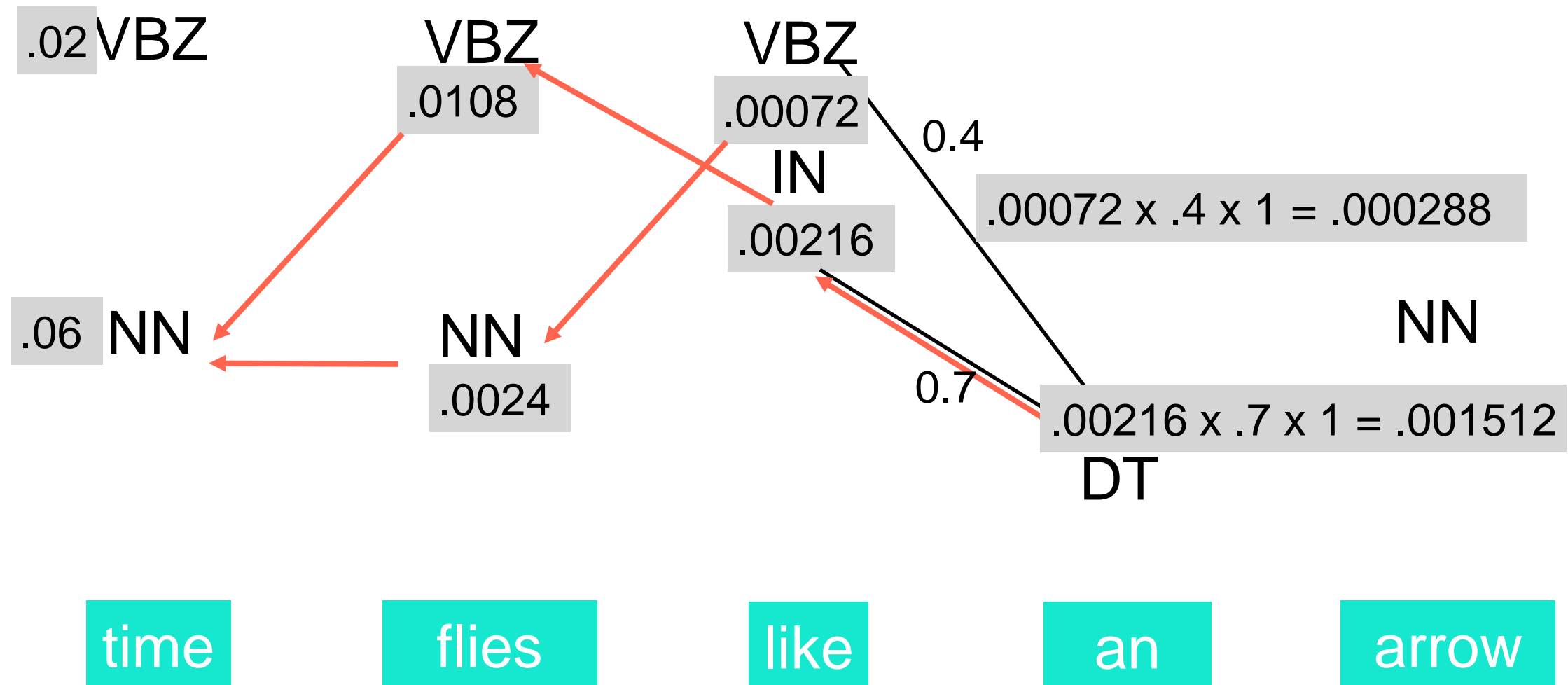
- $P(\text{time} \mid \text{VB}) = 0.2$   
 $P(\text{flies} \mid \text{VB}) = 0.3$   
 $P(\text{like} \mid \text{VB}) = 0.5$
- $P(\text{time} \mid \text{NN}) = 0.3$   
 $P(\text{flies} \mid \text{NN}) = 0.2$   
 $P(\text{arrow} \mid \text{NN}) = 0.5$
- $P(\text{like} \mid \text{P}) = 1.0$
- $P(\text{an} \mid \text{DT}) = 1.0$



- Idea: Because of the Markov assumption, we only need the probabilities for  $X_n$  to compute the probabilities for  $X_{n+1}$ . This suggests a **dynamic programming** algorithm.

# Viterbi Algorithm

- $P(\text{time} \mid \text{VB}) = 0.2$   
 $P(\text{flies} \mid \text{VB}) = 0.3$   
 $P(\text{like} \mid \text{VB}) = 0.5$
- $P(\text{time} \mid \text{NN}) = 0.3$   
 $P(\text{flies} \mid \text{NN}) = 0.2$   
 $P(\text{arrow} \mid \text{NN}) = 0.5$
- $P(\text{like} \mid \text{P}) = 1.0$
- $P(\text{an} \mid \text{DT}) = 1.0$

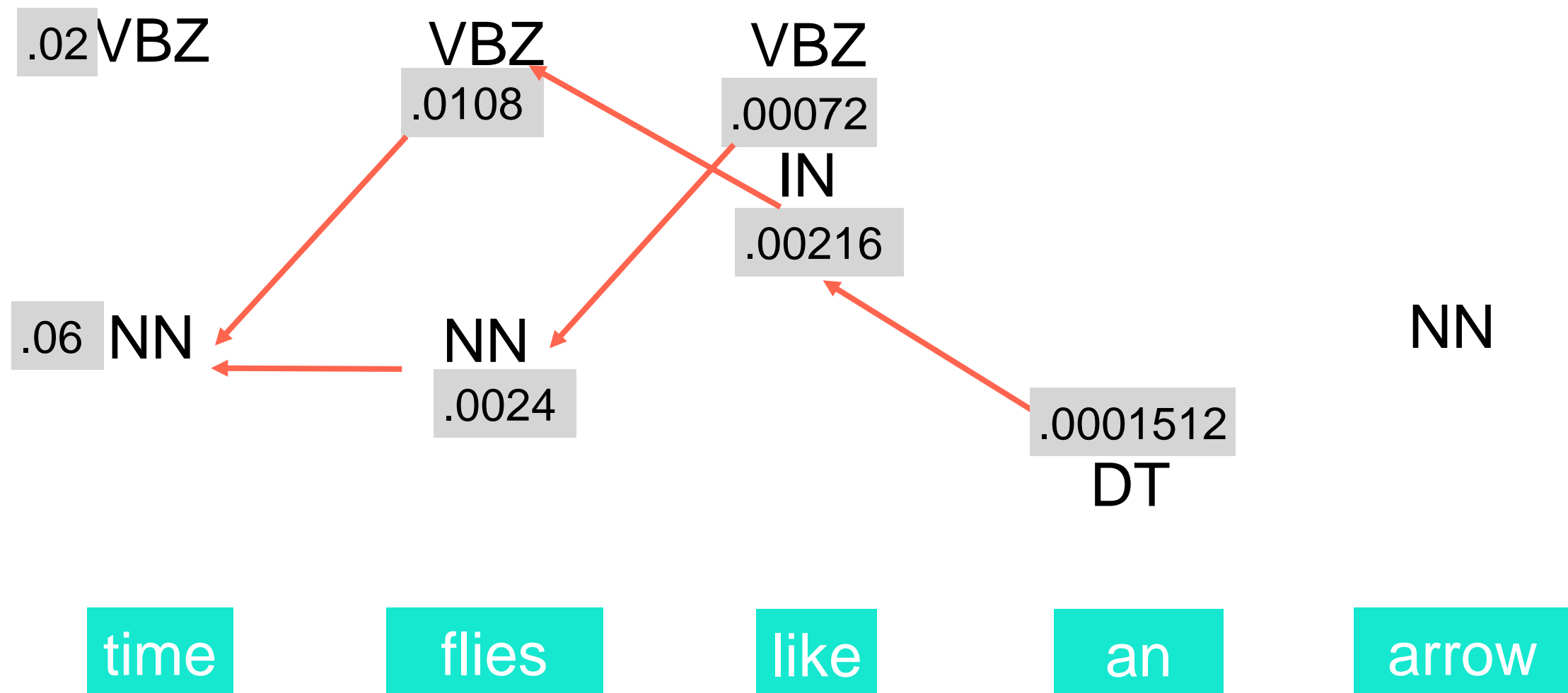


- Idea: Because of the Markov assumption, we only need the probabilities for  $X_n$  to compute the probabilities for  $X_{n+1}$ . This suggests a **dynamic programming** algorithm.



# Viterbi Algorithm

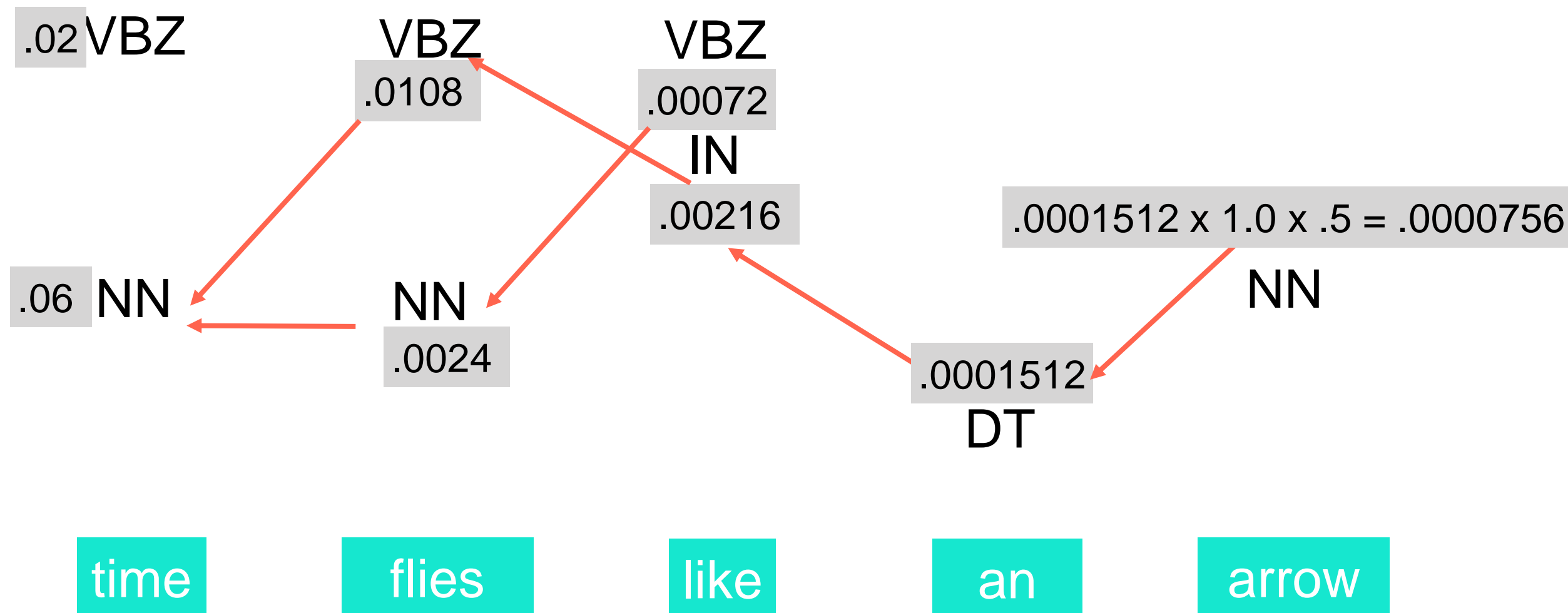
- $P(\text{time} \mid \text{VB}) = 0.2$   
 $P(\text{flies} \mid \text{VB}) = 0.3$   
 $P(\text{like} \mid \text{VB}) = 0.5$
- $P(\text{time} \mid \text{NN}) = 0.3$   
 $P(\text{flies} \mid \text{NN}) = 0.2$   
 $P(\text{arrow} \mid \text{NN}) = 0.5$
- $P(\text{like} \mid \text{P}) = 1.0$
- $P(\text{an} \mid \text{DT}) = 1.0$



- Idea: Because of the Markov assumption, we only need the probabilities for  $X_n$  to compute the probabilities for  $X_{n+1}$ . This suggests a **dynamic programming** algorithm.

# Viterbi Algorithm

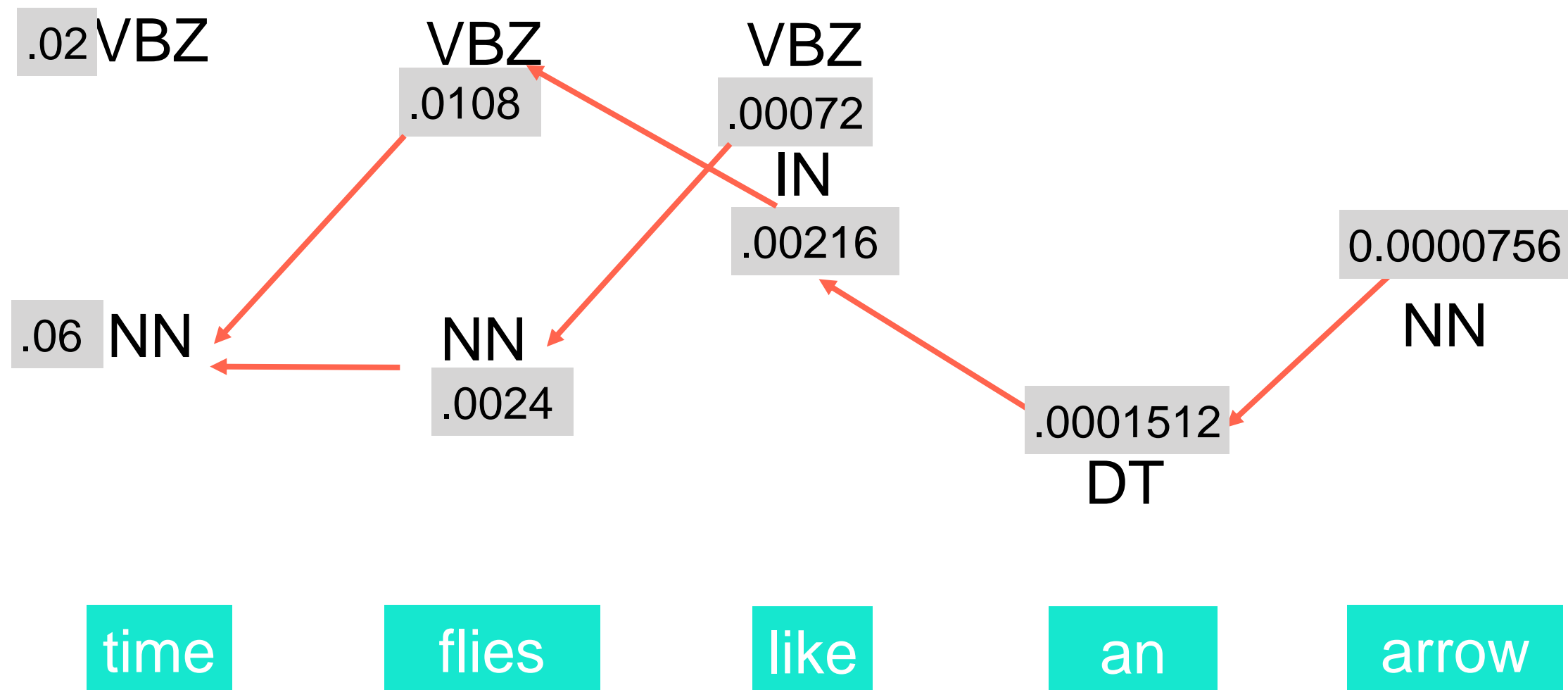
- $P(\text{time} \mid \text{VB}) = 0.2$   
 $P(\text{flies} \mid \text{VB}) = 0.3$   
 $P(\text{like} \mid \text{VB}) = 0.5$
- $P(\text{time} \mid \text{NN}) = 0.3$   
 $P(\text{flies} \mid \text{NN}) = 0.2$   
 $P(\text{arrow} \mid \text{NN}) = 0.5$
- $P(\text{like} \mid \text{P}) = 1.0$
- $P(\text{an} \mid \text{DT}) = 1.0$



- Idea: Because of the Markov assumption, we only need the probabilities for  $X_n$  to compute the probabilities for  $X_{n+1}$ . This suggests a **dynamic programming** algorithm.

# Viterbi Algorithm

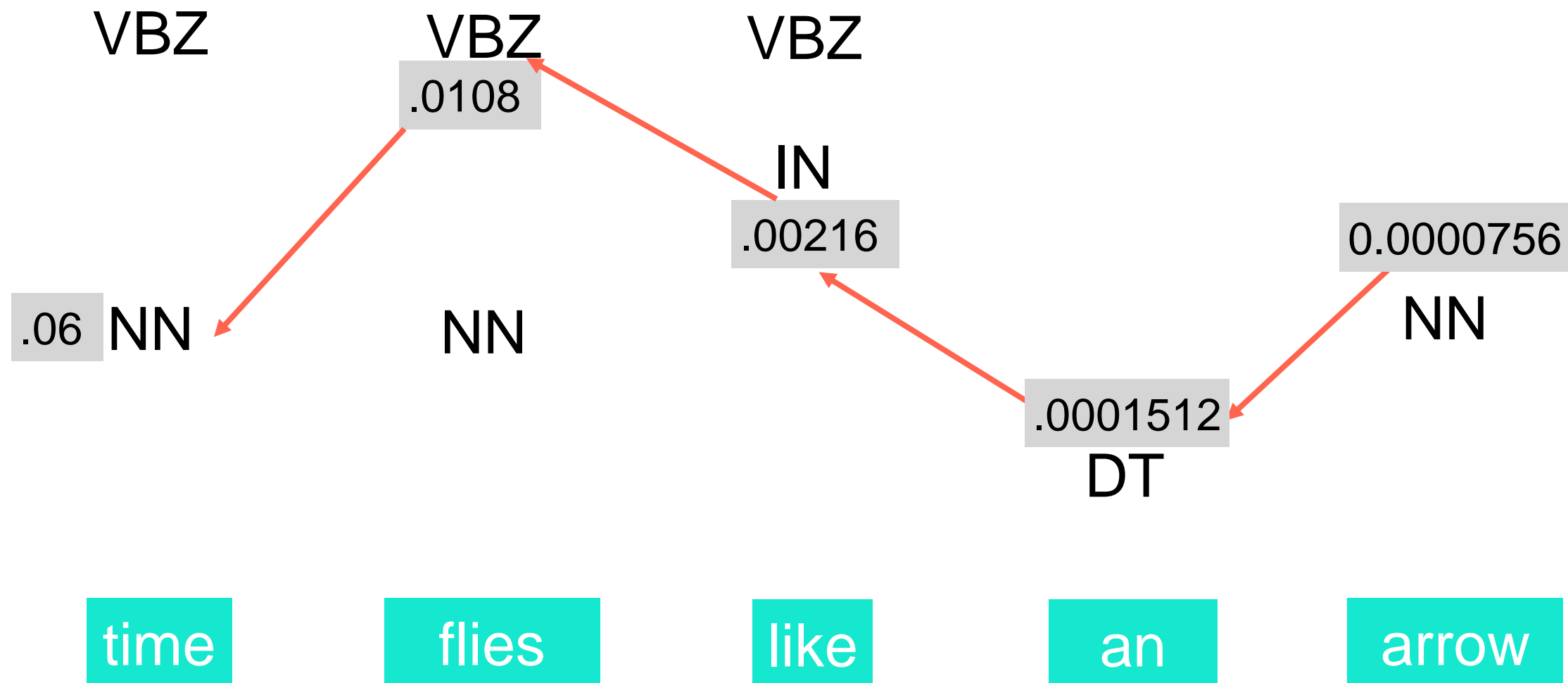
- $P(\text{time} \mid \text{VB}) = 0.2$   
 $P(\text{flies} \mid \text{VB}) = 0.3$   
 $P(\text{like} \mid \text{VB}) = 0.5$
- $P(\text{time} \mid \text{NN}) = 0.3$   
 $P(\text{flies} \mid \text{NN}) = 0.2$   
 $P(\text{arrow} \mid \text{NN}) = 0.5$
- $P(\text{like} \mid \text{P}) = 1.0$
- $P(\text{an} \mid \text{DT}) = 1.0$



- Idea: Because of the Markov assumption, we only need the probabilities for  $X_n$  to compute the probabilities for  $X_{n+1}$ . This suggests a **dynamic programming** algorithm.

# Viterbi Algorithm

- $P(\text{time} \mid \text{VB}) = 0.2$   
 $P(\text{flies} \mid \text{VB}) = 0.3$   
 $P(\text{like} \mid \text{VB}) = 0.5$
- $P(\text{time} \mid \text{NN}) = 0.3$   
 $P(\text{flies} \mid \text{NN}) = 0.2$   
 $P(\text{arrow} \mid \text{NN}) = 0.5$
- $P(\text{like} \mid \text{P}) = 1.0$
- $P(\text{an} \mid \text{DT}) = 1.0$



- Idea: Because of the Markov assumption, we only need the probabilities for  $X_n$  to compute the probabilities for  $X_{n+1}$ . This suggests a **dynamic programming** algorithm.

# Viterbi Algorithm

- **Input:** Sequence of observed words  $w_1, \dots, w_n$
- Create a table  $\pi$ , such that each entry  $\pi[k, t]$  contains the score of the highest-probability sequence ending in tag  $t$  at time  $k$ .
- initialize  $\pi[0, \text{start}] = 1.0$  and  $\pi[0, t] = 0.0$  for all tags  $t \in T$ .
- for  $k = 1$  to  $n$ :
  - for  $t \in T$ :
    - $\pi[k, t] \leftarrow \max_s \pi[k - 1, s] \cdot P(t|s) \cdot P(w_k|t)$ 
      - transition probability
      - emission probability
- **return**  $\max_s \pi[n, s]$

# Trigram Language Model

- Instead of using a unigram context  $P(t_i | t_{i-1})$  use a bigram context  $P(t_i | t_{i-2} t_{i-1})$ 
  - Think of this as having states that represent *pairs of tags*.
- So the HMM probability for a given tag and word sequence is:
$$\prod_{i=1}^n P(t_i | t_{i-2} t_{i-1}) P(w_i | t_i)$$
- Need to handle data sparseness when estimating transition probabilities (for example using backoff or linear interpolation)

# More POS tagging tricks

- It is also often useful in practice to add an end-of-sentence marker (just like we did for n-gram language models).

$$P(t_1, \dots, t_n, w_1, \dots, w_n) = \left[ \prod_{i=1}^n P(t_i | t_{i-2} t_{i-1}) P(w_i | t_i) \right] P(t_{n+1} | t_n)$$

where  $t_{-1} = t_0 = \text{START}$  and  $t_{n+1} = \text{STOP}$ .

- Another useful trick is to replace words with “pseudo words” representing an entire class.
  - For example: replace {"01", "85", "90", ...} with *twoDigitNumber*  
replace {"1985", "2018", ...} with *fourDigitNumber*  
replace {"1", "1.0", "234.3" ...} with *otherNum*  
replace {"IBM", "DNC", ...} with *allCaps etc.*

Using a smoothed trigram HMM model with these tricks, we can build a tagger that is close to the state-of-the art (~97% accuracy on the Penn Treebank).

# HMMs as Language Models

- We can also use an HMM as language models (language generation, MT, ...), i.e. **evaluate**  $P(w_1, \dots, w_n)$  for a given sentence.

What is the advantage over a plain word n-gram model?

- Problem: There are many tag-sequences that could have generated  $w_1, \dots, w_n$ .

$$P(w_1, \dots, w_n, t_1, \dots, t_n) = \prod_{i=1}^n P(t_i | t_{i-1}) P(w_i | t_i)$$

- This is an example of **spurious ambiguity**.

- Need to compute: 
$$P(w_1, \dots, w_n) = \sum_{t_1, \dots, t_n} P(w_1, \dots, w_n, t_1, \dots, t_n)$$
$$= \sum_{t_1, \dots, t_n} \left[ \prod_{i=1}^n P(t_i | t_{i-1}) P(w_i | t_i) \right]$$



# Forward Algorithm

- **Input:** Sequence of observed words  $w_1, \dots, w_n$
- Create a table  $\pi$ , such that each entry  $\pi[k, t]$  contains the score of the highest-probability sequence ending in tag  $t$  at time  $k$ .
- initialize  $\pi[0, \text{start}] = 1.0$  and  $\pi[0, t] = 0.0$  for all tags  $t \in T$ .
- *for*  $k=1$  *to*  $n$ :
  - *for*  $t \in T$ :
    - $\pi[k, t] \leftarrow \sum_s \pi[k-1, s] \cdot P(t|s) \cdot P(w_k|t)$
- **return**  $\sum_s \pi[n, s]$

# Named Entity Recognition as Sequence Labeling

- Use 3 tags:
  - O - outside of named entity
  - I - inside named entity
  - B - first word (beginning) of named entity

...      O                      O                      B                      I                      O                      ...  
identification      of      tetronic      acid      in

- Other encodings are possible (for example, NE-type specific)
- This can also be used for phrase chunking.