

# **Mathematics of POS tagging**

# Argmax Computation

Suppose:

$$x^* = \underset{x}{\operatorname{argmax}} (f(x))$$

Find out value of  $x$  which maximizes  $f(x)$

# Bigram Assumption

Best tag sequence

$$= T^*$$

$$= \operatorname{argmax}_T P(T | W)$$

$$= \operatorname{argmax}_T P(T)P(W | T) \quad (\text{by Bayes Theorem})$$

$$P(T) = P(t_0 = \wedge t_1 t_2 \dots t_{n+1} = .)$$

$$= P(t_0)P(t_1 | t_0)P(t_2 | t_1 t_0)P(t_3 | t_2 t_1 t_0) \dots$$

$$P(t_n | t_{n-1} t_{n-2} \dots t_0)P(t_{n+1} | t_n t_{n-1} \dots t_0)$$

$$= P(t_0)P(t_1 | t_0)P(t_2 | t_1) \dots P(t_n | t_{n-1})P(t_{n+1} | t_n)$$

Bigram Assumption

$$= \prod_{i=0}^{N+1} P(t_i | t_{i-1})$$

# Lexical Probability Assumption

$$P(W|T) = P(w_0|t_0-t_{n+1})P(w_1|w_0t_0-t_{n+1})P(w_2|w_1w_0t_0-t_{n+1}) \dots \\ P(w_n|w_0-w_{n-1}t_0-t_{n+1})P(w_{n+1}|w_0-w_nt_0-t_{n+1})$$

Assumption: A word is determined completely by its tag. This is inspired by speech recognition

$$= P(w_0|t_0)P(w_1|t_1) \dots P(w_{n+1}|t_{n+1})$$

$$= \prod_{i=0}^{n+1} P(w_i|t_i)$$

$$= \prod_{i=0}^{n+1} P(w_i|t_i) \quad (\text{Lexical Probability Assumption})$$

# Best tag sequence

$$T^* = \operatorname{argmax} P(T)P(W|T)$$

$$= \prod_{i=0}^{N+1} P(t_i|t_{i-1}) P(w_i/t_i)$$

# Process

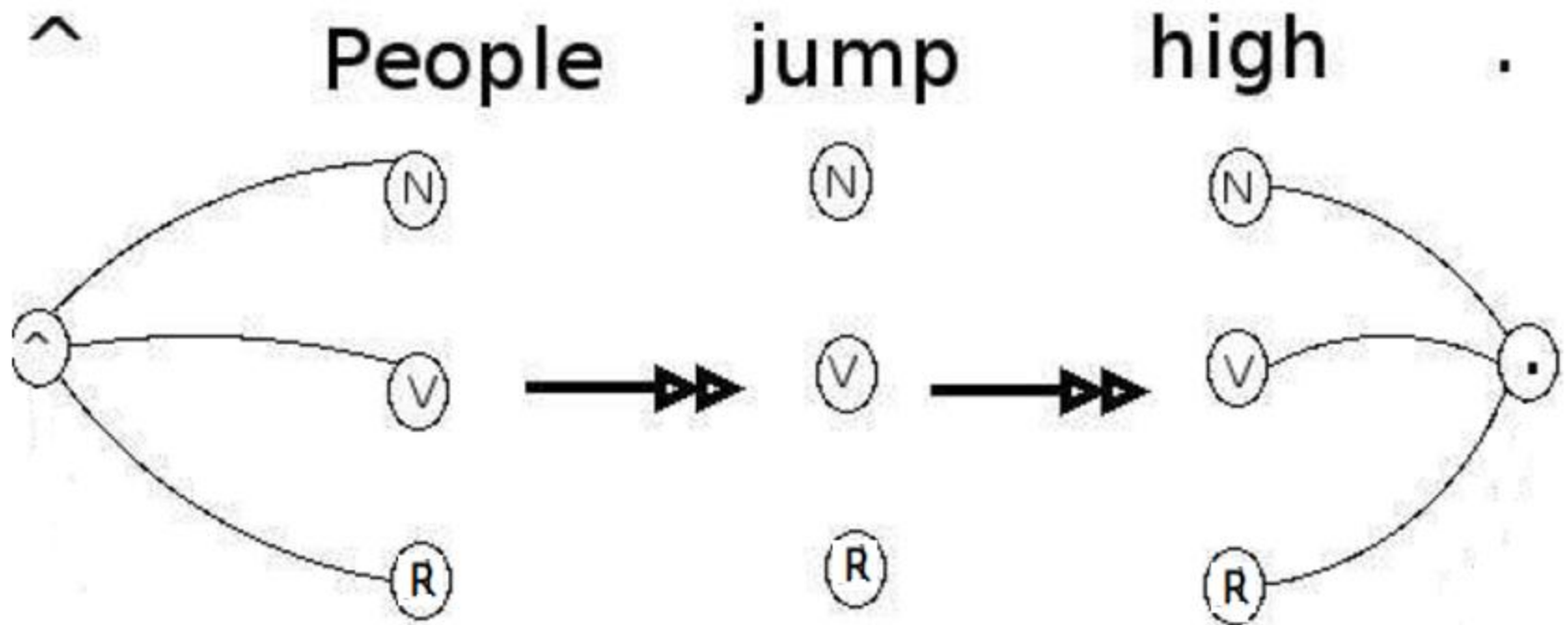
1. List all possible tag for each word in sentence.
2. Choose best suitable tag sequence.

## Example

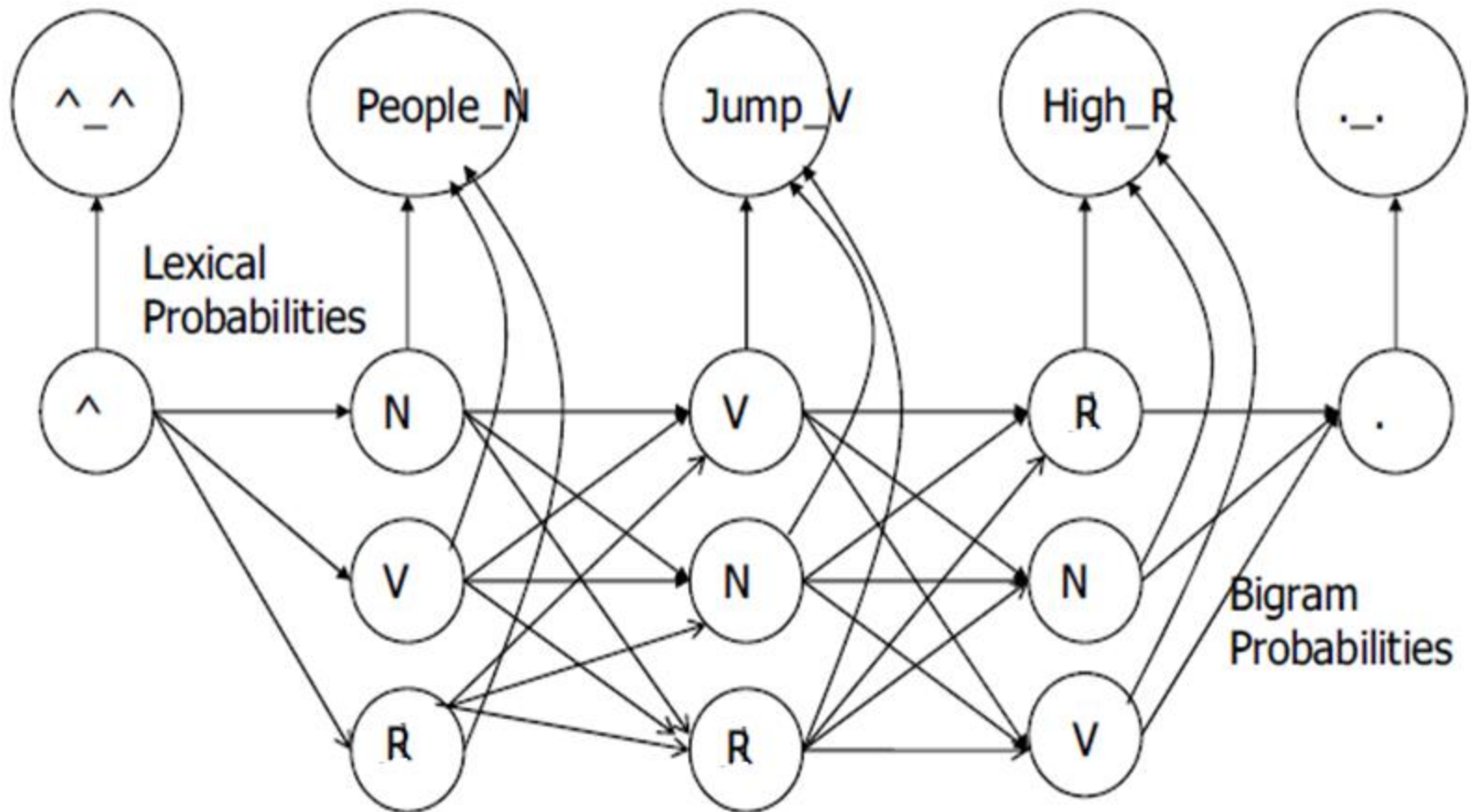
**"People jump high".**

- People : Noun/Verb/Adjective
- jump : Noun/Verb/Adjective
- high : Noun/Verb/Adjective

# Process



# Model



This model is called Generative model.  
Here words are observed from tags as states.  
This is similar to HMM.



# Sequence Labeling

- **Sequence Models:**

- Sequence models are those where there is some sort of dependence through time between the inputs.
- Its job is to assign a label to each unit in a sequence, thus mapping a sequence of observations to a sequence of labels.
- The classical example of a sequence model is the Hidden Markov Model for part-of-speech tagging.

- **Problem:**

- Given a sequence of tokens, infer the most probable sequence of labels(tags) for these tokens

## **Examples :**

- POS
- NER

# HMM

- An HMM is a probabilistic sequence model
  - Given a sequence of units (words, letters, morphemes, sentences, whatever), they compute a probability distribution over possible sequences of labels and choose the best label sequence.

# POS tagging with HMM

## Notations:

$W = w_1, w_2, w_3, \dots, w_n$  -sequence of words(Input)

$T = t_1, t_2, t_3, \dots, t_n$  -sequence of tags(labels)

- It is a process of finding the sequence of tags which is most likely to have generated a given word sequence.

# Markov Processes

- A Markov process is a random process in which the future is independent of the past, given the present
- A Markov process satisfies the **Markov property**
- The markov property is characterized as "**memorylessness**"

# Markov Processes

## Based on Two Properties

1. **Limited Horizon:** Given previous  $t$  states, a state  $i$  is independent of preceding  $0$  to  $t-(k+1)$  states

$$P(X_t=i/X_{t-1}, X_{t-2}, \dots, X_0) = P(X_t=i/X_{t-1}, X_{t-2}, \dots, X_{t-k})$$

- Order  $k$  Markov process

2. **Time invariance (shown for  $k=1$ ):** dependence of  $t$ -state on  $t-1$  state is same everywhere.

$$P(X_t=i/X_{t-1}=j) = P(X_1=i/X_0=j) \dots = P(X_n=i/X_{n-1}=j)$$

# Formal definition of HMM

**An HMM is specified by:**

1. The set  $S = s_1, s_2, s_3, \dots, s_N$  of hidden state
2. The start state  $s_0$
3. The matrix  $A$  of transition probabilities:  
$$a_{ij} = p(s_j/s_i)$$
4. The set  $O$  of possible visible outcomes
5. The matrix  $B$  of output probabilities:  
$$b_{ik} = p(o_k/s_i)$$

# HMM

- In the process, “**tags**” are the “**hidden states**” which produced the “**observable output**” i.e., “**words**”.
- Find out the most probable sequence of tags for the word sequence:

**Best tag sequence:  $T^* = \operatorname{argmax}_T \{P(T/W)\}$**

$$T^* = \operatorname{argmax}_T (P(T) * P(W/T))$$

$$T^* = \operatorname{argmax}_T \prod_i P(w_i | w_1 \dots w_{i-1}, t_1 \dots t_i) P(t_i | t_1 \dots t_{i-1})$$

# HMM

Best tag sequence:  $T^* = \operatorname{argmax}_T \{P(T/W)\}$   
 $T^* = \operatorname{argmax}_T (P(T) * P(W/T))$

**Simplify, using 2 assumptions:**

- 1- **Bigram assumption** : The probability of a tag depends on the previous one (bigram model-  $P(t_i/t_{i-1})$ ).
- 2- **Lexical Generation probability assumption** : The probability of a word appearing depends only on its own POS tag, ( $P(w_i/t_i)$ ).

$$T^* = \operatorname{argmax}_T (\prod_{i=1 \dots T} \operatorname{PROB}(t_i | t_{i-1}) * \operatorname{PROB}(w_i | t_i))$$



# Computing the probability values

## 1. Tag Transition probabilities $p(t_i/t_{i-1})$

$$P(t_i|t_{i-1}) = \frac{C(t_{i-1}, t_i)}{C(t_{i-1})}$$

$$P(NN|DT) = \frac{C(DT, NN)}{C(DT)} = \frac{56,509}{116,454} = 0.49$$

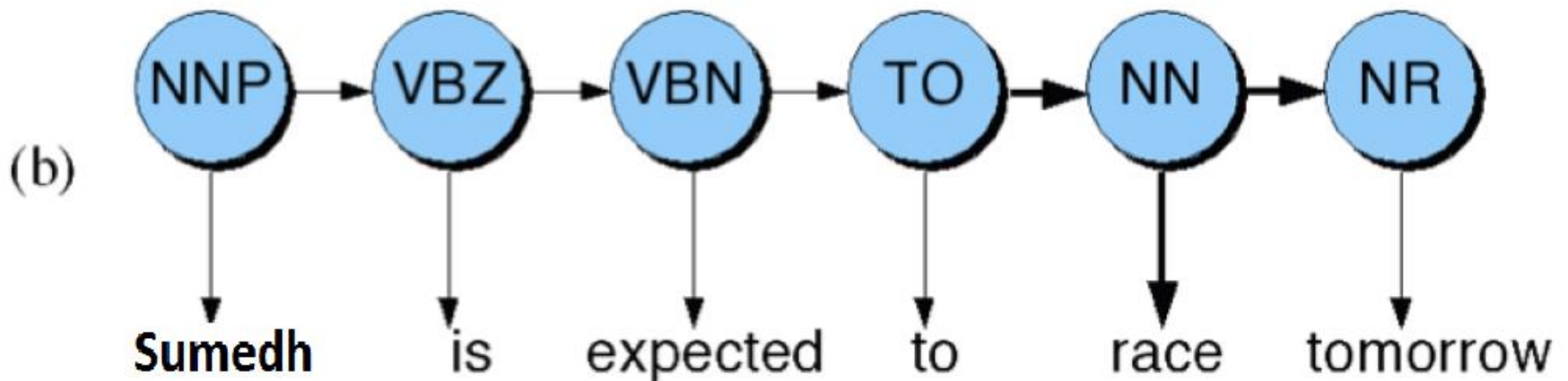
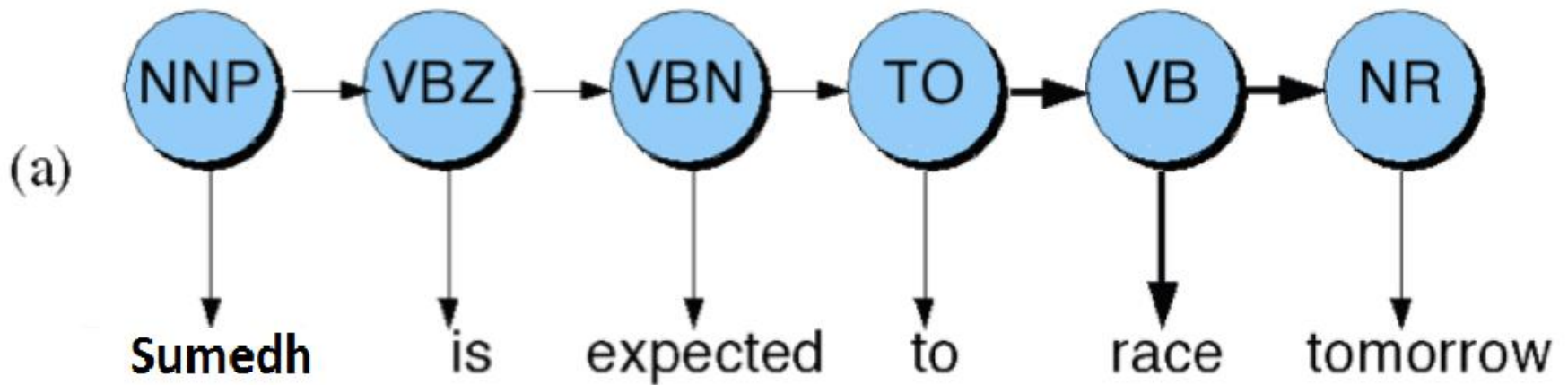
# Computing the probability values

## 2. Word Likelihood probabilities $p(w_i/t_i)$

$$P(w_i|t_i) = \frac{C(t_i, w_i)}{C(t_i)}$$

$$P(is|VBZ) = \frac{C(VBZ, is)}{C(VBZ)} = \frac{10,073}{21,627} = 0.47$$

# Disambiguating race



# Disambiguating race

Difference in probability due to

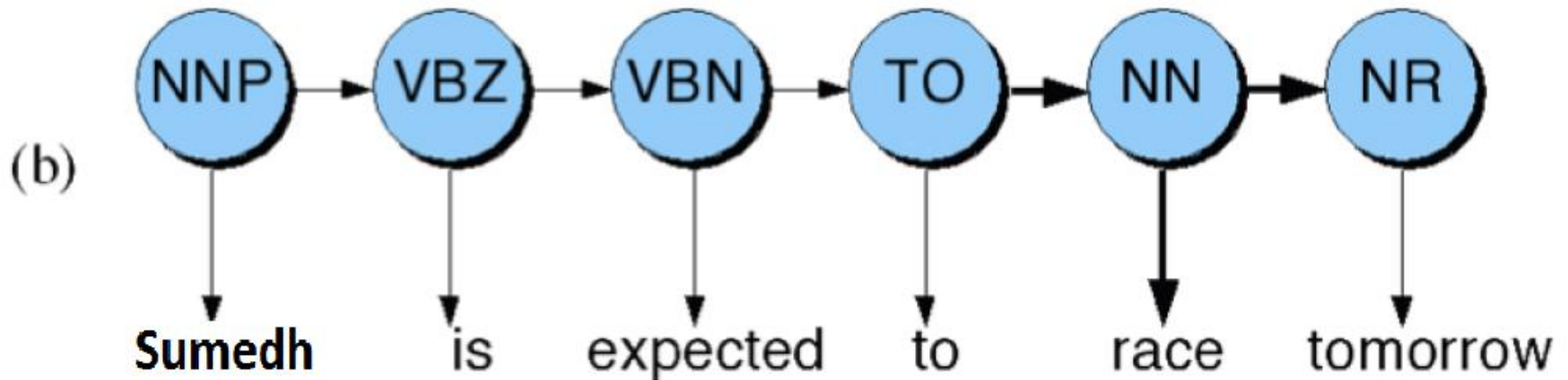
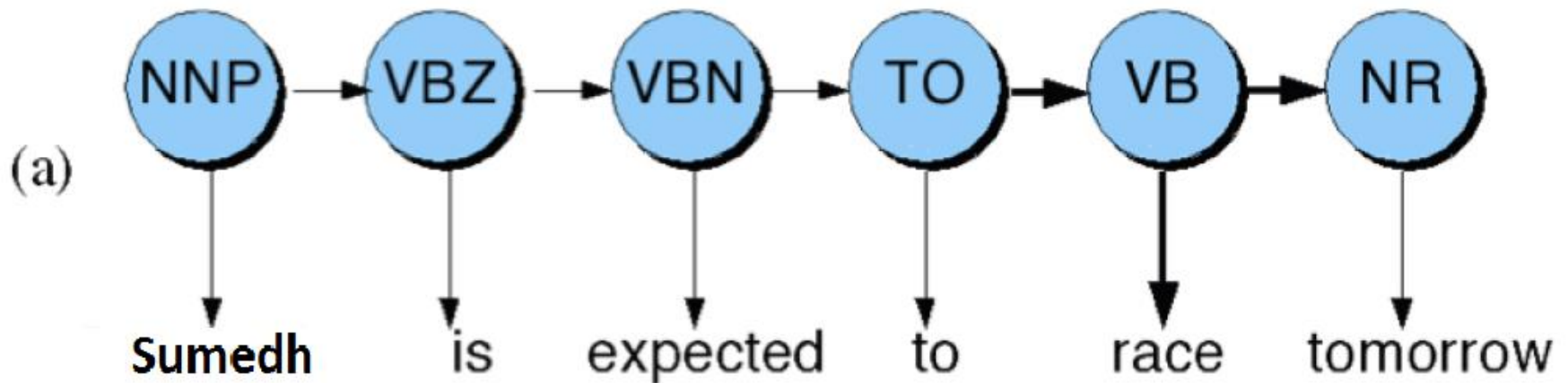
- $P(\text{VB}/\text{TO})$  vs.  $P(\text{NN}/\text{TO})$
- $P(\text{race}/\text{VB})$  vs.  $P(\text{race}/\text{NN})$
- $P(\text{NR}/\text{VB})$  vs.  $P(\text{NR}/\text{NN})$

# Disambiguating race

After computing the probabilities:

- $P(\text{NN}/\text{TO}) * P(\text{NR}/\text{NN}) * P(\text{race}/\text{NN}) =$   
 $0.0047 * 0.0012 * 0.00057 = 0.00000000032$
- $P(\text{VB}/\text{TO}) * P(\text{NR}/\text{VB}) * P(\text{race}/\text{VB}) = 0.83 * 0.0027$   
 $* 0.00012 = 0.00000027$

# What is this model?



**This is a Hidden Markov Model**

# Hidden Markov Models

- Tag Transition probabilities -  $p(t_i/t_{i-1})$
- Word Likelihood probabilities (emissions) -  $p(w_i/t_i)$

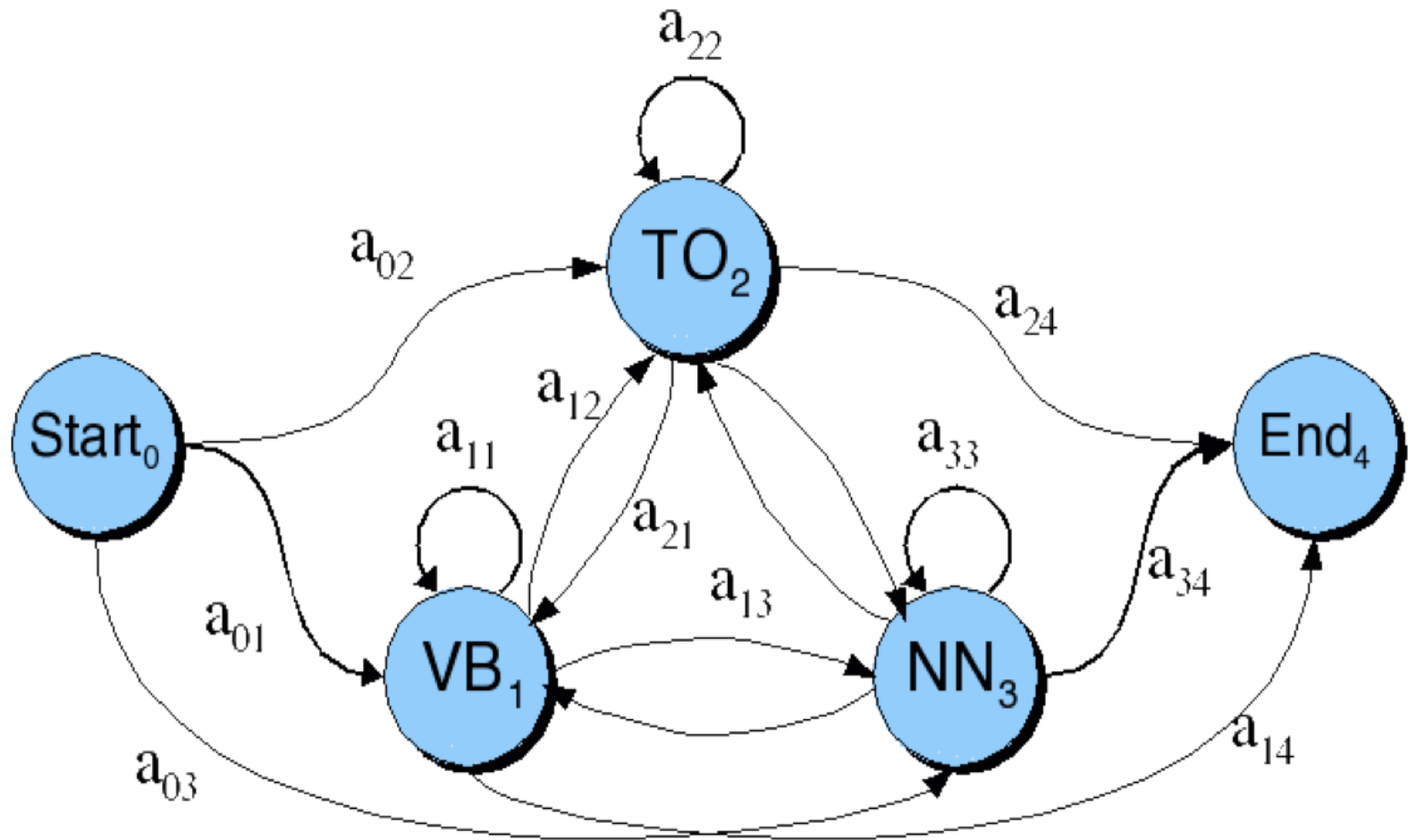
# Hidden Markov Models (HMMs)

## Elements of an HMM model:

- A set of states (here: the tags)
- An output alphabet (here: words)
- Initial state (here: beginning of sentence)
- State transition probabilities (here  $p(t_i/t_{i-1})$ )
- Symbol emission probabilities (here  $p(w_i/t_i)$ )



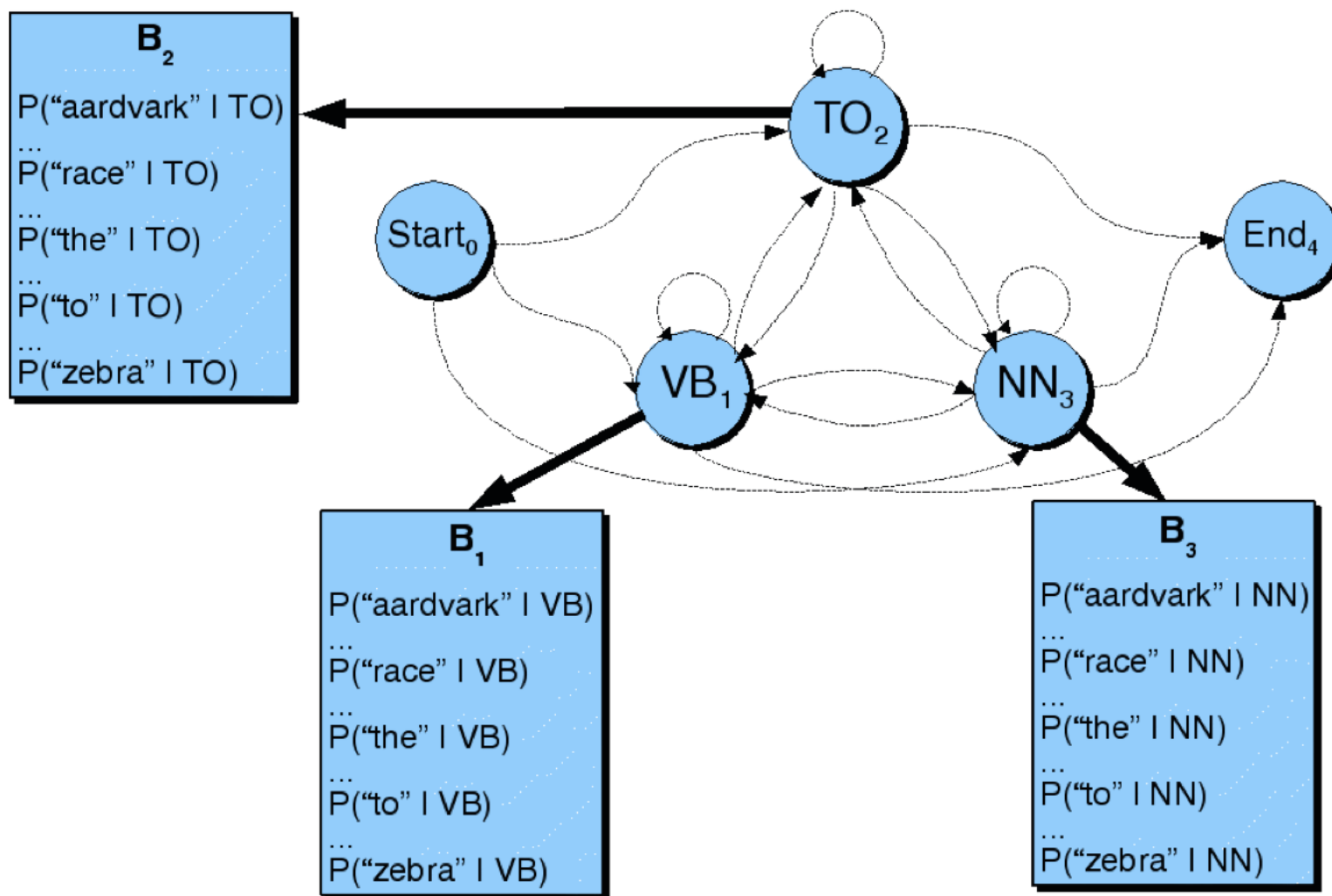
# Graphical Representation



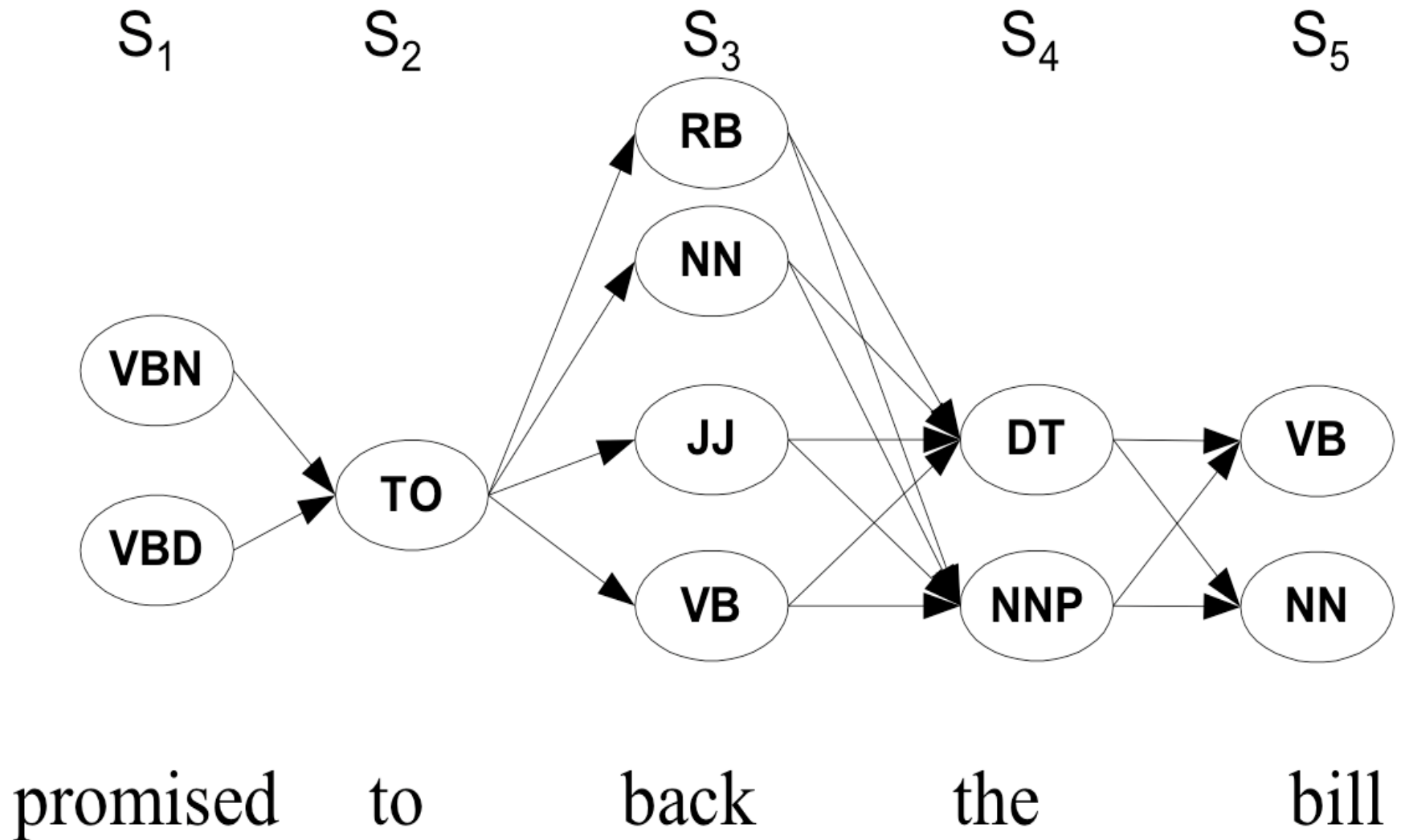
- Edges are labeled with the state transition probabilities:  $P(t_i/t_{i-1})$

# Graphical Representation

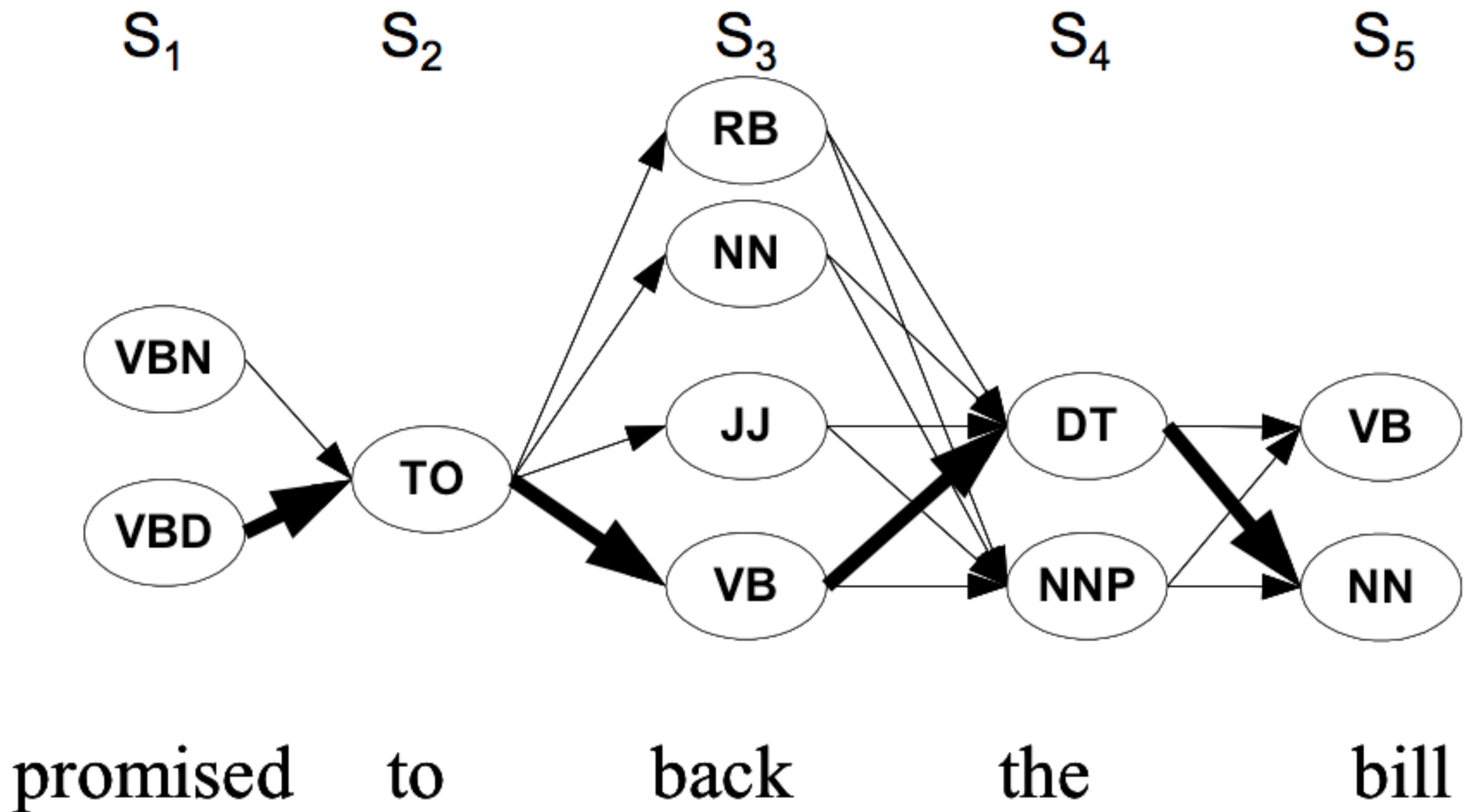
- At each state we emit a word:  $P(w_i/t_i)$



# Walking through the states: best path



# Walking through the states: best path



# Problem

- Consider tag set containing only four tags (ART, N, V, P) then with our independence assumption (the tag  $t_i$  depends only on the tag of the (i-1)th word) we can use a markov chain to capture this bigram probabilities as shown in the figure below. Assume that any bigram probability not in Figure 1 has a value of .0001.

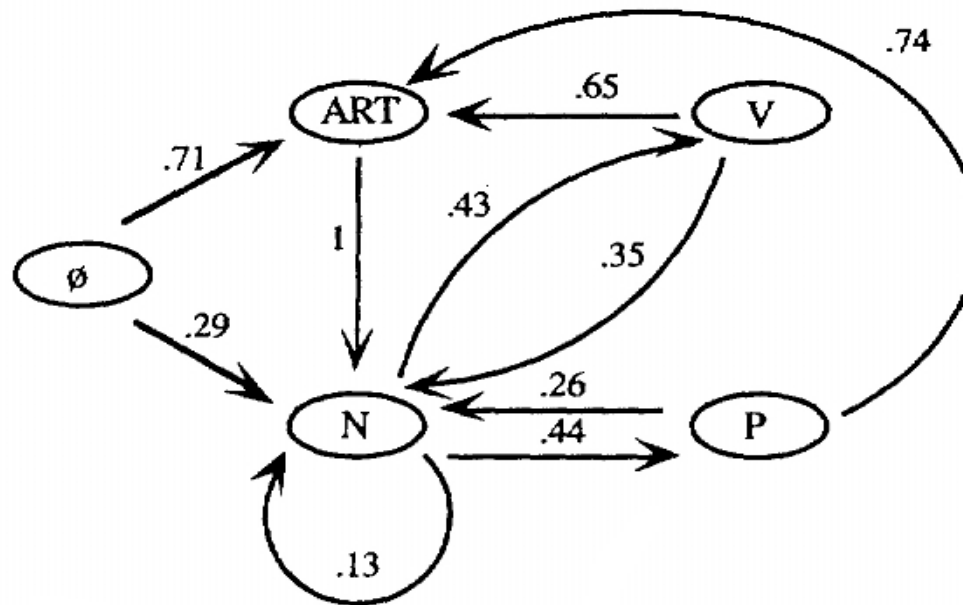


Figure 1

With the help of the above diagram the probability of a sequence can be calculated by multiplying the probabilities along the path.

e.g.

$$P(N, V, ART, N) = 0.29 * 0.43 * 0.65 * 1$$

Such a network representation is called a **Markov model**.

# Viterbi Algorithm

- Uses a dynamic programming approach where we sweep forward through the words one at a time to find the most likely sequence ending in each tag.
- For example, for the word sequence : “**Flies like a flower**”,
  - First, find the four best sequences for the first word each ending in different tag from start.
  - Second, find the four best sequences for the two words "*flies like*": the best ending with "*like*" as a V, the best ending with "*like*" as an N, the best ending with "*like*" as a P, and the best ending with "*like*" as an ART.
  - Using this information we will then find the four best sequences for "*flies like a*", each one ending in a different tag.
  - This process is repeated until all the words in the sentence are considered.
- This algorithm is called the **Viterbi** algorithm.

# Viterbi Algorithm

## Initialization Step

For  $i = 1$  to  $N$  do                      //  $N$  is no of lexical categories and  $T$  is no of words

$SEQSCORE(i,1) = P(W_1 | C_i) * P(C_i | \emptyset)$

$BACKPTR(i,1) = 0$

## Iteration Step

For  $t = 2$  to  $T$  do

For  $i = 1$  to  $N$

$SEQSCORE(i,t) = \max_{j=1,N} (SEQSCORE(j,t-1) * P(C_i | C_j)) * P(w_t | C_i)$

$BACKPTR(i,t) = \text{index of } j \text{ that gave max above}$

## Sequence Identification Step

$C(T) = i \text{ that maximizes } SEQSCORE(i,T)$

For  $i = t-1$  to  $1$  do

$C(i) = BACKPTR(C(i+1), i+1)$  //              Back trace to find the sequence

- The array  $SEQSCORE(i, j)$  stores the probability for the optimal sequence up to word  $j$  ending in the category  $i$ .



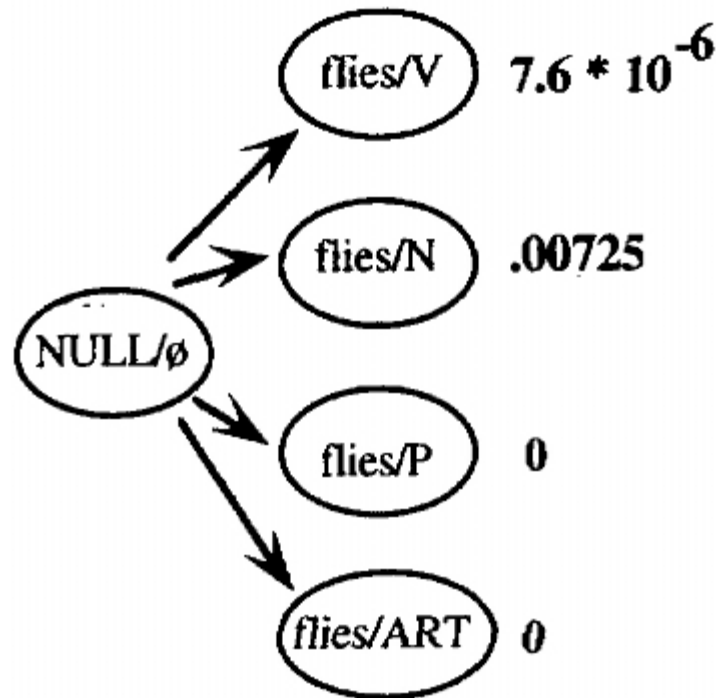
- Consider that after analyzing the corpus we have estimated the bigram probabilities as shown in Figure 1 and lexical-generation probabilities as shown in the table below:-

PROB ( <i>the</i>   ART)	.54	PROB ( <i>a</i>   ART)	.360
PROB ( <i>flies</i>   N)	.025	PROB ( <i>a</i>   N)	.001
PROB ( <i>flies</i>   V)	.076	PROB ( <i>flower</i>   N)	.063
PROB ( <i>like</i>   V)	0.1	PROB ( <i>flower</i>   V)	.05
PROB ( <i>like</i>   P)	0.068	PROB ( <i>birds</i>   N)	.076
PROB ( <i>like</i>   N)	0.012		

Now, using Viterbi Algorithm, the first row is set in the initialization phase using the formula:-

$$\text{SEQSCORE}(i,1) = P(\text{Flies} | C_i) * P(C_i | \emptyset)$$

- The result of the first step of the algorithm, the probability of "flies" in each category has been computed, and is shown below:



$$\text{SEQSCORE}(i,1) = P(\text{Flies} | C_i) * P(C_i | \emptyset)$$

$$P(\text{Flies} | V) * P(V | \emptyset) = 0.076 * 0.0001 = 7.6 * 10^{-6}$$

$$P(\text{Flies} | N) * P(N | \emptyset) = 0.025 * 0.29 = 0.00725$$

$$P(\text{Flies} | P) * P(P | \emptyset) = 0 * 0.0001 = 0$$

$$P(\text{Flies} | \text{ART}) * P(\text{ART} | \emptyset) = 0 * 0.71 = 0$$

The second phase of the algorithm extends the sequences one word at a time, keeping track of the best sequence found so far to each category.

Probability of the state **like/V** is computed as:

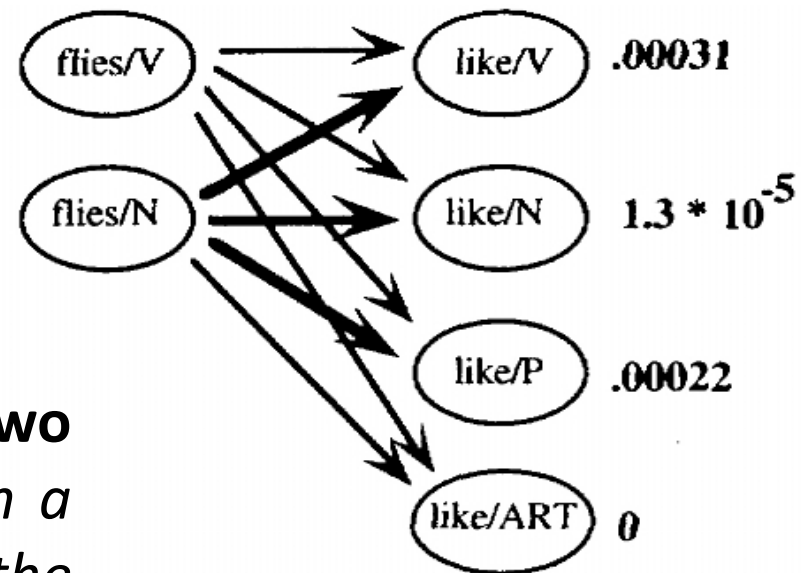
$$\begin{aligned}\text{PROB}(\text{like}/V) &= \text{MAX}(\text{PROB}(\text{flies}/N) * \text{PROB}(V | N), \\ &\quad \text{PROB}(\text{flies}/V) * \text{PROB}(V | V)) * \text{PROB}(\text{like}/V) \\ &= \text{MAX}(.00725 * .43, \quad 7.6 * 10^{-6} * .0001) * 0.1 \\ &= 3.12 * 10^{-4}\end{aligned}$$

$$\text{PROB}(\text{like}/N) = 1.3 * 10^{-5}$$

$$\text{PROB}(\text{like}/P) = 2.2 * 10^{-5}$$

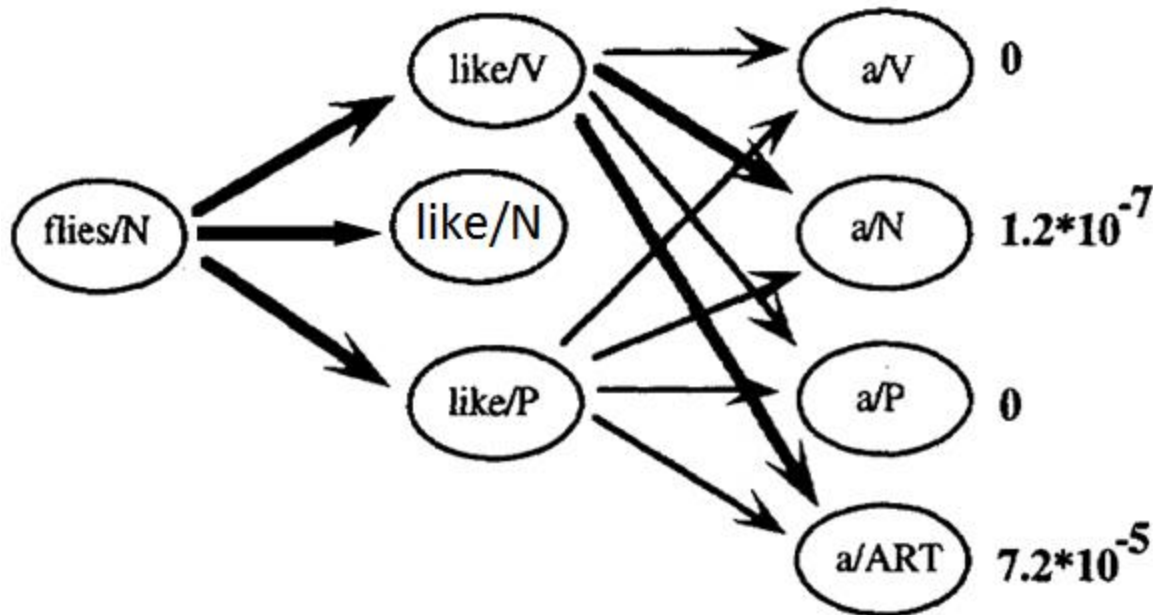
$$\text{PROB}(\text{like}/\text{ART}) = 0$$

The most likely **sequence of length two** generating "*Flies like*" and ending in a *V* has a score of  $3.12 * 10^{-4}$  (and is the sequence *N V*)



Result after first iteration

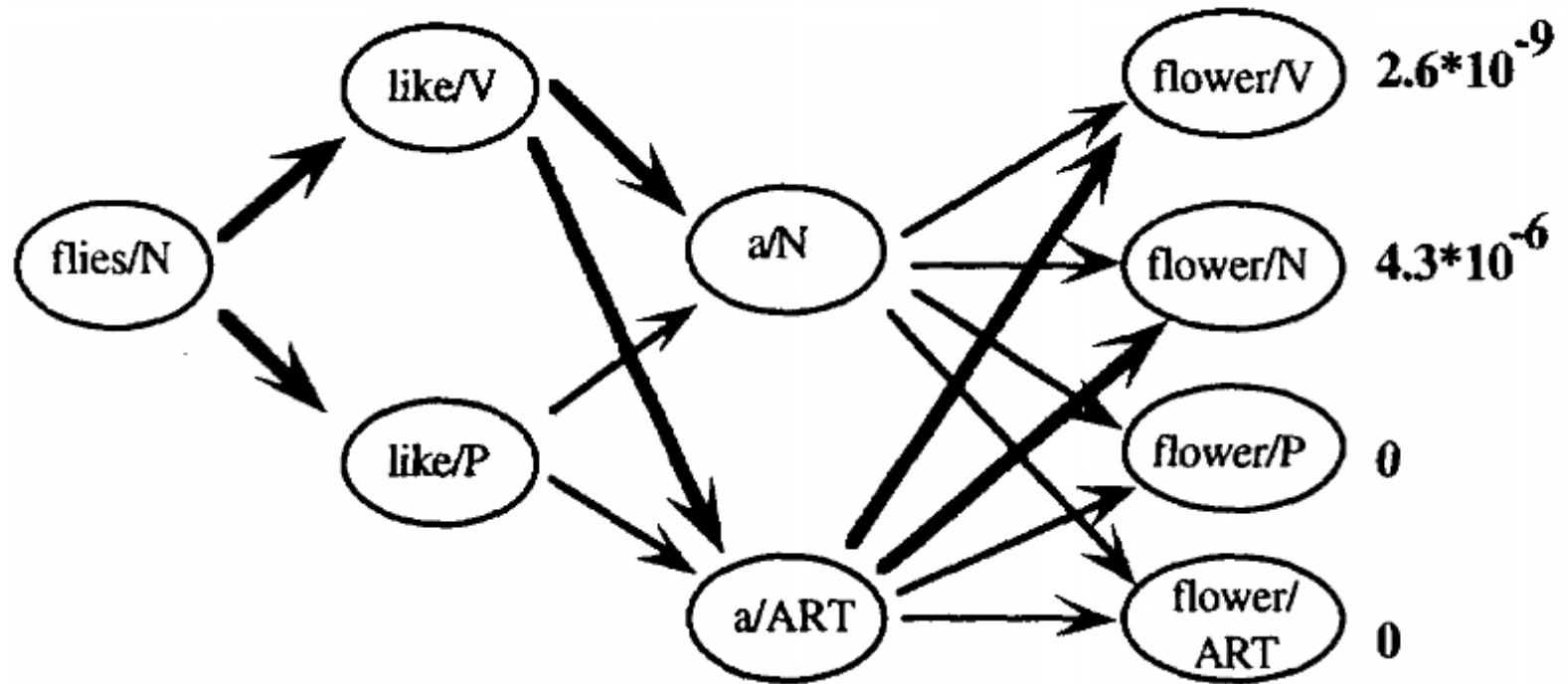
- The computation continues in the same manner until each word has been processed.



Results after second iteration

$$\begin{aligned}
 \text{PROB (a/ART)} &= \text{MAX (PROB (like/V) * PROB (ART | V),} \\
 &\quad \text{PROB(like/P) * PROB (ART|P), PROB(like/N) *} \\
 &\quad \text{PROB (ART|N) ) * PROB (a/ART)} \\
 &= \text{MAX}(0.00031*0.65, 0.00022*0.74, 1.3 \times 10^{-5} * 0.0001) \times 0.360 \\
 &= \mathbf{7.254 \times 10^{-5}}
 \end{aligned}$$

The computation continues in the same manner until each word has been processed.



Results after third iteration [[Allen, 1995](#)]

The highest probability sequence ends in state flower/N.

It is easy to back trace from here to get the complete sequence N, V, ART, N.

# Results and Conclusion

- **Results-** HMM tagger gives an accuracy of about 96% when trained and tested on the Wall Street Journal corpus.
- **Conclusion-**
  - i) HMM tagger is a generative probabilistic model for Part of Speech Tagging which does not assign tags to individual words but selects the best tag sequence for the entire sentence.
  - ii) The use of Viterbi algorithm speeds up the tagging process by reducing the number of computations in each iteration.