

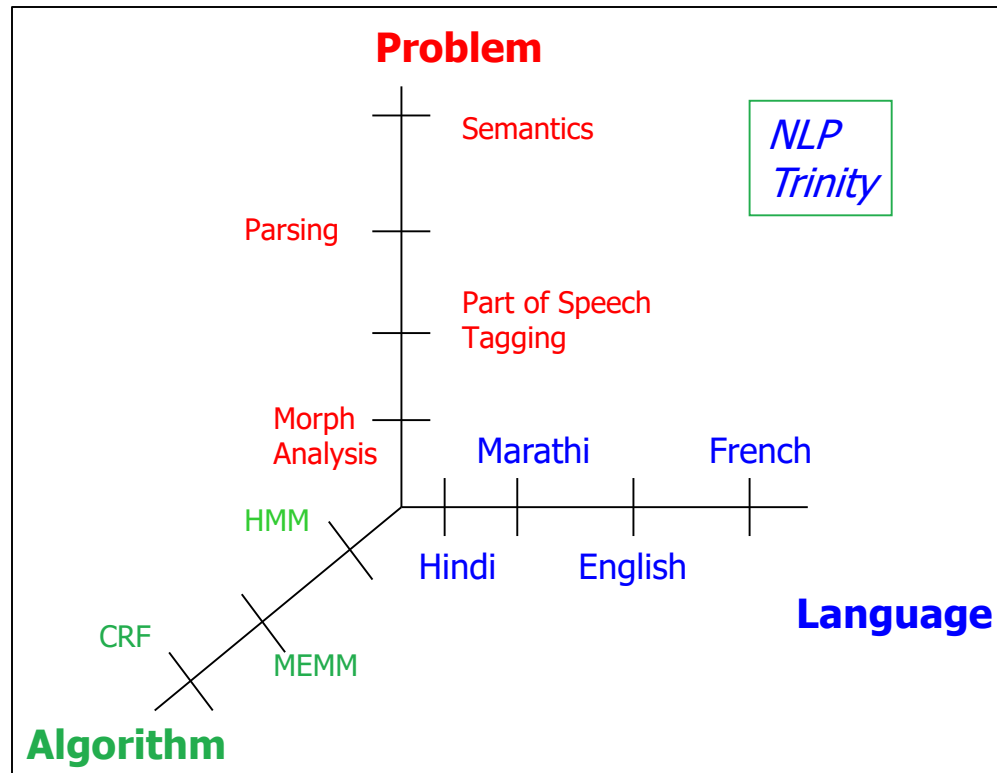
# Part of Speech Tagging

Asif Ekbal  
CSE Dept.,  
IIT Patna

# Part of Speech Tagging

With Hidden Markov Model

# NLP Trinity



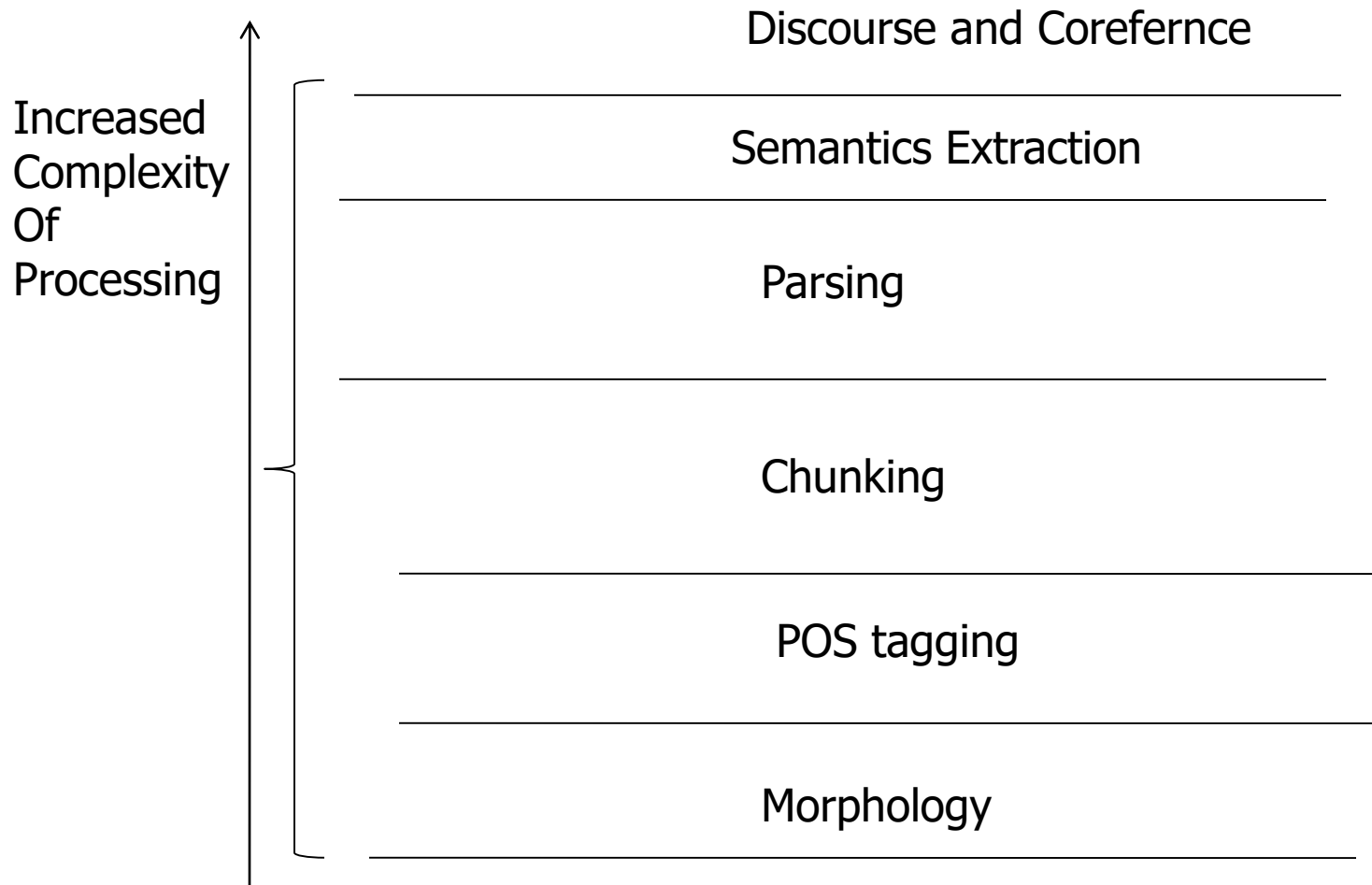
# Part of Speech Tagging

- POS Tagging: attaches to each word in a sentence a part of speech tag from a given set of tags called the **Tag-Set**
- Standard Tag-set : Penn Treebank (for English)- *36 tags*

# Example

- “\_” The\_DT mechanisms\_NNS that\_WDT make\_VBP traditional\_JJ hardware\_NN are\_VBP really\_RB being\_VBG obsoleted\_VBN by\_IN microprocessor-based\_JJ machines\_NNS ,\_, “\_” said\_VBD Mr.\_NNP Benton\_NNP .\_.

# Where does POS tagging fit in



# Example to illustrate complexity of POS tagging

# POS tagging is disambiguation

*N (noun), V (verb), J (adjective), R (adverb) and F (other, i.e., function words).*

***That\_F former\_J Sri\_Lanka\_N skipper\_N and\_F ace\_J batsman\_N Aravinda\_De\_Silva\_N is\_F a\_F man\_N of\_F few\_J words\_N was\_F very\_R much\_R evident\_J on\_F Wednesday\_N when\_F the\_F legendary\_J batsman\_N,\_F who\_F has\_V always\_R let\_V his\_N bat\_N talk\_V,\_F struggled\_V to\_F answer\_V a\_F barrage\_N of\_F questions\_N at\_F a\_F function\_N to\_F promote\_V the\_F cricket\_N league\_N in\_F the\_F city\_N.\_F***



# POS disambiguation

- ***That\_F/N/J*** ('that' can be complementizer (can be put under 'F'), demonstrative (can be put under 'J') or pronoun (can be put under 'N'))
- ***former\_J***
- ***Sri\_N/J Lanka\_N/J*** (Sri Lanka together qualify the skipper)
- ***skipper\_N/V*** ('skipper' can be a verb too)
- ***and\_F***
- ***ace\_J/N*** ('ace' can be both J and N; "Nadal served an ace")
- ***batsman\_N/J*** ('batsman' can be J as it qualifies Aravinda De Silva)
- ***Aravinda\_N De\_N Silva\_N is\_F a\_F***
- ***man\_N/V*** ('man' can verb too as in 'man the boat')
- ***of\_F few\_J***
- ***words\_N/V*** ('words' can be verb too, as in 'he words is speeches beautifully')

# Behaviour of “*That*”

- That
  - *That man is known by the company he keeps.* (Demonstrative)
  - *Man that is known by the company he keeps, gets a good job.* (Pronoun)
  - *That man is known by the company he keeps, is a proverb.* (Complementation)
- Chaotic systems: Systems where a small perturbation in input causes a large change in output

# POS disambiguation

- ***was\_F very\_R much\_R evident\_J on\_F Wednesday\_N***
- ***when\_F/N*** ('when' can be a relative pronoun (put under 'N) as in 'I know the time when he comes')
- ***the\_F legendary\_J batsman\_N***
- ***who\_F/N***
- ***has\_V always\_R let\_V his\_N***
- ***bat\_N/V***
- ***talk\_V/N***
- ***struggle\_V /N***
- ***answer\_V/N***
- ***barrage\_N/V***
- ***question\_N/V***
- ***function\_N/V***
- ***promote\_V cricket\_N league\_N city\_N***

# Mathematics of POS tagging

# Argmax computation (1/2)

Best tag sequence

$$= T^*$$

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

$$= \operatorname{argmax} P(T)P(W|T) \quad (\text{by Baye's Theorem})$$

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

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

$$P(t_n|t_{n-1} t_{n-2} \dots t_0)P(t_{n+1}|t_n t_{n-1} \dots t_0)$$
$$= P(t_0)P(t_1|t_0)P(t_2|t_1) \dots P(t_n|t_{n-1})P(t_{n+1}|t_n)$$

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

Bigram Assumption

# Argmax computation (2/2)

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

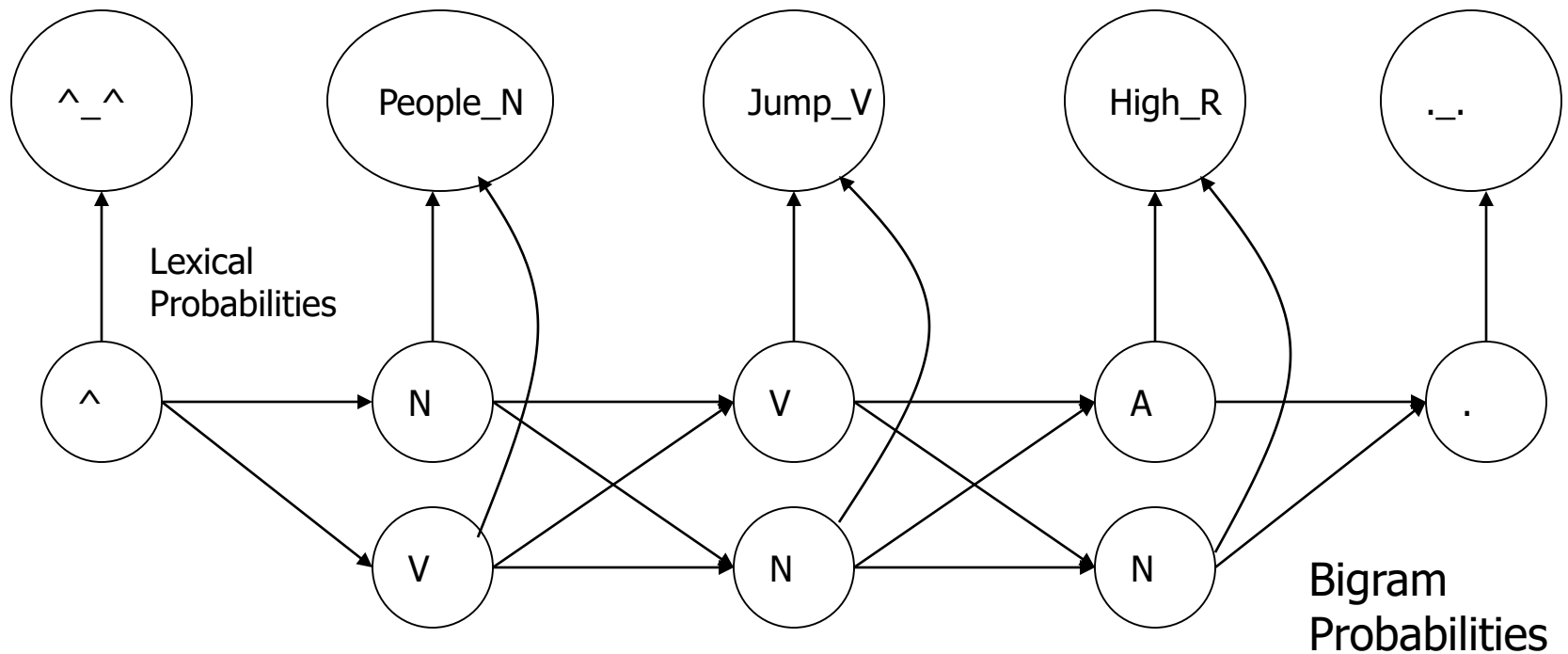
Assumption: A word is determined completely by its tag. This is inspired by speech recognition

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

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

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

# Generative Model



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

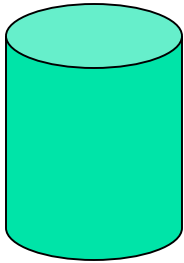
# *Typical POS tag steps*

- Implementation of Viterbi – Unigram, Bigram.
- Evaluation
- Per POS Accuracy
- Confusion Matrix



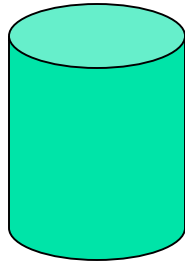
# A Motivating Example

Colored Ball choosing



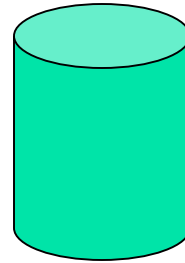
Urn 1

# of Red = 30  
# of Green = 50  
# of Blue = 20



Urn 2

# of Red = 10  
# of Green = 40  
# of Blue = 50



Urn 3

# of Red = 60  
# of Green = 10  
# of Blue = 30

# Example (contd.)

Given :

	$U_1$	$U_2$	$U_3$
$U_1$	0.1	0.4	0.5
$U_2$	0.6	0.2	0.2
$U_3$	0.3	0.4	0.3

Transition probability table

and

	R	G	B
$U_1$	0.3	0.5	0.2
$U_2$	0.1	0.4	0.5
$U_3$	0.6	0.1	0.3

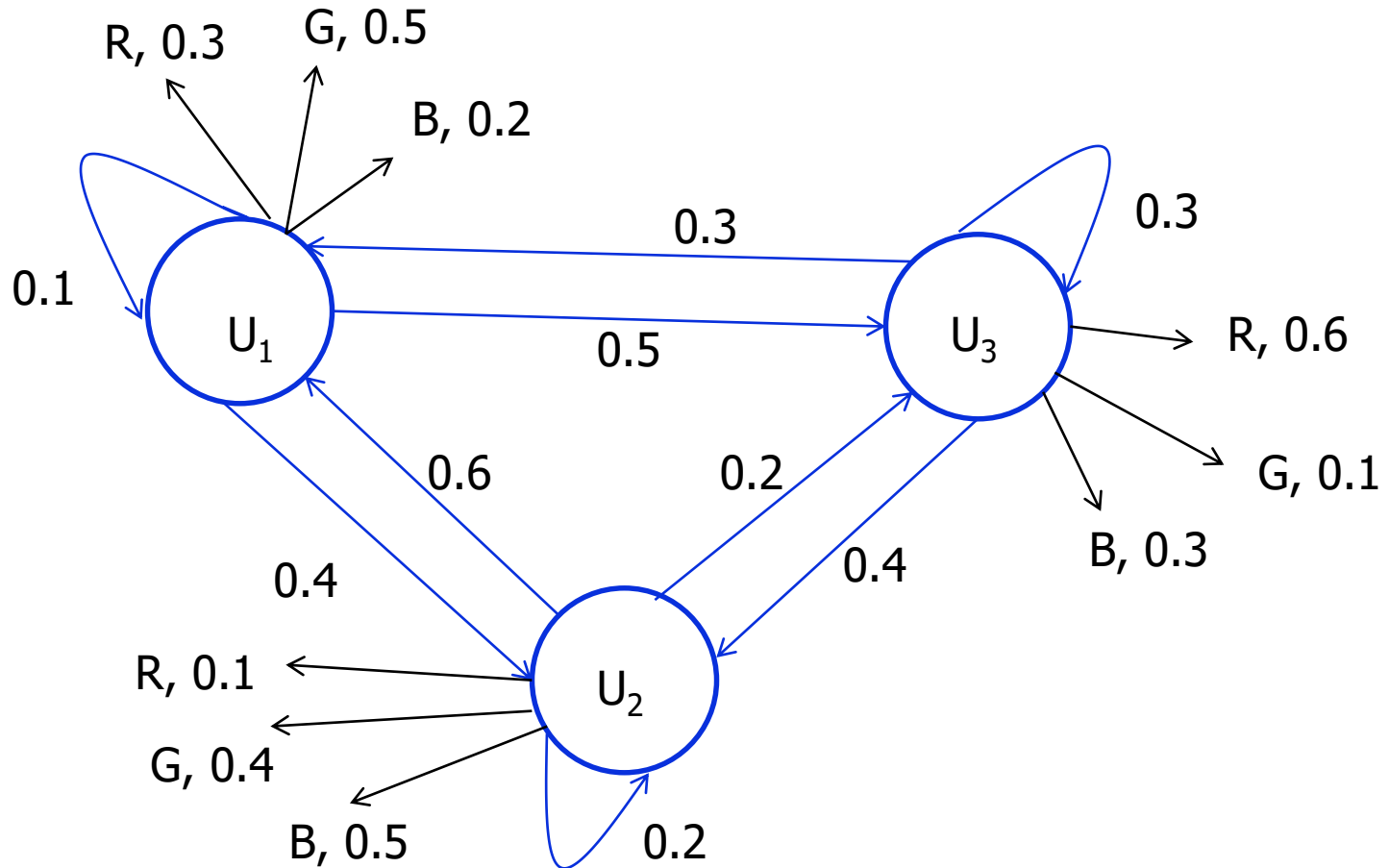
Emission probability table

Observation : RRGGBRGR

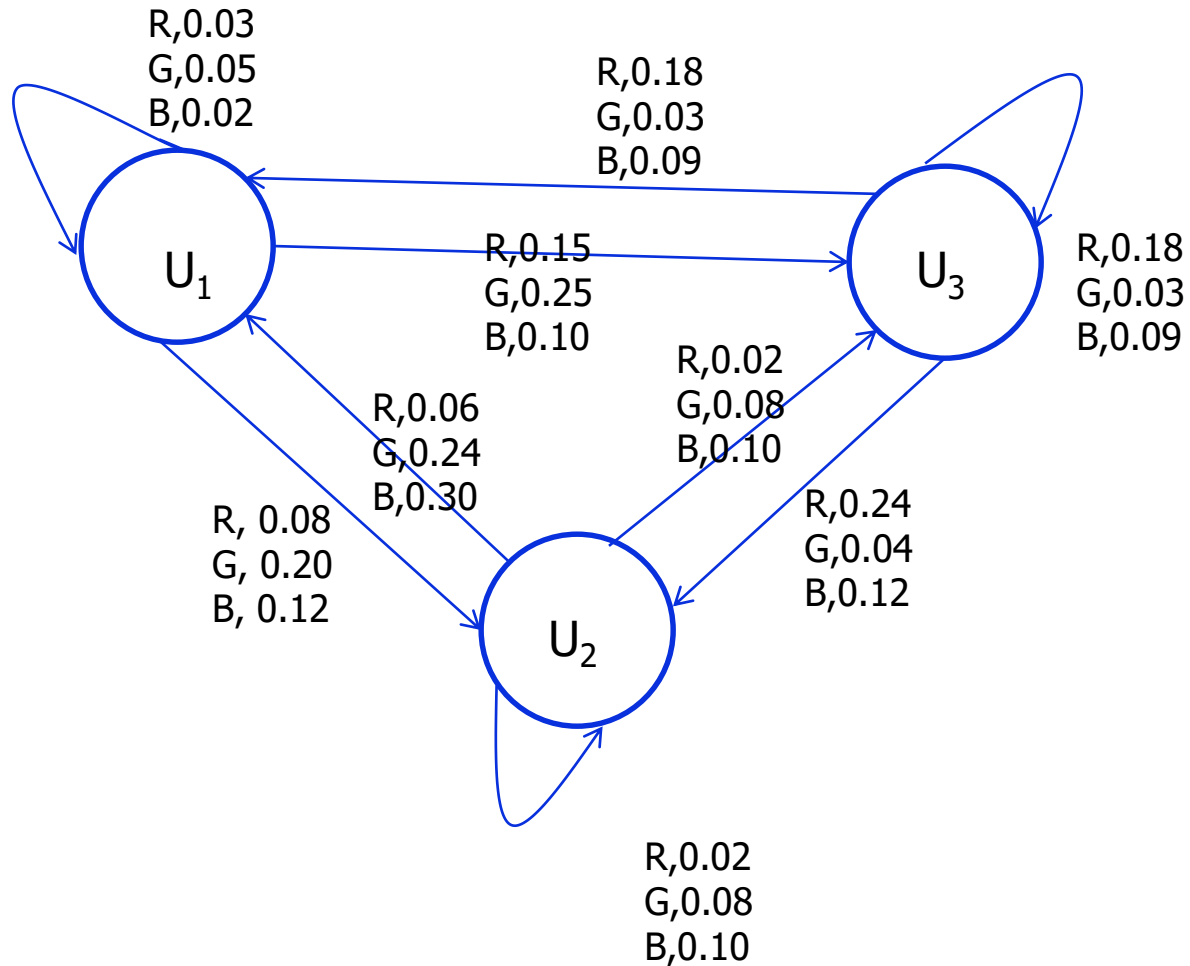
State Sequence : ??

Not so Easily Computable.

# Diagrammatic representation (1/2)



# Diagrammatic representation (2/2)



# Classic problems with respect to HMM

1. Given the observation sequence, find the possible state sequences- *Viterbi*
2. Given the observation sequence, find its probability- forward/backward algorithm
3. Given the observation sequence find the HMM parameters.- Baum-Welch algorithm

# Illustration of Viterbi

- The “start” and “end” are important in a sequence
- Subtrees get eliminated due to the Markov Assumption

## POS Tagset

- N(noun), V(verb), O(other) [simplified]
- ^ (start), . (end) [start & end states]

# Illustration of Viterbi

## Lexicon

people: N, V

laugh: N, V

.

.

.

## Corpora for Training

\_\_\_\_\_  $\wedge w_{11-t_{11}} w_{12-t_{12}} w_{13-t_{13}} \dots w_{1k_1-t_{1k_1}} \cdot$

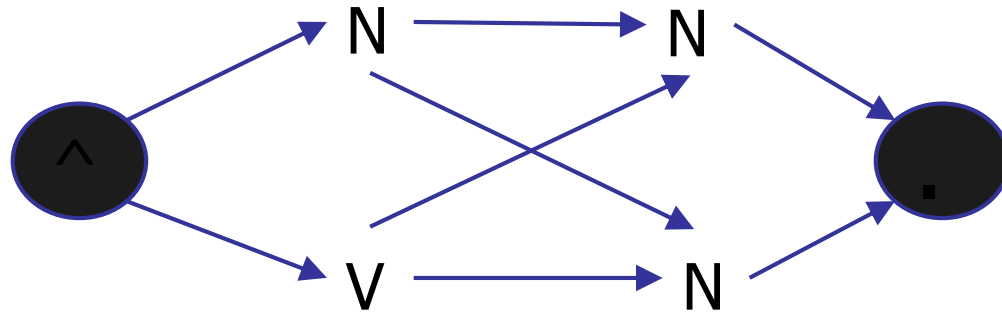
$\wedge w_{21-t_{21}} w_{22-t_{22}} w_{23-t_{23}} \dots w_{2k_2-t_{2k_2}} \cdot$

.

.

$\wedge w_{n1-t_{n1}} w_{n2-t_{n2}} w_{n3-t_{n3}} \dots w_{nk_n-t_{nk_n}} \cdot$

# Inference



Partial sequence graph

	^	N	V	O	.
^	0	0.6	0.2	0.2	0
N	0	0.1	0.4	0.3	0.2
V	0	0.3	0.1	0.3	0.3
O	0	0.3	0.2	0.3	0.2
.	1	0	0	0	0

This transition table will change from language to language due to language divergences.



# Lexical Probability Table

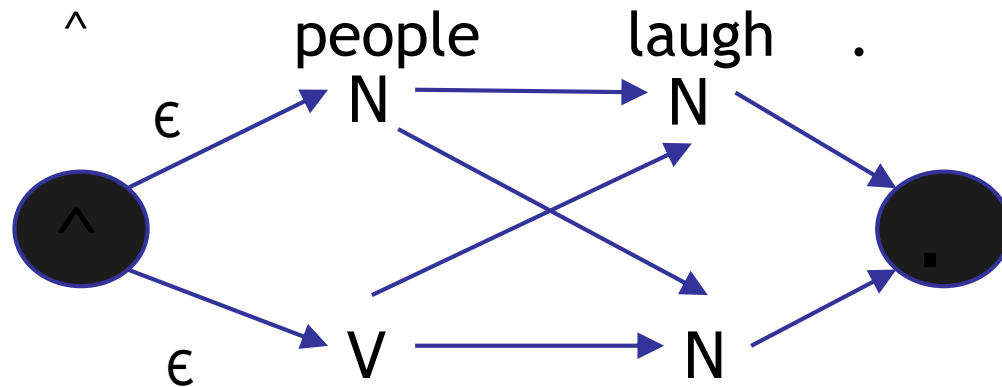
	€	people	laugh	...	...
^	1	0	0	...	0
N	0	$1 \times 10^{-3}$	$1 \times 10^{-5}$	...	...
V	0	$1 \times 10^{-6}$	$1 \times 10^{-3}$	...	...
O	0	0	0	...	...
.	1	0	0	0	0

Size of this table = # pos tags in tagset X vocabulary size

vocabulary size = # unique words in corpus

# Inference

New Sentence:

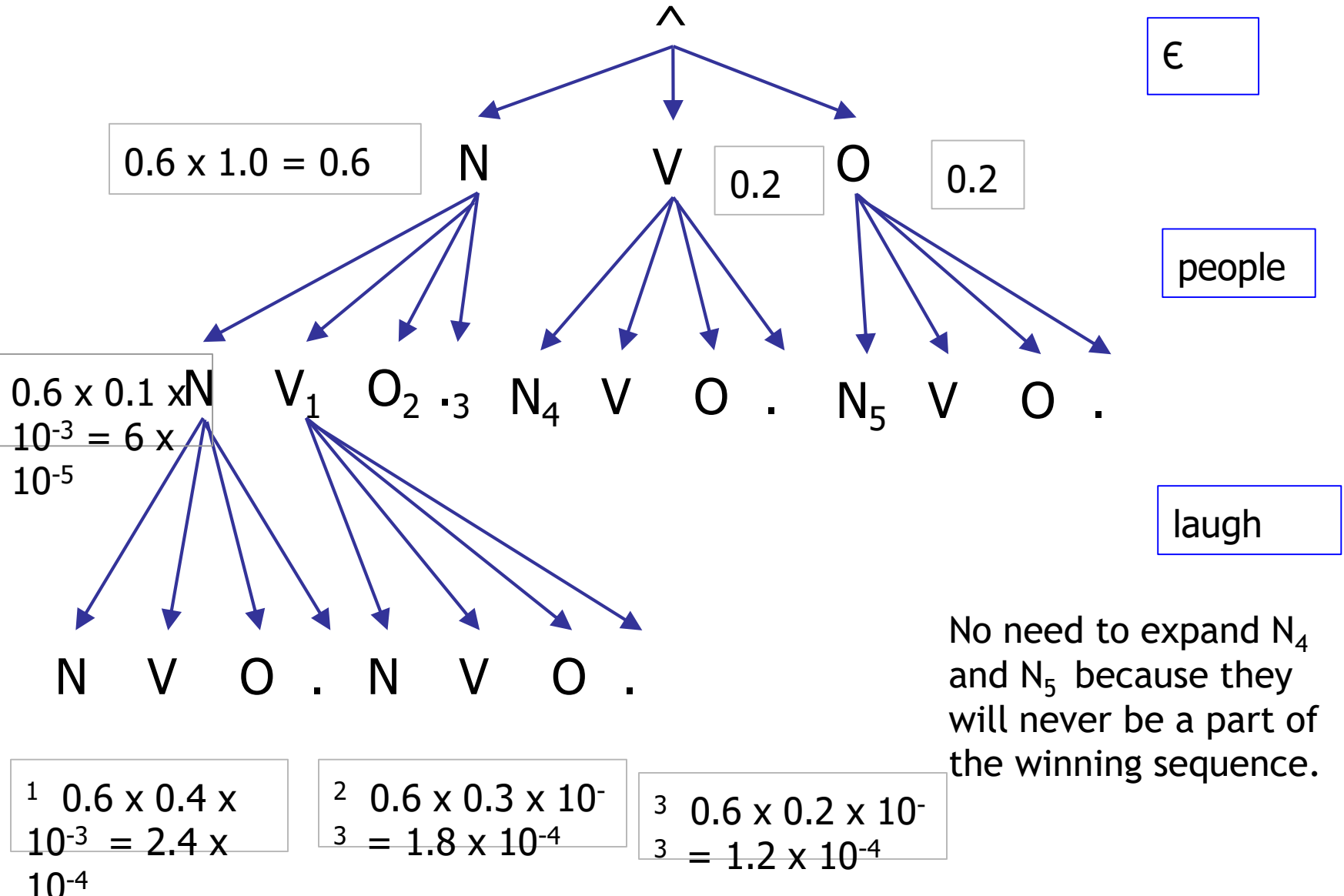


$$p(\text{^ N N .} \mid \text{^ people laugh .})$$
$$= (0.6 \times 1.0) \times (0.1 \times 1 \times 10^{-3}) \times (0.2 \times 1 \times 10^{-5})$$

# Computational Complexity

- If we have to get the probability of each sequence and then find maximum among them, we would run into exponential number of computations
- If  $|s| = \text{\#states (tags + ^ + .)}$   
and  $|o| = \text{length of sentence ( words + ^ + .)}$   
Then,  $\text{\#sequences} = s^{|o|-2}$
- But, a large number of partial computations can be reused using Dynamic Programming

# Dynamic Programming



# Computational Complexity

- Retain only those N / V / O nodes which ends in the highest sequence probability
- Now, complexity reduces from  $|s|^{|o|}$  to  $|s| \cdot |o|$
- Here, we followed the Markov assumption of order 1

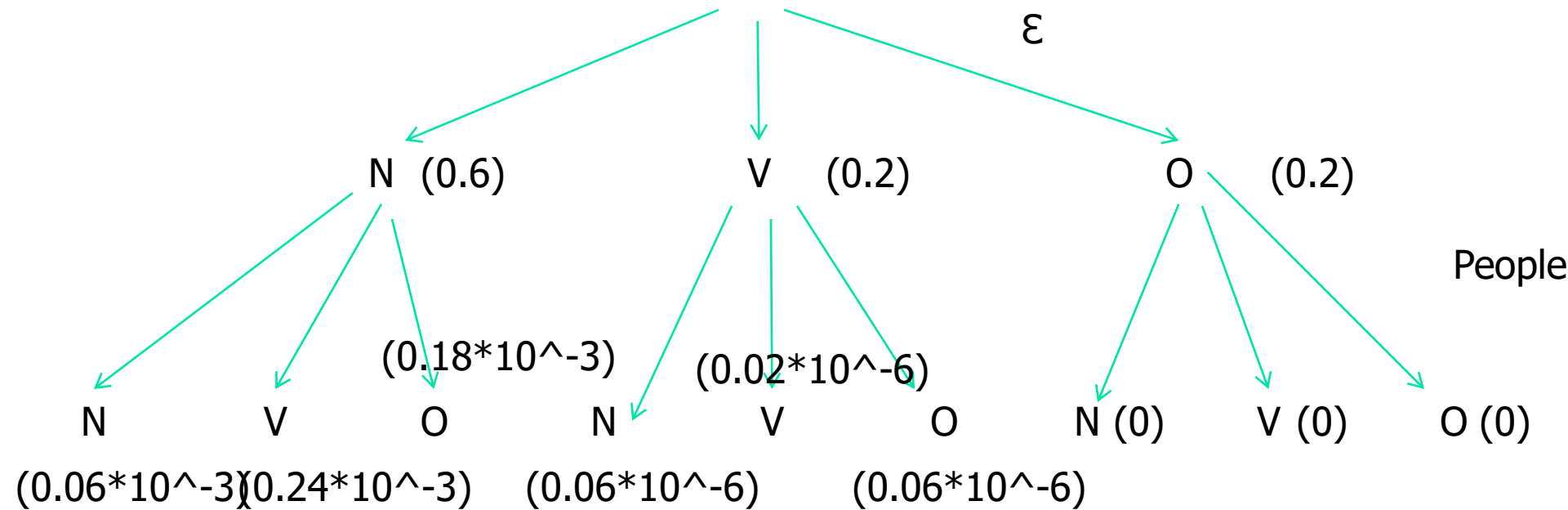
# Points to ponder wrt HMM and Viterbi

# Viterbi Algorithm

- Start with the start state
- Keep advancing sequences that are “maximum” amongst all those ending in the same state

# Viterbi Algorithm

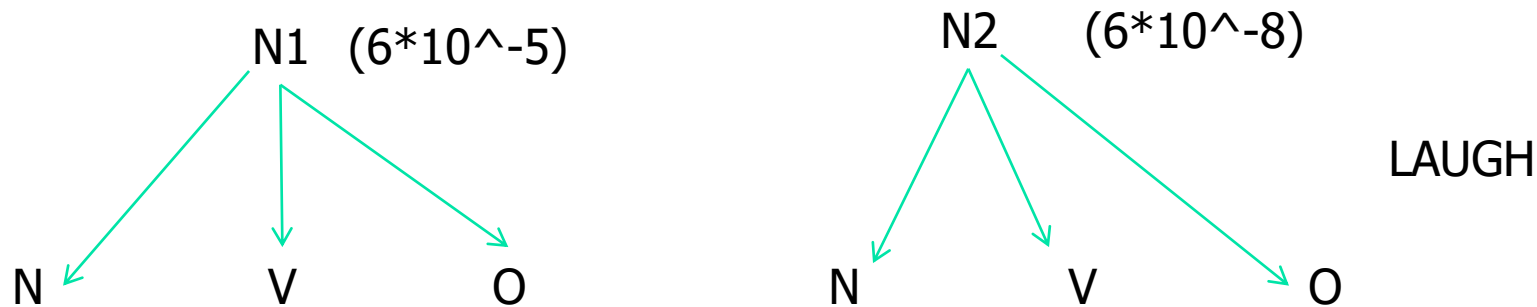
Tree for the sentence: "People laugh"



**Claim:** We do not need to draw all the subtrees in the algorithm



# Viterbi phenomenon (Markov process)



Next step all the probabilities will be multiplied by identical probability (lexical and transition). So children of N2 will have probability less than the children of N1.

# Effect of shifting probability mass

- Will a word be always given the same tag?
- No. Consider the example:
  - ^ people the city with soldiers .
  - ^ quickly people the city . (i.e., 'populate')
- In the first sentence "people" is most likely to be tagged as noun, whereas in the second, probability mass will shift and "people" will be tagged as verb, since it occurs after an adverb.

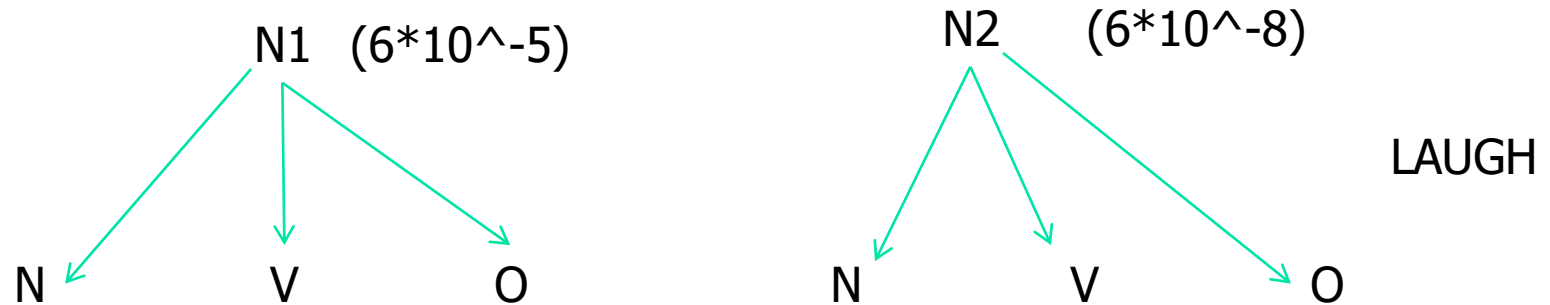
# Tail phenomenon and Language phenomenon

- Long tail Phenomenon: Probability is very low but not zero over a large observed sequence.



- Language Phenomenon:
  - “people” which is predominantly tagged as “Noun” displays a long tail phenomenon.
  - “laugh” is predominantly tagged as “Verb”.

# Viterbi phenomenon (Markov process)



Next step all the probabilities will be multiplied by identical probability (lexical and transition). So children of N2 will have probability less than the children of N1.

# Back to the Urn Example

- Here :
  - $S = \{U_1, U_2, U_3\}$
  - $V = \{R, G, B\}$
- For observation:
  - $O = \{o_1 \dots o_n\}$
- And State sequence
  - $Q = \{q_1 \dots q_n\}$
- $n$  is

$$\pi_i = P(q_1 = U_i)$$

A =

	$U_1$	$U_2$	$U_3$
$U_1$	0.1	0.4	0.5
$U_2$	0.6	0.2	0.2
$U_3$	0.3	0.4	0.3
	R	G	B
$U_1$	0.3	0.5	0.2
$U_2$	0.1	0.4	0.5
$U_3$	0.6	0.1	0.3

B=

# Observations and states

	O <sub>1</sub>	O <sub>2</sub>	O <sub>3</sub>	O <sub>4</sub>	O <sub>5</sub>	O <sub>6</sub>	O <sub>7</sub>	O <sub>8</sub>
OBS:	R	R	G	G	B	R	G	R
State:	S <sub>1</sub>	S <sub>2</sub>	S <sub>3</sub>	S <sub>4</sub>	S <sub>5</sub>	S <sub>6</sub>	S <sub>7</sub>	S <sub>8</sub>

$S_i = U_1/U_2/U_3$ ; A particular state

S: State sequence

O: Observation sequence

$S^*$  = "best" possible state (urn) sequence

Goal: Maximize  $P(S^*|O)$  by choosing "best" S

# Goal

- Maximize  $P(S|O)$  where  $S$  is the State Sequence and  $O$  is the Observation Sequence

$$S^* = \arg \max_S (P(S | O))$$

# False Start

	$O_1$	$O_2$	$O_3$	$O_4$	$O_5$	$O_6$	$O_7$	$O_8$
OBS:	R	R	G	G	B	R	G	R
State:	$S_1$	$S_2$	$S_3$	$S_4$	$S_5$	$S_6$	$S_7$	$S_8$

$$P(S | O) = P(S_{1-8} | O_{1-8})$$

$$P(S | O) = P(S_1 | O).P(S_2 | S_1, O).P(S_3 | S_{1-2}, O)...P(S_8 | S_{1-7}, O)$$

By Markov Assumption (a state depends only on the previous state)

$$P(S | O) = P(S_1 | O).P(S_2 | S_1, O).P(S_3 | S_2, O)...P(S_8 | S_7, O)$$



# Baye's Theorem

$$P(A \mid B) = P(A).P(B \mid A) / P(B)$$

$P(A)$  -: Prior

$P(B|A)$  -: Likelihood

$$\arg \max_s P(S \mid O) = \arg \max_s P(S).P(O \mid S)$$

# State Transitions Probability

$$P(S) = P(S_{1-8})$$

$$P(S) = P(S_1).P(S_2 | S_1).P(S_3 | S_{1-2}).P(S_4 | S_{1-3})...P(S_8 | S_{1-7})$$

**By Markov Assumption (k=1)**

$$P(S) = P(S_1).P(S_2 | S_1).P(S_3 | S_2).P(S_4 | S_3)...P(S_8 | S_7)$$

# Observation Sequence probability

$$P(O | S) = P(O_1 | S_{1-8}).P(O_2 | O_1, S_{1-8}).P(O_3 | O_{1-2}, S_{1-8})...P(O_8 | O_{1-7}, S_{1-8})$$

Assumption that ball drawn depends only on the Urn chosen

$$P(O | S) = P(O_1 | S_1).P(O_2 | S_2).P(O_3 | S_3)...P(O_8 | S_8)$$

$$P(S | O) = P(S).P(O | S)$$

$$P(S | O) = P(S_1).P(S_2 | S_1).P(S_3 | S_2).P(S_4 | S_3)...P(S_8 | S_7).$$

$$P(O_1 | S_1).P(O_2 | S_2).P(O_3 | S_3)...P(O_8 | S_8)$$

# Grouping terms

	$O_0$	$O_1$	$O_2$	$O_3$	$O_4$	$O_5$	$O_6$	$O_7$	$O_8$	
Obs:	$\epsilon$	R	R	G	G	B	R	G	R	
State:	$S_0$	$S_1$	$S_2$	$S_3$	$S_4$	$S_5$	$S_6$	$S_7$	$S_8$	$S_9$

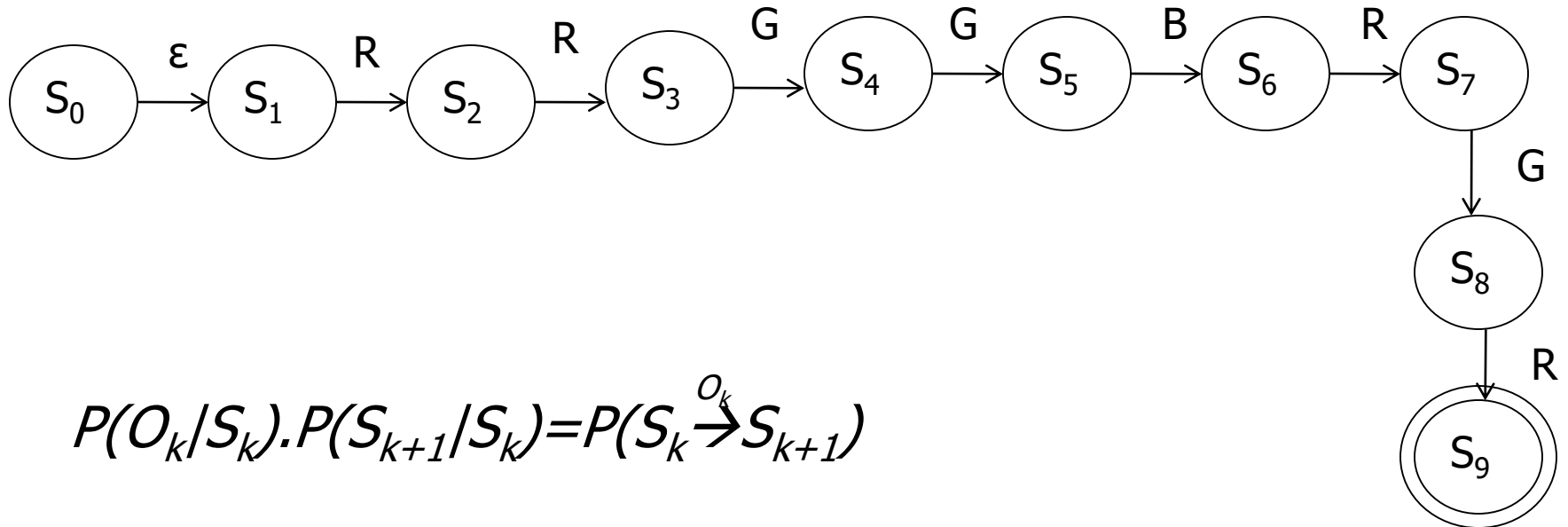
$$\begin{aligned}
 &P(S).P(O|S) \\
 = &[P(O_0|S_0).P(S_1|S_0)]. \\
 &[P(O_1|S_1).P(S_2|S_1)]. \\
 &[P(O_2|S_2).P(S_3|S_2)]. \\
 &[P(O_3|S_3).P(S_4|S_3)]. \\
 &[P(O_4|S_4).P(S_5|S_4)]. \\
 &[P(O_5|S_5).P(S_6|S_5)]. \\
 &[P(O_6|S_6).P(S_7|S_6)]. \\
 &[P(O_7|S_7).P(S_8|S_7)]. \\
 &[P(O_8|S_8).P(S_9|S_8)].
 \end{aligned}$$

We introduce the states  $S_0$  and  $S_9$  as initial and final states respectively.

After  $S_8$  the next state is  $S_9$  with probability 1, i.e.,  
 $P(S_9|S_8)=1$   
 $O_0$  is  $\epsilon$ -transition

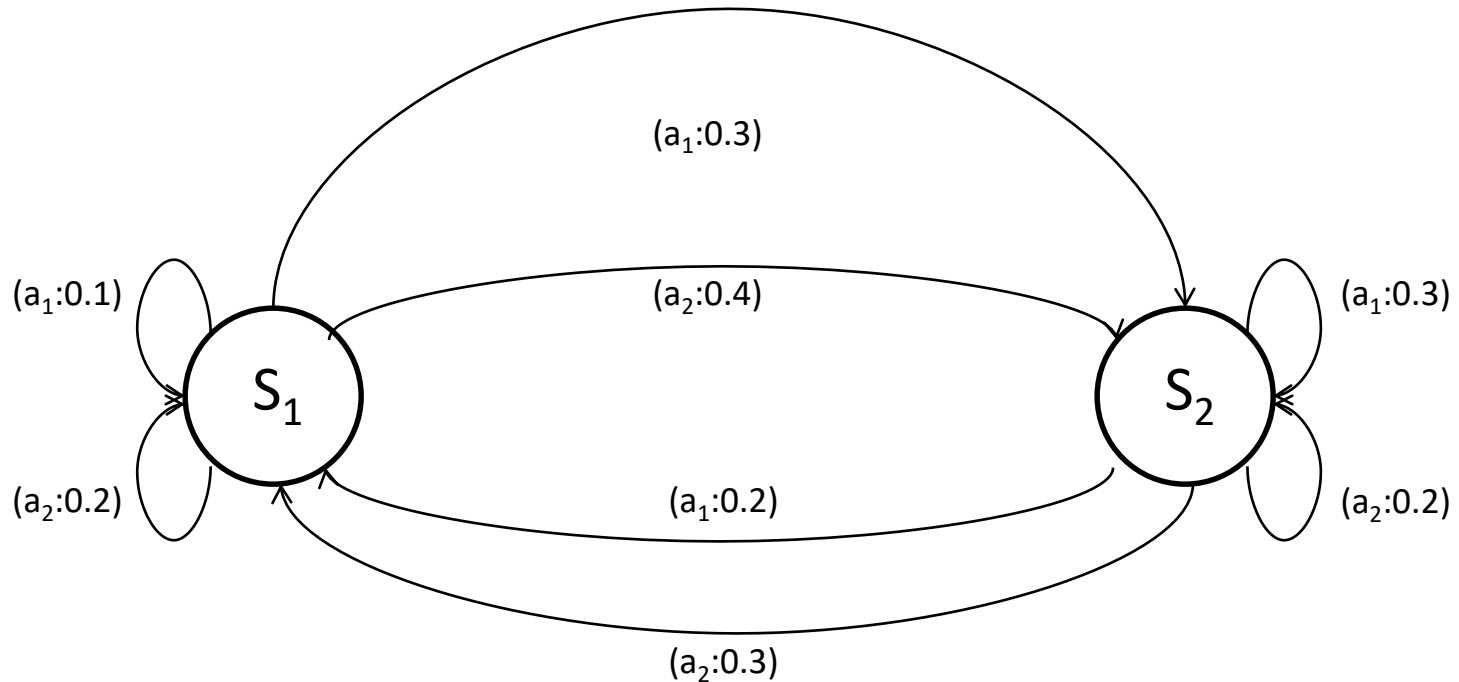
# Introducing useful notation

	$O_0$	$O_1$	$O_2$	$O_3$	$O_4$	$O_5$	$O_6$	$O_7$	$O_8$	
Obs:	$\epsilon$	R	R	G	G	B	R	G	R	
State:	$S_0$	$S_1$	$S_2$	$S_3$	$S_4$	$S_5$	$S_6$	$S_7$	$S_8$	$S_9$



$$P(O_k/S_k) \cdot P(S_{k+1}/S_k) = P(S_k \xrightarrow{O_k} S_{k+1})$$

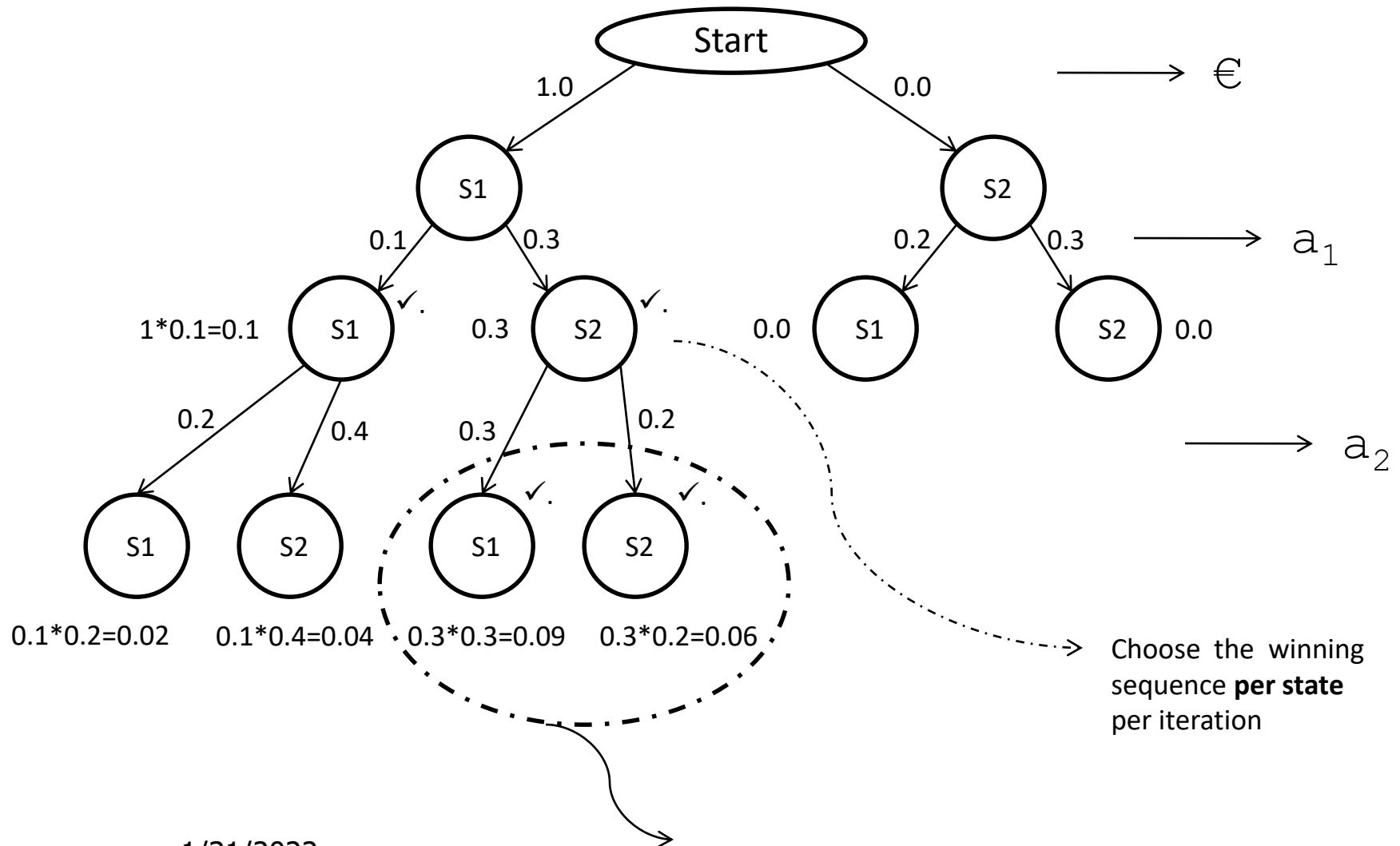
# Probabilistic FSM



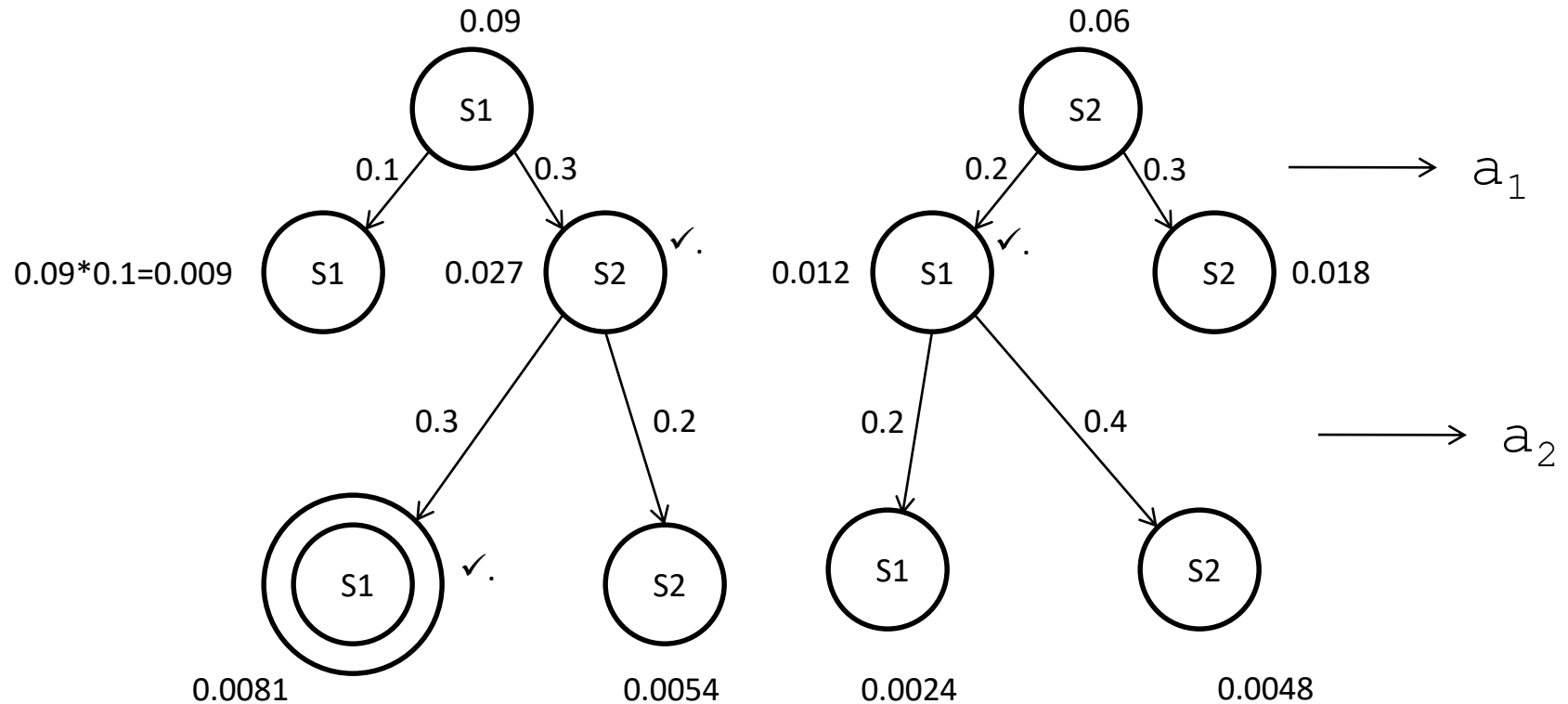
The question here is:

“what is the most likely state sequence given the output sequence seen”

# Developing the tree



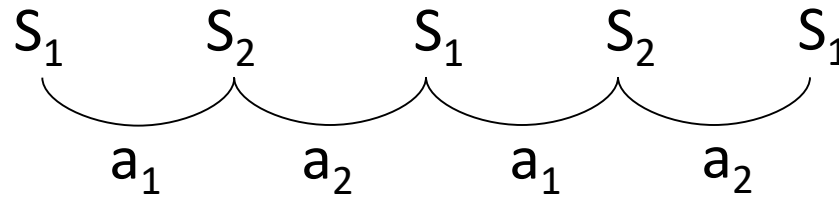
# Tree structure contd...



The problem being addressed by this tree is  $S^* = \arg \max_s P(S \mid a_1 - a_2 - a_1 - a_2, \mu)$   
 $a_1 - a_2 - a_1 - a_2$  is the output sequence and  $\mu$  the model or the machine



Path found:  
(working backward)

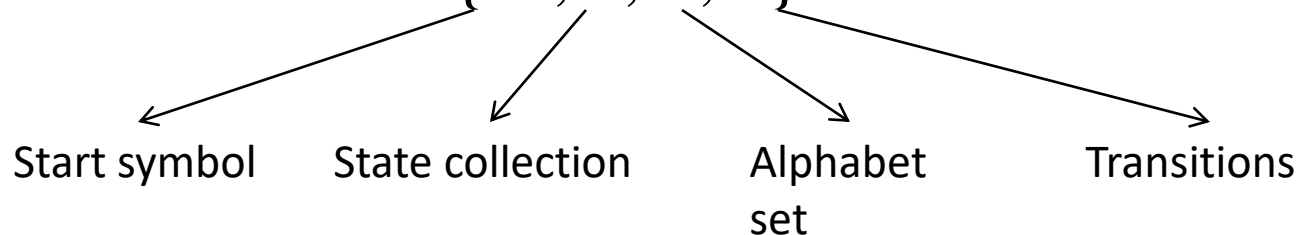


Problem statement: Find the best possible sequence

$$S^* = \arg \max_s P(S | O, \mu)$$

where,  $S \rightarrow$  State Seq,  $O \rightarrow$  Output Seq,  $\mu \rightarrow$  Model or Machine

Model or Machine =  $\{S_0, S, A, T\}$



T is defined as  $P(S_i \xrightarrow{a_k} S_j) \quad \forall i, j, k$

# Probability of observation sequence

# Why probability of observation sequence?: Language modeling problem

Probabilities computed in the context of corpora

1.  $P(\text{"The sun rises in the east"})$
2.  $P(\text{"The sun rise in the east"})$ 
  - Less probable because of grammatical mistake.
3.  $P(\text{The svn rises in the east})$ 
  - Less probable because of lexical mistake.
4.  $P(\text{The sun rises in the west})$ 
  - Less probable because of semantic mistake.

# Uses of language model

## 1. Detect well-formedness

- Lexical, syntactic, semantic, pragmatic, discourse

## 2. Language identification

- Given a piece of text what language does it belong to.

*Good morning* - *English*

*Guten morgen* - *German*

*Bon jour* - *French*

## 3. Automatic speech recognition

## 4. Machine translation

# How to compute $P(o_0 o_1 o_2 o_3 \dots o_m)$ ?

$$P(O) = \sum_S P(O, S) \quad \text{Marginalization}$$

Consider the observation sequence,

$O_0 O_1 O_2 \dots O_m$

$S_0 S_1 S_2 S_3 \dots S_m S_{m+1}$

Where  $S_i$  s represent the state sequences.

# Computing $P(o_0o_1o_2o_3...o_m)$

$$\begin{aligned}P(O, S) &= P(S)P(O | S) \\&= P(S_0S_1S_2...S_{m+1})P(O_0O_1O_2...O_m | S) \\&= P(S_0).P(S_1 | S_0).P(S_2 | S_1)...P(S_{m+1} | S_m). \\&\quad P(O_0 | S_0).P(O_1 | S_1)....P(O_m | S_m) \\&= P(S_0)[P(O_0 | S_0).P(S_1 | S_0)].....[P(O_m | S_m).P(S_{m+1} | S_m)]\end{aligned}$$

# Forward and Backward Probability Calculation

# Forward probability $F(k,i)$

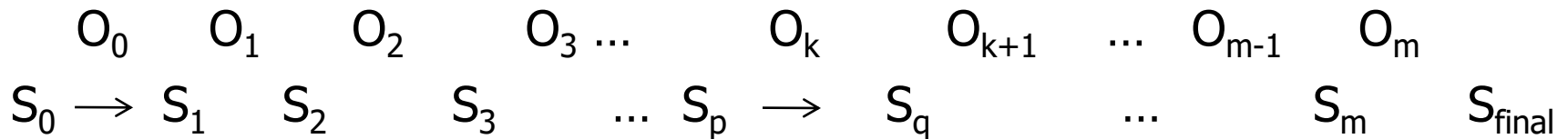
- Define  $F(k,i)$  = Probability of being in state  $S_i$  having seen  $o_0o_1o_2\cdots o_k$
- $F(k,i) = P(o_0o_1o_2\cdots o_k, S_i)$
- With  $m$  as the length of the observed sequence and  $N$  states

$$\begin{aligned} P(\text{observed sequence}) &= P(o_0o_1o_2\cdots o_m) \\ &= \sum_{p=0,N} P(o_0o_1o_2\cdots o_m, S_p) \\ &= \sum_{p=0,N} F(m, p) \end{aligned}$$



# Forward probability (contd.)

$$\begin{aligned}
 F(k, q) &= P(o_0 o_1 o_2 \dots o_k, S_q) \\
 &= P(o_0 o_1 o_2 \dots o_k, S_q) \\
 &= P(o_0 o_1 o_2 \dots o_{k-1}, o_k, S_q) \\
 &= \sum_{p=0, N} P(o_0 o_1 o_2 \dots o_{k-1}, S_p, o_k, S_q) \\
 &= \sum_{p=0, N} P(o_0 o_1 o_2 \dots o_{k-1}, S_p) \cdot \\
 &\quad P(o_k, S_q | o_0 o_1 o_2 \dots o_{k-1}, S_p) \\
 &= \sum_{p=0, N} F(k-1, p) \cdot P(o_k, S_q | S_p) \\
 &= \sum_{p=0, N} F(k-1, p) \cdot P(S_p \xrightarrow{o_k} S_q)
 \end{aligned}$$



# Backward probability $B(k,i)$

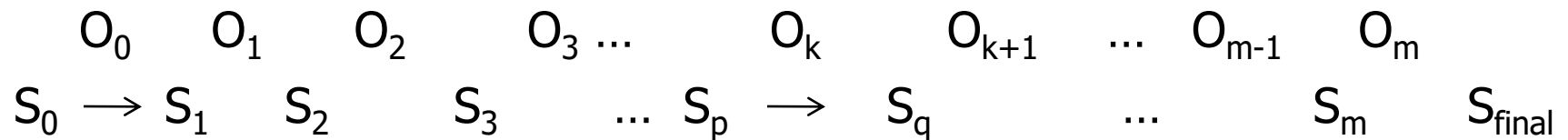
- Define  $B(k,i)$  = Probability of seeing  $o_k o_{k+1} o_{k+2} \dots o_m$  given that the state was  $S_i$

$$B(k,i) = P(o_k o_{k+1} o_{k+2} \dots o_m \mid S_i)$$

- With  $m$  as the length of the whole observed sequence
- $P(\text{observed sequence}) = P(o_0 o_1 o_2 \dots o_m)$   
 $= P(o_0 o_1 o_2 \dots o_m \mid S_0)$   
 $= B(0,0)$

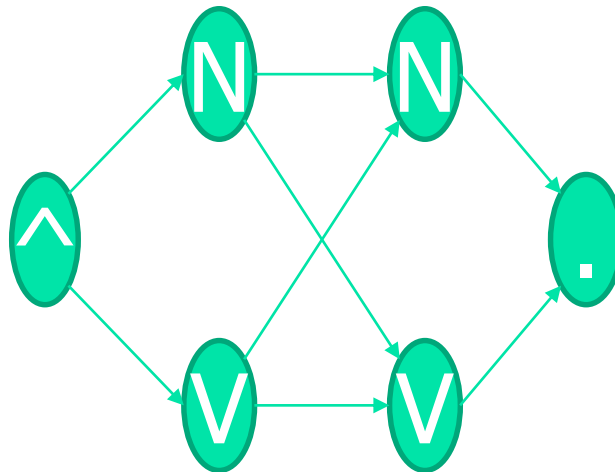
# Backward probability (contd.)

$$\begin{aligned}
 B(k, p) &= P(o_k o_{k+1} o_{k+2} \dots o_m \mid S_p) \\
 &= P(o_{k+1} o_{k+2} \dots o_m, o_k \mid S_p) \\
 &= \sum_{q=0, N} P(o_{k+1} o_{k+2} \dots o_m, o_k, S_q \mid S_p) \\
 &= \sum_{q=0, N} P(o_k, S_q \mid S_p) \\
 &\quad P(o_{k+1} o_{k+2} \dots o_m \mid o_k, S_q, S_p) \\
 &= \sum_{q=0, N} P(o_{k+1} o_{k+2} \dots o_m \mid S_q) \cdot P(o_k, S_q \mid S_p) \\
 &= \sum_{q=0, N} B(k+1, q) \cdot P(S_p \xrightarrow{o_k} S_q)
 \end{aligned}$$



# How Forward Probability Works

- Goal of Forward Probability: To find  $P(O)$  [the probability of Observation Sequence].
- E.g.      ^      People laugh      .



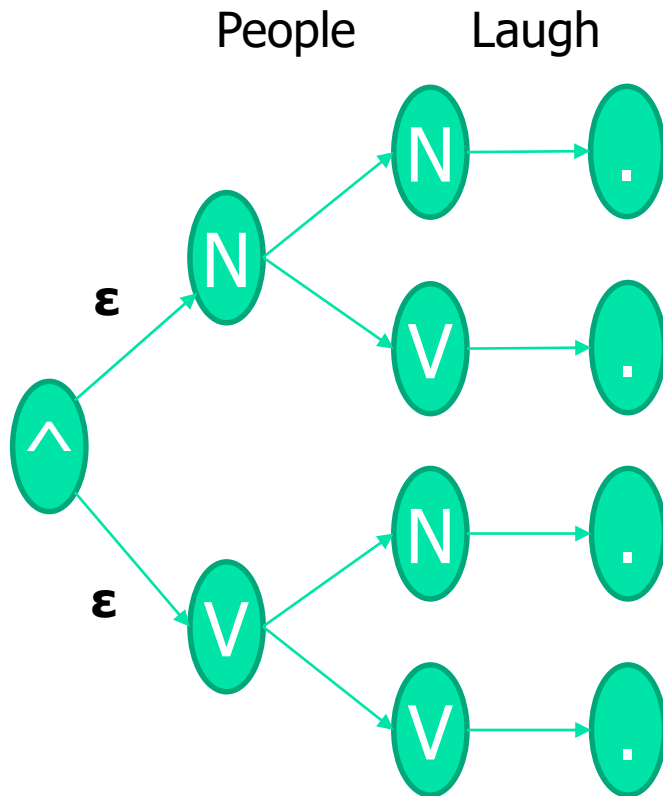
# Transition and Lexical Probability Tables

	<b>^</b>	<b>N</b>	<b>V</b>	<b>.</b>		<b>ε</b>	<b>People</b>	<b>Laugh</b>
<b>^</b>	0	0.7	0.3	0	<b>^</b>	1	0	0
<b>N</b>	0	0.2	0.6	0.2	<b>N</b>	0	0.8	0.2
<b>V</b>	0	0.6	0.2	0.2	<b>V</b>	0	0.1	0.9
<b>.</b>	1	0	0	0	<b>.</b>	1	0	0

Inefficient Computation:

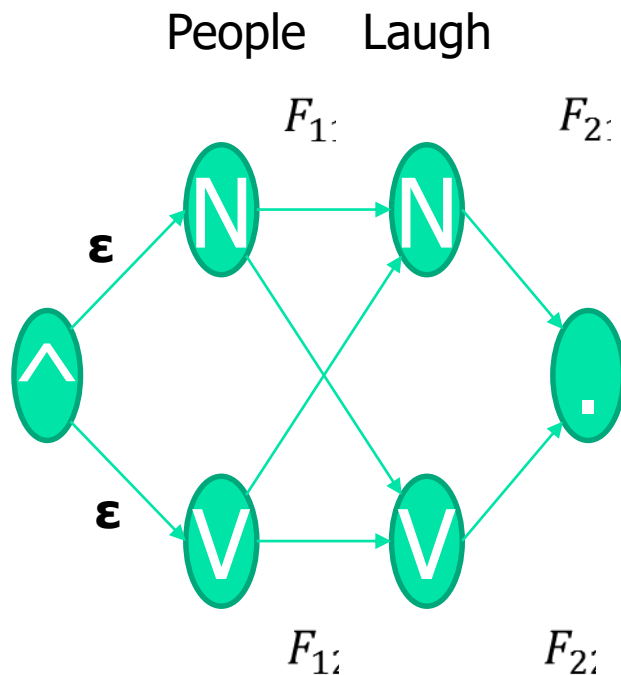
$$P(O) = \sum_S P(O, S) = \sum_S \prod_i P(S_i \xrightarrow{o_j} S_j)$$

# Computation in various paths of the Tree



	$\epsilon$	People	Laugh
Path 1:	$\wedge$	N	N
.			
$P(\text{Path1}) = (1.0 \times 0.7) \times (0.8 \times 0.2) \times (0.2 \times 0.2)$			
	$\epsilon$	People	Laugh
Path 2:	$\wedge$	N	V
.			
$P(\text{Path2}) = (1.0 \times 0.7) \times (0.8 \times 0.6) \times (0.9 \times 0.2)$			
	$\epsilon$	People	Laugh
Path 3:	$\wedge$	V	N
.			
$P(\text{Path3}) = (1.0 \times 0.3) \times (0.1 \times 0.6) \times (0.2 \times 0.2)$			
	$\epsilon$	People	Laugh
Path 4:	$\wedge$	V	V
.			
$P(\text{Path4}) = (1.0 \times 0.3) \times (0.1 \times 0.2) \times (0.9 \times 0.2)$			

# Computations on the Trellis



$F = \text{accumulated } F \times \text{output probability} \times$   
■  $\text{transition probability}$

$$F_{11} = 0.7 \times 1.0$$

$$F_{12} = 0.3 \times 1.0$$

$$F_{21} = F_{11} \times (0.2 \times 0.3) + F_{12} \times (0.6 \times 0.1)$$

$$F_{22} = F_{11} \times (0.6 \times 0.8) + F_{12} \times (0.2 \times 0.1)$$

$$F_{31} = F_{21} \times (0.2 \times 0.2) + F_{22} \times (0.2 \times 0.9)$$

$F_{3i}$

# Number of Multiplications

## Tree

- ❖ Each path has 5 multiplications + 1 addition.
- ❖ There are 4 paths in the tree.
- ❖ Therefore, total of 20 multiplications and 3 additions.

## Trellis

- ❖  $F_{11}$ ,  $\rightarrow$  1 multiplication
- ❖  $F_{12}$ ,  $\rightarrow$  1 multiplication
- ❖  $F_{21} = F_{11} \times (1 \text{ mult}) + F_{12} \times (1 \text{ mult})$   
 $= 4 \text{ multiplications} + 1 \text{ addition}$
- ❖ Similarly, for  $F_{22}$  and  $F_{31}$ , 4 multiplications and 1 addition each.
- ❖ So, total of 14 multiplications and 3 additions.



# Complexity

■ Let  $|S| = \# \text{States}$

And  $|O| = \text{Observation length} - |\{\wedge, \cdot\}|$

- ❖ Stage 1 of Trellis:  $|S|$  multiplications
- ❖ Stage 2 of Trellis:  $|S|$  nodes; each node needs computation over  $|S|$  arcs.
  - ❖ Each Arc = 1 multiplication
  - ❖ Accumulated F = 1 more multiplication
  - ❖ Total  $2|S|^2$  multiplications
- ❖ Same for each stage before reading  $\cdot$
- ❖ At final stage ( $\cdot$ )  $\rightarrow 2|S|$  multiplications

# Summary : Forward Algorithm

1. Accumulate F over each stage of trellis.
2. Take sum of F values multiplied by  $P(S_i \xrightarrow{O_j} S_{i+1})$ .
3. Complexity =  $|S| + 2|S|^2 (|O| - 1) + 2|S|$   
 $= 2|S|^2 |O| - 2|S|^2 + 3|S|$   
 $= O(|S|^2 \cdot |O|)$   
*i.e., linear in the length of input and quadratic in number of states.*

# Reading List

- TnT (<http://www.aclweb.org/anthology-new/A/A00/A00-1031.pdf>)
- Brill Tagger  
([http://delivery.acm.org/10.1145/1080000/1075553/p112-brill.pdf?ip=182.19.16.71&acc=OPEN&CFID=129797466&CFTOKEN=72601926&acm\\_\\_=1342975719\\_082233e0ca9b5d1d67a9997c03a649d1](http://delivery.acm.org/10.1145/1080000/1075553/p112-brill.pdf?ip=182.19.16.71&acc=OPEN&CFID=129797466&CFTOKEN=72601926&acm__=1342975719_082233e0ca9b5d1d67a9997c03a649d1))