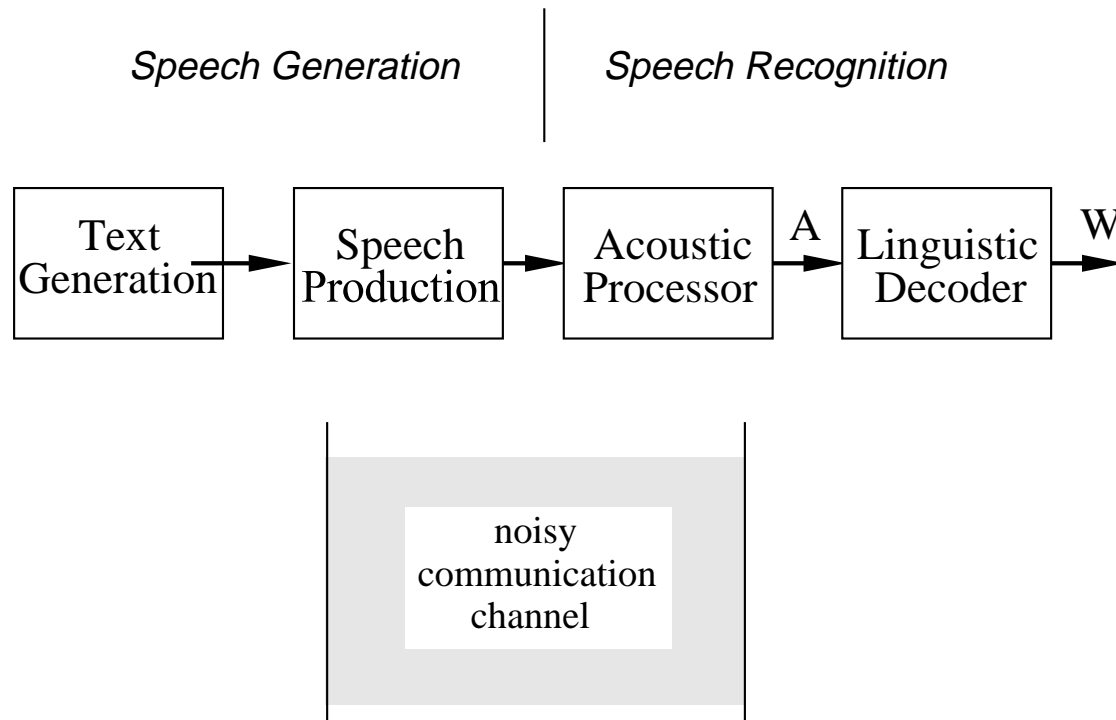


# Hidden Markov Modelling

- Introduction
- Problem formulation
- Forward-Backward algorithm
- Viterbi search
- Baum-Welch parameter estimation
- Other considerations
  - Multiple observation sequences
  - Phone-based models for continuous speech recognition
  - Continuous density HMMs
  - Implementation issues

# Information Theoretic Approach to ASR



Recognition is achieved by maximizing the probability of the linguistic string,  $\mathbf{W}$ , given the acoustic evidence,  $\mathbf{A}$ , i.e., choose the linguistic sequence  $\hat{\mathbf{W}}$  such that

$$P(\hat{\mathbf{W}}|\mathbf{A}) = \max_{\mathbf{W}} P(\mathbf{W}|\mathbf{A})$$

# Information Theoretic Approach to ASR

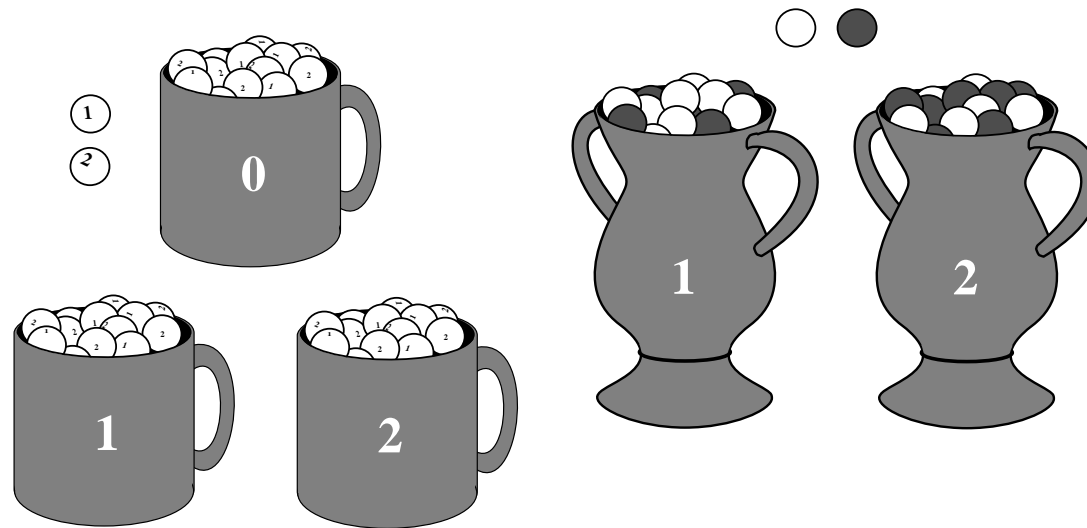
- From Bayes rule:

$$P(\mathbf{W}|\mathbf{A}) = \frac{P(\mathbf{A}|\mathbf{W})P(\mathbf{W})}{P(\mathbf{A})}$$

- Hidden Markov modelling (HMM) deals with the quantity  $P(\mathbf{A}|\mathbf{W})$
- Change in notation:

$$\begin{array}{lcl} \mathbf{A} & \rightarrow & \mathbf{O} \\ \mathbf{W} & \rightarrow & \lambda \\ P(\mathbf{A}|\mathbf{W}) & \rightarrow & P(\mathbf{O}|\lambda) \end{array}$$

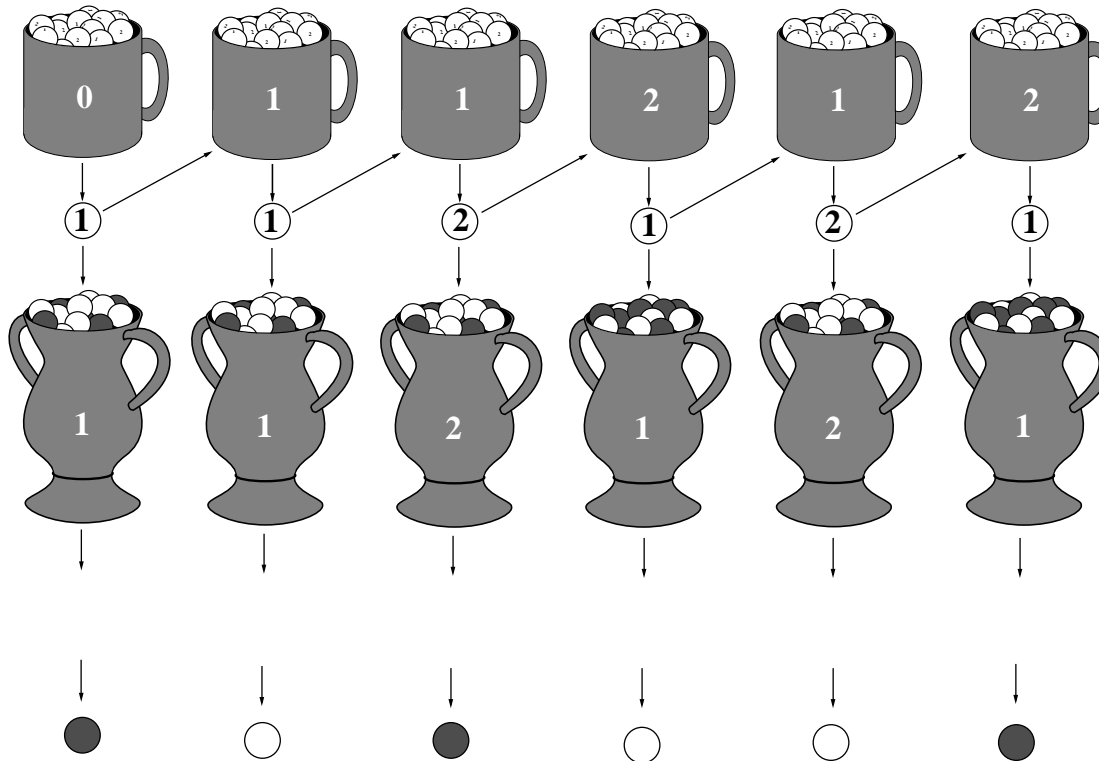
## HMM: An Example



- Consider 3 mugs, each with mixtures of *state* stones, 1 and 2
- The fractions for the  $i^{th}$  mug are  $a_{i1}$  and  $a_{i2}$ , and  $a_{i1} + a_{i2} = 1$
- Consider 2 urns, each with mixtures of black and white balls
- The fractions for the  $i^{th}$  urn are  $b_i(B)$  and  $b_i(W)$ ;  $b_i(B) + b_i(W) = 1$
- The parameter vector for this model is:

$$\lambda = \{a_{01}, a_{02}, a_{11}, a_{12}, a_{21}, a_{22}, b_1(B), b_1(W), b_2(B), b_2(W)\}$$

# HMM: An Example (cont'd)



Observation Sequence:  $\mathbf{O} = \{B, W, B, W, W, B\}$

State Sequence:  $\mathbf{Q} = \{1, 1, 2, 1, 2, 1\}$

**Goal:** Given the model  $\lambda$  and the observation sequence  $\mathbf{O}$ , how can the underlying state sequence  $\mathbf{Q}$  be determined?

# Elements of a Discrete Hidden Markov Model

- $N$ : number of states in the model
  - states,  $\mathbf{s} = \{s_1, s_2, \dots, s_N\}$
  - state at time  $t$ ,  $q_t \in \mathbf{s}$
- $M$ : number of observation symbols (i.e., discrete observations)
  - observation symbols,  $\mathbf{v} = \{v_1, v_2, \dots, v_M\}$
  - observation at time  $t$ ,  $o_t \in \mathbf{v}$
- $\mathbf{A} = \{a_{ij}\}$ : state transition probability distribution
  - $a_{ij} = P(q_{t+1} = s_j | q_t = s_i)$ ,  $1 \leq i, j \leq N$
- $\mathbf{B} = \{b_j(k)\}$ : observation symbol probability distribution in state  $j$ 
  - $b_j(k) = P(v_k \text{ at } t | q_t = s_j)$ ,  $1 \leq j \leq N, 1 \leq k \leq M$
- $\pi = \{\pi_i\}$ : initial state distribution
  - $\pi_i = P(q_1 = s_i)$ ,  $1 \leq i \leq N$

Notationally, an HMM is typically written as:  $\lambda = \{\mathbf{A}, \mathbf{B}, \pi\}$

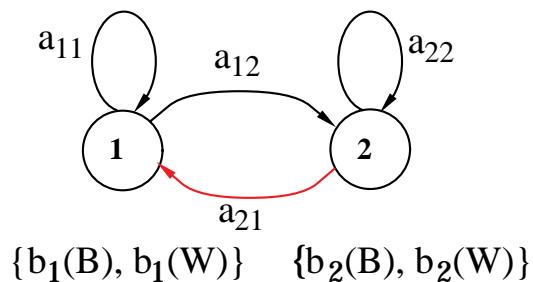
# HMM: An Example (cont'd)

For our simple example:

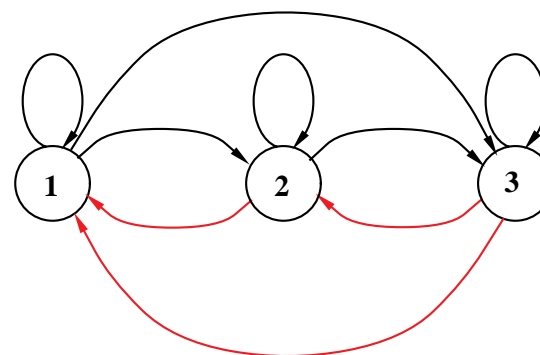
$$\pi = \{a_{01}, a_{02}\}, \quad \mathbf{A} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}, \quad \text{and} \quad \mathbf{B} = \begin{bmatrix} b_1(B) & b_1(W) \\ b_2(B) & b_2(W) \end{bmatrix}$$

## State Diagram

2-state

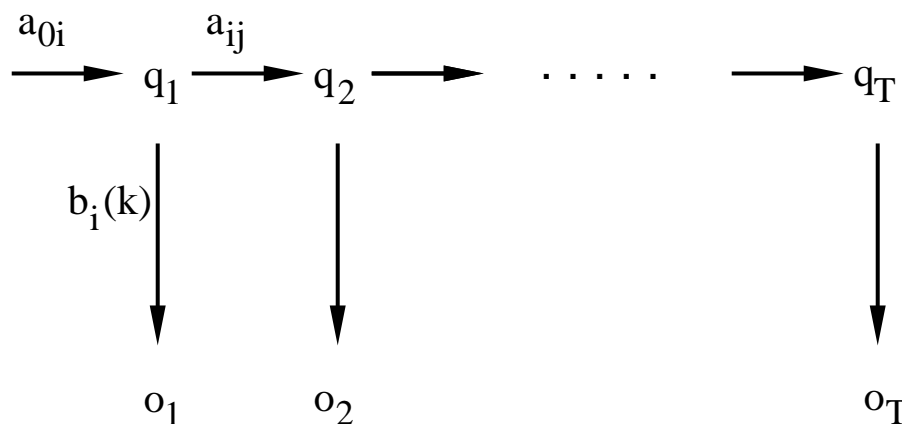


3-state



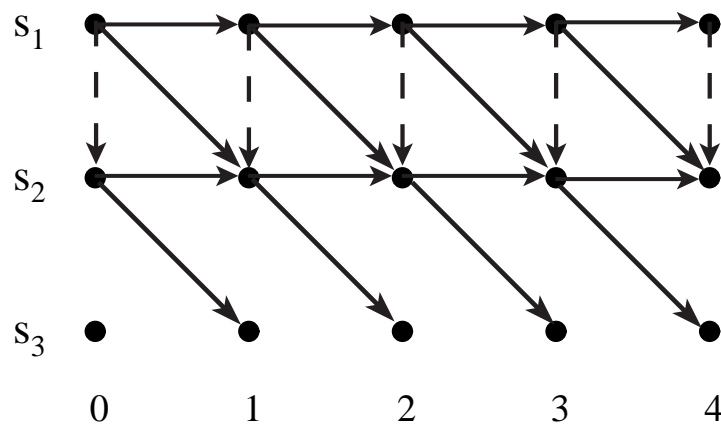
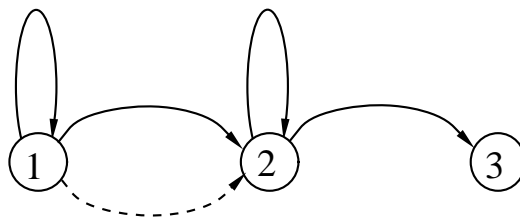
# Generation of HMM Observations

1. Choose an initial state,  $q_1 = s_i$ , based on the initial state distribution,  $\pi$
2. For  $t = 1$  to  $T$ :
  - Choose  $o_t = v_k$  according to the symbol probability distribution in state  $s_i$ ,  $b_i(k)$
  - Transition to a new state  $q_{t+1} = s_j$  according to the state transition probability distribution for state  $s_i$ ,  $a_{ij}$
3. Increment  $t$  by 1, return to step 2 if  $t \leq T$ ; else, terminate





# Representing State Diagram by Trellis



The dashed line represents a *null* transition, where no observation symbol is generated

# Three Basic HMM Problems

1. **Scoring:** Given an observation sequence  $\mathbf{O} = \{o_1, o_2, \dots, o_T\}$  and a model  $\lambda = \{\mathbf{A}, \mathbf{B}, \pi\}$ , how do we compute  $P(\mathbf{O} | \lambda)$ , the probability of the observation sequence?  
==> The Forward-Backward Algorithm
2. **Matching:** Given an observation sequence  $\mathbf{O} = \{o_1, o_2, \dots, o_T\}$ , how do we choose a state sequence  $\mathbf{Q} = \{q_1, q_2, \dots, q_T\}$  which is optimum in some sense?  
==> The Viterbi Algorithm
3. **Training:** How do we adjust the model parameters  $\lambda = \{A, B, \pi\}$  to maximize  $P(\mathbf{O} | \lambda)$ ?  
==> The Baum-Welch Re-estimation Procedures

# Computation of $P(\mathbf{O}|\lambda)$

$$P(\mathbf{O}|\lambda) = \sum_{\text{all } \mathbf{Q}} P(\mathbf{O}, \mathbf{Q}|\lambda)$$

$$P(\mathbf{O}, \mathbf{Q}|\lambda) = P(\mathbf{O}|\mathbf{Q}, \lambda)P(\mathbf{Q}|\lambda)$$

- Consider the *fixed* state sequence:  $\mathbf{Q} = q_1 q_2 \dots q_T$

$$P(\mathbf{O}|\mathbf{Q}, \lambda) = b_{q_1}(o_1)b_{q_2}(o_2)\dots b_{q_T}(o_T)$$

$$P(\mathbf{Q}|\lambda) = \pi_{q_1} a_{q_1 q_2} a_{q_2 q_3} \dots a_{q_{T-1} q_T}$$

Therefore:

$$P(\mathbf{O}|\lambda) = \sum_{q_1, q_2, \dots, q_T} \pi_{q_1} b_{q_1}(o_1) a_{q_1 q_2} b_{q_2}(o_2) \dots a_{q_{T-1} q_T} b_{q_T}(o_T)$$

- Calculation required  $\approx 2T \cdot N^T$  (there are  $N^T$  such sequences)  
For  $N = 5, T = 100 \Rightarrow 2 \cdot 100 \cdot 5^{100} \approx 10^{72}$  computations!

# The Forward Algorithm

- Let us define the forward variable,  $\alpha_t(i)$ , as the probability of the partial observation sequence up to time  $t$  *and* state  $s_i$  at time  $t$ , given the model, i.e.

$$\alpha_t(i) = P(o_1 o_2 \dots o_t, q_t = s_i | \lambda)$$

- It can easily be shown that:

$$\alpha_1(i) = \pi_i b_i(o_1), \quad 1 \leq i \leq N$$

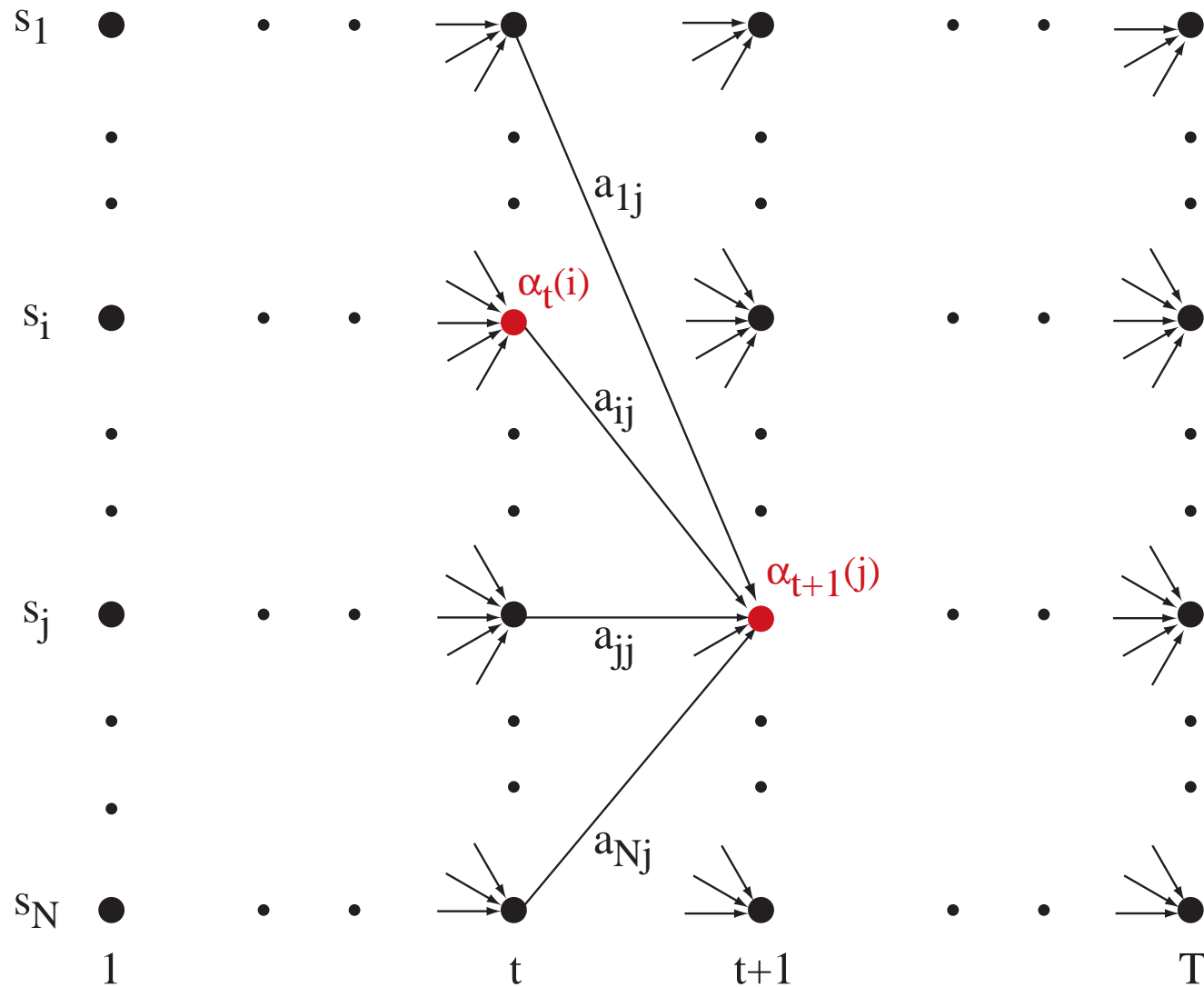
$$P(\mathbf{O} | \lambda) = \sum_{i=1}^N \alpha_T(i)$$

- By induction:

$$\alpha_{t+1}(j) = \left[ \sum_{i=1}^N \alpha_t(i) a_{ij} \right] b_j(o_{t+1}), \quad \begin{array}{l} 1 \leq t \leq T-1 \\ 1 \leq j \leq N \end{array}$$

- Calculation is on the order of  $N^2 T$ .  
For  $N = 5, T = 100 \Rightarrow 100 \cdot 5^2$  computations, instead of  $10^{72}$

# Forward Algorithm Illustration



# The Backward Algorithm

- Similarly, let us define the backward variable,  $\beta_t(i)$ , as the probability of the partial observation sequence from time  $t + 1$  to the end, given state  $s_i$  at time  $t$  and the model, i.e.

$$\beta_t(i) = P(o_{t+1} o_{t+2} \dots o_T | q_t = s_i, \lambda)$$

- It can easily be shown that:

$$\beta_T(i) = 1, \quad 1 \leq i \leq N$$

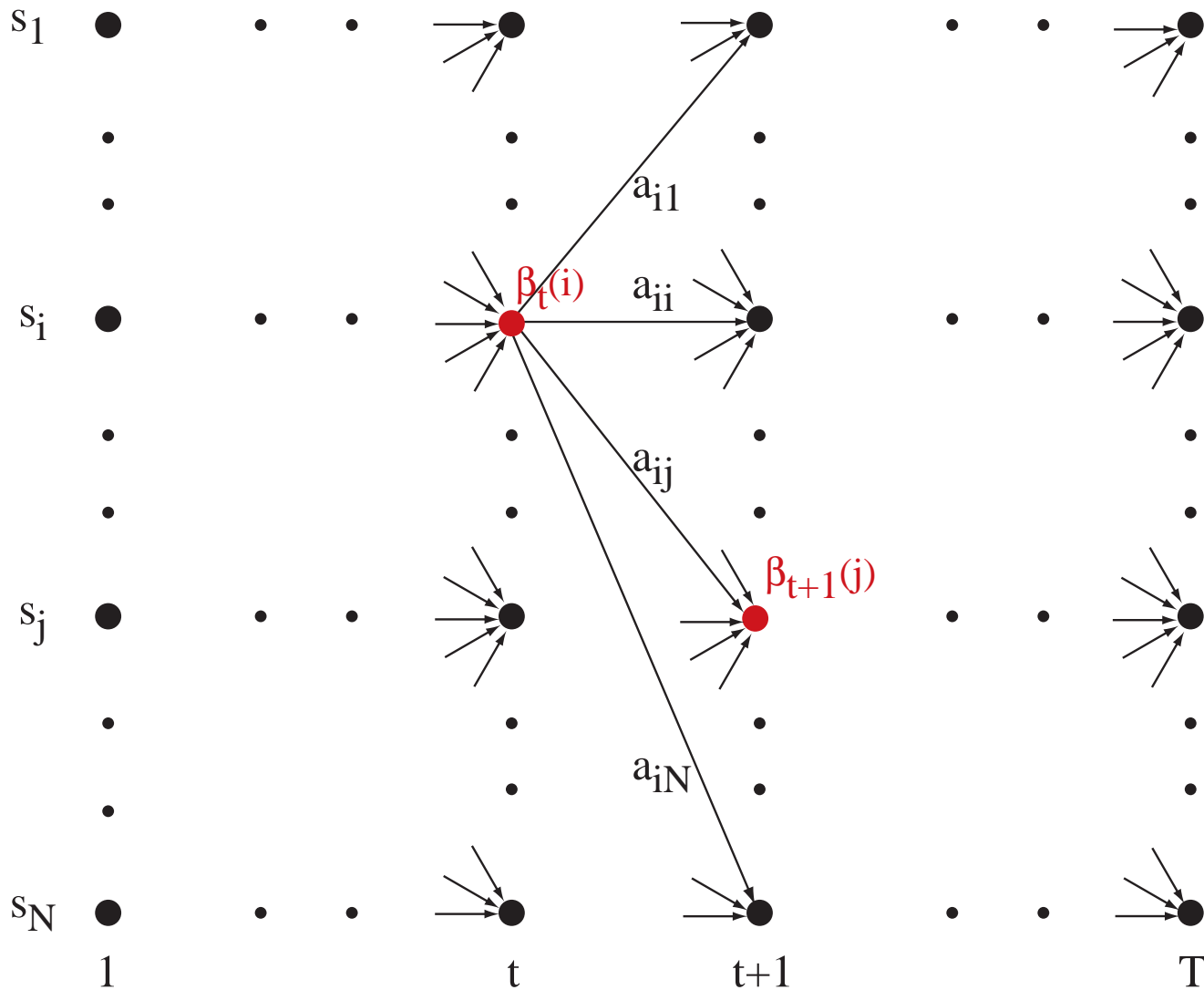
and:

$$P(\mathbf{O} | \lambda) = \sum_{i=1}^N \pi_i b_i(o_1) \beta_1(i)$$

- By induction:

$$\beta_t(i) = \sum_{j=1}^N a_{ij} b_j(o_{t+1}) \beta_{t+1}(j), \quad \begin{array}{l} t = T - 1, T - 2, \dots, 1 \\ 1 \leq i \leq N \end{array}$$

# Backward Procedure Illustration



# Finding Optimal State Sequences

- One criterion chooses states,  $q_t$ , which are *individually* most likely
  - This maximizes the expected number of correct states
- Let us define  $\gamma_t(i)$  as the probability of being in state  $s_i$  at time  $t$ , given the observation sequence and the model, i.e.

$$\gamma_t(i) = P(q_t = s_i | \mathbf{O}, \lambda) \quad \sum_{i=1}^N \gamma_t(i) = 1, \quad \forall t$$

- Then the individually most likely state,  $q_t$ , at time  $t$  is:

$$q_t = \operatorname{argmax}_{1 \leq i \leq N} \gamma_t(i) \quad 1 \leq t \leq T$$

- Note that it can be shown that:

$$\gamma_t(i) = \frac{\alpha_t(i)\beta_t(i)}{P(\mathbf{O}|\lambda)}$$



# Finding Optimal State Sequences

- The individual optimality criterion has the problem that the optimum state sequence may not obey state transition constraints
- Another optimality criterion is to choose the state sequence which maximizes  $P(\mathbf{Q}, \mathbf{O}|\lambda)$ ; This can be found by the *Viterbi* algorithm
- Let us define  $\delta_t(i)$  as the highest probability along a single path, at time  $t$ , which accounts for the first  $t$  observations, i.e.

$$\delta_t(i) = \max_{q_1, q_2, \dots, q_{t-1}} P(q_1 q_2 \dots q_{t-1}, q_t = s_i, o_1 o_2 \dots o_t | \lambda)$$

- By induction:  $\delta_{t+1}(j) = [\max_i \delta_t(i) a_{ij}] b_j(o_{t+1})$
- To retrieve the state sequence, we must keep track of the state sequence which gave the best path, at time  $t$ , to state  $s_i$ 
  - We do this in a separate array  $\psi_t(i)$

# The Viterbi Algorithm

## 1. Initialization:

$$\begin{aligned}\delta_1(i) &= \pi_i b_i(o_1), & 1 \leq i \leq N \\ \psi_1(i) &= 0\end{aligned}$$

## 2. Recursion:

$$\begin{aligned}\delta_t(j) &= \max_{1 \leq i \leq N} [\delta_{t-1}(i) a_{ij}] b_j(o_t), & 2 \leq t \leq T & \quad 1 \leq j \leq N \\ \psi_t(j) &= \operatorname{argmax}_{1 \leq i \leq N} [\delta_{t-1}(i) a_{ij}], & 2 \leq t \leq T & \quad 1 \leq j \leq N\end{aligned}$$

## 3. Termination:

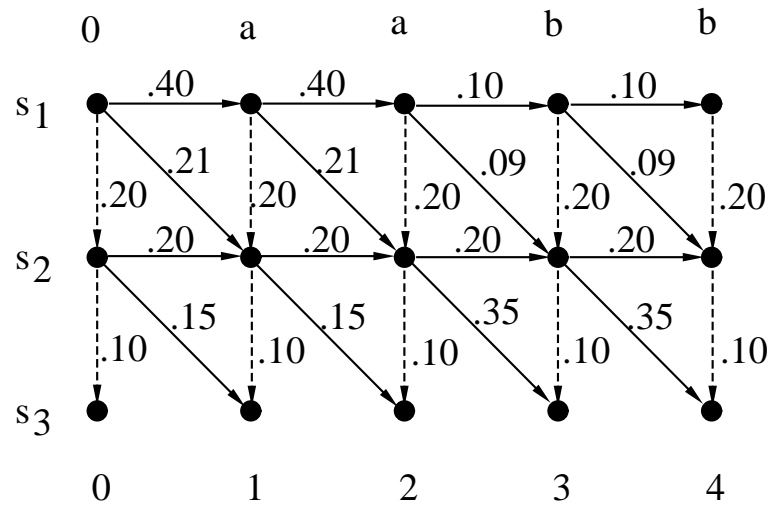
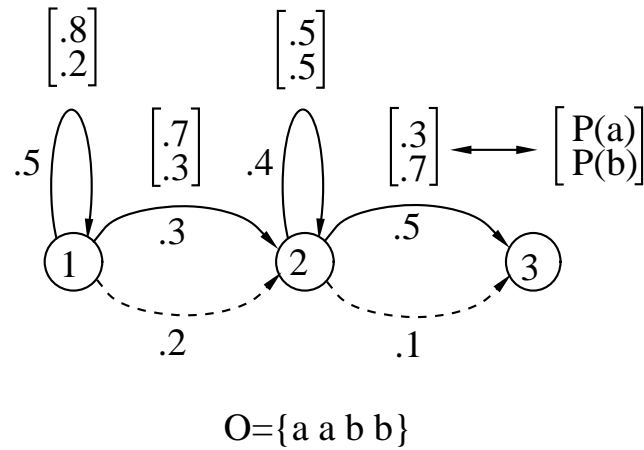
$$\begin{aligned}P^* &= \max_{1 \leq i \leq N} [\delta_T(i)] \\ q_T^* &= \operatorname{argmax}_{1 \leq i \leq N} [\delta_T(i)]\end{aligned}$$

## 4. Path (state-sequence) backtracking:

$$q_t^* = \psi_{t+1}(q_{t+1}^*), \quad t = T - 1, T - 2, \dots, 1$$

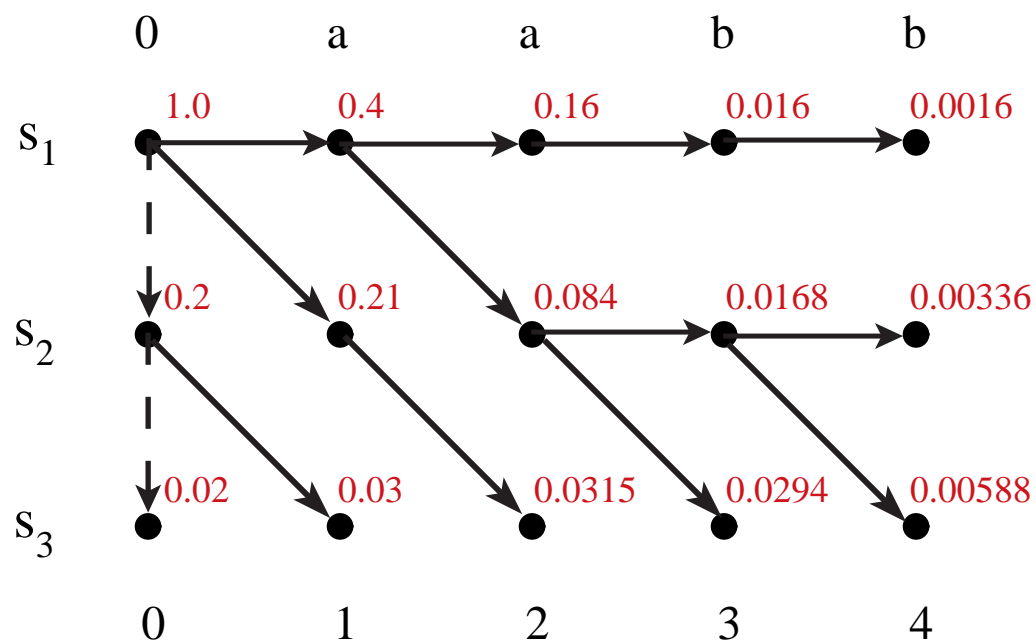
Computation  $\approx N^2 T$

## The Viterbi Algorithm: An Example



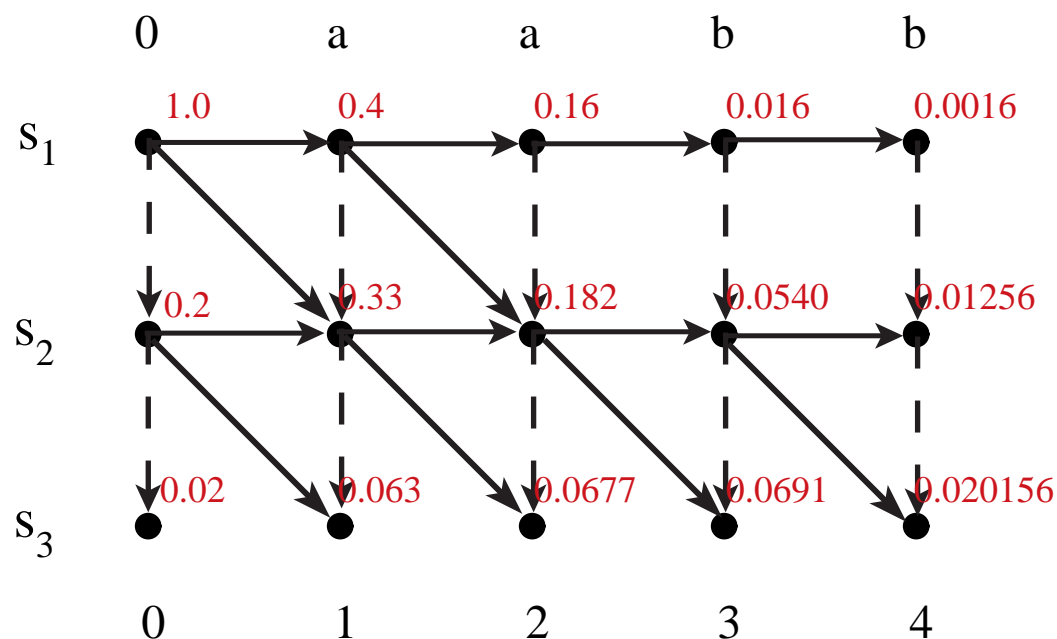
# The Viterbi Algorithm: An Example (cont'd)

	0	<i>a</i>	<i>aa</i>	<i>aab</i>	<i>aabb</i>
$s_1$	1.0	$s_{1,a}$ .4	$s_{1,a}$ .16	$s_{1,b}$ .016	$s_{1,b}$ .0016
$s_2$	$s_{1,0}$ .2	$s_{1,0}$ .08 <b><math>s_{1,a}</math> .21</b> $s_{2,a}$ .04	$s_{1,0}$ .032 <b><math>s_{1,a}</math> .084</b> $s_{2,a}$ .042	$s_{1,0}$ .0032 $s_{1,b}$ .0144 <b><math>s_{2,b}</math> .0168</b>	$s_{1,0}$ .00032 $s_{1,b}$ .00144 <b><math>s_{2,b}</math> .00336</b>
$s_3$	$s_{2,0}$ .02	$s_{2,0}$ .021 <b><math>s_{2,a}</math> .03</b>	$s_{2,0}$ .0084 <b><math>s_{2,a}</math> .0315</b>	$s_{2,0}$ .00168 <b><math>s_{2,b}</math> .0294</b>	$s_{2,0}$ .000336 <b><math>s_{2,b}</math> .00588</b>



# Matching Using Forward-Backward Algorithm

	0	<i>a</i>	<i>aa</i>	<i>aab</i>	<i>aabb</i>
$s_1$	1.0	$s_{1,a}$ .4	$s_{1,a}$ .16	$s_{1,b}$ .016	$s_{1,b}$ .0016
$s_2$	$s_{1,0}$ .2	$s_{1,0}$ .08 $s_{1,a}$ .21 $s_{2,a}$ .04	$s_{1,0}$ .032 $s_{1,a}$ .084 $s_{2,a}$ .066	$s_{1,0}$ .0032 $s_{1,b}$ .0144 $s_{2,b}$ .0364	$s_{1,0}$ .00032 $s_{1,b}$ .00144 $s_{2,b}$ .0108
$s_3$	$s_{2,0}$ .02	$s_{2,0}$ .033 $s_{2,a}$ .03	$s_{2,0}$ .0182 $s_{2,a}$ .0495	$s_{2,0}$ .0054 $s_{2,b}$ .0637	$s_{2,0}$ .001256 $s_{2,b}$ .0189



# Baum-Welch Re-estimation

- Baum-Welch re-estimation uses EM to determine ML parameters
- Define  $\xi_t(i, j)$  as the probability of being in state  $s_i$  at time  $t$  and state  $s_j$  at time  $t + 1$ , given the model and observation sequence

$$\xi_t(i, j) = P(q_t = s_i, q_{t+1} = s_j | \mathbf{O}, \lambda)$$

- Then:

$$\xi_t(i, j) = \frac{\alpha_t(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)}{P(\mathbf{O} | \lambda)}$$

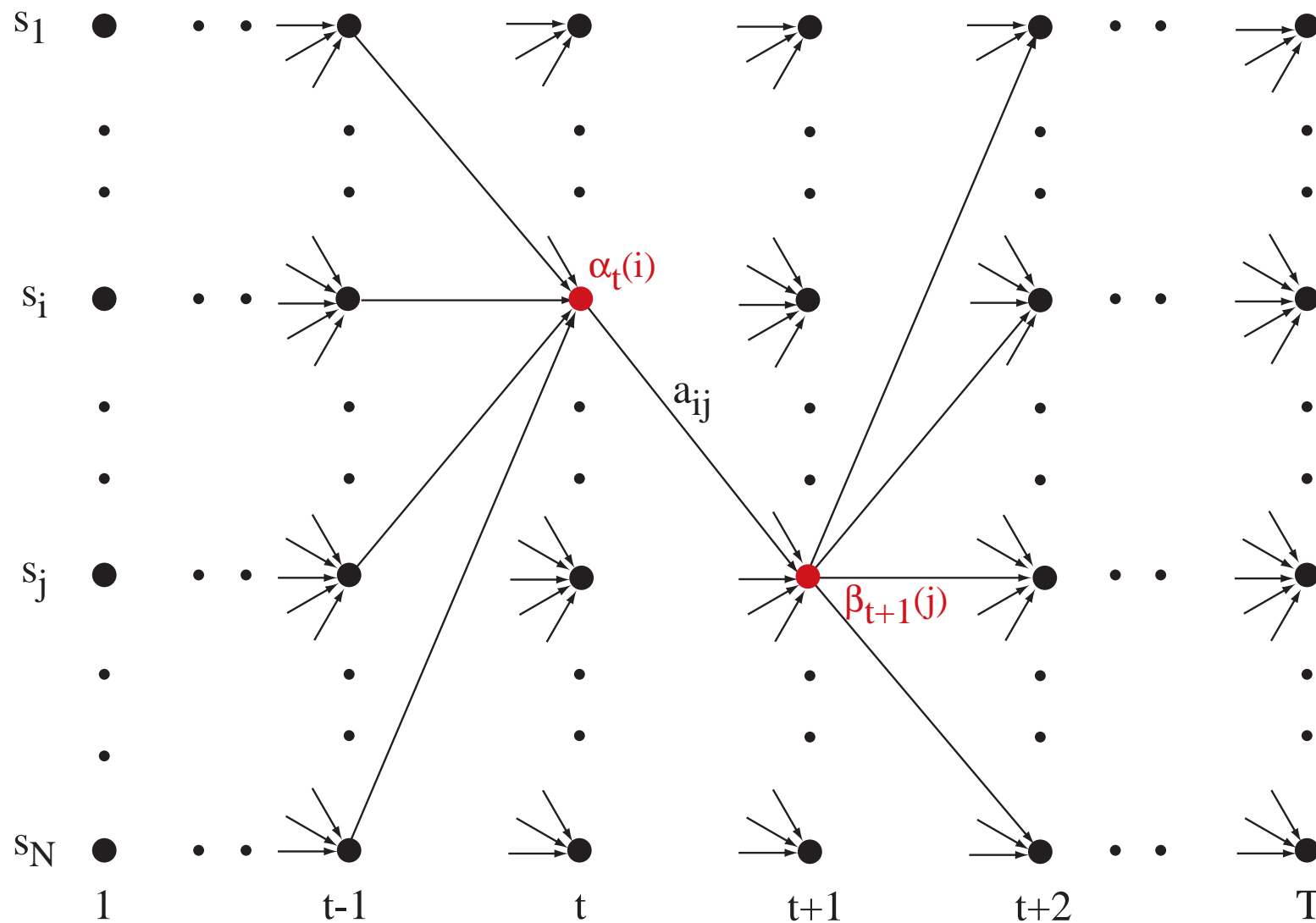
$$\gamma_t(i) = \sum_{j=1}^N \xi_t(i, j)$$

- Summing  $\gamma_t(i)$  and  $\xi_t(i, j)$ , we get:

$$\sum_{t=1}^{T-1} \gamma_t(i) = \text{expected number of transitions from } s_i$$

$$\sum_{t=1}^{T-1} \xi_t(i, j) = \text{expected number of transitions from } s_i \text{ to } s_j$$

# Baum-Welch Re-estimation Procedures



# Baum-Welch Re-estimation Formulas

$$\begin{aligned}\bar{\pi} &= \text{expected number of times in state } s_i \text{ at } t = 1 \\ &= \gamma_1(i)\end{aligned}$$

$$\begin{aligned}\bar{a}_{ij} &= \frac{\text{expected number of transitions from state } s_i \text{ to } s_j}{\text{expected number of transitions from state } s_i} \\ &= \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)}\end{aligned}$$

$$\begin{aligned}\bar{b}_j(k) &= \frac{\text{expected number of times in state } s_j \text{ with symbol } v_k}{\text{expected number of times in state } s_j} \\ &= \frac{\sum_{\substack{t=1 \\ o_t=v_k}}^T \gamma_t(j)}{\sum_{t=1}^T \gamma_t(j)}\end{aligned}$$



# Baum-Welch Re-estimation Formulas

- If  $\lambda = (\mathbf{A}, \mathbf{B}, \pi)$  is the initial model, and  $\bar{\lambda} = (\bar{\mathbf{A}}, \bar{\mathbf{B}}, \bar{\pi})$  is the re-estimated model. Then it can be proved that either:
  1. The initial model,  $\lambda$ , defines a critical point of the likelihood function, in which case  $\bar{\lambda} = \lambda$ , or
  2. Model  $\bar{\lambda}$  is more likely than  $\lambda$  in the sense that  $P(\mathbf{O}|\bar{\lambda}) > P(\mathbf{O}|\lambda)$ , i.e., we have found a new model  $\bar{\lambda}$  from which the observation sequence is more likely to have been produced.
- Thus we can improve the probability of  $\mathbf{O}$  being observed from the model if we iteratively use  $\bar{\lambda}$  in place of  $\lambda$  and repeat the re-estimation until some limiting point is reached. The resulting model is called the maximum likelihood HMM.

# Multiple Observation Sequences

- Speech recognition typically uses left-to-right HMMs. These HMMs can not be trained using a single observation sequence, because only a small number of observations are available to train each state. To obtain reliable estimates of model parameters, one must use multiple observation sequences. In this case, the re-estimation procedure needs to be modified.

- Let us denote the set of  $K$  observation sequences as

$$\mathbf{O} = \{\mathbf{O}^{(1)}, \mathbf{O}^{(2)}, \dots, \mathbf{O}^{(K)}\}$$

where  $\mathbf{O}^{(k)} = \{o_1^{(k)}, o_2^{(k)}, \dots, o_{T_k}^{(k)}\}$  is the  $k$ -th observation sequence.

- Assume that the observations sequences are mutually independent, we want to estimate the parameters so as to maximize

$$P(\mathbf{O} | \lambda) = \prod_{k=1}^K P(\mathbf{O}^{(k)} | \lambda) = \prod_{k=1}^K P_k$$

## Multiple Observation Sequences (cont'd)

- Since the re-estimation formulas are based on frequency of occurrence of various events, we can modify them by adding up the individual frequencies of occurrence for each sequence

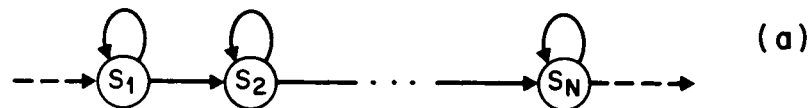
$$\bar{a}_{ij} = \frac{\sum_{k=1}^K \sum_{t=1}^{T_k-1} \xi_t^k(i, j)}{\sum_{k=1}^K \sum_{t=1}^{T_k-1} \gamma_t^k(i)} = \frac{\sum_{k=1}^K \frac{1}{P_k} \sum_{t=1}^{T_k-1} \alpha_t^k(i) a_{ij} b_j(o_{t+1}^{(k)}) \beta_{t+1}^k(j)}{\sum_{k=1}^K \frac{1}{P_k} \sum_{t=1}^{T_k-1} \alpha_t^k(i) \beta_t^k(i)}$$

$$\bar{b}_j(\ell) = \frac{\sum_{k=1}^K \sum_{\substack{t=1 \\ o_t^{(k)} = v_\ell}}^{T_k} \gamma_t^k(j)}{\sum_{k=1}^K \sum_{t=1}^{T_k} \gamma_t^k(j)} = \frac{\sum_{k=1}^K \frac{1}{P_k} \sum_{\substack{t=1 \\ o_t^{(k)} = v_\ell}}^{T_k} \alpha_t^k(i) \beta_t^k(i)}{\sum_{k=1}^K \frac{1}{P_k} \sum_{t=1}^{T_k} \alpha_t^k(i) \beta_t^k(i)}$$

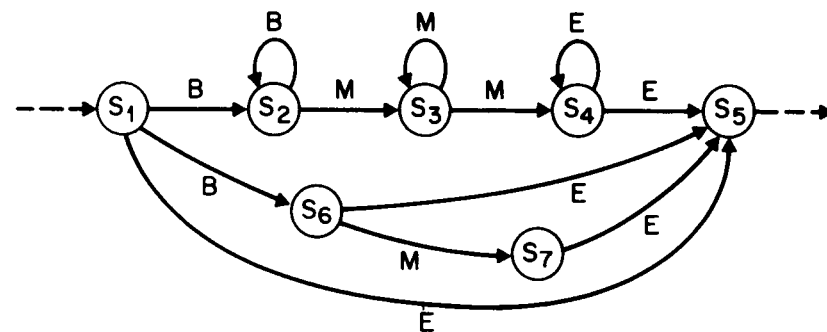
# Phone-based HMMs

- Word-based HMMs are appropriate for small vocabulary speech recognition. For large vocabulary ASR, sub-word-based (e.g., phone-based) models are more appropriate.

WORD MODEL

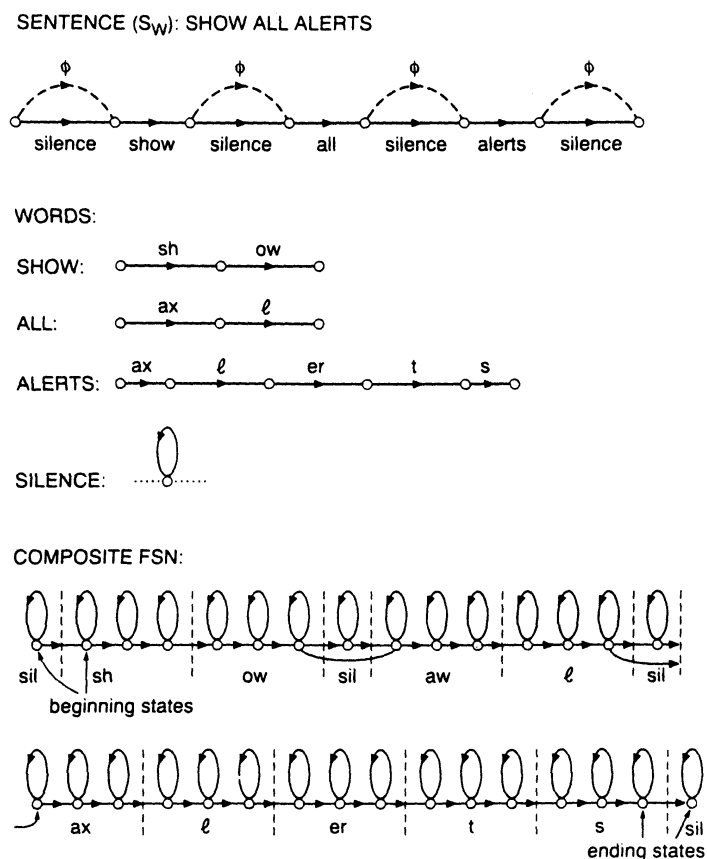


SUB-WORD UNIT



# Phone-based HMMs (cont'd)

- The phone models can have many states, and words are made up from a concatenation of phone models.



# Continuous Density Hidden Markov Models

- A *continuous density* HMM replaces the discrete observation probabilities,  $b_j(k)$ , by a continuous PDF  $b_j(\mathbf{x})$
- A common practice is to represent  $b_j(\mathbf{x})$  as a mixture of Gaussians:

$$b_j(\mathbf{x}) = \sum_{k=1}^M c_{jk} N[\mathbf{x}, \boldsymbol{\mu}_{jk}, \boldsymbol{\Sigma}_{jk}] \quad 1 \leq j \leq N$$

where  $c_{jk}$  is the mixture weight

$$c_{jk} \geq 0 \quad (1 \leq j \leq N, 1 \leq k \leq M, \text{ and } \sum_{k=1}^M c_{jk} = 1, 1 \leq j \leq N),$$

$N$  is the normal density, and

$\boldsymbol{\mu}_{jk}$  and  $\boldsymbol{\Sigma}_{jk}$  are the mean vector and covariance matrix associated with state  $j$  and mixture  $k$ .

# Acoustic Modelling Variations

- *Semi-continuous* HMMs first compute a VQ codebook of size  $M$ 
  - The VQ codebook is then modelled as a family of Gaussian PDFs
  - Each codeword is represented by a Gaussian PDF, and may be used together with others to model the acoustic vectors
  - From the CD-HMM viewpoint, this is equivalent to using the same set of  $M$  mixtures to model all the states
  - It is therefore often referred to as a *Tied Mixture* HMM
- All three methods have been used in many speech recognition tasks, with varying outcomes
- For large-vocabulary, continuous speech recognition with sufficient amount (i.e., tens of hours) of training data, CD-HMM systems currently yield the best performance, but with considerable increase in computation

# Implementation Issues

- Scaling: to prevent underflow
- *Segmental K-means Training*: to train observation probabilities by first performing Viterbi alignment
- Initial estimates of  $\lambda$ : to provide robust models
- Pruning: to reduce search computation



# MIT

## References

- X. Huang, A. Acero, and H. Hon, *Spoken Language Processing*, Prentice-Hall, 2001.
- F. Jelinek, *Statistical Methods for Speech Recognition*. MIT Press, 1997.
- L. Rabiner and B. Juang, *Fundamentals of Speech Recognition*, Prentice-Hall, 1993.