

# **HIDDEN MARKOV MODELS IN SPEECH RECOGNITION**

**Wayne Ward**

**Carnegie Mellon University  
Pittsburgh, PA**

# **Acknowledgements**

**Much of this talk is derived from the paper  
"An Introduction to Hidden Markov Models",  
by Rabiner and Juang**

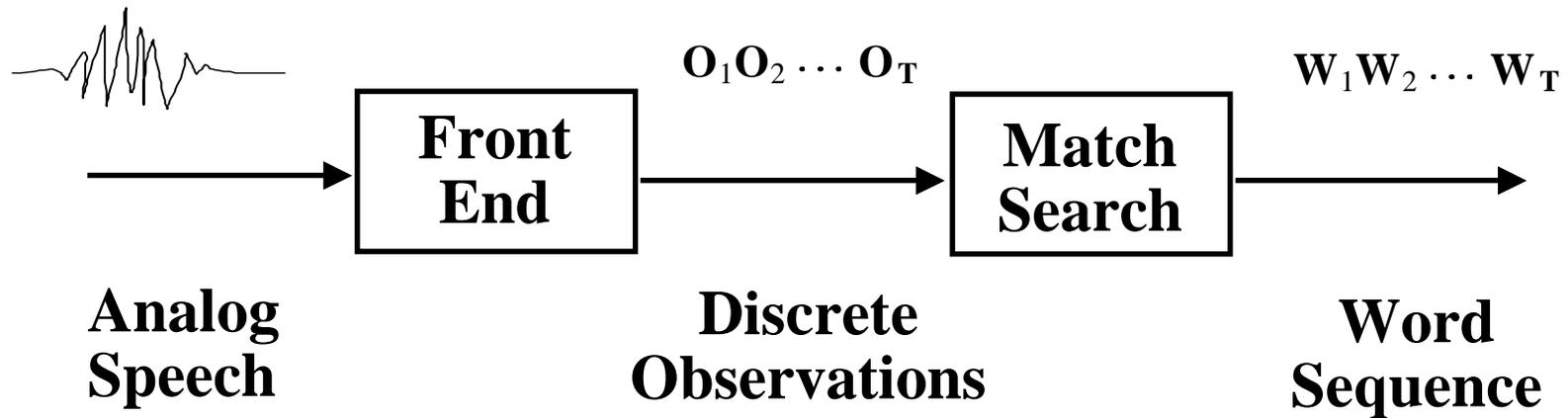
**and from the talk**

**"Hidden Markov Models: Continuous Speech  
Recognition"  
by Kai-Fu Lee**

# Topics

- **MarkovModelsandHiddenMarkovModels**
- **HMMs appliedtospeechrecognition**
  - **Training**
  - **Decoding**

# Speech Recognition



# ML Continuous Speech Recognition

## Goal:

Given acoustic data  $A = a_1, a_2, \dots, a_k$

Find word sequence  $W = w_1, w_2, \dots, w_n$

Such that  $P(W|A)$  is maximized

## Bayes Rule:

$$P(W|A) = \frac{\overset{\text{acoustic model (HMMs)}}{\downarrow} P(A|W) \cdot \overset{\text{language model}}{\swarrow} P(W)}{P(A)}$$

$P(A)$  is a constant for a complete sentence

# Markov Models

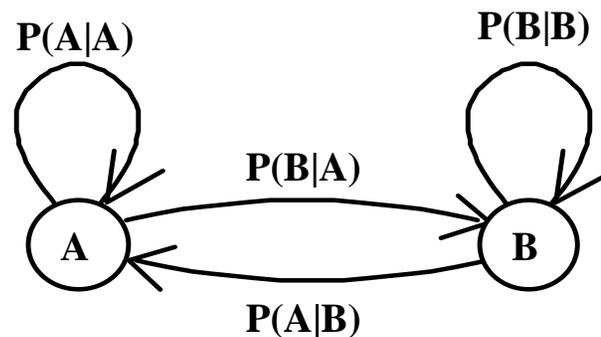
**Elements:**

**States:**

$$\mathbf{S} = \{S_0, S_1, \dots, S_N\}$$

**Transition probabilities:**

$$P(q_t = S_i | q_{t-1} = S_j)$$



**Markov Assumption:**

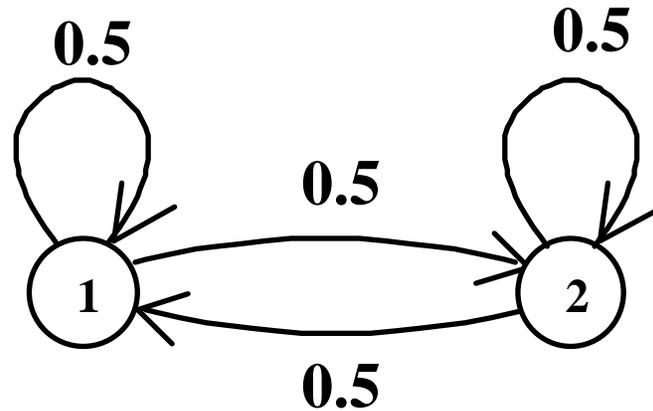
**Transition probability depends only on current state**

$$P(q_t = S_i | q_{t-1} = S_j, q_{t-2} = S_k, \dots) = P(q_t = S_i | q_{t-1} = S_j) = a_{ji}$$

$$a_{ji} \geq 0 \quad \forall j, i$$

$$\sum_{i=0}^N a_{ji} = 1 \quad \forall j$$

# SingleFairCoin



$$P(H)=1.0$$

$$P(H)=0.0$$

$$P(T)=0.0$$

$$P(T)=1.0$$

**Outcome head correspond to state 1, tail to state 2**

**Observation sequence uniquely defines state sequence**

# Hidden Markov Models

## Elements:

**States**

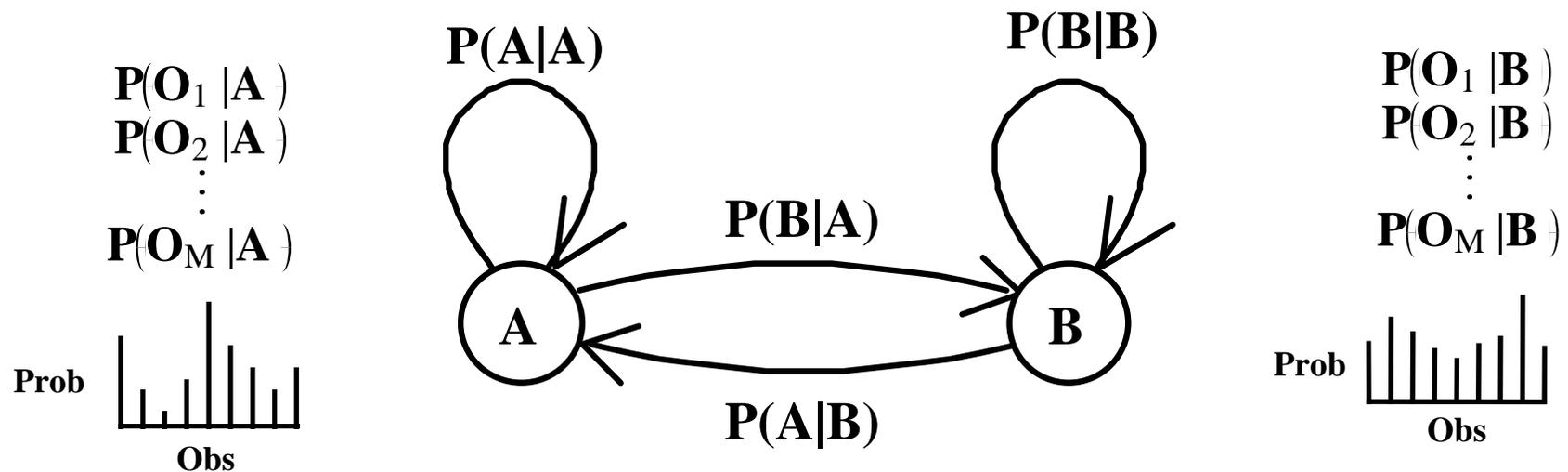
$$S = \{S_0, S_1, \dots, S_N\}$$

**Transition probabilities**

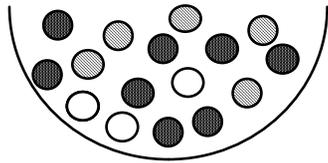
$$P(q_t = S_i | q_{t-1} = S_j) = a_{ji}$$

**Output prob distributions  
(at state  $j$  for symbol  $k$ )**

$$P(y_t = O_k | q_t = S_j) = b_j(k)$$



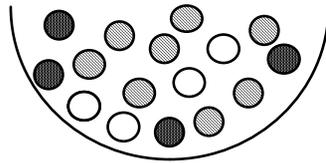
# Discrete Observation HMM



$$P(R)=0.31$$

$$P(B)=0.50$$

$$P(Y)=0.19$$

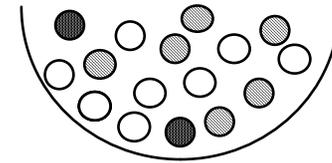


$$P(R)=0.50$$

$$P(B)=0.25$$

$$P(Y)=0.25$$

...



$$P(R)=0.38$$

$$P(B)=0.12$$

$$P(Y)=0.50$$

Observation sequence: **RBYY...R**

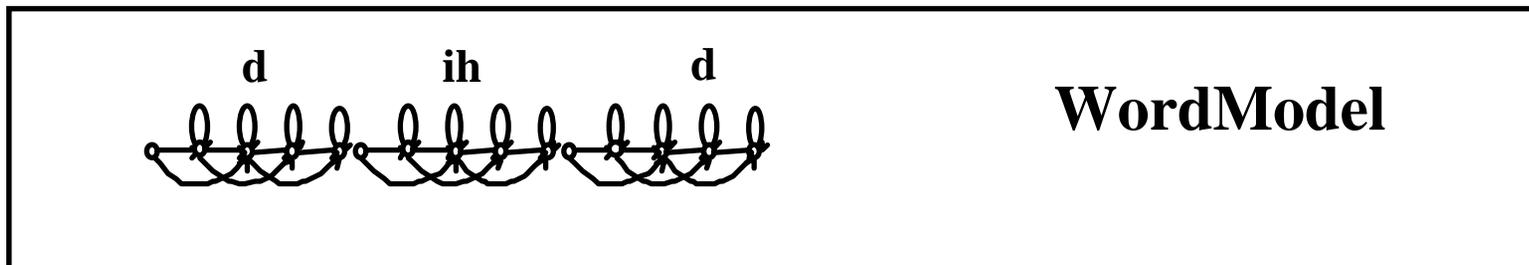
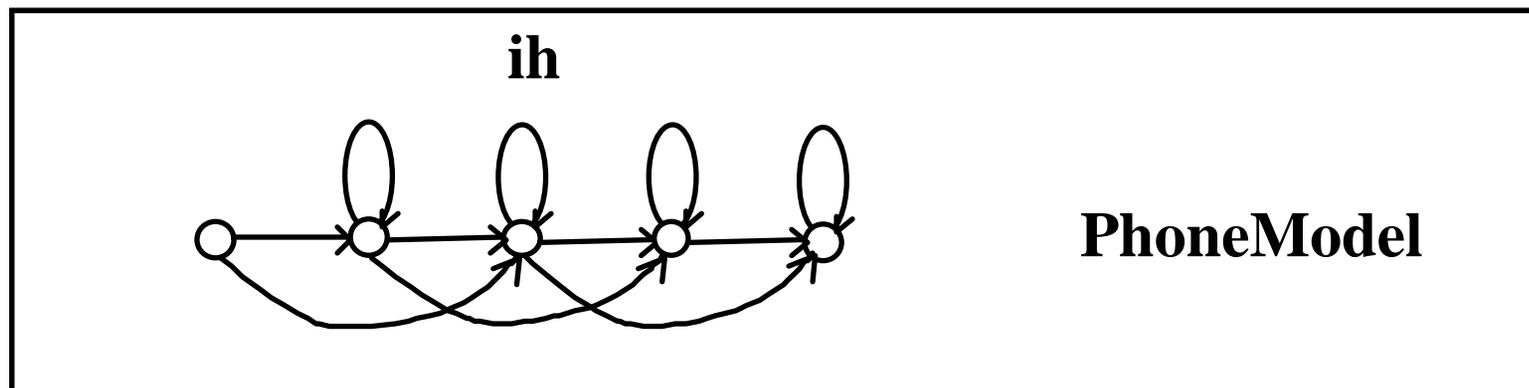
not unique to state sequence

# HMMs InSpeechRecognition

Representspeechasasequenceofobservations

UseHMMtomodelsomeunitofspeech(phone,word)

Concatenateunitsintolargerunits



# HMM Problems And Solutions

## Evaluation:

- Problem- Compute Probability of observation sequence given a model
- Solution- **Forward Algorithm** and **Viterbi Algorithm**

## Decoding:

- Problem- Find state sequence which maximizes probability of observation sequence
- Solution- **Viterbi Algorithm**

## Training:

- Problem- Adjust model parameters to maximize probability of observed sequences
- Solution- **Forward-Backward Algorithm**

# Evaluation

**Probability of observation sequence**  $O = O_1 O_2 \dots O_T$

**given HMM model  $\lambda$  is:**

$$P(O | \lambda) = \sum_{\forall Q} P(Q | \lambda) \quad Q = q_0 q_1 \dots q_T \text{ is a state sequence}$$

$$= \sum a_{q_0 q_1} b_{q_1}(O_1) \times a_{q_1 q_2} b_{q_2}(O_2) \dots a_{q_{T-1} q_T} b_{q_T}(O_T)$$

**Not practical since the number of paths is**  $O(N^T)$

$N$  = number of states in model

$T$  = number of observations in sequence

# The Forward Algorithm

$$\alpha_t(j) = P(O_1 O_2 \cdots O_t, q_t = S_j | \lambda)$$

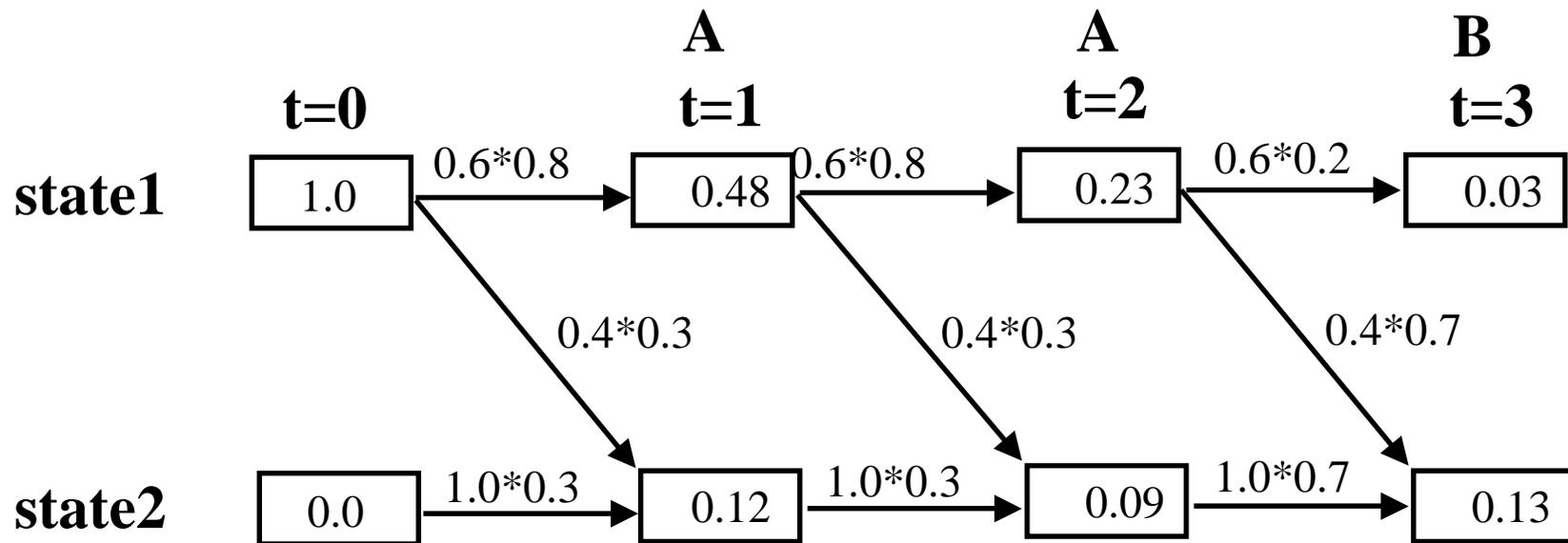
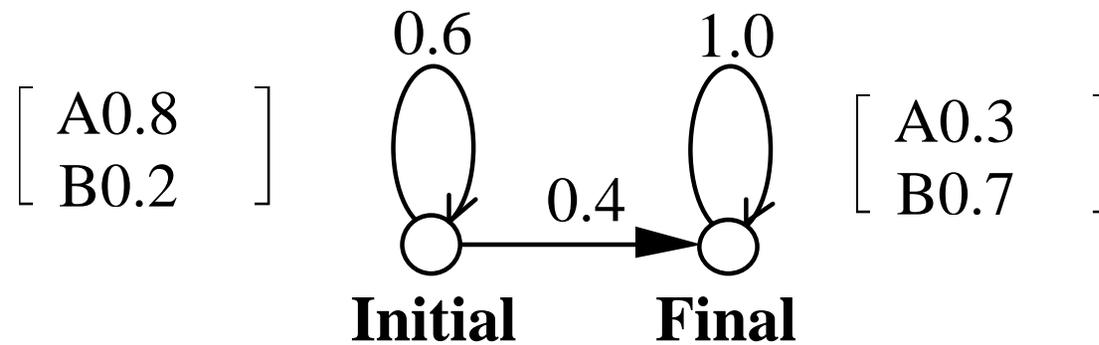
Compute  $\alpha$  recursively:

$$\alpha_0(j) = \begin{cases} \mathbf{1} & \text{if } j \text{ is start state} \\ \mathbf{0} & \text{otherwise} \end{cases}$$

$$\alpha_t(j) = \left[ \sum_{i=0}^N \alpha_{t-1}(i) a_{ij} \right] b_j(O_t) \quad t > 0$$

$$P(O | \lambda) = \alpha_T(S_N) \quad \text{Computation is } O(N^2 T)$$

# ForwardTrellis



# The Backward Algorithm

$$\beta_t(i) = P(O_{t+1} O_{t+2} \cdots O_T | q_t = S_i, \lambda)$$

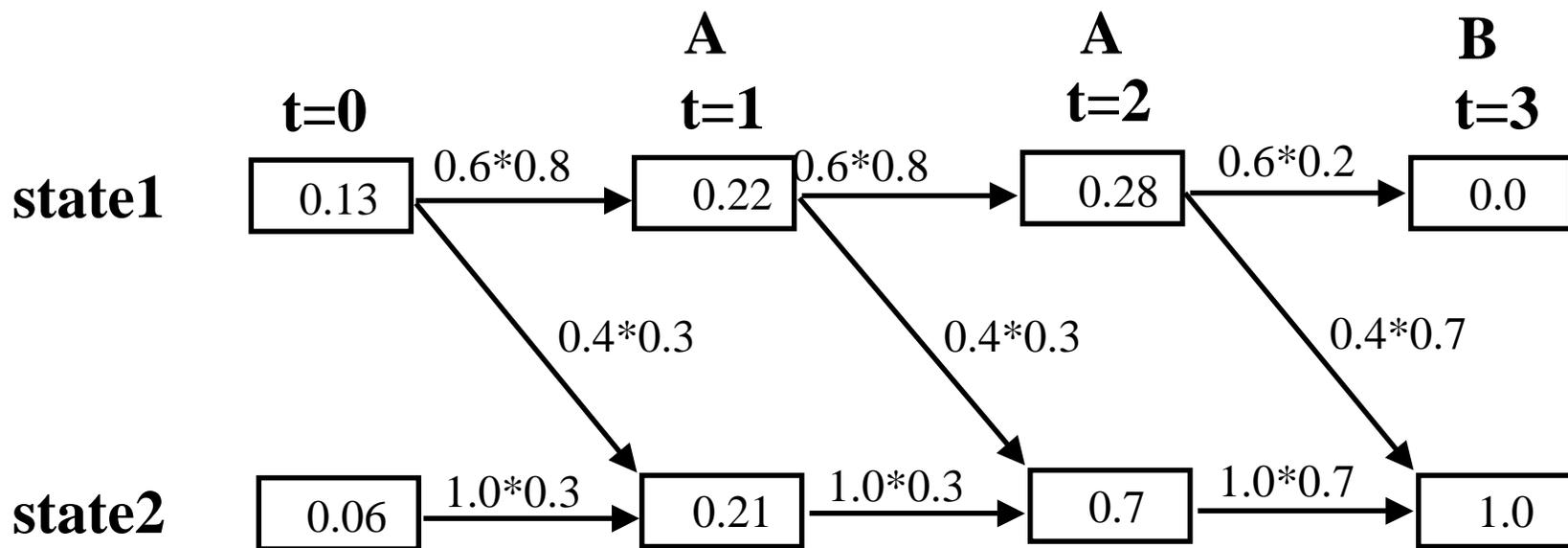
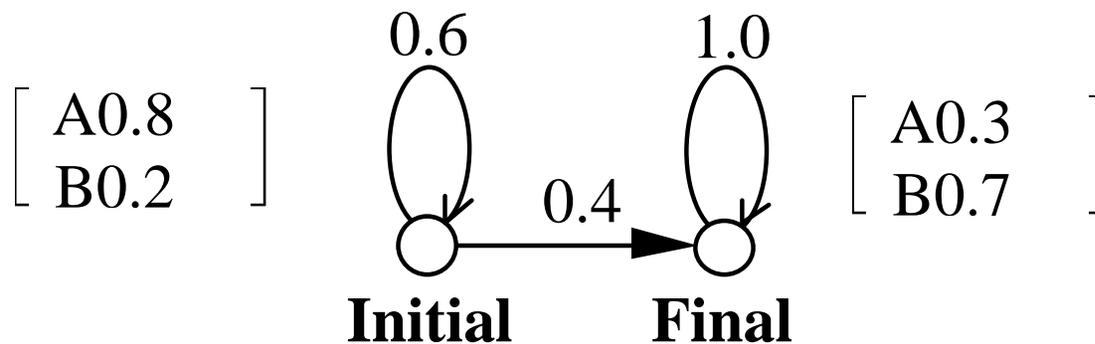
Compute  $\beta$  recursively:

$$\beta_T(i) = \begin{cases} \mathbf{1} & \text{if } i \text{ is end state} \\ \mathbf{0} & \text{otherwise} \end{cases}$$

$$\beta_t(i) = \sum_{j=0}^N a_{ij} b_j(O_{t+1}) \beta_{t+1}(j) \quad t < T$$

$$P(O | \lambda) = \beta_0(S_0) = \alpha_T(S_N) \quad \text{Computation is } O(N^2T)$$

# BackwardTrellis



# The Viterbi Algorithm

For decoding:

Find the state sequence  $Q$  which maximizes  $P(O, Q | \lambda)$

Similar to Forward Algorithm except **MAX** instead of **SUM**

$$VP_t(i) = \text{MAX}_{q_0, \dots, q_{t-1}} P(O_1 O_2 \dots O_t, q_t = i | \lambda)$$

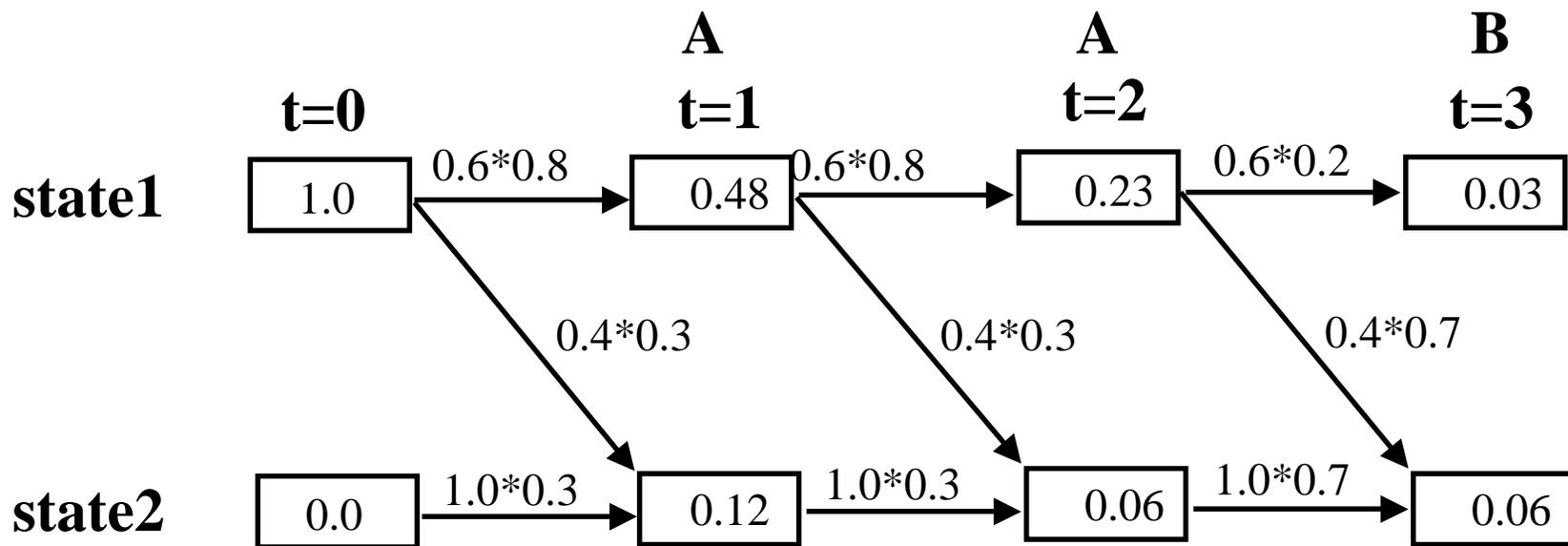
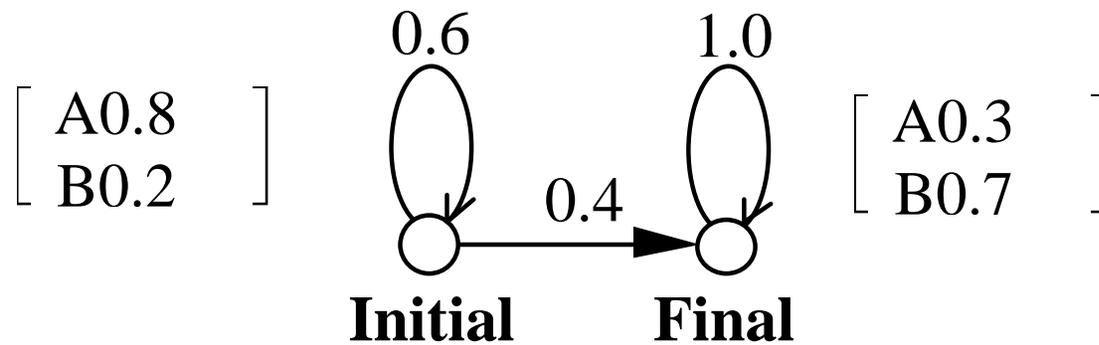
Recursive Computation:

$$VP_t(j) = \text{MAX}_{i=0, \dots, N} VP_{t-1}(i) a_{ij} b_j(O_t) \quad t > 0$$

$$P(O, Q | \lambda) = V \quad P_T(S_N)$$

Save each maximum for backtrace at end

# Viterbi Trellis



# Training HMM Parameters

## Train parameters of HMM

- Tune  $\lambda$  to maximize  $P(O | \lambda)$
- No efficient algorithm for global optimum
- Efficient iterative algorithm finds a local optimum

## Baum-Welch (Forward-Backward) re-estimation

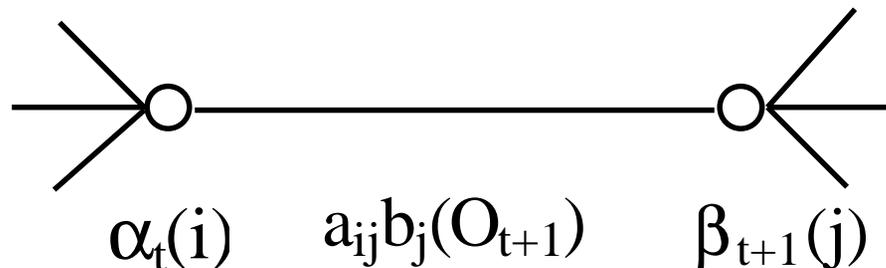
- Compute probabilities using current model  $\lambda$
- Refine  $\lambda \rightarrow \hat{\lambda}$  based on computed values
- Use  $\alpha$  and  $\beta$  from Forward-Backward

# Forward-Backward Algorithm

$\xi_t(i,j)$  = Probability of transiting from  $S_i$  to  $S_j$  at time  $t$  given  $O$

$$= P(q_t = S_i, q_{t+1} = S_j | O, \lambda)$$

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



# Baum-Welch Reestimation

$$\bar{a}_{ij} = \frac{\text{expected number of transitions from } S_i \text{ to } S_j}{\text{expected number of transitions from } S_i}$$

$$= \frac{\sum_{t=0}^{T-1} \xi_t(i, j)}{\sum_{t=0}^{T-1} \sum_{j=0}^N \xi_t(i, j)}$$

$$\bar{b}_j(k) = \frac{\text{expected number of times in state } j \text{ with symbol } k}{\text{expected number of times in state } j}$$

$$= \frac{\sum_{t: O_t = k} \sum_{i=0}^N \xi_t(i, j)}{\sum_{t=0}^{T-1} \sum_{i=0}^N \xi_t(i, j)}$$

# Convergence of FB Algorithm

1. Initialize  $\lambda = (A, B)$
2. Compute  $\alpha$ ,  $\beta$ , and  $\xi$
3. Estimate  $\bar{\lambda} = (\bar{A}, \bar{B})$  from  $\xi$
4. Replace  $\lambda$  with  $\bar{\lambda}$
5. If not converged goto 2

It can be shown that  $P(O | \bar{\lambda}) > P(O | \lambda)$  unless  $\bar{\lambda} = \lambda$

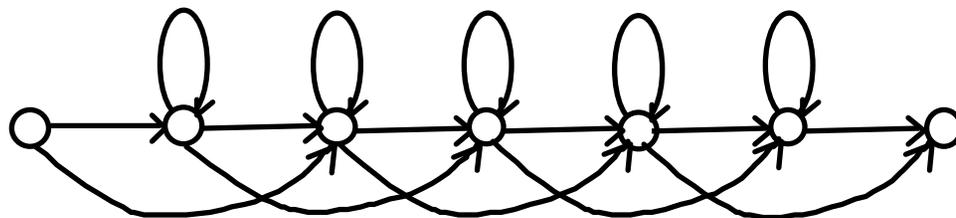
# HMMs InSpeechRecognition

**Representspeechasasequenceofsymbols**

**UseHMMtomodelsomeunitofspeech(phone,word)**

**OutputProbabilities- Prob ofobservingsymbolinastate**

**Transition Prob - Prob ofstayinginorskippingstate**

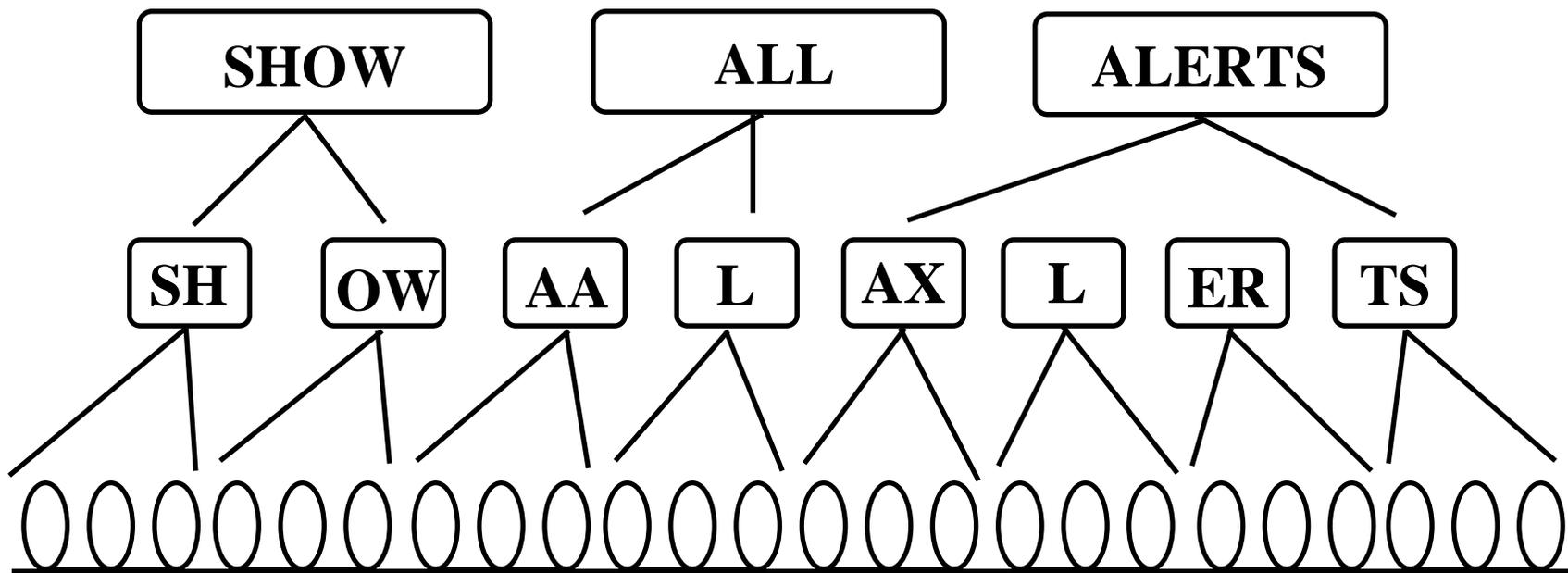


**PhoneModel**

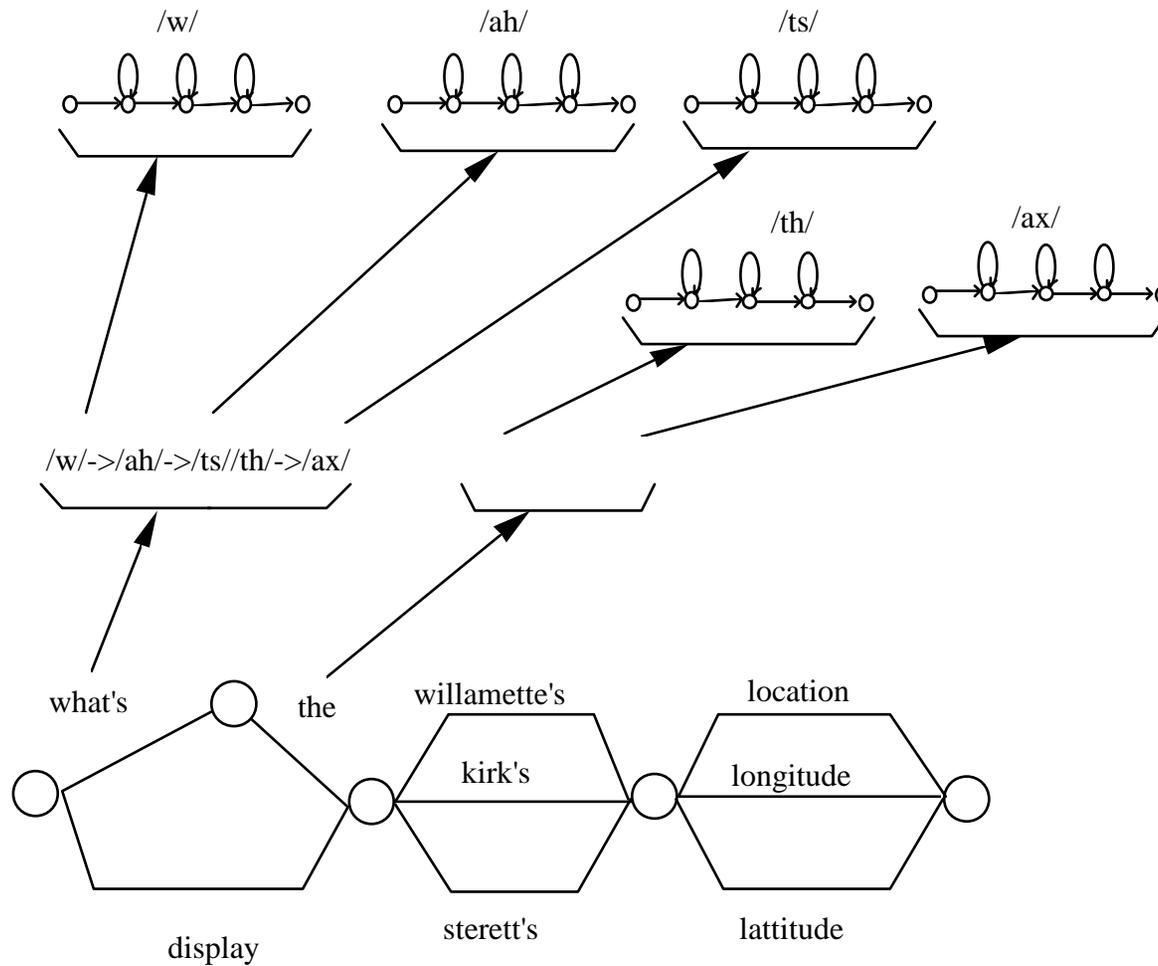
# Training HMMs for Continuous Speech

- Use only orthograph transcription of sentence
  - no need for segmented/labelled data
- Concatenate phoneme models to give word model
- Concatenate word models to give sentence model
- Train entire sentence model on entire spoken sentence

# Forward-Backward Training for Continuous Speech



# RecognitionSearch



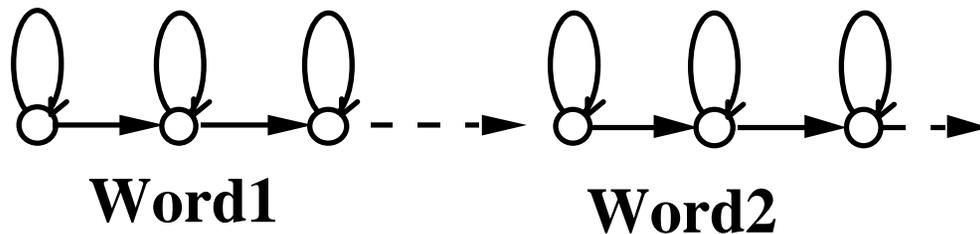
# Viterbi Search

- **Uses Viterbi decoding**
  - **Takes MAX, not SUM**
  - **Finds optimal state sequence  $P(O, Q | \lambda)$**   
**not optimal word sequence  $P(O | \lambda)$**
- **Time synchronous**
  - **Extends all paths by 1 timestep**
  - **All paths have same length (no need to normalize to compare scores)**

# Viterbi Search Algorithm

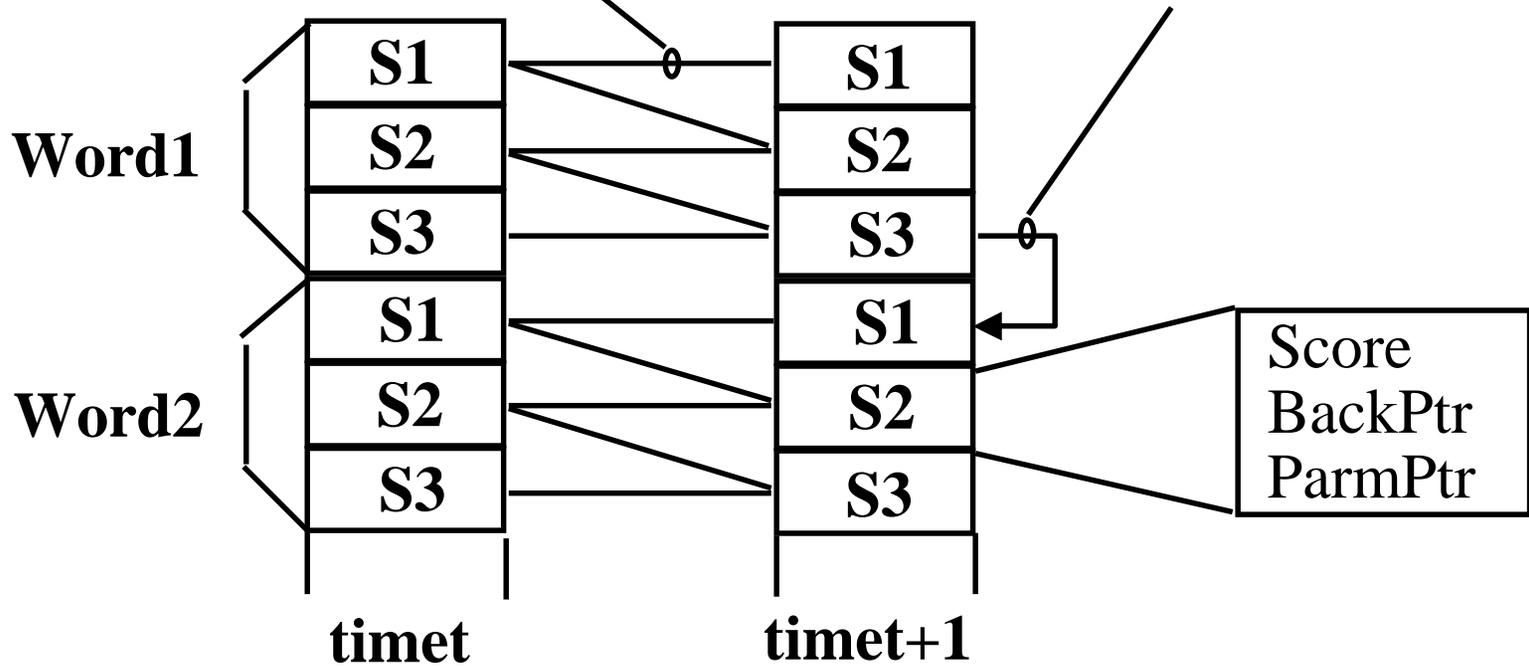
0. Create statelist with one cell for each state in system
1. Initialize statelist with initial states for  $t=0$
2. Clear statelist for  $t+1$
3. Compute within-word transitions from  $t$  to  $t+1$ 
  - If new state reached, update score and BackPtr
  - If better score for state, update score and BackPtr
4. Compute between word transitions at  $t+1$ 
  - If new state reached, update score and BackPtr
  - If better score for state, update score and BackPtr
5. If end of utterance, print backtrace and quit
6. Else increment  $t$  and go to step 2

# Viterbi Search Algorithm



OldProb(S1) • OutProb • Transprob

OldProb(S3) • P(W2|W1)



# Viterbi BeamSearch

## Viterbi Search

All states enumerated

Not practical for large grammars

Most states inactive at any given time

## Viterbi BeamSearch- prune less likely paths

States worse than threshold range from best are pruned

From and To structures created dynamically- list of active states

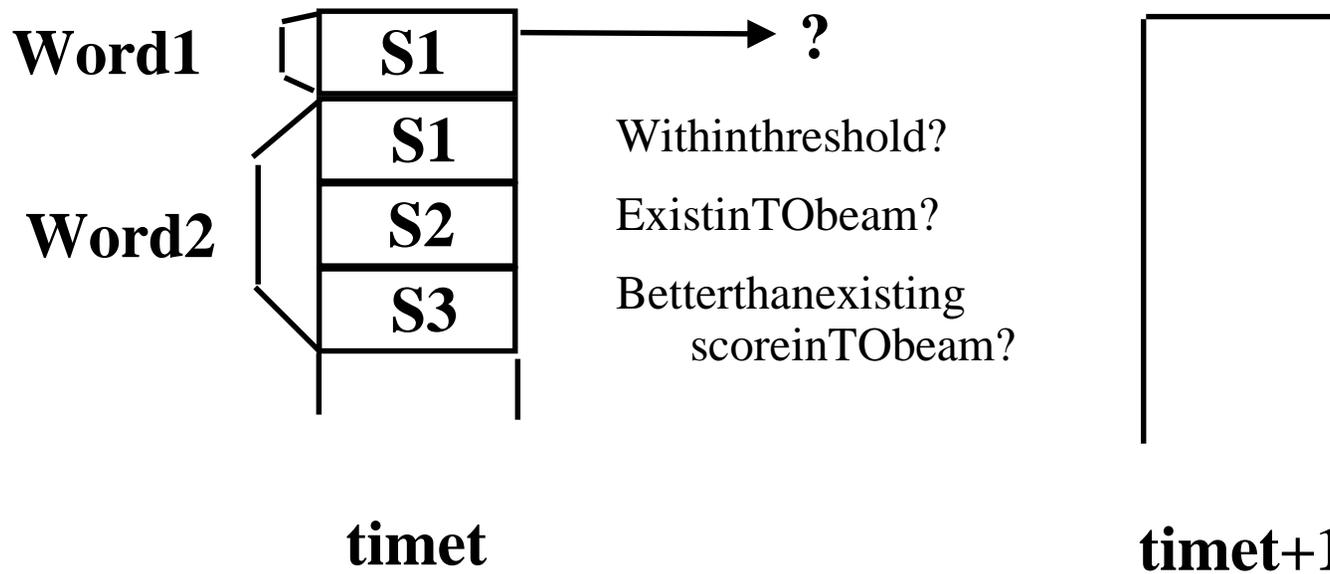
# Viterbi BeamSearch

**FROM  
BEAM**

States within threshold  
from best state

**TO  
BEAM**

Dynamically constructed



# Continuous Density HMMs

Model so far has assumed discrete observations, each observation in a sequence was one of a set of  $M$  discrete symbols

Speech input must be Vector Quantized in order to provide discrete input.

VQ leads to quantization error

The discrete probability density  $b_j(k)$  can be replaced with the continuous probability density  $b_j(\mathbf{x})$  where  $\mathbf{x}$  is the observation vector

Typically Gaussian densities are used

A single Gaussian is not adequate, so a weighted sum of Gaussians is used to approximate actual PDF

# Mixture Density Functions

$b_j(\mathbf{x})$  is the probability density function for state  $j$

$$b_j(x) = \sum_{m=1}^M c_{jm} N[x, \mu_{jm}, U_{jm}]$$

$\mathbf{x}$  = Observation vector  $x_1, x_2, \dots, x_D$

$M$  = Number of mixtures (Gaussians)

$c_{jm}$  = Weight of mixture  $m$  in state  $j$  where

$$\sum_{m=1}^M c_{jm} = 1$$

$N$  = Gaussian density function

$\mu_{jm}$  = Mean vector for mixture  $m$ , state  $j$

$U_{jm}$  = Covariance matrix for mixture  $m$ , state  $j$

# DiscreteHmmvs.ContinuousHMM

## □ ProblemswithDiscrete:

- quantization errors
- Codebookand HMMsmodelled separately

## • ProblemswithContinuousMixtures:

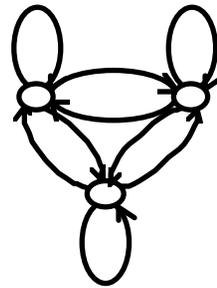
- Smallnumberofmixturesperforms poorly
- Largenumberofmixturesincreasescomputation andparameterstobeestimated

$$c_{jm}, \mu_{jm}, U_{jm} \text{ for } j=1, \dots, N \text{ and } m=1, \dots, M$$

- ContinuousmakesmoreassumptionsthanDiscrete, especiallyifdiagonalcovariance pdf
- Discreteprobabilityisatablelookup,continuous mixturesrequiremanymultiplications

# Model Topologies

**Ergodic-** Fully connected, each state has transition to every other state



**Left-to-Right-** Transitions only to states with higher index than current state. Inherently impose temporal order. These most often used for speech.

