# HIDDEN MARKOV MODELS IN SPEECH RECOGNITION

**Wayne Ward**

**Carnegie Mellon University**

**Pittsburgh, PA**

# Acknowledgements

Much of this talk is derived from the paper

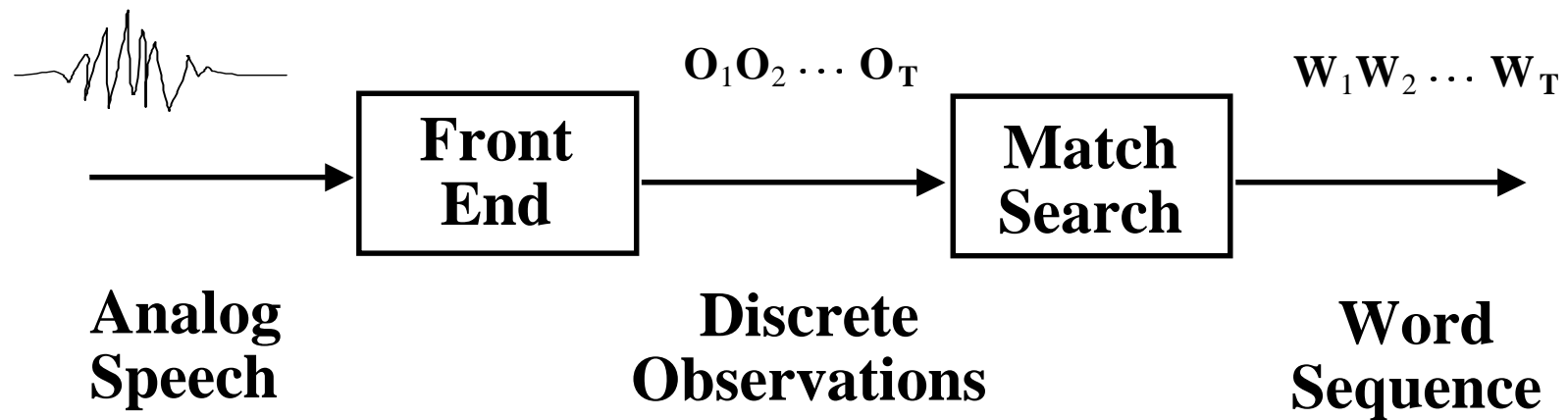"An Introduction to Hidden Markov Models",

by Rabiner and Juang

and from the talk

"Hidden Markov Models: Continuous Speech
    Recognition"

by Kai-Fu Lee

# Topics

- **Markov Models and Hidden Markov Models**

- **HMMs applied to speech recognition**

  - **Training**

  - **Decoding**

# Speech Recognition



$$O_1 O_2 \cdots O_T \qquad W_1 W_2 \cdots W_T$$

**Front End** → **Match Search**

**Analog Speech**     **Discrete Observations**     **Word Sequence**

# ML Continuous Speech Recognition

**Goal:**

> Given acoustic data $A = a_1, a_2, ..., a_k$
>
> Find word sequence $W = w_1, w_2, ... w_n$
>
> Such that $P(W \mid A)$ is maximized

**Bayes Rule:**

acoustic model (HMMs)　　　　　　　　language model

$$P(W \mid A) = \frac{P(A \mid W) \cdot P(W)}{P(A)}$$

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

# **Markov Models**

**Elements:**
 **States :** $\quad S = \{S_0, S_1, \cdots S_N\}$
 **Transition probabilities :** $\quad P(q_t = S_i \mid q_{t-1} = S_j)$



**P(A | A)**  **P(B | B)**

**P(B | A)**
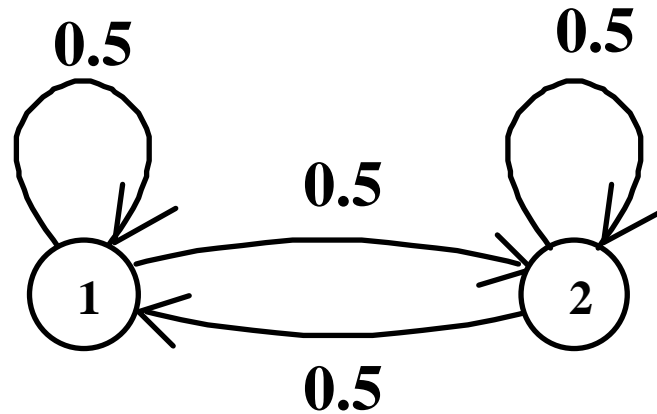
**A**  **B**

**P(A | B)**

**Markov Assumption:**
 **Transition probability depends only on current state**

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

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

# Single Fair Coin



0.5             0.5

0.5

1          2

0.5

P(H) = 1.0          P(H) = 0.0

P(T) = 0.0          P(T) = 1.0

Outcome head corresponds to state 1, tail to state 2

Observation sequence uniquely defines state sequence
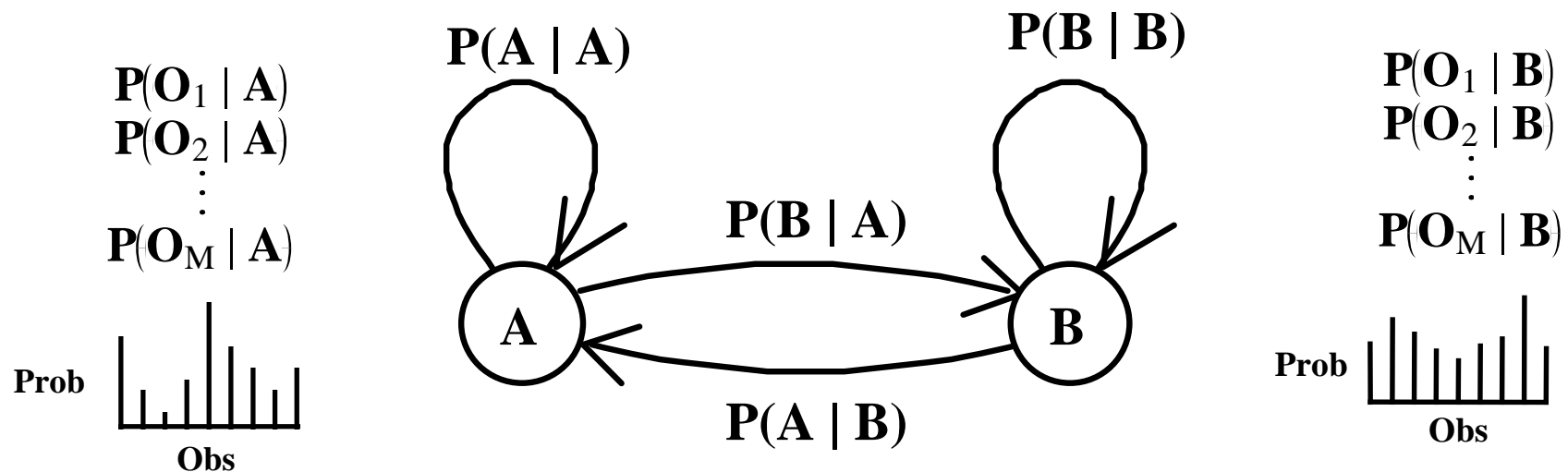
# Hidden Markov Models

**Elements:**

   **States**                                     $S = \{S_0, S_1, \ldots S_N\}$
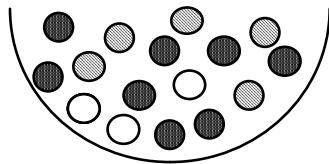
   **Transition probabilities**           $P(q_t = S_i \mid q_{t-1} = S_j) = a_{ji}$

   **Output prob distributions**        $P(y_t = O_k \mid q_t = S_j) = b_j(k)$
      **(at state j for symbol k)**

$P(O_1 \mid A)$
$P(O_2 \mid A)$
$\vdots$
$P(O_M \mid A)$

**Prob**

**Obs**

**P(A | A)**

**P(B | A)**

**A**

**B**

**P(A | B)**

**P(B | B)**

$P(O_1 \mid B)$
$P(O_2 \mid B)$
$\vdots$
$P(O_M \mid B)$

**Prob**

**Obs**

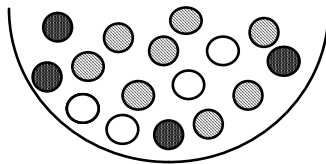# Discrete Observation HMM

P(R) = 0.31

P(B) = 0.50

P(Y) = 0.19
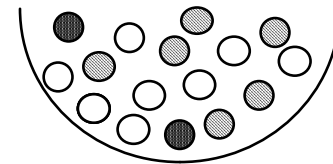
P(R) = 0.50

P(B) = 0.25

P(Y) = 0.25

P(R) = 0.38

P(B) = 0.12

P(Y) = 0.50

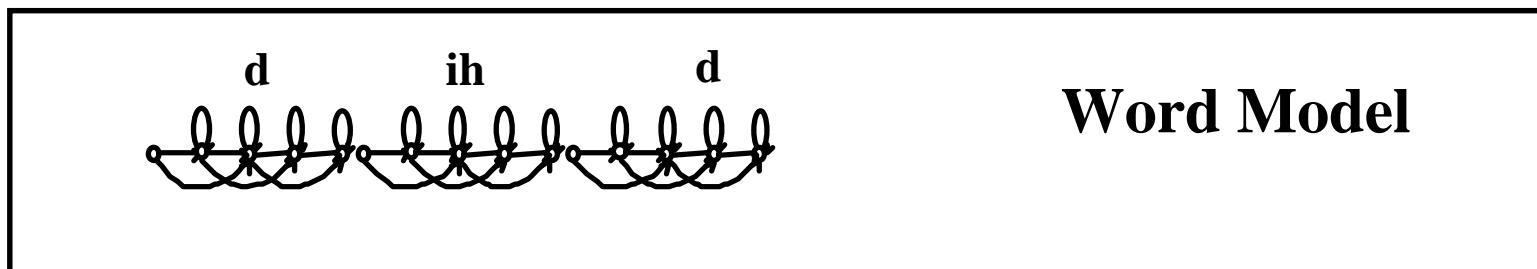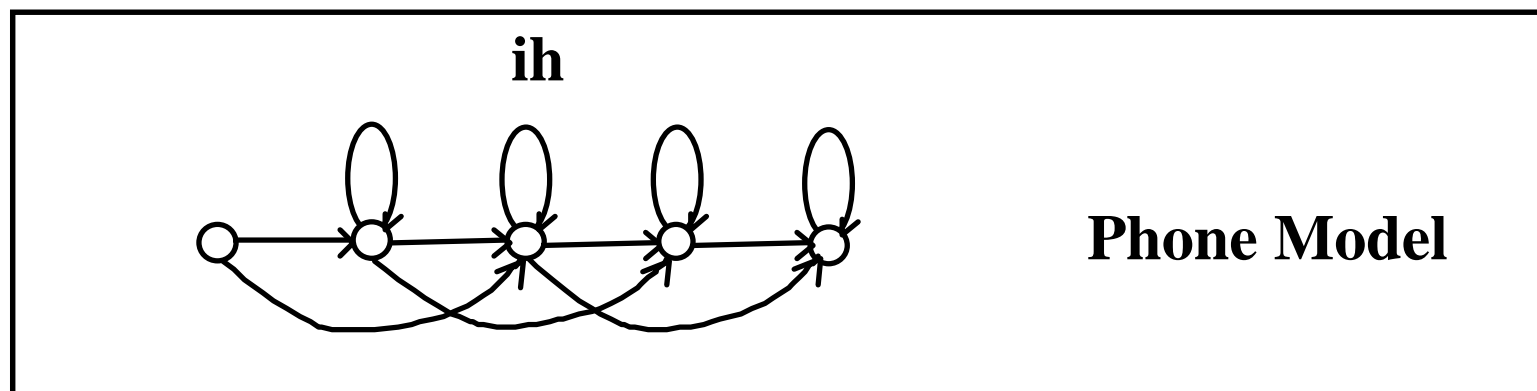**Observation sequence: R B Y Y ••• R**

**not unique to state sequence**

# HMMs In Speech Recognition

**Represent speech as a sequence of observations**

**Use HMM to model some unit of speech (phone, word)**

**Concatenate units into larger units**

ih

**Phone Model**

d          ih          d

**Word Model**

# HMM Problems And Solutions

**Evaluation:**
   - Problem - Compute Probabilty 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** $O = O_1 \, O_2 \cdots O_T$

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

$$P(O \mid \lambda) = \sum_{\forall Q} P(O, Q \mid \lambda) \qquad Q = q_0 q_1 \ldots q_T \text{ is a state sequence}$$

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

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

$\quad$ N = number of states in model
$\quad$ T = number of observations in sequence

# The Forward Algorithm

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

**Compute $\alpha$ recursively:**

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

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

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

# **Forward Trellis**

$$\begin{bmatrix} A & 0.8 \\ B & 0.2 \end{bmatrix}$$

0.6

1.0

0.4

**Initial**    **Final**

$$\begin{bmatrix} A & 0.3 \\ B & 0.7 \end{bmatrix}$$

|  | **t=0** | **A** **t=1** | **A** **t=2** | **B** **t=3** |
|---|---|---|---|---|
| **state 1** | 1.0 | 0.48 | 0.23 | 0.03 |
| **state 2** | 0.0 | 0.12 | 0.09 | 0.13 |

0.6 * 0.8        0.6 * 0.8        0.6 * 0.2

0.4 * 0.3        0.4 * 0.3        0.4 * 0.7

1.0 * 0.3        1.0 * 0.3        1.0 * 0.7

**14**

# The Backward Algorithm

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

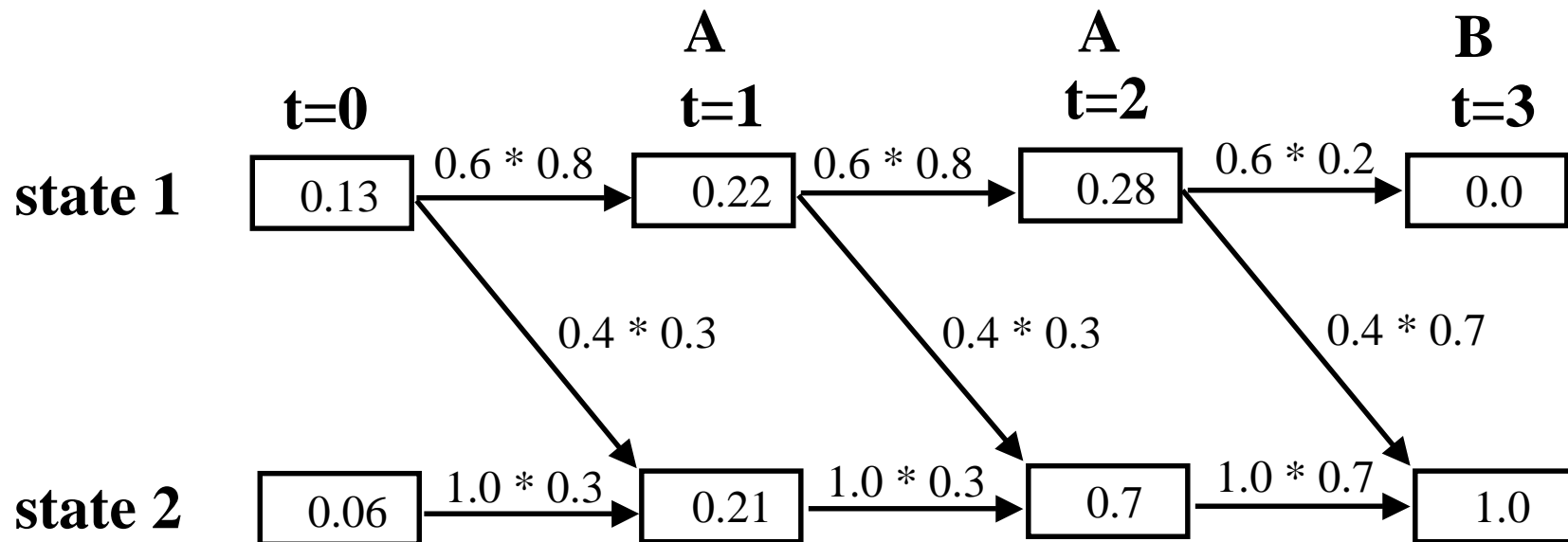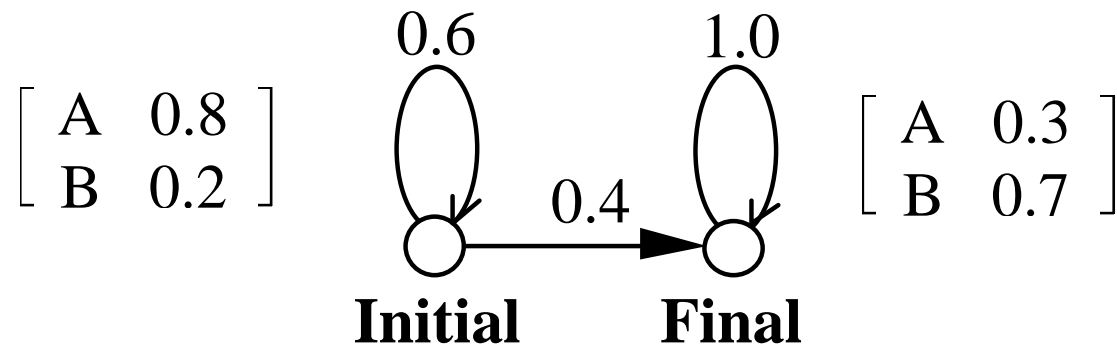**Compute $\beta$ recursively:**

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

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

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

# **Backward Trellis**

$$\begin{bmatrix} A & 0.8 \\ B & 0.2 \end{bmatrix}$$

0.6    1.0

0.4

**Initial**    **Final**

$$\begin{bmatrix} A & 0.3 \\ B & 0.7 \end{bmatrix}$$

|  | **t=0** | | **A** **t=1** | | **A** **t=2** | | **B** **t=3** |
|---|---|---|---|---|---|---|---|
| **state 1** | 0.13 | 0.6 * 0.8 | 0.22 | 0.6 * 0.8 | 0.28 | 0.6 * 0.2 | 0.0 |
| | | 0.4 * 0.3 | | 0.4 * 0.3 | | 0.4 * 0.7 | |
| **state 2** | 0.06 | 1.0 * 0.3 | 0.21 | 1.0 * 0.3 | 0.7 | 1.0 * 0.7 | 1.0 |

# **The Viterbi Algorithm**

For decoding:

Find the state sequence **Q** which maximizes **P(O, Q | λ )**

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

$$VP_t(i) = MAX_{q_0, \cdots q_{t-1}} \ P(O_1 O_2 \cdots O_t, q_t=i \mid \lambda \ )$$
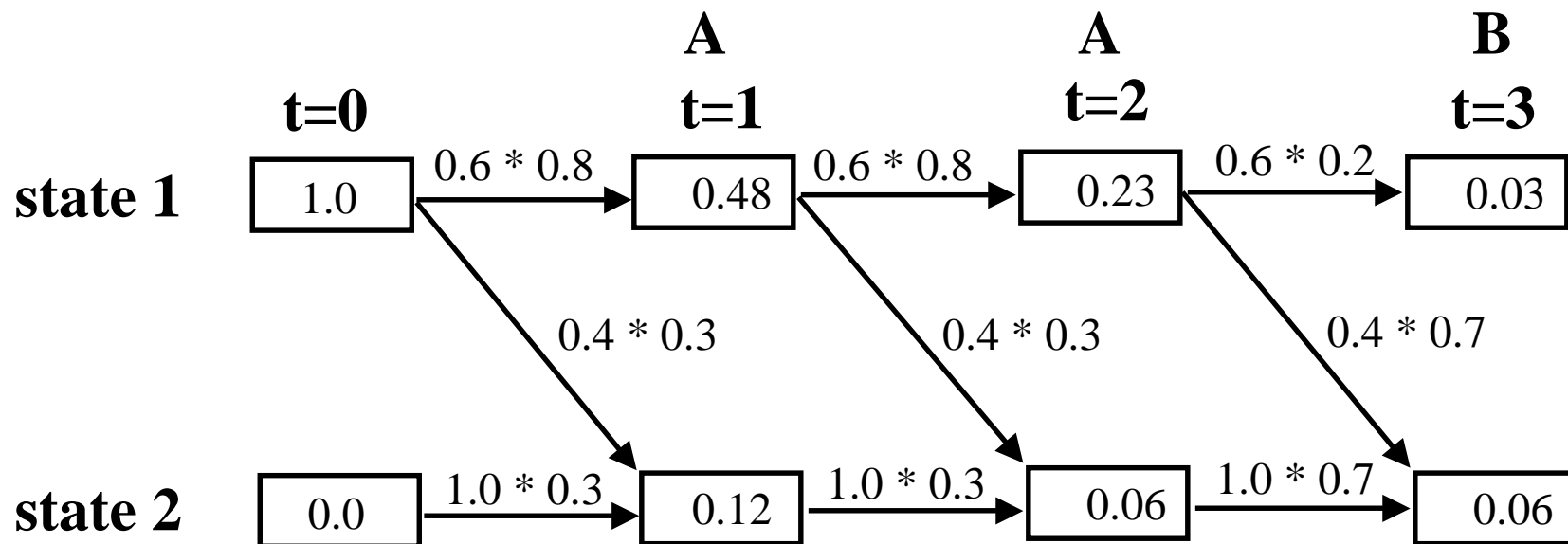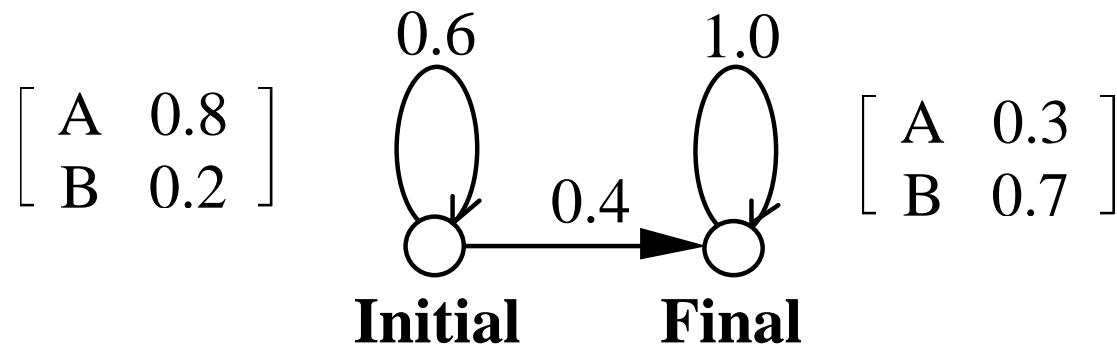
Recursive Computation:

$$VP_t(j) = MAX_{i=0, \ldots, N} \ VP_{t-1}(i) \ a_{ij} b_j(O_t) \qquad t > 0$$

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

Save each maximum for backtrace at end

# **Viterbi Trellis**

$$\begin{bmatrix} A & 0.8 \\ B & 0.2 \end{bmatrix}$$

0.6   1.0

0.4

**Initial**     **Final**

$$\begin{bmatrix} A & 0.3 \\ B & 0.7 \end{bmatrix}$$

|  |  | **A** | **A** | **B** |
|---|---|---|---|---|
|  | **t=0** | **t=1** | **t=2** | **t=3** |
| **state 1** | 1.0 | 0.48 | 0.23 | 0.03 |
| **state 2** | 0.0 | 0.12 | 0.06 | 0.06 |

0.6 * 0.8    0.6 * 0.8    0.6 * 0.2

0.4 * 0.3    0.4 * 0.3    0.4 * 0.7

1.0 * 0.3    1.0 * 0.3    1.0 * 0.7

# Training HMM Parameters

**Train parameters of HMM**

- Tune $\lambda$ to maximize $P(O \mid \lambda)$

- No efficient algorithm for global optimum

- Efficient iterative algorithm finds a local optimum
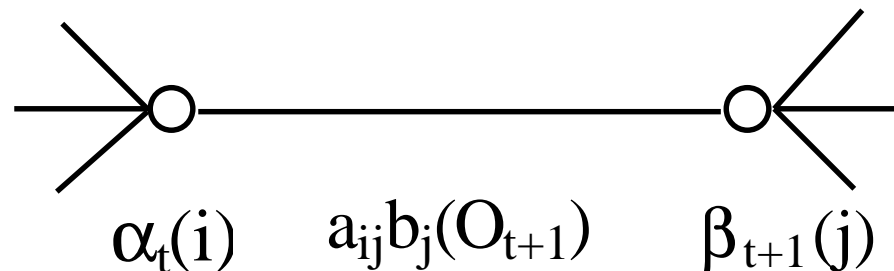
**Baum-Welch (Forward-Backward) re-estimation**

- Compute probabilities using current model $\lambda$

- Refine $\lambda \longrightarrow \lambda$ based on computed values

- Use $\alpha$ and $\beta$ from Forward-Backward

# Forward-Backward Algorithm

$$\xi_t(i,j) = \quad \begin{array}{l}\textbf{Probability of transiting from } S_i \textbf{ to } S_j \\ \textbf{at time } t \textbf{ given O}\end{array}$$

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

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

$$\alpha_t(i) \qquad a_{ij}b_j(O_{t+1}) \qquad \beta_{t+1}(j)$$

20

# Baum-Welch Reestimation

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

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

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

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

# 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 go to 2**

**It can be shown that $P(O \mid \bar{\lambda}) > P(O \mid \lambda)$ unless $\bar{\lambda} = \lambda$**
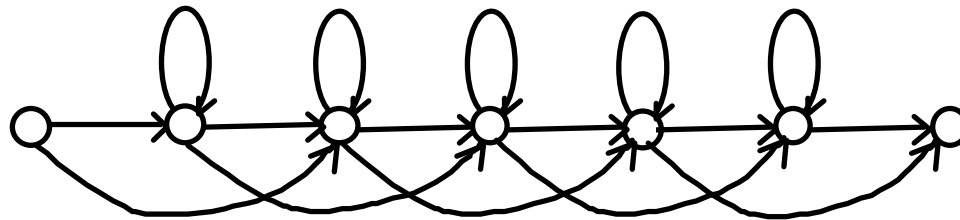
# HMMs In Speech Recognition

Represent speech as a sequence of symbols

Use HMM to model some unit of speech (phone, word)

Output Probabilities - Prob of observing symbol in a state

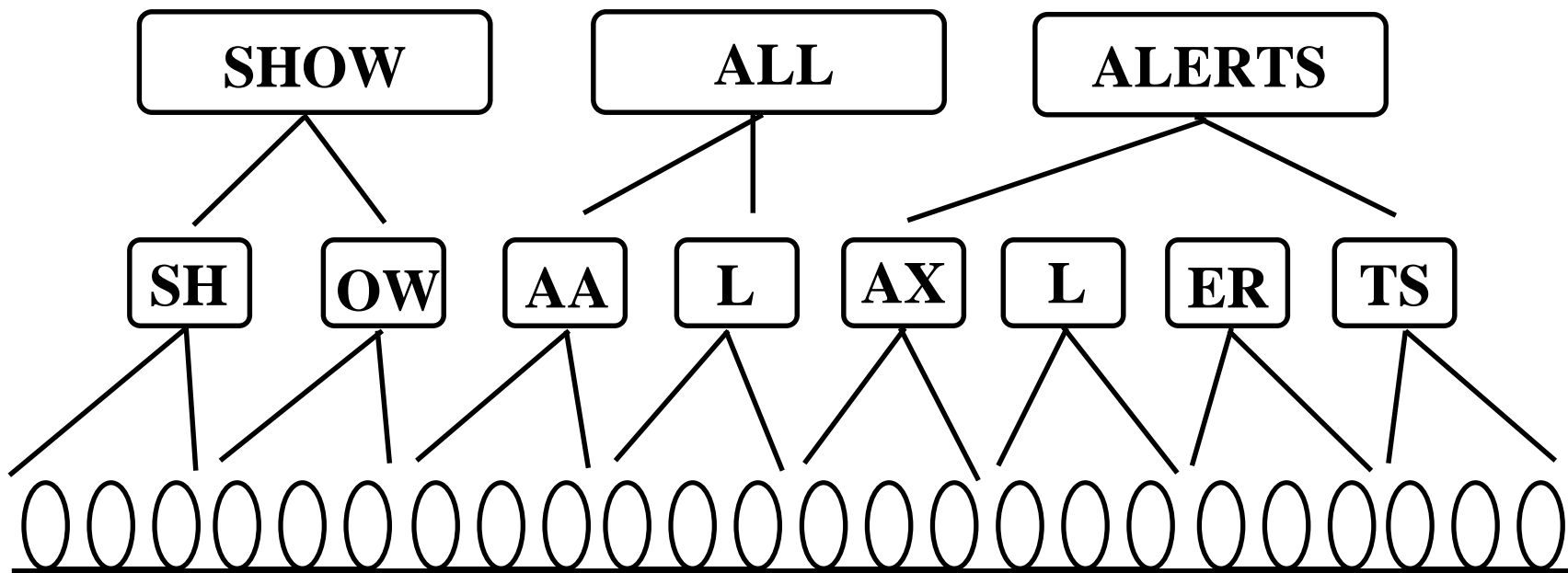Transition Prob - Prob of staying in or skipping state
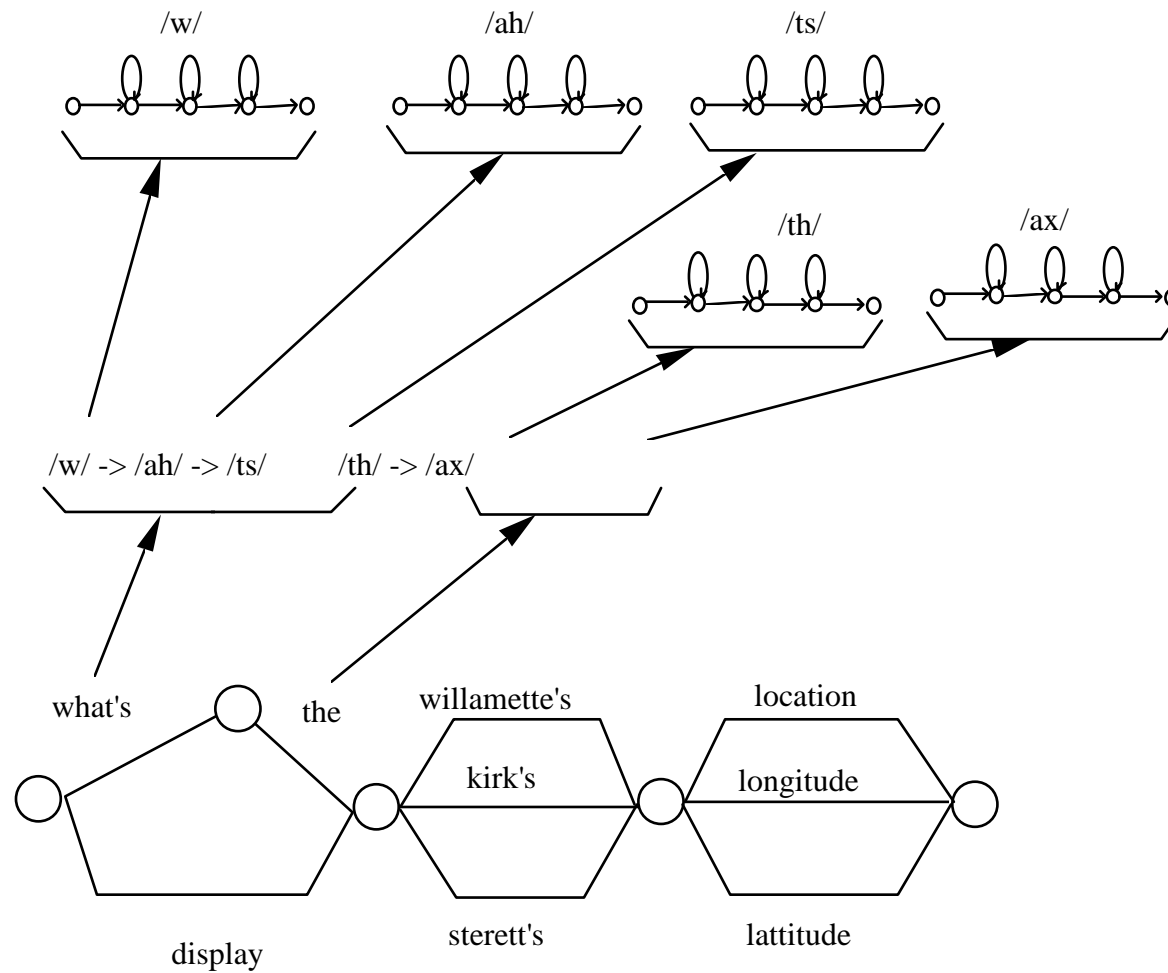
**Phone Model**

# Training HMMs for Continuous Speech

- Use only orthograph transcription of sentence

  - no need for segmented/labelled data

- Concatenate phone 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

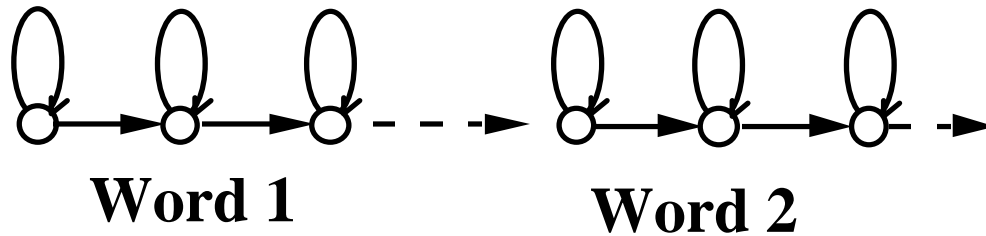# Recognition Search

/w/             /ah/             /ts/

/th/             /ax/

/w/ -> /ah/ -> /ts/       /th/ -> /ax/

what's          the      willamette's       location

kirk's         longitude

display         sterett's        lattitude

**26**

# Viterbi Search

- Uses Viterbi decoding

  - Takes MAX, not SUM

  - Finds optimal state sequence $P(O, Q \mid \lambda)$
    not optimal word sequence $P(O \mid \lambda)$

- Time synchronous

  - Extends all paths by 1 time step

  - All paths have same length (no need to normalize to compare scores)
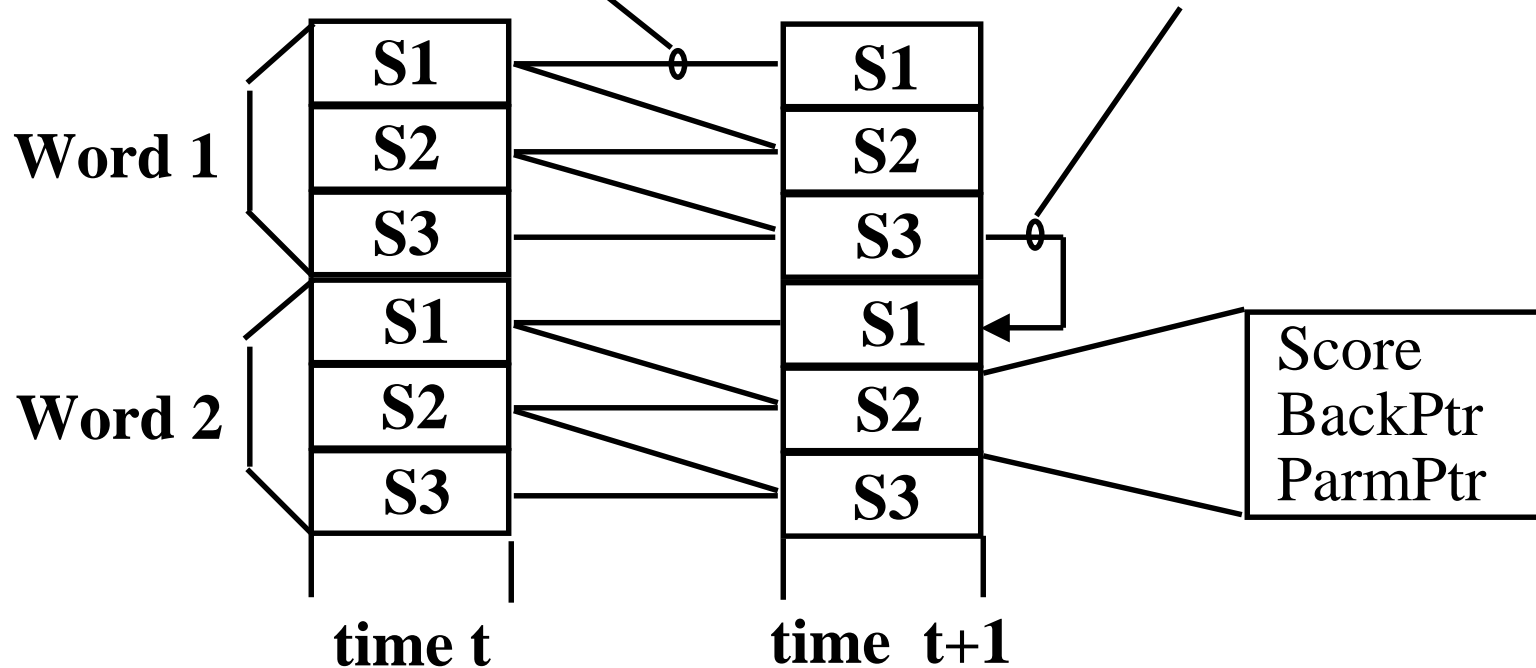
# Viterbi Search Algorithm

0. Create state list with one cell for each state in system

1. Initialize state list with initial states for time  t= 0

2. Clear state list for time t+1

3. Compute within-word transitions from time 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 time 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



Word 1          Word 2

OldProb(S1) • OutProb • Transprob          OldProb(S3) • P(W2 | W1)

|          |          |          |
|----------|----------|----------|----------|
| Word 1   | S1       |          | S1 |
|          | S2       |          | S2 |
|          | S3       |          | S3 |
| Word 2   | S1       |          | S1 |
|          | S2       |          | S2 |
|          | S3       |          | S3 |

Score
BackPtr
ParmPtr

time t          time  t+1

29

# Viterbi Beam Search

**Viterbi Search**

    All states enumerated

    Not practical for large grammars

    Most states inactive at any given time

**Viterbi Beam Search  -  prune less likely paths**

    States worse than threshold range from best are pruned

    From and To structures created dynamically - list of active
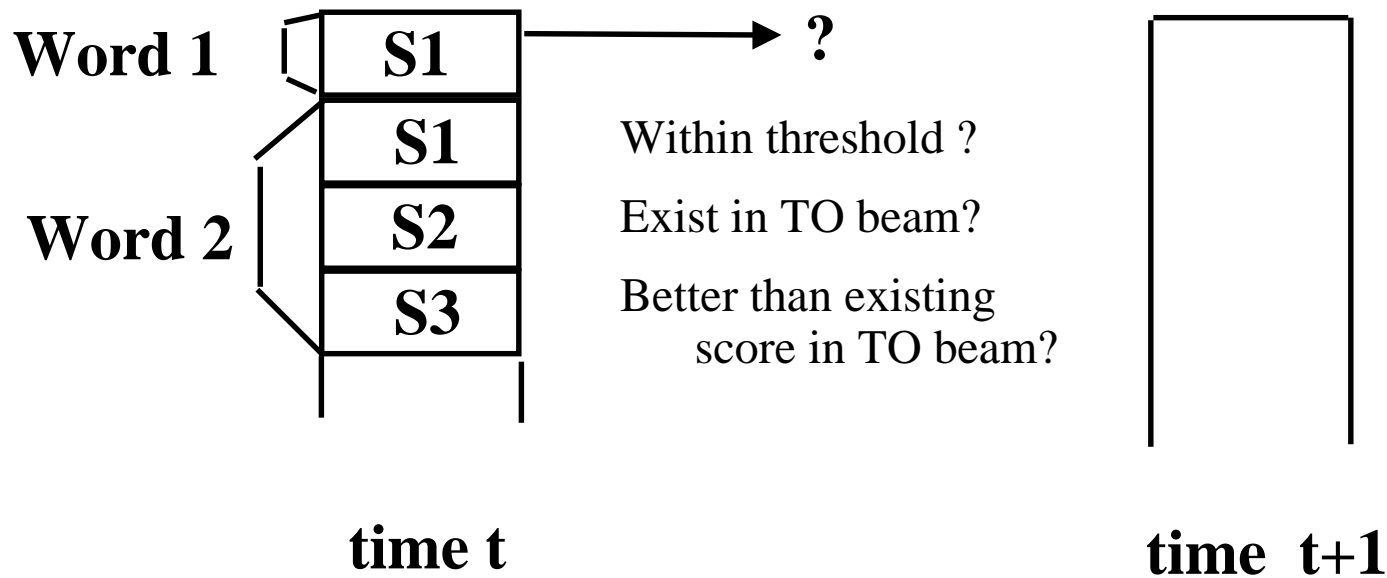       states

# Viterbi Beam Search

### FROM BEAM
States within threshold from best state

### TO BEAM
Dynamically constructed

**Word 1**   S1 ——————→ **?**

S1

**Word 2**   S2

S3

Within threshold ?

Exist in TO beam?

Better than existing score in TO beam?

**time t**                    **time  t+1**

# Continuous Density HMMs

Model so far has assumed discete 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

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

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

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

# Discrete Hmm vs. Continuous HMM

☐ **Problems with Discrete:**

- **quantization errors**

- **Codebook and HMMs modelled separately**

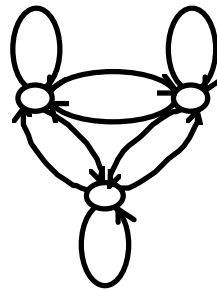- **Problems with Continuous Mixtures:**

  - **Small number of mixtures performs poorly**

  - **Large number of mixtures increases computation and parameters to be estimated**

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

- **Continuous makes more assumptions than Discrete, especially if diagonal covariance pdf**

- **Discrete probability is a table lookup, continuous mixtures require many multiplications**

# **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.