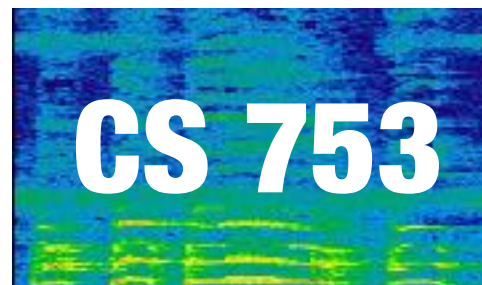# HMMs for Acoustic Modeling (Part I)
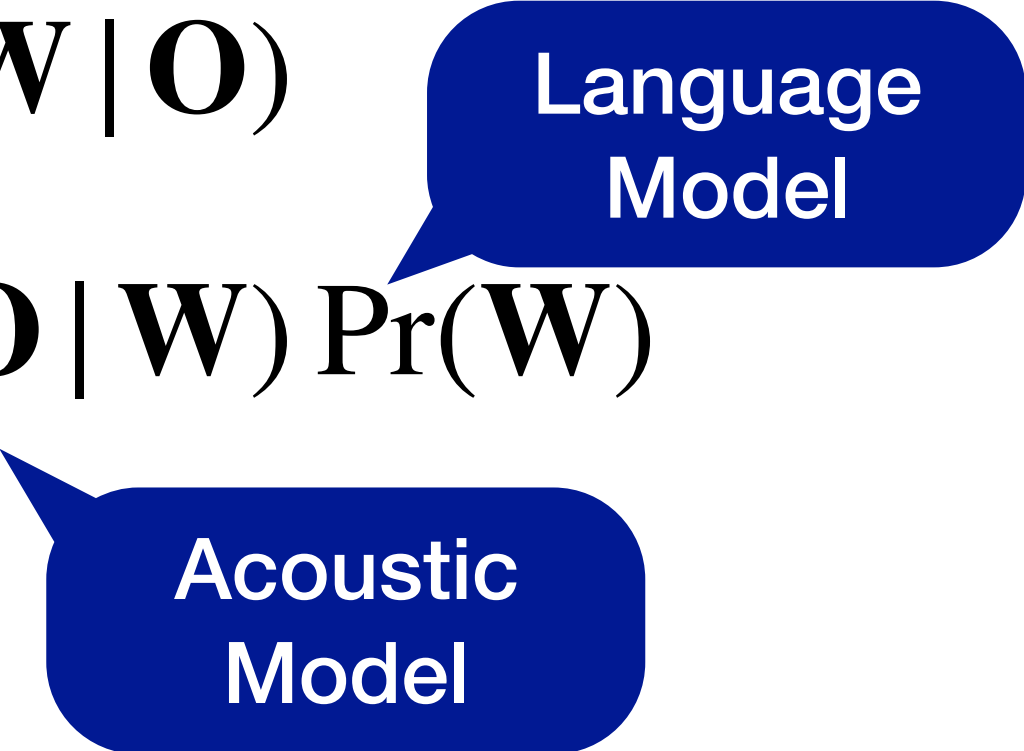
## Lecture 2

**CS 753**
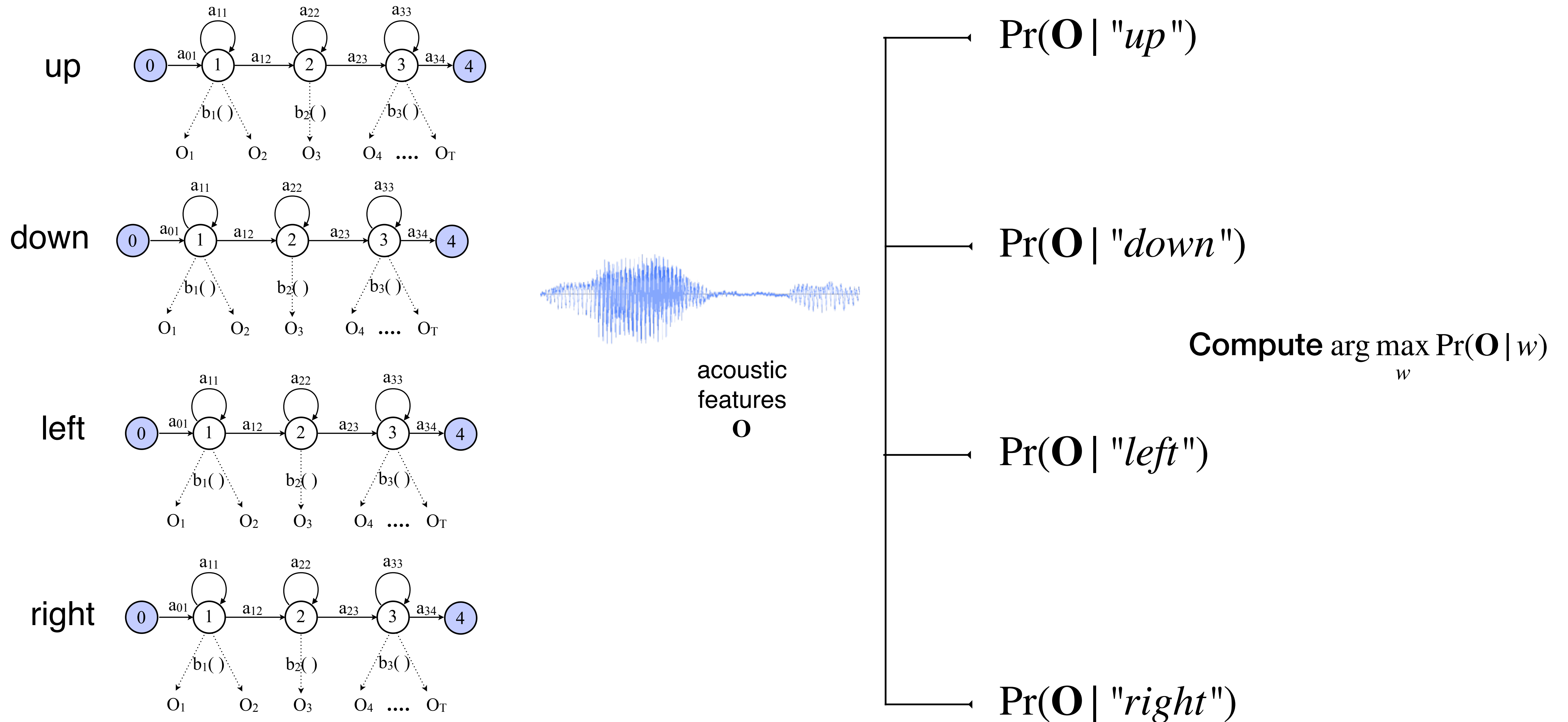
Instructor: Preethi Jyothi

# Recall: Statistical ASR

Let $\mathbf{O}$ be a sequence of acoustic features corresponding to a speech signal. That is, $\mathbf{O} = \{O_1, \ldots, O_T\}$, where $O_i \in \mathbb{R}^d$ refers to a d-dimensional acoustic feature vector and $T$ is the length of the sequence.

Let $\mathbf{W}$ denote a word sequence. An ASR decoder solves the foll. problem:

$$\mathbf{W}* = \arg\max_{W} \Pr(\mathbf{W} \,|\, \mathbf{O})$$

$$= \arg\max_{W} \Pr(\mathbf{O} \,|\, \mathbf{W}) \, \Pr(\mathbf{W})$$

Language Model

Acoustic Model

# Isolated word recognition



up

$a_{11}$ $a_{22}$ $a_{33}$

$0$ $\xrightarrow{a_{01}}$ $1$ $\xrightarrow{a_{12}}$ $2$ $\xrightarrow{a_{23}}$ $3$ $\xrightarrow{a_{34}}$ $4$

$b_1()$ $b_2()$ $b_3()$

$O_1$ $O_2$ $O_3$ $O_4$ .... $O_T$

down

$a_{11}$ $a_{22}$ $a_{33}$

$0$ $\xrightarrow{a_{01}}$ $1$ $\xrightarrow{a_{12}}$ $2$ $\xrightarrow{a_{23}}$ $3$ $\xrightarrow{a_{34}}$ $4$

$b_1()$ $b_2()$ $b_3()$

$O_1$ $O_2$ $O_3$ $O_4$ .... $O_T$

left

$a_{11}$ $a_{22}$ $a_{33}$

$0$ $\xrightarrow{a_{01}}$ $1$ $\xrightarrow{a_{12}}$ $2$ $\xrightarrow{a_{23}}$ $3$ $\xrightarrow{a_{34}}$ $4$

$b_1()$ $b_2()$ $b_3()$

$O_1$ $O_2$ $O_3$ $O_4$ .... $O_T$

right

$a_{11}$ $a_{22}$ $a_{33}$

$0$ $\xrightarrow{a_{01}}$ $1$ $\xrightarrow{a_{12}}$ $2$ $\xrightarrow{a_{23}}$ $3$ $\xrightarrow{a_{34}}$ $4$

$b_1()$ $b_2()$ $b_3()$

$O_1$ $O_2$ $O_3$ $O_4$ .... $O_T$

acoustic features $\mathbf{O}$

$\Pr(\mathbf{O} \mid \textit{"up"})$

$\Pr(\mathbf{O} \mid \textit{"down"})$

**Compute** $\arg\max_{w} \Pr(\mathbf{O} \mid w)$

$\Pr(\mathbf{O} \mid \textit{"left"})$

$\Pr(\mathbf{O} \mid \textit{"right"})$
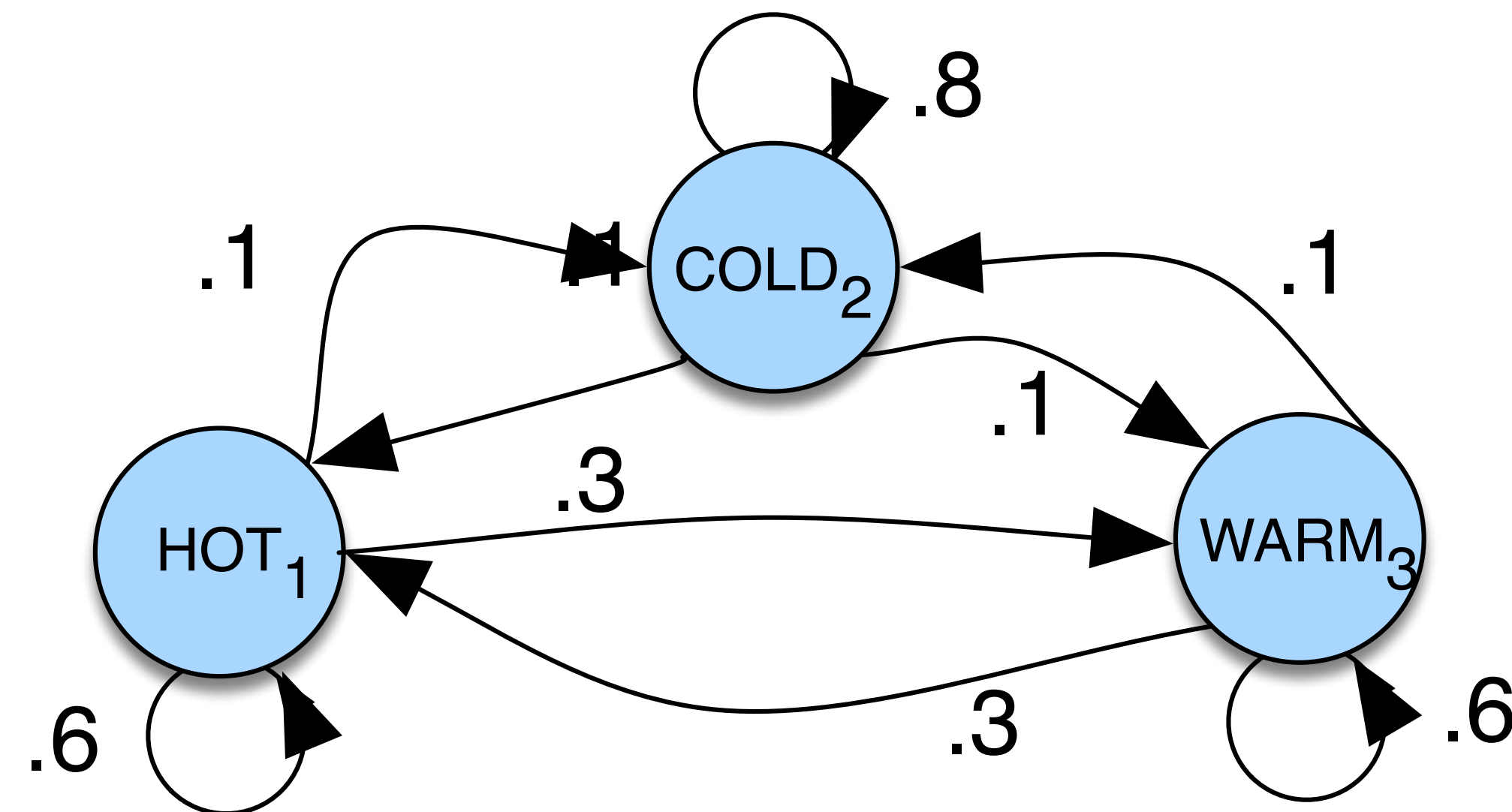
# What are Hidden Markov Models (HMMs)?

Following slides contain figures/material from "Hidden Markov Models", "Speech and Language Processing", D. Jurafsky and J. H. Martin, 2019. (https://web.stanford.edu/~jurafsky/slp3/A.pdf)

.6 .3 .1 .6

$$\pi = [0.1, 0.7, 0.2]$$

ing the probabili

$$\sum_{j=1}^{n} a_{ij} = 1 \quad \forall i$$

$\pi = \pi_1, \pi_2, ..., \pi_N$ an **initial probability distrit**
probability that the Markov
Some states $j$ may have $\pi_j =$
be initial states. Also, $\sum_{i=1}^{n}$

# HMM Assumptions



**Markov Assumption:** $P(q_i|q_1...q_{i-1}) = P(q_i|q_{i-1})$

**Output I**

# Hidden Markov Model

$Q = q_1 q_2 \ldots q_N$     a set of $N$ **states**

$A = a_{11} \ldots a_{ij} \ldots a_{NN}$     a **transition probability matrix** $A$, each $a_{ij}$ representing the probability of moving from state $i$ to state $j$, s.t. $\sum_{j=1}^{N} a_{ij} = 1 \quad \forall i$

$O = o_1 o_2 \ldots o_T$     a sequence of $T$ **observations**, each one drawn from a vocabulary $V = v_1, v_2, ..., v_V$

$B = b_i(o_t)$     a sequence of **observation likelihoods**, also called **emission probabilities**, each expressing the probability of an observation $o_t$ being generated from a state $i$

$\pi = \pi_1, \pi_2, ..., \pi_N$     an **initial probability distribution** over states. $\pi_i$ is the probability that the Markov chain will start in state $i$. Some states $j$ may have $\pi_j = 0$, meaning that they cannot be initial states. Also, $\sum_{i=1}^{n} \pi_i = 1$

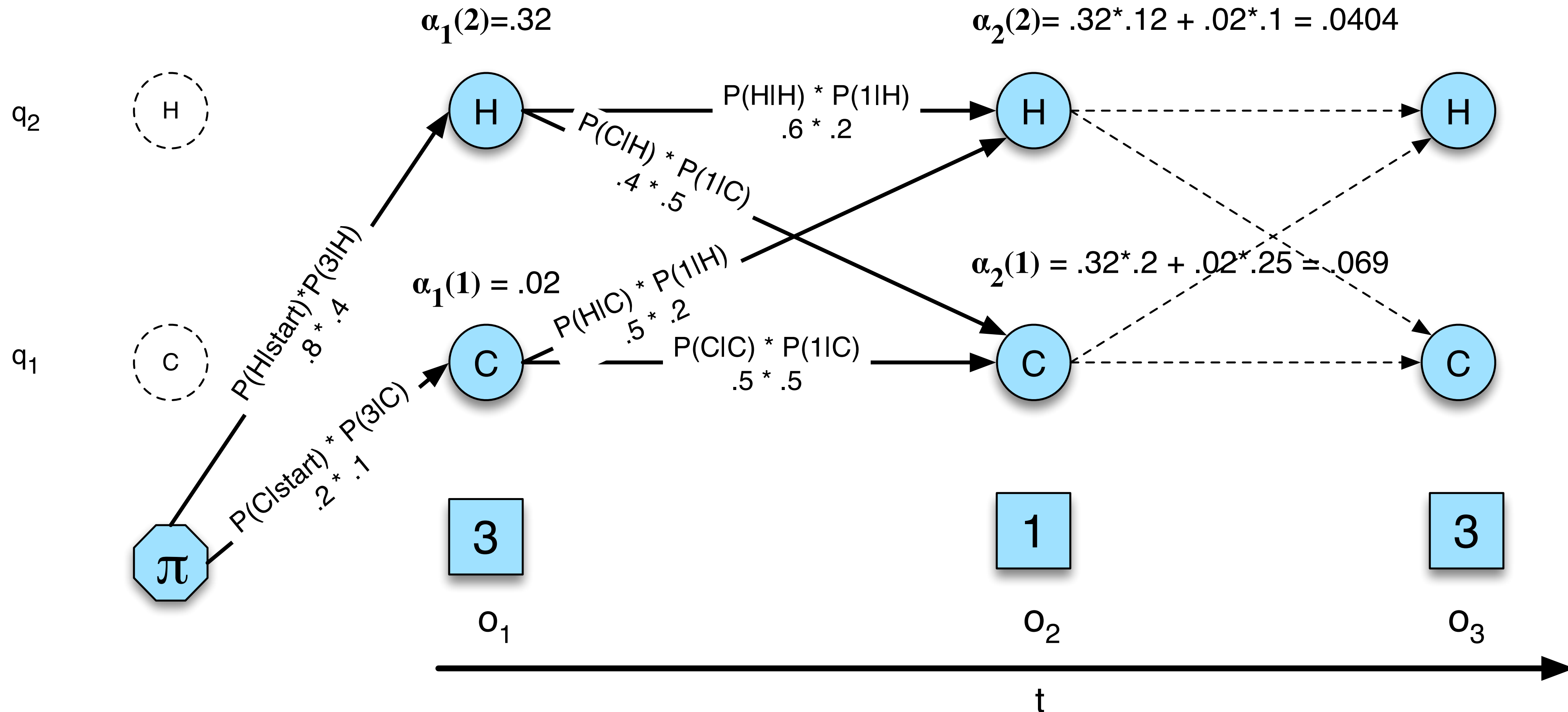| | |
|---|---|
| **Problem 1 (Likelihood):** | Given an HMM $\lambda = (A, B)$ and an observation sequence $O$, determine the likelihood $P(O\|\lambda)$. |
| **Problem 2 (Decoding):** | Given an observation sequence $O$ and an HMM $\lambda = (A, B)$, discover the best hidden state sequence $Q$. |
| **Problem 3 (Learning):** | Given an observation sequence $O$ and the set of states in the HMM, learn the HMM parameters $A$ and $B$. |

**Computing Likelihood:** Given an HMM $\lambda = (A, B)$ and an observation sequence $O$, determine the likelihood $P(O\|\lambda)$.
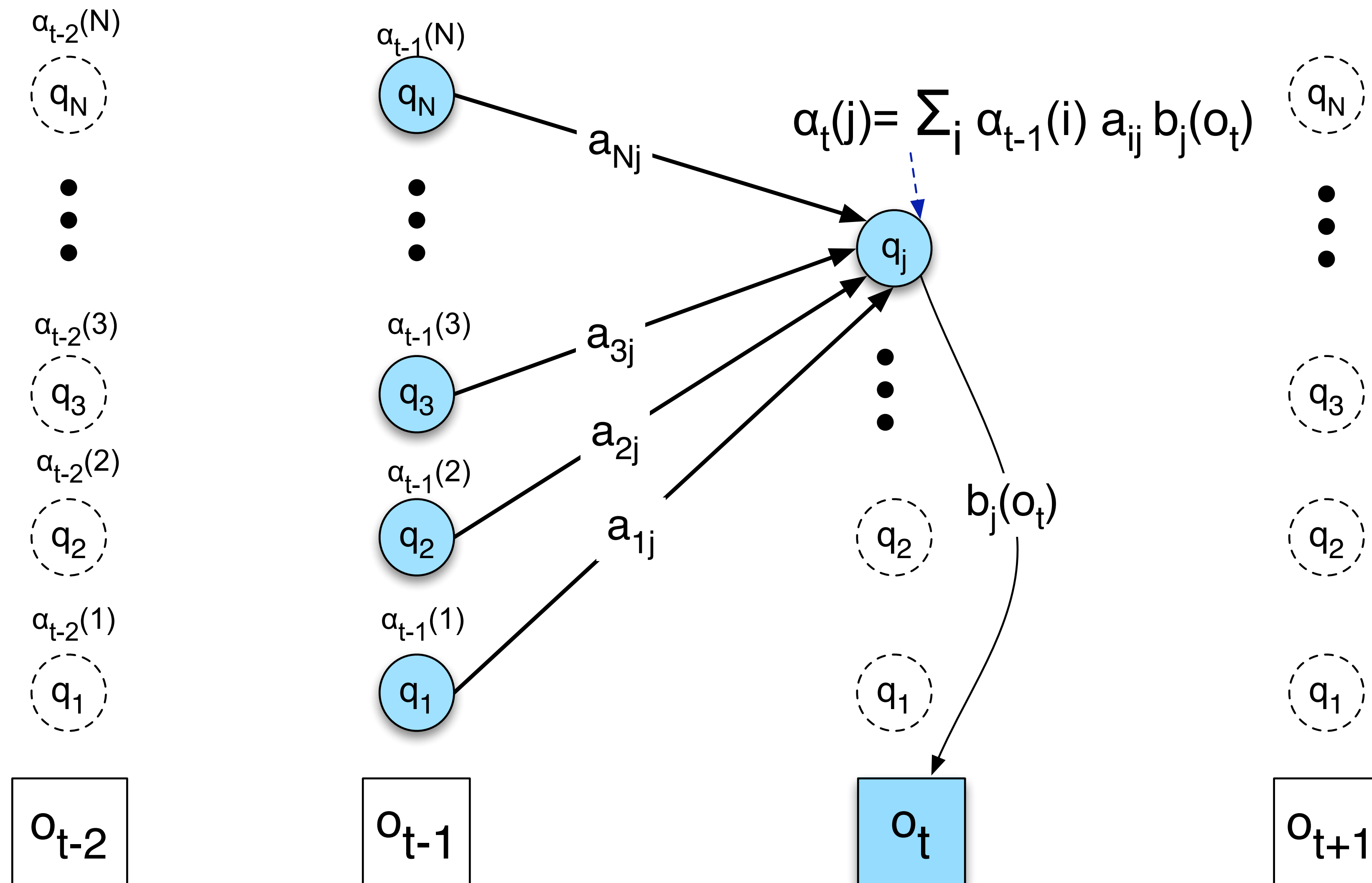
# Forward Algorithm

$$\alpha_t(j) = P(o_1, o_2 \ldots o_t, q_t = j | \lambda)$$

$$\alpha_t(j) = \sum_{i=1}^{N} \alpha_{t-1}(i) a_{ij} b_j(o_t)$$

# Visualizing the forward recursion



$$\alpha_t(j) = \sum_i \alpha_{t-1}(i) \, a_{ij} \, b_j(o_t)$$

# Forward Algorithm

1. Initialization:

$$\alpha_1(j) \;=\; \pi_j b_j(o_1) \;\; 1 \le j \le N$$

2. Recursion:

$$\alpha_t(j) = \sum_{i=1}^{N} \alpha_{t-1}(i) a_{ij} b_j(o_t); \;\; 1 \le j \le N, 1 < t \le T$$

3. Termination:

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

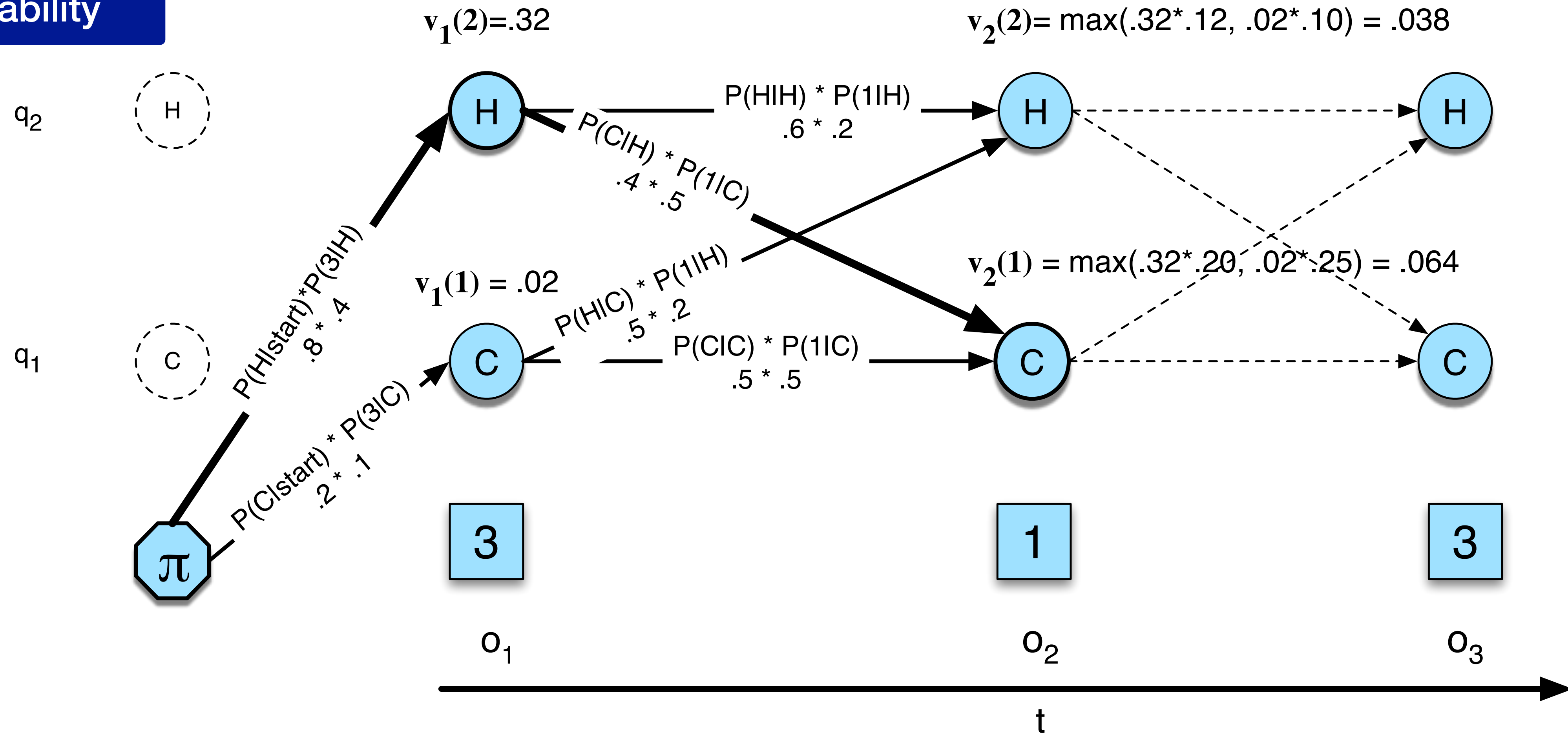| | |
|---|---|
| **Problem 1 (Likelihood):** | Given an HMM $\lambda = (A, B)$ and an observation sequence $O$, determine the likelihood $P(O|\lambda)$. |
| **Problem 2 (Decoding):** | Given an observation sequence $O$ and an HMM $\lambda = (A, B)$, discover the best hidden state sequence $Q$. |
| **Problem 3 (Learning):** | Given an observation sequence $O$ and the set of states in the HMM, learn the HMM parameters $A$ and $B$. |

**Decoding**: Given as input an HMM $\lambda = (A, B)$ and a sequence of observations $O = o_1, o_2, \ldots, o_T$, find the most probable sequence of states $Q = q_1 q_2 q_3 \ldots q_T$.

# Viterbi Trellis

$$v_t(j) = \max_{q_1,\dots,q_{t-1}} P(q_1\dots q_{t-1}, o_1, o_2 \dots o_t, q_t = j | \lambda)$$

$$v_t(j) = \max_{i=1}^{N} v_{t-1}(i)\, a_{ij}\, b_j(o_t)$$

**Viterbi Path Probability**

**$v_1(2)$=.32**

**$v_2(2)$= max(.32*.12, .02*.10) = .038**

$q_2$   ( H )

P(H|H) * P(1|H)
.6 * .2

H     H     H

P(C|H) * P(1|C)
.4 * .5

P(H|start)*P(3|H)
.8 * .4

**$v_1(1)$ = .02**

P(H|C) * P(1|H)
.5 * .2

**$v_2(1)$ = max(.32*.20, .02*.25) = .064**

$q_1$   ( C )

C

P(C|C) * P(1|C)
.5 * .5

C     C

P(C|start) * P(3|C)
.2 * .1

π

**3**        **1**        **3**

$o_1$        $o_2$        $o_3$

t

1. **Initialization:**

$$v_1(j) = \pi_j b_j(o_1) \qquad 1 \leq j \leq N$$
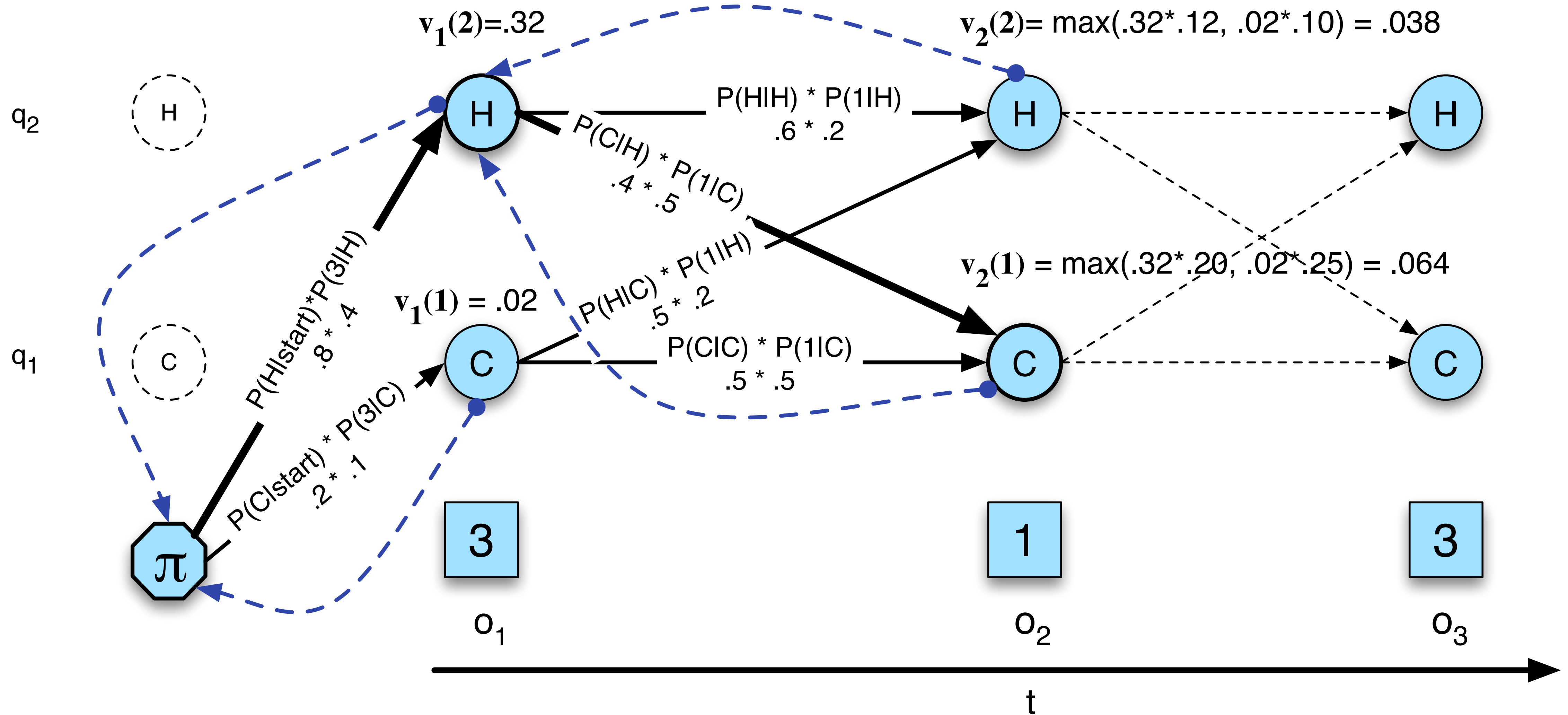$$bt_1(j) = 0 \qquad\qquad 1 \leq j \leq N$$

2. **Recursion**

$$v_t(j) = \max_{i=1}^{N} v_{t-1}(i)\, a_{ij}\, b_j(o_t); \quad 1 \leq j \leq N, 1 < t \leq T$$
$$bt_t(j) = \operatorname*{argmax}_{i=1}^{N} v_{t-1}(i)\, a_{ij}\, b_j(o_t); \quad 1 \leq j \leq N, 1 < t \leq T$$

3. **Termination:**

$$\text{The best score:} \quad P* = \max_{i=1}^{N} v_T(i)$$

$$\text{The start of backtrace:} \quad q_T* = \operatorname*{argmax}_{i=1}^{N} v_T(i)$$

# Viterbi backtrace

# Gaussian Observation Model

- So far, we considered HMMs with discrete outputs

- In acoustic models, HMMs output real valued vectors

- Hence, observation probabilities are defined using probability density functions

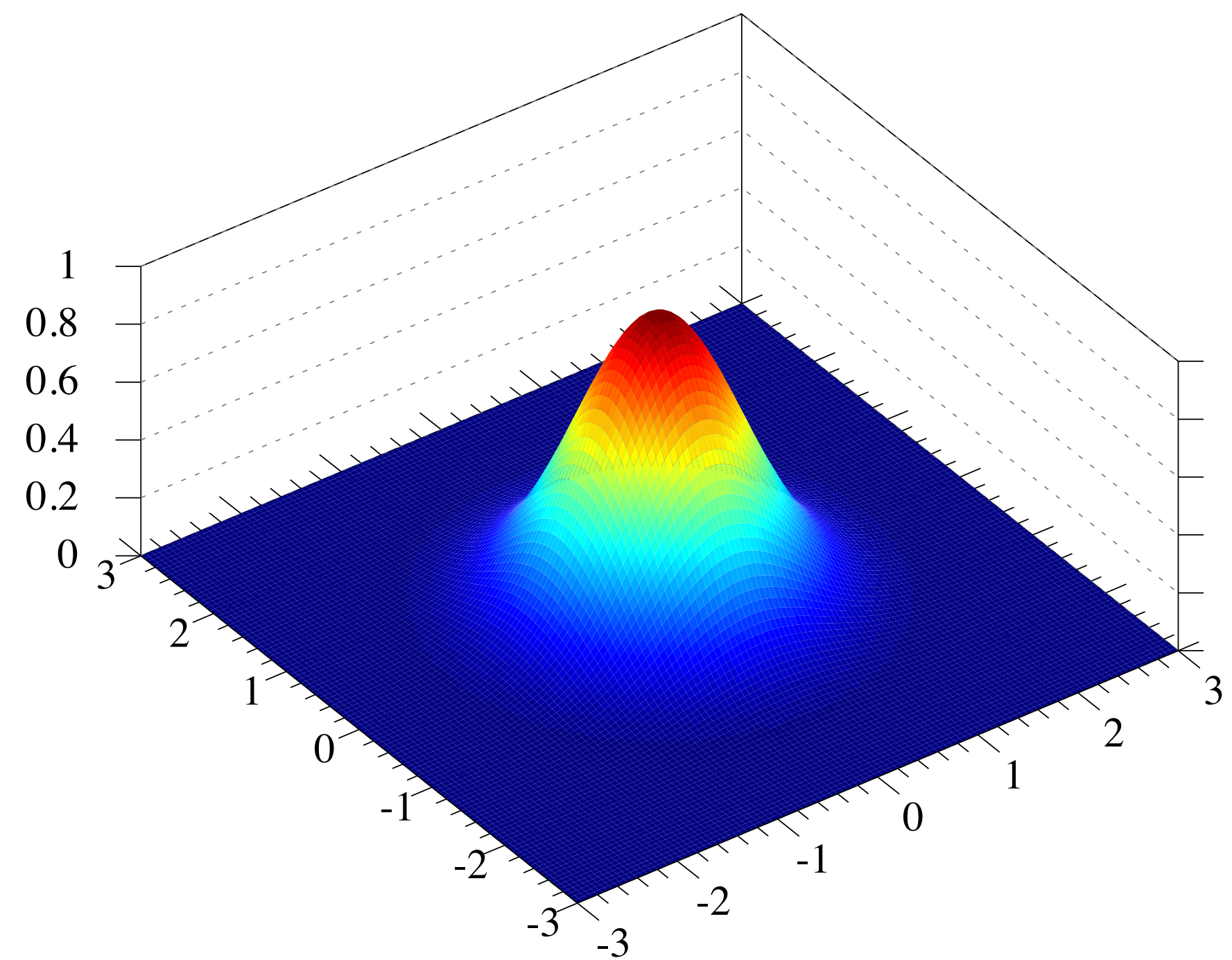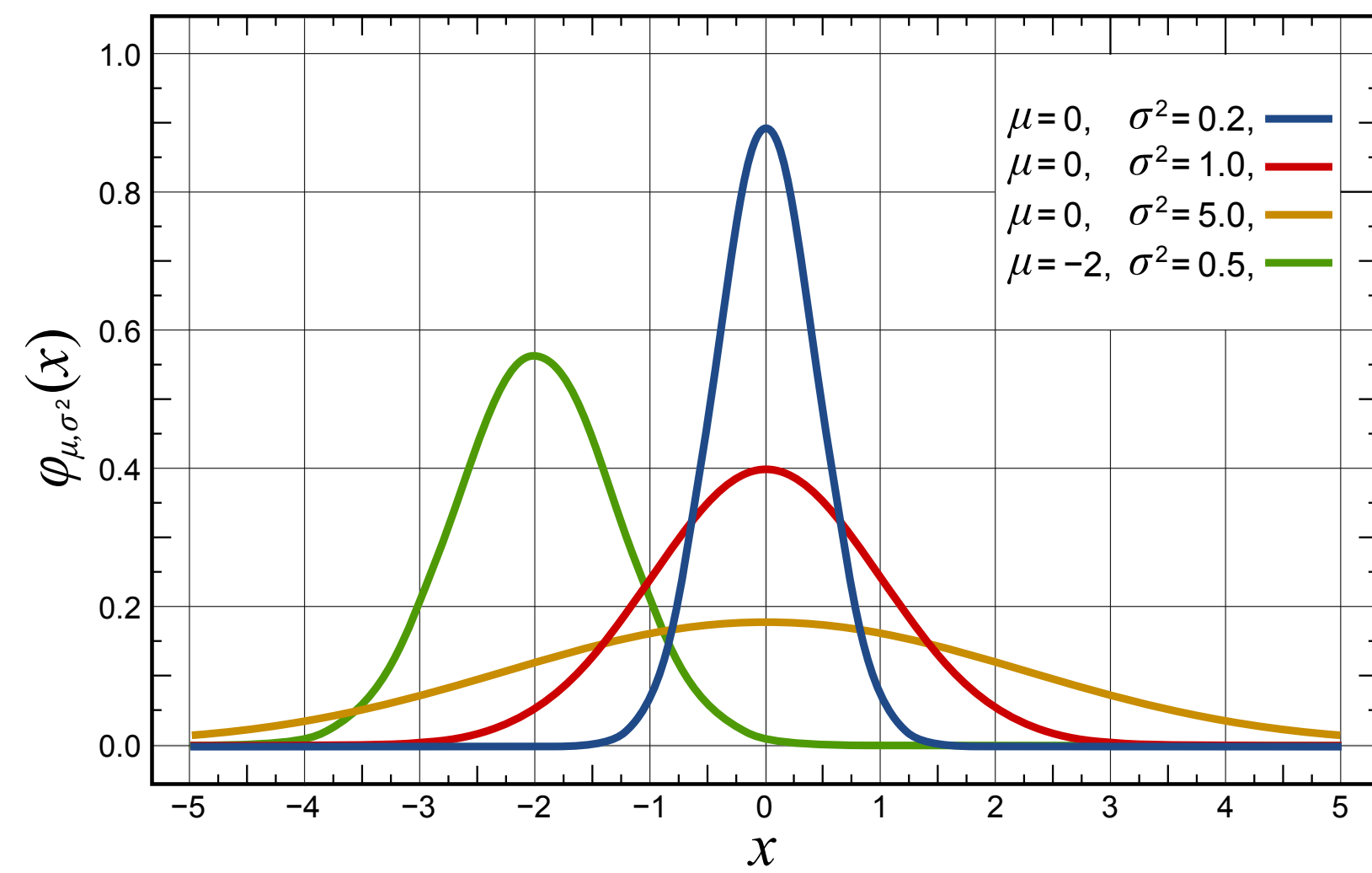- A widely used model: Gaussian distribution

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$$

- HMM emission/observation probabilities $b_j(x) = \mathcal{N}(x|\mu_j, \sigma_j^2)$ where $\mu_j$ is the mean associated with state $j$ and $\sigma_j^2$ is its variance

- For multivariate Gaussians, $b_j(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \Sigma_j)$ where $\Sigma_j$ is the covariance matrix associated with state j

# Gaussian Mixture Model

- A single Gaussian observation model assumes that the observed acoustic feature vectors are unimodal
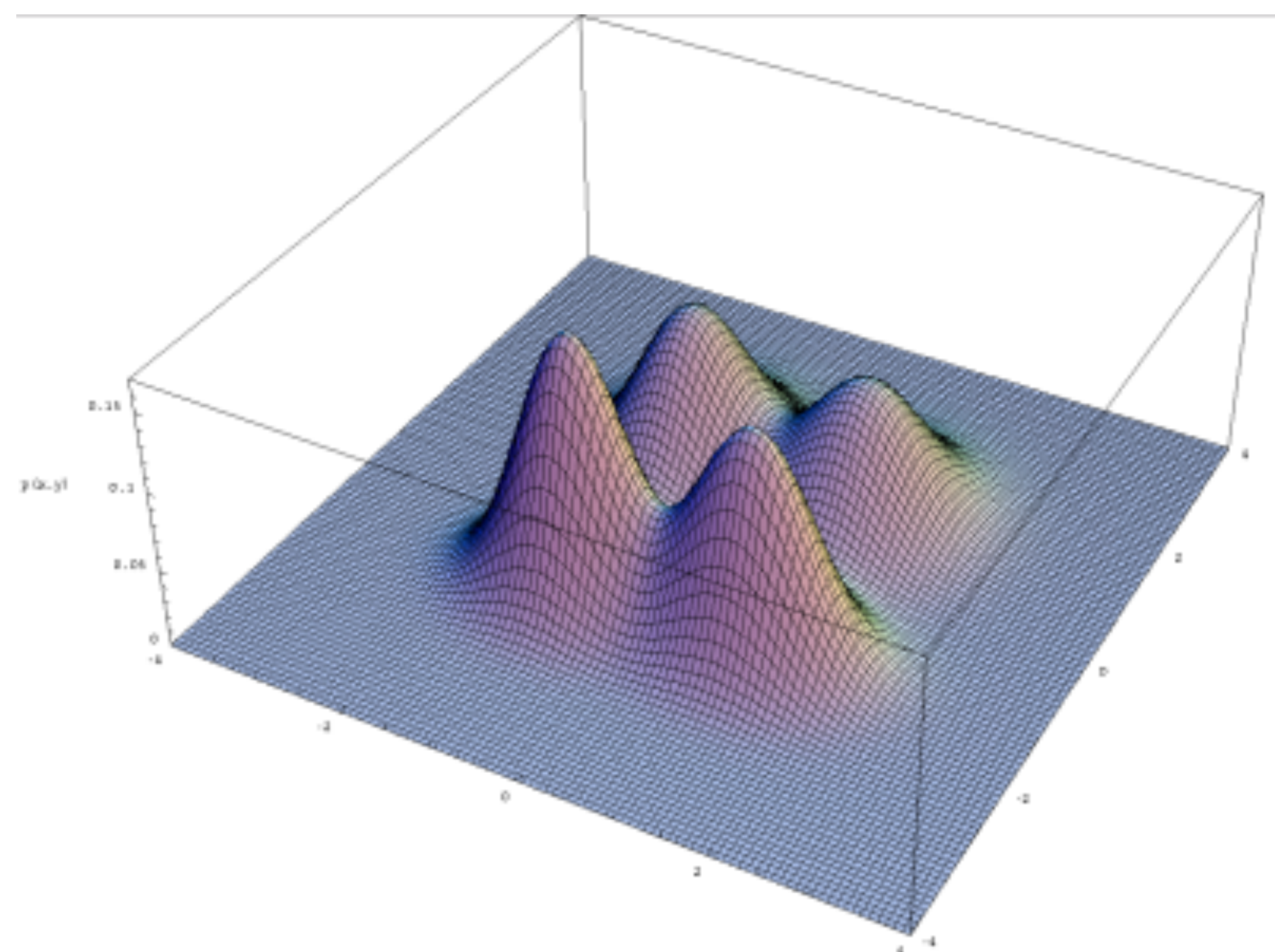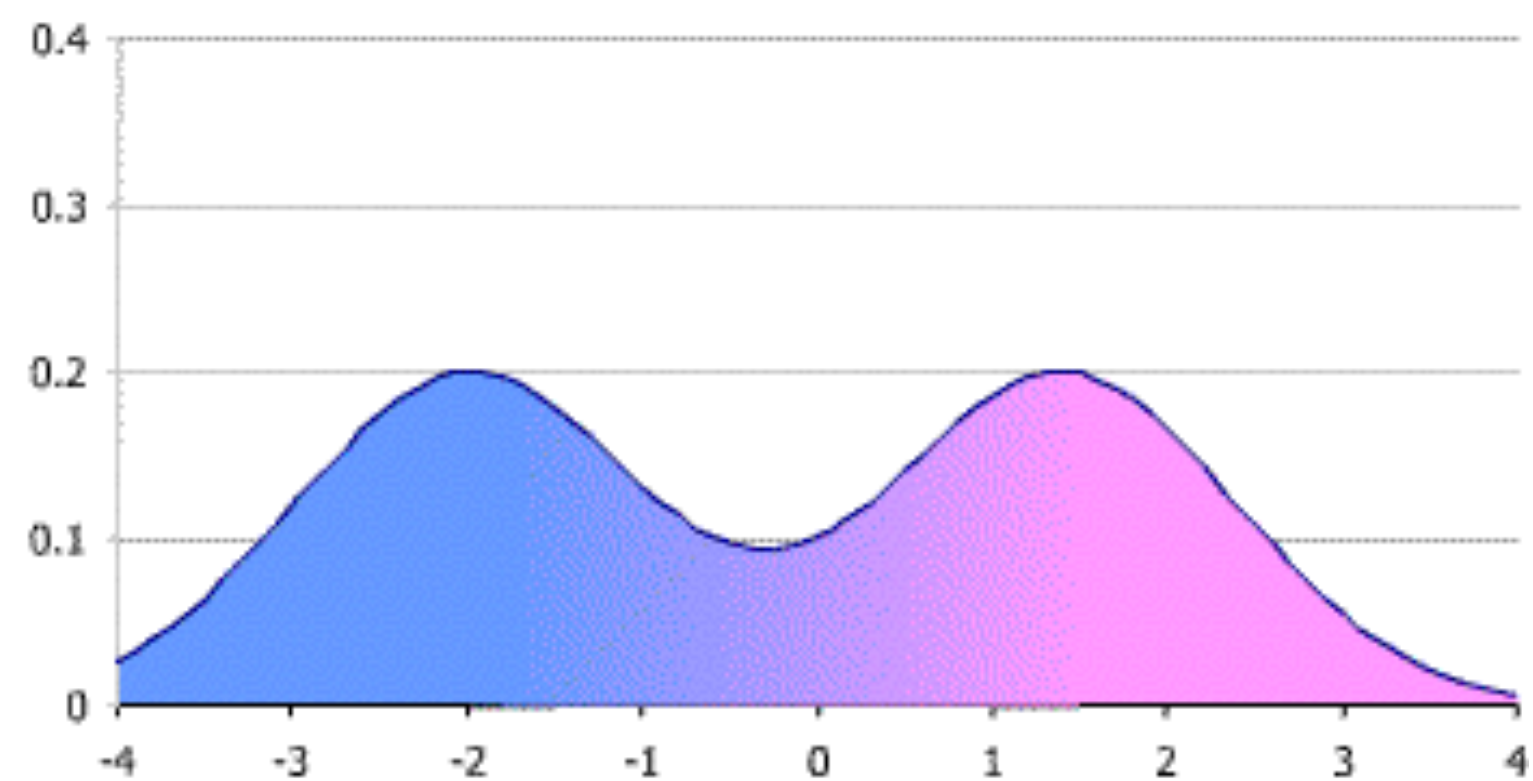
# Unimodal

# Gaussian Mixture Model

- A single Gaussian observation model assumes that the observed acoustic feature vectors are unimodal

- More generally, we use a "mixture of Gaussians" to model multiple modes in the data

# Mixture Models

# Gaussian Mixture Model

- A single Gaussian observation model assumes that the observed acoustic feature vectors are unimodal

- More generally, we use a "mixture of Gaussians" to model multiple modes in the data

- Instead of $b_j(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \Sigma_j)$ in the single Gaussian case, $b_j(\mathbf{x})$ now becomes:

$$b_j(\mathbf{x}) = \sum_{m=1}^{M} c_{jm} \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_{jm}, \Sigma_{jm})$$

where $c_{jm}$ is the mixing probability for Gaussian component $m$ of state $j$

$$\sum_{m=1}^{M} c_{jm} = 1, \quad c_{jm} \geq 0$$