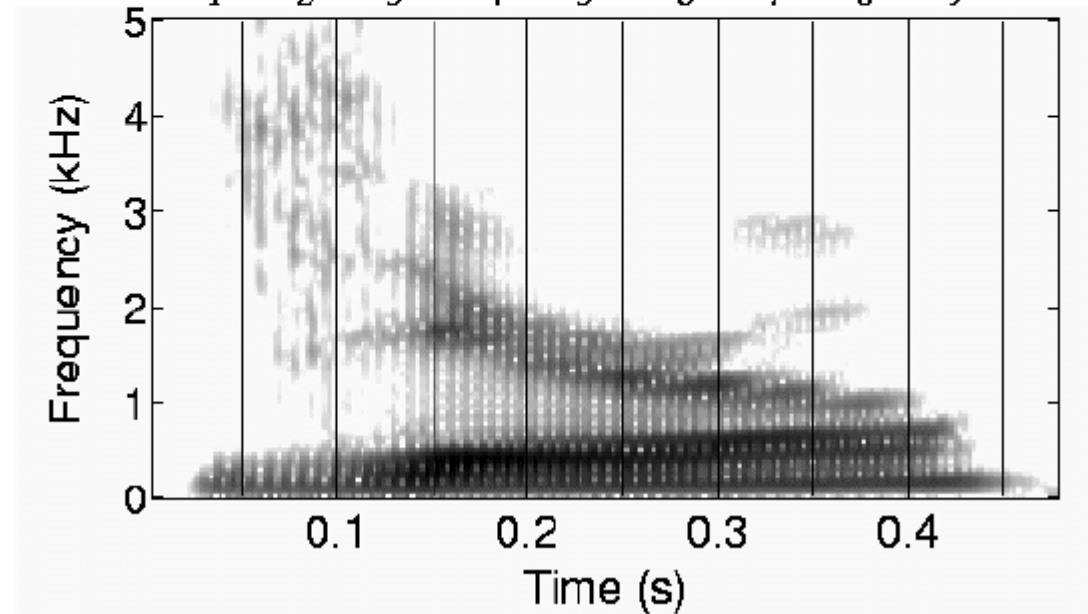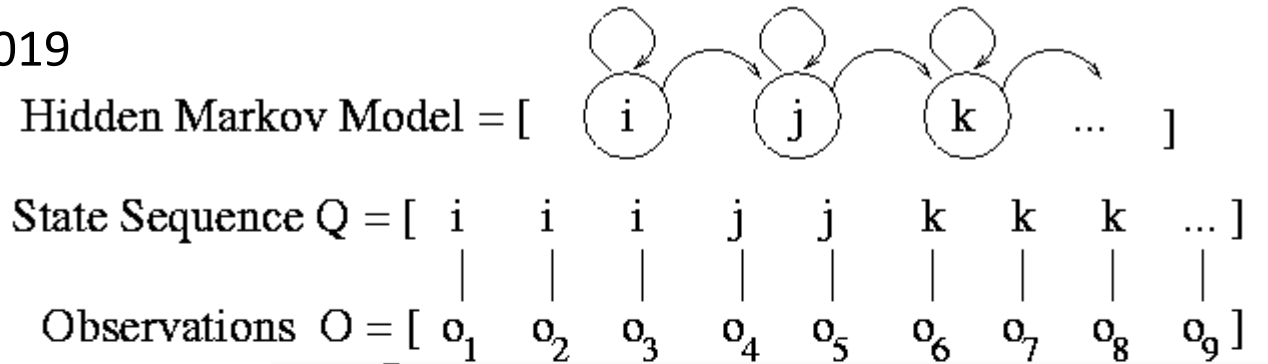# CS440/ECE448 Lecture 20: Hidden Markov Models

Slides by Svetlana Lazebnik, 11/2016

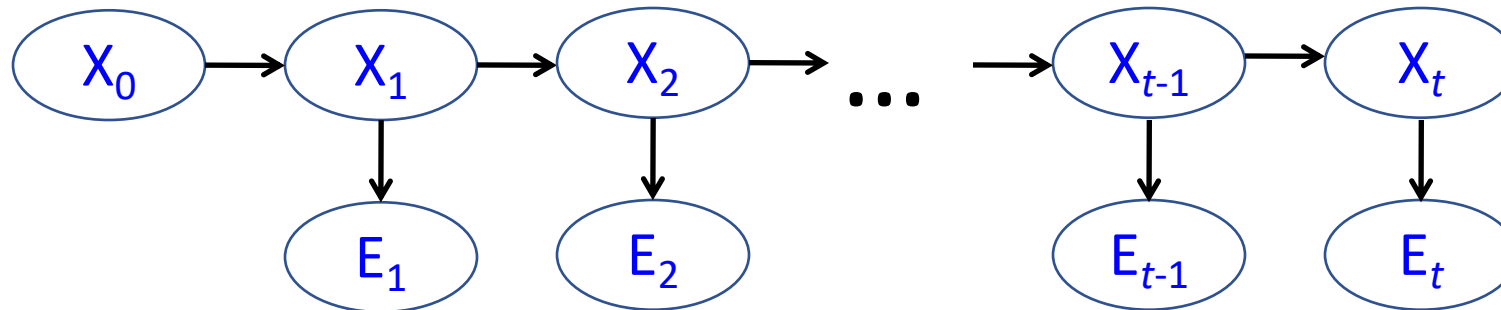Modified by Mark Hasegawa-Johnson, 3/2019

# Probabilistic reasoning over time

- So far, we've mostly dealt with *episodic* environments
  - Exceptions: games with multiple moves, planning
- In particular, the Bayesian networks we've seen so far describe static situations
  - Each random variable gets a single fixed value in a single problem instance
- Now we consider the problem of describing probabilistic environments that evolve over time
  - Examples: robot localization, human activity detection, tracking, speech recognition, machine translation,

# Hidden Markov Models

- At each time slice $t$, the state of the world is described by an **unobservable (hidden) variable** $X_t$ and an **observable *evidence* variable** $E_t$

- **Transition model:** distribution over the current state given the whole past history:
  $P(X_t \mid X_0, ..., X_{t-1}) = P(X_t \mid \mathbf{X}_{0:t-1})$

- **Observation model:** $P(E_t \mid \mathbf{X}_{0:t}, \mathbf{E}_{1:t-1})$

# Hidden Markov Models
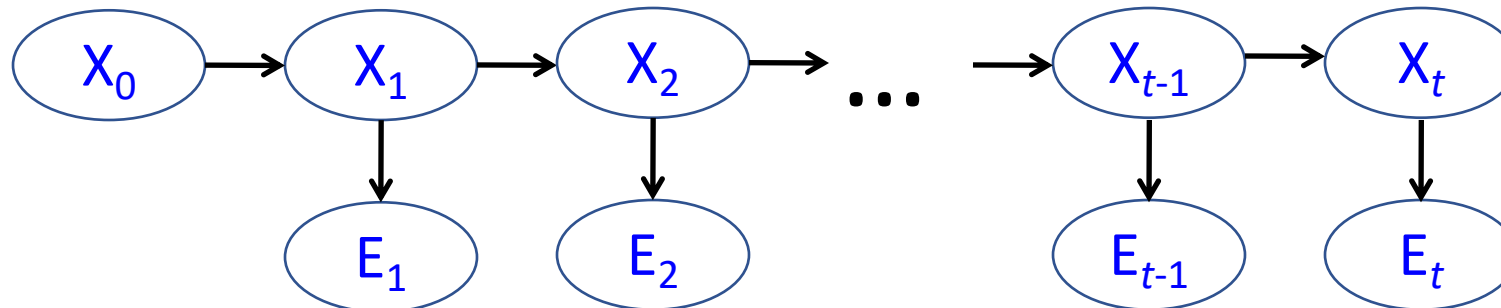
- **Markov assumption** (first order)
  - The current state is conditionally independent of all the other states given the state in the previous time step
  - What does $P(X_t \mid \mathbf{X}_{0:t-1})$ simplify to?

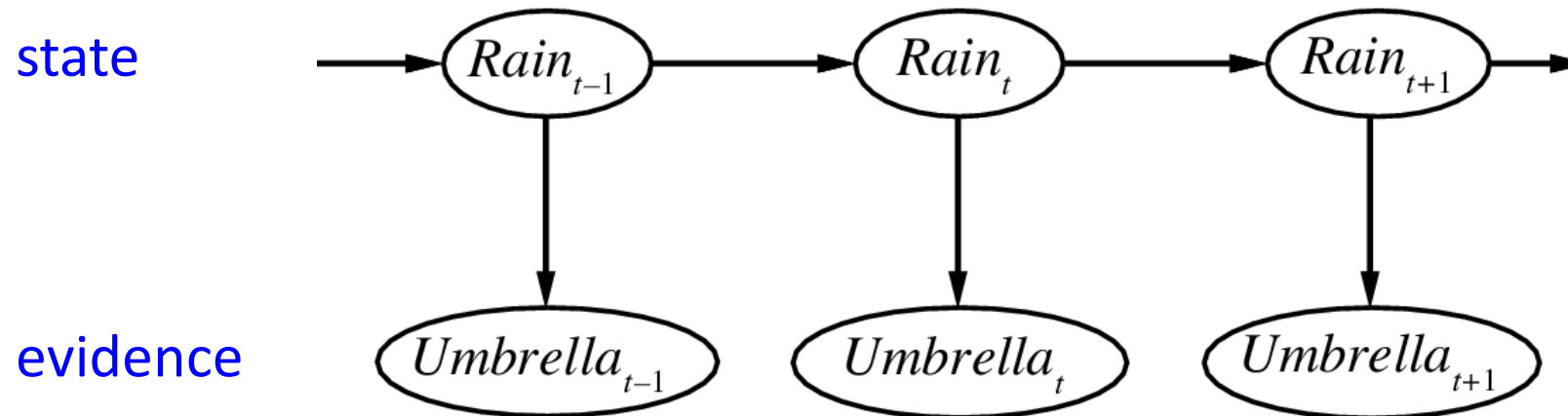  $P(X_t \mid \mathbf{X}_{0:t-1}) = P(X_t \mid X_{t-1})$

- Markov assumption for observations
  - The evidence at time $t$ depends only on the state at time $t$
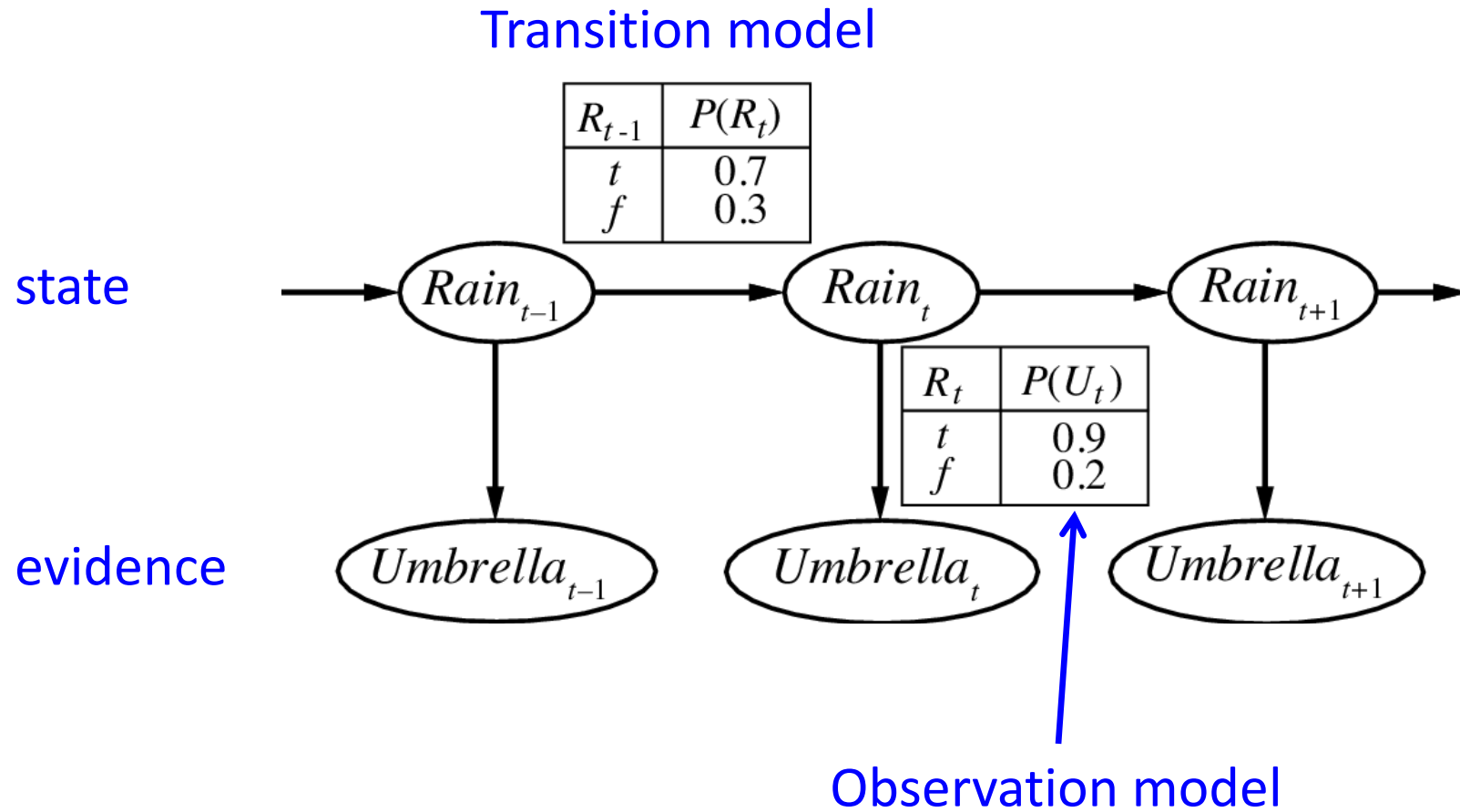  - What does $P(E_t \mid \mathbf{X}_{0:t}, \mathbf{E}_{1:t-1})$ simplify to?

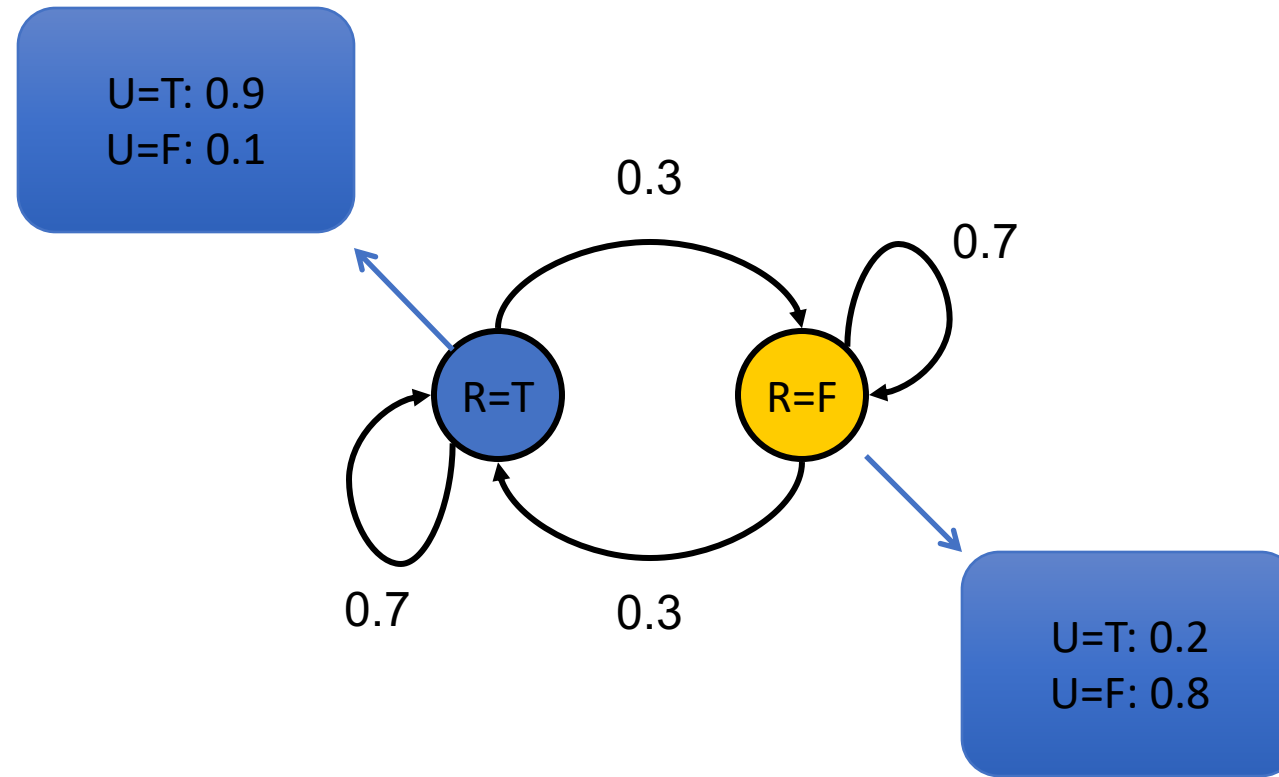  $P(E_t \mid \mathbf{X}_{0:t}, \mathbf{E}_{1:t-1}) = P(E_t \mid X_t)$

# Example

state



evidence

# Example

# An alternative visualization



U=T: 0.9
U=F: 0.1

0.3

0.7

R=T    R=F

0.7    0.3

U=T: 0.2
U=F: 0.8

Transition probabilities

|            | $R_t = T$ | $R_t = F$ |
|------------|-----------|-----------|
| $R_{t-1} = T$ | 0.7    | 0.3       |
| $R_{t-1} = F$ | 0.3    | 0.7       |

Observation (emission) probabilities

|          | $U_t = T$ | $U_t = F$ |
|----------|-----------|-----------|
| $R_t = T$ | 0.9      | 0.1       |
| $R_t = F$ | 0.2      | 0.8       |

# Another example

- **States:** X = {home, office, cafe}

- **Observations:** E = {sms, facebook, email}



30%   SMS
50%   FACEBOOK
20%   EMAIL

HOME

OFFICE        CAFE

10%   SMS
10%   FACEBOOK
80%   EMAIL

80%   SMS
05%   FACEBOOK
10%   EMAIL

## Transition Probabilities

|        | home | office | cafe |
|--------|------|--------|------|
| home   | 0.2  | 0.6    | 0.2  |
| office | 0.5  | 0.2    | 0.3  |
| cafe   | 0.2  | 0.8    | 0.0  |

## Emission Probabilities

|        | sms | facebook | email |
|--------|-----|----------|-------|
| home   | 0.3 | 0.5      | 0.2   |
| office | 0.1 | 0.1      | 0.8   |
| cafe   | 0.8 | 0.1      | 0.1   |

# The Joint Distribution

- Transition model: $P(X_t \mid \mathbf{X}_{0:t-1}) = P(X_t \mid X_{t-1})$

- Observation model: $P(E_t \mid \mathbf{X}_{0:t}, \mathbf{E}_{1:t-1}) = P(E_t \mid X_t)$

- How do we compute the full joint $P(\mathbf{X}_{0:t}, \mathbf{E}_{1:t})$?

$$P(\mathbf{X}_{0:t}, \mathbf{E}_{1:t}) = P(X_0) \prod_{i=1}^{t} P(X_i \mid X_{i-1}) P(E_i \mid X_i)$$

# Review: Bayes net inference

- Inference:
  - Trees: Sum-Product Algorithm (Textbook: "Variable Elimination" Algorithm)
  - Other Nets: Junction Tree Algorithm (Textbook: "Join Tree" Algorithm)
  - In General: NP-Complete, because clique size = graph size in general
- Parameter learning
  - Fully observed: Count # times each event occurs
  - Partially observed: Expectation-Maximization algorithm
    - Estimate Probability of each event at each time
    - E[# times event occurs] = sum_t(Probability event occurs at time t)

# Sum-Product Algorithm for HMMs (Forward algorithm)

- An HMM is a tree!
- Let's say we want to find $P(X_3|E_1,E_2,E_3) = P(X_3,E_1,E_2,E_3)/P(E_1,E_2,E_3)$

$P(X_3,E_1,E_2,E_3) = \sum_{X0} \sum_{X1} \sum_{X2} P(X_0,X_1,X_2 X_3,E_1,E_2,E_3)$

$\qquad = \sum_{X0} \sum_{X1} \sum_{X2} P(X_0)P(X_1|X_0)P(E_1|X_1)P(X_2|X_1)P(E_2|X_2) P(X_3|X_2)P(E_3|X_3)$

Let's rearrange the sums:

$\qquad = [\sum_{X2} [\sum_{X1} [\sum_{X0} P(X_0)P(X_1|X_0)] P(E_1|X_1)P(X_2|X_1)] P(E_2|X_2) P(X_3|X_2) ] P(E_3|X_3)$

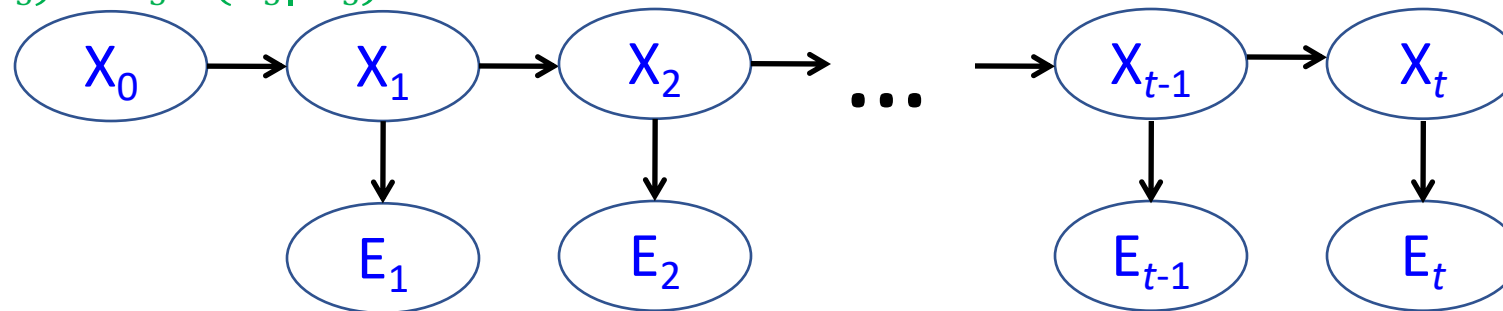Let's compute $F_1 = [\sum_{X0} P(X_0)P(X_1|X_0)]$ for any value of $X_1$ (so we marginalize out $X_0$)

$\qquad = [\sum_{X2} [\sum_{X1} F_1 P(E_1|X_1)P(X_2|X_1)] P(E_2|X_2) P(X_3|X_2) ] P(E_3|X_3)$

Now let's compute $F_2 = [\sum_{X1} F_1 P(E_1|X_1)P(X_2|X_1)]$ for any value of $X_2$ (so we marginalize out $X_1$)

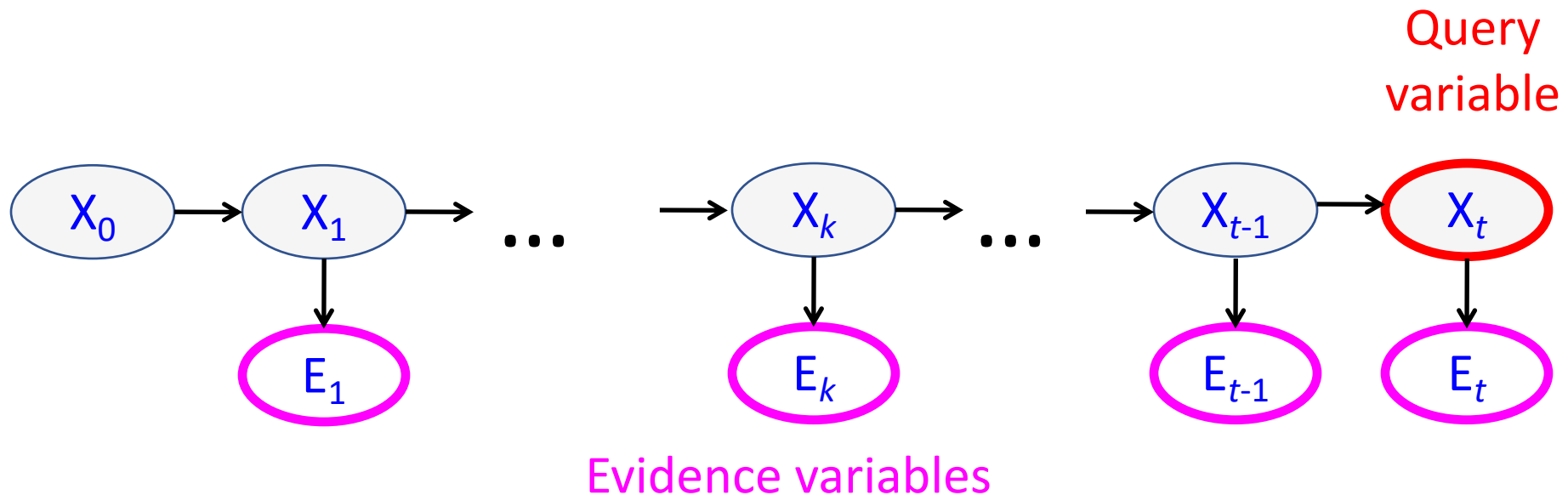$\qquad = [\sum_{X2} F_2 P(E_2|X_2) P(X_3|X_2) ] P(E_3|X_3)$

And $F_3 = [\sum_{X2} F_2 P(E_2|X_2) P(X_3|X_2) ]$ (so we marginalize out $X_2$)

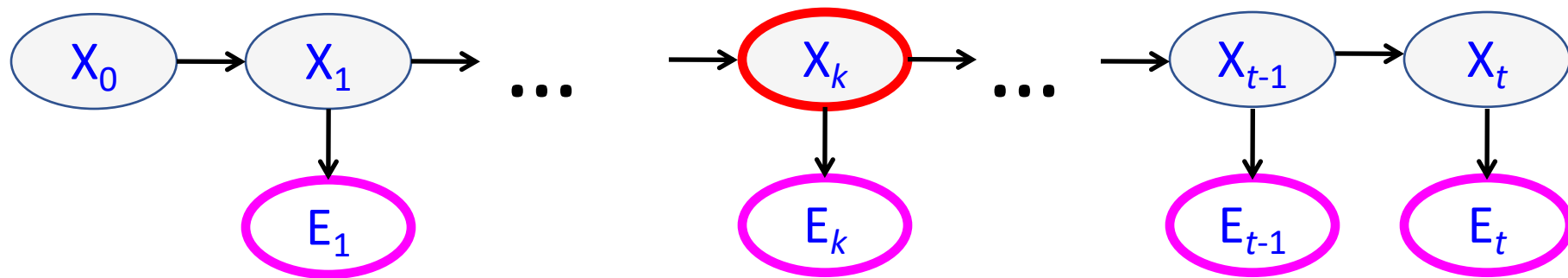$\Rightarrow P(X_3,E_1,E_2,E_3) = F_3 P(E_3|X_3)$

# HMM inference tasks

- **Filtering:** what is the distribution over the current state $X_t$ given all the evidence so far, $\mathbf{e}_{1:t}$ ?
  - The forward algorithm = sum-product algorithm for Xt given e1:t
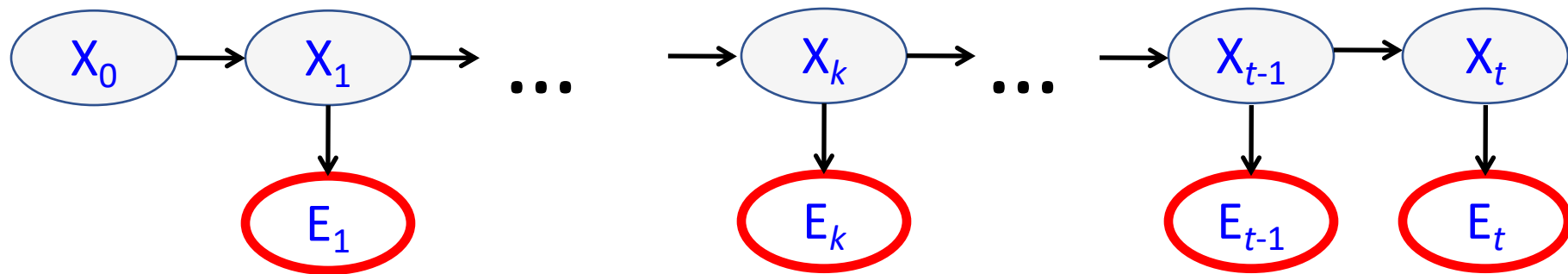


Query variable

Evidence variables

# HMM inference tasks

- **Filtering:** what is the distribution over the current state $X_t$ given all the evidence so far, $\mathbf{e}_{1:t}$ ?

- **Smoothing:** what is the distribution of some state $X_k$ given the entire observation sequence $\mathbf{e}_{1:t}$?
  - The forward-backward algorithm = sum-product algorithm for Xk given e1:t, when 1 < k < t
  - Xk = query variable, unknown, need to consider all its possible values
  - E1:t = evidence variables, known, only need to consider the given values

# HMM inference tasks

- **Filtering:** what is the distribution over the current state $X_t$ given all the evidence so far, $\mathbf{e}_{1:t}$ ?

- **Smoothing:** what is the distribution of some state $X_k$ given the entire observation sequence $\mathbf{e}_{1:t}$?

- **Evaluation:** compute the probability of a given observation sequence $\mathbf{e}_{1:t}$
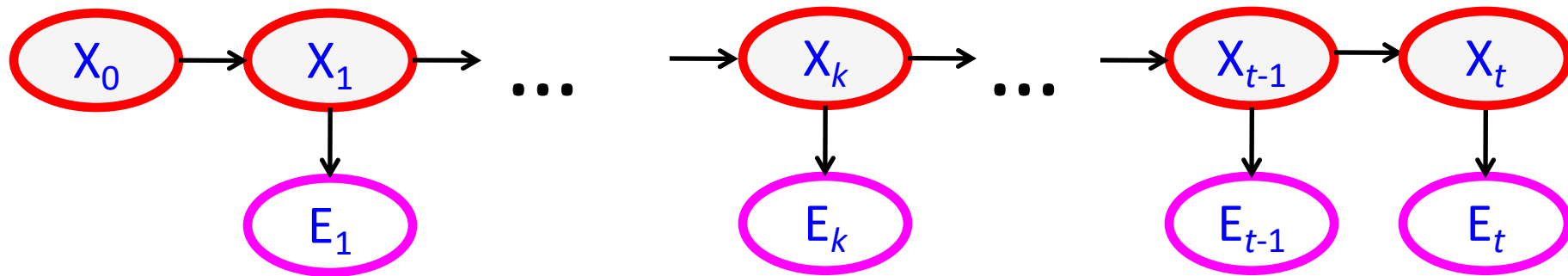
# HMM inference tasks

- **Filtering:** what is the distribution over the current state $X_t$ given all the evidence so far, $\mathbf{e}_{1:t}$

- **Smoothing:** what is the distribution of some state $X_k$ given the entire observation sequence $\mathbf{e}_{1:t}$?

- **Evaluation:** compute the probability of a given observation sequence $\mathbf{e}_{1:t}$

- **Decoding:** what is the most likely state sequence $\mathbf{X}_{0:t}$ given the observation sequence $\mathbf{e}_{1:t}$?
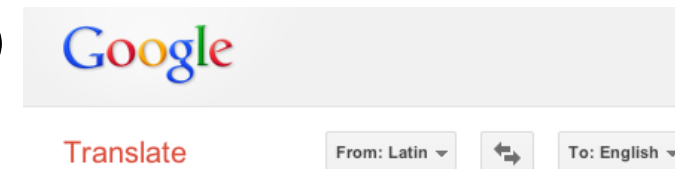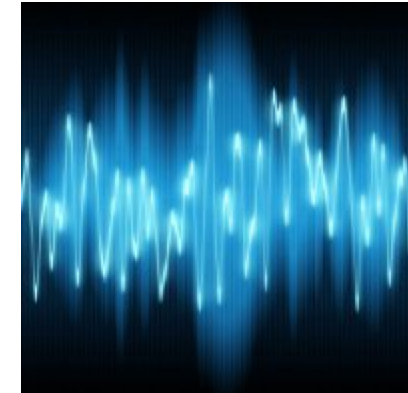  - The Viterbi algorithm

# HMM Learning and Inference

- Inference tasks
  - **Filtering:** what is the distribution over the current state $X_t$ given all the evidence so far, $\mathbf{e}_{1:t}$
  - **Smoothing:** what is the distribution of some state $X_k$ given the entire observation sequence $\mathbf{e}_{1:t}$?
  - **Evaluation:** compute the probability of a given observation sequence $\mathbf{e}_{1:t}$
  - **Decoding:** what is the most likely state sequence $\mathbf{X}_{0:t}$ given the observation sequence $\mathbf{e}_{1:t}$?
- Learning
  - Given a training sample of sequences, learn the model parameters (transition and emission probabilities)
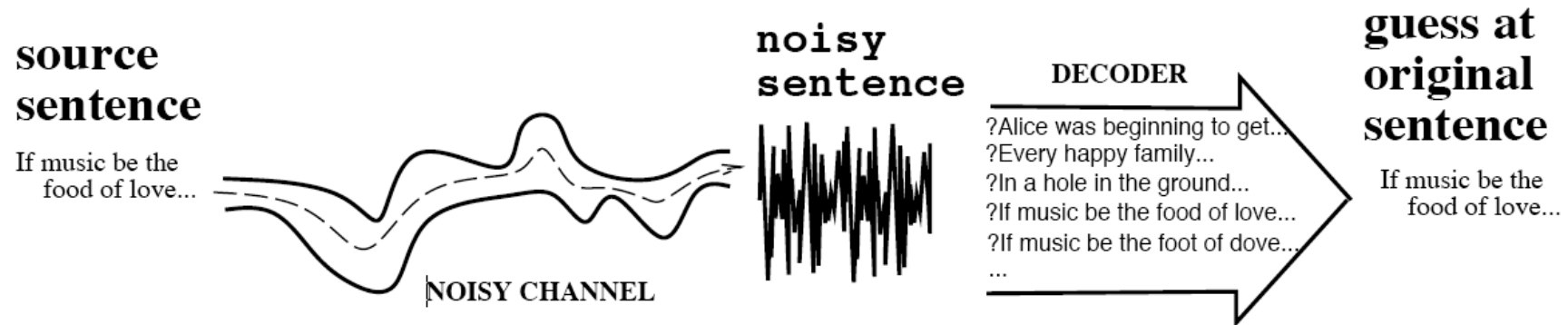    - EM algorithm

# Applications of HMMs

- Speech recognition HMMs:
  - Observations are acoustic signals (continuous valued)
  - States are specific positions in specific words (so, tens of thousands)

- Machine translation HMMs:
  - Observations are words (tens of thousands)
  - States are translation options

- Robot tracking:
  - Observations are range readings (continuous)
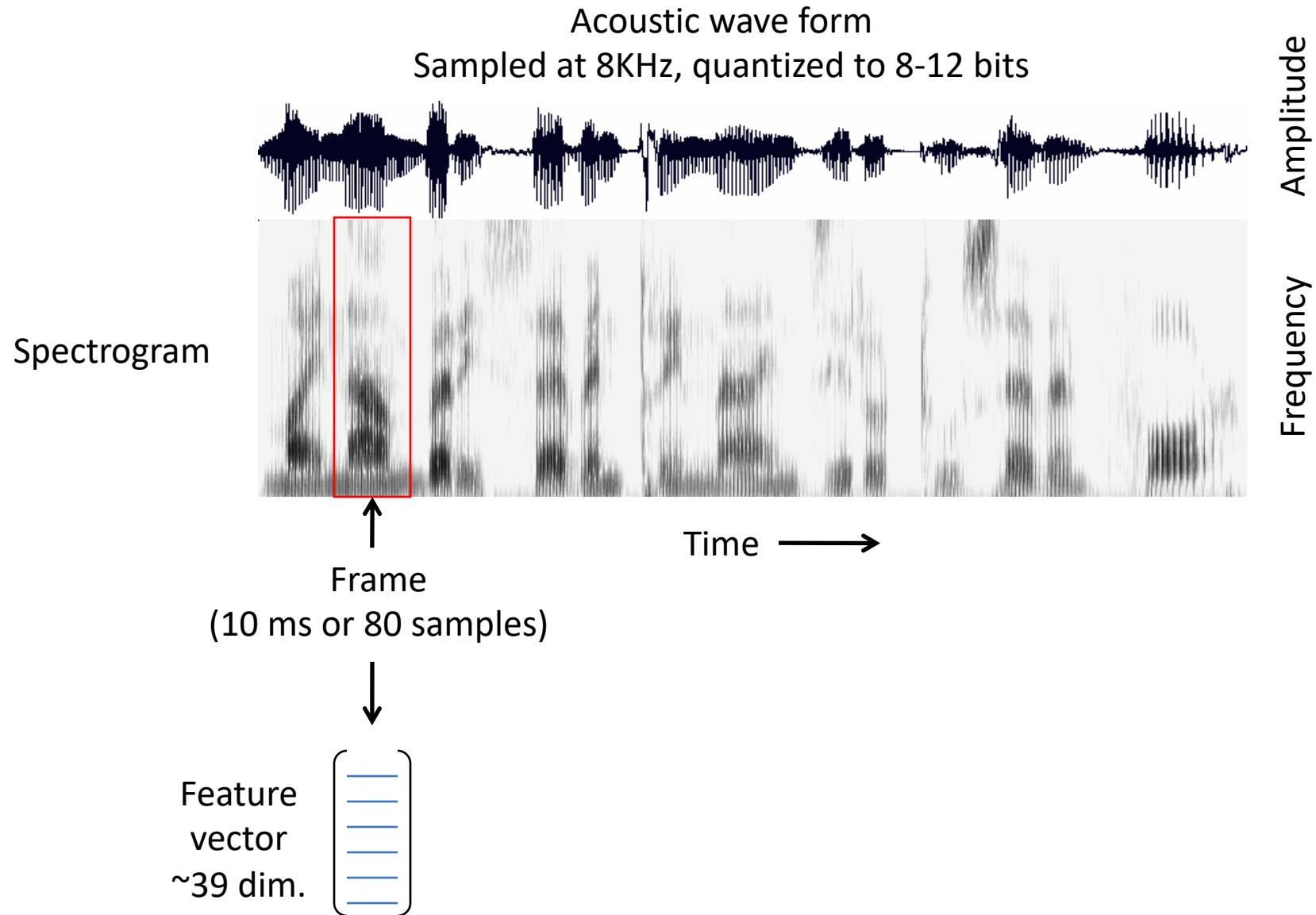  - States are positions on a map (continuous)

Source: Tamara Berg

# Application of HMMs: Speech recognition

- "Noisy channel" model of speech

# Speech feature extraction

# Speech feature extraction

Acoustic wave form
Sampled at 8KHz, quantized to 8-12 bits

Amplitude

Spectrogram

Frequency

Time

Frame
(10 ms or 80 samples)

Feature
vector
~39 dim.

# Phonetic model

- **Phones:** speech sounds
- **Phonemes:** groups of speech sounds that have a unique meaning/function in a language (e.g., there are several different ways to pronounce "t")

# Phonetic model

| IPA Symbol | ARPAbet Symbol | Word | IPA Transcription | ARPAbet Transcription |
|---|---|---|---|---|
| [p] | [p] | parsley | [ˈparsli] | [p aa r s l iy] |
| [t] | [t] | tarragon | [ˈtærəgɑn] | [t ae r ax g aa n] |
| [k] | [k] | catnip | [ˈkætnip] | [k ae t n ix p] |
| [b] | [b] | bay | [beɪ] | [b ey] |
| [d] | [d] | dill | [dɪl] | [d ih l] |
| [g] | [g] | garlic | [ˈgɑrlik] | [g aa r l ix k] |
| [m] | [m] | mint | [mɪnt] | [m ih n t] |
| [n] | [n] | nutmeg | [ˈnʌtmɛg] | [n ah t m eh g] |
| [ŋ] | [ng] | ginseng | [ˈdʒɪnsɪŋ] | [jh ih n s ix ng] |
| [f] | [f] | fennel | [ˈfɛnl] | [f eh n el] |
| [v] | [v] | clove | [kloʊv] | [k l ow v] |
| [θ] | [th] | thistle | [ˈθɪsl] | [th ih s el] |
| [ð] | [dh] | heather | [ˈhɛðɚ] | [h eh dh axr] |
| [s] | [s] | sage | [seɪdʒ] | [s ey jh] |
| [z] | [z] | hazelnut | [ˈheɪzlnʌt] | [h ey z el n ah t] |
| [ʃ] | [sh] | squash | [skwɑʃ] | [s k w a sh] |
| [ʒ] | [zh] | ambrosia | [æmˈbroʊʒə] | [ae m b r ow zh ax] |
| [tʃ] | [ch] | chicory | [ˈtʃɪkɚi] | [ch ih k axr iy ] |
| [dʒ] | [jh] | sage | [seɪdʒ] | [s ey jh] |
| [l] | [l] | licorice | [ˈlɪkɚɪʃ] | [l ih k axr ix sh] |
| [w] | [w] | kiwi | [ˈkiwi] | [k iy w iy] |
| [r] | [r] | parsley | [ˈpɑrsli] | [p aa r s l iy] |
| [j] | [y] | yew | [yu] | [y uw] |
| [h] | [h] | horseradish | [ˈhɔrsrædɪʃ] | [h ao r s r ae d ih sh] |
| [ʔ] | [q] | uh-oh | [ʔʌʔoʊ] | [q ah q ow] |
| [ɾ] | [dx] | butter | [ˈbʌɾɚ] | [b ah dx axr ] |
| [ɾ̃] | [nx] | wintergreen | [wɪɾ̃ɚgrin] | [w ih nx axr g r i n ] |
| [l̩] | [el] | thistle | [ˈθɪsl] | [th ih s el] |

**Figure 4.1** IPA and ARPAbet symbols for transcription of English consonants.

| IPA Symbol | ARPAbet Symbol | Word | IPA Transcription | ARPAbet Transcription |
|---|---|---|---|---|
| [i] | [iy] | lily | [ˈlɪli] | [l ih l iy] |
| [ɪ] | [ih] | lily | [ˈlɪli] | [l ih l iy] |
| [eɪ] | [ey] | daisy | [ˈdeɪzi] | [d ey z i] |
| [ɛ] | [eh] | poinsettia | [pɔmˈsɛriə] | [p oy n s eh dx iy ax] |
| [æ] | [ae] | aster | [ˈæstɚ] | [ae s t axr] |
| [ɑ] | [aa] | poppy | [ˈpɑpi] | [p aa p i] |
| [ɔ] | [ao] | orchid | [ˈɔrkɪd] | [ao r k ix d] |
| [ʊ] | [uh] | woodruff | [ˈwʊdrʌf] | [w uh d r ah f] |
| [oʊ] | [ow] | lotus | [ˈloʊɾəs] | [l ow dx ax s] |
| [u] | [uw] | tulip | [ˈtulɪp] | [t uw l ix p] |
| [ʌ] | [uh] | buttercup | [ˈbʌɾɚˌkʌp] | [b uh dx axr k uh p] |
| [ɝ] | [er] | bird | [ˈbɝd] | [b er d] |
| [aɪ] | [ay] | iris | [ˈaɪrɪs] | [ay r ix s] |
| [aʊ] | [aw] | sunflower | [ˈsʌnflaʊɚ] | [s ah n f l aw axr] |
| [ɔɪ] | [oy] | poinsettia | [pɔmˈsɛriə] | [p oy n s eh dx iy ax] |
| [ju] | [y uw] | feverfew | [fivɚfju] | [f iy v axr f y u] |
| [ə] | [ax] | woodruff | [ˈwʊdrəf] | [w uh d r ax f] |
| [ɚ] | [axr] | heather | [ˈhɛðɚ] | [h eh dh axr] |
| [ɨ] | [ix] | tulip | [ˈtulɪp] | [t uw l ix p] |
| [ʉ] | [ux] | | [] | [] |

**Figure 4.2** IPA and ARPAbet symbols for transcription of English vowels

# HMM models for phones

HMM states in most speech recognition systems correspond to **_subsegments_** _of_ **_triphones_**

- **_Triphone_**: the /b/ in "about" (ax-b+aw) sounds different from the /b/ in "Abdul" (ae-b+d). There are around 60 phones and as many as $60^3$ context-dependent _triphones._
- **_Subsegments_**: /b/ has three subsegments: the closure, the silence, and the release. There are $3\times60^3$ subsegments of triphones.
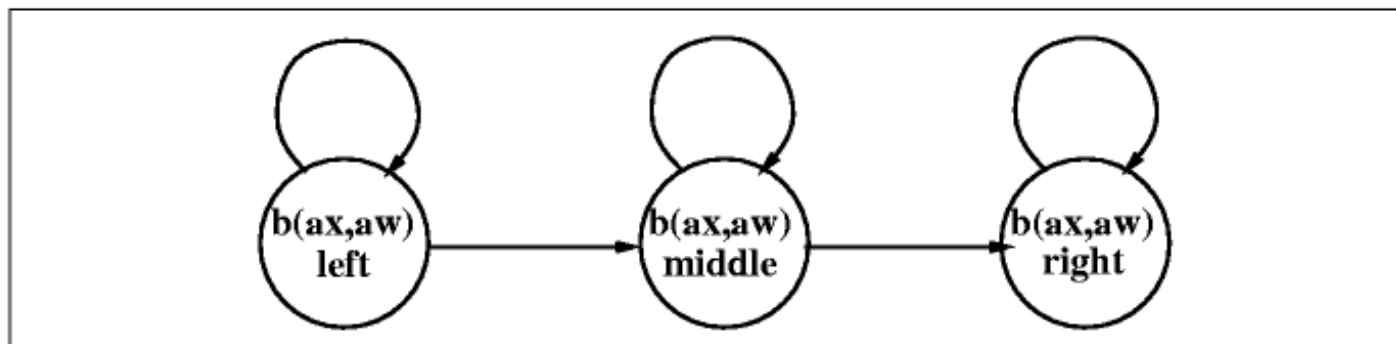


**Figure 7.11** An example of the context-dependent triphone b(ax,aw) (the phone [b] preceded by a [ax] and followed by a [aw], as in the beginning of *about*, showing its left, middle, and right subphones.
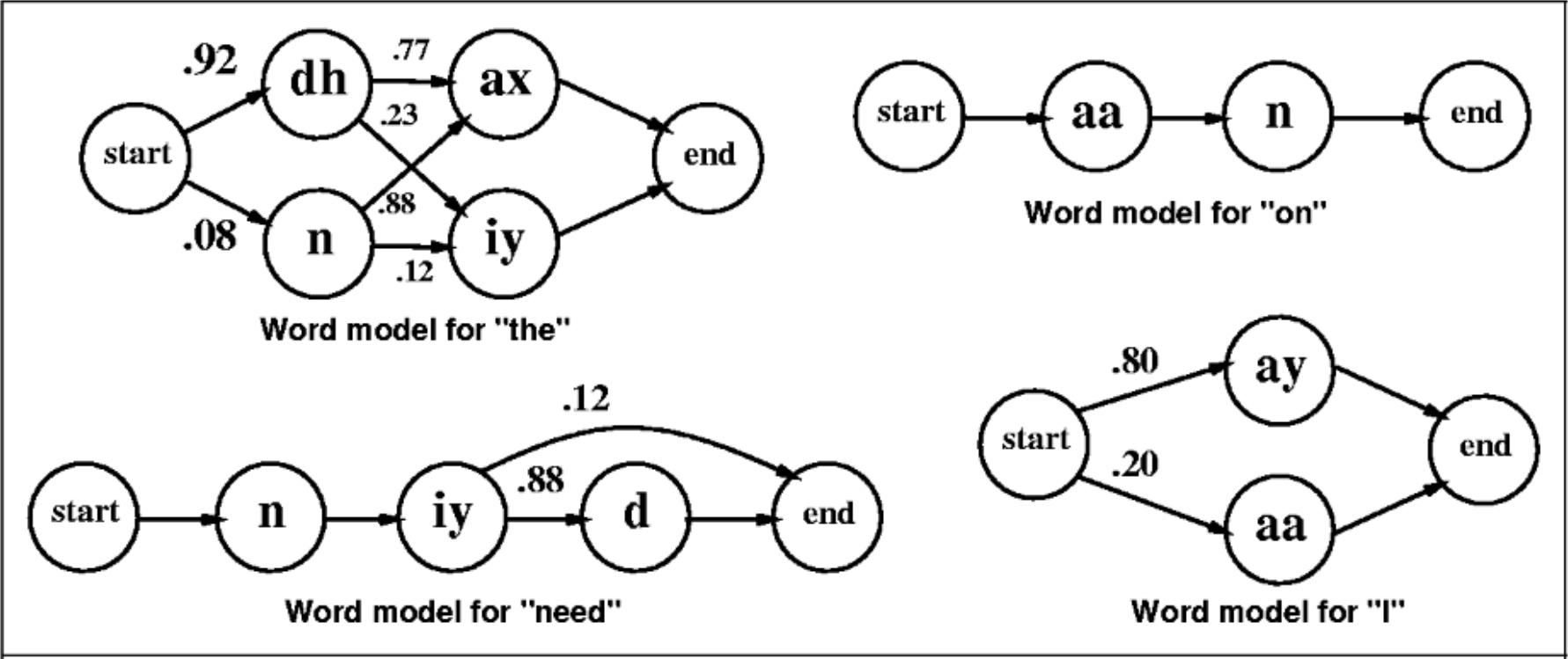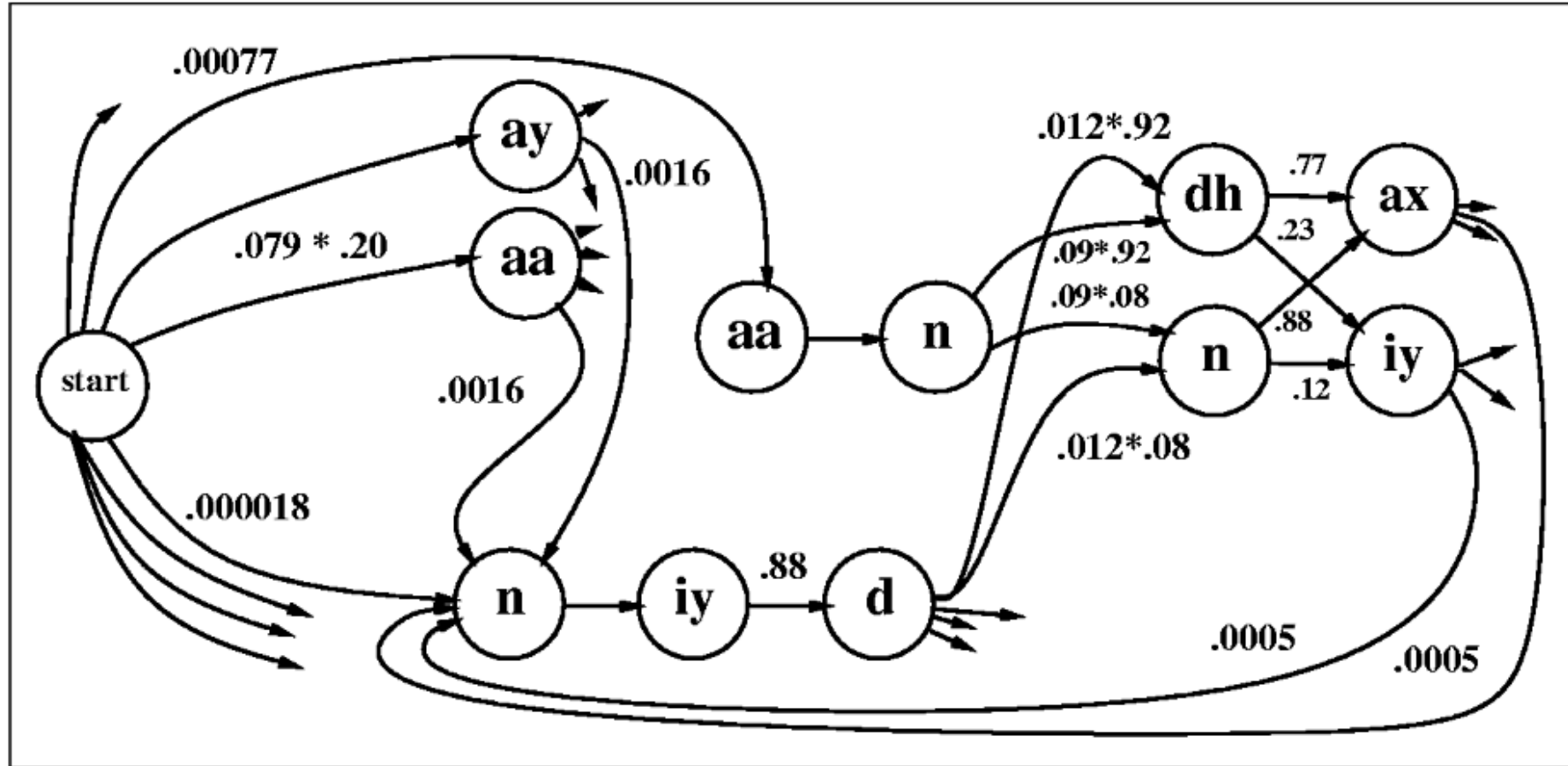
# HMM models for words



**Figure 7.5** Pronunciation networks for the words *I*, *on*, *need*, and *the*. All networks (especially *the*) are significantly simplified.
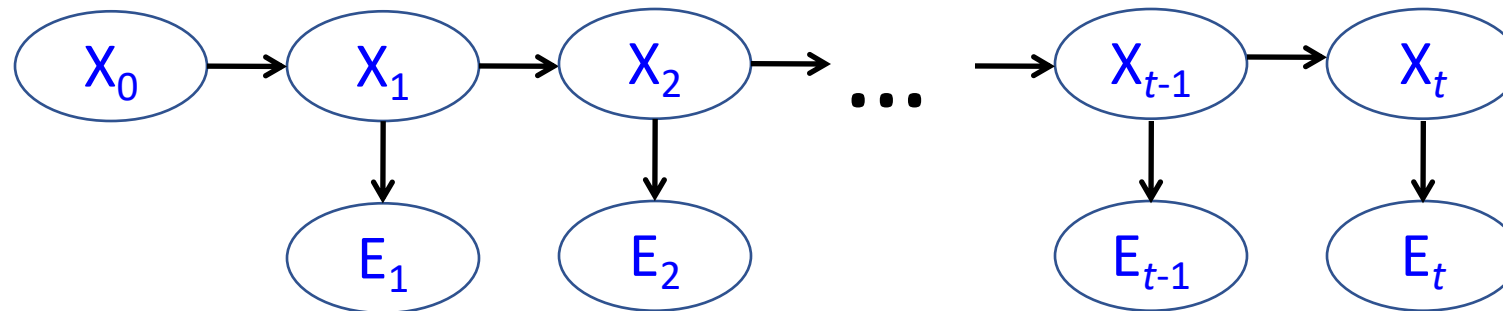
# Putting words together



- Given a sequence of acoustic features, how do we find the corresponding word sequence?

# The Viterbi Algorithm

$$\max_{\boldsymbol{X}_{0:t}} P(\boldsymbol{X}_{0:t}, \boldsymbol{E}_{0:t})$$

$$= \max_{X_t} P(E_t|X_t) \max_{X_{t-1}} P(X_t|X_{t-1})P(E_{t-1}|X_{t-1}) \max_{X_{t-2}} \dots$$

Complexity changes from O{N^T} to O{TN^2}
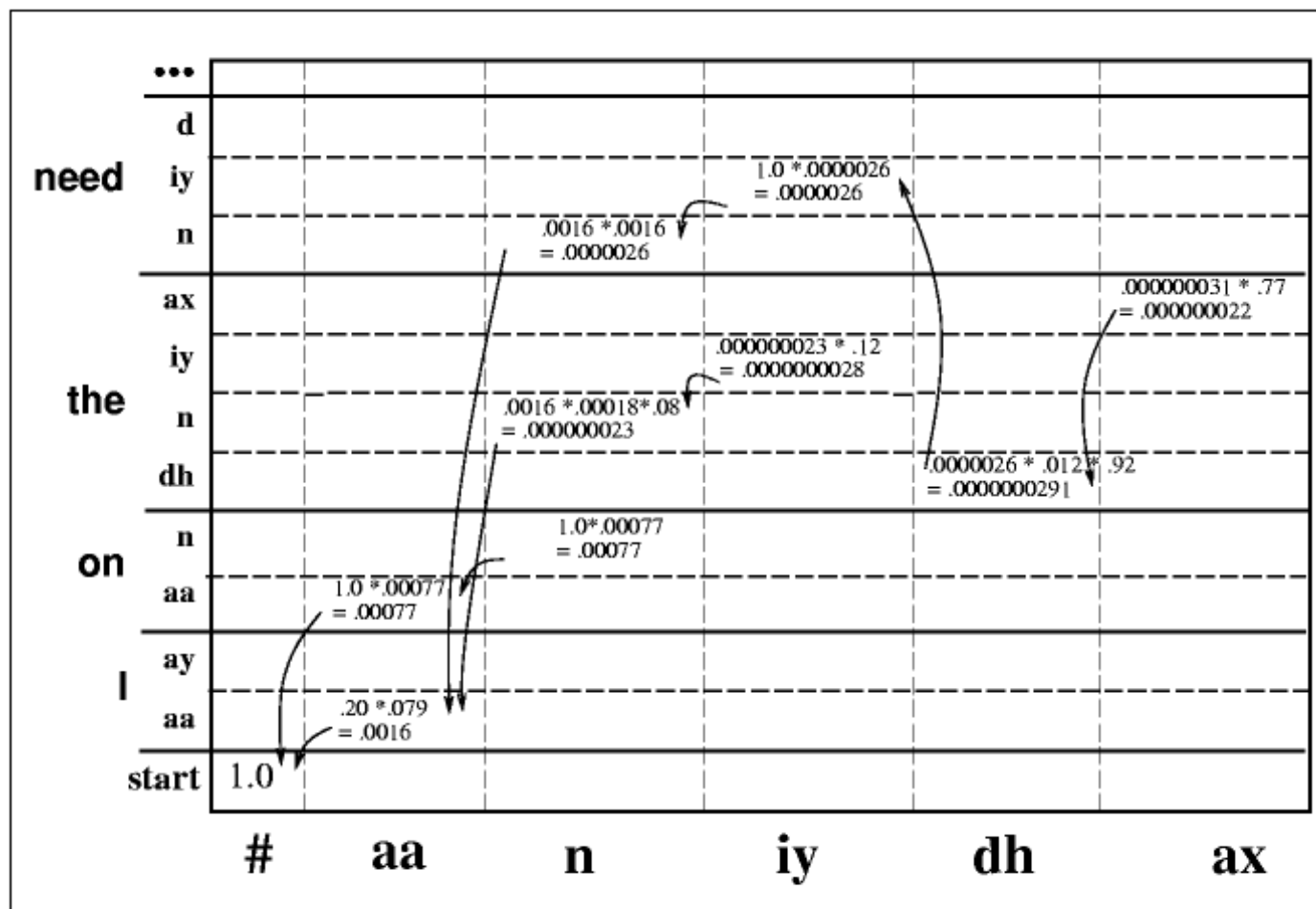
# Decoding with the Viterbi algorithm



**Figure 7.10**  The entries in the individual state columns for the Viterbi algorithm. Each cell keeps the probability of the best path so far and a pointer to the previous cell along that path. Backtracing from the successful last word (*the*), we can reconstruct the word sequence *I need the*.

# For more information

- CS 447: Natural Language Processing
- ECE 417: Multimedia Signal Processing
- ECE 594: Mathematical Models of Language
- Linguistics 506: Computational Linguistics
- D. Jurafsky and J. Martin, "Speech and Language Processing," 2nd ed., Prentice Hall, 2008