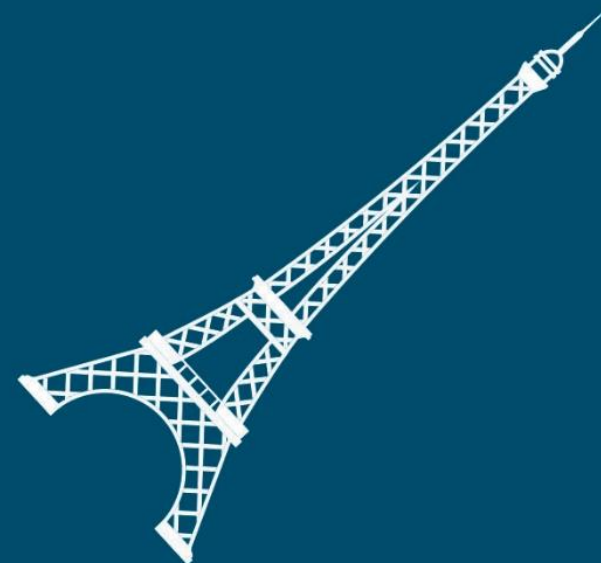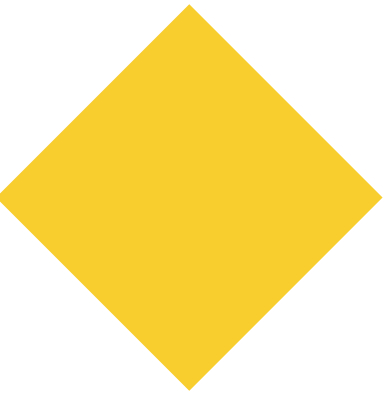# A.I. IN AUDIO & SIGNAL PROCESSING

Session 3: HMM for speech processing

# COURSE STRUCTURE

Quick Summary

**Audio processing for AI**
- Signal, audio, speech encoding (4h)
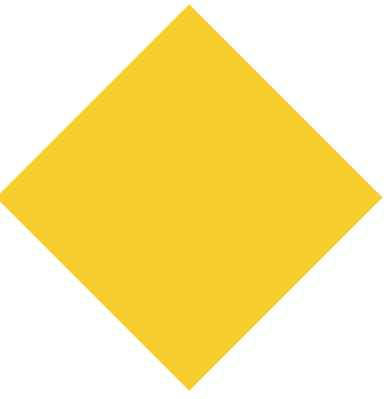- Deep learning for audio processing (4h+4h)

**Automata for language modelling**
- HMM for speech processing (4h)
- Automata and transducer (4h)

**Towards speaking with an AI-bot**
- Speech synthesis (4h)
- Automatic speech recognition (4h)
- Speaker and emotion recognition (4h)
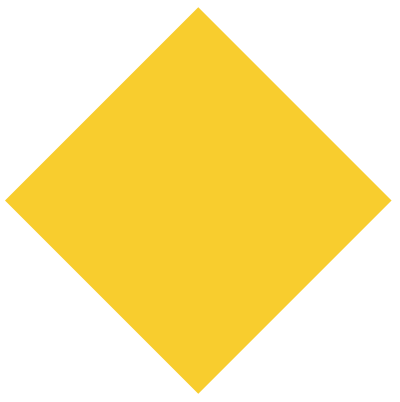
# SESSION 3: HMM FOR SPEECH PROCESSING

Quick Summary

1. **Markov models & HMM**

2. **Scoring a sentence**

3. **Decoding a sequence of phonems**

4. **Training a language model**

# HMM FOR SPEECH PROCESSING.

Markov models and HMM

# MARKOV MODELS & HMM

## Markov property defining a Markov Model

$$\forall n \geq 0, (i_0, \ldots, i_{n-1}, i, j) \in \boldsymbol{E}^{n+2},$$
$$P(X_{n+1} = j \mid X_0 = i_0, X_1 = i_1, \ldots, X_{n-1} = i_{n-1}, X_n = i) = P(X_{n+1} = j \mid X_n = i)$$

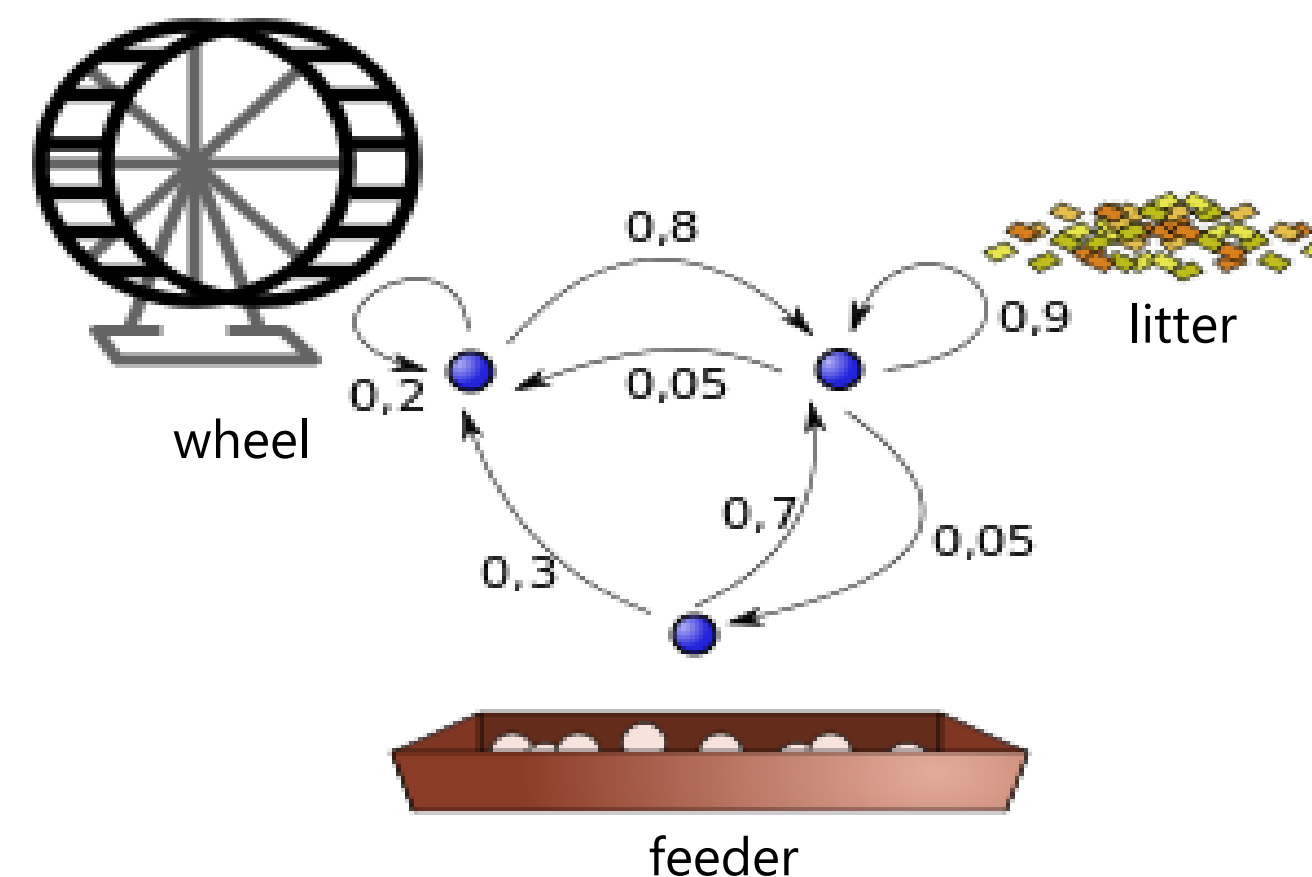We consider homogeneous models ($p_{i,j}$ is constant over time).
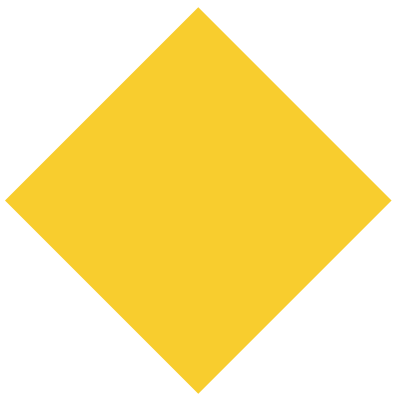
| $E$ | states space |
|---|---|
| $X_0, X_1, \ldots, X_{n-1}, X_n$ | random variable sequence of successive states |
| $p_{i,j}$ | transition probability from state i to state j |
| $n$ | time index (noted t further) |

## Transition probability

$$p_{i,j} = P(X_1 = j \mid X_0 = i) \qquad \text{with} \qquad \forall i \in \boldsymbol{E}, \ \sum_{j \in E} p_{i,j} = 1$$

## Example of Markov process

- Hamster pet

  $\rightarrow$ hamster activity at $t_n$ is predictable,
  knowing its activity at $t_0$



0,8
0,9  litter
wheel
0,2    0,05
0,7    0,05
0,3
feeder

# MARKOV MODELS & HMM

## Hidden Markov Model

Markov model with "partially observable" states
Usually, part only of the model is known:
 → either the sequence of observations O is unknown
 → either the sequence of states Q is unknown
 → either the transition probabilities are unknown

## Elements of a discrete HMM

$S = \{S_0, S_1, \dots, S_N\}$     set of possible states
$V = \{V_0, V_1, \dots, V_M\}$     set of possible observations

$Q = (q_0, q_1, \dots, q_T)$     sequence of states with $t$ from 0 to $T$
$\mathcal{O} = (\sigma_0, \sigma_1, \dots, \sigma_T)$     sequence of observations with $t$ from 0 to $T$

$a_{i,j} = P(q_{t+1} = S_j | q_t = S_i)$     state transition probability (matrix $A$)
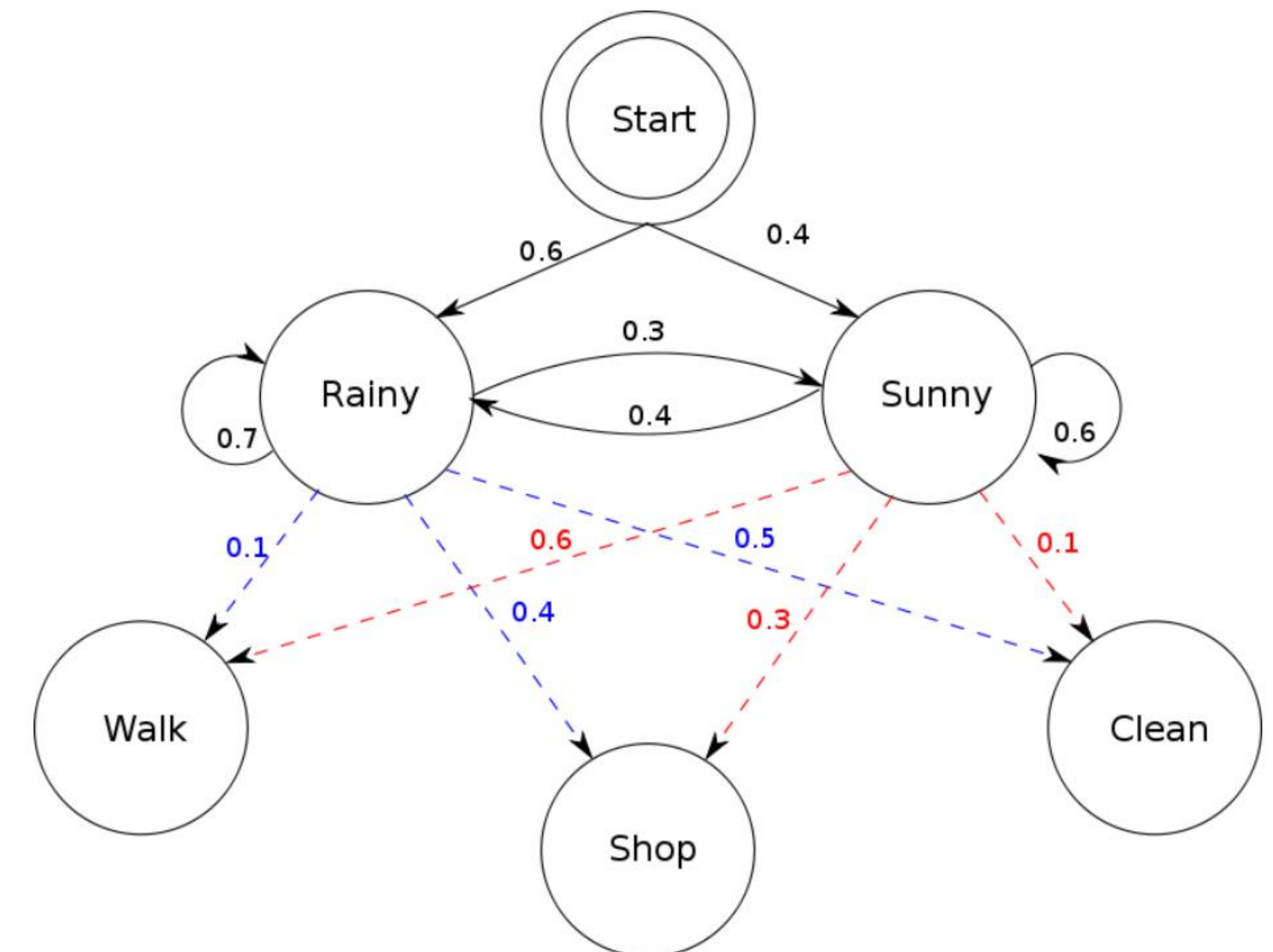$b_j(k) = P(\sigma_t = V_k | q_t = S_j)$     observation probability (matrix $B$)

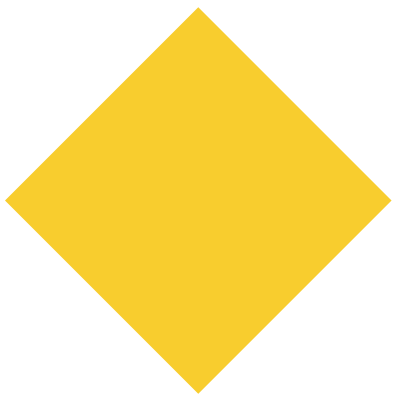$\pi = \{\pi_0, \pi_2, \dots, \pi_N\}$     initial state distribution, with $\pi_i = P(q_0 = S_i)$
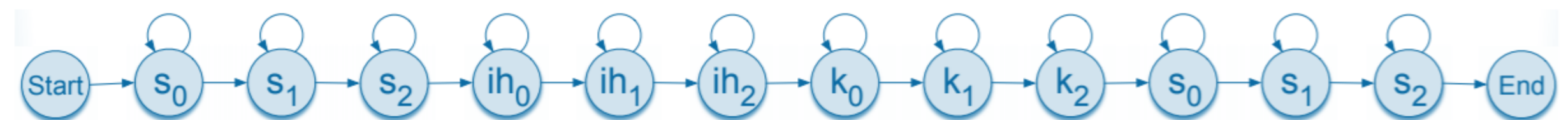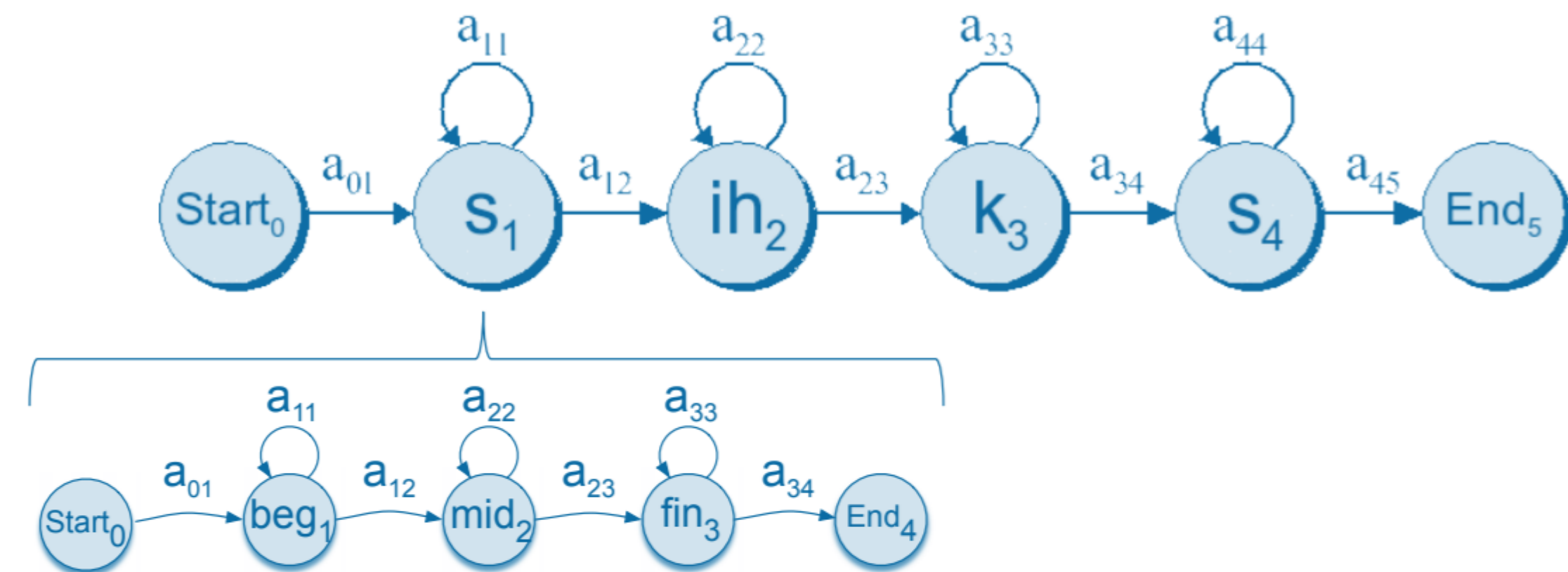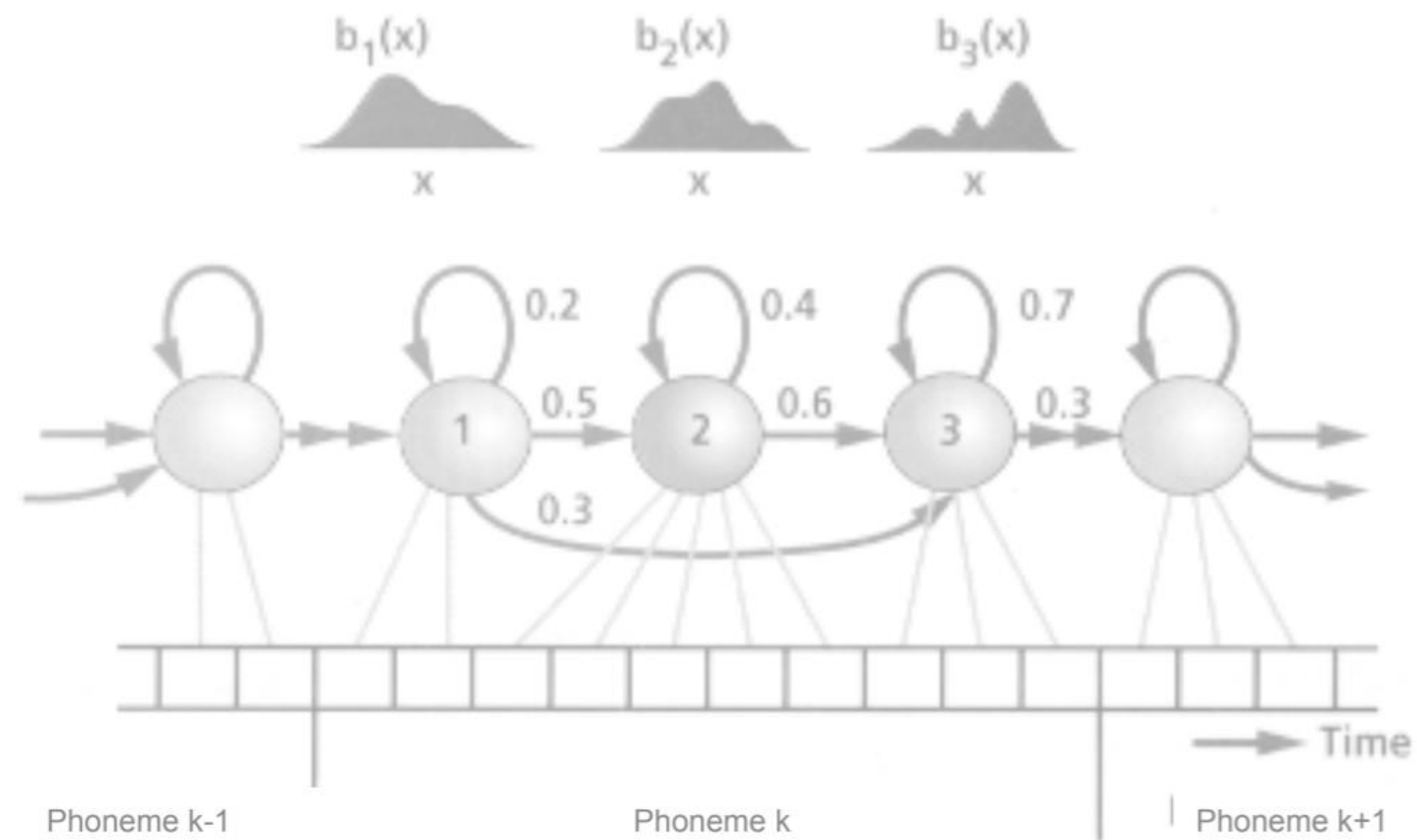
## Example of Hidden Markov Model
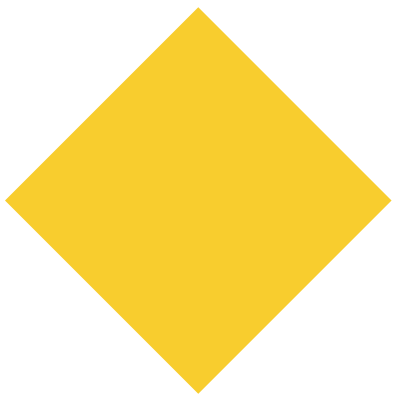
- activity depending on the weather



L. Baum, J. A. Eagon, T. Petrie. 1966-1970.

# MARKOV MODELS & HMM

HMM application to speech

Basic problems for HMMs

- Scoring
    Given the state sequence $Q = (q_0, q_1, …, q_T)$, and a model $\lambda = (A, B, \pi)$, how do we efficiently compute $\mathrm{P}(\mathcal{O}|\lambda, Q)$, the probability of the observation sequence, given the model?

    → **Forward algorithm**

- Matching/Decoding
    Given the observation sequence $\mathcal{O} = (\sigma_0, \sigma_1, …, \sigma_T)$, and a model λ, how do we choose a corresponding state sequence $Q = (q_0, q_1, …, q_T)$ which is optimal in some meaningful sense (i.e., best "explains" the observations). $\mathrm{P}(Q|\lambda, \mathcal{O})$?

    → **Viterbi algorithm**

- Training
    How do we adjust the model parameters model $\lambda = (A, B, \pi)$ to maximize $\mathrm{P}(\lambda|\mathrm{Q}, \mathcal{O})$?
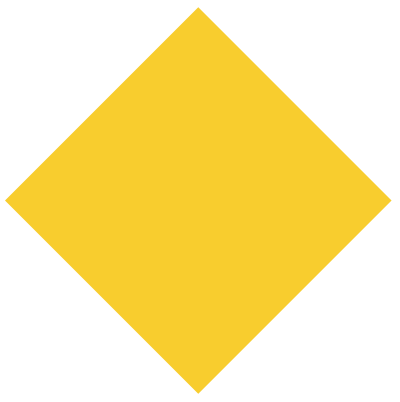
    → **Baum-Welch re-estimation procedures**
        (known as forward-backward algorithm)

# HMM FOR SPEECH PROCESSING.

Scoring a sentence

# SCORING A SENTENCE

## Goal

Find $P(\mathcal{O}|\lambda)$,

$P(\mathcal{O}|\lambda)$: probability to observe $\mathcal{O} = (\sigma_0, \sigma_1, \ldots, \sigma_n)$, knowing the model $\lambda = (A, B, \pi)$

## Analytical solving

law of total probability

$$(1) \qquad P(\mathcal{O}|\lambda) = \sum_{all\ Q} P(\mathcal{O}|Q, \lambda)\, P(Q|\lambda)$$

Indépendance of observations knowing Q

$\sigma_t$ depends on $q_t$ and $q_0, q_1, \ldots, q_{t-1}$

besides, as $Q$ follow Markov property

$$(2) \qquad P(\mathcal{O}|Q, \lambda) = \prod_{t=0}^{T} P(\sigma_t|Q, \lambda)$$

$$P(\sigma_t|Q, \lambda) = P(\sigma_t|q_t, \lambda)$$
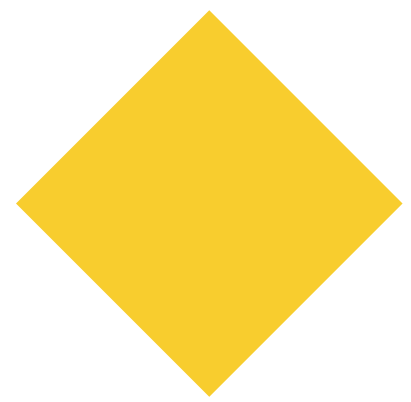$$= b_{q_t}(\sigma_t) \qquad \text{by definition}$$

initial state and transition probabilities

$$(3) \qquad P(Q|\lambda) = \pi_{q_0} \prod_{t=1}^{T} a_{q_{t-1}, q_t}$$

(1), (2) and (3) give the result

$$(4) \qquad P(\mathcal{O}|\lambda) = \sum_{all\ Q} \left[ \pi_{q_0} . b_{q_0}(\sigma_0) . \prod_{t=1}^{T} a_{q_{t-1}, q_t} . b_{q_t}(\sigma_t) \right]$$

# MARKOV MODELS & HMM

Computational solving: Forward algorithm

- Initialization

$$\alpha_0(i) = \pi_i \cdot b_i(\sigma_0) \quad \text{for } i \in [\![0, N]\!]$$

- Induction

$$\alpha_t(j) = \left[ \sum_{i=0}^{N} \alpha_{t-1}(i) \cdot a_{i,j} \right] \cdot b_j(\sigma_t) \quad \text{for t} \in [\![1, T]\!], \ j \in [\![0, N]\!]$$

- Termination

$$P(\mathcal{O}|\lambda) = \sum_{i=0}^{N} \alpha_T(i)$$



$a_{i,j} = P(q_{t+1} = S_j | q_t = S_i)$     state transition probability (matrix $A$)
$b_j(k) = P(\sigma_t = V_k | q_t = S_j)$     observation probability (matrix $B$)

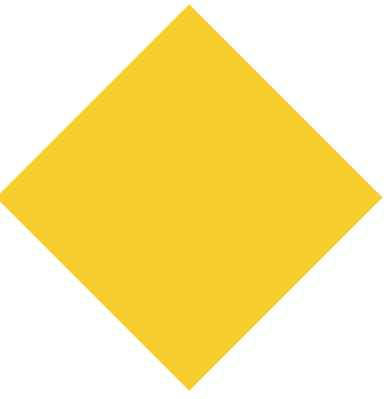$\pi = \{\pi_0, \pi_2, \dots, \pi_N\}$     initial state distribution, with $\pi_i = P(q_0 = S_i)$

# HMM FOR SPEECH PROCESSING.

Decoding a sequence of phonems
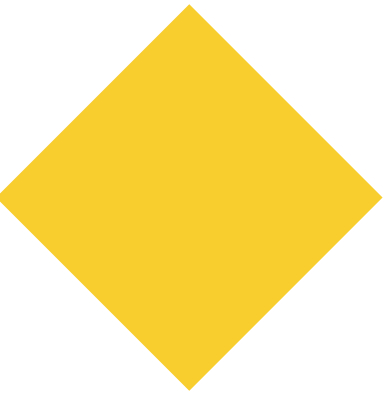
# DECODING A SEQUENCE OF PHONEMS

Goal

Find most probable sequence of state $Q = (q_0, q_1, \dots, q_T)$, given observations $\mathcal{O}$ and model $\lambda$.
$\rightarrow$ find $Q$ maximizing $\text{P}(Q|\mathcal{O}, \lambda)$

Forward algorithm provides a probability through all path sequence $Q$
$\rightarrow$ find the optimum path sequence

Solving approaches

- Consider the path sequence maximizing successively each $a_{i,j}$
$\rightarrow$ possibly not optimal

- Consider the path sequence maximizing $\text{P}(Q|\mathcal{O}, \lambda)$ with respect to the whole sequence
$\rightarrow$ Viterbi algorithm

# DECODING A SEQUENCE OF PHONEMS

Analytical solving

From equations (1) and (4), (see scoring previous chapter)

$$P(Q|\mathcal{O}, \lambda) = \pi_{q_0} . b_{q_0}(\sigma_0) . \prod_{t=1}^{T} a_{q_{t-1}, q_t} . b_{q_t}(\sigma_t)$$

$$\underbrace{\phantom{b_{q_0}(\sigma_0) . \prod_{t=1}^{T} a_{q_{t-1}, q_t} . b_{q_t}(\sigma_t)}}_{\delta_T}$$

$$P(Q|\mathcal{O}, \lambda) = \pi_{q_0} . b_{q_0}(\sigma_0) . \prod_{t=1}^{T} a_{q_{t-1}, q_t} . b_{q_t}(\sigma_t)$$

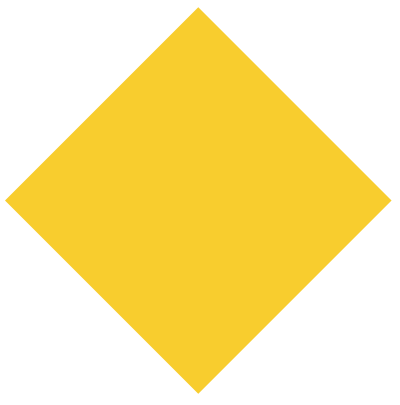Idea to compute iteratively overtime the probability $\delta_t$ for t$\in \llbracket 1, T \rrbracket$

$$\delta_t(j) = \max_{0 \le i \le N} (\delta_{t-1} \, a_{i,j}) . b_j(\sigma_t)$$

And thus, compute at each time step t, the most likely state transition

---

Viterbi algorithm assumptions
- $\mathcal{O}$ and $Q$ are both in sequences
- $\mathcal{O}$ and $Q$ are isomorphic (one observed event per hidden event)
- $Q$ verifies Markov property

# DECODING A SEQUENCE OF PHONEMS

Viterbi algorithm

1. Initialization:

$$\delta_1(i) = \pi_i b_i(o_1), \qquad 1 \leq i \leq N$$
$$\psi_1(i) = 0$$

2. Recursion:

$$\delta_t(j) = \max_{1 \leq i \leq N}[\delta_{t-1}(i)a_{ij}]b_j(o_t), \qquad 2 \leq t \leq T \qquad 1 \leq j \leq N$$
$$\psi_t(j) = \operatorname*{argmax}_{1 \leq i \leq N}[\delta_{t-1}(i)a_{ij}], \qquad 2 \leq t \leq T \qquad 1 \leq j \leq N$$
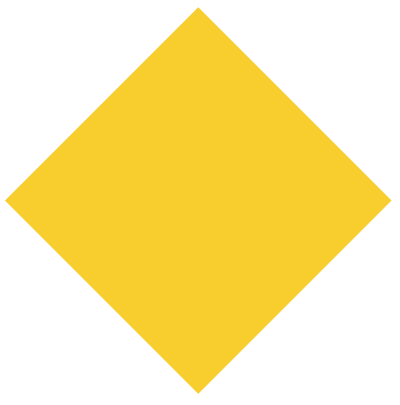
3. Termination:

$$P^* = \max_{1 \leq i \leq N}[\delta_T(i)]$$
$$q_T^* = \operatorname*{argmax}_{1 \leq i \leq N}[\delta_T(i)]$$
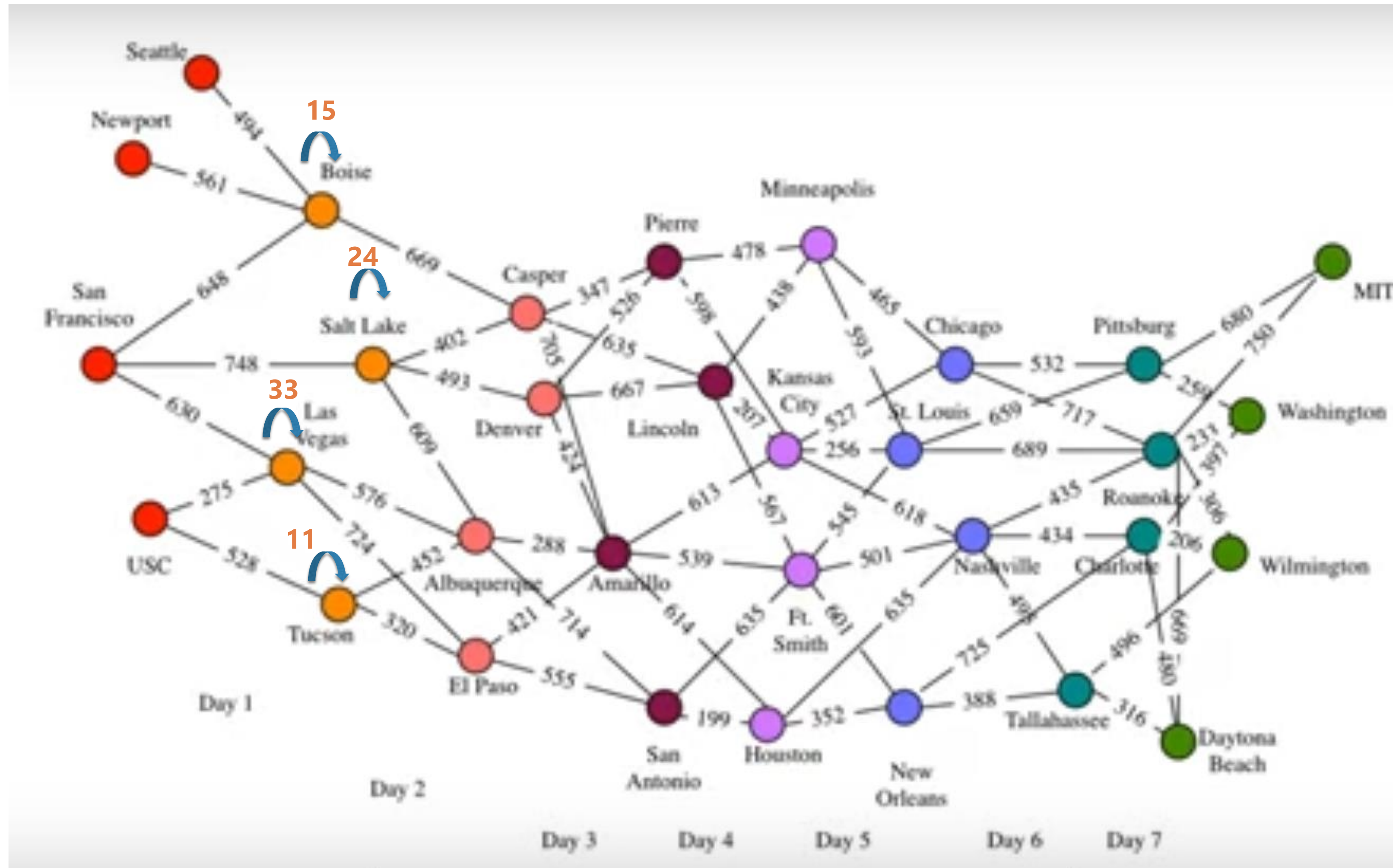
4. Path (state-sequence) backtracking:

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

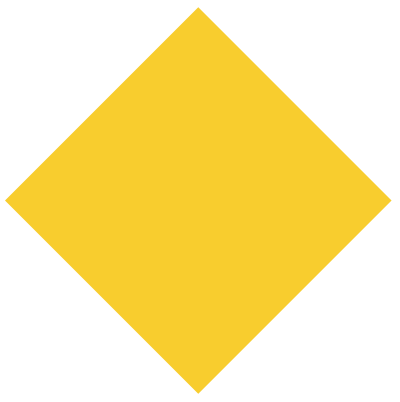Example

# DECODING A SEQUENCE OF PHONEMS

Viterbi algorithm

# HMM FOR SPEECH PROCESSING.

Training a language model

# TRAINING A LANGUAGE MODEL

<u>Goal</u>

Adjusting model parameters to maximize $P(Q, \mathcal{O}|\lambda)$.
    $\mathcal{O} = (\sigma_0, \sigma_1, \dots, \sigma_T)$  is one of the training sequence
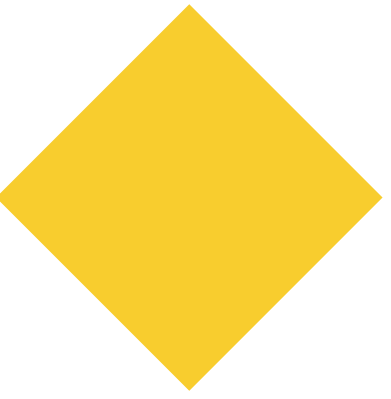
<u>Analytical solving</u>

    $\rightarrow$ none

<u>Baum-Welch re-estimation procedures</u>

    Iterative algorithm that:
        - Compute statistics on the current model given the training data
        - Adapt the model given the previous statistics
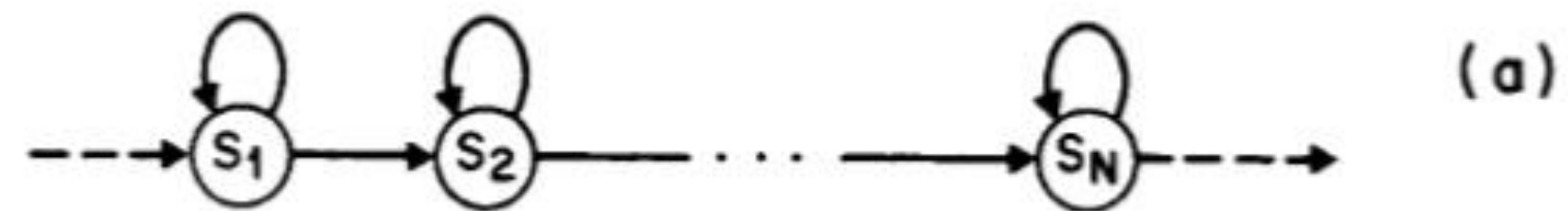        - Return to 1$^{st}$ step until convergence

    Also known as forward-backward algorithm

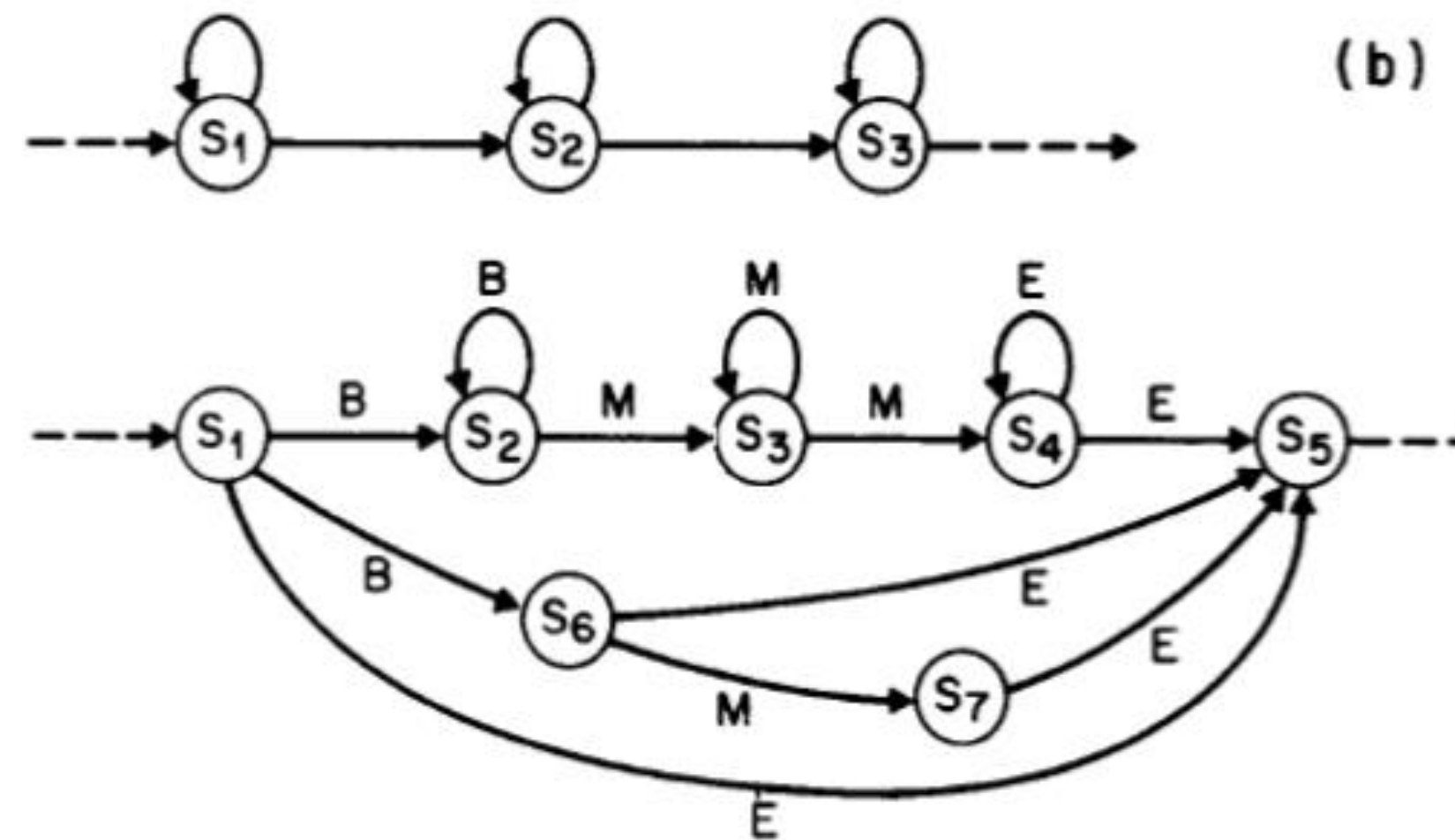Language model using HMM



WORD MODEL
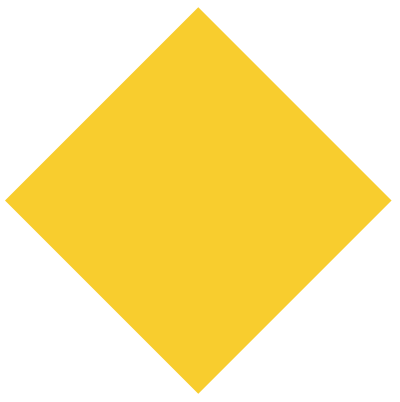
(a)

SUB-WORD UNIT

(b)

# HMM FOR SPEECH PROCESSING.

Thank you for your attention.

References:
- Xavier Anguera

# PRACTICAL EXERCISE

1. Modelize Rainy-sunny model with hmmlearn

Example of Hidden Markov Model

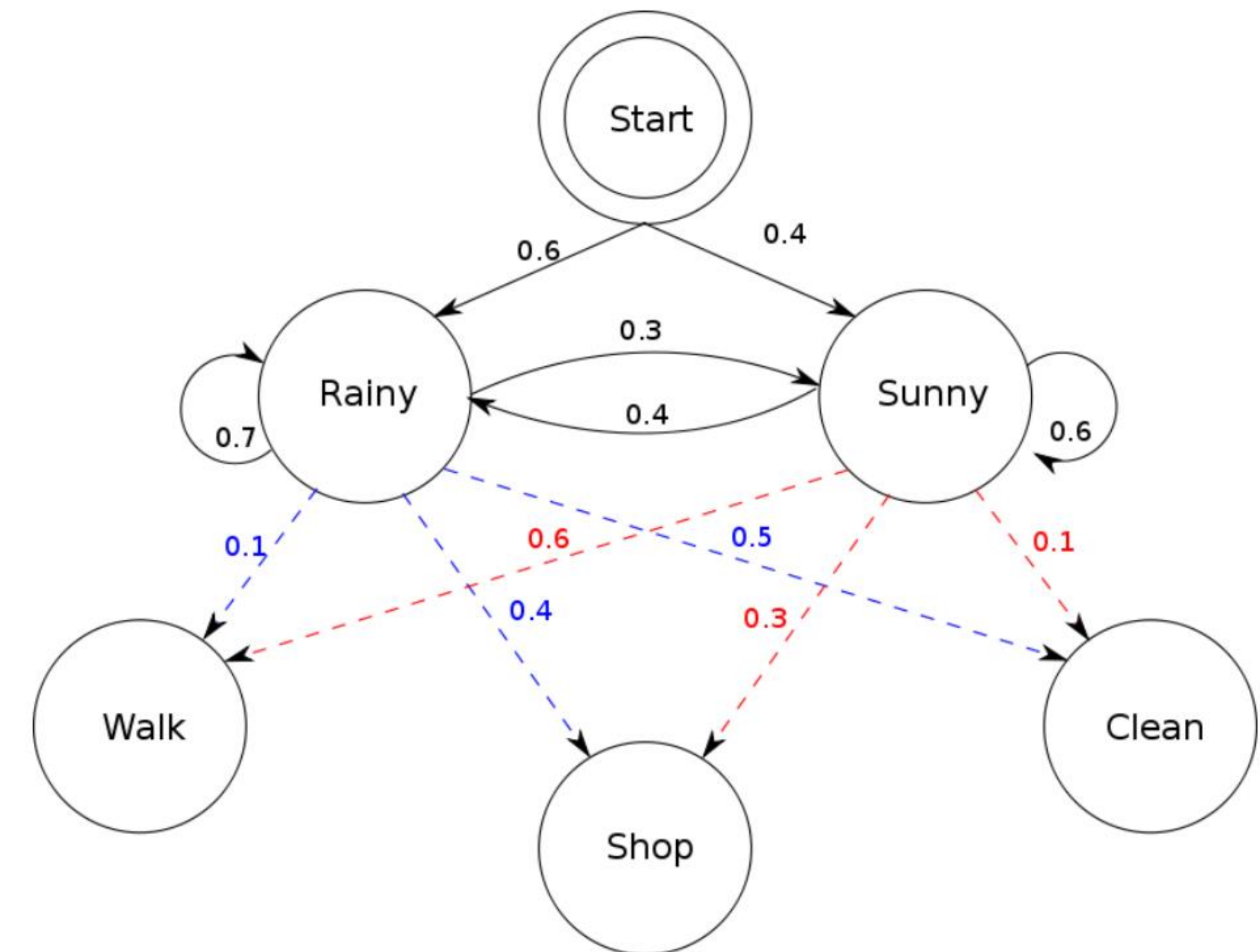- activity depending on the weather

Use the following items:
- from hmmlearn import hmm
- MultinomialHMM
- startprob_
- transmat_
- emissionprob_

TO DO:
- write starting probability
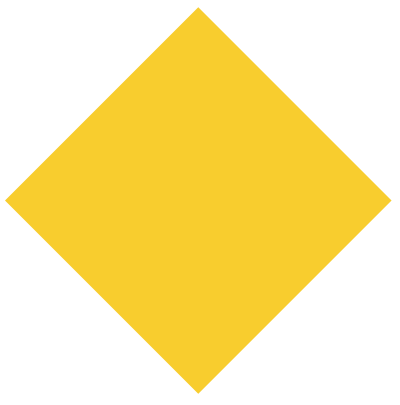- transition matrix
- emission probability

# PRACTICAL EXERCISE

2. Solve scoring problem

    Find probability of observations for the following
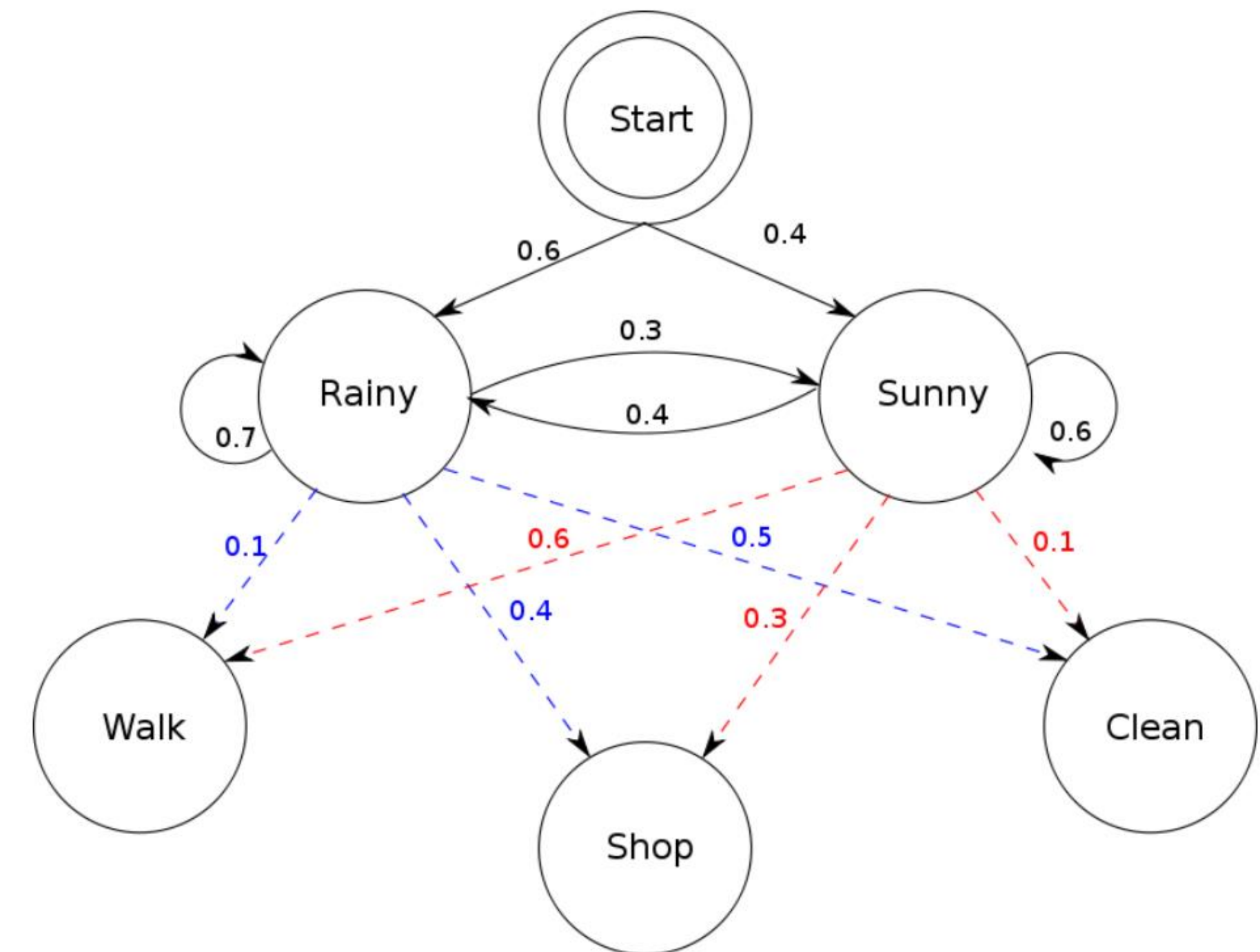sequences of states:
- (Start)
- (Rainy)
- (Sunny)
- (Sunny, Sunny, Sunny)
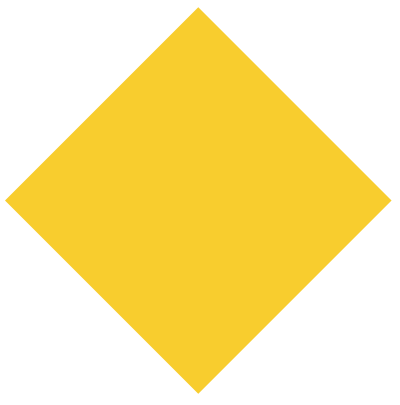
Use the following items:
- model.score

Example of Hidden Markov Model

- activity depending on the weather

# PRACTICAL EXERCISE

## 2. Solve scoring problem

Find the seqence of states for the following observations:
- (Walk)
- (Shop)
- (Clean)
-(Clean, Clean, Clean)

## Use the following items:
- model.decode

Example of Hidden Markov Model

- activity depending on the weather