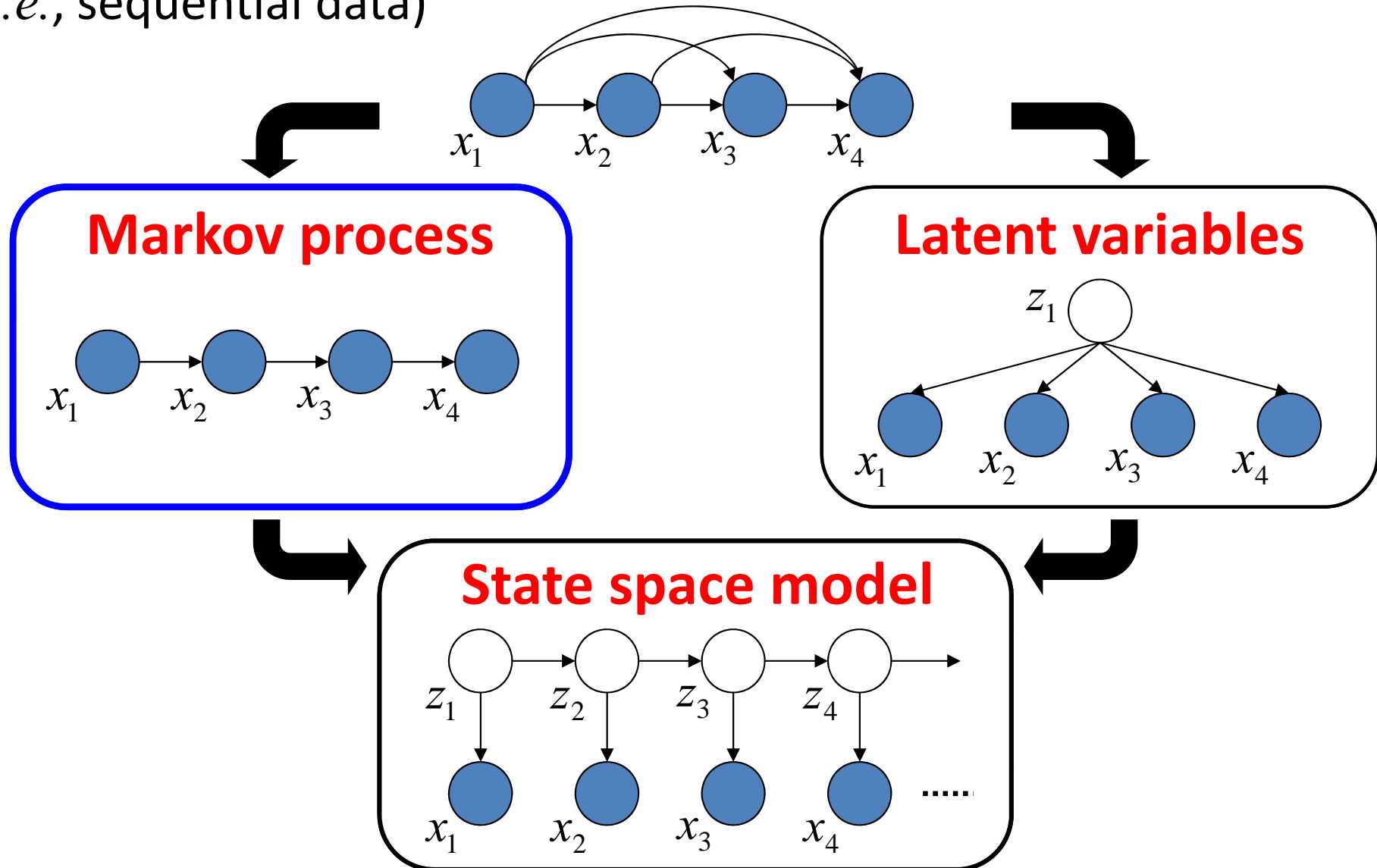


Sequential Data Modeling

Hidden Markov Models

Basic Approaches

How to efficiently model joint probability of high-dimensional data (*i.e.*, sequential data)



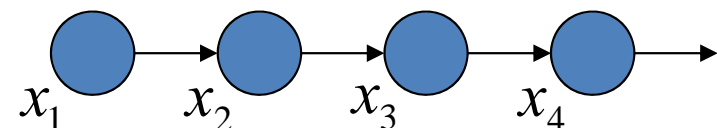
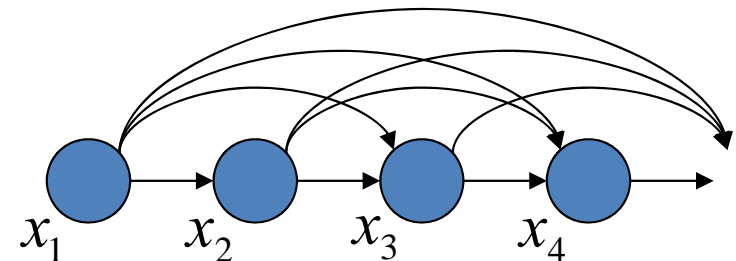
Review: Markov Process

- Assume that the conditional probability distribution of the present states depends only on a few past states

$$p(x_1, \dots, x_N) = p(x_1) \prod_{n=2}^N p(x_n | x_1, \dots, x_{n-1})$$

1st order Markov chain

$$p(x_1, \dots, x_N) = p(x_1) \prod_{n=2}^N p(x_n | x_{n-1})$$



Show probabilistic graphical model

- Nodes: random variables
- Edges: causality relationship between variables

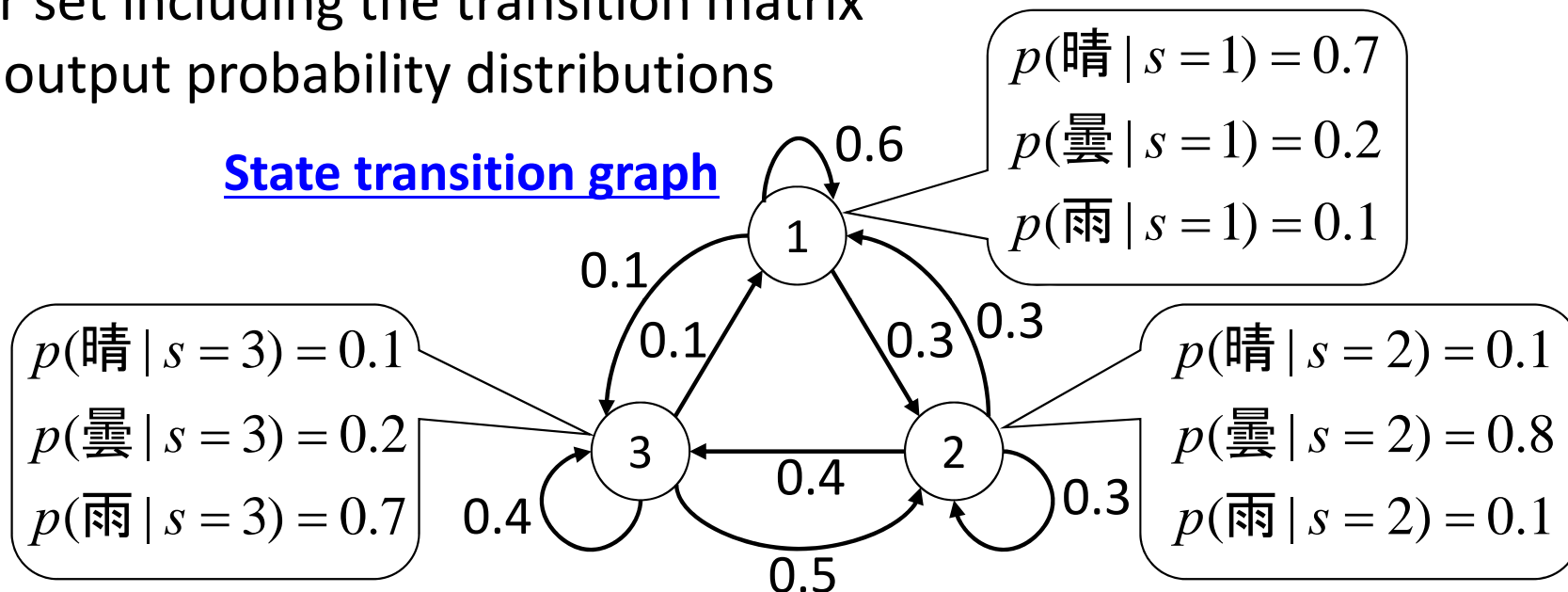
Easy to see conditional dependency between variables!

Hidden Markov Model (HMM)

- Use of **discrete latent variables**
- Parameter set including the transition matrix and state output probability distributions

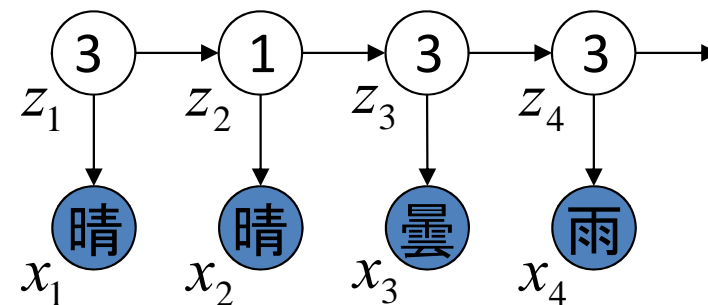
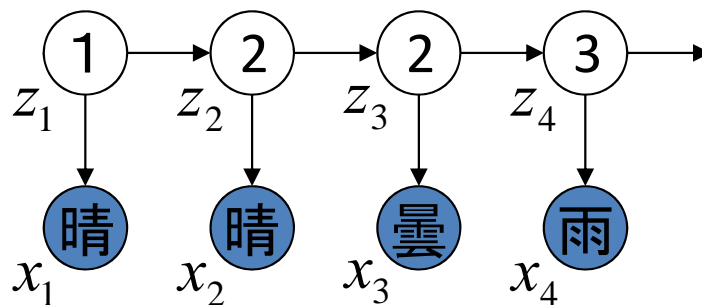
State transition graph

晴: fine
雨: rain
曇: cloud



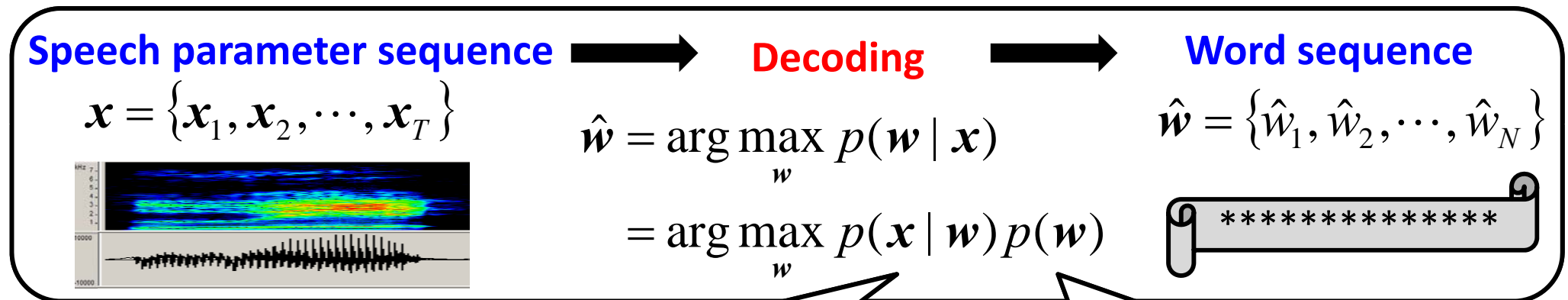
Even if 晴 晴 曇 雨 ... is observed, **a state sequence is not observed** (*i.e.*, **state sequence = latent variable**).

Possible state sequences...



Example of Application

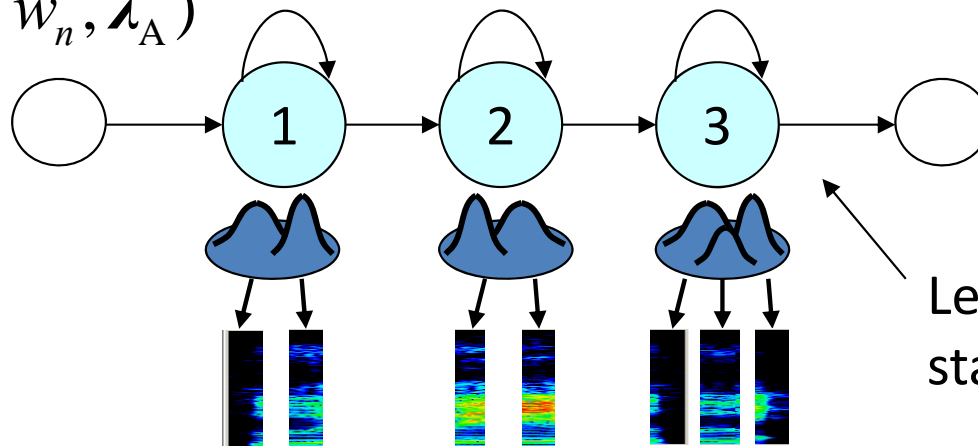
- Automatic speech recognition (*i.e.*, conversion from speech into text)



Acoustic model $p(\mathbf{x} | \mathbf{w}, \lambda_A)$
e.g., **HMM** with **GMM**

Language model $p(\mathbf{w} | \lambda_L)$
e.g., **Markov model**

$$p(\mathbf{x} | w_n, \lambda_A)$$

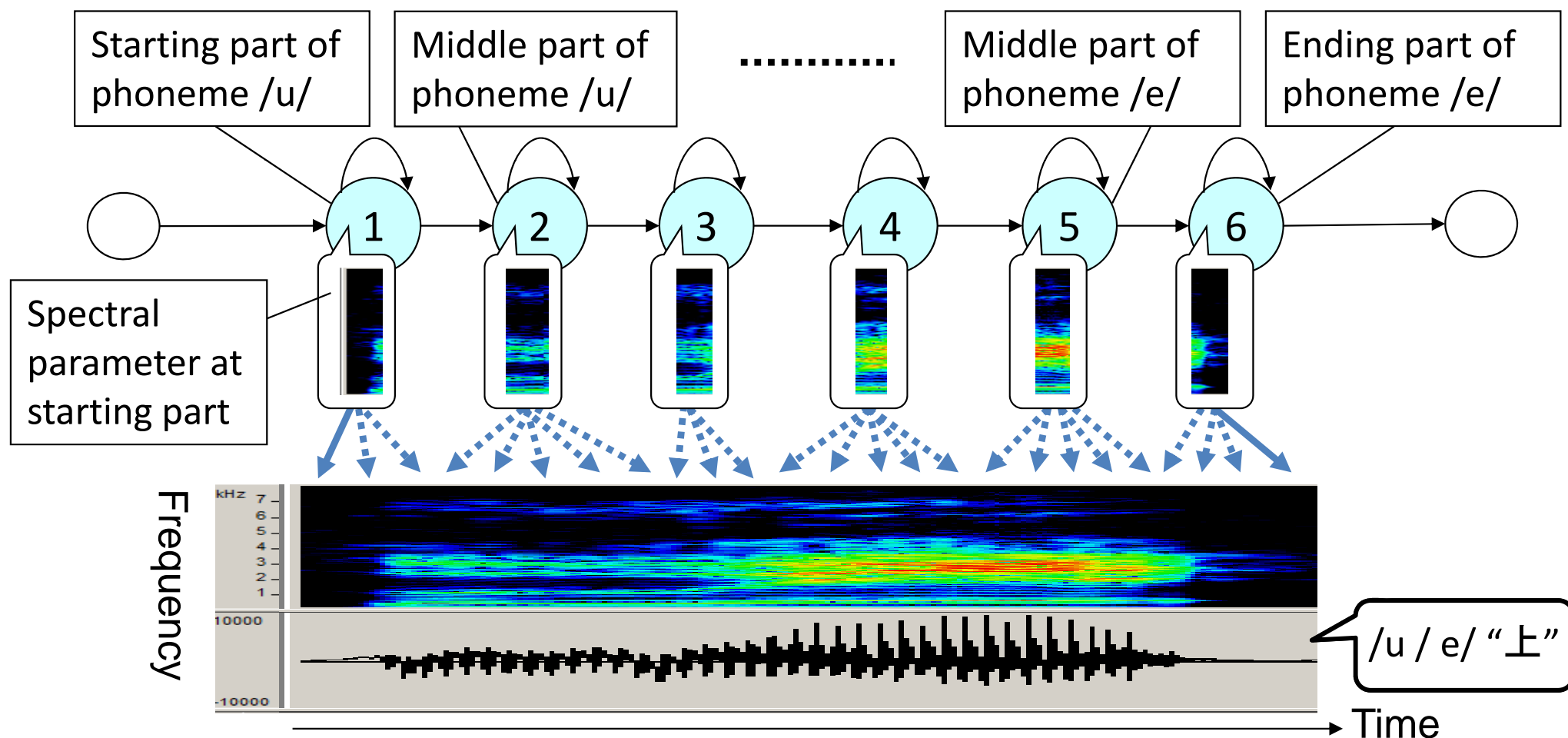


$$p(w_n | w_{n-1}, w_{n-2}, \lambda_L)$$

Left-to-right HMM using GMM as state output probability density

Example: Acoustic Model

- HMM effectively models a time sequence of speech parameters
 - ✓ Handle fluctuation of speaking speed with self-loop transition
 - ✓ Handle fluctuation of articulation with output probability density function



Speech waveform and its spectral sequence

HMM as Mixture Model

- Observation sequence of length N

$$\mathbf{x} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$$

- Corresponding state sequence of length N

$$\mathbf{z} = \{z_1, z_2, \dots, z_N\}$$

- Likelihood of HMM

$$p(\mathbf{x} | \lambda) = \sum_{\text{all } \mathbf{z}} p(\mathbf{x}_1, \dots, \mathbf{x}_N, z_1, \dots, z_N | \lambda)$$

$$= \sum_{\text{all } \mathbf{z}} \left\{ \underbrace{p(z_1 | \lambda) \left[\prod_{n=2}^N p(z_n | z_{n-1}, \lambda) \right]}_{\text{Prior probability of } \mathbf{z} : p(\mathbf{z})} \underbrace{\prod_{n=1}^N p(\mathbf{x}_n | z_n, \lambda)}_{\text{Output probability of } \mathbf{x} \text{ given } \mathbf{z} : p(\mathbf{x} | \mathbf{z})} \right\}$$

of mixture components : S^N

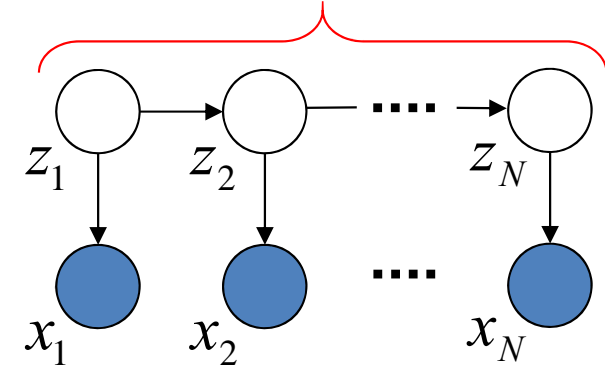
Prior probability of $\mathbf{z} : p(\mathbf{z})$

Output probability of \mathbf{x} given $\mathbf{z} : p(\mathbf{x} | \mathbf{z})$

$$= \sum_{\text{all } \mathbf{z}} p(\mathbf{z} | \lambda) p(\mathbf{x} | \mathbf{z}, \lambda)$$

← Mixture model with discrete latent variables \mathbf{z}

Markov chain of latent variables



x_4 depends on x_1 , x_2 , and x_3 .

Elements of HMM

- HMM parameter set $\lambda = \{\boldsymbol{\pi}, \boldsymbol{A}, \boldsymbol{B}(\cdot)\}$

- Set of S finite states: $s = \{1, 2, \dots, S\}$

- Initial state distribution

$$\boldsymbol{\pi} = \{\pi_1, \pi_2, \dots, \pi_S\}$$

$$\pi_i = p(s = i \mid \lambda)$$

- Transition probabilities

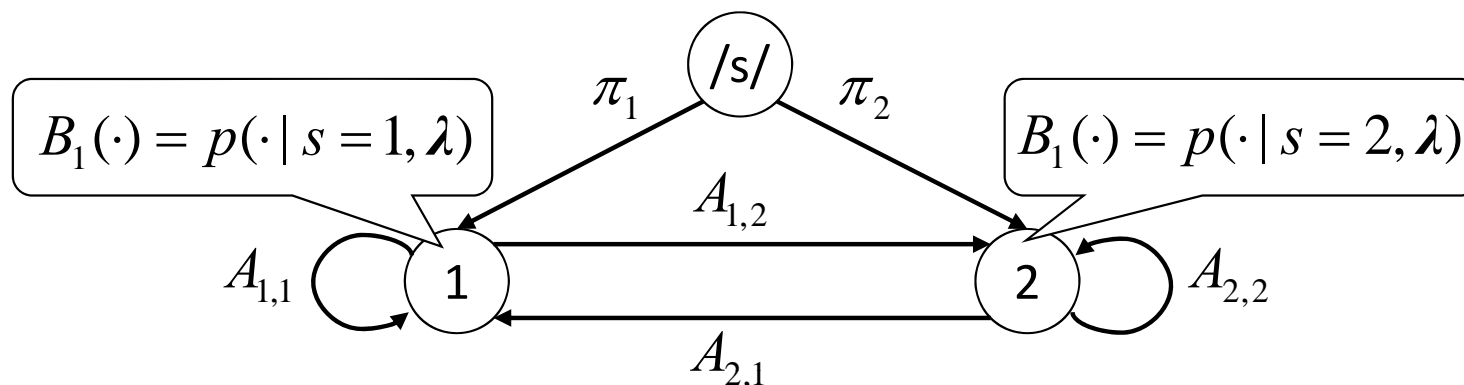
$$\boldsymbol{A} = \{A_{1,1}, A_{1,2}, \dots, A_{S,S}\}$$

$$A_{i,j} = p(s = j \mid s = i, \lambda)$$

- Output probabilities (or probability densities) of \mathbf{x}_n

$$\boldsymbol{B}(\mathbf{x}_n) = \{B_1(\mathbf{x}_n), B_2(\mathbf{x}_n), \dots, B_S(\mathbf{x}_n)\}$$

$$B_i(\mathbf{x}_n) = p(\mathbf{x}_n \mid s = i, \lambda)$$



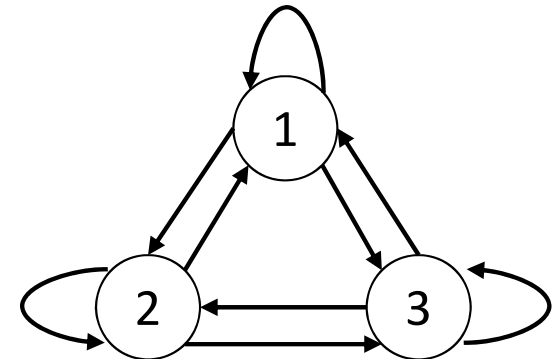
Type of HMM: Ergodic HMM

- Initial probabilities

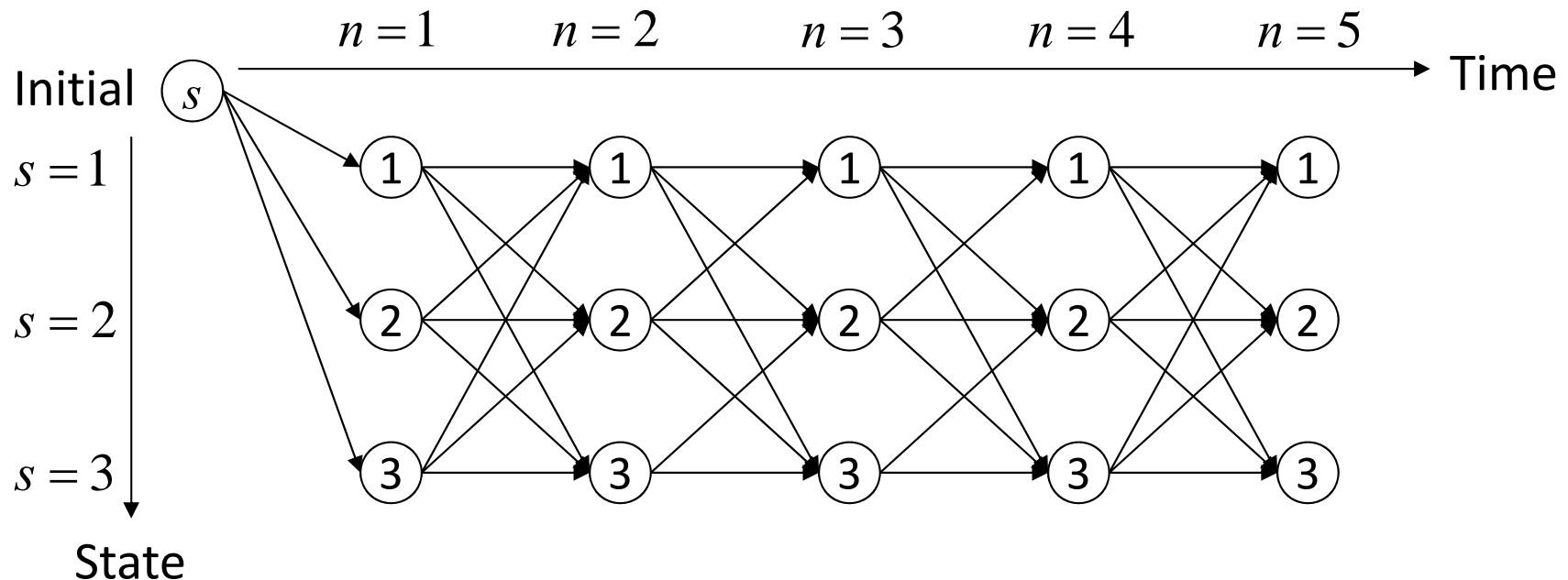
$$\pi_s = [\pi_1 \quad \pi_2 \quad \pi_3]$$

- Transition probabilities

$$A = \begin{bmatrix} A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,1} & A_{3,2} & A_{3,3} \end{bmatrix}$$



- Trellis graph

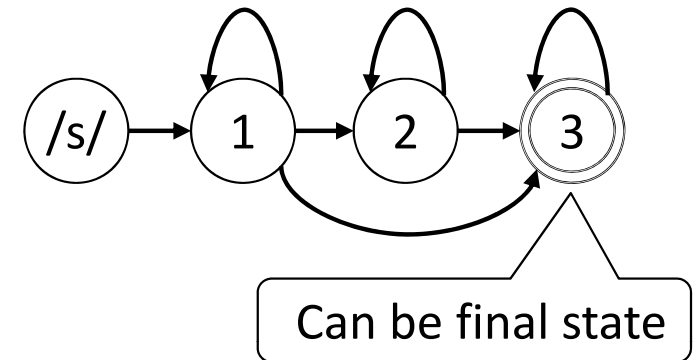


Type of HMM: Left-to-Right HMM

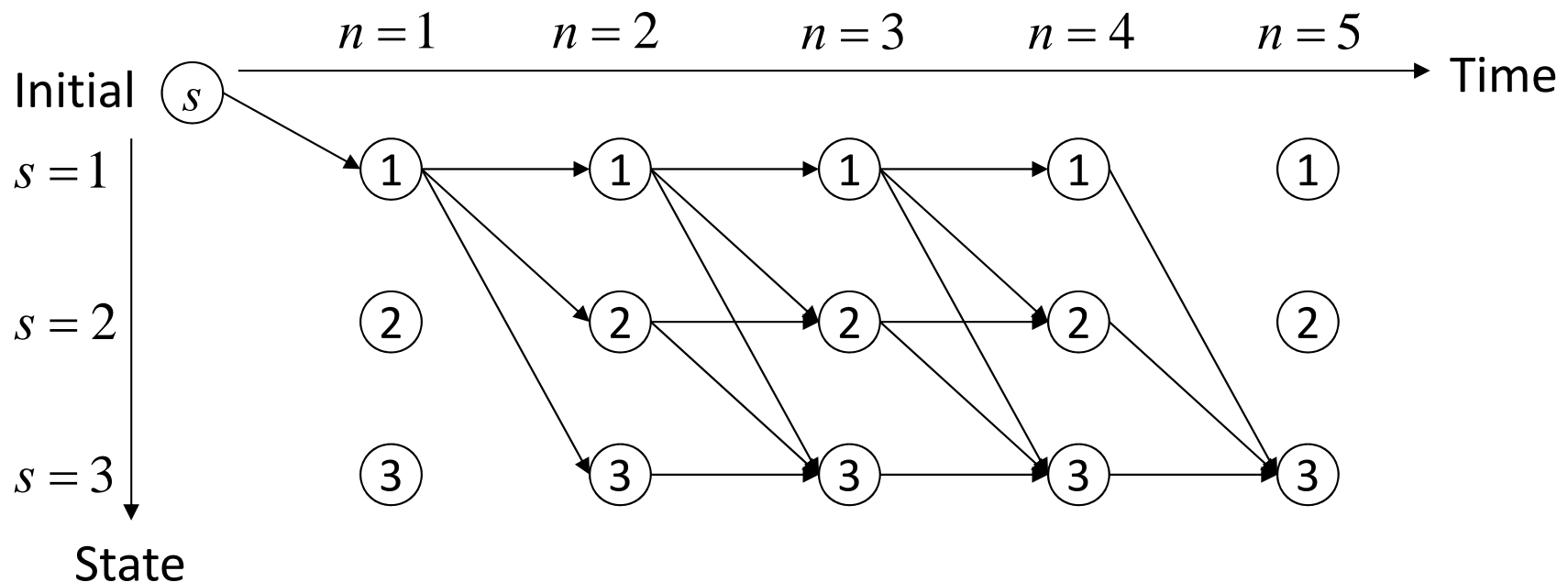
- Initial probabilities $\pi_s = \begin{cases} 1 & s = 1 \\ 0 & s \neq 1 \end{cases}$

- Transition probabilities $A = \begin{bmatrix} A_{1,1} & A_{1,2} & A_{1,3} \\ 0 & A_{2,2} & A_{2,3} \\ 0 & 0 & A_{3,3} \end{bmatrix}$

Suitable for speech model



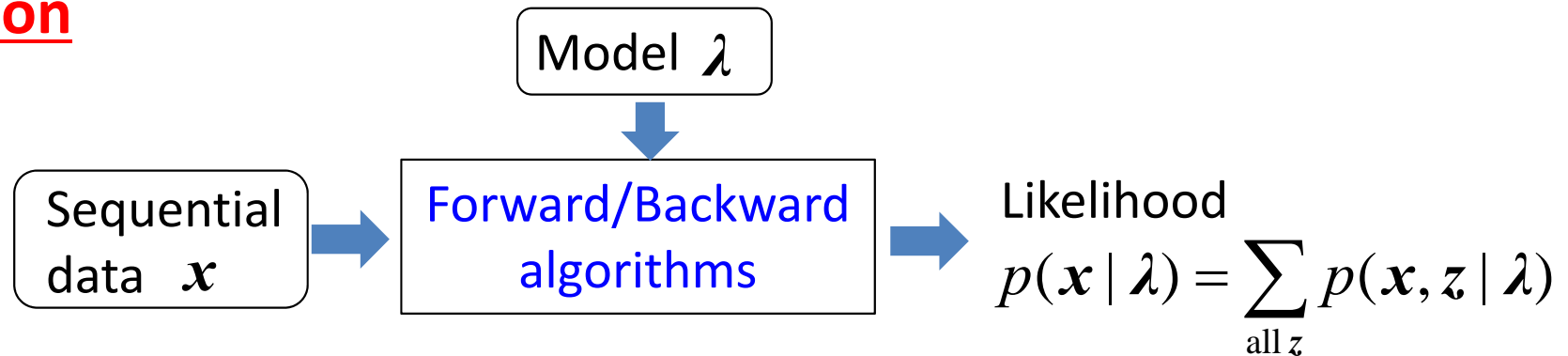
- Trellis graph



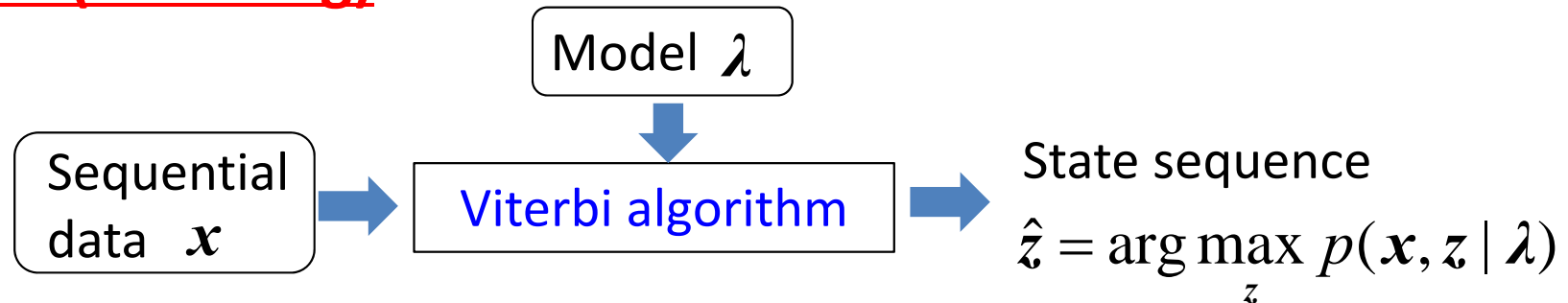
Strong constraints on temporal structure can be used!

Evaluation/Alignment/Training

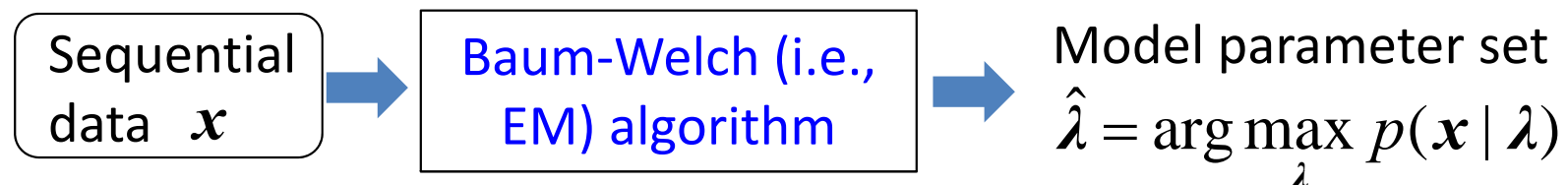
Evaluation



Alignment (Decoding)



Training

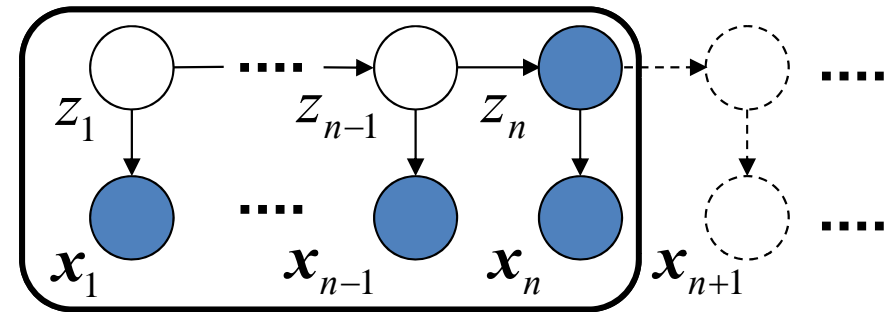


Forward Algorithm

- Recursively calculate forward probabilities
 - Forward probability that HMM generates $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ and ends in state s

$$\alpha_n(s) = p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n, z_n = s \mid \lambda)$$

Marginalizing out
all possible previous states!



- Recursive calculation

➤ Initialization ($n = 1; 1 \leq s \leq S$)

$$\alpha_1(s) = \pi_s B_s(\mathbf{x}_1)$$

➤ Recursion ($2 \leq n \leq N; 1 \leq s \leq S$)

$$\alpha_n(s) = \left[\sum_{s'=1}^S \alpha_{n-1}(s') A_{s',s} \right] B_s(\mathbf{x}_n)$$

➤ Termination

$$p(\mathbf{x}_1, \dots, \mathbf{x}_N \mid \lambda) = \sum_{s=1}^S \alpha_N(s)$$

$$\begin{aligned} \alpha(z_n) &= p(\mathbf{x}_1, \dots, \mathbf{x}_n, z_n) \\ &= p(\mathbf{x}_1, \dots, \mathbf{x}_n | z_n) p(z_n) \\ &= p(\mathbf{x}_n | z_n) p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1} | z_n) p(z_n) \\ &= p(\mathbf{x}_n | z_n) p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1}, z_n) \\ &= p(\mathbf{x}_n | z_n) \sum_{z_{n-1}} p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1}, z_{n-1}, z_n) \\ &= p(\mathbf{x}_n | z_n) \sum_{z_{n-1}} p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1}, z_n | z_{n-1}) p(z_{n-1}) \\ &= p(\mathbf{x}_n | z_n) \sum_{z_{n-1}} p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1} | z_{n-1}) p(z_n | z_{n-1}) p(z_{n-1}) \\ &= p(\mathbf{x}_n | z_n) \sum_{z_{n-1}} p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1}, z_{n-1}) p(z_n | z_{n-1}) \end{aligned}$$

By use of the definition (13.34) for $\alpha(z_n)$, we then obtain

$$\alpha(z_n) = p(\mathbf{x}_n | z_n) \sum_{z_{n-1}} \alpha(z_{n-1}) p(z_n | z_{n-1}).$$

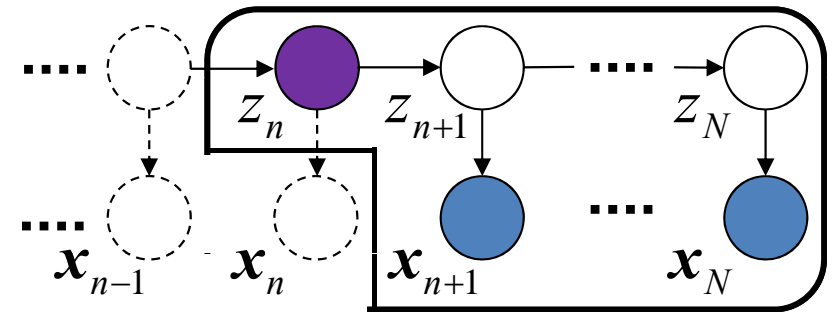
Backward Algorithm

- Recursively calculate backward probabilities

- Backward probability** that HMM generates $\{\mathbf{x}_{n+1}, \dots, \mathbf{x}_N\}$ when starting in state s at time n

$$\beta_n(s) = p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N \mid z_n = s, \lambda)$$

Marginalizing out all possible succeeding states!



- Recursive calculation**

- Initialization ($n = N; 1 \leq s \leq S$)

$$\beta_N(s) = 1 \quad p(\mathbf{z}_N \mid \mathbf{X}) = \frac{p(\mathbf{X}, \mathbf{z}_N) \beta(\mathbf{z}_N)}{p(\mathbf{X})}$$

- Recursion ($1 \leq n \leq N-1; 1 \leq s \leq S$)

$$\beta_n(s) = \sum_{s'=1}^S A_{s,s'} B_{s'}(\mathbf{x}_{n+1}) \beta_{n+1}(s')$$

- Termination

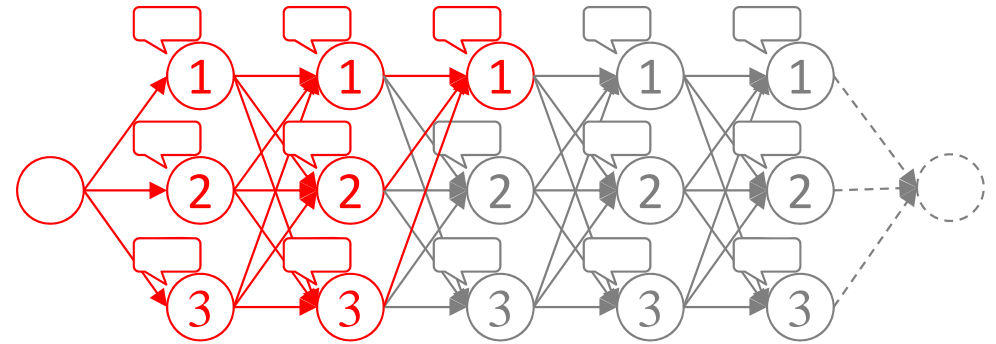
$$p(\mathbf{x}_1, \dots, \mathbf{x}_N \mid \lambda) = \sum_{s=1}^S \pi_s B_s(\mathbf{x}_1) \beta_1(s)$$

$$\begin{aligned} \beta(z_n) &= p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N \mid z_n) \\ &= \sum_{z_{n+1}} p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N, z_{n+1} \mid z_n) \\ &= \sum_{z_{n+1}} p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N \mid z_n, z_{n+1}) p(z_{n+1} \mid z_n) \\ &= \sum_{z_{n+1}} p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N \mid z_{n+1}) p(z_{n+1} \mid z_n) \\ &= \sum_{z_{n+1}} p(\mathbf{x}_{n+2}, \dots, \mathbf{x}_N \mid z_{n+1}) p(\mathbf{x}_{n+1} \mid z_{n+1}) p(z_{n+1} \mid z_n) \\ &\downarrow \\ \beta(z_n) &= \sum_{z_{n+1}} \beta(z_{n+1}) p(\mathbf{x}_{n+1} \mid z_{n+1}) p(z_{n+1} \mid z_n). \end{aligned}$$

Product of Forward/Backward Probs

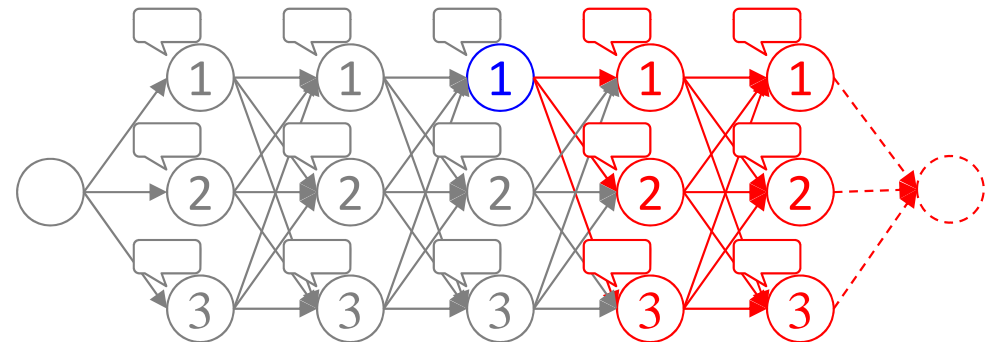
- Forward probability

$$\alpha_n(s) = p(\mathbf{x}_1, \dots, \mathbf{x}_n, z_n = s \mid \lambda)$$



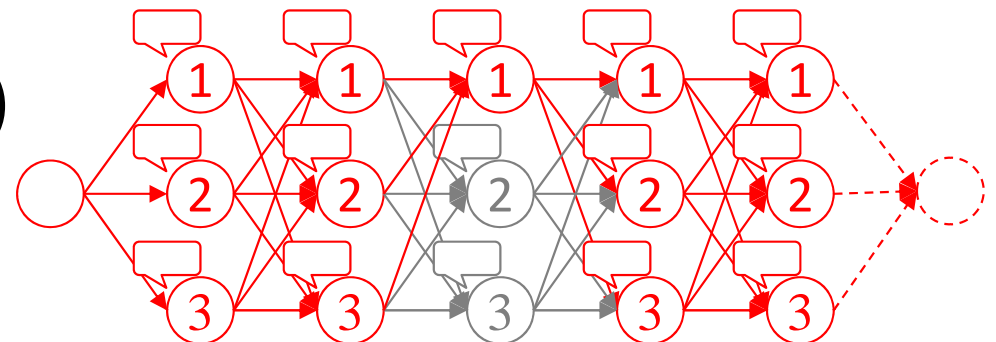
- Backward probability

$$\beta_n(s) = p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N \mid z_n = s, \lambda)$$



- Their product

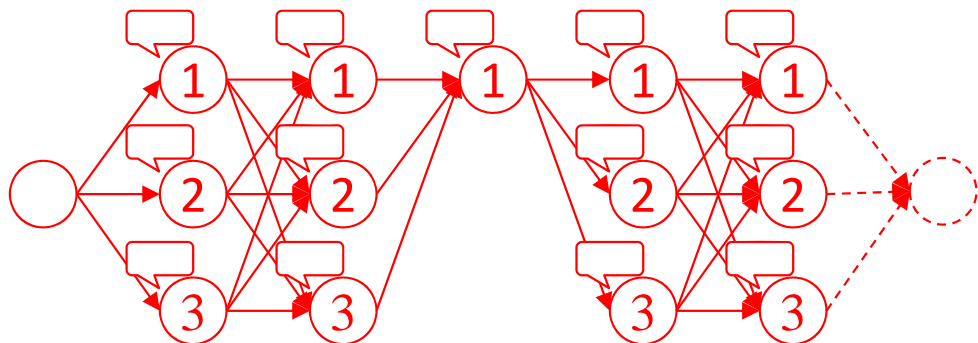
$$\alpha_n(s)\beta_n(s) = p(\mathbf{x}_1, \dots, \mathbf{x}_N, z_n = s \mid \lambda)$$



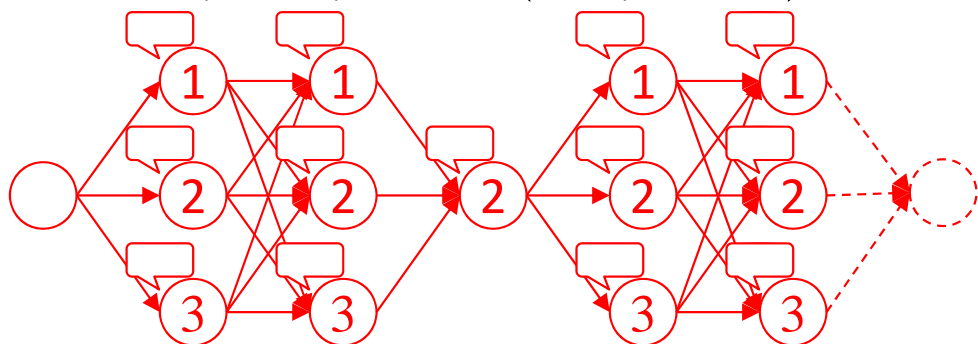
Considering all possible paths
passing through state s at time n

Likelihood from Forward/Backward Probs

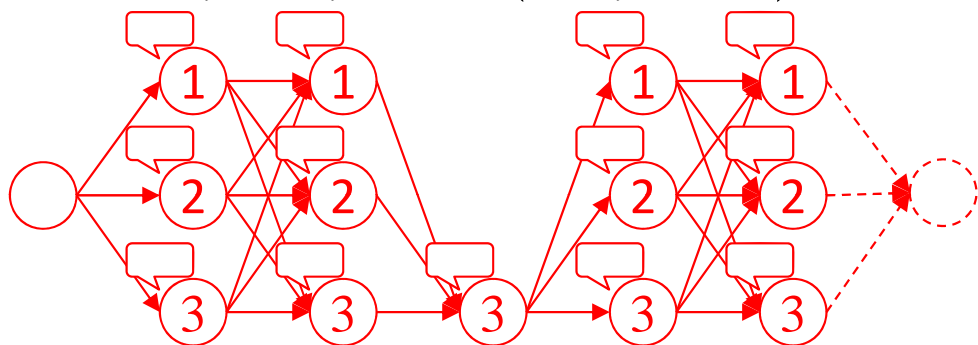
$$\alpha_n(1)\beta_n(1) = p(\mathbf{x}, z_n = 1 | \lambda)$$



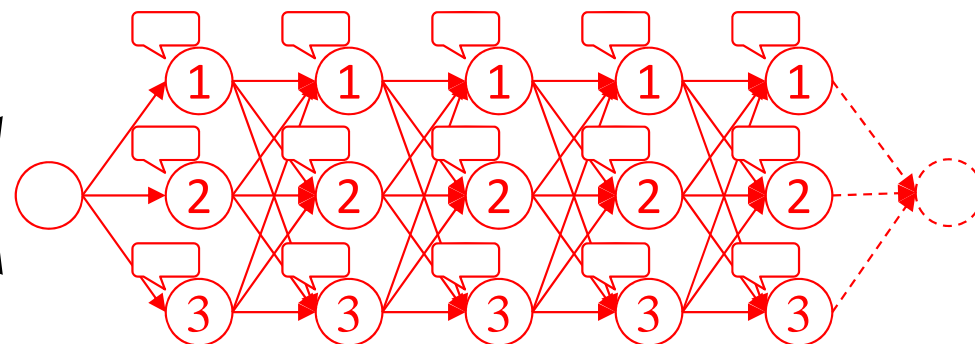
$$\alpha_n(2)\beta_n(2) = p(\mathbf{x}, z_n = 2 | \lambda)$$



$$\alpha_n(3)\beta_n(3) = p(\mathbf{x}, z_n = 3 | \lambda)$$



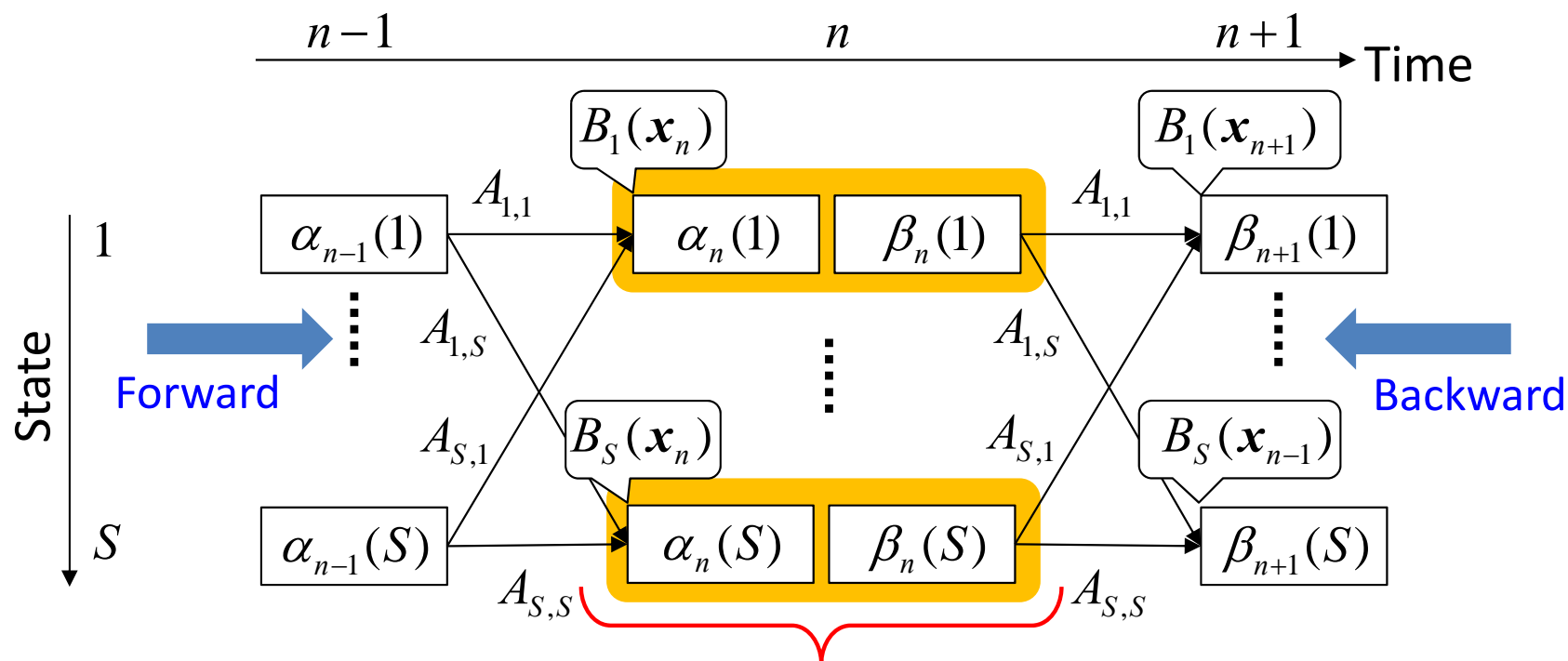
$$\sum_{s=1}^S \alpha_n(s)\beta_n(s) = \sum_{s=1}^S p(\mathbf{x}, z_n = s | \lambda) = p(\mathbf{x} | \lambda)$$



Sum of the products of forward/backward probabilities at time n is equal to likelihood over all possible paths!

Likelihood Calculation

- Forward probability : $\alpha_n(s) = p(\mathbf{x}_1, \dots, \mathbf{x}_n, z_n = s \mid \lambda)$
- Backward probability : $\beta_n(s) = p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N \mid z_n = s, \lambda)$
- Their product: $\alpha_n(s)\beta_n(s) = p(\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{x}_{n+1}, \dots, \mathbf{x}_N, z_n = s \mid \lambda)$



Sum of their products at time n = Likelihood

$$\sum_{s=1}^S \alpha_n(s) \beta_n(s) = p(\mathbf{x}_1, \dots, \mathbf{x}_N \mid \lambda)$$

App.: Derivation of Forward Algorithm

$$\begin{aligned}
 p(\mathbf{x}_1, \dots, \mathbf{x}_N | \lambda) &= \sum_{\text{all } \mathbf{z}} \left\{ p(z_1 | \lambda) \left[\prod_{n=2}^N p(z_n | z_{n-1}, \lambda) \right] \prod_{n=1}^N p(\mathbf{x}_n | z_n, \lambda) \right\} \\
 &= \sum_{s_N=1}^S \left(B_{s_N}(\mathbf{x}_N) \underbrace{\sum_{\text{all } \mathbf{z} \setminus z_N} \left\{ p(z_N = s_N | z_{N-1}, \lambda) p(z_1 | \lambda) \left[\prod_{n=2}^{N-1} p(z_n | z_{n-1}, \lambda) \right] \prod_{n=1}^{N-1} p(\mathbf{x}_n | z_n, \lambda) \right\}}_{\substack{\text{Factor out the common factor at } n=N \text{ in each state } s=s_N \\ \text{All possible state sequences from } n=1 \text{ to } n=N-1}} \right) \\
 &= \sum_{s_N=1}^S \left(B_{s_N}(\mathbf{x}_N) \sum_{s_{N-1}=1}^S A_{s_{N-1}, s_N} \left(B_{s_{N-1}}(\mathbf{x}_{N-1}) \sum_{s_{N-2}=1}^S A_{s_{N-2}, s_{N-1}} \cdots \left(B_{s_2}(\mathbf{x}_2) \sum_{s_1=1}^S A_{s_1, s_2} \underbrace{\left(\pi_{s_1} B_{s_1}(\mathbf{x}_1) \right)}_{\alpha_1(s_1)} \right) \right) \right) \\
 &\quad \underbrace{\qquad \qquad \qquad}_{\alpha_2(s_2)} \qquad \qquad \qquad \underbrace{\qquad \qquad \qquad}_{\alpha_{N-1}(s_{N-1})} \\
 &= \sum_{s=1}^S \underbrace{\qquad \qquad \qquad}_{\alpha_N(s)}
 \end{aligned}$$

App.: Derivation of Backward Algorithm

$$\begin{aligned}
 p(\mathbf{x}_1, \dots, \mathbf{x}_N | \lambda) &= \sum_{\text{all } \mathbf{z}} \left\{ p(z_1 | \lambda) \left[\prod_{n=2}^N p(z_n | z_{n-1}, \lambda) \right] \prod_{n=1}^N p(\mathbf{x}_n | z_n, \lambda) \right\} \\
 &\quad \text{Factor out the common factor at } n=1 \text{ in each state } s=s_1 \\
 &= \sum_{s_1=1}^S \pi_{s_1} B_{s_1}(\mathbf{x}_1) \sum_{\text{all } \mathbf{z} \setminus z_1} \left\{ p(z_2 | z_1 = s, \lambda) \left[\prod_{n=3}^N p(z_n | z_{n-1}, \lambda) \right] \prod_{n=2}^N p(\mathbf{x}_n | z_n, \lambda) \right\} \\
 &\quad \text{All possible state sequences from } n=2 \text{ to } n=N \\
 &\quad \text{Iteratively factor out in the same manner from } n=2 \text{ to } n=N-1 \\
 &= \sum_{s_1=1}^S \pi_{s_1} B_{s_1}(\mathbf{x}_1) \left(\sum_{s_2=1}^S A_{s_1, s_2} B_{s_2}(\mathbf{x}_2) \cdots \left(\sum_{s_{N-1}=1}^S A_{s_{N-2}, s_{N-1}} B_{s_{N-1}}(\mathbf{x}_{N-1}) \left(\sum_{s_N=1}^S A_{s_{N-1}, s_N} B_{s_N}(\mathbf{x}_N) \right) \right) \right) \\
 &\quad \underbrace{\hspace{15em}}_{\beta_{N-1}(s_{N-1})} \\
 &\quad \underbrace{\hspace{10em}}_{\beta_{N-2}(s_{N-2})} \\
 &\quad \underbrace{\hspace{5em}}_{\beta_1(s_1)} \\
 &= \sum_{s=1}^S \pi_s B_s(\mathbf{x}_1) \beta_1(s)
 \end{aligned}$$

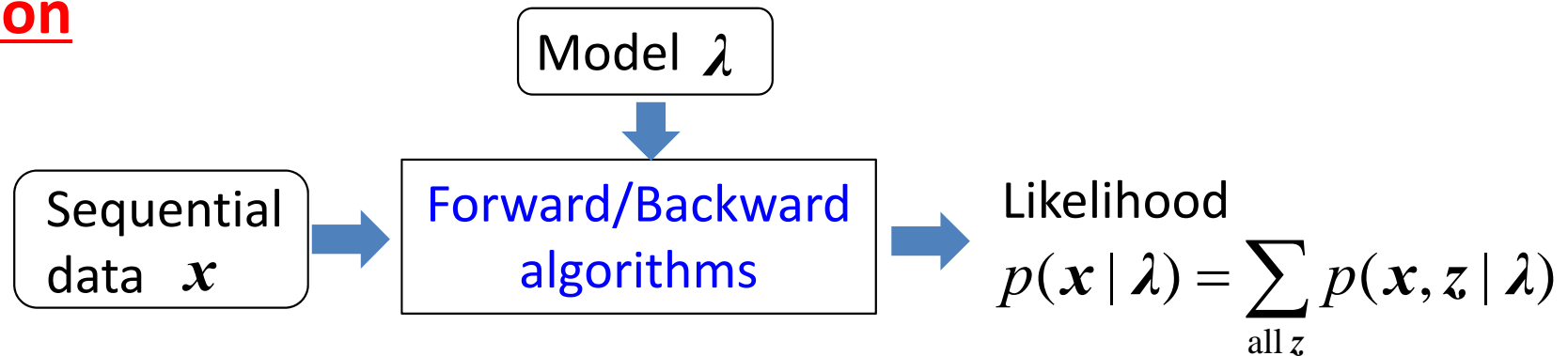
* Note that $\beta_N(s_N)=1$

Sequential Data Modeling

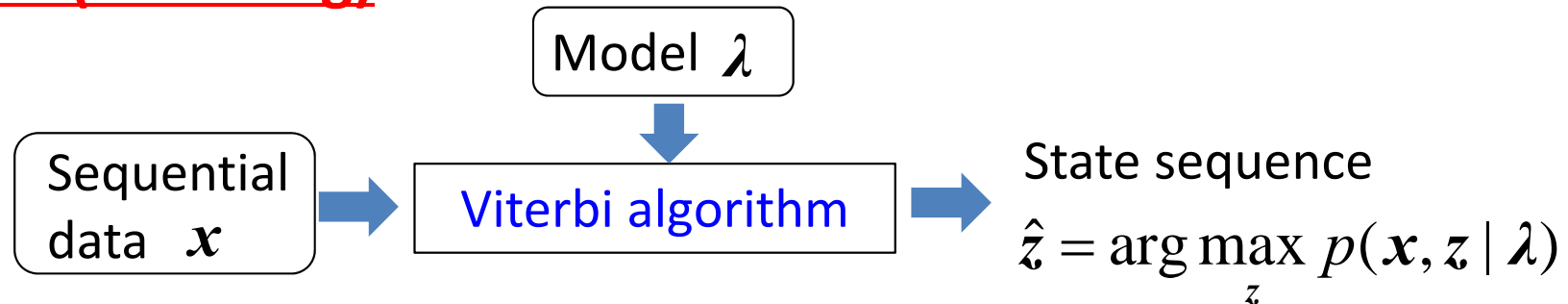
“Discrete Latent Variable Models 1”

Evaluation/Alignment/Training

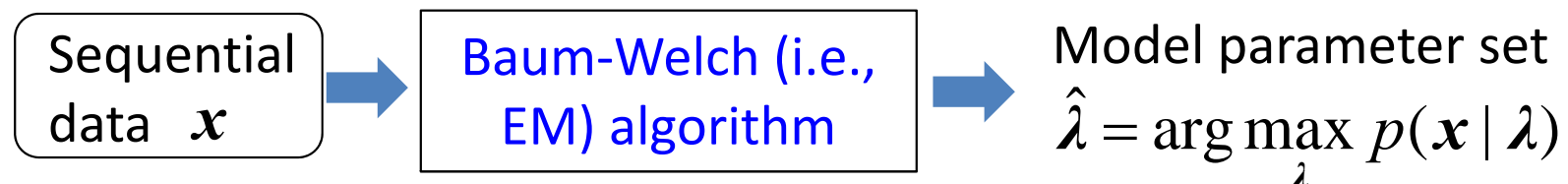
Evaluation



Alignment (Decoding)

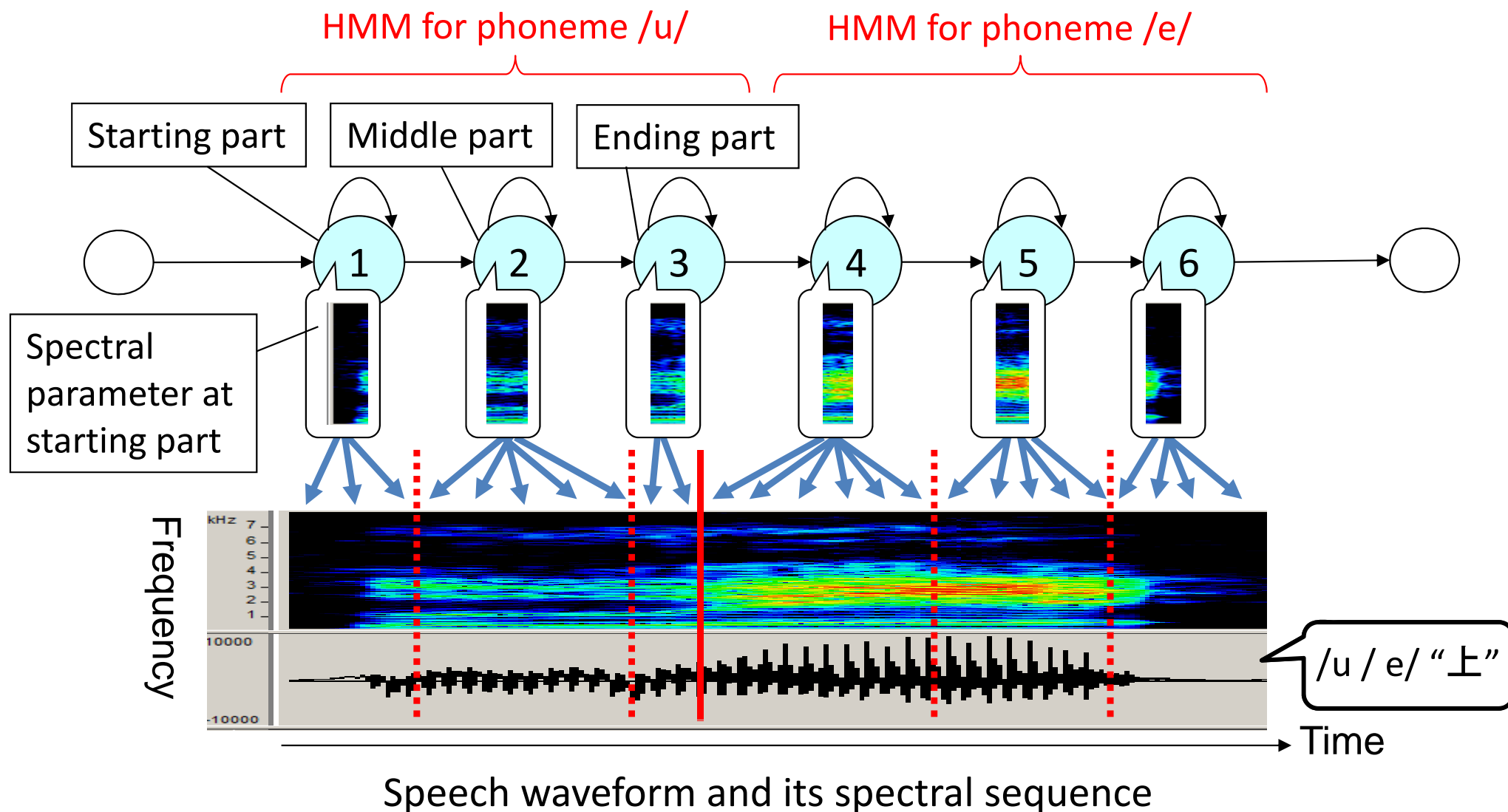


Training



Alignment

- Find the best state sequence for a given sequential data sample
e.g., find phoneme boundaries over a speech signal



Viterbi Algorithm

- Recursively calculate the best path probability
 - The most likely path probability that generates $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ and ends in state s

$$\chi_n(s) = p(\mathbf{x}_1, \dots, \mathbf{x}_n, \hat{z}_1, \dots, \hat{z}_{n-1}, \hat{z}_n = s \mid \lambda)$$

- Recursive calculation

➤ **Initialization** ($n = 1; 1 \leq s \leq S$)

$$\chi_1(s) = \pi_s B_s(\mathbf{x}_1) \quad \phi_1(s) = s$$

➤ **Recursion** ($2 \leq n \leq N; 1 \leq s \leq S$)

$$\chi_n(s) = B_s(\mathbf{x}_n) \max_{s'} \chi_{n-1}(s') A_{s',s}$$

$$\phi_n(s) = \arg \max_{s'} \chi_{n-1}(s') A_{s',s}$$

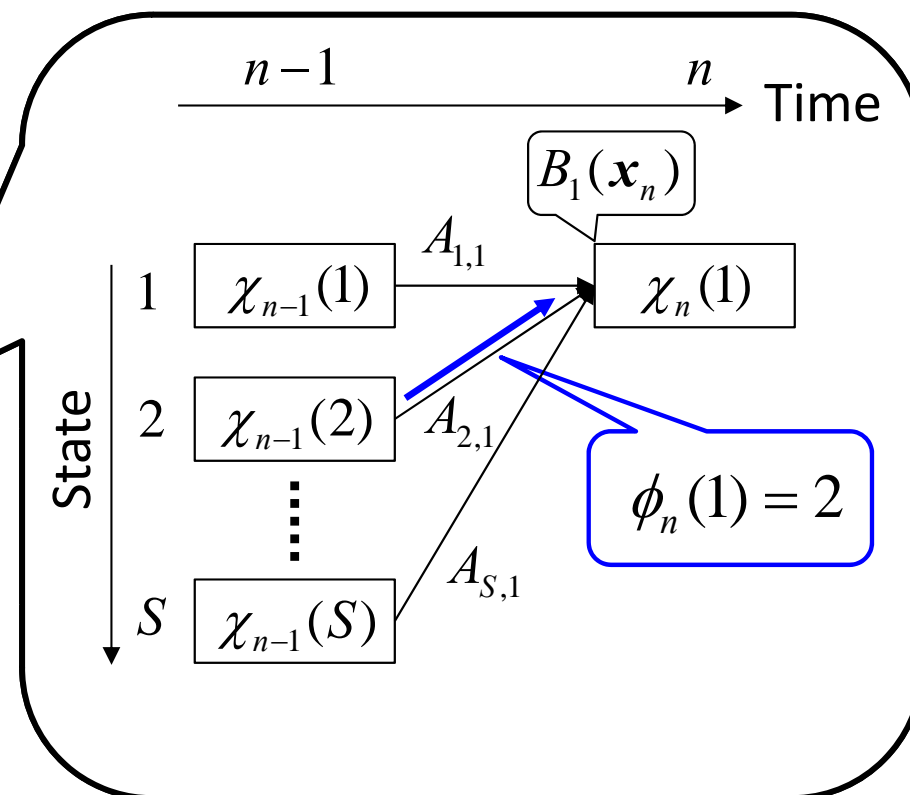
➤ **Termination**

$$p(\mathbf{x}, \hat{\mathbf{z}} \mid \lambda) = \max_s \chi_N(s)$$

$$\hat{z}_N = \arg \max_s \chi_N(s)$$

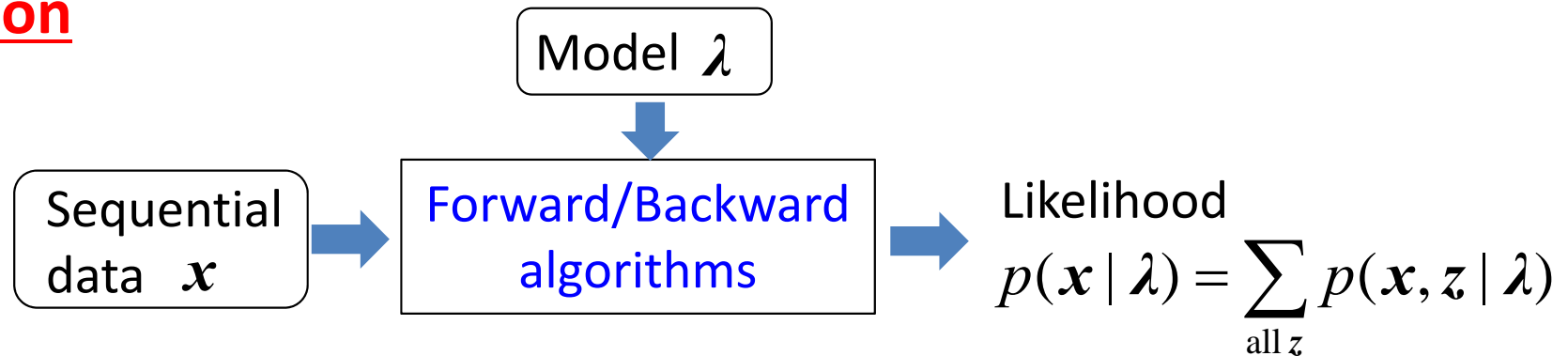
➤ **Path backtracking** ($2 \leq n \leq N$)

$$\hat{z}_{n-1} = \phi_n(\hat{z}_n)$$

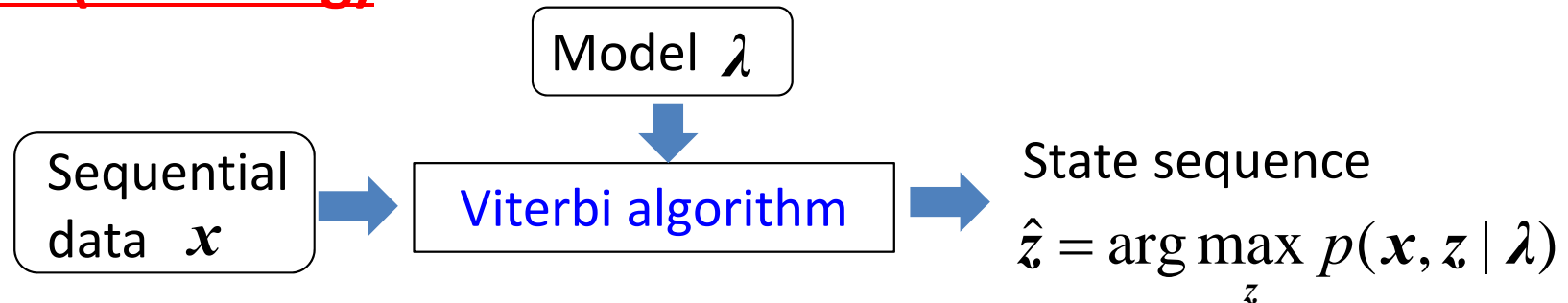


Evaluation/Alignment/Training

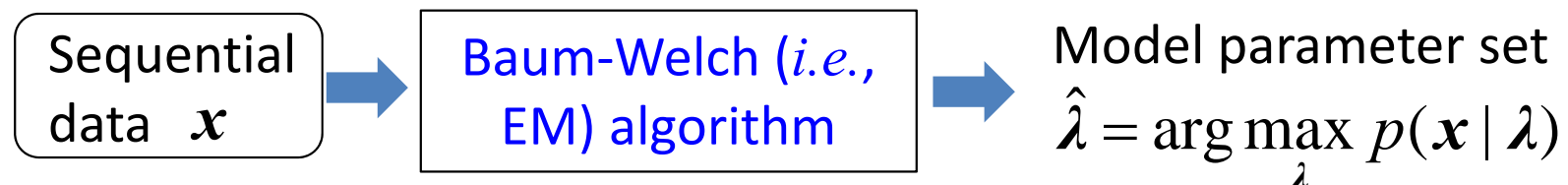
Evaluation



Alignment (Decoding)



Training

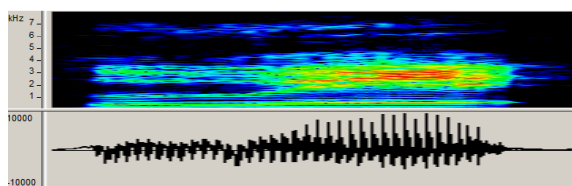


Training

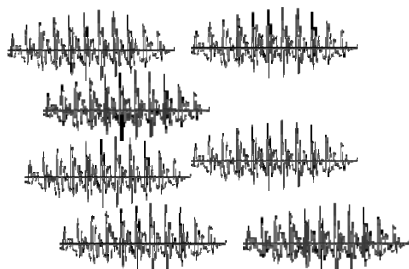
- Optimize HMM parameters so that the HMM appropriately models given observation sequences (*e.g.*, training the HMM from data)



“Keyword”

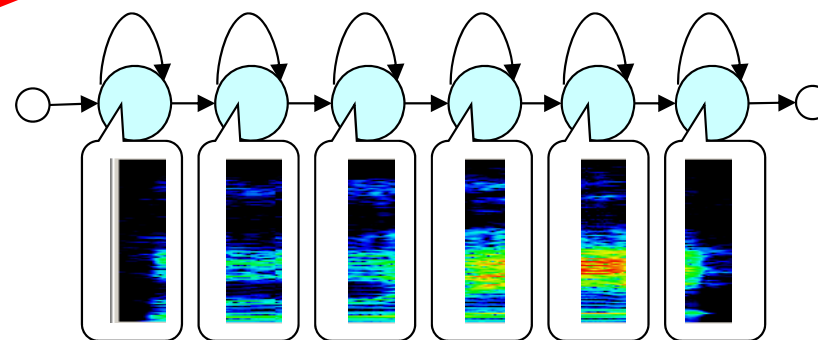


Speech waveform and
its spectral sequence



Training data samples

Optimize HMM parameters
by maximizing likelihood for
training data!



HMM to be optimized

Training with EM Algorithm

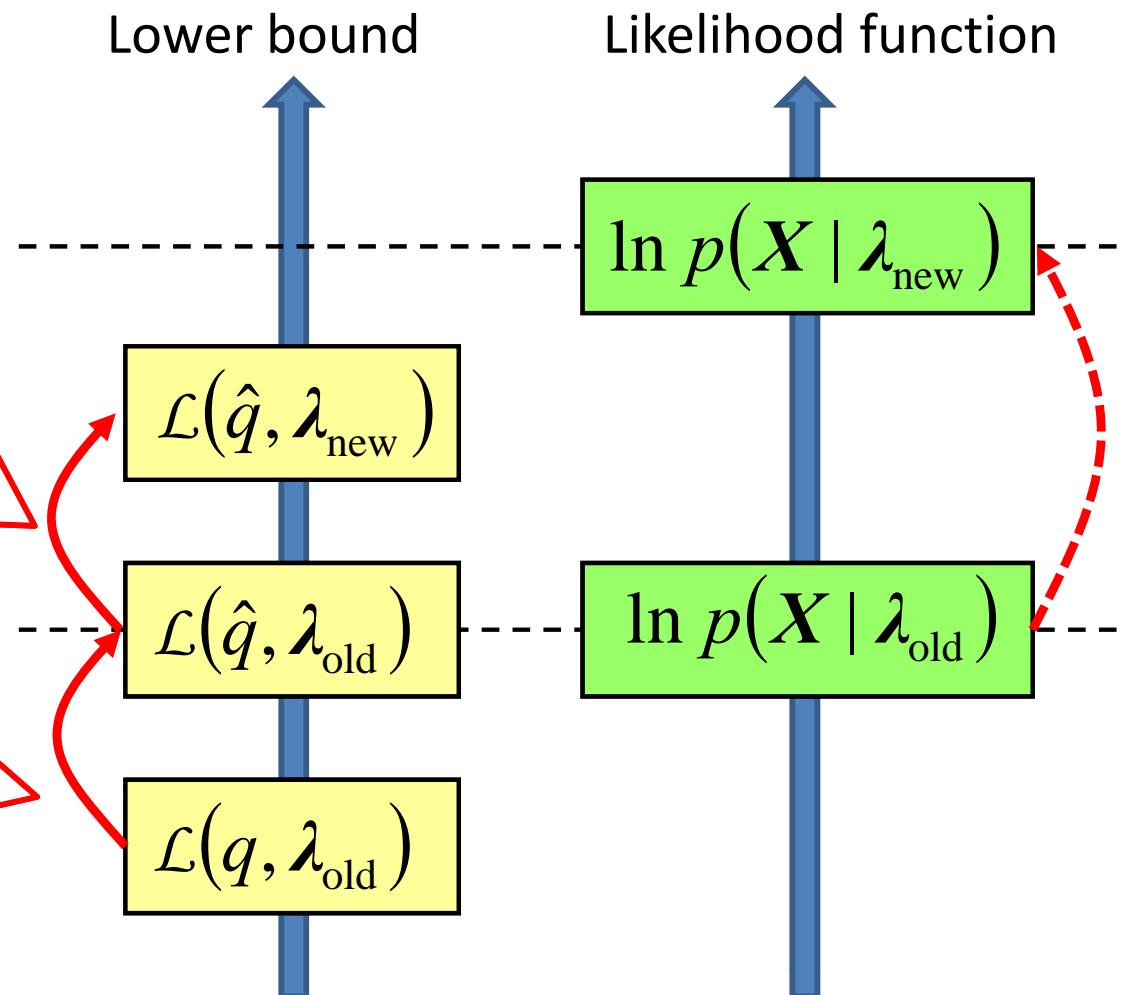
- Iteratively update lower bound of likelihood function through two steps: **E-step** and **M-step**

M-step: maximize lower bound with respect to λ while fixing q

$$\hat{\lambda} = \arg \max_{\lambda} \sum_{u=1}^U \sum_{\text{all } \mathbf{z}^{(u)}} \left\{ \hat{q}(\mathbf{z}^{(u)}) \cdot \ln p(\mathbf{x}^{(u)}, \mathbf{z}^{(u)} | \lambda) \right\}$$

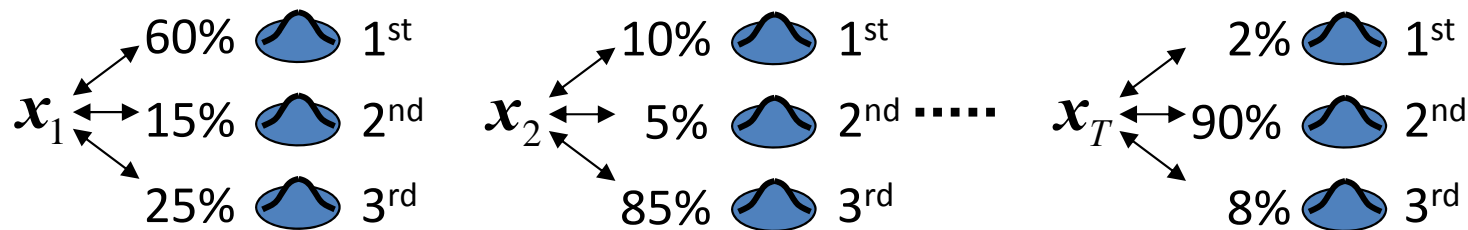
E-step: maximize lower bound with respect to q while fixing λ

$$\hat{q}(\mathbf{z}^{(u)}) = p(\mathbf{z}^{(u)} | \mathbf{x}^{(u)}, \lambda)$$

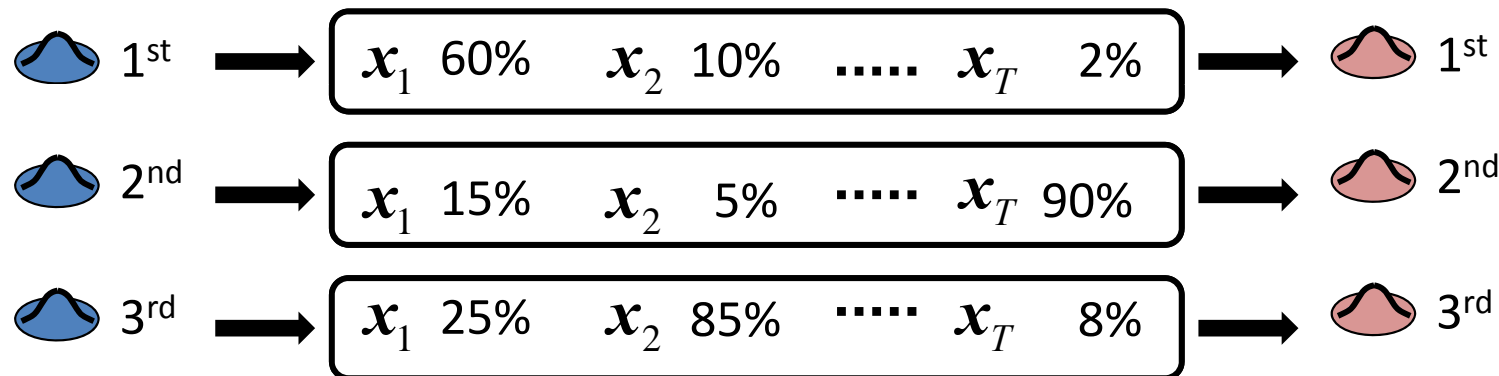


Review: Why Do We Use EM Algorithm?

- In a **mixture model**, we don't know which data can be used to update parameters of each mixture component due to **latent variables**.
- EM algorithm** is capable of addressing this issue!
 - E-step**: assign each data to individual mixture components (*i.e.*, estimate latent variables) based on current model parameters



- M-step**: update parameters of each mixture component using the assigned data



Lower Bound of HMM Likelihood

Log-scaled likelihood function for U samples of sequential data

$$\begin{aligned}\underline{\ln p(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(U)} \mid \boldsymbol{\lambda})} &= \sum_{u=1}^U \ln \sum_{\text{all } \mathbf{z}^{(u)}} p(\mathbf{x}^{(u)}, \mathbf{z}^{(u)} \mid \boldsymbol{\lambda}) \\ &\geq \sum_{u=1}^U \sum_{\text{all } \mathbf{z}^{(u)}} q(\mathbf{z}^{(u)}) \ln \frac{p(\mathbf{x}^{(u)}, \mathbf{z}^{(u)} \mid \boldsymbol{\lambda})}{q(\mathbf{z}^{(u)})} = \underline{\mathcal{L}(q, \boldsymbol{\lambda})}\end{aligned}$$

Lower bound

E-step: calculate posterior probabilities of latent variables (*i.e.*, state sequences)

$$\hat{q}(\mathbf{z}^{(u)}) = p(\mathbf{z}^{(u)} \mid \mathbf{x}^{(u)}, \boldsymbol{\lambda}_{\text{old}}) = \frac{p(\mathbf{x}^{(u)}, \mathbf{z}^{(u)} \mid \boldsymbol{\lambda}_{\text{old}})}{\sum_{\text{all } \mathbf{z}^{(u)}} p(\mathbf{x}^{(u)}, \mathbf{z}^{(u)} \mid \boldsymbol{\lambda}_{\text{old}})}$$

M-step: maximize auxiliary function with respect to model parameters

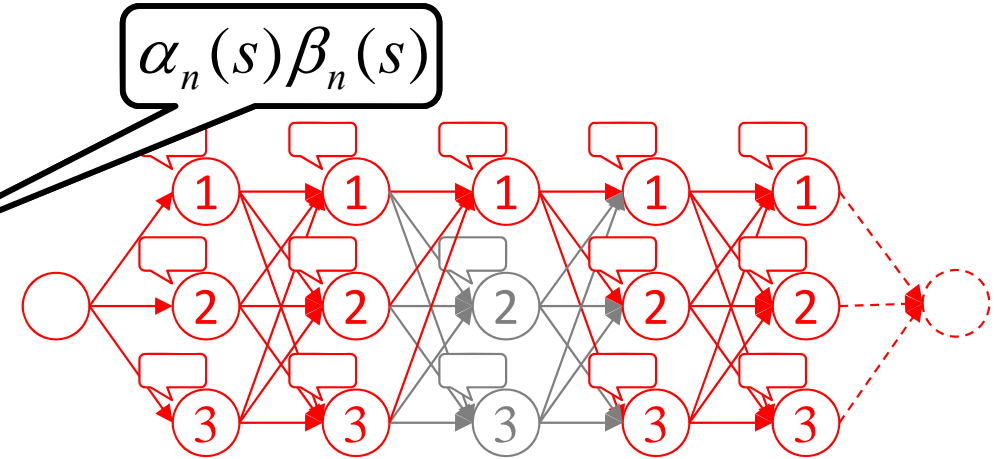
$$Q(\boldsymbol{\lambda}_{\text{new}}, \boldsymbol{\lambda}_{\text{old}}) = \sum_{u=1}^U \sum_{\text{all } \mathbf{z}^{(u)}} \hat{q}(\mathbf{z}^{(u)}) \ln p(\mathbf{x}^{(u)}, \mathbf{z}^{(u)} \mid \boldsymbol{\lambda}_{\text{new}})$$

E-Step

- Calculate posterior probabilities of latent variables

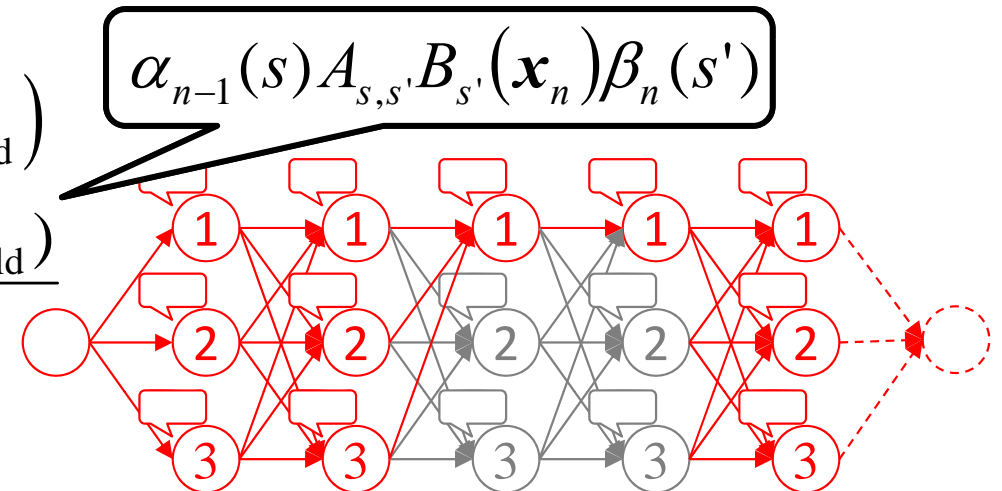
Expected # of samples observed in state s at time n in sample u

$$\begin{aligned}\gamma_s^{(u)}(n) &= \hat{q}(z_n^{(u)} = s) \\ &= p(z_n^{(u)} = s \mid \mathbf{x}^{(u)}, \lambda_{\text{old}}) \\ &= \frac{p(\mathbf{x}^{(u)}, z_n^{(u)} = s \mid \lambda_{\text{old}})}{p(\mathbf{x}^{(u)} \mid \lambda_{\text{old}})}\end{aligned}$$



Expected # of samples from state s' at time $n - 1$ to state s at time n in sample u

$$\begin{aligned}\xi_{s,s'}^{(u)}(n-1) &= \hat{q}(z_{n-1}^{(u)} = s, z_n^{(u)} = s') \\ &= p(z_{n-1}^{(u)} = s, z_n^{(u)} = s' \mid \mathbf{x}^{(u)}, \lambda_{\text{old}}) \\ &= \frac{p(\mathbf{x}^{(u)}, z_{n-1}^{(u)} = s, z_n^{(u)} = s' \mid \lambda_{\text{old}})}{p(\mathbf{x}^{(u)} \mid \lambda_{\text{old}})}\end{aligned}$$



Sufficient Statistics

Auxiliary function

$$\begin{aligned}
 Q(\lambda, \lambda_{\text{old}}) &= \sum_{u=1}^U \sum_{\text{all } \mathbf{z}^{(u)}} \hat{q}(\mathbf{z}^{(u)}) \ln p(\mathbf{x}^{(u)}, \mathbf{z}^{(u)} | \lambda) \\
 &= \sum_{s=1}^S \sum_{u=1}^U \gamma_s^{(u)} (n=1) \ln \pi_s + \sum_{s=1}^S \sum_{s'=1}^S \sum_{u=1}^U \sum_{n=2}^{N^{(u)}} \xi_{s,s'}^{(u)} (n-1) \ln A_{s,s'} \\
 &\quad + \sum_{s=1}^S \sum_{u=1}^U \sum_{\text{all } \mathbf{o}} \sum_{n \in \{\mathbf{x}_n^{(u)} = \text{"o"}\}} \gamma_s^{(u)} (n) \ln B_s(\text{"o"}) \quad \text{Sum of posterior probabilities over all sequences} \\
 &= \sum_{s=1}^S \gamma_s (n=1) \ln \pi_s + \sum_{s=1}^S \sum_{s'=1}^S \xi_{s,s'} \ln A_{s,s'} + \sum_{s=1}^S \sum_{\text{all "o"}} \gamma_s(\text{"o"}) \ln B_s(\text{"o"})
 \end{aligned}$$

Sufficient statistics

Expected # of samples in state s at time $n=1$

$$\gamma_s(n=1) = \sum_{u=1}^U \gamma_s^{(u)}(n=1)$$

Expected # of samples from state s to state s'

$$\xi_{s,s'} = \sum_{u=1}^U \sum_{n=2}^{N^{(u)}} \xi_{s,s'}^{(u)}(n-1)$$

Expected # of samples of observing "o" in state s

$$\gamma_s(\text{"o"}) = \sum_{u=1}^U \sum_{n \in \{\mathbf{x}_n^{(u)} = \text{"o"}\}} \gamma_s^{(u)}(n)$$

M-Step

Auxiliary function

$$Q(\lambda, \lambda_{\text{old}}) = \sum_{s=1}^S \gamma_s(n=1) \ln \pi_s + \sum_{s=1}^S \sum_{s'=1}^S \xi_{s,s'} \ln A_{s,s'} + \sum_{s=1}^S \sum_{\text{all "o"}} \gamma_s(\text{"o"}) \ln B_s(\text{"o"})$$

ML estimates

$$\left. \frac{\partial \left\{ Q(\lambda, \lambda_{\text{old}}) + \varepsilon \left(1 - \sum_{s=1}^S \pi_s \right) \right\}}{\partial \pi_s} \right|_{\lambda=\hat{\lambda}} = 0$$
$$\left. \frac{\partial \left\{ Q(\lambda, \lambda_{\text{old}}) + \varepsilon \left(1 - \sum_{s'=1}^S A_{s,s'} \right) \right\}}{\partial A_{s,s'}} \right|_{\lambda=\hat{\lambda}} = 0$$
$$\left. \frac{\partial \left\{ Q(\lambda, \lambda_{\text{old}}) + \varepsilon \left(1 - \sum_{\text{all "o"}} B_s(\text{"o"}) \right) \right\}}{\partial B_s(\text{"o"})} \right|_{\lambda=\hat{\lambda}} = 0$$

For each state,

Initial state probability $\hat{\pi}_s = \frac{\gamma_s(n=1)}{\sum_{s=1}^S \gamma_s(n=1)}$

Transition probability $\hat{A}_{s,s'} = \frac{\xi_{s,s'}}{\sum_{s'=1}^S \xi_{s,s'}}$

Output probability $\hat{B}_s(\text{"o"}) = \frac{\gamma_s(\text{"o"})}{\sum_{\text{all "o"}} \gamma_s(\text{"o"})}$

App.: Auxiliary Function

$$Q(\lambda, \lambda_{\text{old}}) = \sum_{u=1}^U \sum_{\text{all } \mathbf{z}^{(u)}} \hat{q}(\mathbf{z}^{(u)}) \ln p(\mathbf{x}^{(u)}, \mathbf{z}^{(u)} | \lambda)$$

$$= \sum_{u=1}^U \sum_{\text{all } \mathbf{z}^{(u)}} \hat{q}(\mathbf{z}^{(u)}) \ln \left\{ p(z_1^{(u)} | \lambda) \left[\prod_{n=2}^{N^{(u)}} p(z_n^{(u)} | z_{n-1}^{(u)}, \lambda) \right] \prod_{n=1}^{N^{(u)}} p(\mathbf{x}_n^{(u)} | z_n^{(u)}, \lambda) \right\}$$

$$= \sum_{u=1}^U \left\{ \sum_{\text{all } \mathbf{z}^{(u)}} \hat{q}(\mathbf{z}^{(u)}) \ln p(z_1^{(u)} | \lambda) + \sum_{n=2}^{N^{(u)}} \sum_{\text{all } \mathbf{z}^{(u)}} \hat{q}(\mathbf{z}^{(u)}) \ln p(z_n^{(u)} | z_{n-1}^{(u)}, \lambda) \right.$$

$$\left. + \sum_{n=1}^{N^{(u)}} \sum_{\text{all } \mathbf{z}^{(u)}} \hat{q}(\mathbf{z}^{(u)}) \ln p(\mathbf{x}_n^{(u)} | z_n^{(u)}, \lambda) \right\}$$

$$\sum_{s=1}^S \hat{q}(z_1^{(u)} = s) \ln p(z_1^{(u)} = s | \lambda) = \gamma_s^{(u)}(n=1) = \pi_s$$

$$\sum_{n=2}^{N^{(u)}} \sum_{s=1}^S \sum_{s'=1}^S \hat{q}(z_n^{(u)} = s, z_{n-1}^{(u)} = s') \cdot \ln p(z_n^{(u)} = s | z_{n-1}^{(u)} = s', \lambda) = \xi_{s,s'}^{(u)}(n-1) = A_{s,s'}$$

$$\sum_{n=1}^{N^{(u)}} \sum_{s=1}^S \hat{q}(z_n^{(u)} = s) \ln p(\mathbf{x}_n^{(u)} | z_n^{(u)} = s, \lambda) = \gamma_s^{(u)}(n) = \sum_{\text{all "o"} n \in \{\mathbf{x}_n^{(u)} = \text{"o"}\}} \sum \ln B_s(\text{"o"})$$

$$\begin{aligned} & \sum_{\text{all } \mathbf{z}} p(\mathbf{z}) f(z_1) \\ &= \sum_{\text{all } \mathbf{z}} p(z_2, \dots, z_N | z_1) p(z_1) f(z_1) \\ &= \sum_{\text{all } z_1} p(z_1) f(z_1) \end{aligned}$$

Hidden Markov Model Summary

Expectation Maximization

Maximum likelihood for the HMM

- Marginal likelihood

$$p(\mathbf{X}|\boldsymbol{\theta}) = \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$$

$$\begin{aligned}\gamma(\mathbf{z}_n) &= p(\mathbf{z}_n|\mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \\ \xi(\mathbf{z}_{n-1}, \mathbf{z}_n) &= p(\mathbf{z}_{n-1}, \mathbf{z}_n|\mathbf{X}, \boldsymbol{\theta}^{\text{old}})\end{aligned}$$

- EM

- $\max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) = \sum_{k=1}^K \gamma(z_{1k}) \ln \pi_k + \sum_{n=2}^N \sum_{j=1}^K \sum_{k=1}^K \xi(z_{n-1,j}, z_{nk}) \ln A_{jk}$$

the same as GMM

$$+ \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \ln p(\mathbf{x}_n|\boldsymbol{\phi}_k).$$

(

Calculation

- E-step

$$\begin{aligned}\gamma(\mathbf{z}_n) &= p(\mathbf{z}_n | \mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \\ \xi(\mathbf{z}_{n-1}, \mathbf{z}_n) &= p(\mathbf{z}_{n-1}, \mathbf{z}_n | \mathbf{X}, \boldsymbol{\theta}^{\text{old}})\end{aligned}$$

- M-step

$$\begin{aligned}\pi_k &= \frac{\gamma(z_{1k})}{\sum_{j=1}^K \gamma(z_{1j})} & A_{jk} &= \frac{\sum_{n=2}^N \xi(z_{n-1,j}, z_{nk})}{\sum_{l=1}^K \sum_{n=2}^N \xi(z_{n-1,j}, z_{nl})} \\ \boldsymbol{\mu}_k &= \frac{\sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n}{\sum_{n=1}^N \gamma(z_{nk})} & \boldsymbol{\Sigma}_k &= \frac{\sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T}{\sum_{n=1}^N \gamma(z_{nk})}\end{aligned}$$

Key quantities

$$\begin{aligned}\gamma(\mathbf{z}_n) &= p(\mathbf{z}_n | \mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \\ \xi(\mathbf{z}_{n-1}, \mathbf{z}_n) &= p(\mathbf{z}_{n-1}, \mathbf{z}_n | \mathbf{X}, \boldsymbol{\theta}^{\text{old}})\end{aligned}$$

$$\gamma(\mathbf{z}_n) = \frac{p(\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{z}_n) p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_n)}{p(\mathbf{X})} = \frac{\alpha(\mathbf{z}_n) \beta(\mathbf{z}_n)}{p(\mathbf{X})}$$

$$\xi(\mathbf{z}_{n-1}, \mathbf{z}_n) = \frac{\alpha(\mathbf{z}_{n-1}) p(\mathbf{x}_n | \mathbf{z}_n) p(\mathbf{z}_n | \mathbf{z}_{n-1}) \beta(\mathbf{z}_n)}{p(\mathbf{X})}$$

$$p(\mathbf{X}) = \sum_{\mathbf{z}_n} \alpha(\mathbf{z}_n) \beta(\mathbf{z}_n) \quad p(\mathbf{X}) = \sum_{\mathbf{z}_N} \alpha(\mathbf{z}_N)$$

Inference

- Marginal likelihood (Forward-Backward / sum-product): $P(X)$
- Best state path (Viterbi / max-sum): $\max_z p(X, Z)$
- Prediction: $p(x_{N+1} | X)$

Prediction

$$\begin{aligned} p(\mathbf{x}_{N+1}|\mathbf{X}) &= \sum_{\mathbf{z}_{N+1}} p(\mathbf{x}_{N+1}, \mathbf{z}_{N+1}|\mathbf{X}) \\ &= \sum_{\mathbf{z}_{N+1}} p(\mathbf{x}_{N+1}|\mathbf{z}_{N+1})p(\mathbf{z}_{N+1}|\mathbf{X}) \\ &= \sum_{\mathbf{z}_{N+1}} p(\mathbf{x}_{N+1}|\mathbf{z}_{N+1}) \sum_{\mathbf{z}_N} p(\mathbf{z}_{N+1}, \mathbf{z}_N|\mathbf{X}) \\ &= \sum_{\mathbf{z}_{N+1}} p(\mathbf{x}_{N+1}|\mathbf{z}_{N+1}) \sum_{\mathbf{z}_N} p(\mathbf{z}_{N+1}|\mathbf{z}_N)p(\mathbf{z}_N|\mathbf{X}) \\ &= \sum_{\mathbf{z}_{N+1}} p(\mathbf{x}_{N+1}|\mathbf{z}_{N+1}) \sum_{\mathbf{z}_N} p(\mathbf{z}_{N+1}|\mathbf{z}_N) \frac{p(\mathbf{z}_N, \mathbf{X})}{p(\mathbf{X})} \\ &= \frac{1}{p(\mathbf{X})} \sum_{\mathbf{z}_{N+1}} p(\mathbf{x}_{N+1}|\mathbf{z}_{N+1}) \sum_{\mathbf{z}_N} p(\mathbf{z}_{N+1}|\mathbf{z}_N) \alpha(\mathbf{z}_N) \end{aligned}$$

Thinking



- Recover: $p(x_n | x_1, \dots, x_{n-1}, x_{n+1}, \dots, x_N)$

Scaling factors for HMM

- $\alpha(z_n)$ is obtained from the previous value $\alpha(z_{n-1})$ by multiplying by quantities $p(z_n/z_{n-1})$ and $p(x_n/z_n)$.
- $p(z_n/z_{n-1}) < 1 \Rightarrow$ the values of $\alpha(z_n)$ can go to zero exponentially quickly.
- Rescale $\alpha(z_n)$ as

$$\hat{\alpha}(z_n) = p(z_n | x_1, \dots, x_n) = \frac{\alpha(z_n)}{p(x_1, \dots, x_n)} \quad c_n = p(x_n | x_1, \dots, x_{n-1})$$

- Then

$$\alpha(z_n) = p(z_n | x_1, \dots, x_n) p(x_1, \dots, x_n) = \left(\prod_{m=1}^n c_m \right) \hat{\alpha}(z_n)$$

Rewrite forward-backward recursion

- Forward (c_n the coefficient that normalizes the right side)

$$c_n \hat{\alpha}(\mathbf{z}_n) = p(\mathbf{x}_n | \mathbf{z}_n) \sum_{\mathbf{z}_{n-1}} \hat{\alpha}(\mathbf{z}_{n-1}) p(\mathbf{z}_n | \mathbf{z}_{n-1}) \quad \alpha(\mathbf{z}_n) = p(\mathbf{x}_n | \mathbf{z}_n) \sum_{\mathbf{z}_{n-1}} \alpha(\mathbf{z}_{n-1}) p(\mathbf{z}_n | \mathbf{z}_{n-1})$$

- Backward

$$c_{n+1} \hat{\beta}(\mathbf{z}_n) = \sum_{\mathbf{z}_{n+1}} \hat{\beta}(\mathbf{z}_{n+1}) p(\mathbf{x}_{n+1} | \mathbf{z}_{n+1}) p(\mathbf{z}_{n+1} | \mathbf{z}_n) \quad \beta(\mathbf{z}_n) = \sum_{\mathbf{z}_{n+1}} \beta(\mathbf{z}_{n+1}) p(\mathbf{x}_{n+1} | \mathbf{z}_{n+1}) p(\mathbf{z}_{n+1} | \mathbf{z}_n)$$

- Marginal

$$\begin{aligned} \gamma(\mathbf{z}_n) &= \hat{\alpha}(\mathbf{z}_n) \hat{\beta}(\mathbf{z}_n) \\ \xi(\mathbf{z}_{n-1}, \mathbf{z}_n) &= c_n \hat{\alpha}(\mathbf{z}_{n-1}) p(\mathbf{x}_n | \mathbf{z}_n) p(\mathbf{z}_n | \mathbf{z}_{n-1}) \hat{\beta}(\mathbf{z}_n) \end{aligned}$$

HMM Extensions

1. Generalize model of state duration:

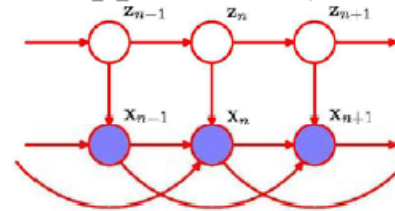
- vanilla HMM restricted in model of how long stay in state - prob. that model will spend D steps in state k and then transition out:

$$P(D) = (A_{kk})^D (1 - A_{kk}) \propto \exp(-D \log A_{kk})$$

- instead associate distribution with time spent in state k : $P(t|k)$ (see *semi-Markov* models for sequence segmentation applications)

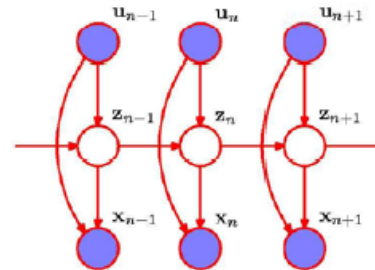
2. Combine with auto-regressive Markov model:

- include long-range relationships
- directly model relations between observations



3. Supervised setting:

- include additional observations
- *input-output* HMM



$$z_{n+1} \perp\!\!\!\perp z_{n-1} \mid z_n$$

4. *factorial hidden Markov model* (Ghahramani and Jordan, 1997)

- reduce latent states
- add Markov chains
- no conditional independence

