

Special Topics in Natural Language Processing

CS6980

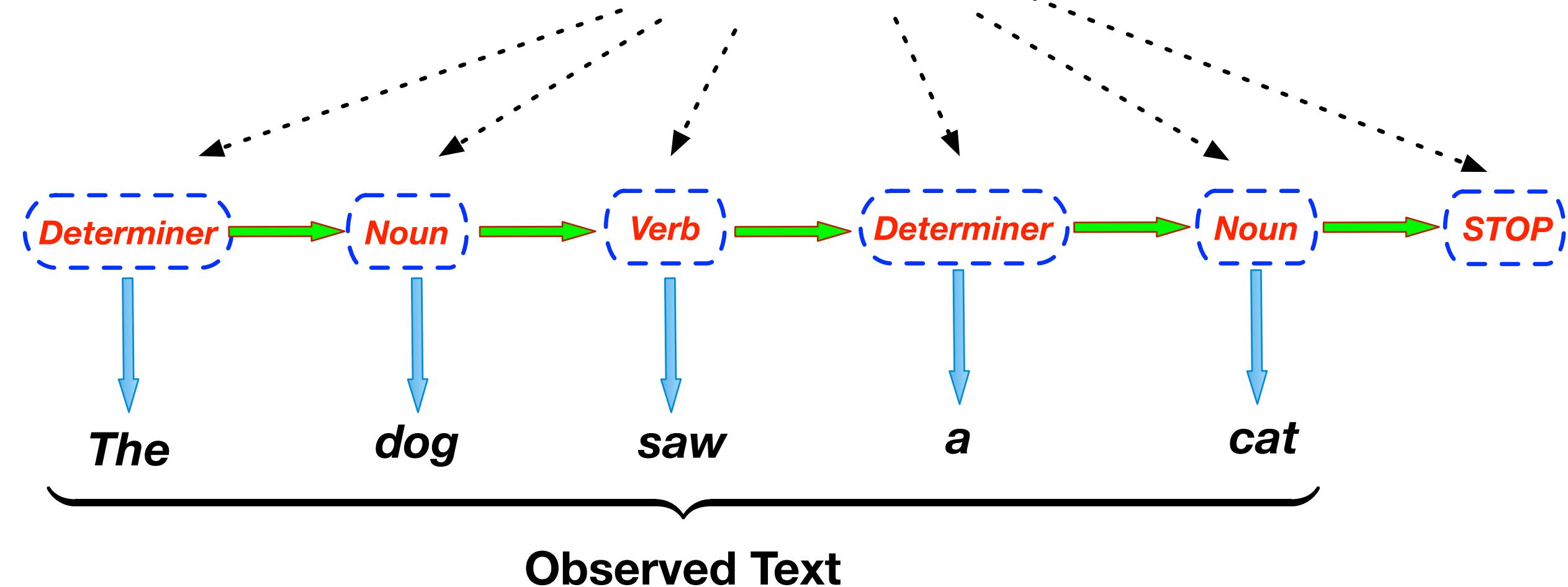
Ashutosh Modi
CSE Department, IIT Kanpur



Lecture 11: Sequence Prediction 2
Jan 29, 2020

Hidden Markov Models (HMM)

Latent or Hidden States



HMM Setting

Set of States (Λ)

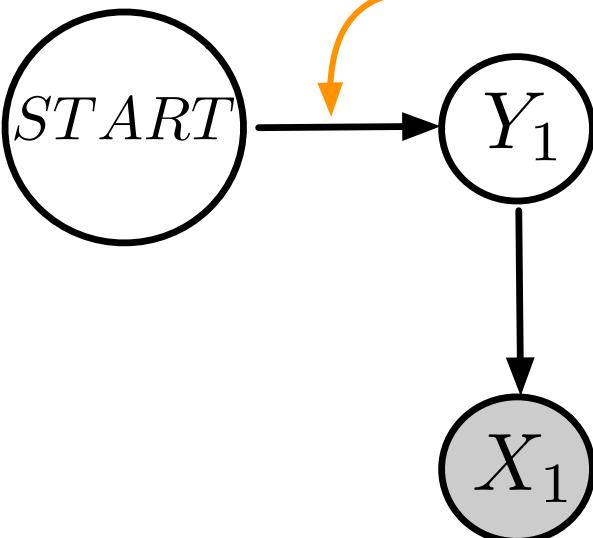
$$\{c_0 = \text{START}, c_1, c_2, \dots, c_K, c_{K+1} = \text{STOP}\}$$

Set of Observations

$$\{w_1, w_2, \dots w_J\}$$

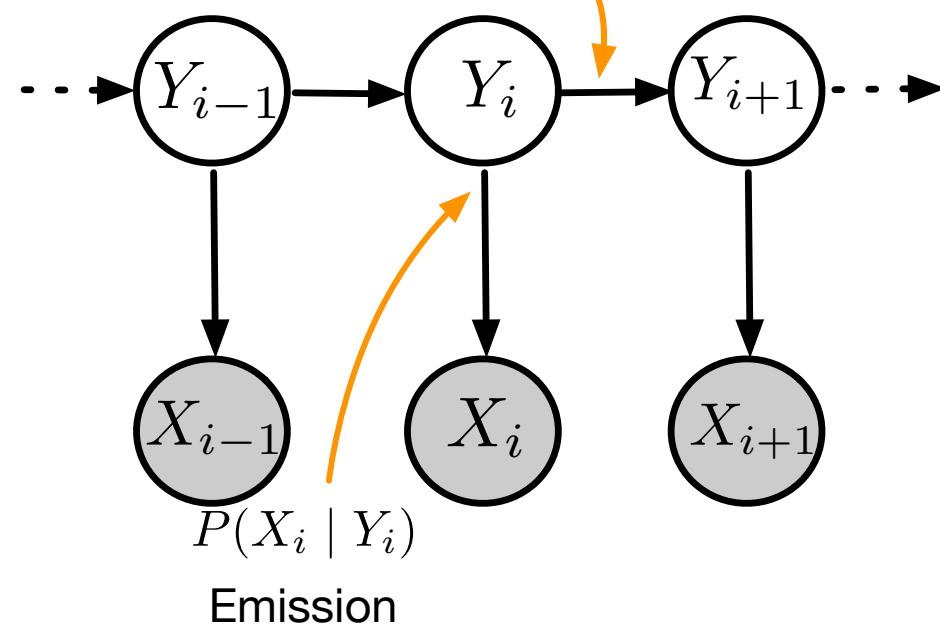
Start
Probability

$$P(Y_1 | Y_0 = \text{START})$$



Transition
Probability

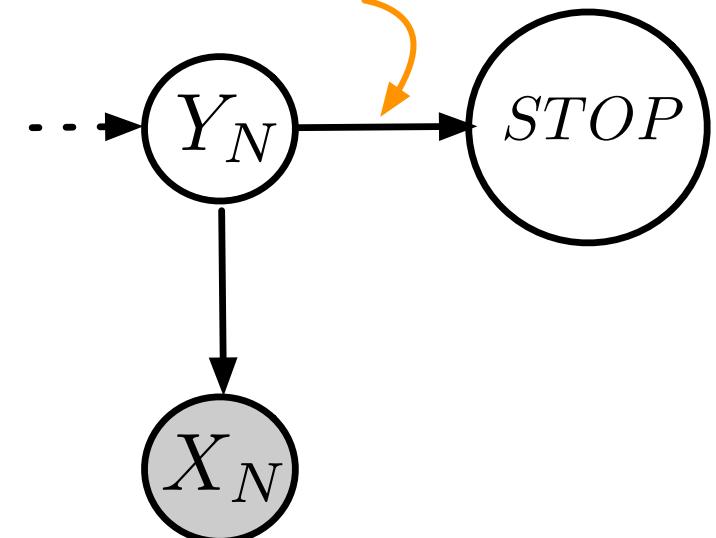
$$P(Y_{i+1} | Y_i)$$



Emission
Probability

End
Probability

$$P(Y_{N+1} = \text{STOP} | Y_N)$$



HMM Joint Distribution

- We have a generative model and are interested in the joint distribution

$$P(X, Y) = P(X_1 = x_1, \dots, X_N = x_N, Y_1 = y_1, \dots, Y_N = y_N)$$

$$P(X_1 = x_1, \dots, X_N = x_N, Y_1 = y_1, \dots, Y_N = y_N) =$$

$$P(y_1|START) \times \left(\prod_{i=1}^{N-1} P(y_{i+1}|y_i) \right) \times \left(\prod_{i=1}^N P(x_i|y_i) \right) \times P(STOP|y_N)$$

Initial / Start
Probability

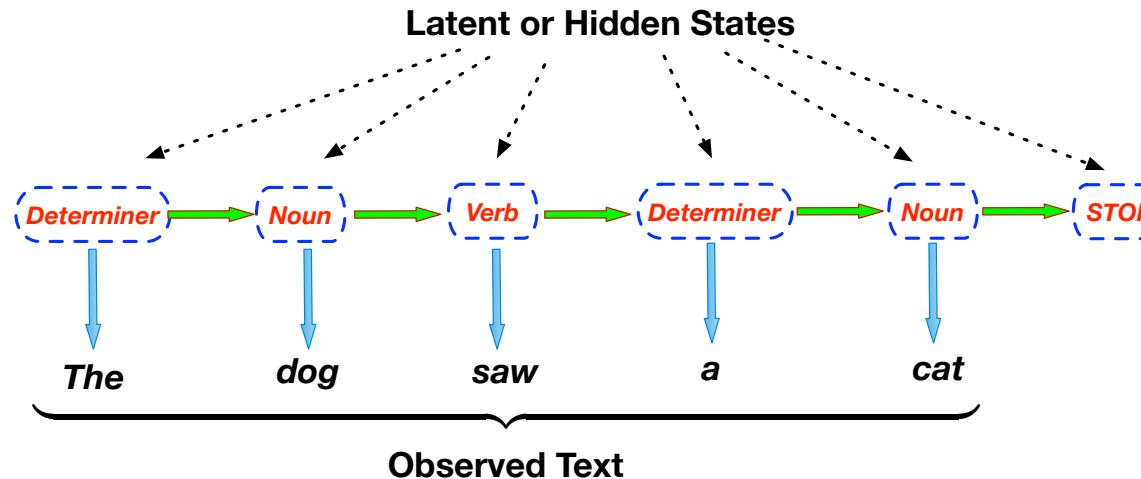
Transition
Probability

Emission
Probability

End / Final
Probability



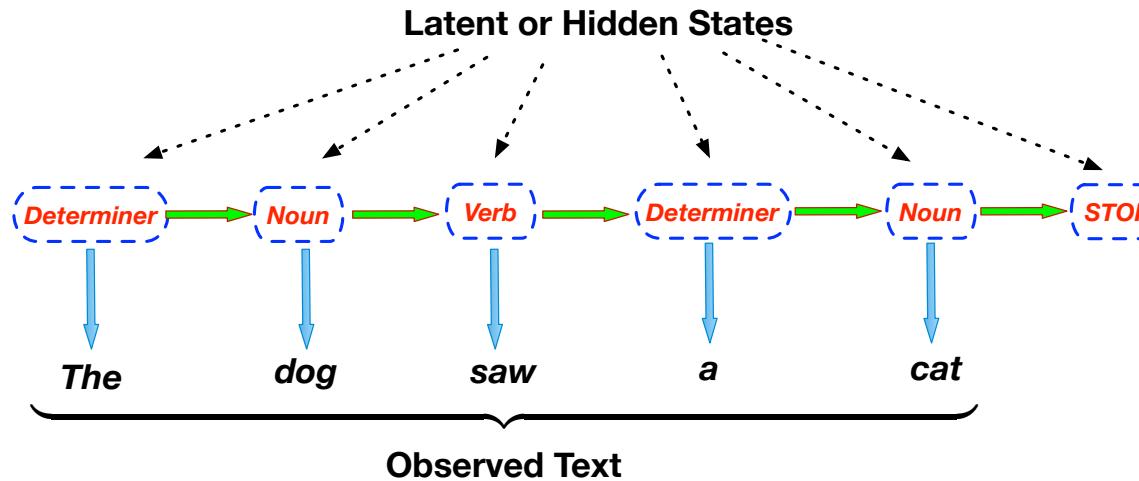
Hidden Markov Models (HMM)



Three Questions:

1. **Learning Problem:** How do we estimate the parameters of distributions?
2. **Decoding Problem:** Given the observed text what is the hidden POS sequence that best explains the observation?
3. **Likelihood Problem:** Given the parameters of the distributions what is the probability of the observed sequence?

Hidden Markov Models (HMM)



Three Questions:

1. **Learning Problem:** How do we estimate the parameters of distributions?
2. **Decoding Problem:** Given the observed text what is the hidden POS sequence that best explains the observation?
3. **Likelihood Problem:** Given the parameters of the distributions what is the probability of the observed sequence?

Learning Problem

$$P_{\text{init}}(y_1 = c_k | \text{START}) = \frac{\text{Count}(y_1 = c_k)}{\left(\sum_{l=1}^K \text{Count}(y_1 = c_l)\right)} = M$$

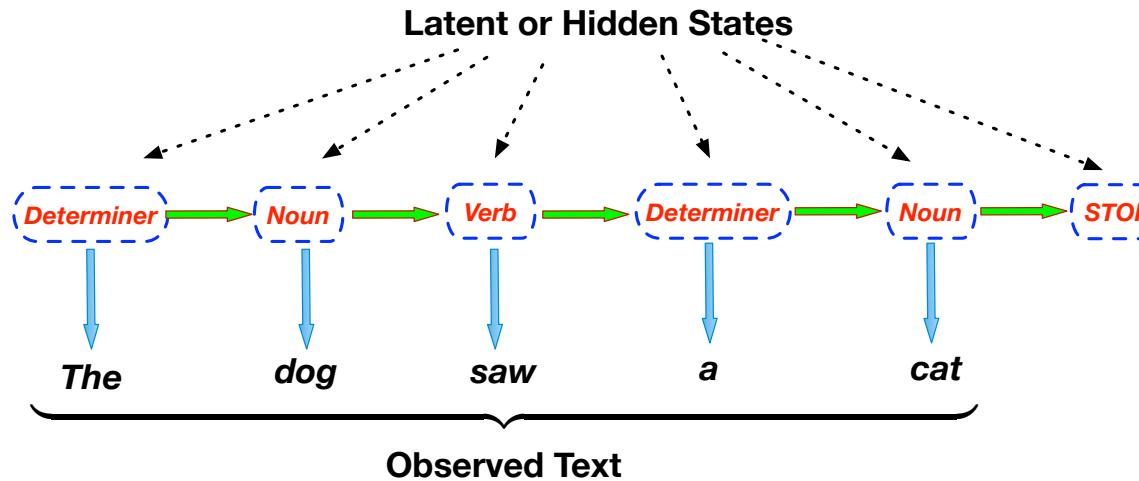
$$P_{\text{trans}}(y_{i+1} = c_k | y_i = c_l) = \frac{\text{Count}(c_k, c_l)}{\text{Count}(c_l)}$$

$$P_{\text{emiss}}(x_i = w_j | y_i = c_k) = \frac{\text{Count}(w_j, c_k)}{\text{Count}(c_k)}$$

$$P_{\text{final}}(\text{STOP} | y_N = c_k) = \frac{\text{Count}(y_N = c_k)}{\left(\sum_{l=1}^K \text{Count}(y_N = c_l)\right)} = M$$



Hidden Markov Models (HMM)



Three Questions:

1. **Learning Problem:** How do we estimate the parameters of distributions?
2. **Decoding Problem:** Given the observed text what is the hidden POS sequence that best explains the observation?
3. **Likelihood Problem:** Given the parameters of the distributions what is the probability of the observed sequence?



Decoding Problem

- Two ways possible



Decoding Problem

- Two ways possible
 - **Posterior Decoding:** Maximize probability of each state

$$y_i^* = \arg \max_{y_i \in \Lambda} P(Y_i = y_i | X_1 = x_1, \dots, X_N = x_N)$$



Decoding Problem

- Two ways possible
 - **Posterior Decoding:** Maximize probability of each state

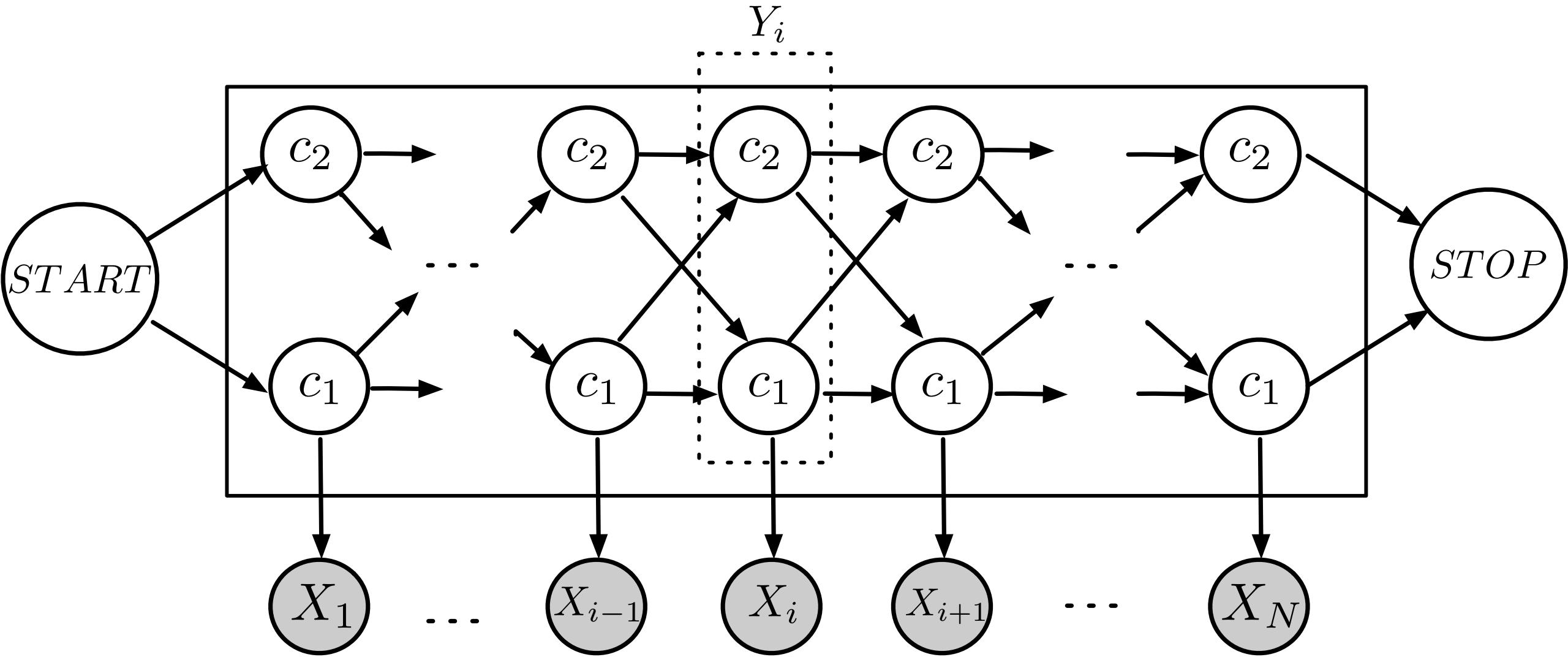
$$y_i^* = \arg \max_{y_i \in \Lambda} P(Y_i = y_i | X_1 = x_1, \dots, X_N = x_N)$$

- **Viterbi Decoding:** Maximize the probability of the entire sequence

$$\begin{aligned} y^* &= \arg \max_{y=y_1 \dots y_N} P(Y_1 = y_1, \dots, Y_N = y_N | X_1 = x_1, \dots, X_N = x_N) \\ &= \arg \max_{y=y_1 \dots y_N} P(Y_1 = y_1, \dots, Y_N = y_N, X_1 = x_1, \dots, X_N = x_N) \end{aligned}$$



Trellis Diagram



Decoding Problem

- Decoding requires the following calculation

$$P(Y | X) = \frac{P(X, Y)}{P(X)}$$

$$P(Y | X) = \frac{P(X, Y)}{\sum_Y P(X, Y)}$$



Decoding Problem

- Decoding requires the following calculation

$$P(Y | X) = \frac{P(X, Y)}{P(X)}$$

$$P(Y | X) = \frac{P(X, Y)}{\sum_Y P(X, Y)}$$

But what is the sum in the denominator over?

Over all possible sequences $Y = y_1, \dots, y_N$



Decoding Problem

In general, number of sequences = K^N



Decoding Problem

In general, number of sequences = K^N

**EXPONENTIAL
IN LENGTH**



Decoding Problem

In general, number of sequences = K^N

**EXPONENTIAL
IN LENGTH**

What do we do?



Decoding Problem

In general, number of sequences = K^N

EXponential
in length

What do we do?

DYNAMIC PROGRAMMING



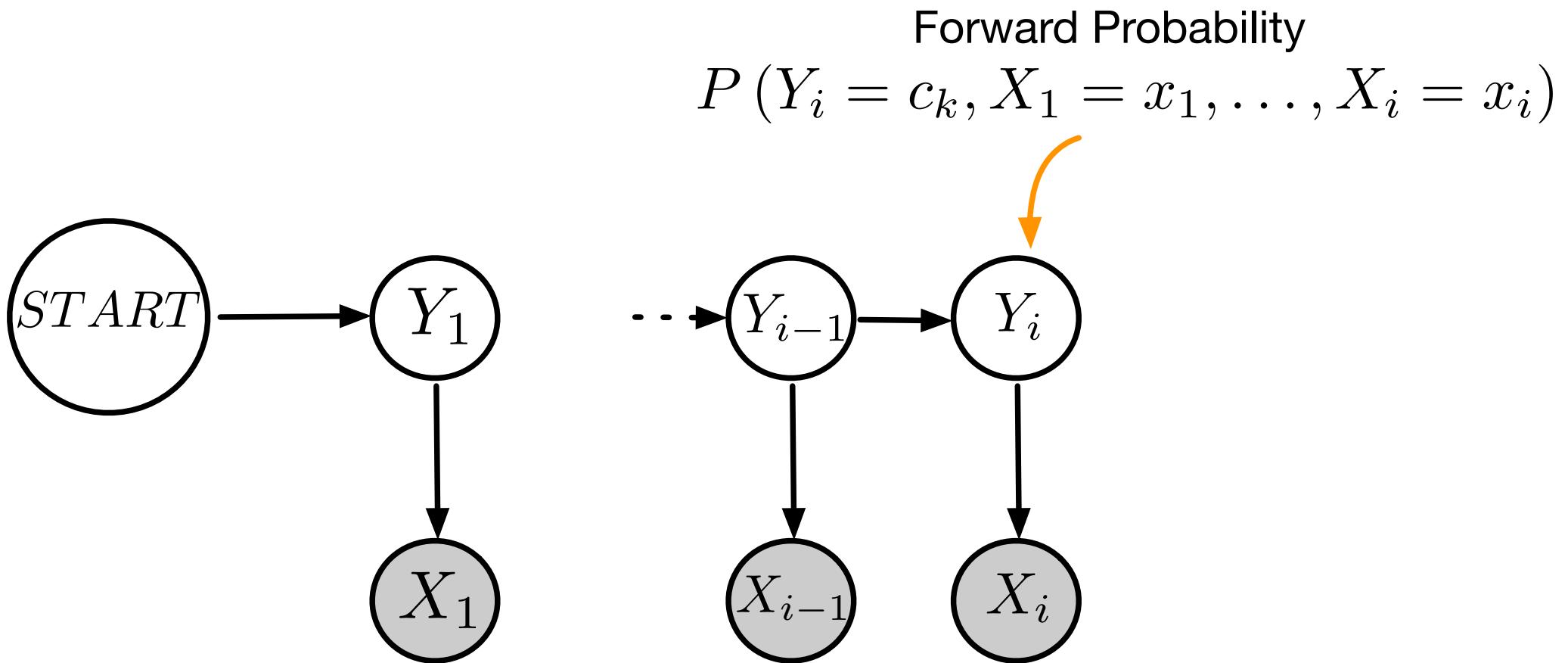
Forward Probability

$\text{forward}(i, c_k) := P(Y_i = c_k, X_1 = x_1, \dots, X_i = x_i)$



Forward Probability

$\text{forward}(i, c_k) := P(Y_i = c_k, X_1 = x_1, \dots, X_i = x_i)$



Forward Probability

$$\text{forward}(i, c_k) := P(Y_i = c_k, X_1 = x_1, \dots, X_i = x_i)$$

$$\text{forward}(i, c_k) = \left(\sum_{c_l} P(Y_i = c_k \mid Y_{i-1} = c_l) \times \text{forward}(i-1, c_l) \right) \times P(X_i \mid Y_i = c_k)$$



Forward Probability

$\text{forward}(i, c_k) := P(Y_i = c_k, X_1 = x_1, \dots, X_i = x_i)$



Forward Probability

$$\text{forward}(i, c_k) := P(Y_i = c_k, X_1 = x_1, \dots, X_i = x_i)$$

$$P(Y_i = c_k, X_1 = x_1, \dots, X_i = x_i) = P(Y_i = c_k, X_{1:i})$$



Forward Probability

$$\text{forward}(i, c_k) := P(Y_i = c_k, X_1 = x_1, \dots, X_i = x_i)$$

$$P(Y_i = c_k, X_1 = x_1, \dots, X_i = x_i) = P(Y_i = c_k, X_{1:i})$$

$$P(Y_i = c_k, X_{1:i})$$

$$= \sum_{Y_{i-1}} P(Y_i = c_k, Y_{i-1}, X_{1:i})$$

Marginalization



Forward Probability

$$\text{forward}(i, c_k) := P(Y_i = c_k, X_1 = x_1, \dots, X_i = x_i)$$

$$P(Y_i = c_k, X_1 = x_1, \dots, X_i = x_i) = P(Y_i = c_k, X_{1:i})$$

$$P(Y_i = c_k, X_{1:i})$$

$$= \sum_{Y_{i-1}} P(Y_i = c_k, Y_{i-1}, X_{1:i}) \quad \text{Marginalization}$$

$$= \sum_{Y_{i-1}} P(X_i \mid Y_i, Y_{i-1}, X_{1:i-1}) \times P(Y_i, Y_{i-1}, X_{1:i-1}) \quad \text{Bayes Rule}$$



Forward Probability

$$\text{forward}(i, c_k) := P(Y_i = c_k, X_1 = x_1, \dots, X_i = x_i)$$

$$P(Y_i = c_k, X_1 = x_1, \dots, X_i = x_i) = P(Y_i = c_k, X_{1:i})$$

$$P(Y_i = c_k, X_{1:i})$$

$$= \sum_{Y_{i-1}} P(Y_i = c_k, Y_{i-1}, X_{1:i}) \quad \text{Marginalization}$$

$$= \sum_{Y_{i-1}} P(X_i \mid Y_i, Y_{i-1}, X_{1:i-1}) \times P(Y_i, Y_{i-1}, X_{1:i-1}) \quad \text{Bayes Rule}$$

$$= \sum_{Y_{i-1}} P(X_i \mid Y_i, Y_{i-1}, X_{1:i-1}) \times P(Y_i \mid Y_{i-1}, X_{1:i-1}) \times P(Y_{i-1}, X_{1:i-1})$$



Forward Probability

$$\text{forward}(i, c_k) := P(Y_i = c_k, X_1 = x_1, \dots, X_i = x_i)$$

$$P(Y_i = c_k, X_1 = x_1, \dots, X_i = x_i) = P(Y_i = c_k, X_{1:i})$$

$$P(Y_i = c_k, X_{1:i})$$

$$= \sum_{Y_{i-1}} P(Y_i = c_k, Y_{i-1}, X_{1:i}) \quad \text{Marginalization}$$

$$= \sum_{Y_{i-1}} P(X_i \mid Y_i, Y_{i-1}, X_{1:i-1}) \times P(Y_i, Y_{i-1}, X_{1:i-1}) \quad \text{Bayes Rule}$$

$$= \sum_{Y_{i-1}} P(X_i \mid Y_i, Y_{i-1}, X_{1:i-1}) \times P(Y_i \mid Y_{i-1}, X_{1:i-1}) \times P(Y_{i-1}, X_{1:i-1})$$

$$= \sum_{Y_{i-1}} P(X_i \mid Y_i) \times P(Y_i \mid Y_{i-1}) \times \text{forward}(i - 1, y_{i-1}) \quad \text{Independence}$$



Forward Probability

$\text{forward}(i, c_k) := P(Y_i = c_k, X_1 = x_1, \dots, X_i = x_i)$

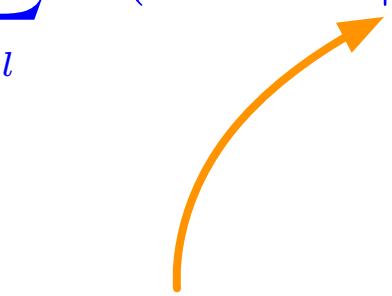
$$\text{forward}(i, c_k) = \left(\sum_{c_l} P(Y_i = c_k \mid Y_{i-1} = c_l) \times \text{forward}(i-1, c_l) \right) \times P(X_i \mid Y_i = c_k)$$



Forward Probability

$$\text{forward}(i, c_k) := P(Y_i = c_k, X_1 = x_1, \dots, X_i = x_i)$$

$$\text{forward}(i, c_k) = \left(\sum_{c_l} P(Y_i = c_k \mid Y_{i-1} = c_l) \times \text{forward}(i-1, c_l) \right) \times P(X_i \mid Y_i = c_k)$$



Forward Probability

$$\text{forward}(i, c_k) := P(Y_i = c_k, X_1 = x_1, \dots, X_i = x_i)$$

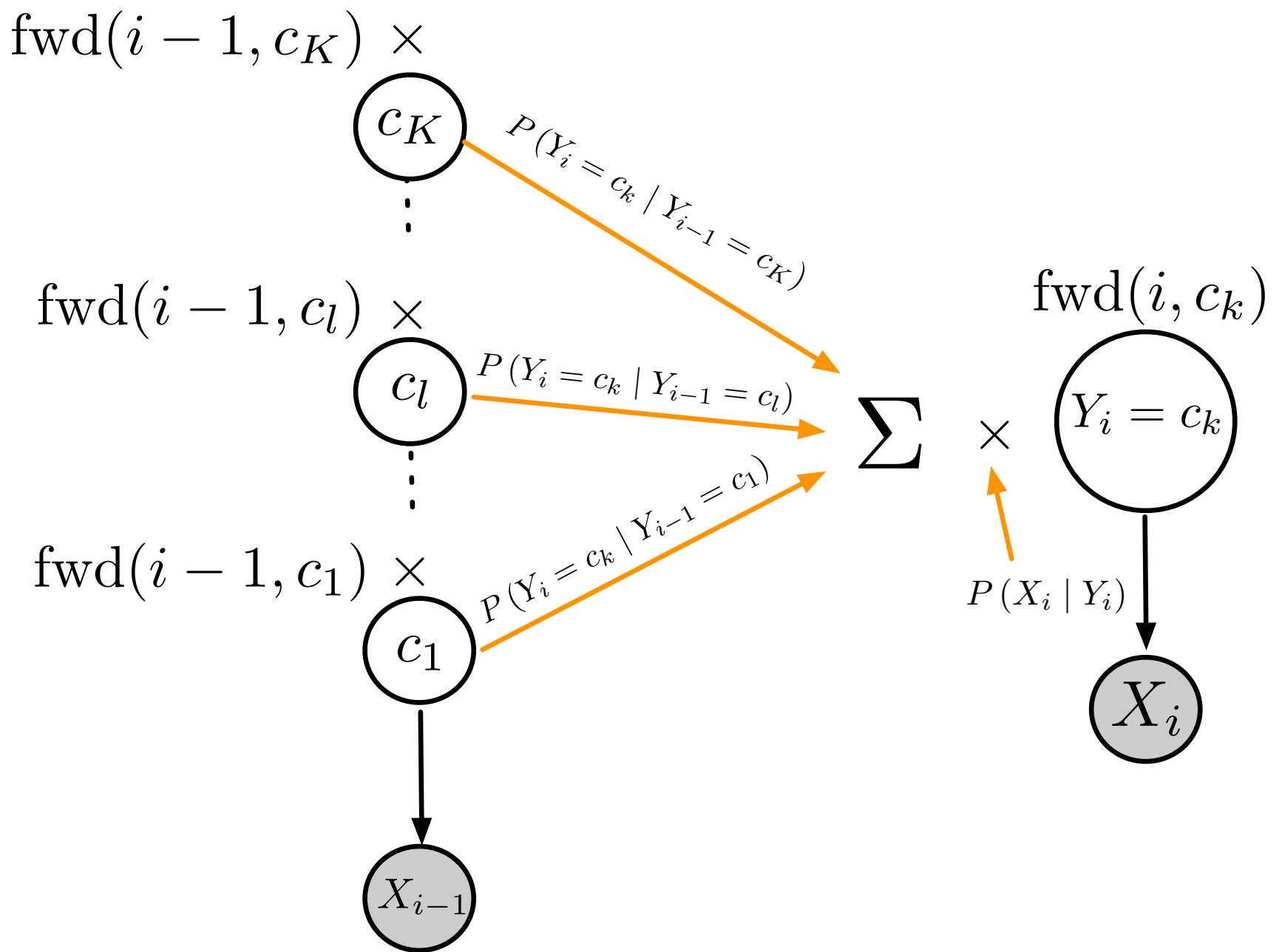
$$\text{forward}(i, c_k) = \left(\sum_{c_l} P(Y_i = c_k \mid Y_{i-1} = c_l) \times \text{forward}(i-1, c_l) \right) \times P(X_i \mid Y_i = c_k)$$

Transition Probability

Emission Probability



Forward Probability



Forward Probability

$\text{forward}(i, c_k) := P(Y_i = c_k, X_1 = x_1, \dots, X_i = x_i)$

$$\text{forward}(i, c_k) = \left(\sum_{c_l} P(Y_i = c_k \mid Y_{i-1} = c_l) \times \text{forward}(i-1, c_l) \right) \times P(X_i \mid Y_i = c_k)$$

$\text{forward}(N+1, STOP) = ?$



Forward Probability

$$\text{forward}(i, c_k) := P(Y_i = c_k, X_1 = x_1, \dots, X_i = x_i)$$

$$\text{forward}(i, c_k) = \left(\sum_{c_l} P(Y_i = c_k \mid Y_{i-1} = c_l) \times \text{forward}(i-1, c_l) \right) \times P(X_i \mid Y_i = c_k)$$

$$\text{forward}(N+1, STOP)$$

$$= P(Y_{N+1} = STOP, X_1 = x_1, \dots, X_N = x_N, X_{N+1})$$



Forward Probability

$$\text{forward}(i, c_k) := P(Y_i = c_k, X_1 = x_1, \dots, X_i = x_i)$$

$$\text{forward}(i, c_k) = \left(\sum_{c_l} P(Y_i = c_k \mid Y_{i-1} = c_l) \times \text{forward}(i-1, c_l) \right) \times P(X_i \mid Y_i = c_k)$$

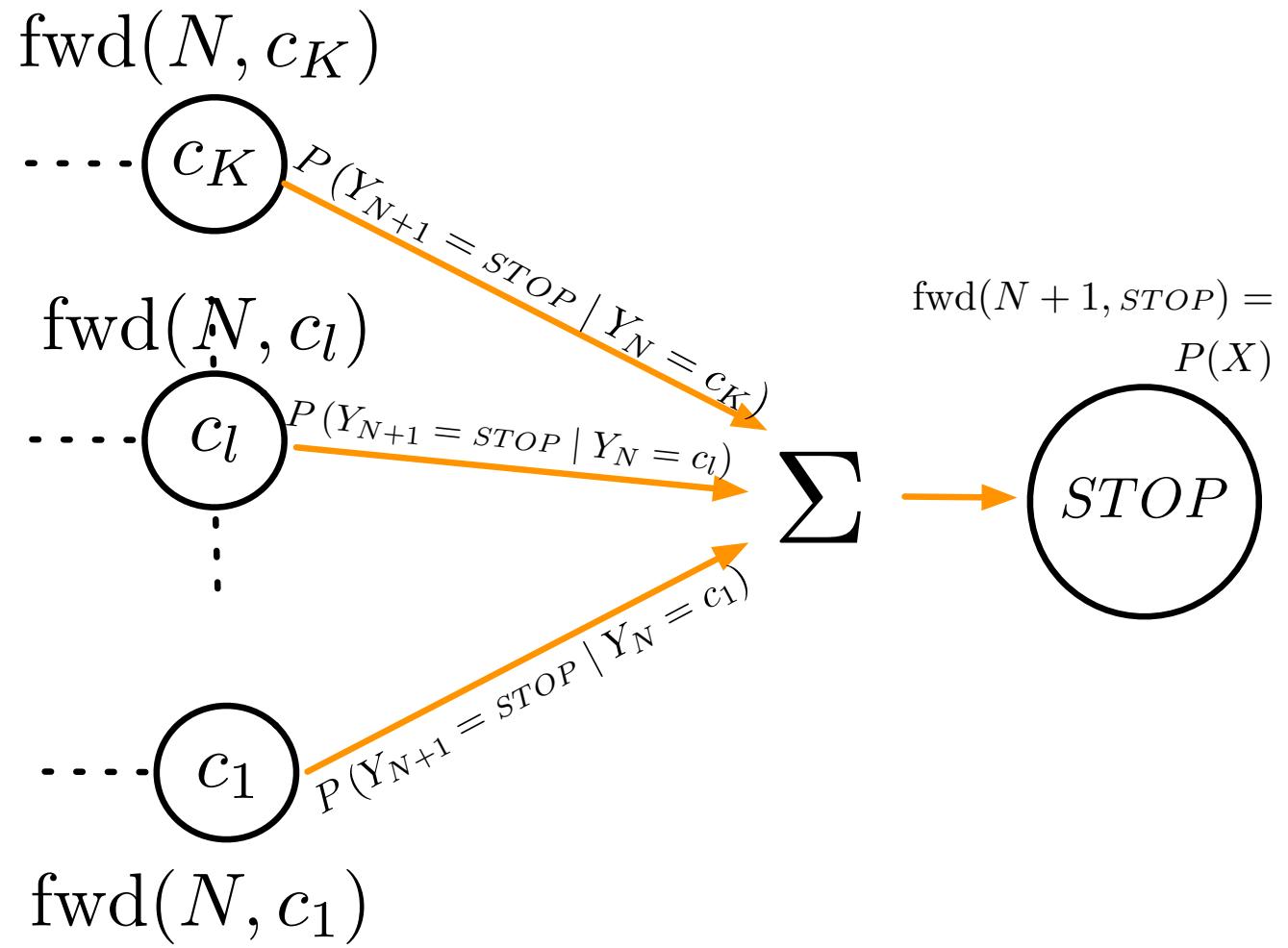
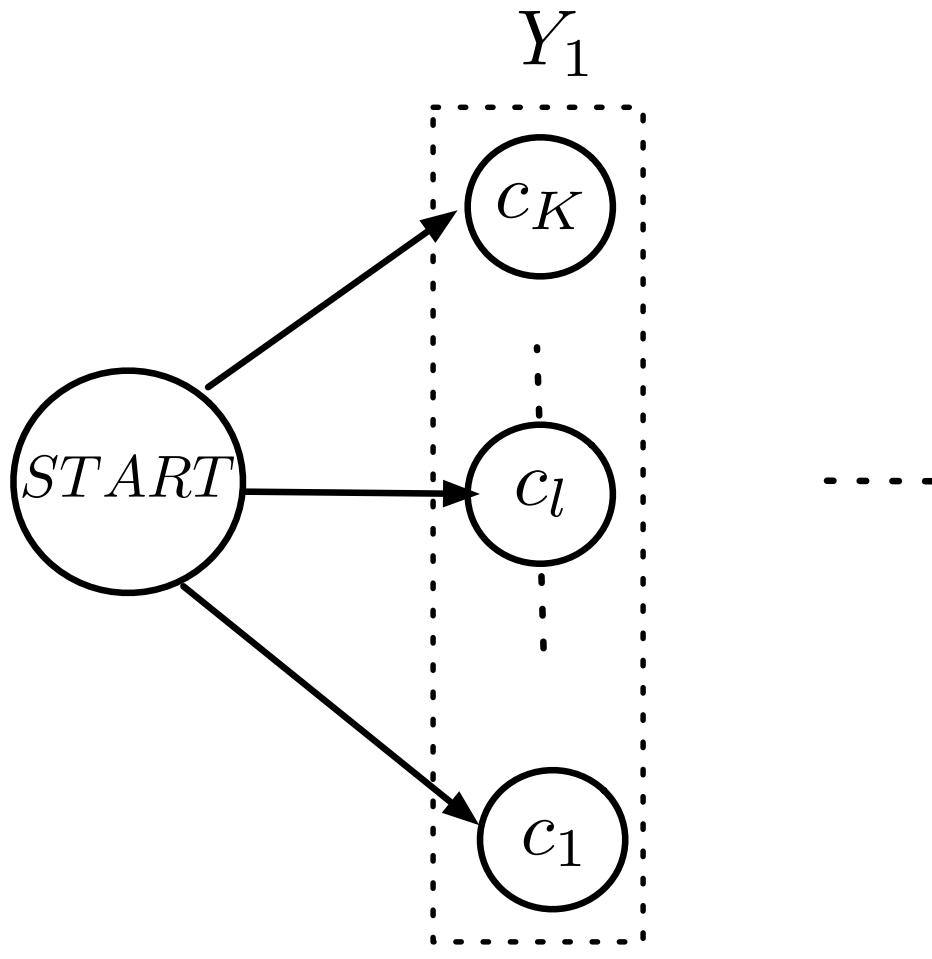
$$\text{forward}(N+1, STOP)$$

$$= P(Y_{N+1} = STOP, X_1 = x_1, \dots, X_N = x_N, X_{N+1})$$

$$= P(X)$$



Forward Probability



Forward Probability

$\text{forward}(i, c_k) := P(Y_i = c_k, X_1 = x_1, \dots, X_i = x_i)$

$$\text{forward}(i, c_k) = \left(\sum_{c_l} P(Y_i = c_k \mid Y_{i-1} = c_l) \times \text{forward}(i-1, c_l) \right) \times P(X_i \mid Y_i = c_k)$$

$\text{forward}(N+1, STOP)$

$= P(Y_{N+1} = STOP, X_1 = x_1, \dots, X_N = x_N, X_{N+1})$

$= P(X)$

$$P(X) = \text{forward}(N+1, STOP)$$



Decoding Problem

Brute force approach $\sim O(K^N)$

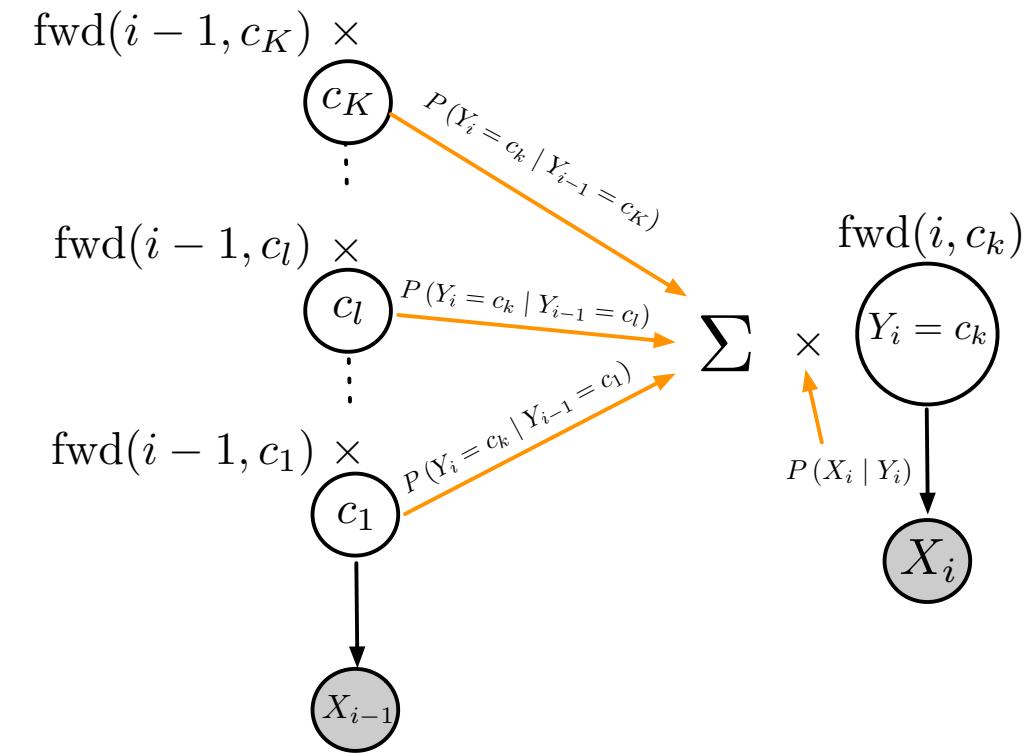
What is the complexity of forward probability calculations?



Decoding Problem

Brute force approach $\sim O(K^N)$

What is the complexity of forward probability calculations?

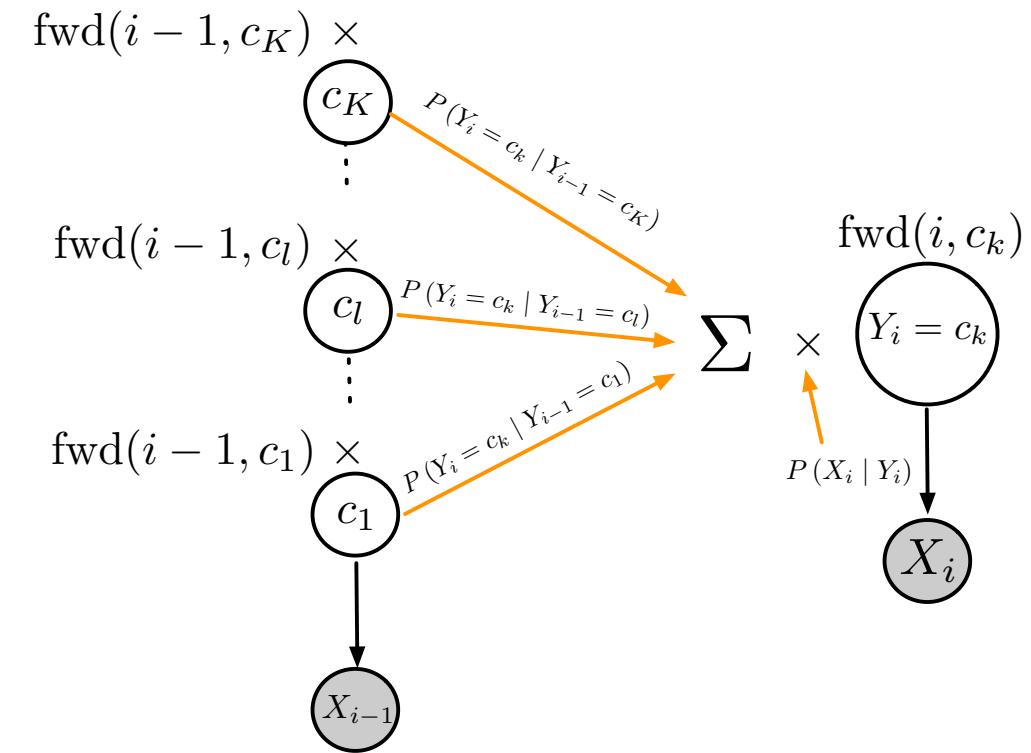


At any given time step, for every state we have information coming from K previous states. This accounts for $O(K)$ operations

Decoding Problem

Brute force approach $\sim O(K^N)$

What is the complexity of forward probability calculations?



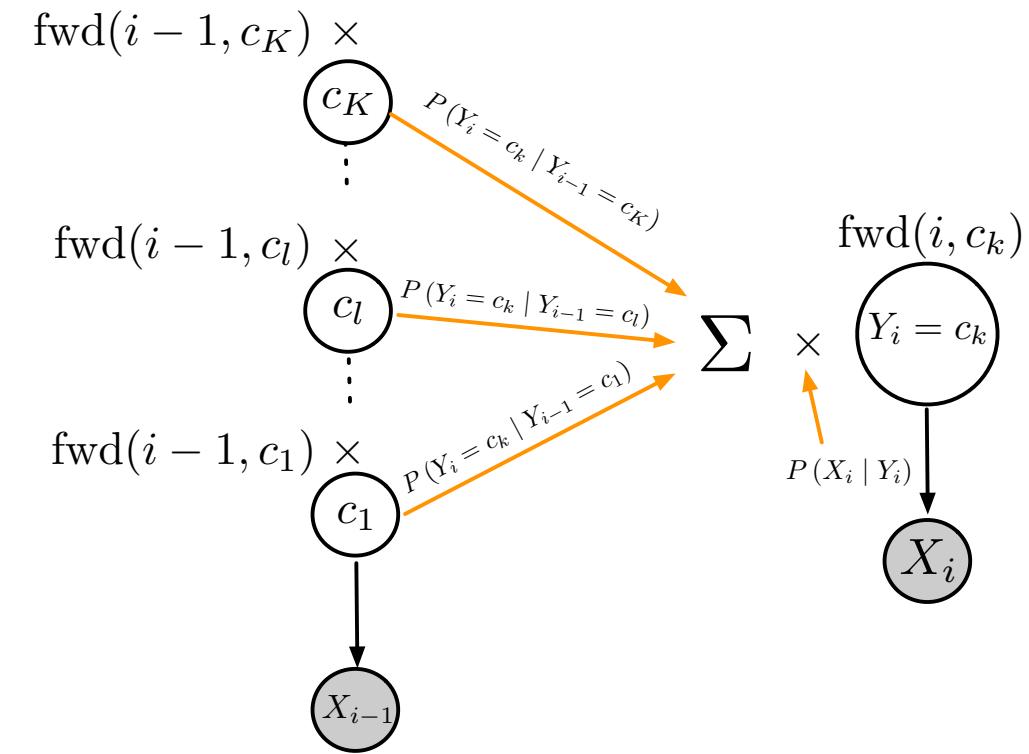
At any given time step, for every state we have information coming from K previous states. This accounts for $O(K)$ operations

We do this for all states at a given time step. This makes it $O(K^2)$ operations.

Decoding Problem

Brute force approach $\sim O(K^N)$

What is the complexity of forward probability calculations?



At any given time step, for every state we have information coming from K previous states. This accounts for $O(K)$ operations

We do this for all states at a given time step. This makes it $O(K^2)$ operations.

Repeating this over the whole sequence (length = N) leads to $O(N K^2)$ operations.



Decoding Problem

Brute force approach $\sim O(K^N)$

What is the complexity of forward probability calculations?

Complexity of forward probability calculations $\sim O(NK^2)$



Decoding Problem

Brute force approach $\sim O(K^N)$

What is the complexity of forward probability calculations?

Complexity of forward probability calculations $\sim O(NK^2)$

We have drastically reduced number of operations:
exponential \rightarrow linear



Backward Probability

$\text{backward}(i, c_l) := P(X_{i+1} = x_{i+1}, \dots, X_N = x_N \mid Y_i = c_l)$

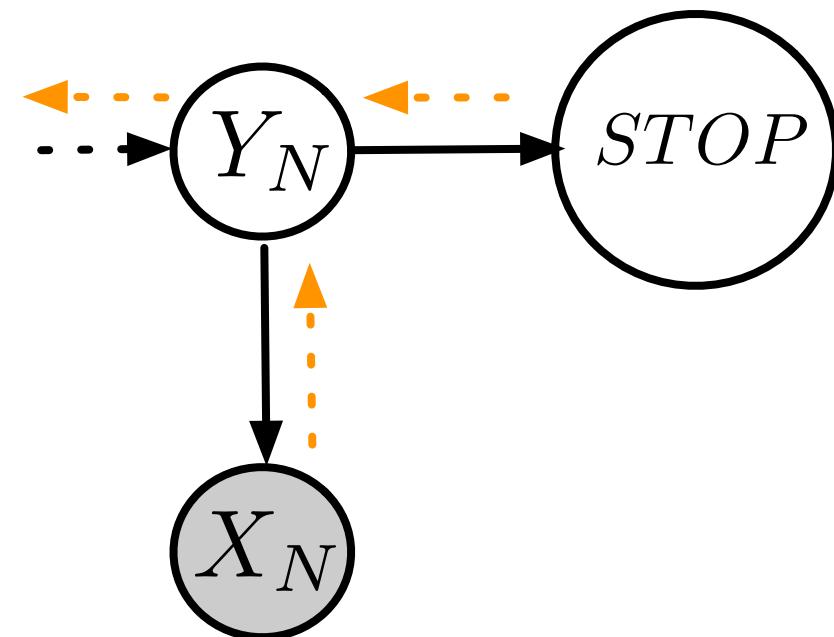
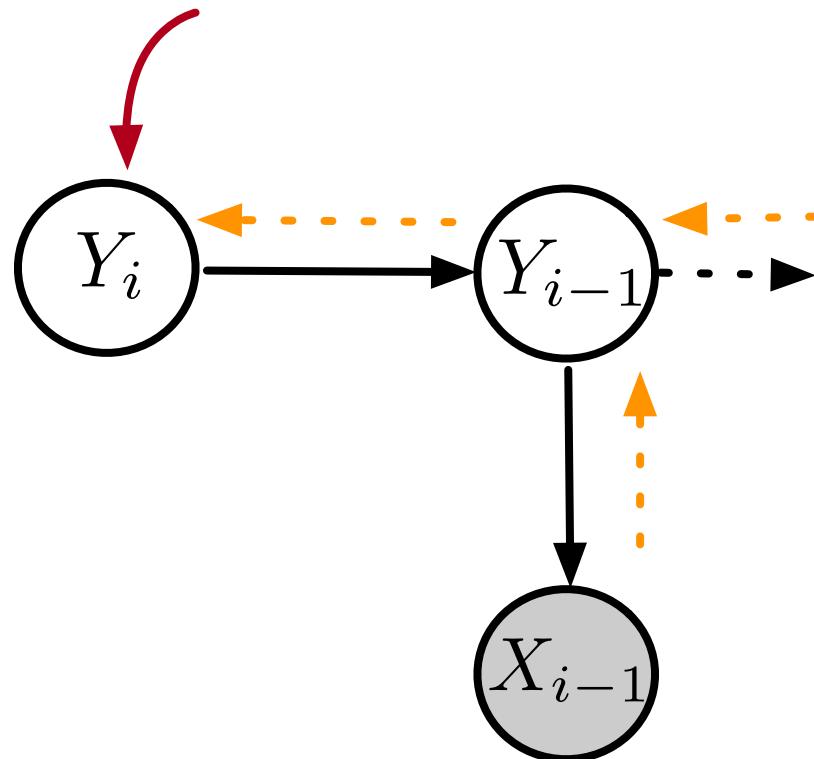


Backward Probability

$$\text{backward}(i, c_l) := P(X_{i+1} = x_{i+1}, \dots, X_N = x_N \mid Y_i = c_l)$$

Backward Probability

$$P(X_{i+1} = x_{i+1}, \dots, X_N = x_N \mid Y_i = c_l)$$



Backward Probability

$$\text{backward}(i, c_l) := P(X_{i+1} = x_{i+1}, \dots, X_N = x_N \mid Y_i = c_l)$$

$$\text{backward}(i, c_l) = \sum_{c_k} P(X_{i+1} = x_{i+1} \mid Y_{i+1} = c_k) \times \text{backward}(i+1, c_k) \times P(Y_{i+1} = c_k \mid Y_i = c_l)$$



Backward Probability

$\text{backward}(i, c_l) := P(X_{i+1} = x_{i+1}, \dots, X_N = x_N \mid Y_i = c_l)$

$P(X_{i+1} = x_{i+1}, \dots, X_N = x_N \mid Y_i = c_l) = P(X_{i+1:N} \mid Y_i)$



Backward Probability

$\text{backward}(i, c_l) := P(X_{i+1} = x_{i+1}, \dots, X_N = x_N \mid Y_i = c_l)$

$P(X_{i+1} = x_{i+1}, \dots, X_N = x_N \mid Y_i = c_l) = P(X_{i+1:N} \mid Y_i)$

$P(X_{i+1:N} \mid Y_i)$



Backward Probability

$$\text{backward}(i, c_l) := P(X_{i+1} = x_{i+1}, \dots, X_N = x_N \mid Y_i = c_l)$$

$$P(X_{i+1} = x_{i+1}, \dots, X_N = x_N \mid Y_i = c_l) = P(X_{i+1:N} \mid Y_i)$$

$$P(X_{i+1:N} \mid Y_i)$$

$$= \sum_{Y_{i+1}} P(X_{i+1:N}, Y_{i+1} \mid Y_i)$$



Backward Probability

$$\text{backward}(i, c_l) := P(X_{i+1} = x_{i+1}, \dots, X_N = x_N \mid Y_i = c_l)$$

$$P(X_{i+1} = x_{i+1}, \dots, X_N = x_N \mid Y_i = c_l) = P(X_{i+1:N} \mid Y_i)$$

$$P(X_{i+1:N} \mid Y_i)$$

$$= \sum_{Y_{i+1}} P(X_{i+1:N}, Y_{i+1} \mid Y_i)$$

$$= \sum_{Y_{i+1}} P(X_{i+1}, Y_{i+1}, X_{i+2:N} \mid Y_i)$$



Backward Probability

$$\text{backward}(i, c_l) := P(X_{i+1} = x_{i+1}, \dots, X_N = x_N \mid Y_i = c_l)$$

$$P(X_{i+1} = x_{i+1}, \dots, X_N = x_N \mid Y_i = c_l) = P(X_{i+1:N} \mid Y_i)$$

$$P(X_{i+1:N} \mid Y_i)$$

$$= \sum_{Y_{i+1}} P(X_{i+1:N}, Y_{i+1} \mid Y_i)$$

$$= \sum_{Y_{i+1}} P(X_{i+1}, Y_{i+1}, X_{i+2:N} \mid Y_i)$$

$$= \sum_{Y_{i+1}} P(X_{i+1} \mid Y_{i+1}, X_{i+2:N}, Y_i) \times P(Y_{i+1}, X_{i+2:N}, Y_i)$$



Backward Probability

$$\text{backward}(i, c_l) := P(X_{i+1} = x_{i+1}, \dots, X_N = x_N \mid Y_i = c_l)$$

$$P(X_{i+1} = x_{i+1}, \dots, X_N = x_N \mid Y_i = c_l) = P(X_{i+1:N} \mid Y_i)$$

$$P(X_{i+1:N} \mid Y_i)$$

$$= \sum_{Y_{i+1}} P(X_{i+1:N}, Y_{i+1} \mid Y_i)$$

$$= \sum_{Y_{i+1}} P(X_{i+1}, Y_{i+1}, X_{i+2:N} \mid Y_i)$$

$$= \sum_{Y_{i+1}} P(X_{i+1} \mid Y_{i+1}, X_{i+2:N}, Y_i) \times P(Y_{i+1}, X_{i+2:N} \mid Y_i)$$

$$= \sum_{Y_{i+1}} P(X_{i+1} \mid Y_{i+1}) \times P(X_{i+2:N} \mid Y_{i+1}) \times P(Y_{i+1} \mid Y_i)$$



Backward Probability

$$\text{backward}(i, c_l) := P(X_{i+1} = x_{i+1}, \dots, X_N = x_N \mid Y_i = c_l)$$

$$\text{backward}(i, c_l) = \sum_{c_k} P(X_{i+1} = x_{i+1} \mid Y_{i+1} = c_k) \times \text{backward}(i + 1, c_k) \times P(Y_{i+1} = c_k \mid Y_i = c_l)$$



Backward Probability

$$\text{backward}(i, c_l) := P(X_{i+1} = x_{i+1}, \dots, X_N = x_N \mid Y_i = c_l)$$

$$\text{backward}(i, c_l) = \sum_{c_k} P(X_{i+1} = x_{i+1} \mid Y_{i+1} = c_k) \times \text{backward}(i + 1, c_k) \times P(Y_{i+1} = c_k \mid Y_i = c_l)$$

Emission
Probability

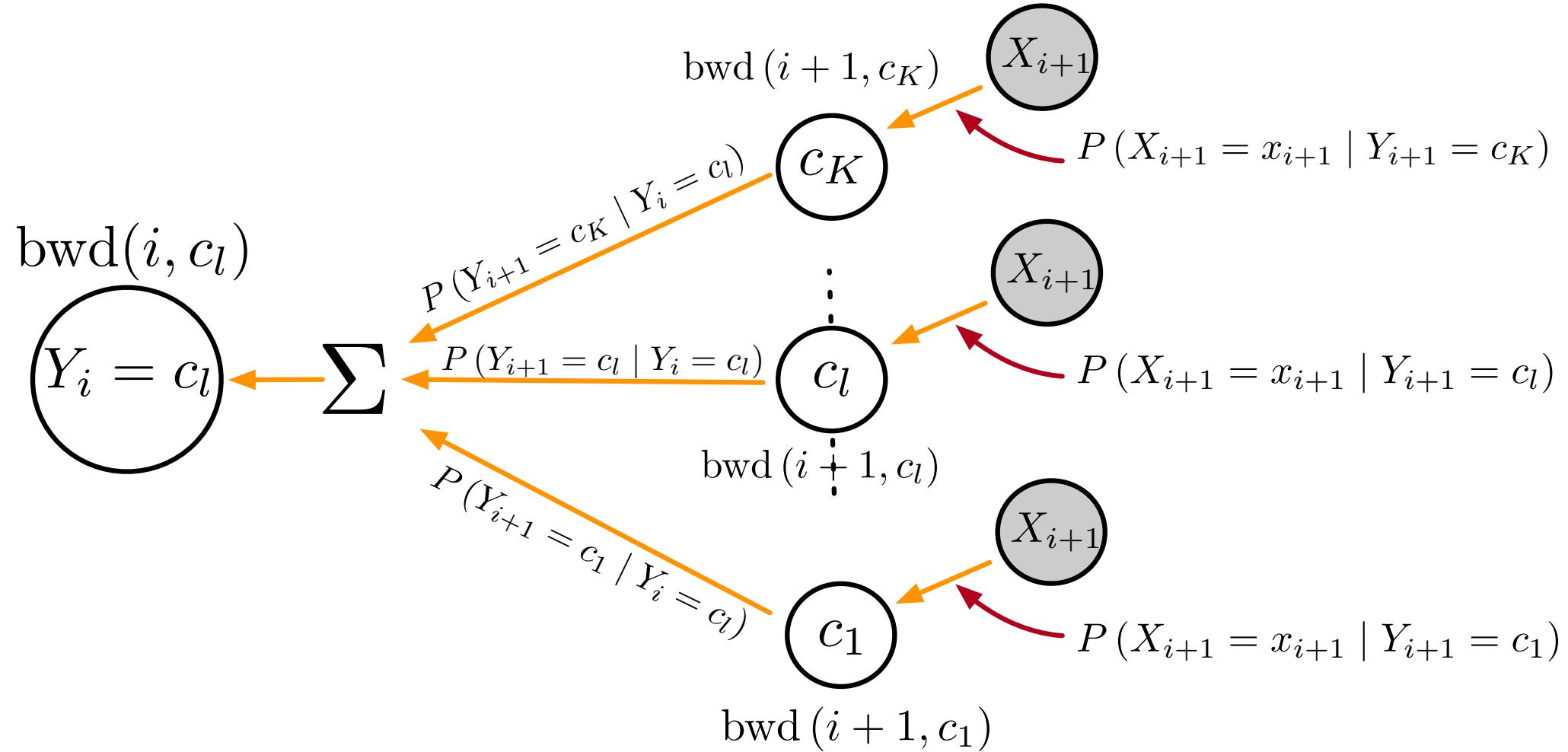


Transition
Probability



Backward Probability

$$\text{backward}(i, c_l) = \sum_{c_k} P(X_{i+1} = x_{i+1} \mid Y_{i+1} = c_k) \times \text{backward}(i + 1, c_k) \times P(Y_{i+1} = c_k \mid Y_i = c_l)$$



Backward Probability

$\text{backward}(i, c_l) := P(X_{i+1} = x_{i+1}, \dots, X_N = x_N \mid Y_i = c_l)$

$\text{backward}(i, c_l) = \sum_{c_k} P(X_{i+1} = x_{i+1} \mid Y_{i+1} = c_k) \times \text{backward}(i + 1, c_k) \times P(Y_{i+1} = c_k \mid Y_i = c_l)$

$\text{backward}(0, \text{START}) = ?$



Backward Probability

$\text{backward}(i, c_l) := P(X_{i+1} = x_{i+1}, \dots, X_N = x_N \mid Y_i = c_l)$

$$\text{backward}(i, c_l) = \sum_{c_k} P(X_{i+1} = x_{i+1} \mid Y_{i+1} = c_k) \times \text{backward}(i + 1, c_k) \times P(Y_{i+1} = c_k \mid Y_i = c_l)$$

$$\text{backward}(0, \text{START}) = P(X)$$



Summary so far

$$\text{forward}(i, c_k) := P(Y_i = c_k, X_1 = x_1, \dots, X_i = x_i)$$

$$\text{forward}(i, c_k) = \left(\sum_{c_l} P(Y_i = c_k \mid Y_{i-1} = c_l) \times \text{forward}(i-1, c_l) \right) \times P(X_i \mid Y_i = c_k)$$



Summary so far

$$\text{forward}(i, c_k) := P(Y_i = c_k, X_1 = x_1, \dots, X_i = x_i)$$

$$\text{forward}(i, c_k) = \left(\sum_{c_l} P(Y_i = c_k \mid Y_{i-1} = c_l) \times \text{forward}(i-1, c_l) \right) \times P(X_i \mid Y_i = c_k)$$

$$\text{backward}(i, c_l) := P(X_{i+1} = x_{i+1}, \dots, X_N = x_N \mid Y_i = c_l)$$

$$\text{backward}(i, c_l) = \sum_{c_k} P(X_{i+1} = x_{i+1} \mid Y_{i+1} = c_k) \times \text{backward}(i+1, c_k) \times P(Y_{i+1} = c_k \mid Y_i = c_l)$$



Summary so far

$$\text{forward}(i, c_k) := P(Y_i = c_k, X_1 = x_1, \dots, X_i = x_i)$$

$$\text{forward}(i, c_k) = \left(\sum_{c_l} P(Y_i = c_k \mid Y_{i-1} = c_l) \times \text{forward}(i-1, c_l) \right) \times P(X_i \mid Y_i = c_k)$$

$$\text{backward}(i, c_l) := P(X_{i+1} = x_{i+1}, \dots, X_N = x_N \mid Y_i = c_l)$$

$$\text{backward}(i, c_l) = \sum_{c_k} P(X_{i+1} = x_{i+1} \mid Y_{i+1} = c_k) \times \text{backward}(i+1, c_k) \times P(Y_{i+1} = c_k \mid Y_i = c_l)$$

$$P(X) = \text{forward}(N+1, STOP)$$

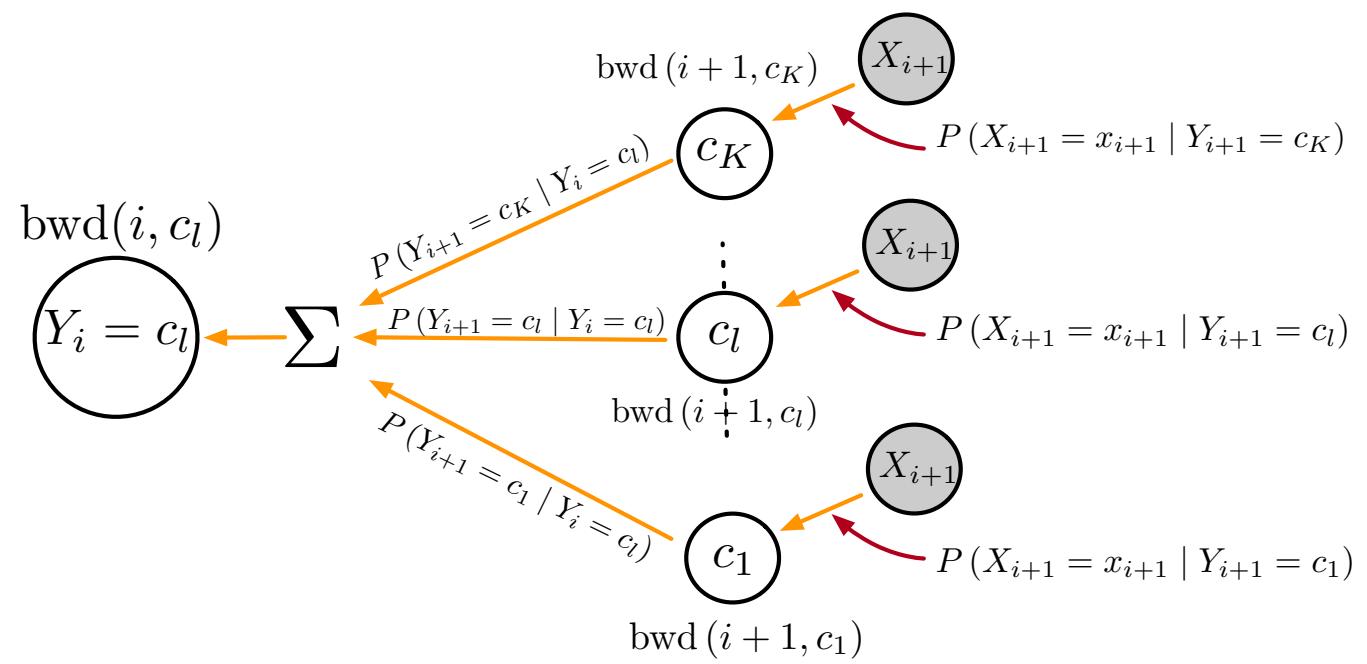
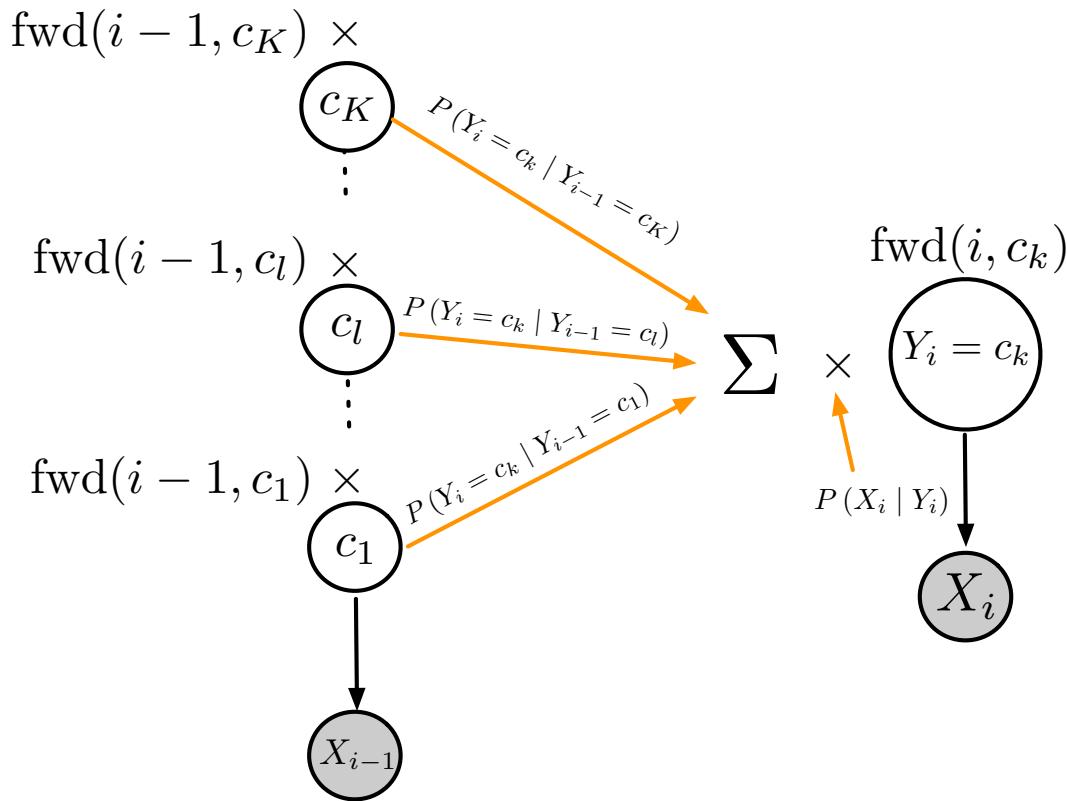
$$\text{backward}(0, START) = P(X)$$



Summary so far

forward $(i, c_k) := P(Y_i = c_k, X_1 = x_1, \dots, X_i = x_i)$

backward $(i, c_l) := P(X_{i+1} = x_{i+1}, \dots, X_N = x_N \mid Y_i = c_l)$



Forward Backward

$$P(X) = P(X_1, \dots, X_i, X_{i+1}, \dots, X_N)$$



Forward Backward

$$\begin{aligned} P(X) &= P(X_1, \dots, X_i, X_{i+1}, \dots, X_N) \\ &= \sum_{c_l} P(X_1, \dots, X_i, X_{i+1}, \dots, X_N, Y_i = c_l) \end{aligned}$$



Forward Backward

$$\begin{aligned} P(X) &= P(X_1, \dots, X_i, X_{i+1}, \dots, X_N) \\ &= \sum P(X_1, \dots, X_i, X_{i+1}, \dots, X_N, Y_i = c_l) \\ &= \sum_{c_l} P(X_{i+1}, \dots, X_N \mid Y_i = c_l) \times P(Y_i = c_l, X_1, \dots, X_i) \end{aligned}$$



Forward Backward

$$\begin{aligned} P(X) &= P(X_1, \dots, X_i, X_{i+1}, \dots, X_N) \\ &= \sum_{c_l} P(X_1, \dots, X_i, X_{i+1}, \dots, X_N, Y_i = c_l) \\ &= \sum_{c_l} P(X_{i+1}, \dots, X_N \mid Y_i = c_l) \times P(Y_i = c_l, X_1, \dots, X_i) \\ &= \sum_{c_l} \text{backward}(i, c_l) \times \text{forward}(i, c_l) \end{aligned}$$



Forward Backward

$$\begin{aligned} P(X) &= P(X_1, \dots, X_i, X_{i+1}, \dots, X_N) \\ &= \sum_{c_l} P(X_1, \dots, X_i, X_{i+1}, \dots, X_N, Y_i = c_l) \\ &= \sum_{c_l} P(X_{i+1}, \dots, X_N \mid Y_i = c_l) \times P(Y_i = c_l, X_1, \dots, X_i) \\ &= \sum_{c_l} \text{backward}(i, c_l) \times \text{forward}(i, c_l) \end{aligned}$$

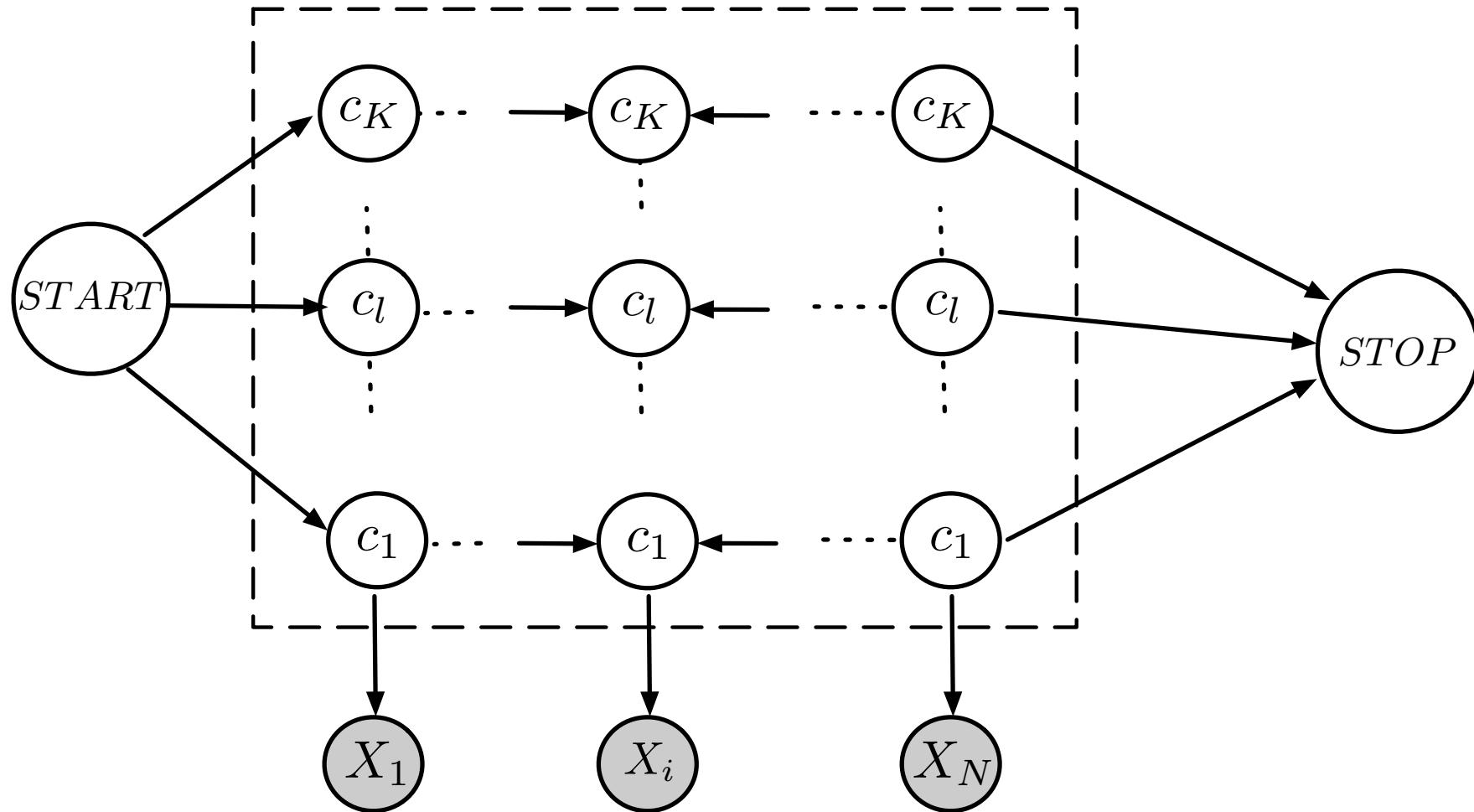
Does not depend on i

$$P(X) = \sum_{c_l} \text{backward}(i, c_l) \times \text{forward}(i, c_l)$$



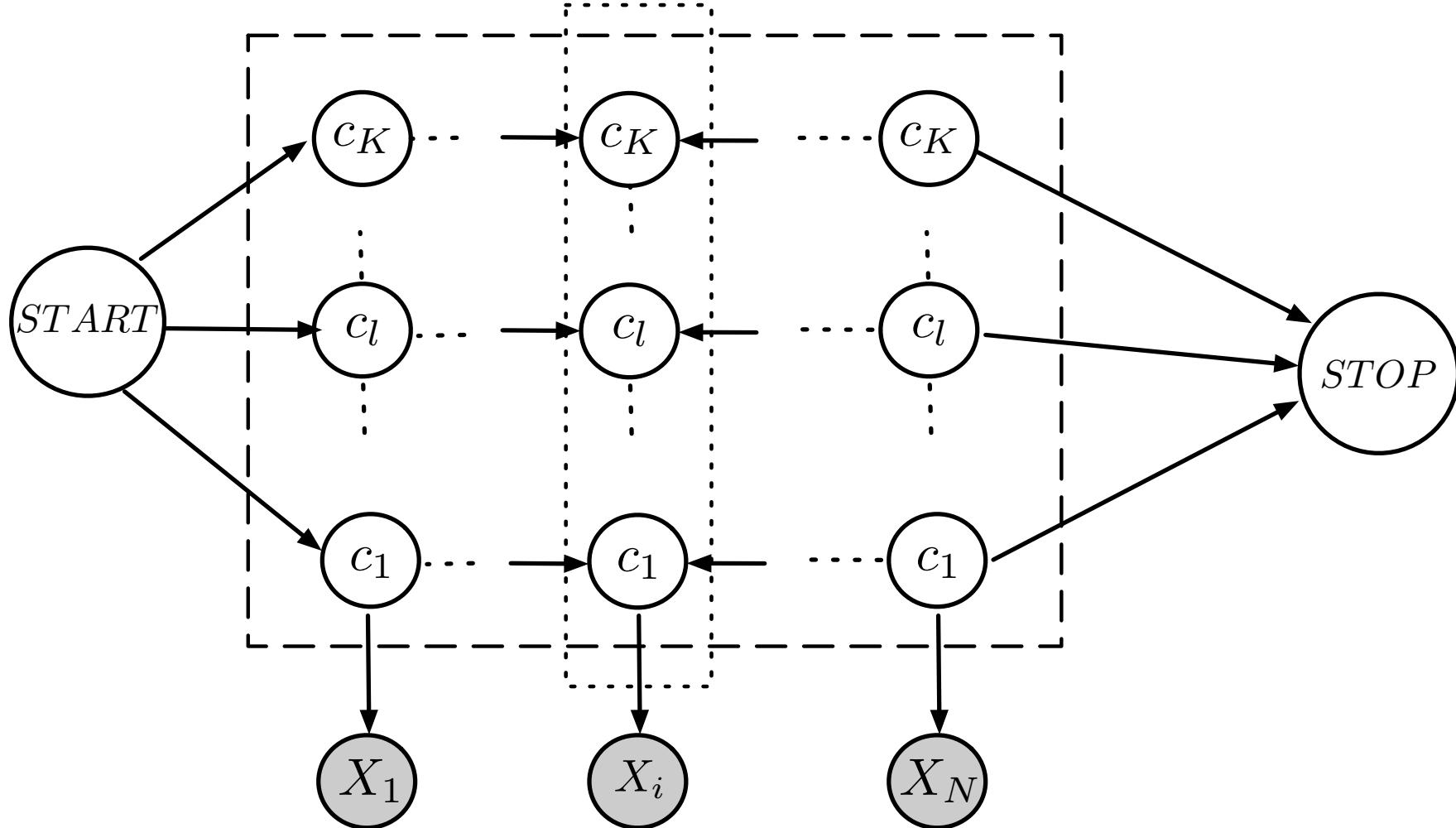
Forward Backward

$$P(X) = \sum_{c_l} \text{backward}(i, c_l) \times \text{forward}(i, c_l)$$



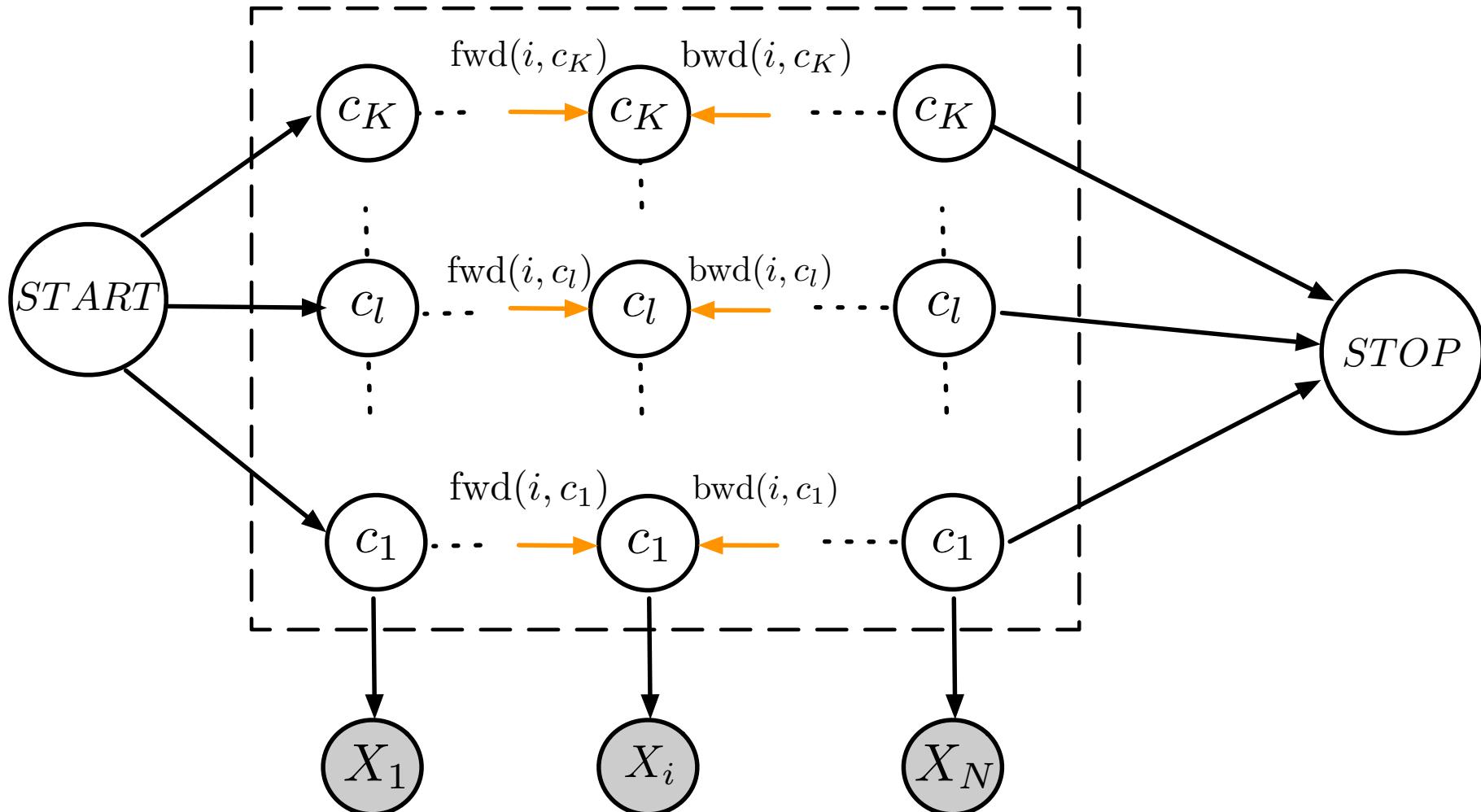
Forward Backward

$$P(X) = \sum_{c_l} \text{backward } (i, c_l) \times \text{forward } (i, c_l)$$



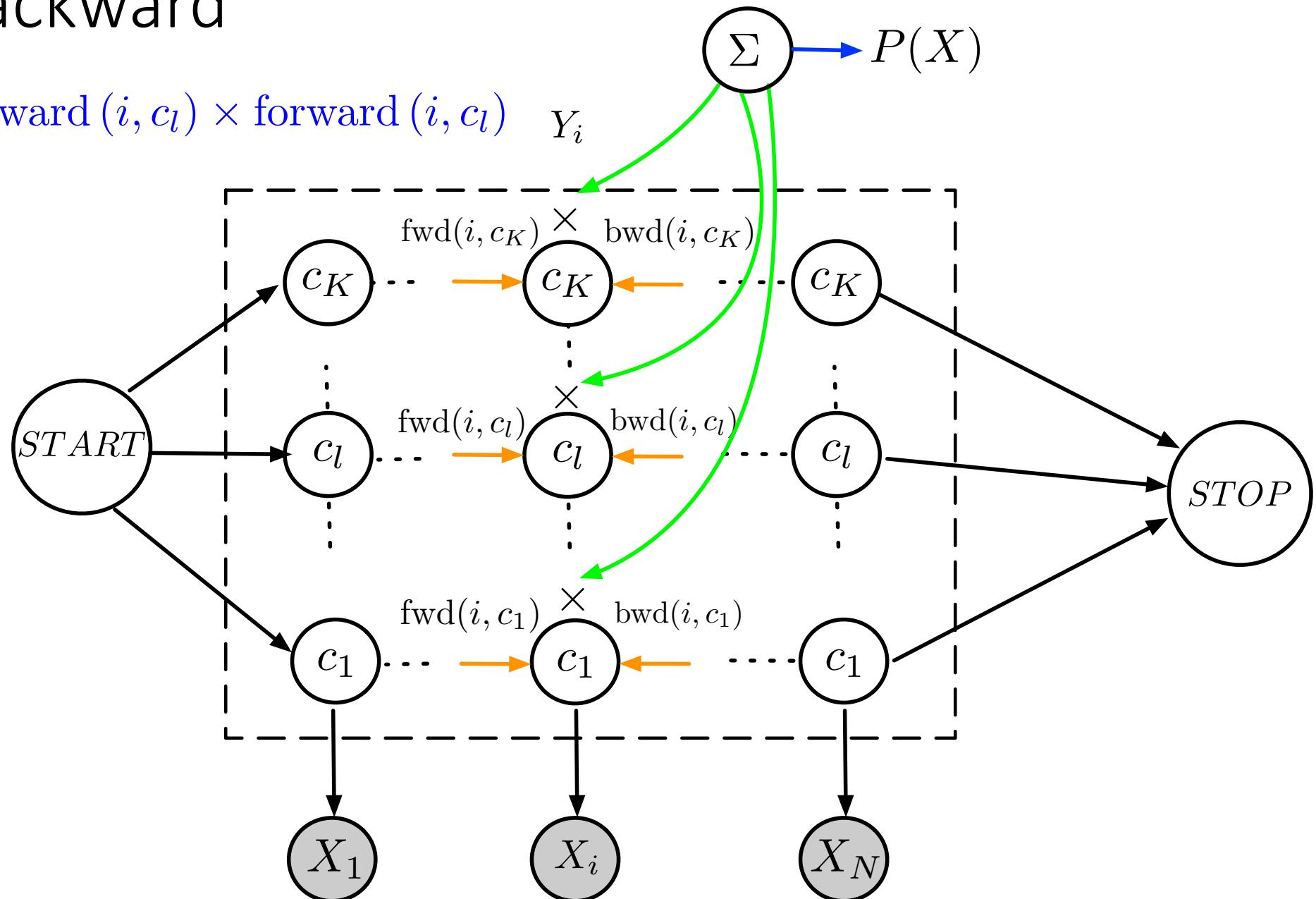
Forward Backward

$$P(X) = \sum_{c_l} \text{backward}(i, c_l) \times \text{forward}(i, c_l) Y_i$$



Forward Backward

$$P(X) = \sum_{c_l} \text{backward}(i, c_l) \times \text{forward}(i, c_l)$$



Computation

Brute force approach $\sim O(K^N)$

What is the complexity of $P(X)$ calculation?



Computation

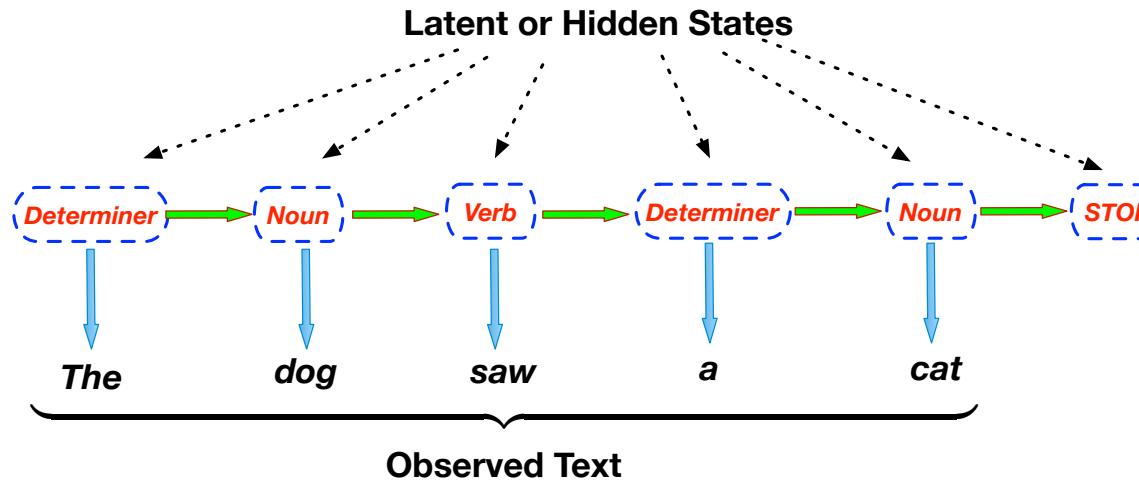
Brute force approach $\sim O(K^N)$

What is the complexity of $P(X)$ calculation?

Complexity of $P(X)$ calculation: $\sim O(NK^2)$



Hidden Markov Models (HMM)



Three Questions:

1. **Learning Problem:** How do we estimate the parameters of distributions?
2. **Decoding Problem:** Given the observed text what is the hidden POS sequence that best explains the observation?
3. **Likelihood Problem:** Given the parameters of the distributions what is the probability of the observed sequence?



Decoding Problem

- Two ways possible

- **Posterior Decoding:** Maximize probability of each state

$$y_i^* = \arg \max_{c_l} P(Y_i = c_l | X_1 = x_1, \dots, X_N = x_N)$$



Decoding Problem

- Two ways possible

- **Posterior Decoding:** Maximize probability of each state

$$y_i^* = \arg \max_{c_l} P(Y_i = c_l | X_1 = x_1, \dots, X_N = x_N)$$

- **Viterbi Decoding:** Maximize the probability of the entire sequence

$$\begin{aligned} y^* &= \arg \max_{y=y_1 \dots y_N} P(Y_1 = y_1, \dots, Y_N = y_N | X_1 = x_1, \dots, X_N = x_N) \\ &= \arg \max_{y=y_1 \dots y_N} P(Y_1 = y_1, \dots, Y_N = y_N, X_1 = x_1, \dots, X_N = x_N) \end{aligned}$$



Posterior Decoding

$$y_i^* = \arg \max_{c_l} P(Y_i = c_l | X_1 = x_1, \dots, X_N = x_N)$$



Posterior Decoding

$$\begin{aligned}y_i^* &= \arg \max_{c_l} P(Y_i = c_l | X_1 = x_1, \dots, X_N = x_N) \\&= \arg \max_{c_l} \frac{P(Y_i = c_l, X_1 = x_1, \dots, X_N = x_N)}{P(X)}\end{aligned}$$



Posterior Decoding

$$y_i^* = \arg \max_{c_l} P(Y_i = c_l | X_1 = x_1, \dots, X_N = x_N)$$

$$= \arg \max_{c_l} \frac{P(Y_i = c_l, X_1 = x_1, \dots, X_N = x_N)}{P(X)}$$

$$= \arg \max_{c_l} \frac{P(X_N = x_n, \dots, X_{i+1} = x_{i+1} | Y_i = c_l) \times (Y_i = c_l, X_1 = x_1, \dots, X_i = x_i)}{P(X)}$$



Posterior Decoding

$$y_i^* = \arg \max_{c_l} P(Y_i = c_l | X_1 = x_1, \dots, X_N = x_N)$$

$$= \arg \max_{c_l} \frac{P(Y_i = c_l, X_1 = x_1, \dots, X_N = x_N)}{P(X)}$$

$$= \arg \max_{c_l} \frac{P(X_N = x_n, \dots, X_{i+1} = x_{i+1} | Y_i = c_l) \times (Y_i = c_l, X_1 = x_1, \dots, X_i = x_i)}{P(X)}$$

$$= \arg \max_{c_l} \frac{\text{forward}(i, c_l) \times \text{backward}(i, c_l)}{P(X)}$$



Posterior Decoding

$$y_i^* = \arg \max_{c_l} P(Y_i = c_l | X_1 = x_1, \dots, X_N = x_N)$$

$$= \arg \max_{c_l} \frac{P(Y_i = c_l, X_1 = x_1, \dots, X_N = x_N)}{P(X)}$$

$$= \arg \max_{c_l} \frac{P(X_N = x_n, \dots, X_{i+1} = x_{i+1} | Y_i = c_l) \times (Y_i = c_l, X_1 = x_1, \dots, X_i = x_i)}{P(X)}$$

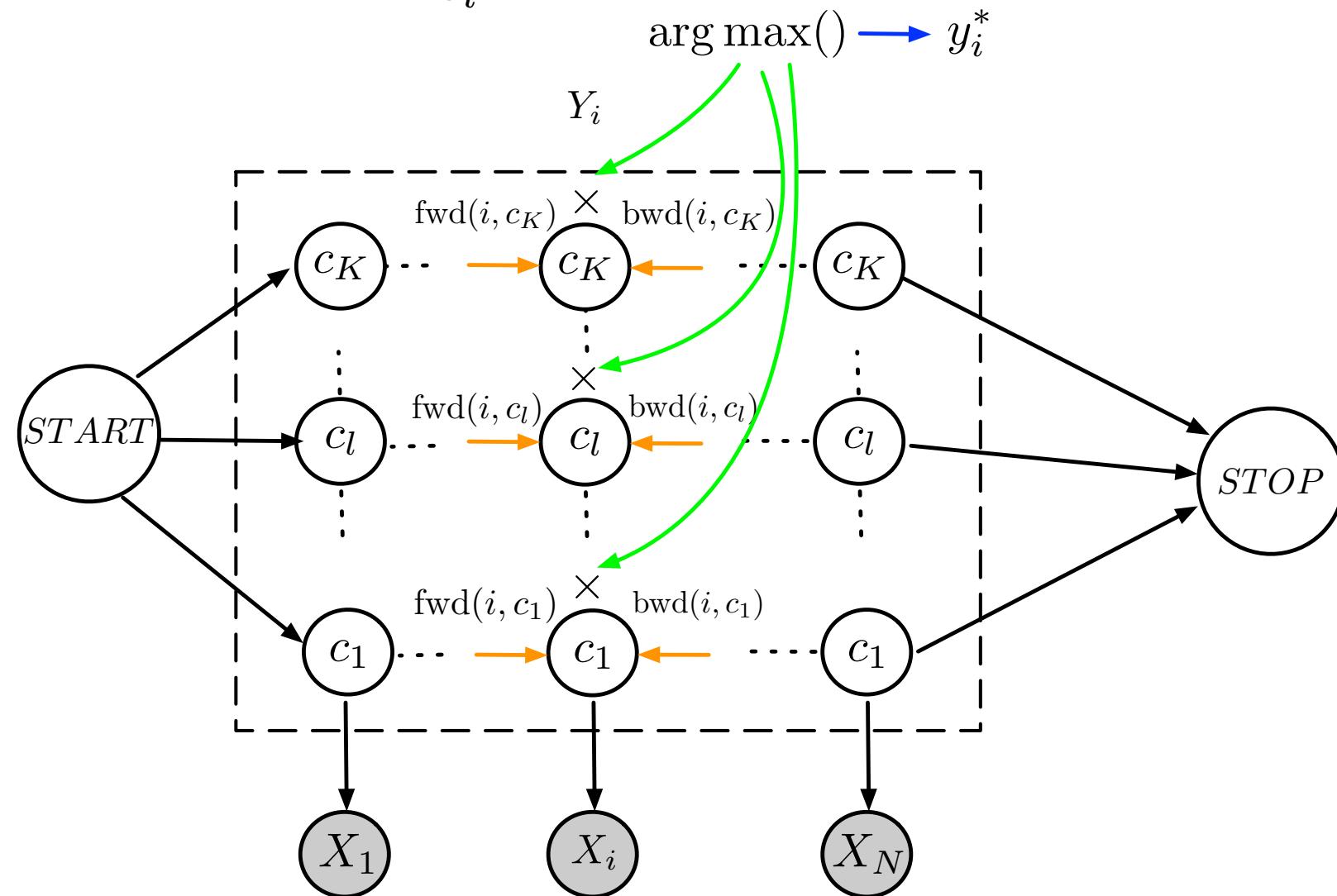
$$= \arg \max_{c_l} \frac{\text{forward}(i, c_l) \times \text{backward}(i, c_l)}{P(X)}$$

$$y_i^* = \arg \max_{c_l} \text{forward}(i, c_l) \times \text{backward}(i, c_l)$$



Posterior Decoding

$$y_i^* = \arg \max_{c_l} \text{forward}(i, c_l) \times \text{backward}(i, c_l)$$



References

1. Michael Collin's NLP Lecture Notes:
<http://www.cs.columbia.edu/~mcollins/hmms-spring2013.pdf>
2. Chapter 6, Speech and Language Processing, Dan Jurafsky and James Martin
3. LxMLS Lab Guide: <http://lxmbs.it.pt/2016/LxMLS2016.pdf>



QUIZ

- Suppose, you are given data about the activities that some person did over different days.
- We know that the activity that a person does on a given day is dependent on the weather on that day.
- We would like to predict the most likely weather pattern by observing the activity pattern.
- Just to simplify our problem, let us assume:
 - We have only 3 different activities: STUDY, SHOP, PLAY
 - We have only 2 possible weather conditions: SUNNY, RAIN



QUIZ

	STUDY (x_1)	SHOP (x_2)	PLAY (x_3)
SUNNY (c_1)	0.2	0.4	0.4
RAIN (c_2)	0.5	0.4	0.1

	SUNNY (c_1)	RAIN (c_2)
SUNNY (c_1)	0.6	0.4
RAIN (c_2)	0.5	0.5

	SUNNY (c_1)	RAIN (c_2)
START	0.8	0.2

QUIZ

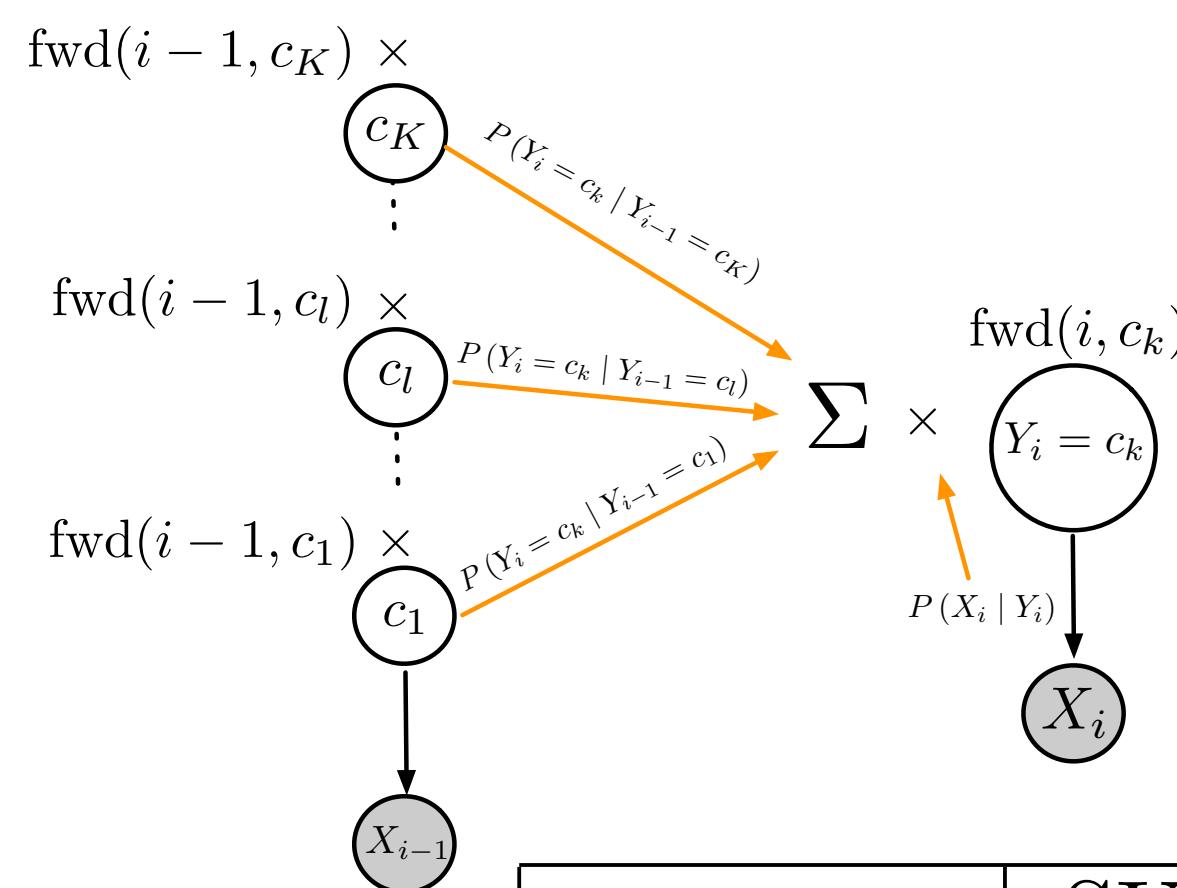
	STUDY (x_1)	SHOP (x_2)	PLAY (x_3)
SUNNY (c_1)	0.2	0.4	0.4
RAIN (c_2)	0.5	0.4	0.1

	SUNNY (c_1)	RAIN (c_2)
SUNNY (c_1)	0.6	0.4
RAIN (c_2)	0.5	0.5

	SUNNY (c_1)	RAIN (c_2)
START	0.8	0.2

Calculate the forward probabilities for the sequence PLAY, STUDY, PLAY





	STUDY (x_1)	SHOP (x_2)	PLAY (x_3)
SUNNY (c_1)	0.2	0.4	0.4
RAIN (c_2)	0.5	0.4	0.1

	SUNNY (c_1)	RAIN (c_2)
SUNNY (c_1)	0.6	0.4
RAIN (c_2)	0.5	0.5

	SUNNY (c_1)	RAIN (c_2)
START	0.8	0.2

Calculate the forward probabilities for the sequence: PLAY, STUDY, PLAY

	SUNNY (c_1)	RAIN (c_2)
fwd(1, c_l)	?	?
fwd(2, c_l)	?	?
fwd(3, c_l)	?	?