



# 10-601 Introduction to Machine Learning

Machine Learning Department  
School of Computer Science  
Carnegie Mellon University

## Hidden Markov Models

Matt Gormley  
Lecture 20  
Nov. 7, 2018

# Reminders

- **Homework 6: PAC Learning / Generative Models**
  - Out: Wed, Oct 31
  - Due: Wed, Nov 7 at 11:59pm (1 week)
- **Homework 7: HMMs**
  - Out: Wed, Nov 7
  - Due: Mon, Nov 19 at 11:59pm

# HMM Outline

- **Motivation**
  - Time Series Data
- **Hidden Markov Model (HMM)**
  - Example: Squirrel Hill Tunnel Closures [courtesy of Roni Rosenfeld]
  - Background: Markov Models
  - From Mixture Model to HMM
  - History of HMMs
  - Higher-order HMMs
- **Training HMMs**
  - (Supervised) Likelihood for HMM
  - Maximum Likelihood Estimation (MLE) for HMM
  - EM for HMM (aka. Baum-Welch algorithm)
- **Forward-Backward Algorithm**
  - Three Inference Problems for HMM
  - Great Ideas in ML: Message Passing
  - Example: Forward-Backward on 3-word Sentence
  - Derivation of Forward Algorithm
  - Forward-Backward Algorithm
  - Viterbi algorithm



Last Lecture



This Lecture

# **SUPERVISED LEARNING FOR HMMS**

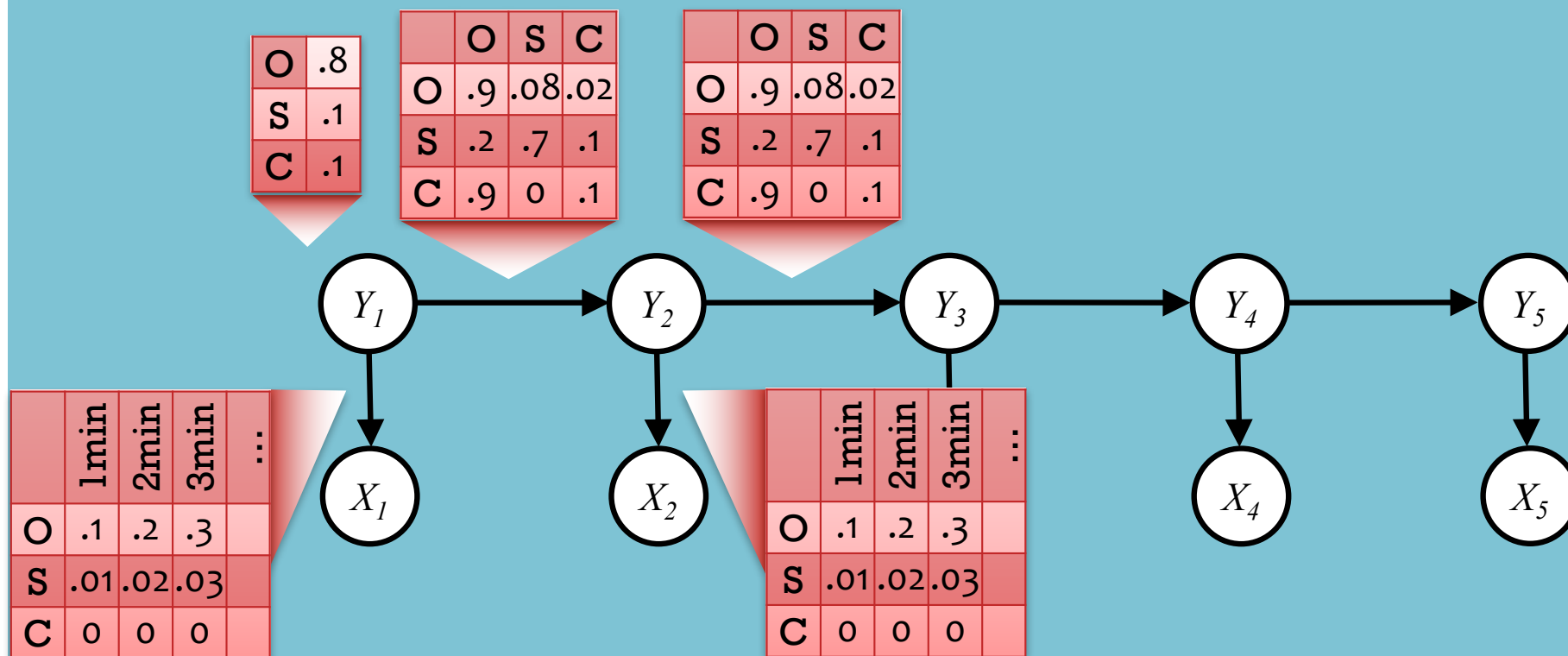
# Hidden Markov Model

## HMM Parameters:

Emission matrix,  $\mathbf{A}$ , where  $P(X_t = k | Y_t = j) = A_{j,k}, \forall t, k$

Transition matrix,  $\mathbf{B}$ , where  $P(Y_t = k | Y_{t-1} = j) = B_{j,k}, \forall t, k$

Initial probs,  $\mathbf{C}$ , where  $P(Y_1 = k) = C_k, \forall k$



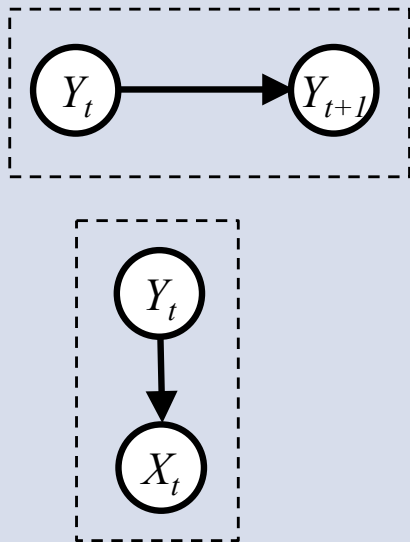
# Training HMMs

## *Whiteboard*

- (Supervised) Likelihood for an HMM
- Maximum Likelihood Estimation (MLE) for HMM

# Supervised Learning for HMMs

Learning an HMM decomposes into solving two (independent) Mixture Models



Data:  $D = \{(\vec{x}^{(i)}, \vec{y}^{(i)})\}_{i=1}^N$        $\vec{x} = [x_1, \dots, x_T]^T$   
Likelihood:  $\vec{y} = [y_1, \dots, y_T]^T$

$$\ell(A, B, C) = \sum_{i=1}^N \log p(\vec{x}^{(i)}, y^{(i)} | A, B, C)$$

$$= \sum_{i=1}^N \left[ \underbrace{\log p(y_1^{(i)} | C)}_{\text{initial}} + \underbrace{\left( \sum_{t=2}^T \log p(y_t^{(i)} | y_{t-1}^{(i)}, B) \right)}_{\text{transition}} + \underbrace{\left( \sum_{t=1}^T \log p(x_t^{(i)} | y_t^{(i)}, A) \right)}_{\text{emission}} \right]$$

MLE:

$$\hat{A}, \hat{B}, \hat{C} = \underset{A, B, C}{\operatorname{argmax}} \ell(A, B, C)$$

$$\Rightarrow \hat{C} = \underset{C}{\operatorname{argmax}} \sum_{i=1}^N \log p(y_1^{(i)} | C)$$

$$\hat{B} = \underset{B}{\operatorname{argmax}} \sum_{i=1}^N \sum_{t=2}^T \log p(y_t^{(i)} | y_{t-1}^{(i)}, B)$$

$$\hat{A} = \underset{A}{\operatorname{argmax}} \sum_{i=1}^N \sum_{t=1}^T \log p(x_t^{(i)} | y_t^{(i)}, A)$$

Can solve in closed form, which yields...

$$\hat{C}_k = \frac{\#(y_1^{(i)} = k)}{N} \quad \forall i, k$$

$$\hat{B}_{jk} = \frac{\#(y_t^{(i)} = k \text{ and } y_{t-1}^{(i)} = j)}{\#(y_{t-1}^{(i)} = j)} \quad \forall i, t > 1, j, k$$

$$\hat{A}_{jk} = \frac{\#(x_t^{(i)} = k \text{ and } y_t^{(i)} = j)}{\#(y_t^{(i)} = j)} \quad \forall i, t, j, k$$

# Hidden Markov Model

## HMM Parameters:

Emission matrix,  $\mathbf{A}$ , where  $P(X_t = k | Y_t = j) = A_{j,k}, \forall t, k$

Transition matrix,  $\mathbf{B}$ , where  $P(Y_t = k | Y_{t-1} = j) = B_{j,k}, \forall t, k$

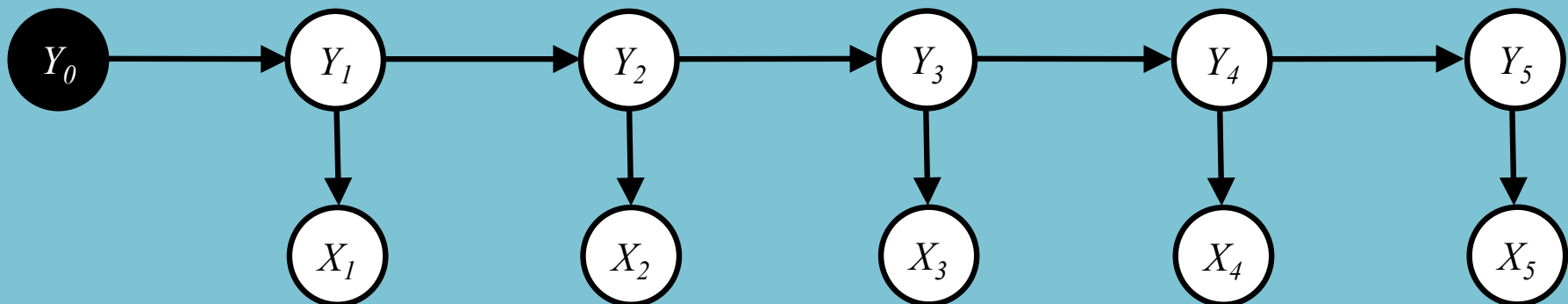
**Assumption:**  $y_0 = \text{START}$

**Generative Story:**

$$Y_t \sim \text{Multinomial}(\mathbf{B}_{Y_{t-1}}) \quad \forall t$$

$$X_t \sim \text{Multinomial}(\mathbf{A}_{Y_t}) \quad \forall t$$

For notational convenience, we fold the *initial probabilities*  $\mathbf{C}$  into the *transition matrix*  $\mathbf{B}$  by our assumption.



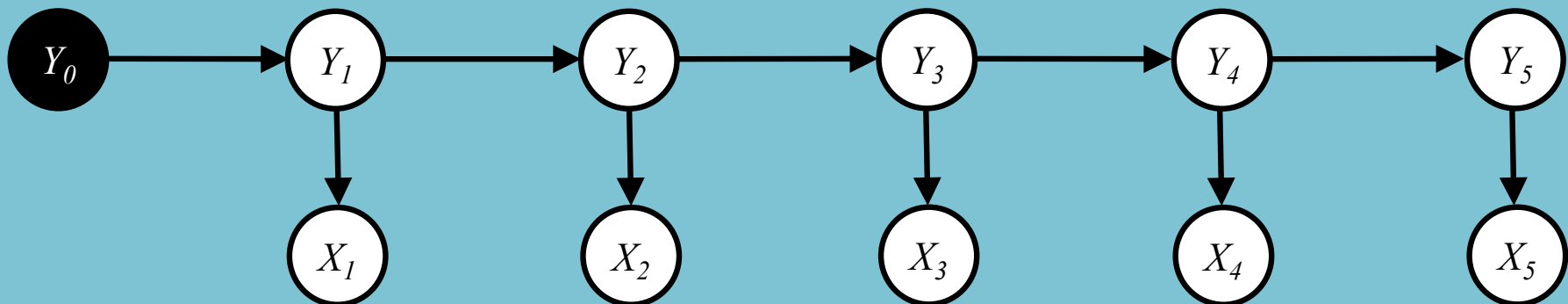


# Hidden Markov Model

**Joint Distribution:**

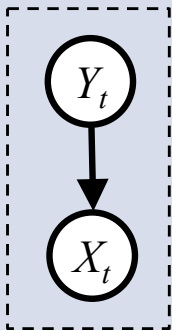
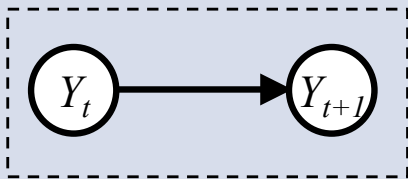
$y_0 = \text{START}$

$$\begin{aligned} p(\mathbf{x}, \mathbf{y} | y_0) &= \prod_{t=1}^T p(x_t | y_t) p(y_t | y_{t-1}) \\ &= \prod_{t=1}^T A_{y_t, x_t} B_{y_{t-1}, y_t} \end{aligned}$$



# Supervised Learning for HMMs

Learning an HMM decomposes into solving two (independent) Mixture Models



$$D = \{(\vec{x}^{(i)}, \vec{y}^{(i)})\}_{i=1}^N$$

Likelihood:  $\ell(A, B) = \sum_{i=1}^N \log p(\vec{x}^{(i)}, \vec{y}^{(i)})$

$$= \sum_{i=1}^N \left[ \sum_{t=1}^T \log p(y_t^{(i)} | y_{t-1}^{(i)}, B) + \log p(x_t^{(i)} | y_t^{(i)}, A) \right]$$

MLE:  $\hat{A}, \hat{B} = \operatorname{argmax} \ell(A, B)$

$$\hat{A} = \operatorname{argmax} \sum_{i=1}^N \left[ \sum_{t=1}^T \log p(x_t^{(i)} | y_t^{(i)}, A) \right]$$

$$\hat{B} = \operatorname{argmax} \sum_{i=1}^N \left[ \sum_{t=1}^T \log p(y_t^{(i)} | y_{t-1}^{(i)}, B) \right]$$

↑ can solve in closed form to get...

$$\hat{B}_{jk} = \frac{\#(y_t^{(i)} = k \text{ and } y_{t-1}^{(i)} = j)}{\#(y_{t-1}^{(i)} = j)}$$

$$\hat{A}_{jk} = \frac{\#(x_t^{(i)} = k \text{ and } y_t^{(i)} = j)}{\#(y_t^{(i)} = j)}$$

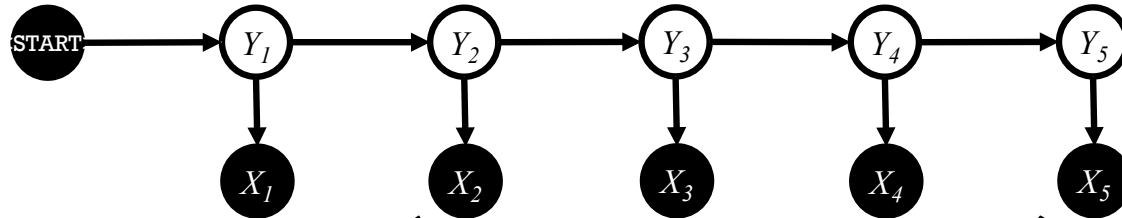
# HMMs: History

- Markov chains: Andrey Markov (1906)
  - Random walks and Brownian motion
- Used in Shannon's work on information theory (1948)
- Baum-Welsh learning algorithm: late 60's, early 70's.
  - Used mainly for speech in 60s-70s.
- Late 80's and 90's: David Haussler (major player in learning theory in 80's) began to use HMMs for modeling biological sequences
- Mid-late 1990's: Dayne Freitag/Andrew McCallum
  - Freitag thesis with Tom Mitchell on IE from Web using logic programs, grammar induction, etc.
  - McCallum: multinomial Naïve Bayes for text
  - With McCallum, IE using HMMs on CORA
- ...

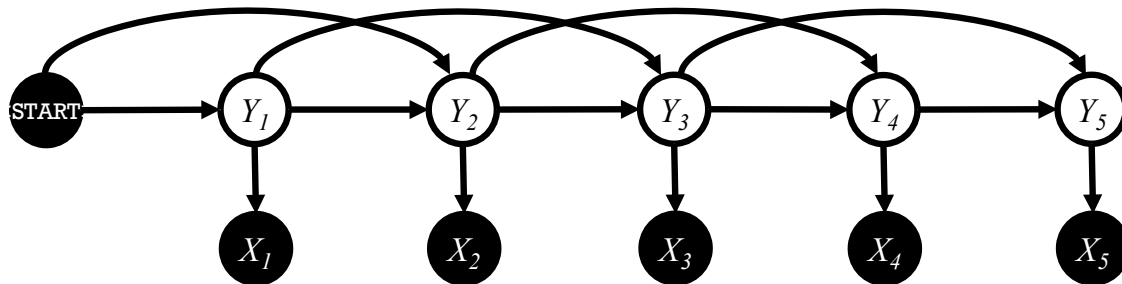


# Higher-order HMMs

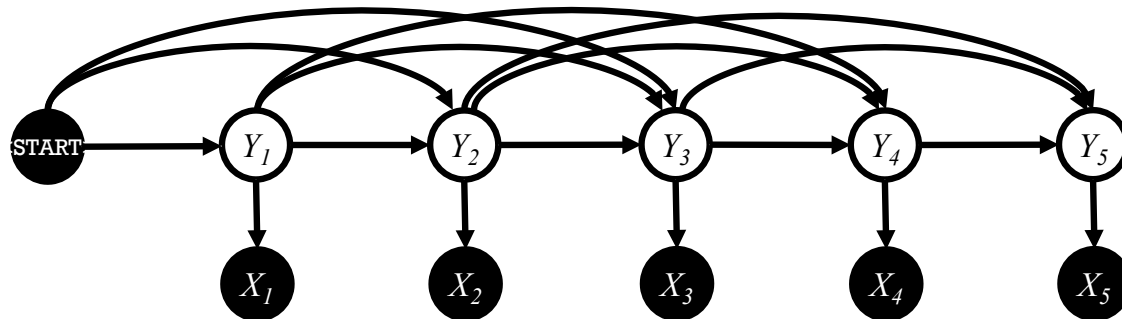
- 1<sup>st</sup>-order HMM (i.e. bigram HMM)



- 2<sup>nd</sup>-order HMM (i.e. trigram HMM)



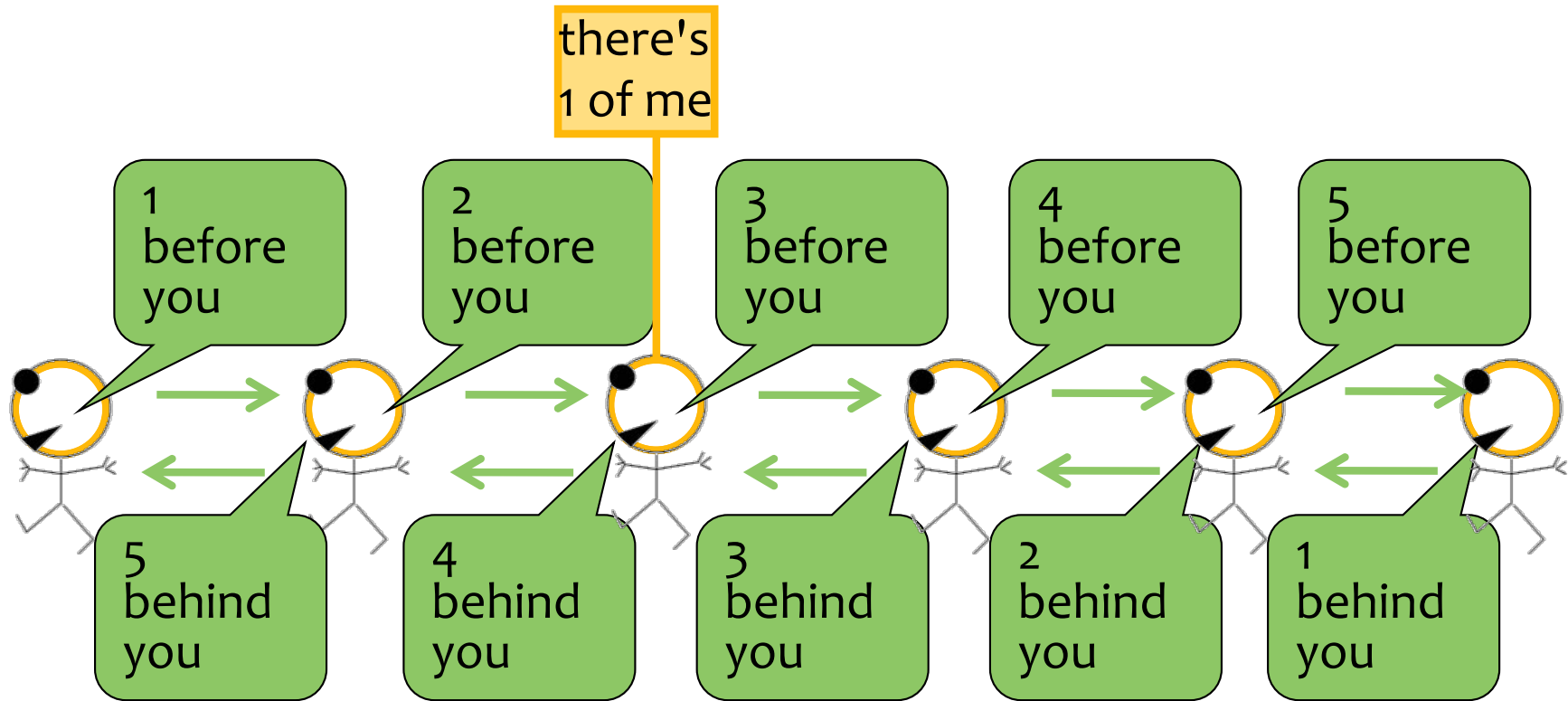
- 3<sup>rd</sup>-order HMM



# **BACKGROUND: MESSAGE PASSING**

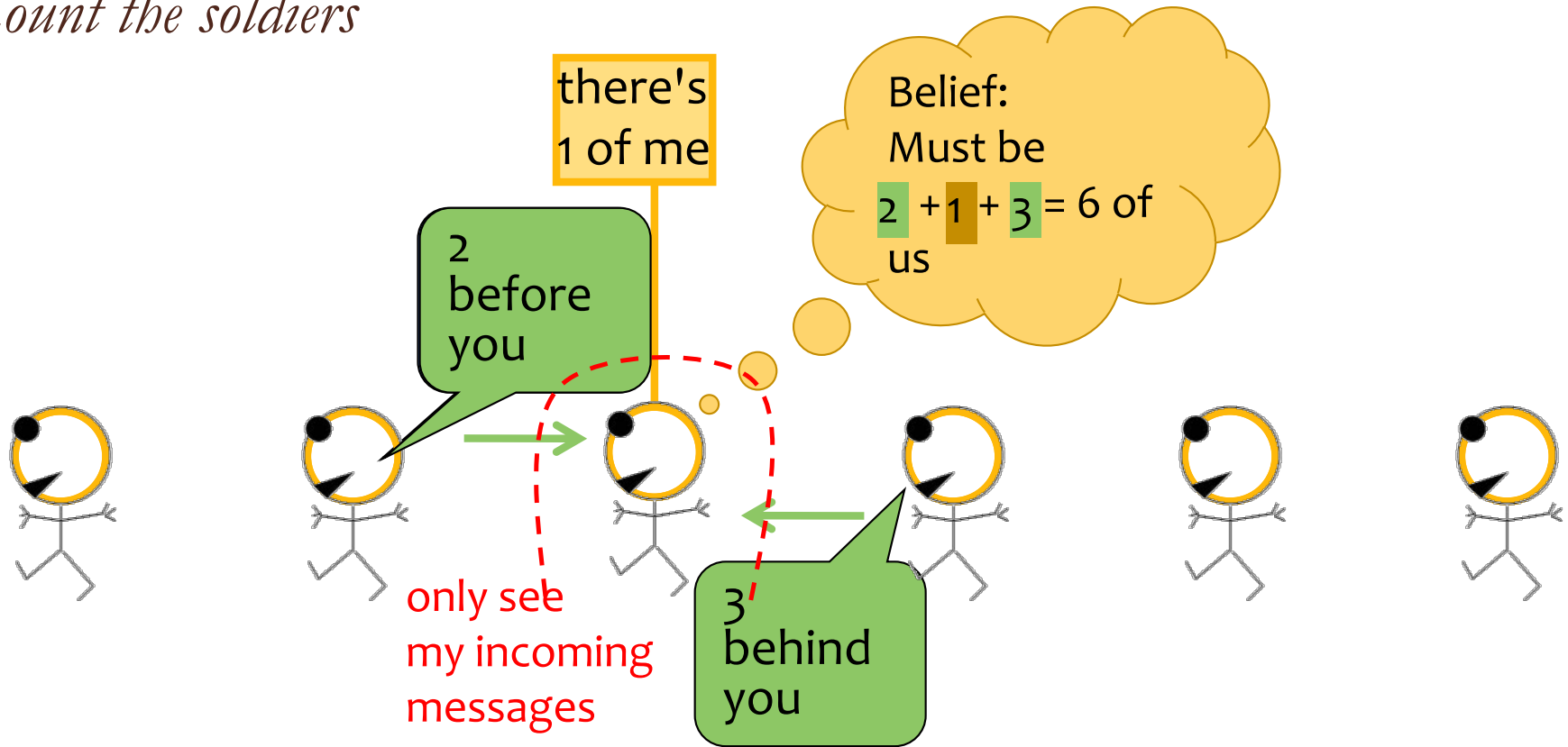
# Great Ideas in ML: Message Passing

*Count the soldiers*



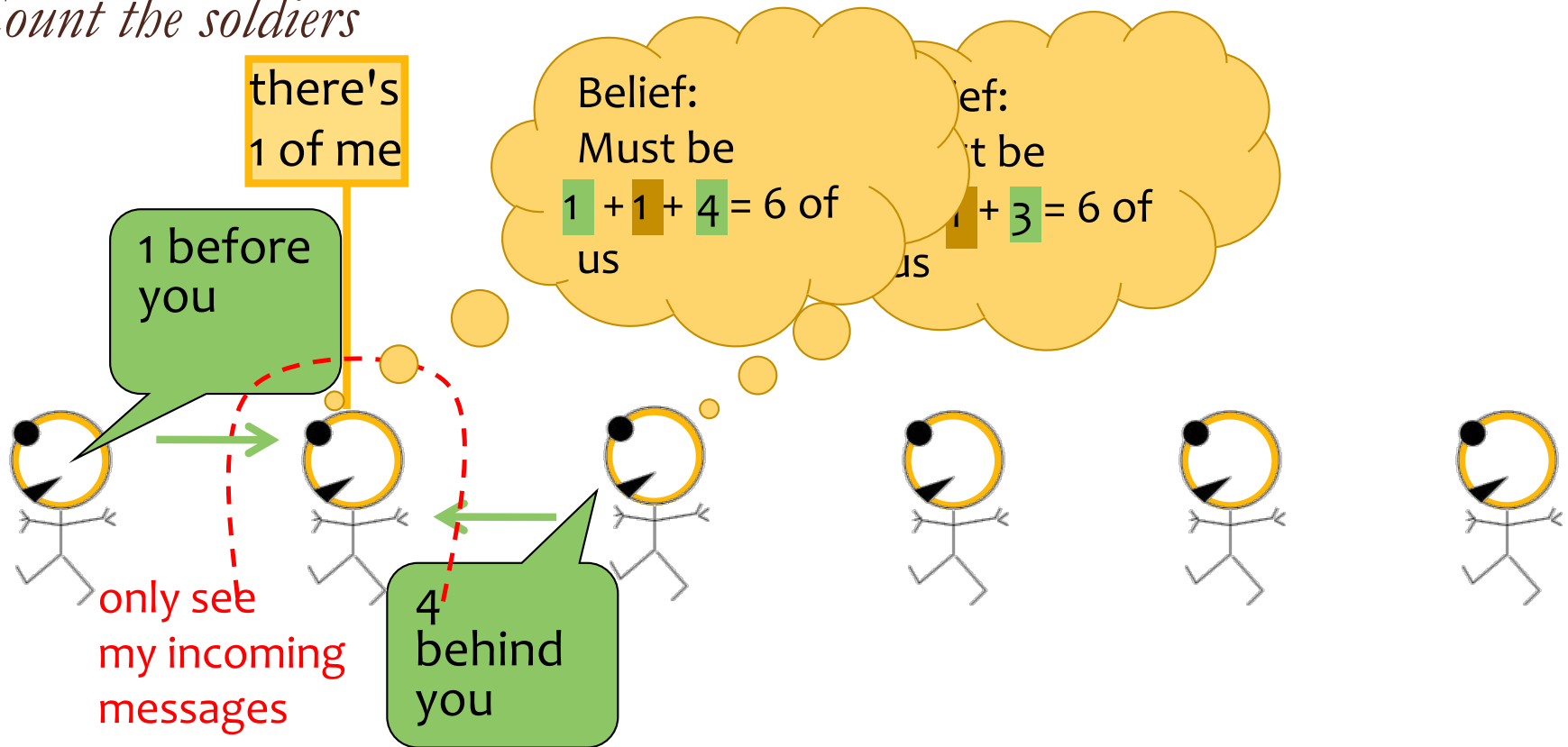
# Great Ideas in ML: Message Passing

*Count the soldiers*



# Great Ideas in ML: Message Passing

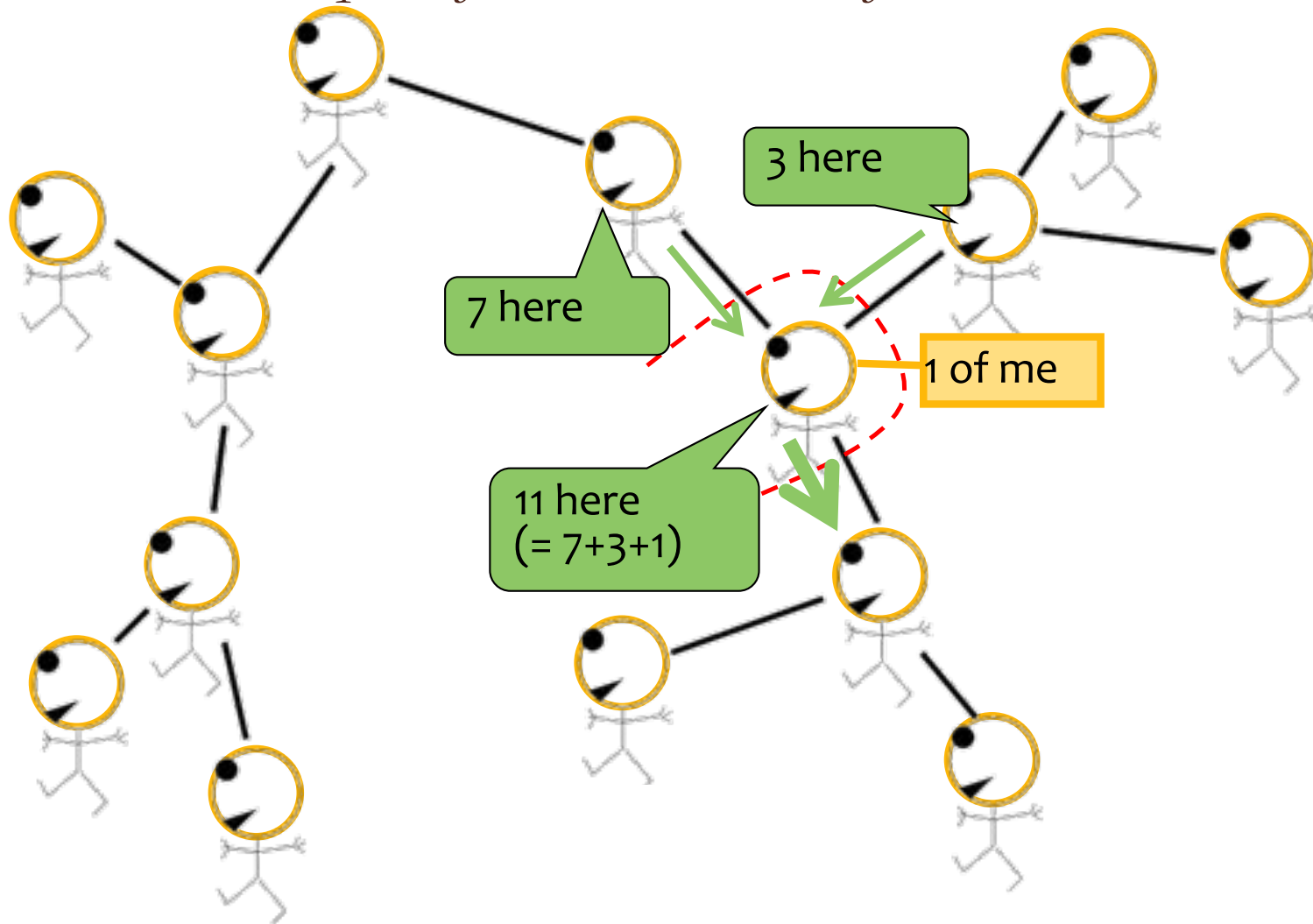
*Count the soldiers*





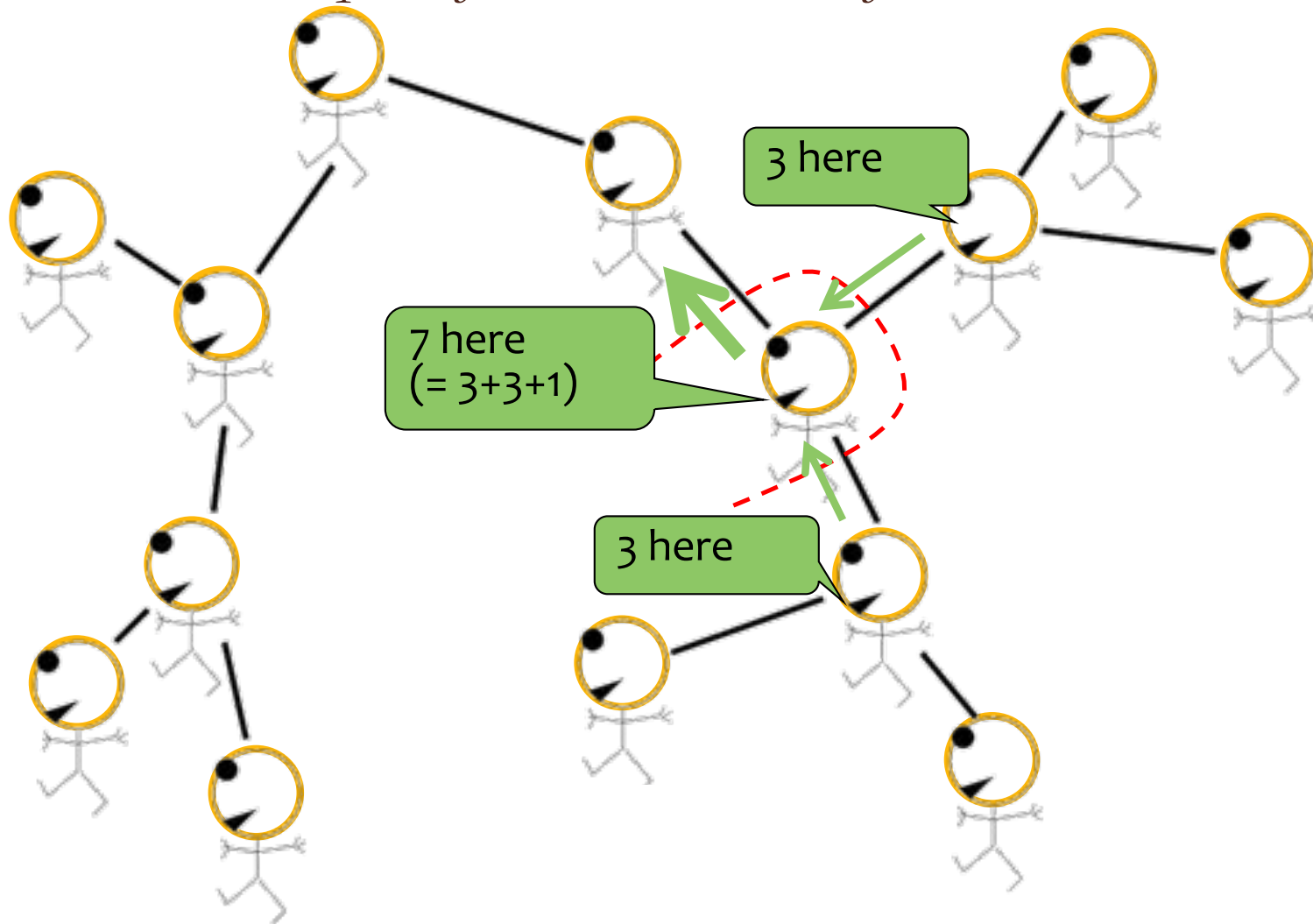
# Great Ideas in ML: Message Passing

*Each soldier receives reports from all branches of tree*



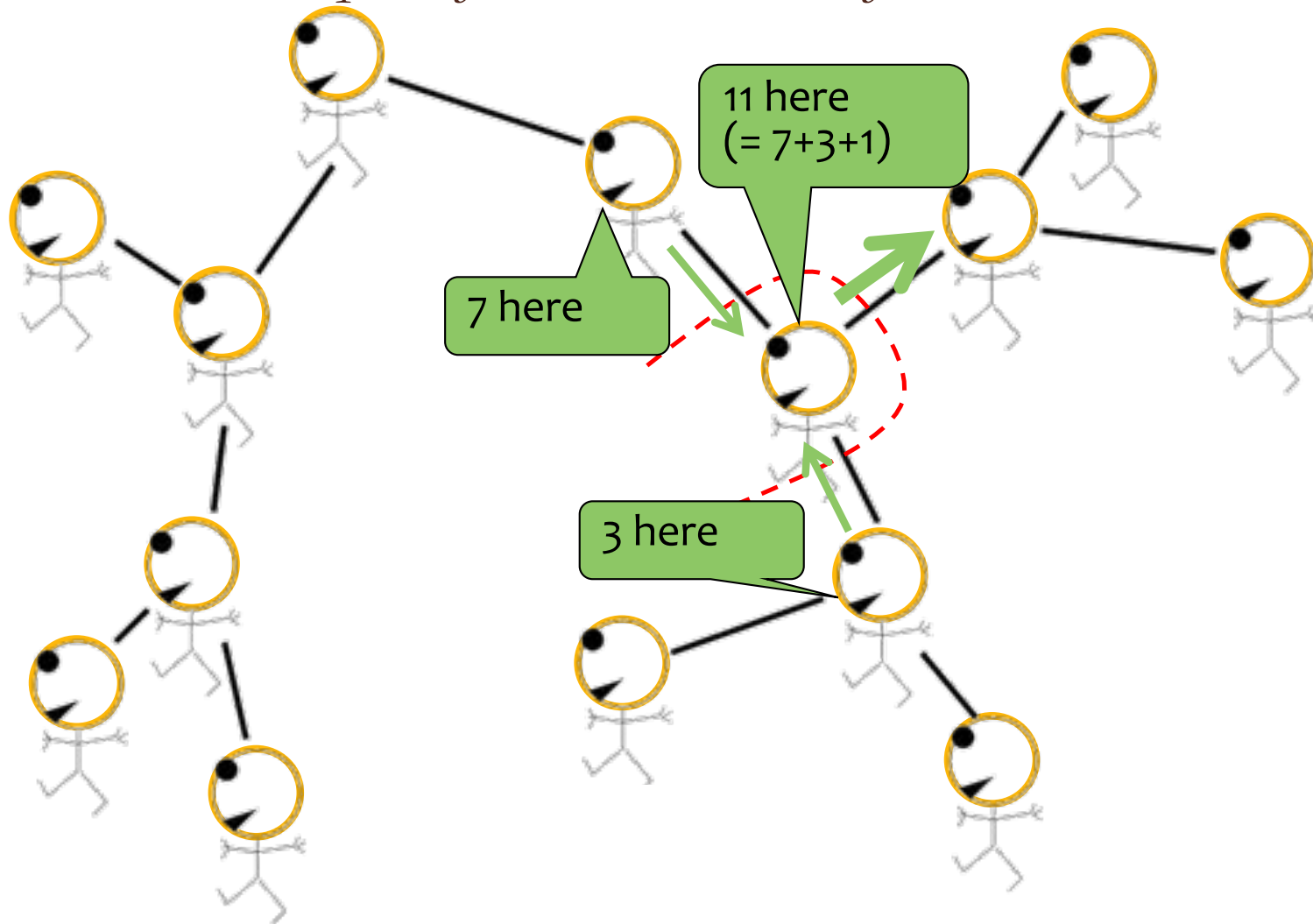
# Great Ideas in ML: Message Passing

*Each soldier receives reports from all branches of tree*



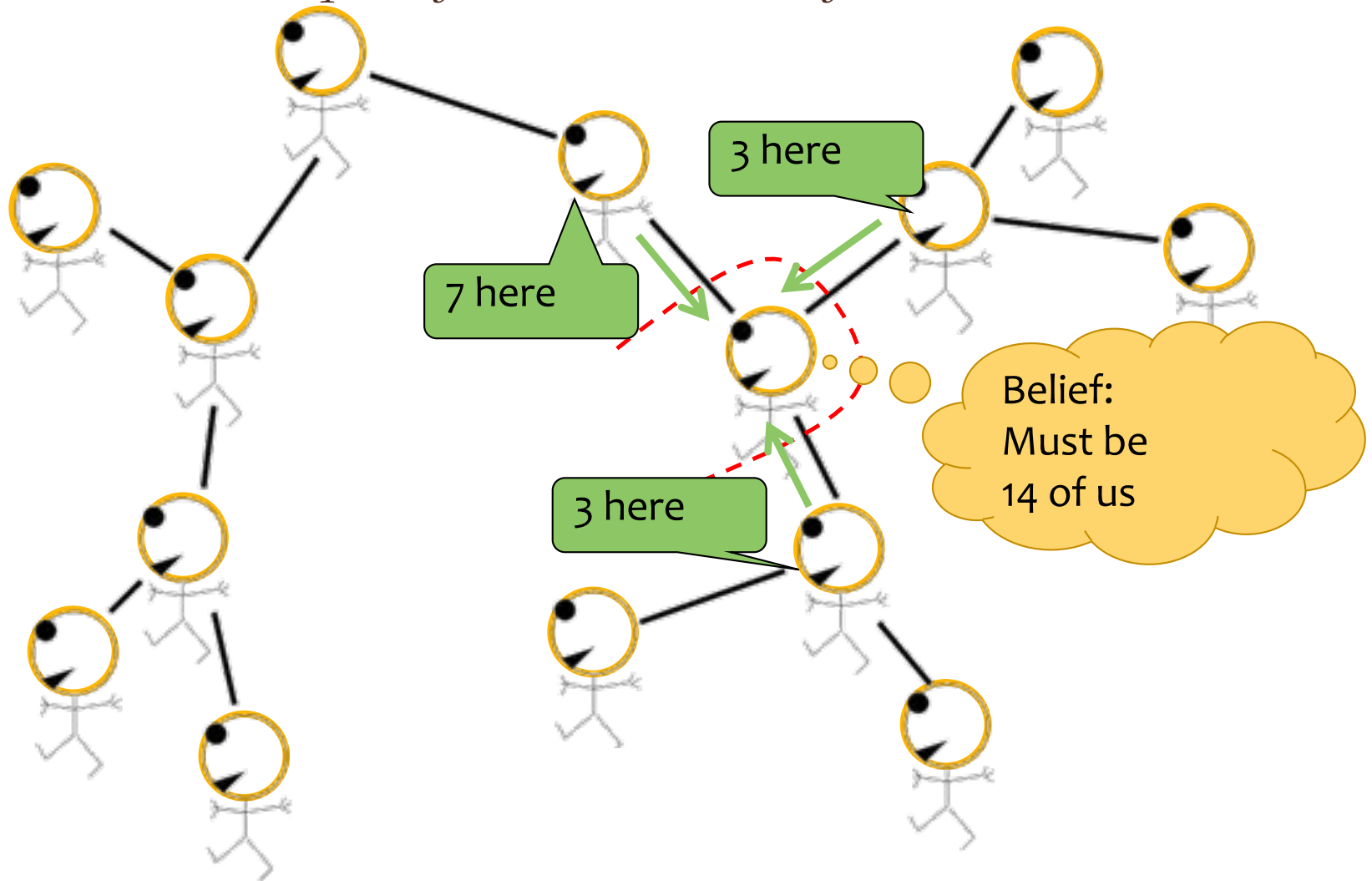
# Great Ideas in ML: Message Passing

*Each soldier receives reports from all branches of tree*



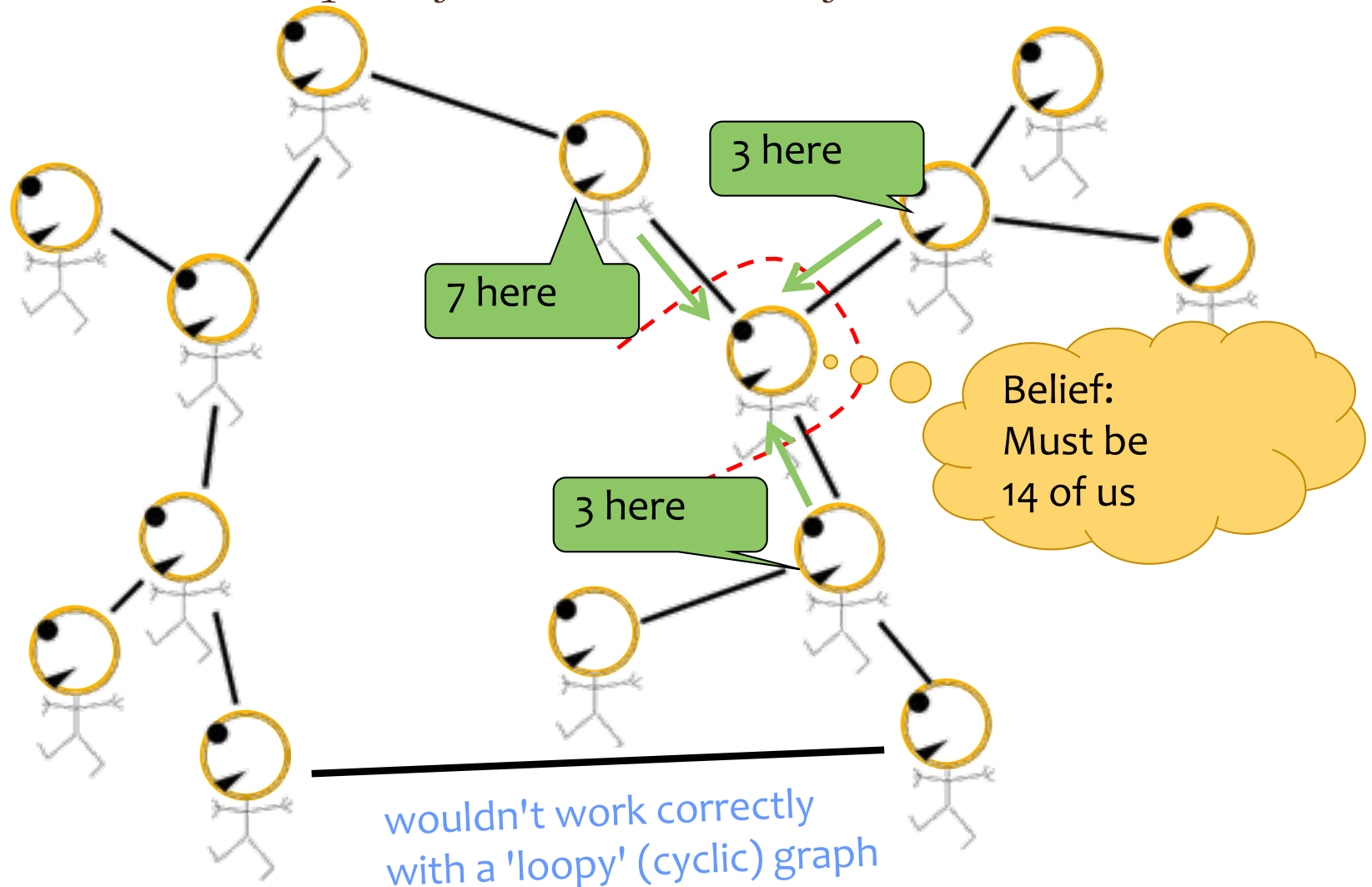
# Great Ideas in ML: Message Passing

*Each soldier receives reports from all branches of tree*



# Great Ideas in ML: Message Passing

*Each soldier receives reports from all branches of tree*



# **THE FORWARD-BACKWARD ALGORITHM**

# Inference for HMMs

















































## *Whiteboard*

### – Three Inference Problems for an HMM

1. Evaluation: Compute the probability of a given sequence of observations
2. Viterbi Decoding: Find the most-likely sequence of hidden states, given a sequence of observations
3. Marginals: Compute the marginal distribution for a hidden state, given a sequence of observations

# Dataset for Supervised Part-of-Speech (POS) Tagging

Data:  $\mathcal{D} = \{x^{(n)}, y^{(n)}\}_{n=1}^N$

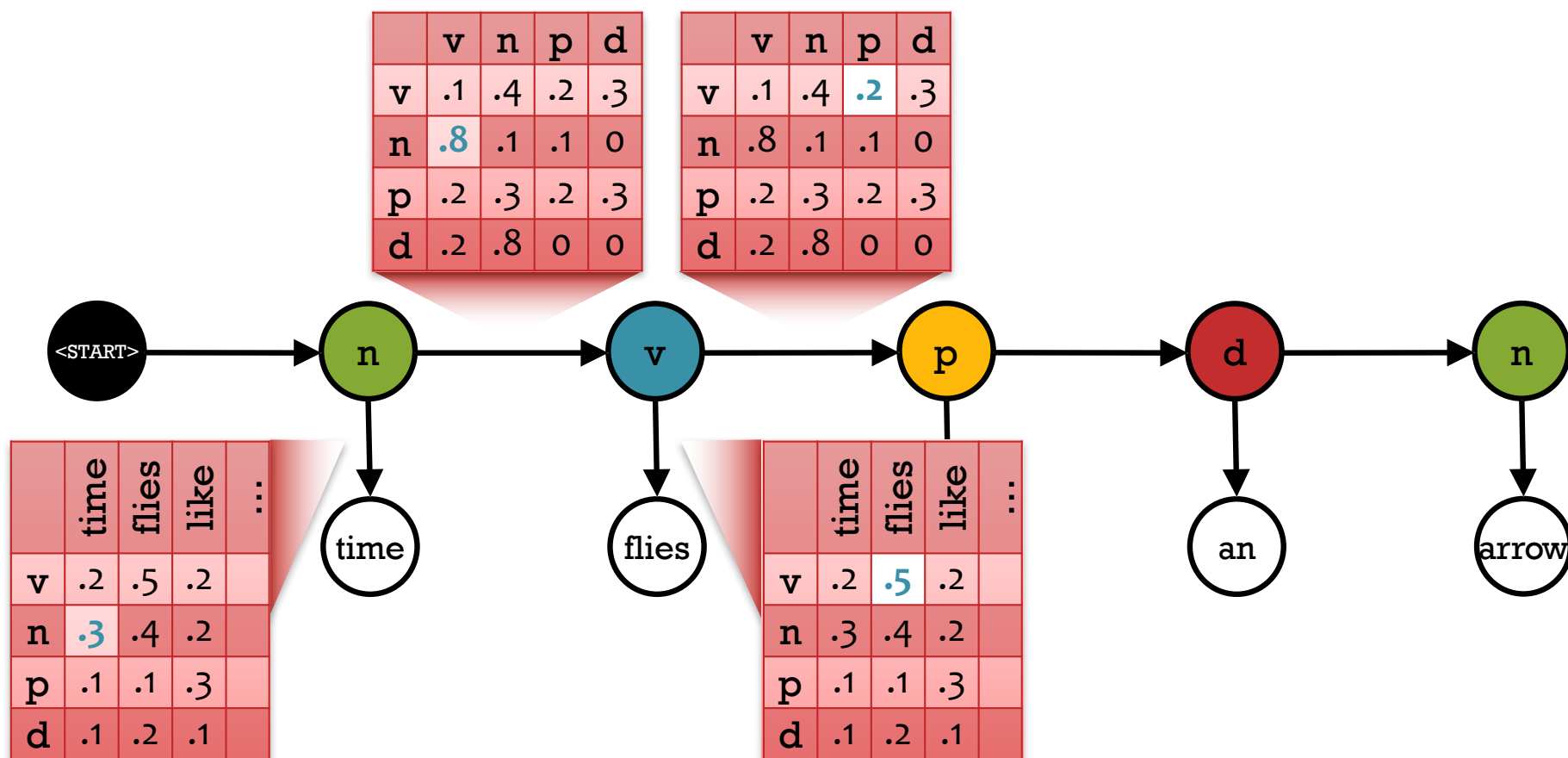
Sample 1:							$y^{(1)}$
							$x^{(1)}$
Sample 2:							$y^{(2)}$
							$x^{(2)}$
Sample 3:							$y^{(3)}$
							$x^{(3)}$
Sample 4:							$y^{(4)}$
							$x^{(4)}$



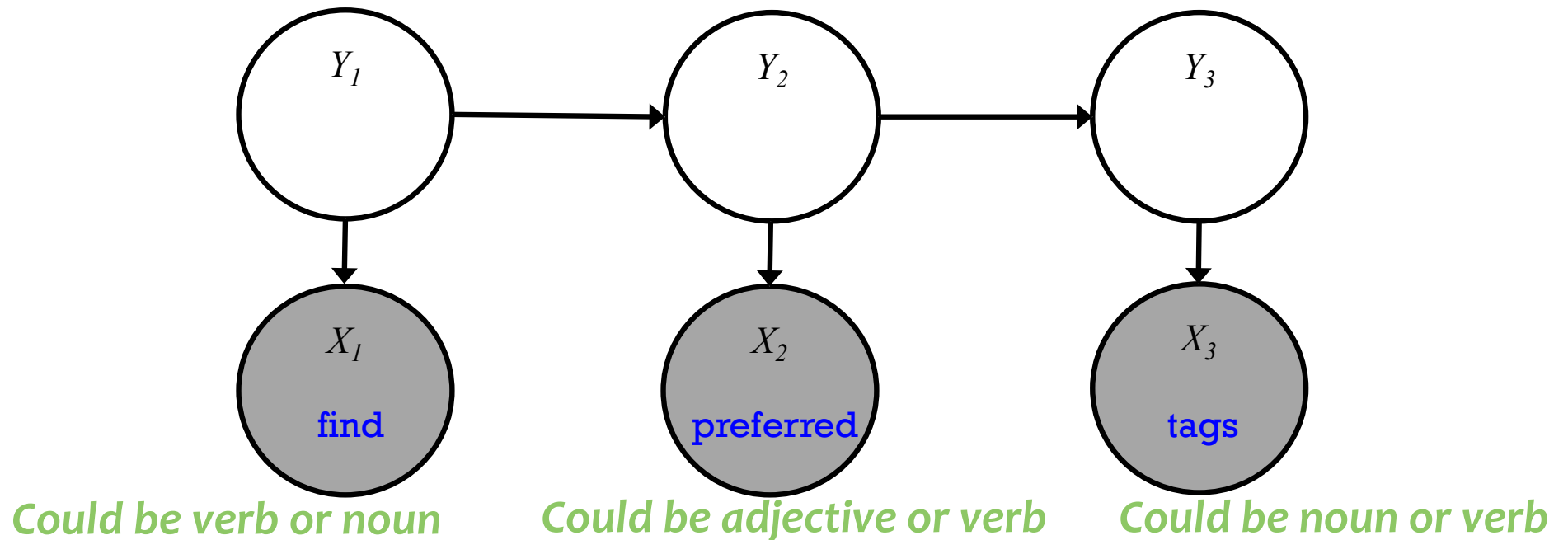
# Hidden Markov Model

A Hidden Markov Model (HMM) provides a joint distribution over the sentence/tags with an assumption of dependence between adjacent tags.

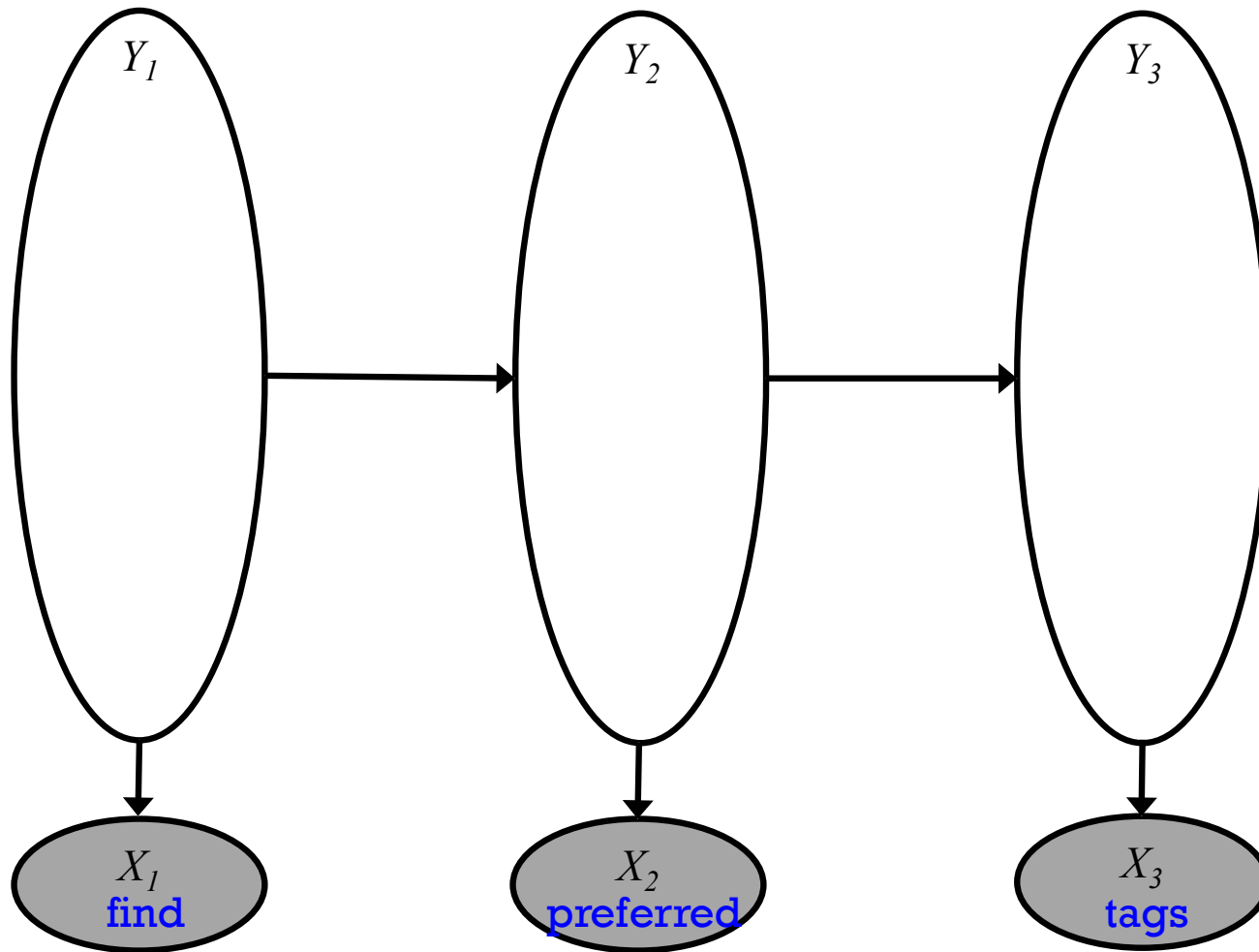
$$p(n, v, p, d, n, \text{time, flies, like, an, arrow}) = (.3 * .8 * .2 * .5 * \dots)$$



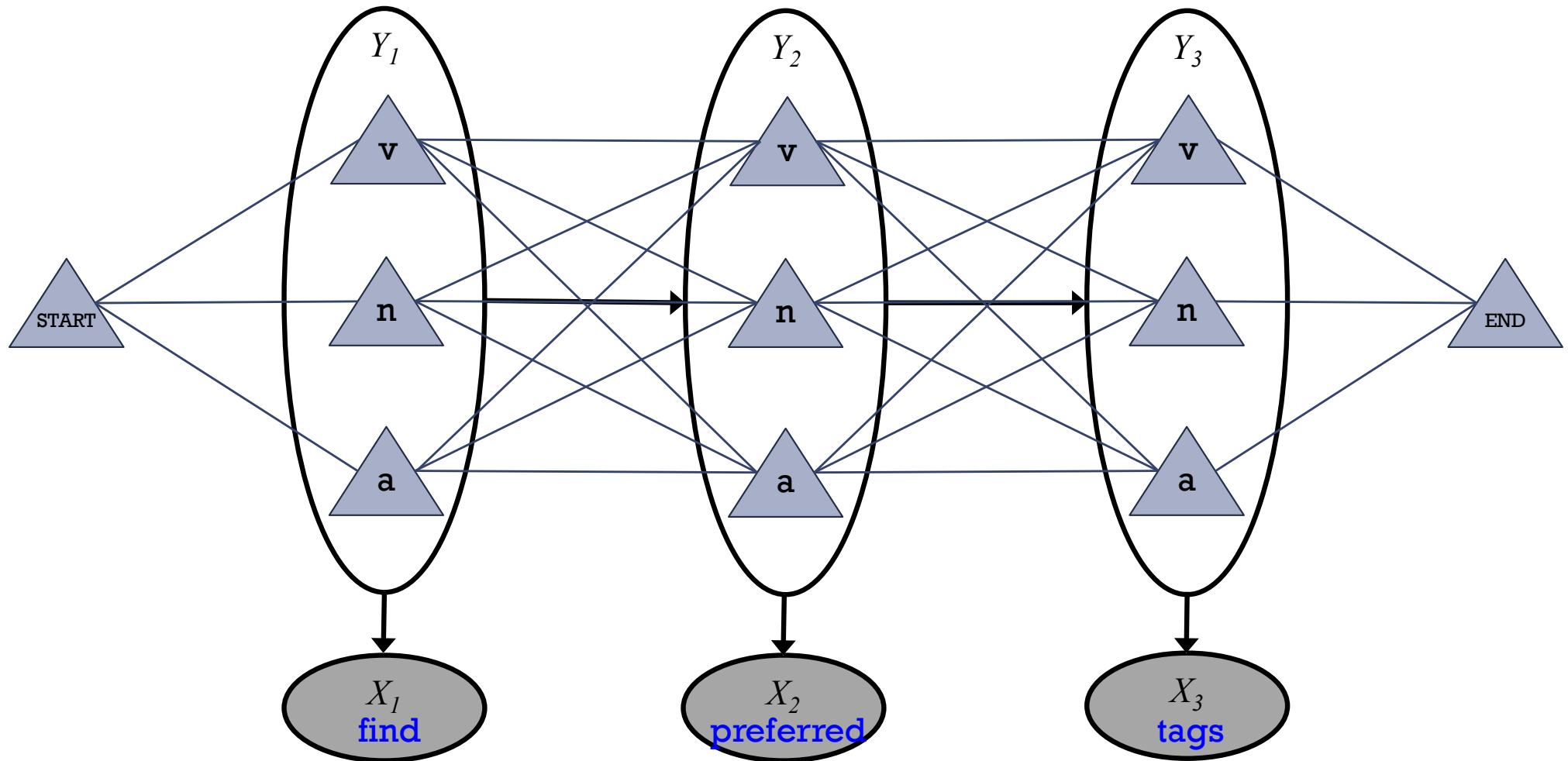
# Forward-Backward Algorithm



# Forward-Backward Algorithm

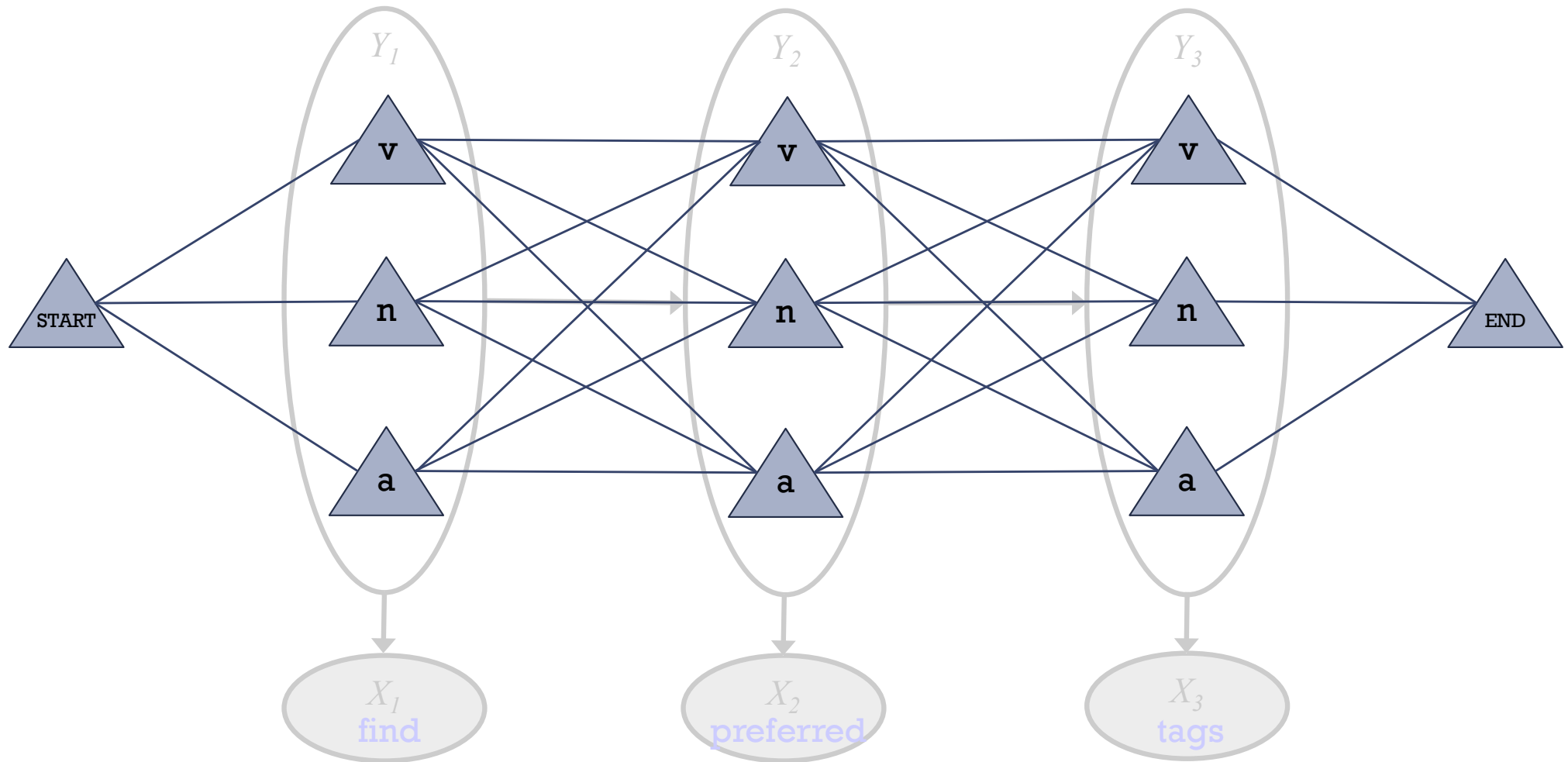


# Forward-Backward Algorithm



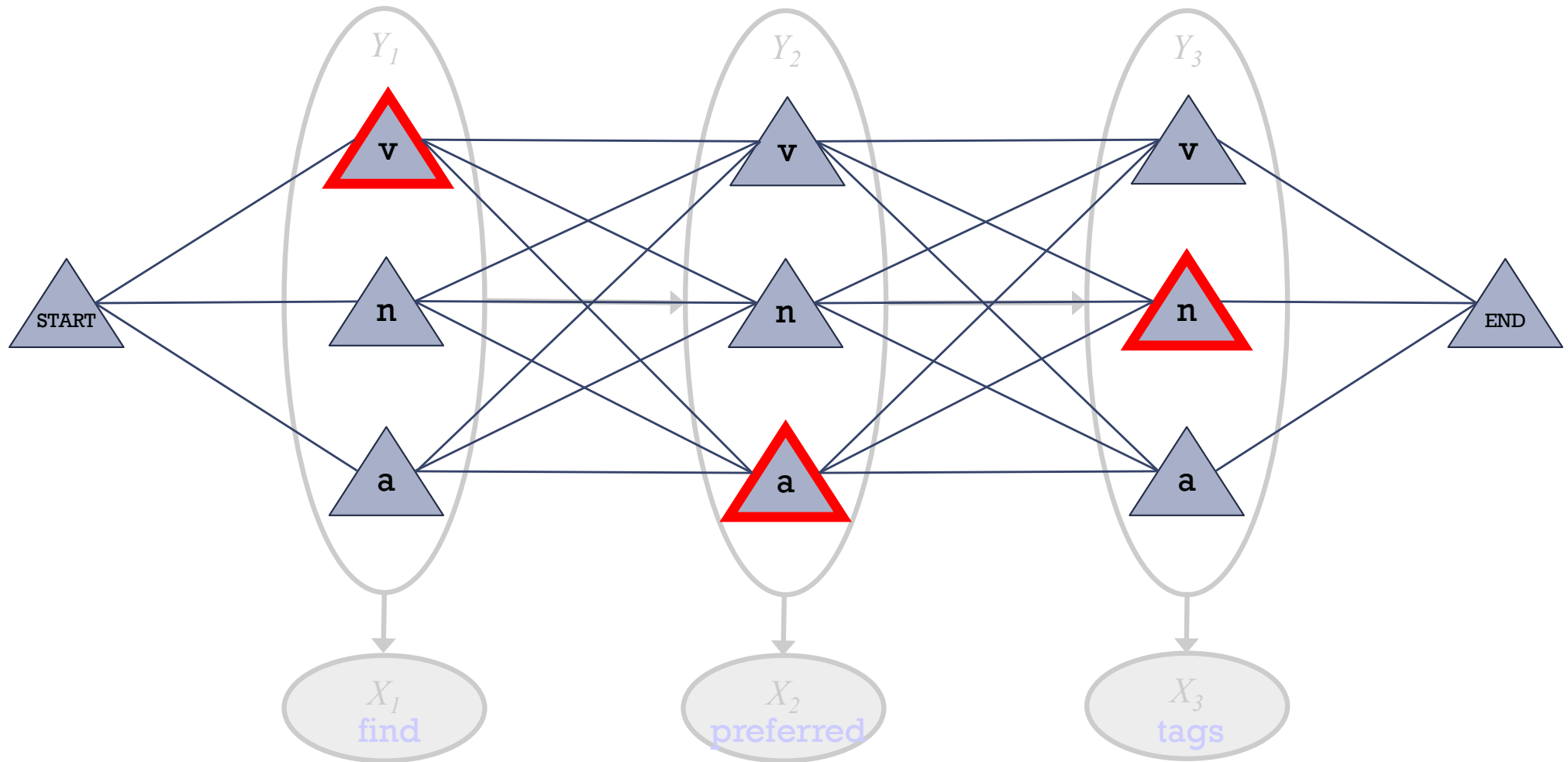
- Let's show the possible values for each variable

# Forward-Backward Algorithm



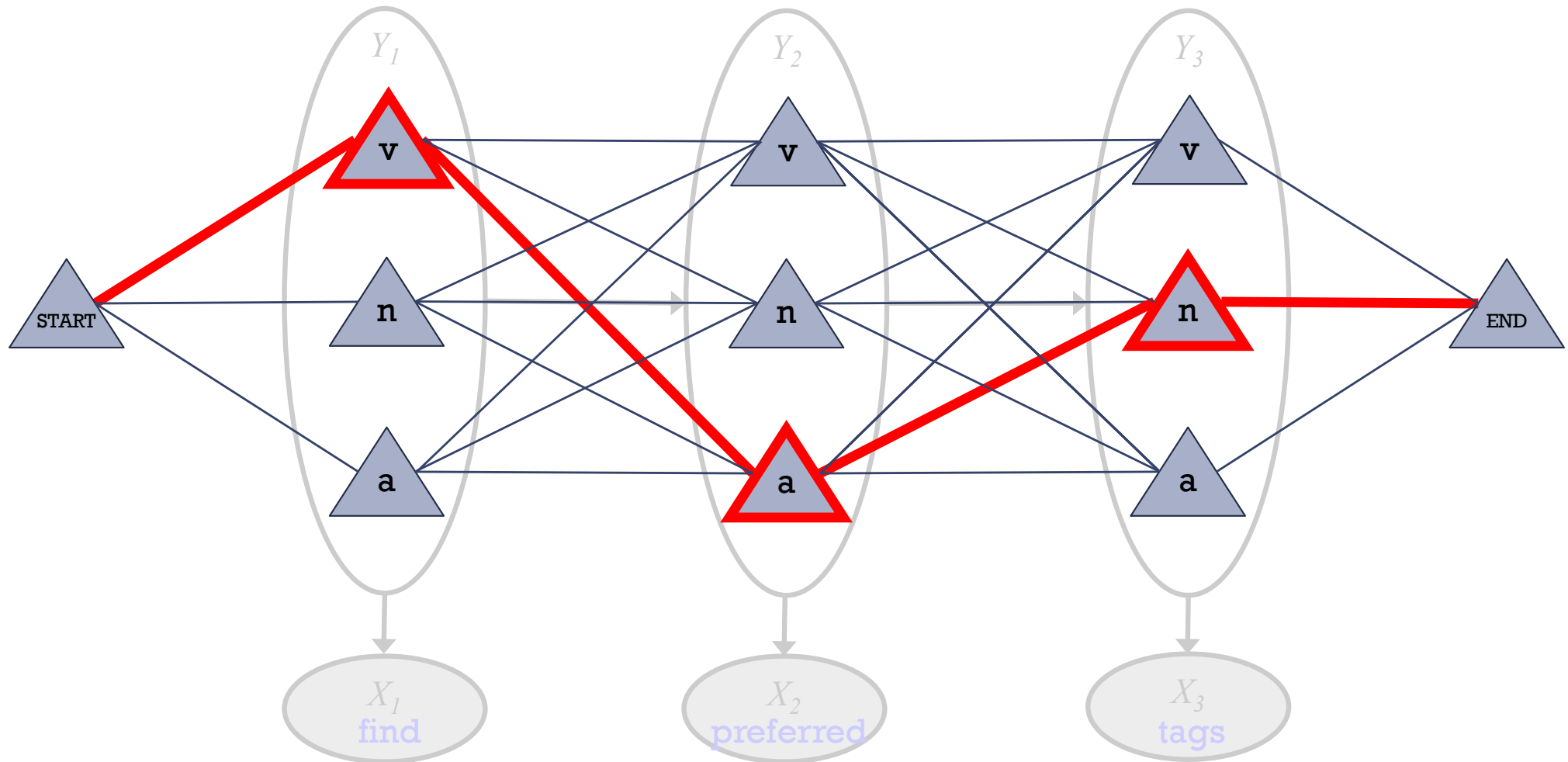
- Let's show the possible values for each variable

# Forward-Backward Algorithm



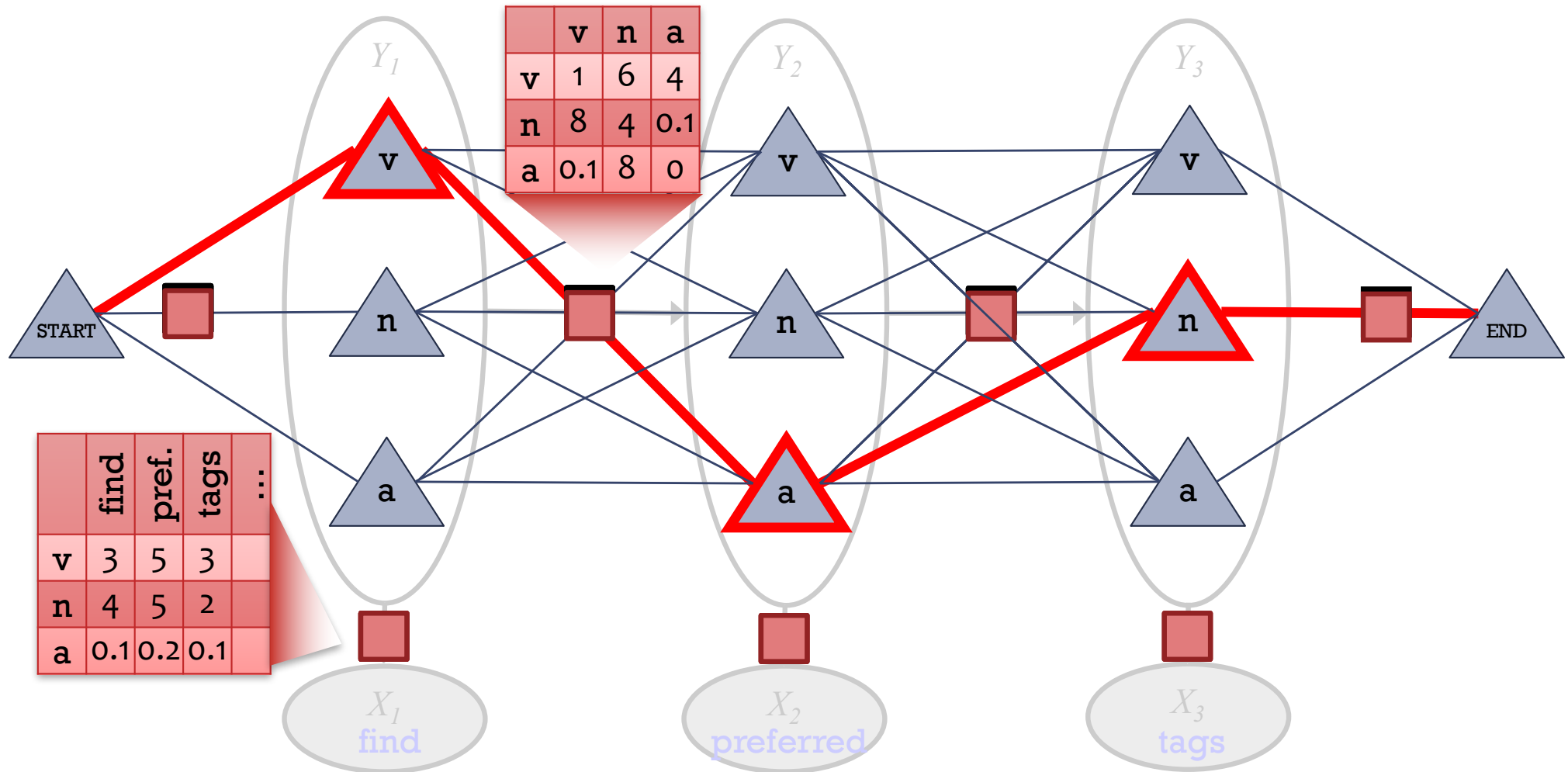
- Let's show the possible *values* for each variable
- One possible assignment

# Forward-Backward Algorithm



- Let's show the possible *values* for each variable
- One possible assignment
- And what the 7 transition / emission factors **think of it** ...

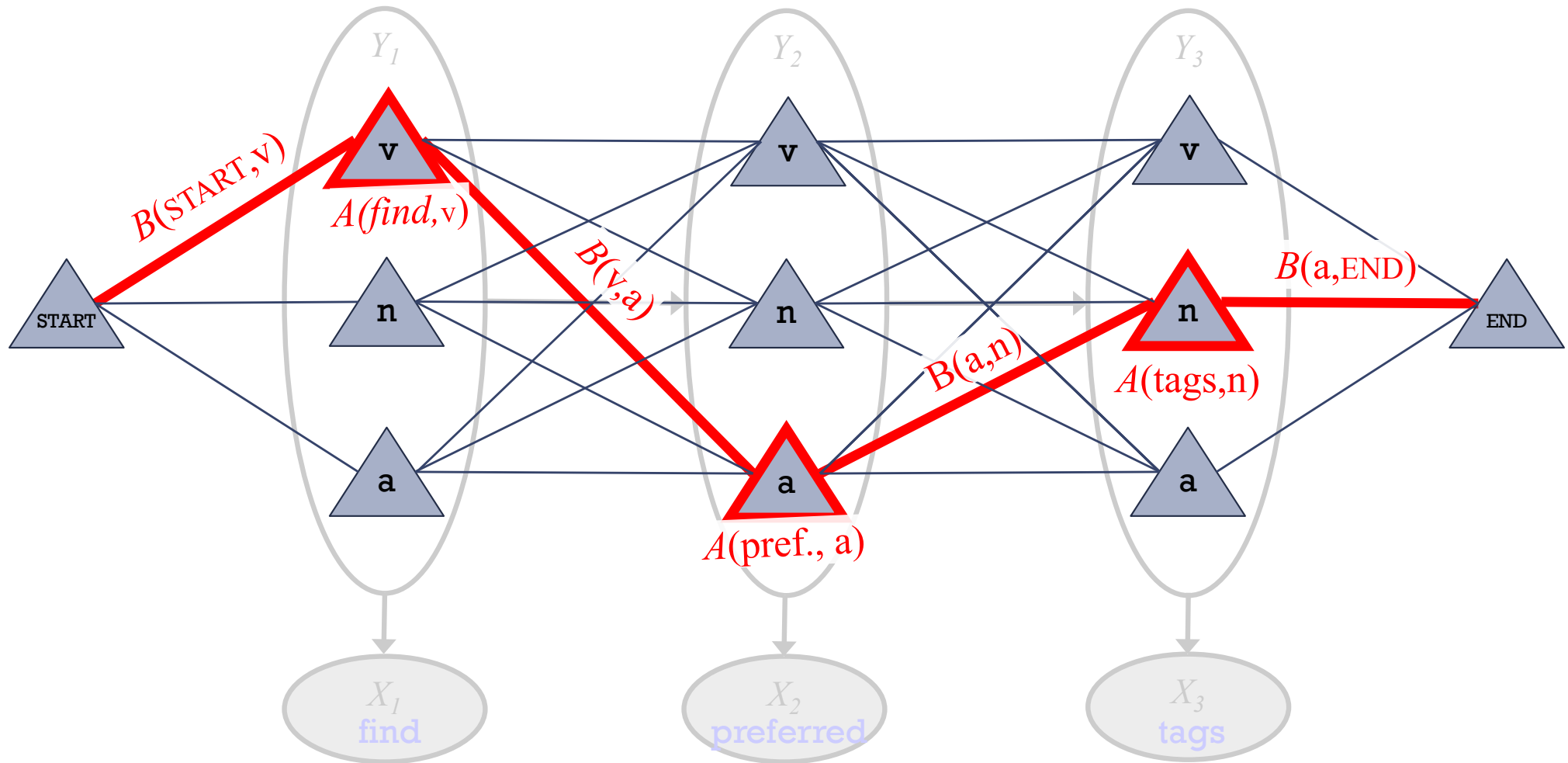
# Forward-Backward Algorithm



- Let's show the possible *values* for each variable
- One possible assignment
- And what the 7 transition / emission factors **think of it** ...

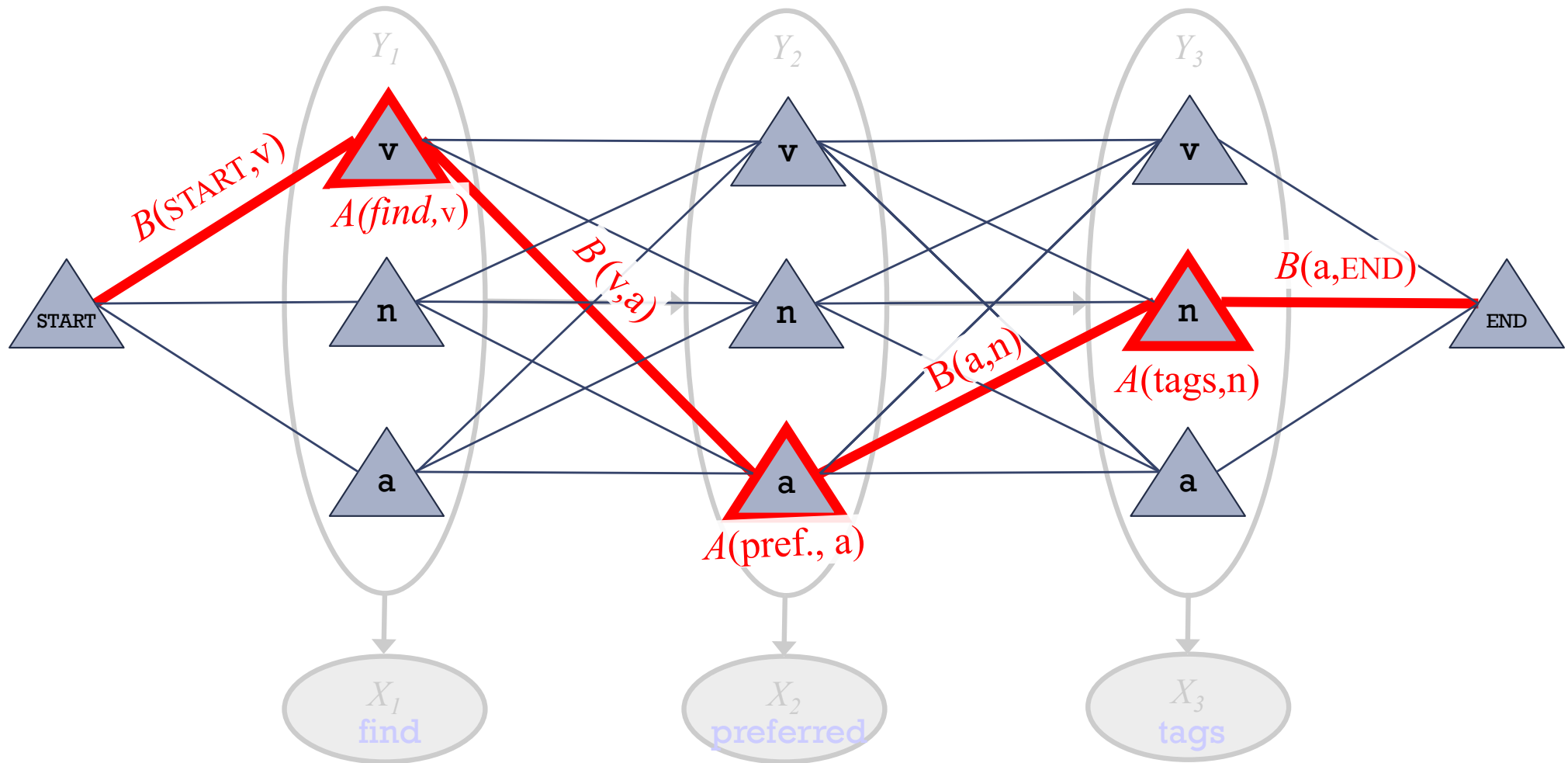


# Viterbi Algorithm: Most Probable Assignment



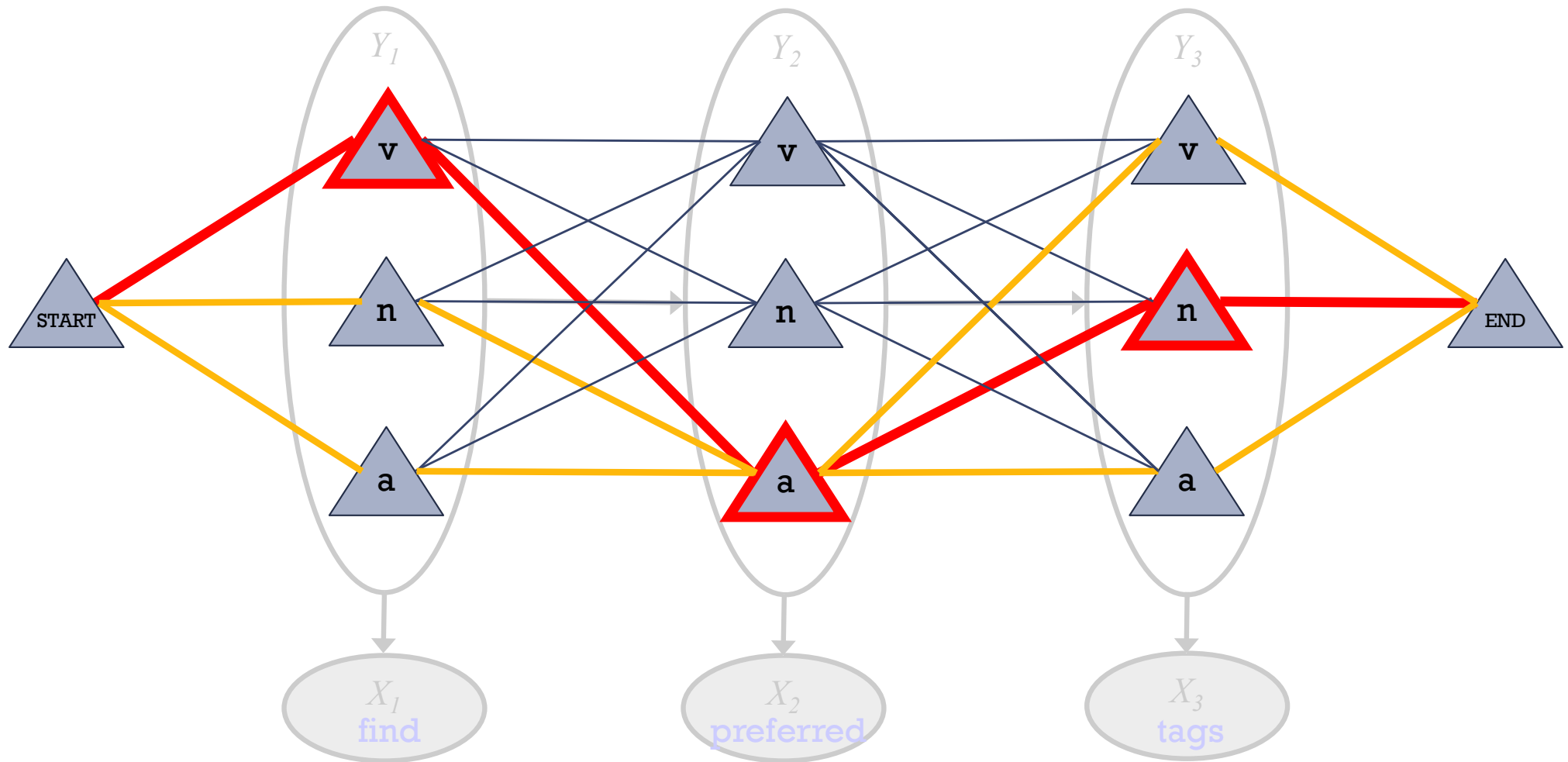
- So  $p(\mathbf{v} \mathbf{a} \mathbf{n}) = (1/Z) * \text{product of 7 numbers}$
- Numbers associated with edges and nodes of path
- Most probable assignment = **path with highest product**

# Viterbi Algorithm: Most Probable Assignment

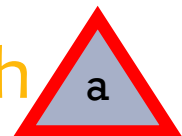


- So  $p(\mathbf{v} \mathbf{a} \mathbf{n}) = (1/Z) * \text{product weight of one path}$

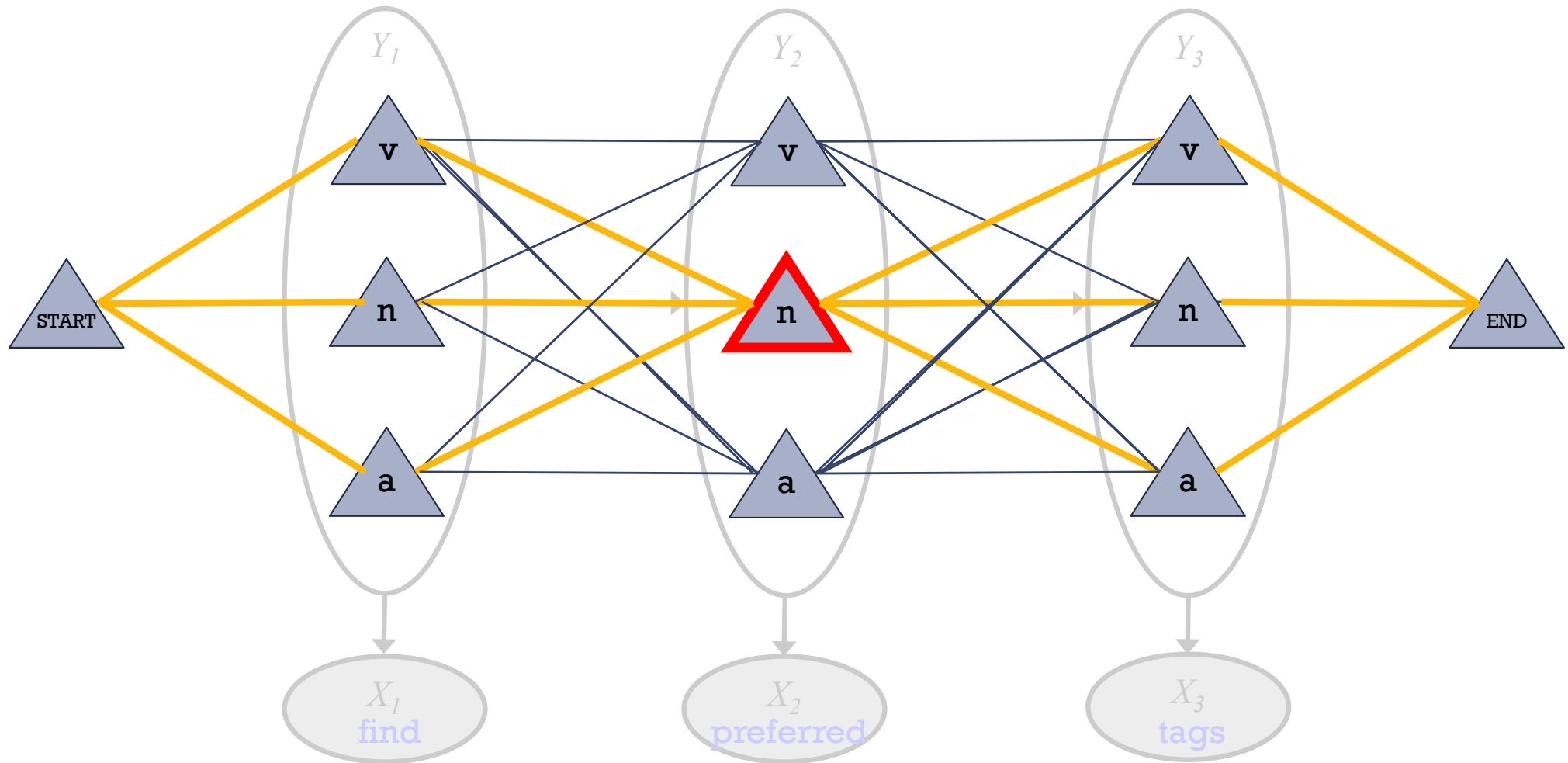
# Forward-Backward Algorithm: Finds Marginals



- So  $p(\mathbf{v} \mathbf{a} \mathbf{n}) = (1/Z) * \text{product weight of one path}$
- Marginal probability  $p(Y_2 = a)$   
 $= (1/Z) * \text{total weight of all paths through } \mathbf{a}$

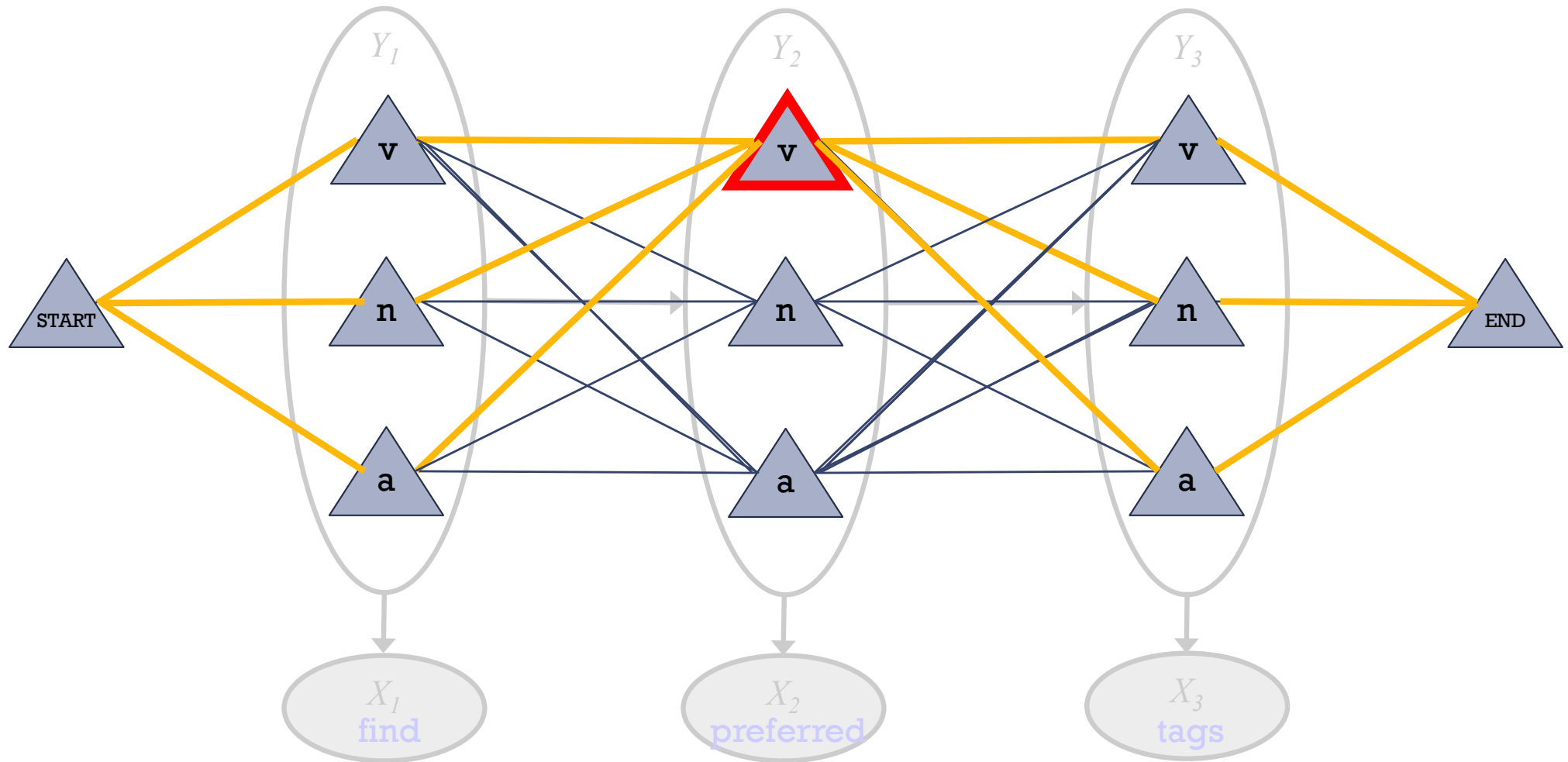


# Forward-Backward Algorithm: Finds Marginals

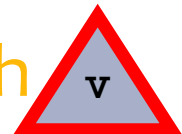


- So  $p(v \ a \ n) = (1/Z) * \text{product weight of one path}$
- Marginal probability  $p(Y_2 = n)$   
 $= (1/Z) * \text{total weight of all paths through } \triangle n$

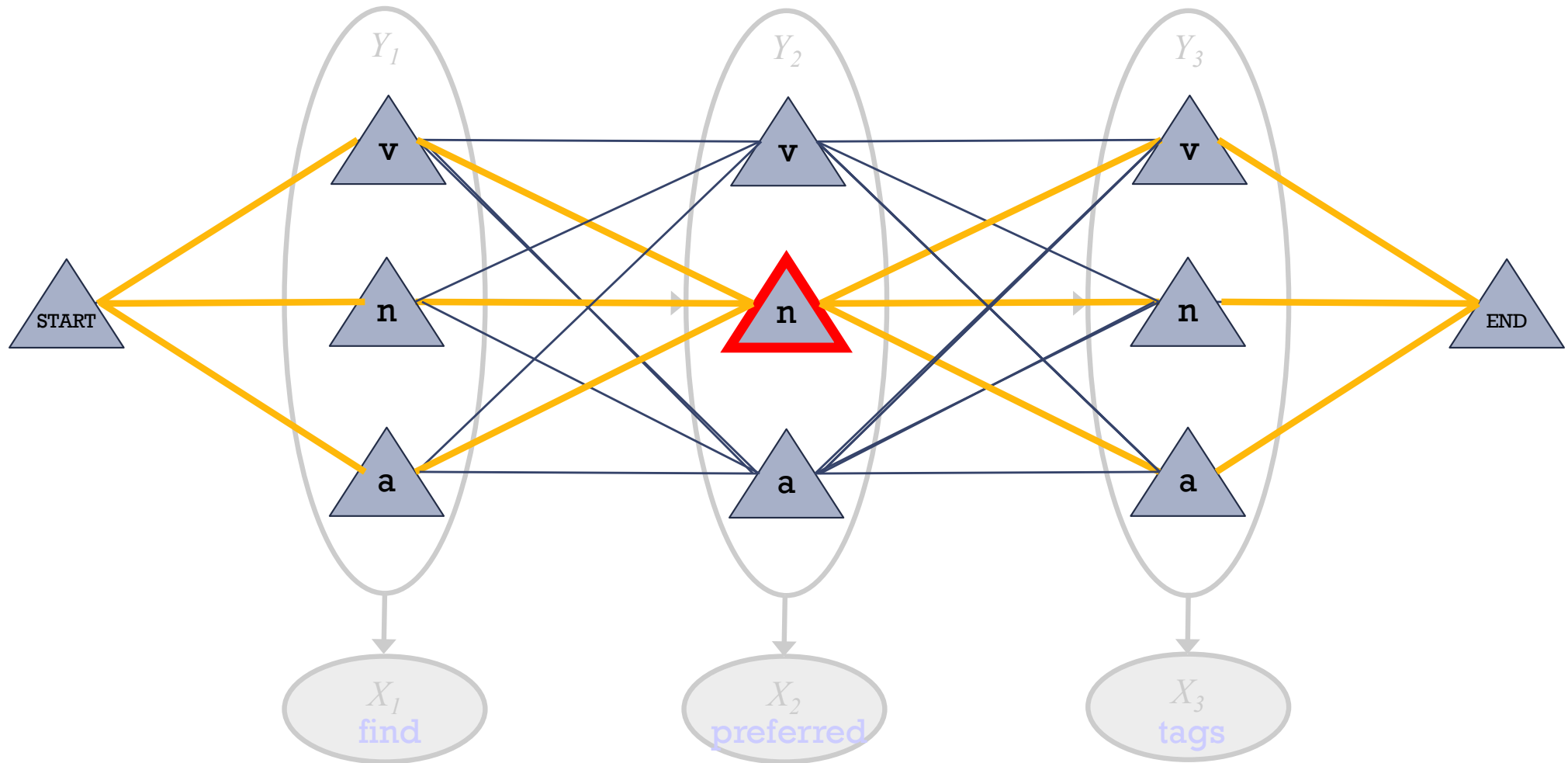
# Forward-Backward Algorithm: Finds Marginals



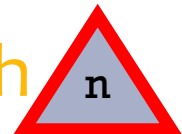
- So  $p(\mathbf{v} \mathbf{a} \mathbf{n}) = (1/Z) * \text{product weight of one path}$
- Marginal probability  $p(Y_2 = v)$   
 $= (1/Z) * \text{total weight of all paths through } \mathbf{v}$



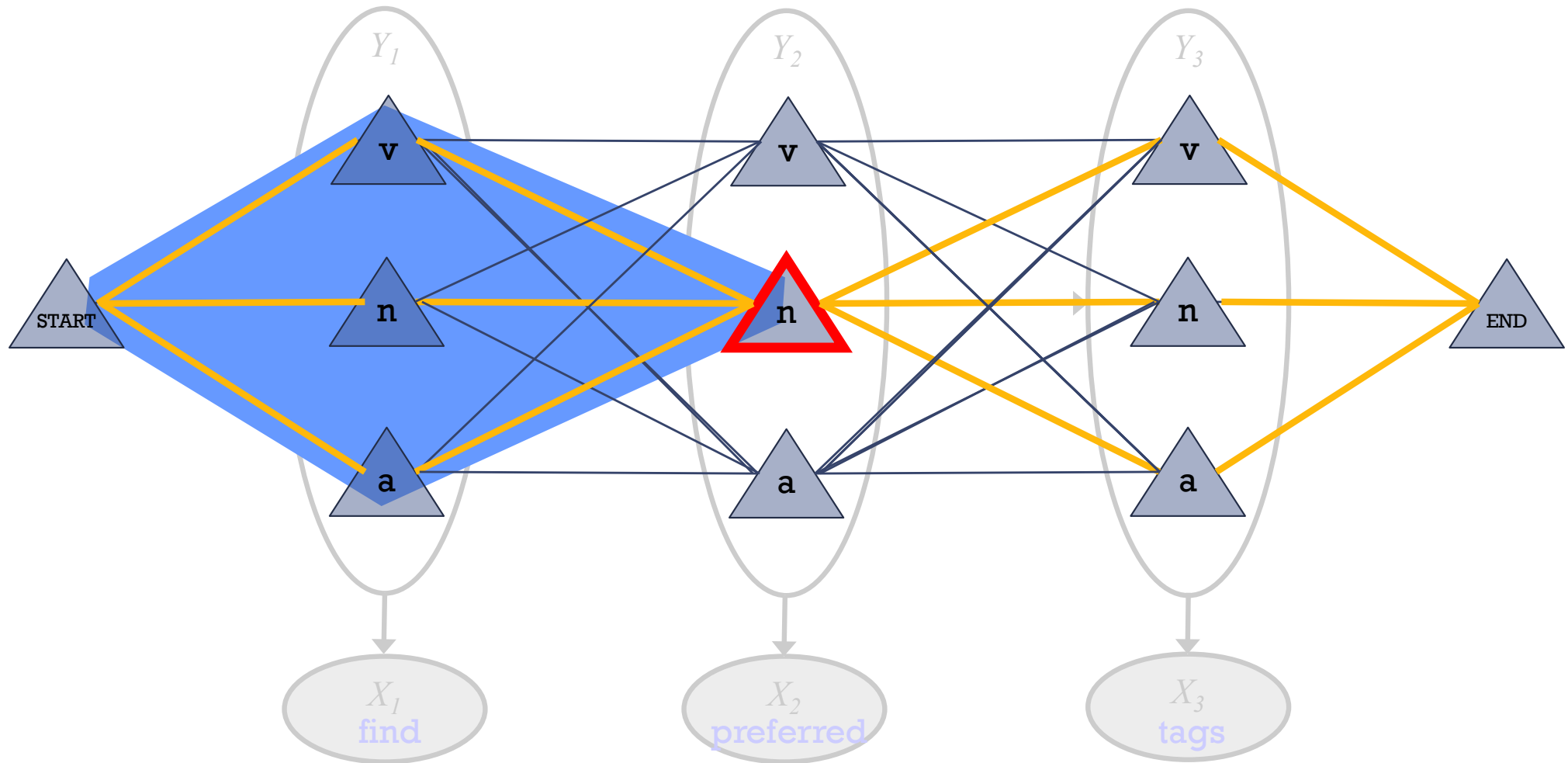
# Forward-Backward Algorithm: Finds Marginals



- So  $p(\mathbf{v} \mathbf{a} \mathbf{n}) = (1/Z) * \text{product weight of one path}$
- Marginal probability  $p(Y_2 = \mathbf{n})$   
 $= (1/Z) * \text{total weight of all paths through } \mathbf{n}$



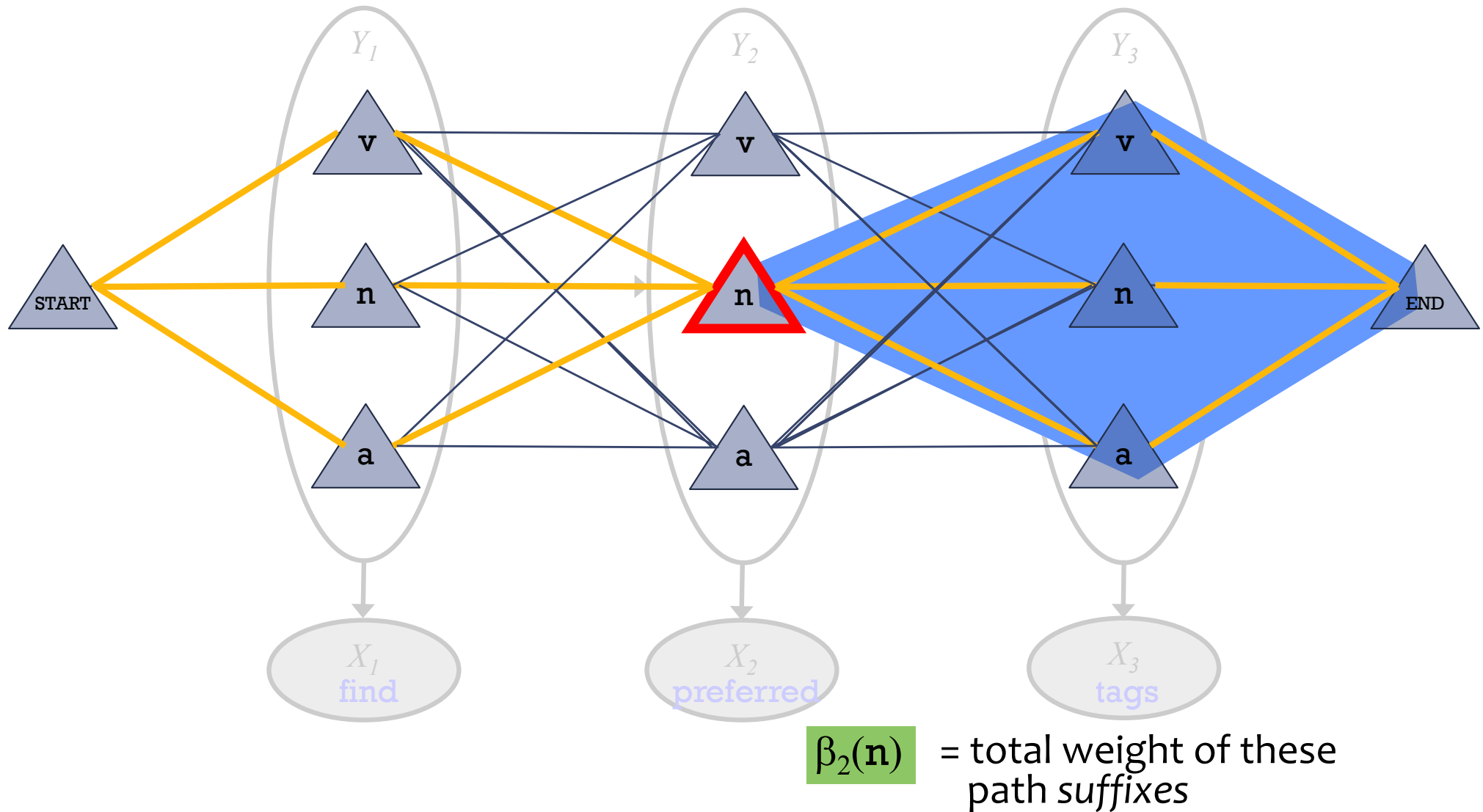
# Forward-Backward Algorithm: Finds Marginals



$\alpha_2(\mathbf{n})$  = total weight of these path *prefixes*

(found by dynamic programming: matrix-vector products)

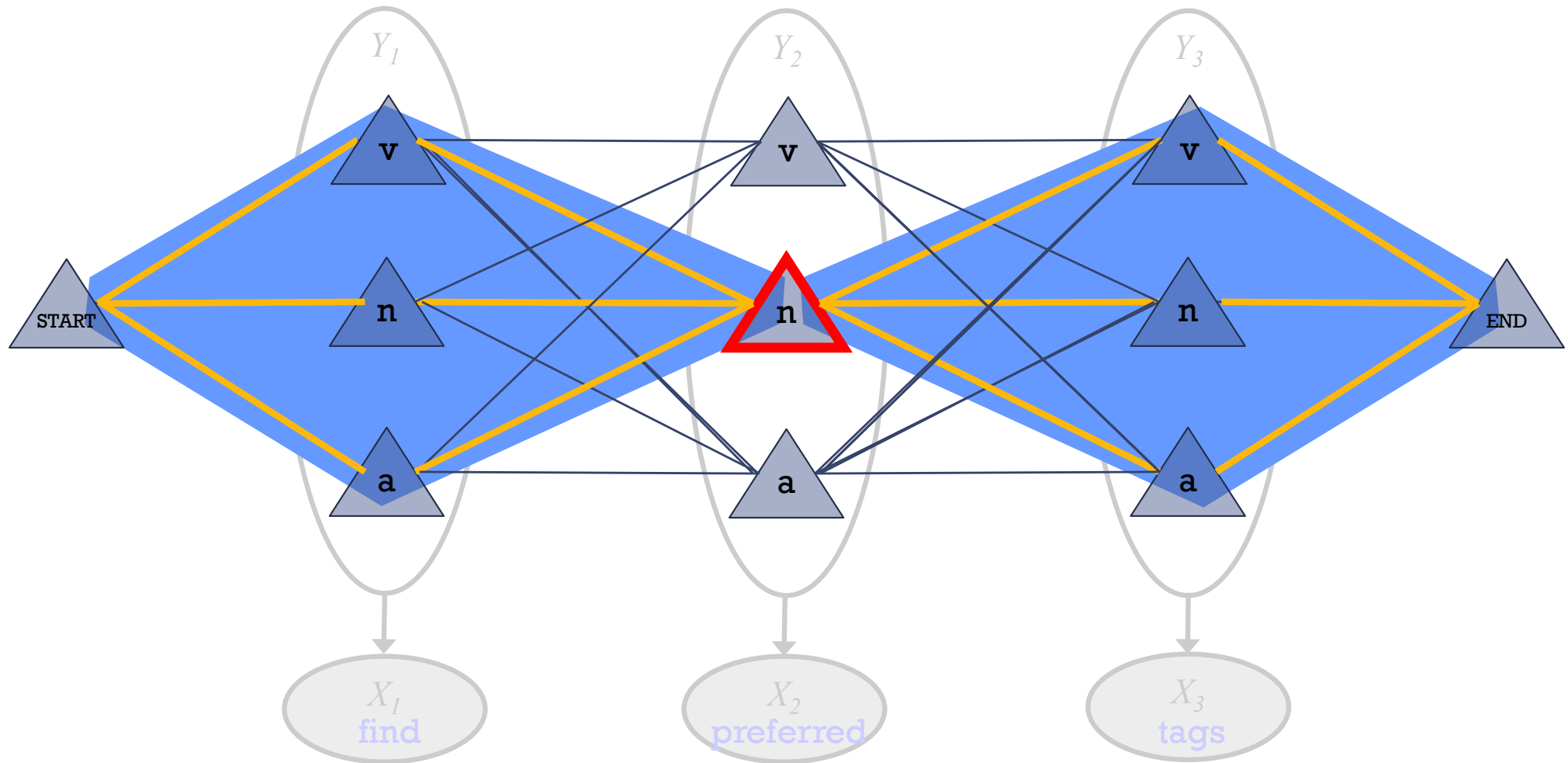
# Forward-Backward Algorithm: Finds Marginals



(found by dynamic programming: matrix-vector products)



# Forward-Backward Algorithm: Finds Marginals



$\alpha_2(\mathbf{n})$  = total weight of these path prefixes  $(a + b + c)$

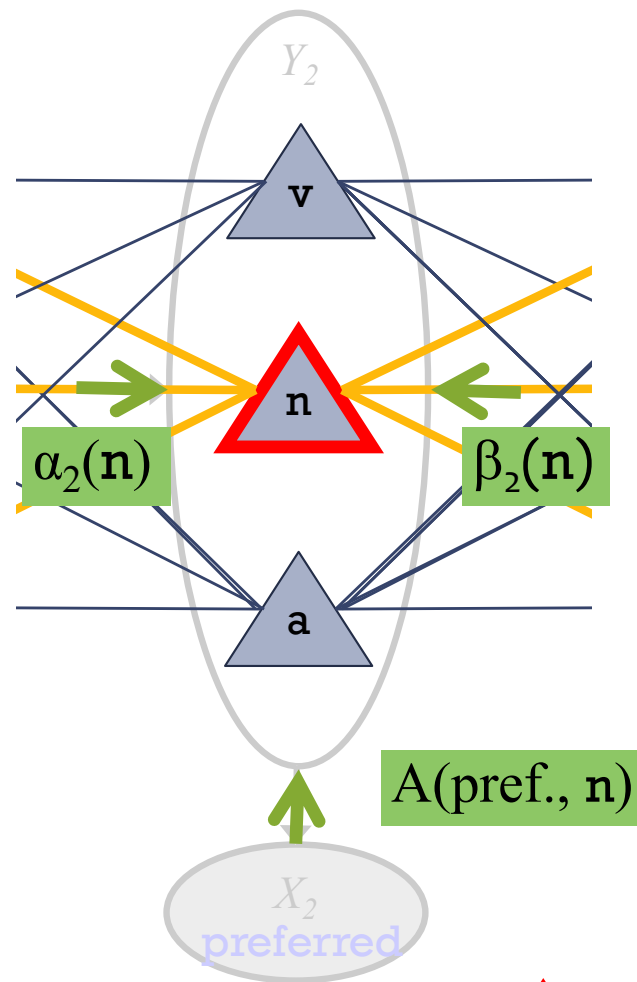
$\beta_2(\mathbf{n})$  = total weight of these path suffixes  $(x + y + z)$

Product gives  $ax+ay+az+bx+by+bz+cx+cy+cz$  = total weight of paths

# Forward-Backward Algorithm: Finds Marginals

Oops! The weight of a path through a state also includes a weight at that state.  
So  $\alpha(n) \cdot \beta(n)$  isn't enough.

The extra weight is the opinion of the emission probability at this variable.

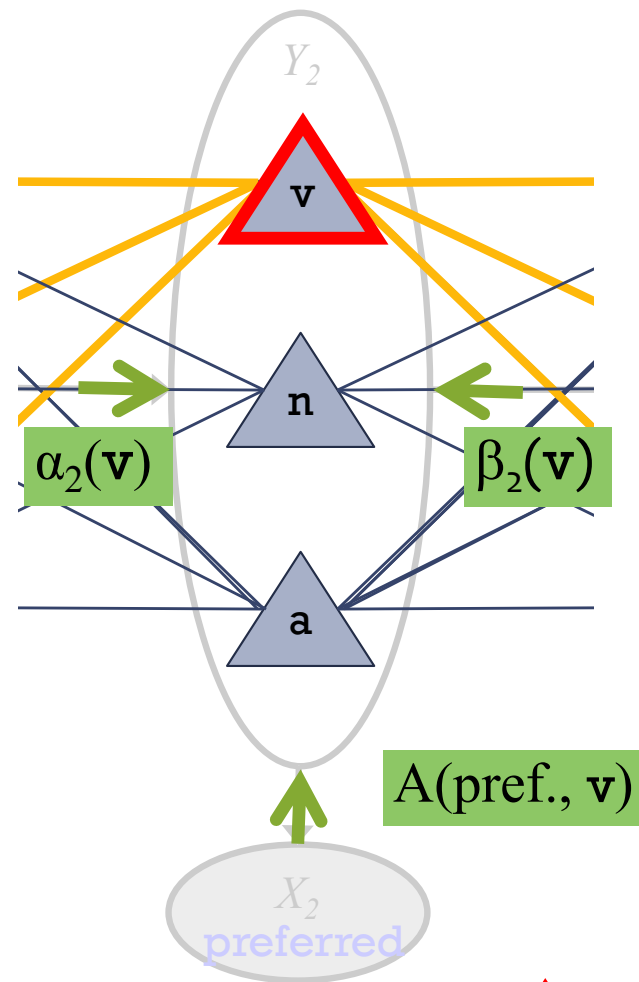


“belief that  $Y_2 = n$ ”

total weight of *all* paths through 

$$= \alpha_2(n) A(\text{pref.}, n) \beta_2(n)$$

# Forward-Backward Algorithm: Finds Marginals



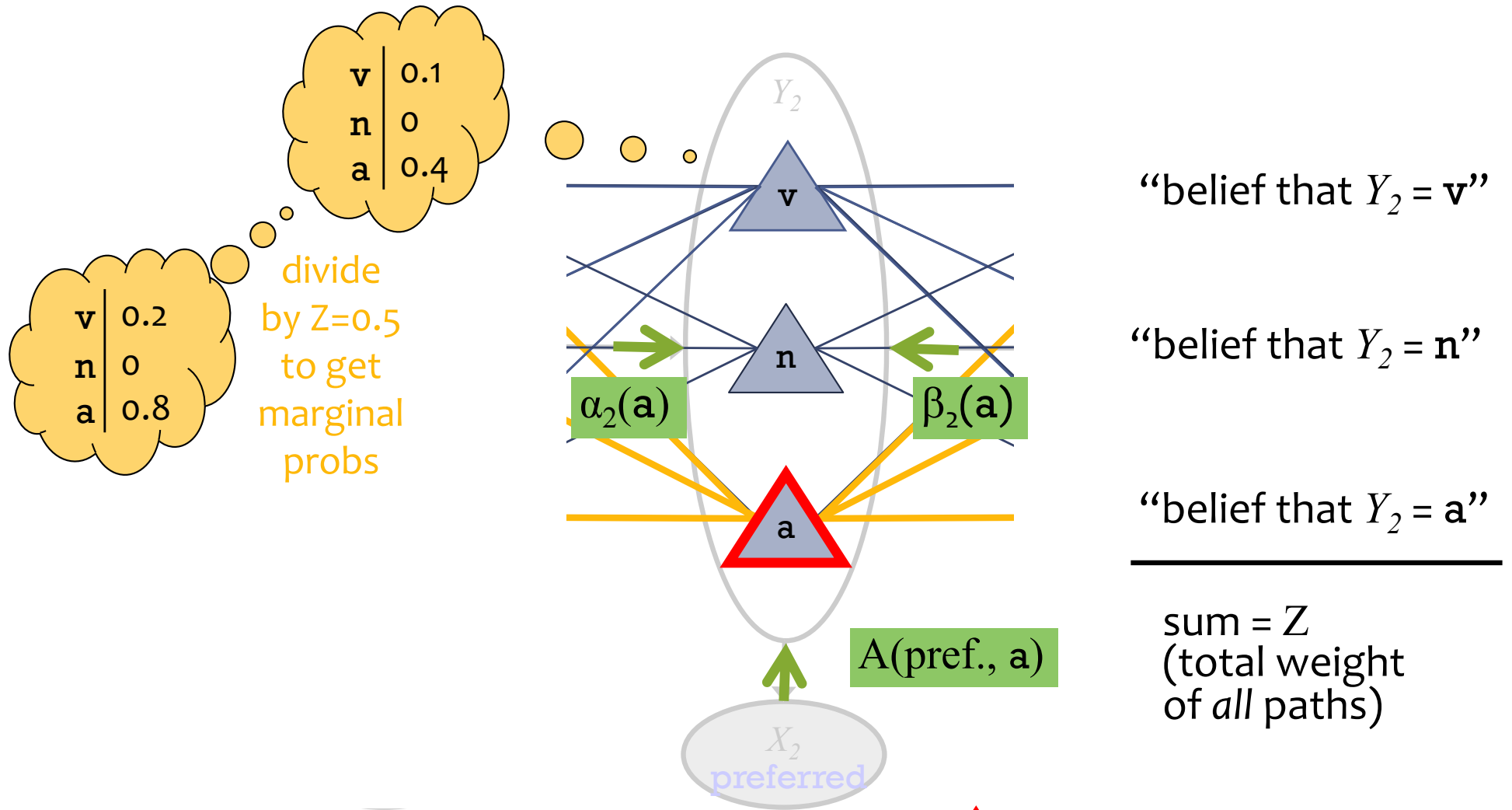
“belief that  $Y_2 = \mathbf{v}$ ”

“belief that  $Y_2 = \mathbf{n}$ ”

total weight of *all* paths through 

$$= \alpha_2(\mathbf{v}) \ A(\text{pref.}, \mathbf{v}) \ \beta_2(\mathbf{v})$$

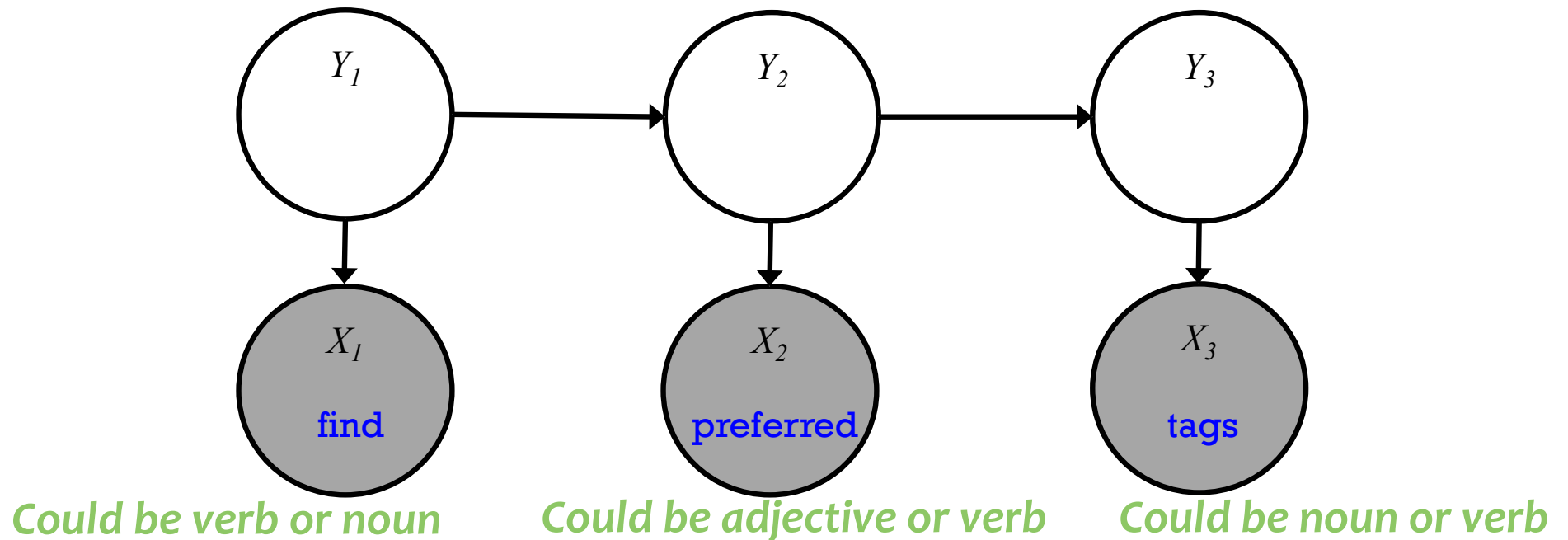
# Forward-Backward Algorithm: Finds Marginals



total weight of *all* paths through 

$$= \alpha_2(\mathbf{a}) \quad A(\text{pref.}, \mathbf{a}) \quad \beta_2(\mathbf{a})$$

# Forward-Backward Algorithm



# Inference for HMMs

## *Whiteboard*

- Derivation of Forward algorithm
- Forward-backward algorithm
- Viterbi algorithm

# Derivation of Forward Algorithm

Definition:  $\alpha_t(k) \triangleq p(x_1, \dots, x_t, y_t = k)$

Derivation:

$$\begin{aligned}
 \alpha_T(\text{END}) &= p(x_1, \dots, x_T, y_T = \text{END}) \\
 &= p(x_1, \dots, x_T | y_T) p(y_T) \quad \leftarrow \text{by def. of joint} \\
 &= p(x_T | y_T) p(x_1, \dots, x_{T-1} | y_T) p(y_T) \quad \leftarrow \text{by cond. indep. of HMM} \\
 &= p(x_T | y_T) p(x_1, \dots, x_{T-1}, y_T) \quad \leftarrow \text{by def. of joint} \\
 &= p(x_T | y_T) \sum_{y_{T-1}} p(x_1, \dots, x_{T-1}, y_{T-1}, y_T) \quad \leftarrow \text{by def. of marginal} \\
 &= p(x_T | y_T) \sum_{y_{T-1}} p(x_1, \dots, x_{T-1}, y_T | y_{T-1}) p(y_{T-1}) \quad \leftarrow \text{by def. of joint} \\
 &= p(x_T | y_T) \sum_{y_{T-1}} \underbrace{p(x_1, \dots, x_{T-1} | y_{T-1}) p(y_T | y_{T-1})}_{\leftarrow \text{by cond. indep. of HMM}} p(y_{T-1}) \\
 &= p(x_T | y_T) \sum_{y_{T-1}} p(x_1, \dots, x_{T-1}, y_{T-1}) p(y_T | y_{T-1}) \quad \leftarrow \text{by def. of joint} \\
 &= p(x_T | y_T) \sum_{y_{T-1}} \alpha_{T-1}(y_{T-1}) p(y_T | y_{T-1}) \quad \leftarrow \text{by def. of } \alpha_t(k)
 \end{aligned}$$

Herein using " $y_T$ " as shorthand for " $y_T = \text{END}$ "

# Forward-Backward Algorithm

Define:

$$\alpha_t(k) \triangleq p(x_1, \dots, x_t, y_t = k)$$

$$\beta_t(k) \triangleq p(x_{t+1}, \dots, x_T | y_t = k)$$

Assume

$$y_0 = \text{START}$$

$$y_{T+1} = \text{END}$$

① Initialize

$$\alpha_0(\text{START}) = 1 \quad \alpha_0(k) = 0 \quad \forall k \neq \text{START}$$

$$\beta_{T+1}(\text{END}) = 1 \quad \beta_{T+1}(k) = 0 \quad \forall k \neq \text{END}$$

② For  $t = 1, \dots, T$ :

For  $k = 1, \dots, K$ :

$$\alpha_t(k) = p(x_t | y_t = k) \sum_{j=1}^K \alpha_{t-1}(j) p(y_t = k | y_{t-1} = j)$$

the alphas include the emission probabilities  
so we don't multiply them in separately

③ For  $t = T, \dots, T$ :

For  $k = 1, \dots, K$ :

$$\beta_t(k) = \sum_{j=1}^K p(x_{t+1} | y_{t+1} = j) \beta_{t+1}(j) p(y_{t+1} = j | y_t = k)$$

④ Compute  $p(\vec{x}) = \alpha_{T+1}(\text{END})$

[Evaluation]

⑤ Compute  $p(y_t = k | \vec{x}) = \frac{\alpha_t(k) \beta_t(k)}{p(\vec{x})}$

[Marginals]



# Viterbi Algorithm

Define:  $\omega_t(k) \triangleq \max_{y_1, \dots, y_{t-1}} p(x_1, \dots, x_t, y_1, \dots, y_{t-1}, y_t = k)$

"backpointers"  $\rightarrow b_t(k) \triangleq \arg \max_{y_1, \dots, y_{t-1}} p(x_1, \dots, x_t, y_1, \dots, y_{t-1}, y_t = k)$

Assume  $y_0 = \text{START}$

① Initialize  $\omega_0(\text{START}) = 1$   $\omega_0(k) = 0 \ \forall k \neq \text{START}$

② For  $t = 1, \dots, T$ :

For  $k = 1, \dots, K$ :

$$\omega_t(k) = \max_{j \in \{1, \dots, K\}} p(x_t | y_t = k) \omega_{t-1}(j) p(y_t = k | y_{t-1} = j)$$

$$b_t(k) = \arg \max_{j \in \{1, \dots, K\}} p(x_t | y_t = k) \omega_{t-1}(j) p(y_t = k | y_{t-1} = j)$$

③ Compute Most Probable Assignment

$$\hat{y}_T = b_{T+1}(\text{END})$$

For  $t = T-1, \dots, 1$

$$\hat{y}_t = b_{t+1}(\hat{y}_{t+1})$$

[Decoding]

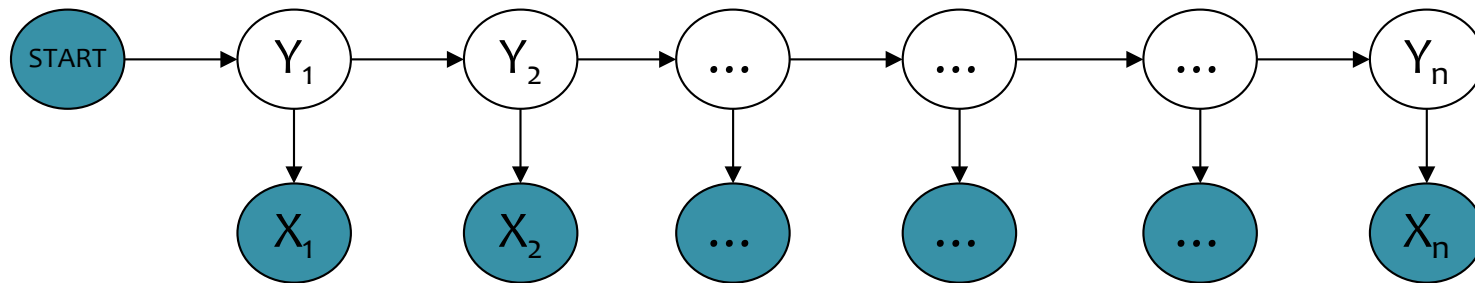
follow the  
"backpointers"

# Inference in HMMs

What is the **computational complexity** of inference for HMMs?

- The **naïve** (brute force) computations for *Evaluation*, *Decoding*, and *Marginals* take **exponential time**,  $O(K^T)$
- The **forward-backward** algorithm and **Viterbi** algorithm run in **polynomial time**,  $O(T \cdot K^2)$ 
  - Thanks to dynamic programming!

# Shortcomings of Hidden Markov Models



- HMM models capture dependences between each state and **only** its corresponding observation
  - NLP example: In a sentence segmentation task, each segmental state may depend not just on a single word (and the adjacent segmental stages), but also on the (non-local) features of the whole line such as line length, indentation, amount of white space, etc.
- Mismatch between learning objective function and prediction objective function
  - HMM learns a joint distribution of states and observations  $P(\mathbf{Y}, \mathbf{X})$ , but in a prediction task, we need the conditional probability  $P(\mathbf{Y}|\mathbf{X})$

# **MBR DECODING**

# Inference for HMMs

- ~~Four~~
- ~~Three~~ Inference Problems for an HMM
    1. Evaluation: Compute the probability of a given sequence of observations
    2. Viterbi Decoding: Find the most-likely sequence of hidden states, given a sequence of observations
    3. Marginals: Compute the marginal distribution for a hidden state, given a sequence of observations
    4. MBR Decoding: Find the lowest loss sequence of hidden states, given a sequence of observations (Viterbi decoding is a special case)

# Minimum Bayes Risk Decoding

- Suppose we given a loss function  $l(\mathbf{y}', \mathbf{y})$  and are asked for a single tagging
- How should we choose just one from our probability distribution  $p(\mathbf{y}|\mathbf{x})$ ?
- A minimum Bayes risk (MBR) decoder  $h(\mathbf{x})$  returns the variable assignment with minimum **expected** loss under the model's distribution

$$\begin{aligned} h_{\theta}(\mathbf{x}) &= \operatorname{argmin}_{\hat{\mathbf{y}}} \mathbb{E}_{\mathbf{y} \sim p_{\theta}(\cdot|\mathbf{x})} [\ell(\hat{\mathbf{y}}, \mathbf{y})] \\ &= \operatorname{argmin}_{\hat{\mathbf{y}}} \sum_{\mathbf{y}} p_{\theta}(\mathbf{y} | \mathbf{x}) \ell(\hat{\mathbf{y}}, \mathbf{y}) \end{aligned}$$

# Minimum Bayes Risk Decoding

$$h_{\theta}(\mathbf{x}) = \operatorname{argmin}_{\hat{\mathbf{y}}} \mathbb{E}_{\mathbf{y} \sim p_{\theta}(\cdot | \mathbf{x})} [\ell(\hat{\mathbf{y}}, \mathbf{y})]$$

Consider some example loss functions:

The **0-1 loss function** returns 1 only if the two assignments are identical and 0 otherwise:

$$\ell(\hat{\mathbf{y}}, \mathbf{y}) = 1 - \mathbb{I}(\hat{\mathbf{y}}, \mathbf{y})$$

The MBR decoder is:

$$\begin{aligned} h_{\theta}(\mathbf{x}) &= \operatorname{argmin}_{\hat{\mathbf{y}}} \sum_{\mathbf{y}} p_{\theta}(\mathbf{y} | \mathbf{x}) (1 - \mathbb{I}(\hat{\mathbf{y}}, \mathbf{y})) \\ &= \operatorname{argmax}_{\hat{\mathbf{y}}} p_{\theta}(\hat{\mathbf{y}} | \mathbf{x}) \end{aligned}$$

which is exactly the Viterbi decoding problem!

# Minimum Bayes Risk Decoding

$$h_{\theta}(\mathbf{x}) = \operatorname{argmin}_{\hat{\mathbf{y}}} \mathbb{E}_{\mathbf{y} \sim p_{\theta}(\cdot | \mathbf{x})} [\ell(\hat{\mathbf{y}}, \mathbf{y})]$$

Consider some example loss functions:

The **Hamming loss** corresponds to accuracy and returns the number of incorrect variable assignments:

$$\ell(\hat{\mathbf{y}}, \mathbf{y}) = \sum_{i=1}^V (1 - \mathbb{I}(\hat{y}_i, y_i))$$

The MBR decoder is:

$$\hat{y}_i = h_{\theta}(\mathbf{x})_i = \operatorname{argmax}_{\hat{y}_i} p_{\theta}(\hat{y}_i | \mathbf{x})$$

This decomposes across variables and requires the variable marginals.



# Learning Objectives

## Hidden Markov Models

*You should be able to...*

1. Show that structured prediction problems yield high-computation inference problems
2. Define the first order Markov assumption
3. Draw a Finite State Machine depicting a first order Markov assumption
4. Derive the MLE parameters of an HMM
5. Define the three key problems for an HMM: evaluation, decoding, and marginal computation
6. Derive a dynamic programming algorithm for computing the marginal probabilities of an HMM
7. Interpret the forward-backward algorithm as a message passing algorithm
8. Implement supervised learning for an HMM
9. Implement the forward-backward algorithm for an HMM
10. Implement the Viterbi algorithm for an HMM
11. Implement a minimum Bayes risk decoder with Hamming loss for an HMM