



# Belief Propagation + Learning fully observable MRFs and CRFs



Matt Gormley  
Lecture 9  
Sep. 28, 2022

# Reminders

- **Homework 2: Learning to Search for RNNs**
  - Out: Sun, Sep 18
  - **Written (except for Empirical Questions)**
    - **Due: Thu, Sep 29 at 11:59pm**
  - **Programming + Empirical Questions**
    - **Due: NEVER?**
- **Homework 3: General Graph CRF Module**
  - Out: Thu, Sep 29
  - Due: Mon, Oct 10 at 11:59pm

# Reminders

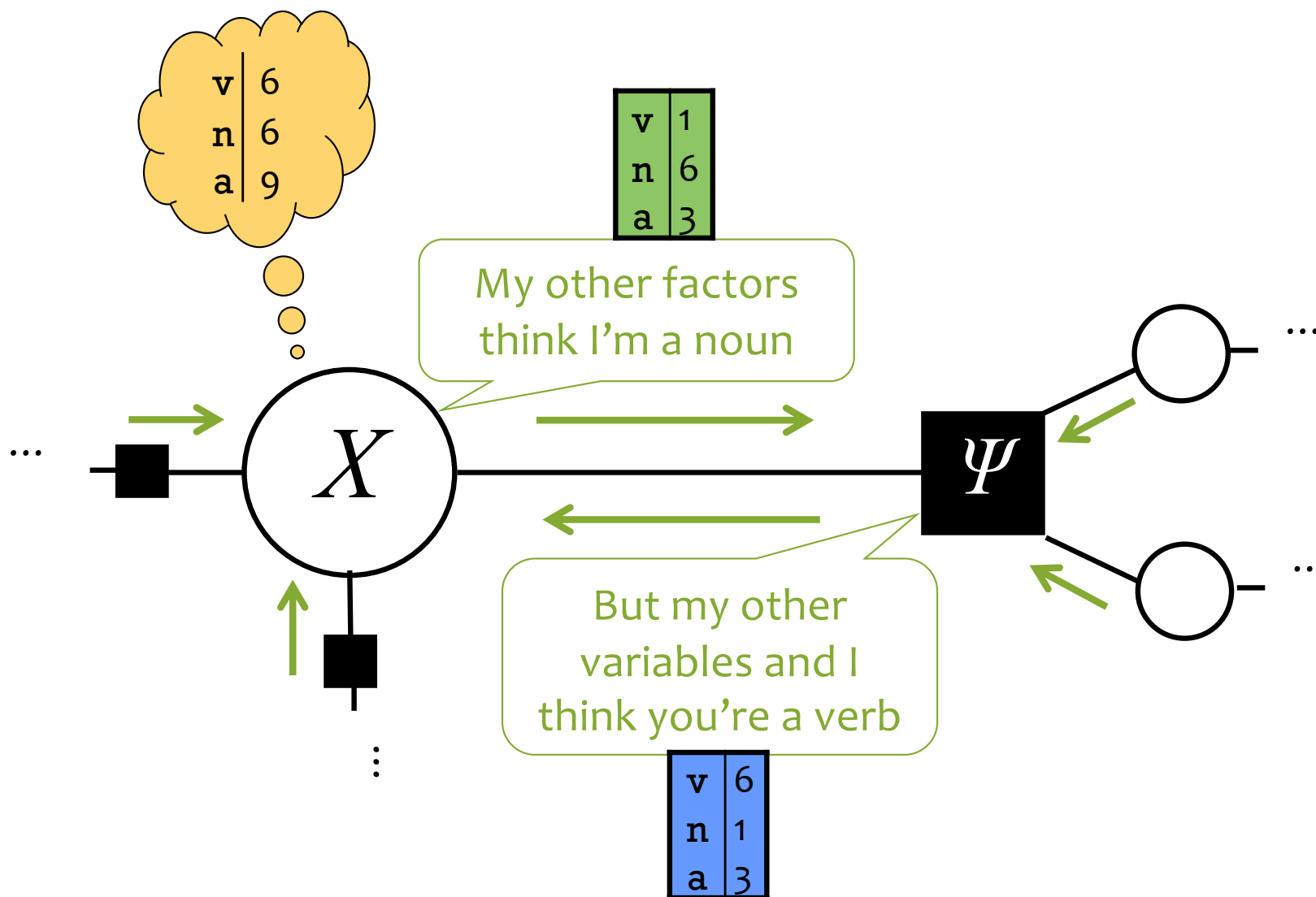
- **Homework 2: Learning to Search for RNNs**
  - Out: Sun, Sep 18
  - **Written (except for Empirical Questions)**
    - Due: Thu, Sep 29 at 11:59pm
  - **Programming + Empirical Questions**
    - Due: Mon, Oct 24 at 9:00am
- **Homework 3: General Graph CRF Module**
  - Out: Thu, Sep 29
  - Due: Mon, Oct 10 at 11:59pm

Exact marginal inference for factor trees

# **SUM-PRODUCT BELIEF PROPAGATION**

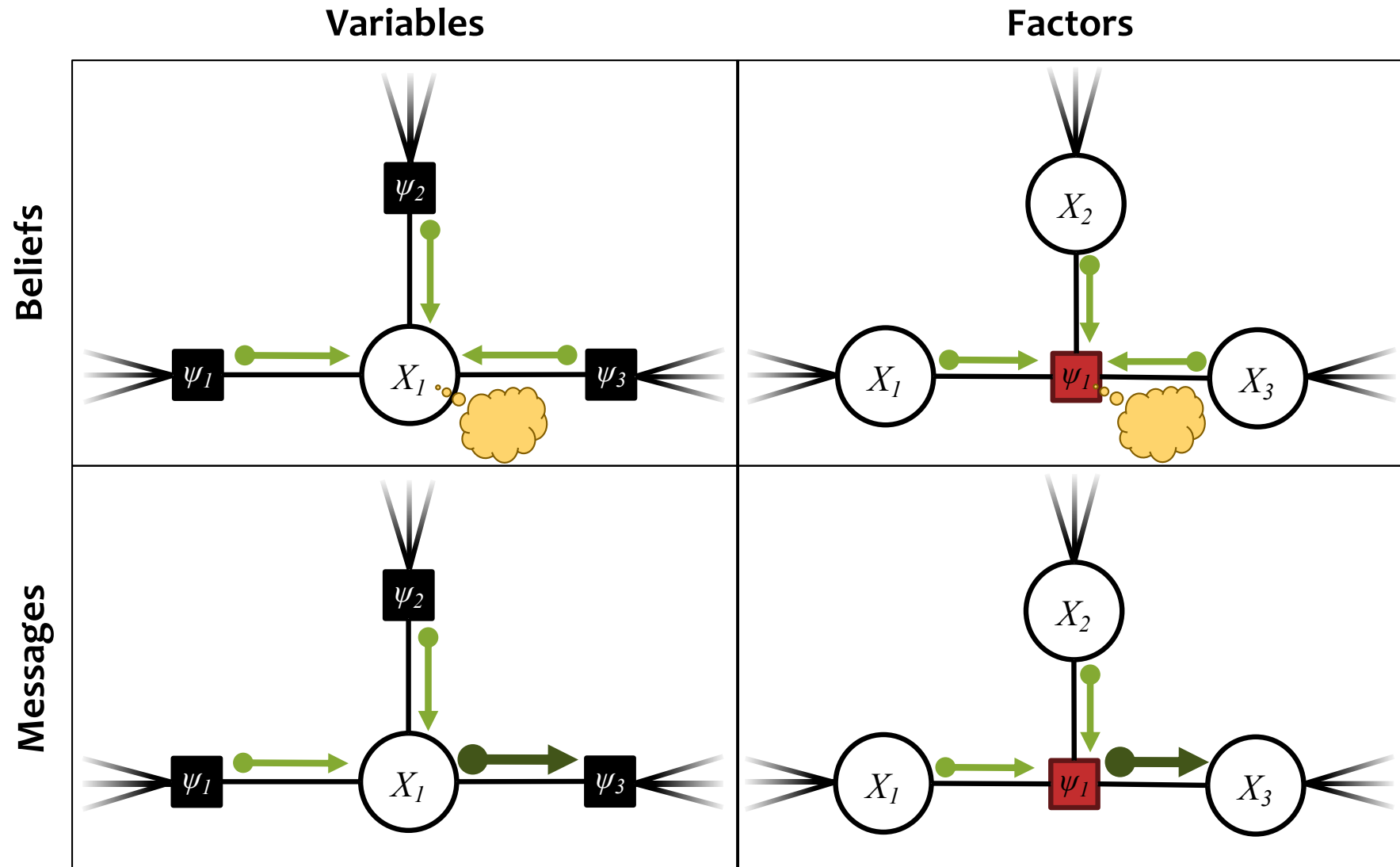


# Message Passing in Belief Propagation



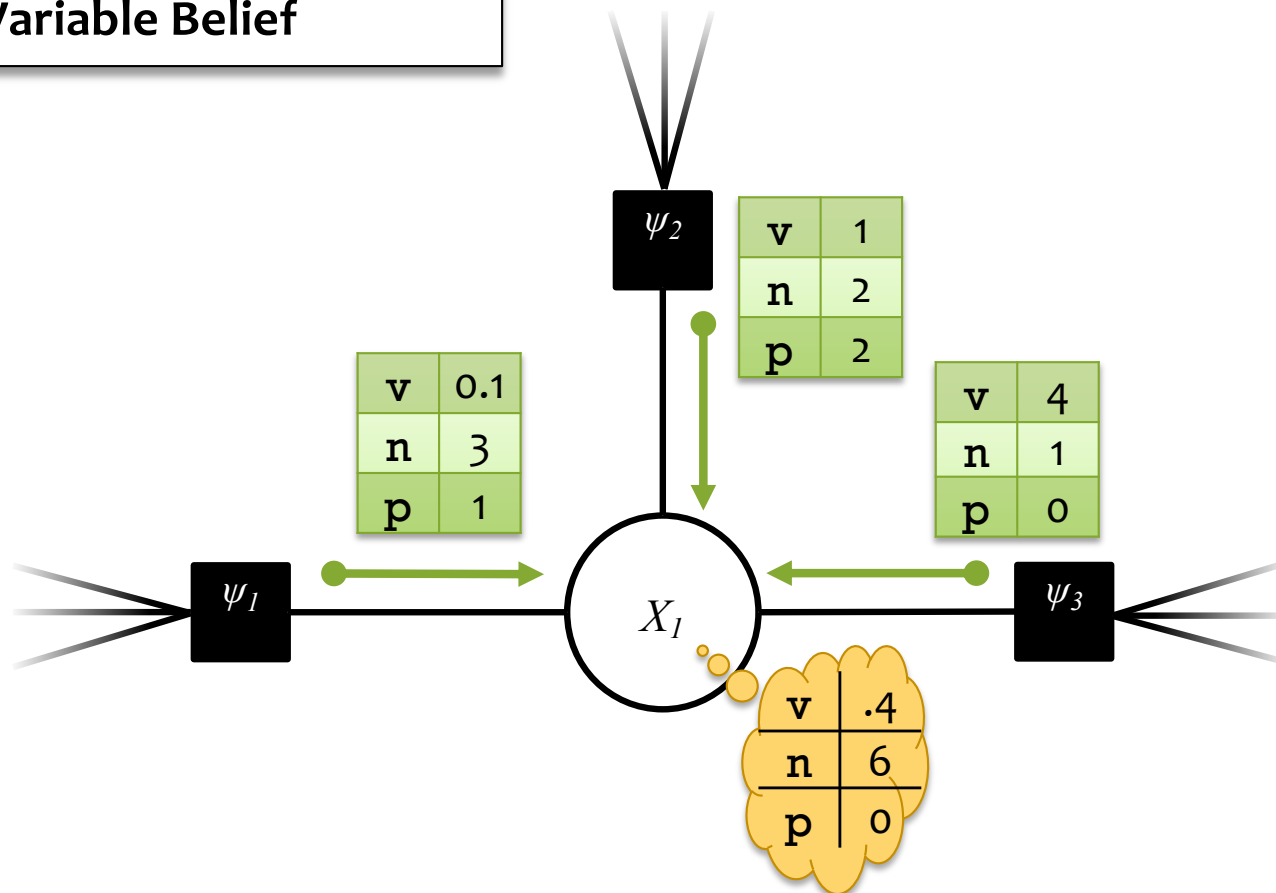
Both of these messages judge the possible values of variable  $X$ .  
Their product = belief at  $X$  = product of all 3 messages to  $X$ .

# Sum-Product Belief Propagation



# Sum-Product Belief Propagation

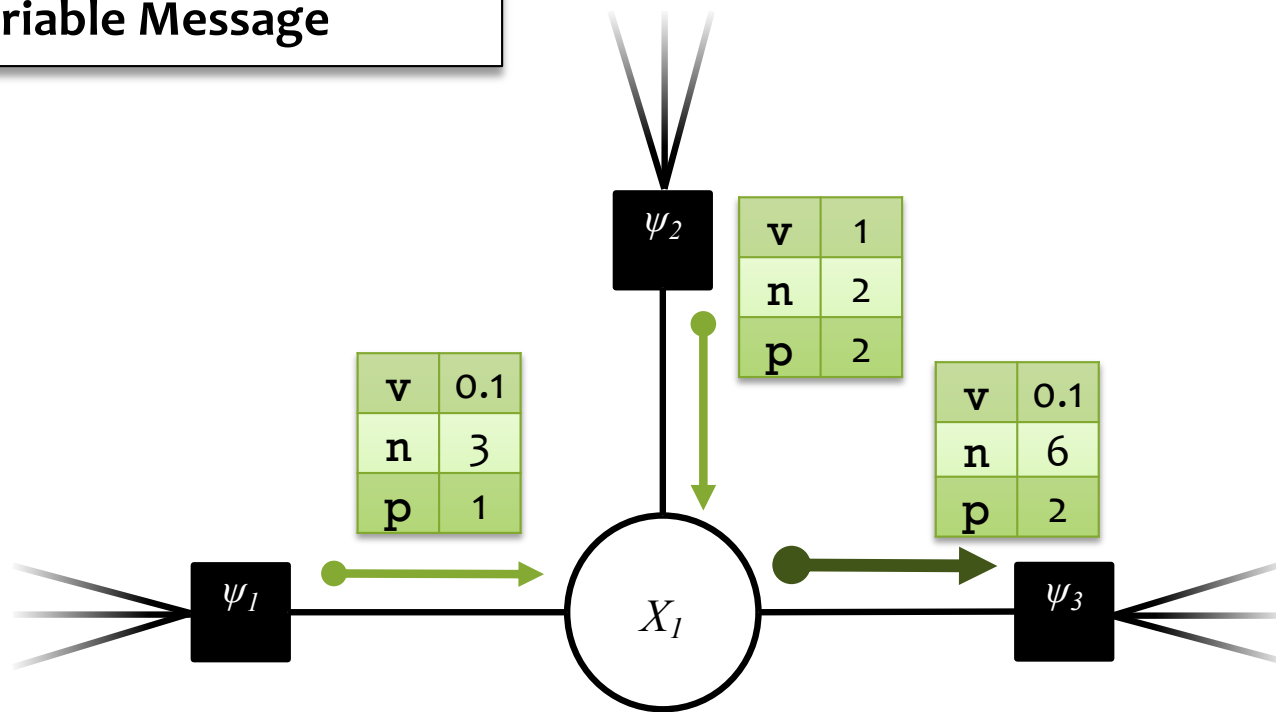
Variable Belief



$$b_i(x_i) = \prod_{\alpha \in \mathcal{N}(i)} \mu_{\alpha \rightarrow i}(x_i)$$

# Sum-Product Belief Propagation

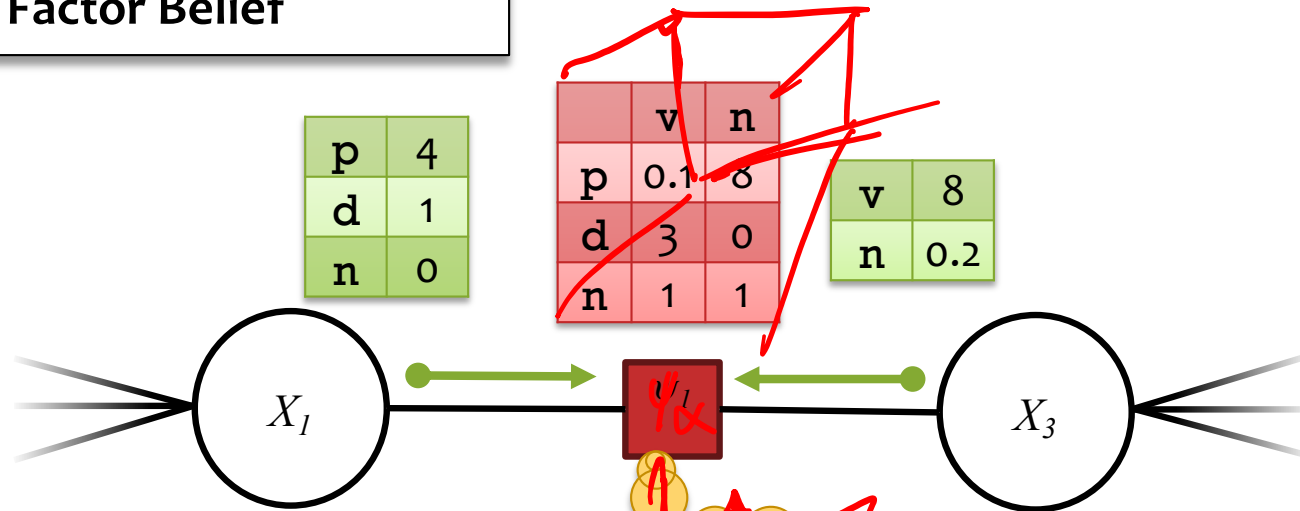
Variable Message



$$\mu_{i \rightarrow \alpha}(x_i) = \prod_{\alpha \in \mathcal{N}(i) \setminus \alpha} \mu_{\alpha \rightarrow i}(x_i)$$

# Sum-Product Belief Propagation

Factor Belief



$$X_\alpha = [p, a, v]$$

$$X_\alpha[1] = p$$

$$X_\alpha[2] = a$$

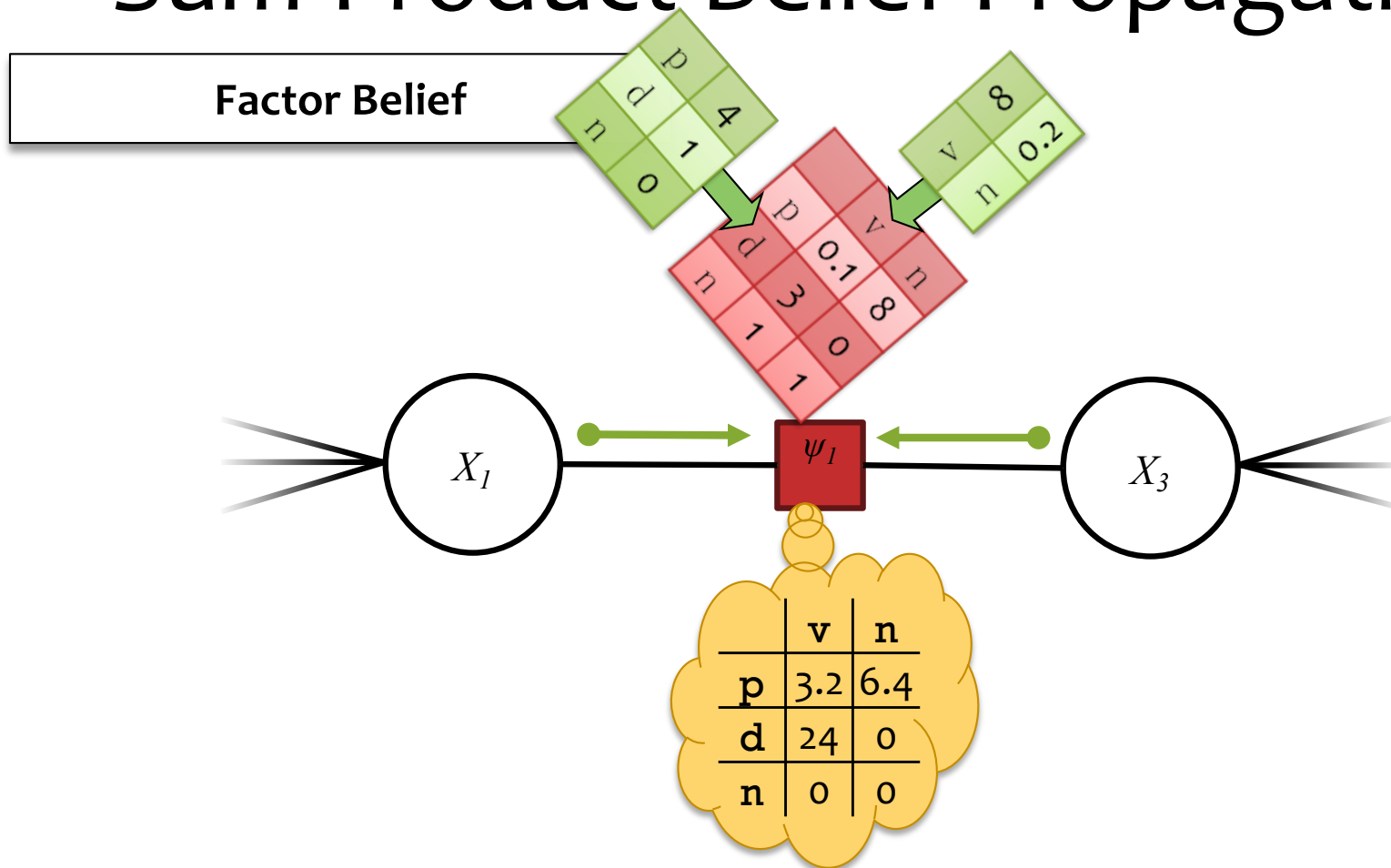
$$X_\alpha[3] = v$$

$$\alpha = \{1, 2, 3\}$$

$$X_\alpha = \{X_1, X_2, X_3\}$$

$$b_\alpha(x_\alpha) = \psi_\alpha(x_\alpha) \prod_{i \in \mathcal{N}(\alpha)} \mu_{i \rightarrow \alpha}(x_\alpha[i])$$

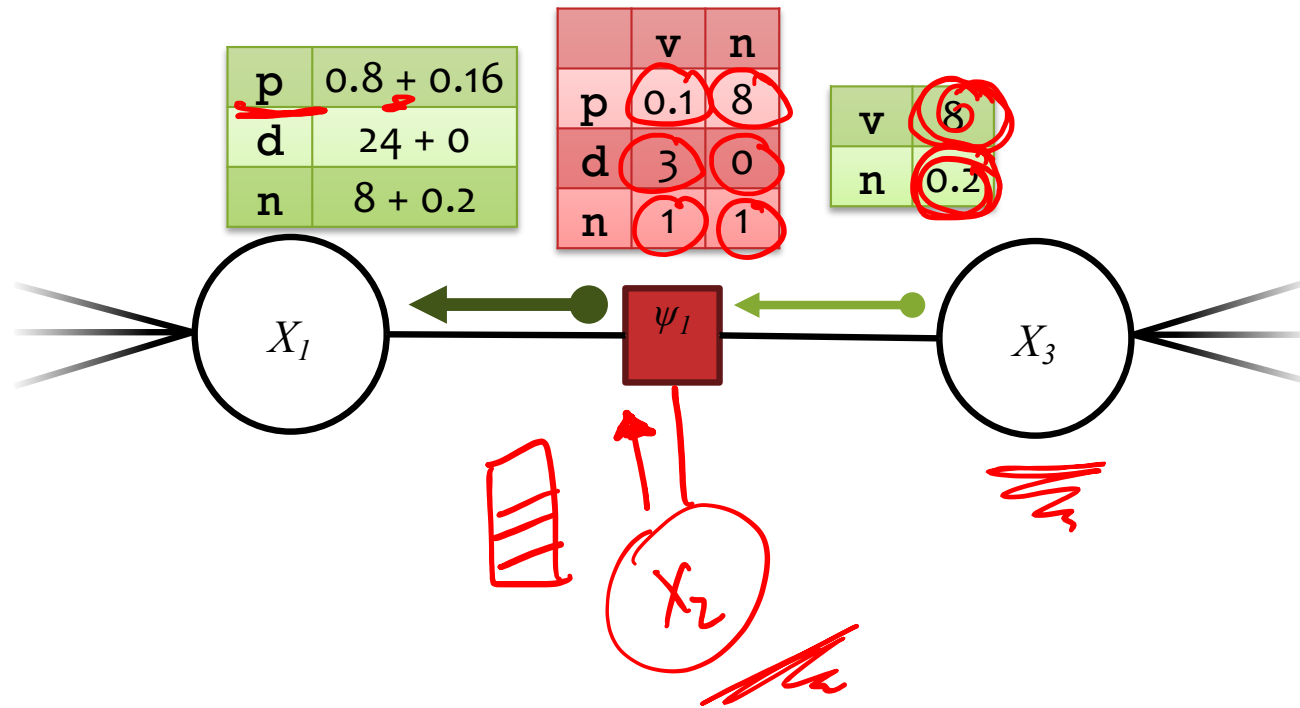
# Sum-Product Belief Propagation



$$b_{\alpha}(\mathbf{x}_{\alpha}) = \psi_{\alpha}(\mathbf{x}_{\alpha}) \prod_{i \in \mathcal{N}(\alpha)} \mu_{i \rightarrow \alpha}(\mathbf{x}_{\alpha}[i])$$

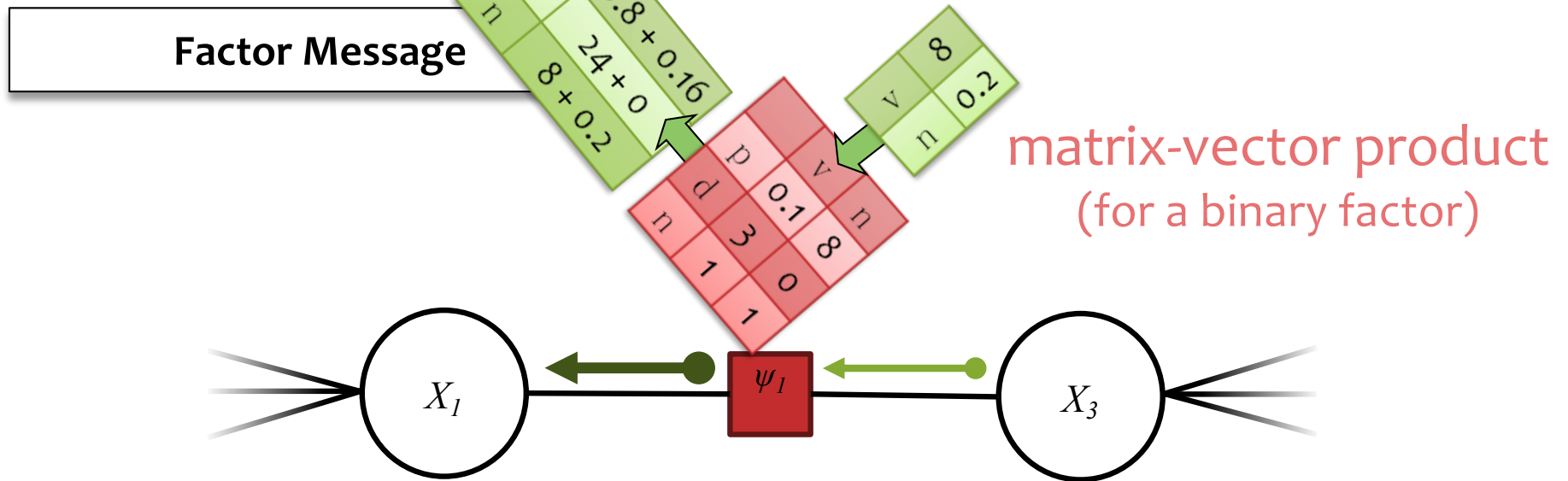
# Sum-Product Belief Propagation

Factor Message



$$\mu_{\alpha \rightarrow i}(x_i) = \sum_{\mathbf{x}_\alpha : \mathbf{x}_\alpha[i] = x_i} \psi_\alpha(\mathbf{x}_\alpha) \prod_{j \in \mathcal{N}(\alpha) \setminus i} \mu_{j \rightarrow \alpha}(\mathbf{x}_\alpha[j])$$

# Sum-Product Belief Propagation



$$\mu_{\alpha \rightarrow i}(x_i) = \sum_{\mathbf{x}_{\alpha} : \mathbf{x}_{\alpha}[i] = x_i} \psi_{\alpha}(\mathbf{x}_{\alpha}) \prod_{j \in \mathcal{N}(\alpha) \setminus i} \mu_{j \rightarrow \alpha}(\mathbf{x}_{\alpha}[j])$$



# Sum-Product Belief Propagation

**Input:** a factor graph with no cycles

**Output:** exact marginals for each variable and factor

## Algorithm:

1. Initialize the messages to the uniform distribution.

$$\mu_{i \rightarrow \alpha}(x_i) = 1 \quad \mu_{\alpha \rightarrow i}(x_i) = 1$$

1. Choose a root node.

2. Send messages from the **leaves** to the **root**.  
Send messages from the **root** to the **leaves**.

$$\mu_{i \rightarrow \alpha}(x_i) = \prod_{\alpha \in \mathcal{N}(i) \setminus \alpha} \mu_{\alpha \rightarrow i}(x_i) \quad \mu_{\alpha \rightarrow i}(x_i) = \sum_{\mathbf{x}_{\alpha} : \mathbf{x}_{\alpha}[i] = x_i} \psi_{\alpha}(\mathbf{x}_{\alpha}) \prod_{j \in \mathcal{N}(\alpha) \setminus i} \mu_{j \rightarrow \alpha}(\mathbf{x}_{\alpha}[j])$$

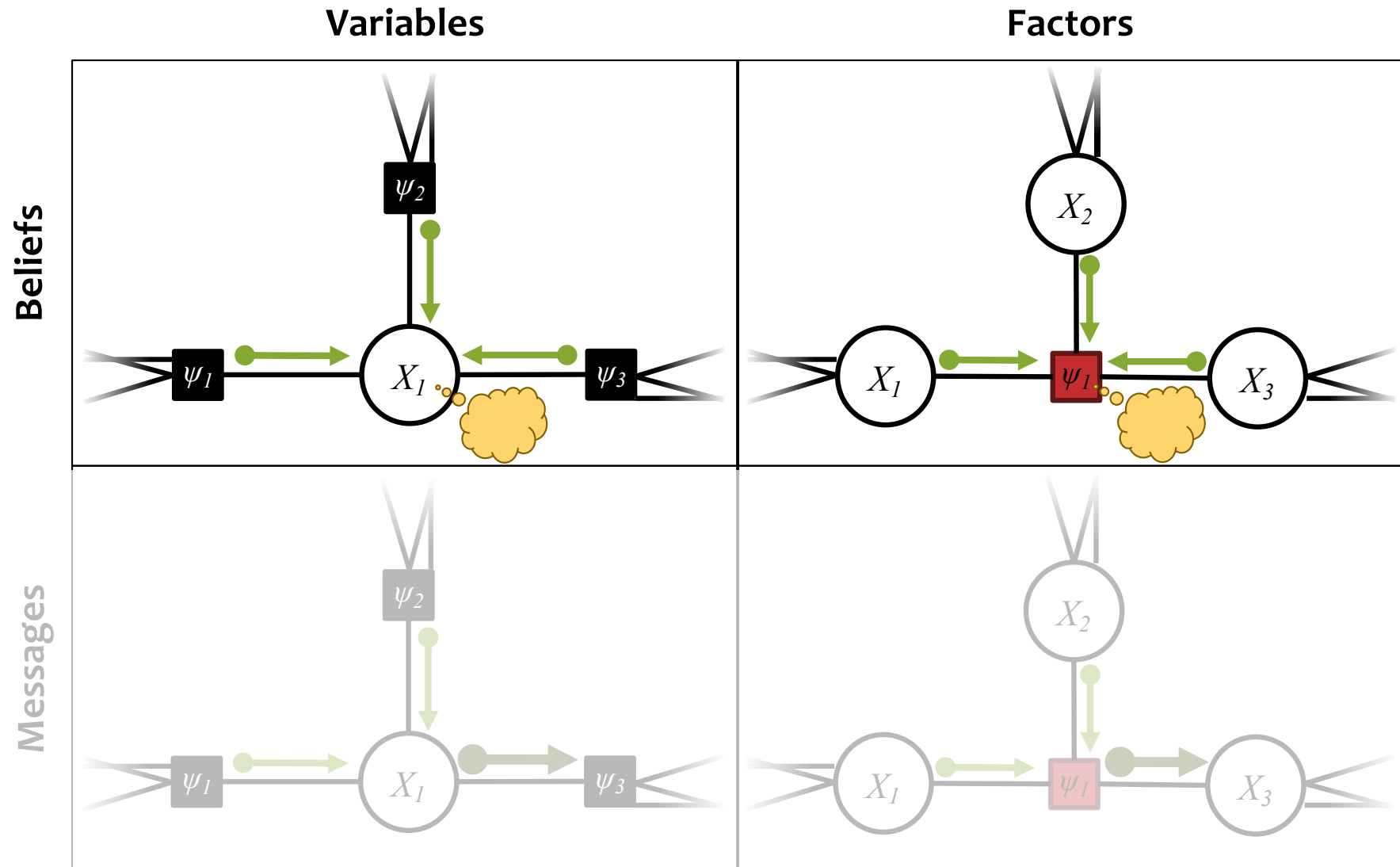
1. Compute the beliefs (unnormalized marginals).

$$b_i(x_i) = \prod_{\alpha \in \mathcal{N}(i)} \mu_{\alpha \rightarrow i}(x_i) \quad b_{\alpha}(\mathbf{x}_{\alpha}) = \psi_{\alpha}(\mathbf{x}_{\alpha}) \prod_{i \in \mathcal{N}(\alpha)} \mu_{i \rightarrow \alpha}(\mathbf{x}_{\alpha}[i])$$

2. Normalize beliefs and return the **exact** marginals.

$$p_i(x_i) \propto b_i(x_i) \quad p_{\alpha}(\mathbf{x}_{\alpha}) \propto b_{\alpha}(\mathbf{x}_{\alpha})$$

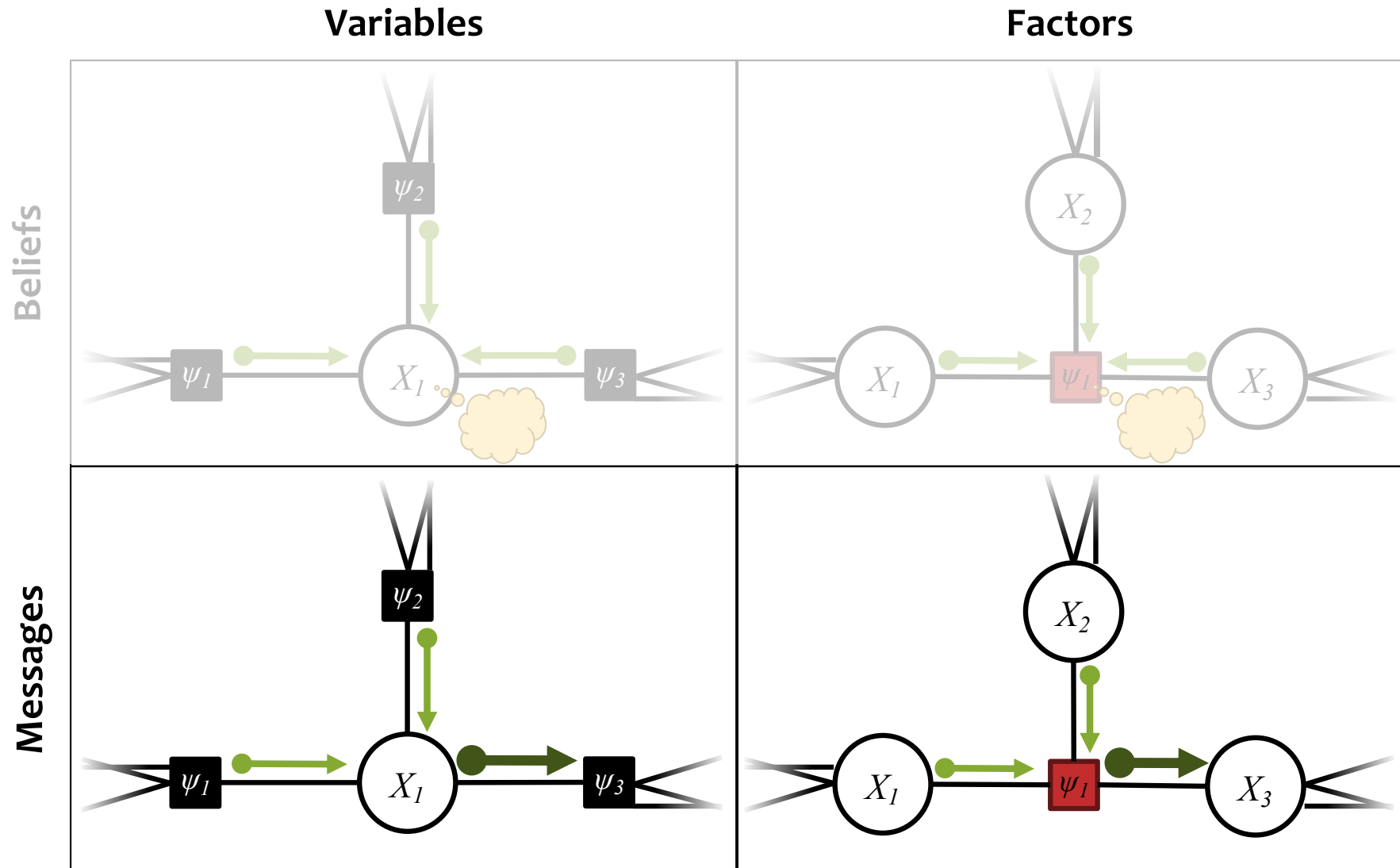
# Sum-Product Belief Propagation



$$b_i(x_i) = \prod_{\alpha \in \mathcal{N}(i)} \mu_{\alpha \rightarrow i}(x_i)$$

$$b_{\alpha}(\mathbf{x}_{\alpha}) = \psi_{\alpha}(\mathbf{x}_{\alpha}) \prod_{i \in \mathcal{N}(\alpha)} \mu_{i \rightarrow \alpha}(\mathbf{x}_{\alpha}[i])$$

# Sum-Product Belief Propagation

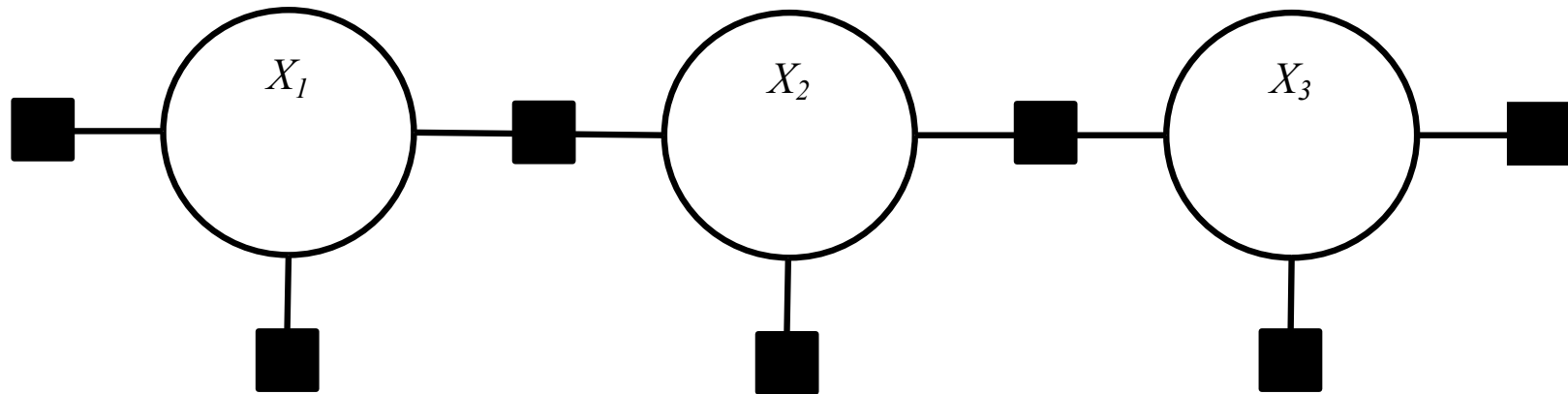


$$\mu_{i \rightarrow \alpha}(x_i) = \prod_{\alpha \in \mathcal{N}(i) \setminus \alpha} \mu_{\alpha \rightarrow i}(x_i)$$

$$\mu_{\alpha \rightarrow i}(x_i) = \sum_{\mathbf{x}_{\alpha} : \mathbf{x}_{\alpha}[i] = x_i} \psi_{\alpha}(\mathbf{x}_{\alpha}) \prod_{j \in \mathcal{N}(\alpha) \setminus i} \mu_{j \rightarrow \alpha}(\mathbf{x}_{\alpha}[j])$$

**FORWARD BACKWARD AS  
SUM-PRODUCT BP**

# CRF Tagging Model



find

preferred

tags

*Could be verb or noun*

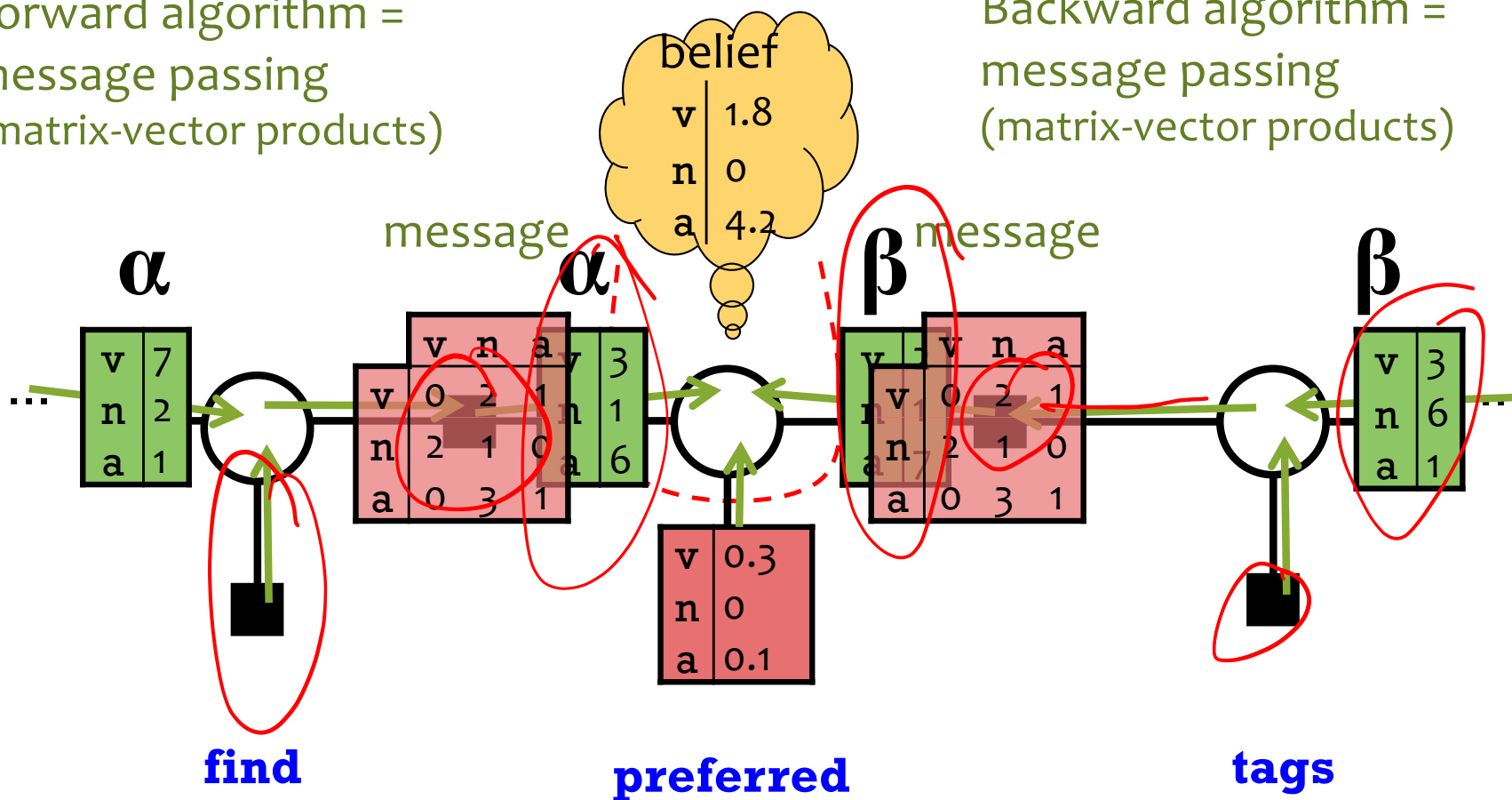
*Could be adjective or verb*

*Could be noun or verb*

# CRF Tagging by Belief Propagation

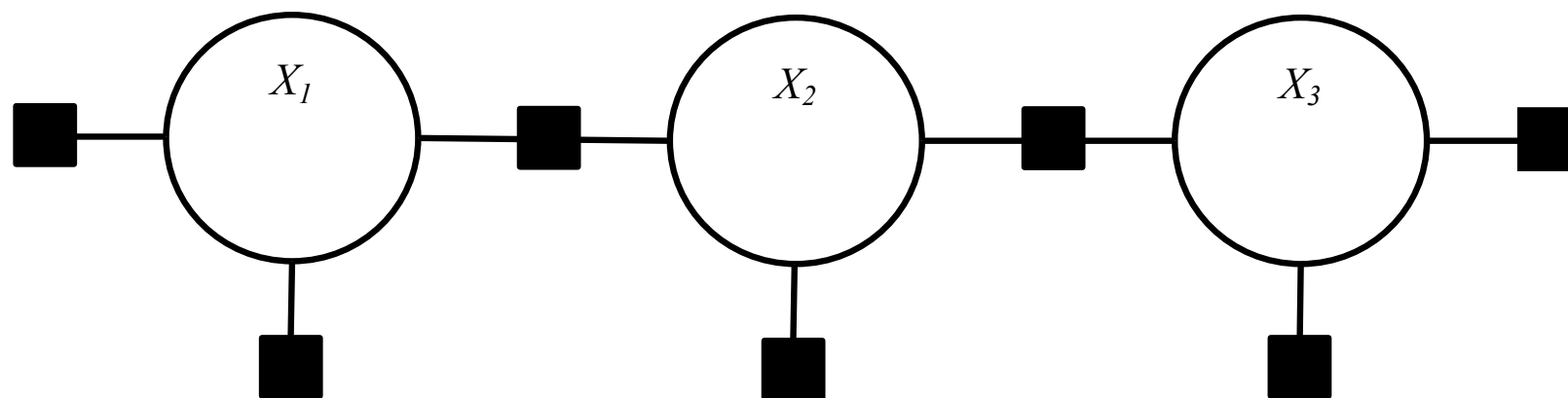
Forward algorithm =  
message passing  
(matrix-vector products)

Backward algorithm =  
message passing  
(matrix-vector products)



- Forward-backward is a message passing algorithm.
- It's the simplest case of belief propagation.

# So Let's Review Forward-Backward ...



find

preferred

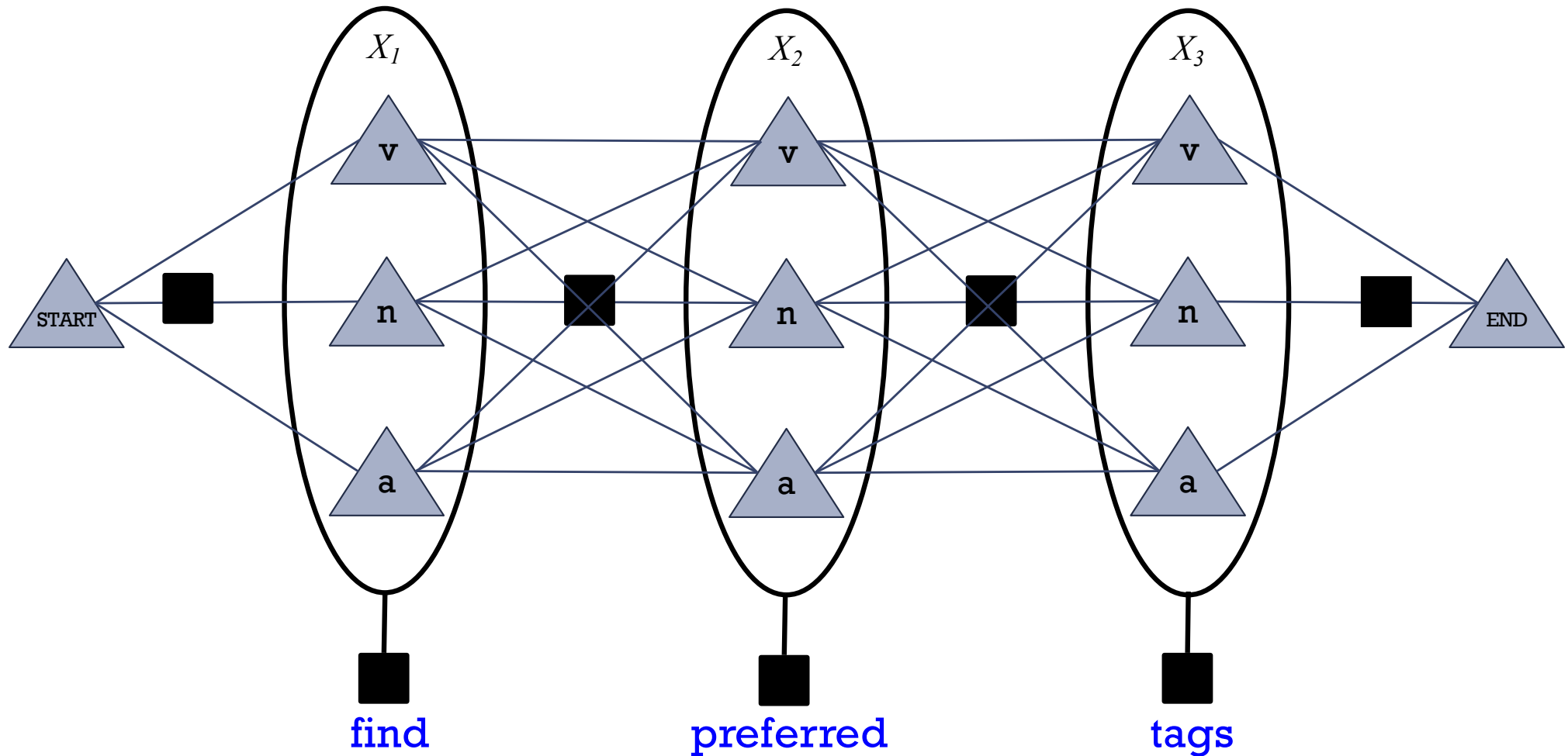
tags

*Could be verb or noun*

*Could be adjective or verb*

*Could be noun or verb*

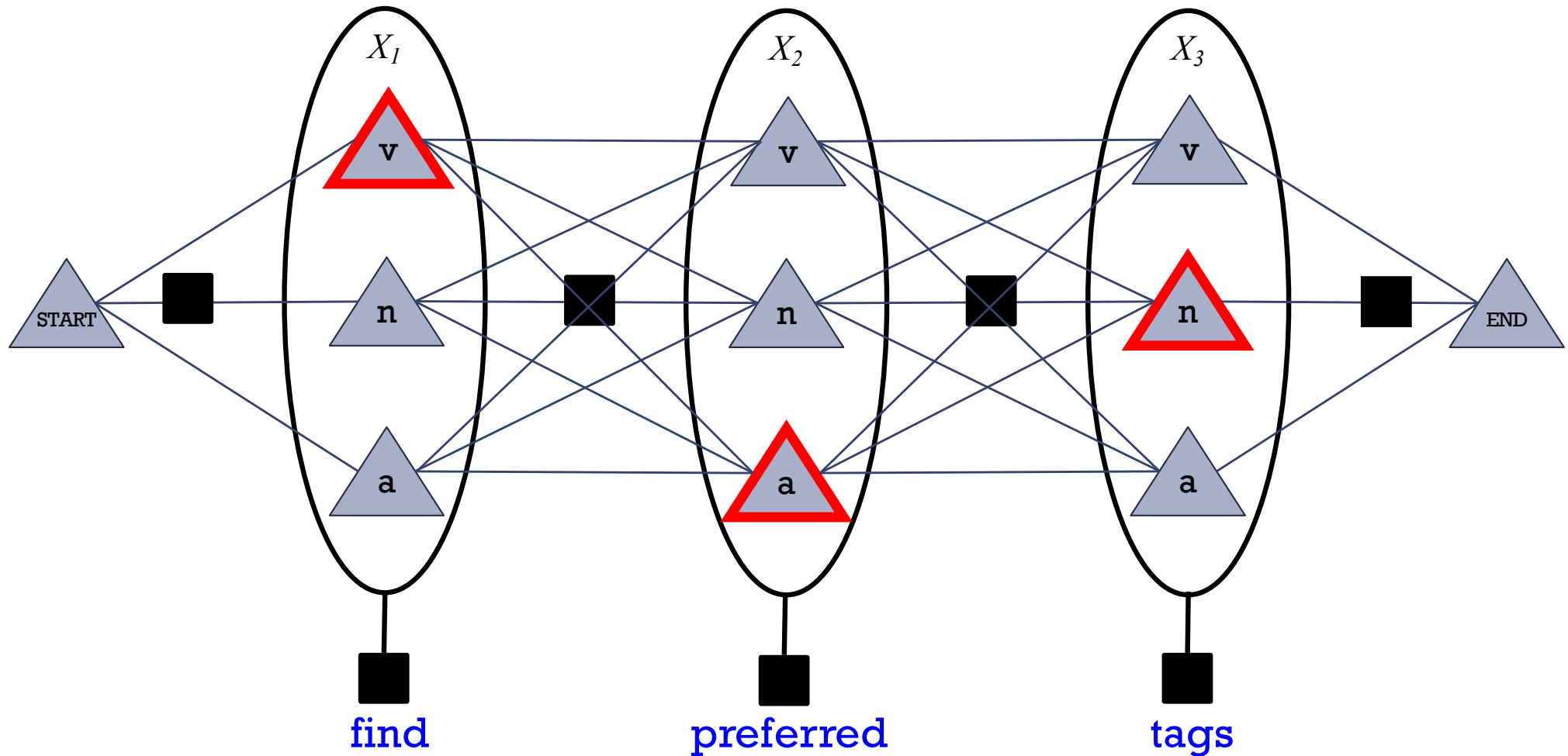
# So Let's Review Forward-Backward ...



- Show the possible *values* for each variable

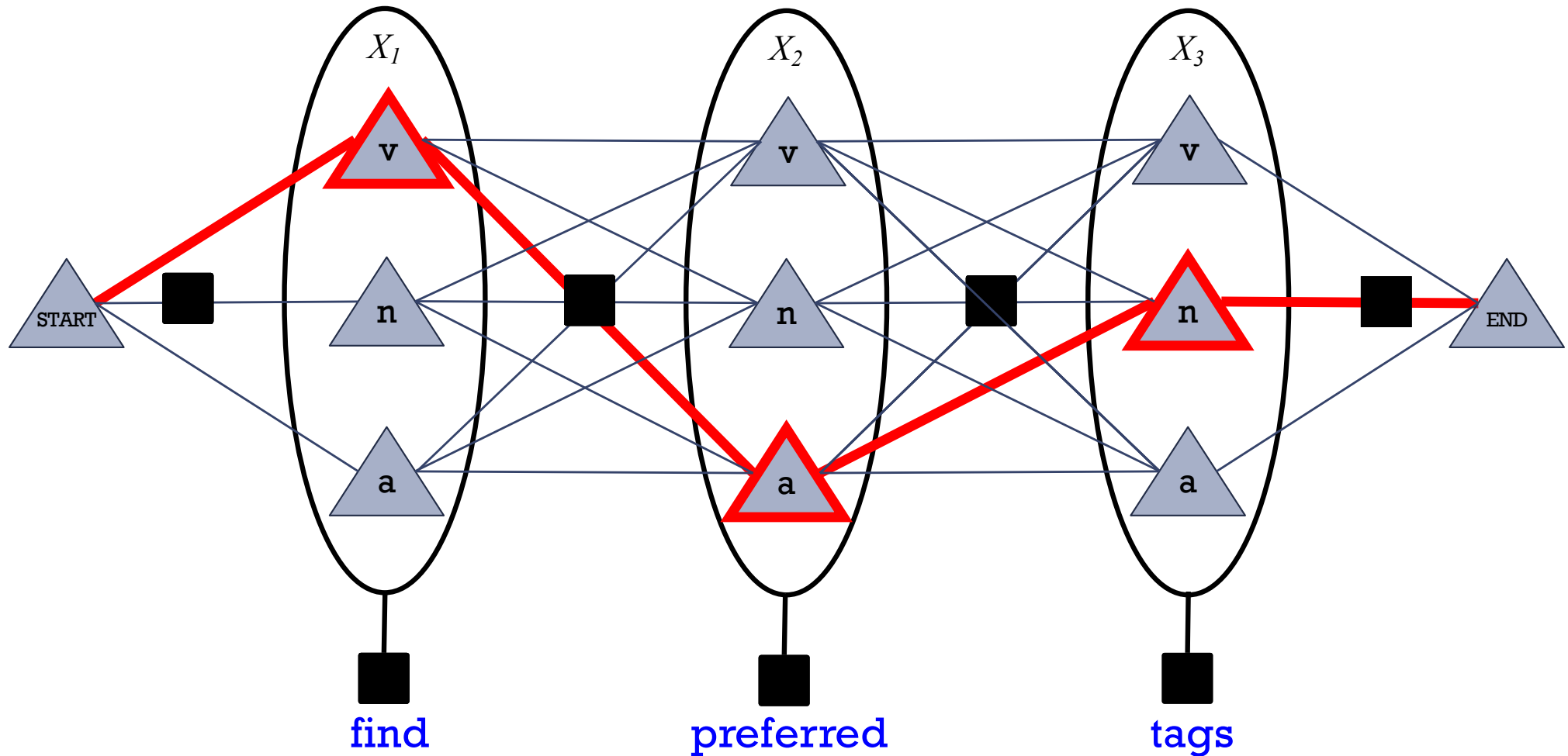


# So Let's Review Forward-Backward ...



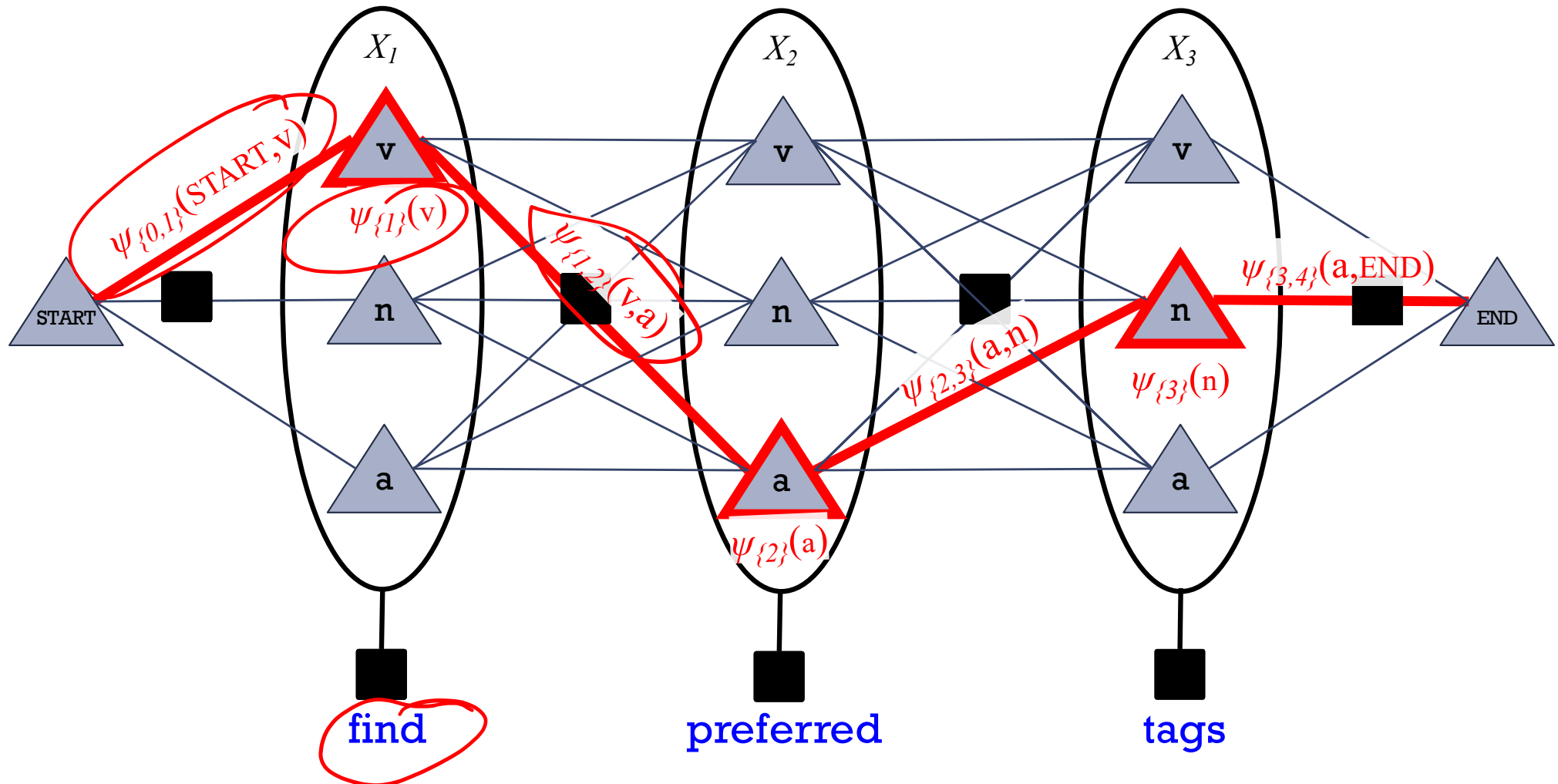
- Let's show the possible *values* for each variable
- One possible assignment

# So Let's Review Forward-Backward ...



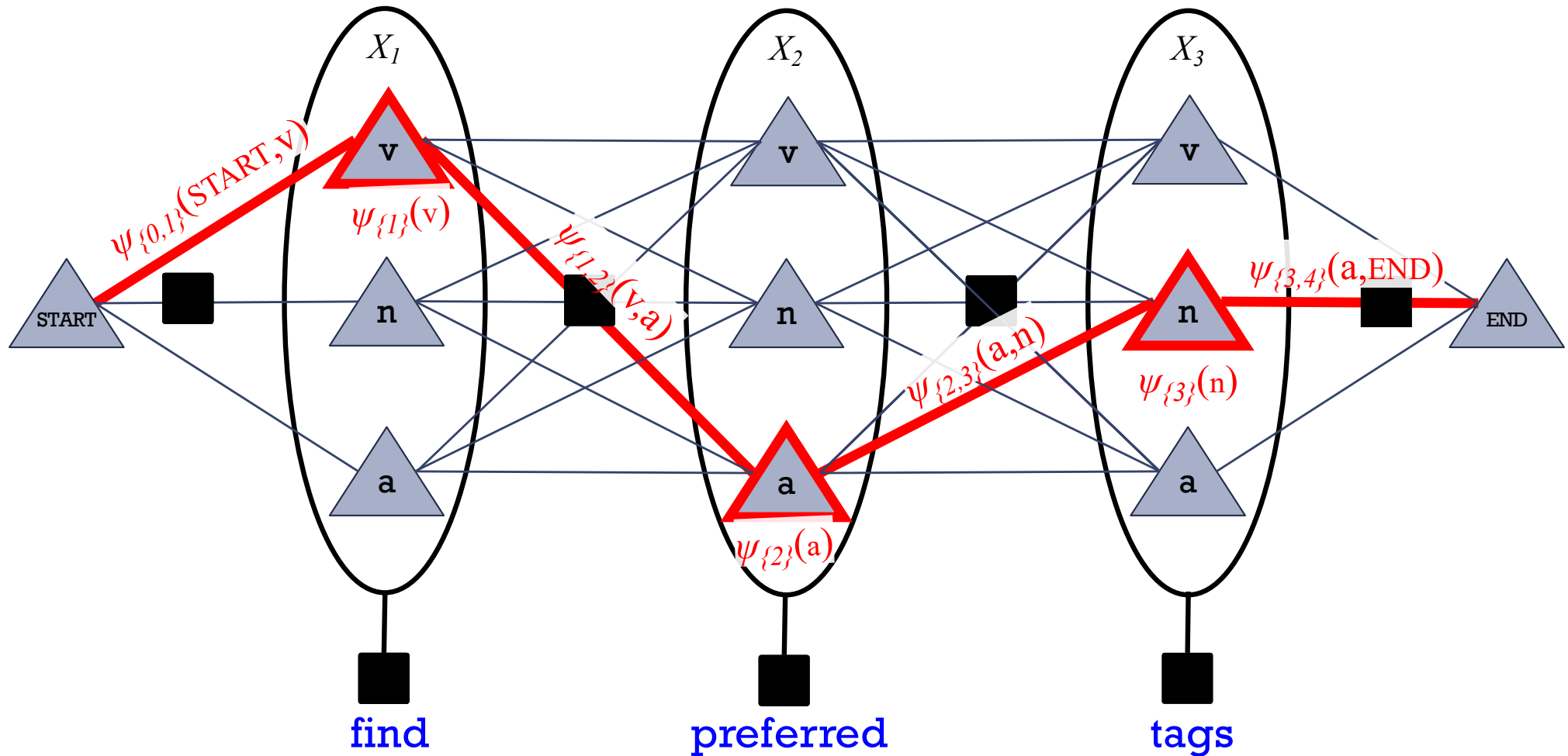
- Let's show the possible *values* for each variable
- One possible assignment
- And what the 7 factors **think of it** ...

# Viterbi Algorithm: Most Probable Assignment



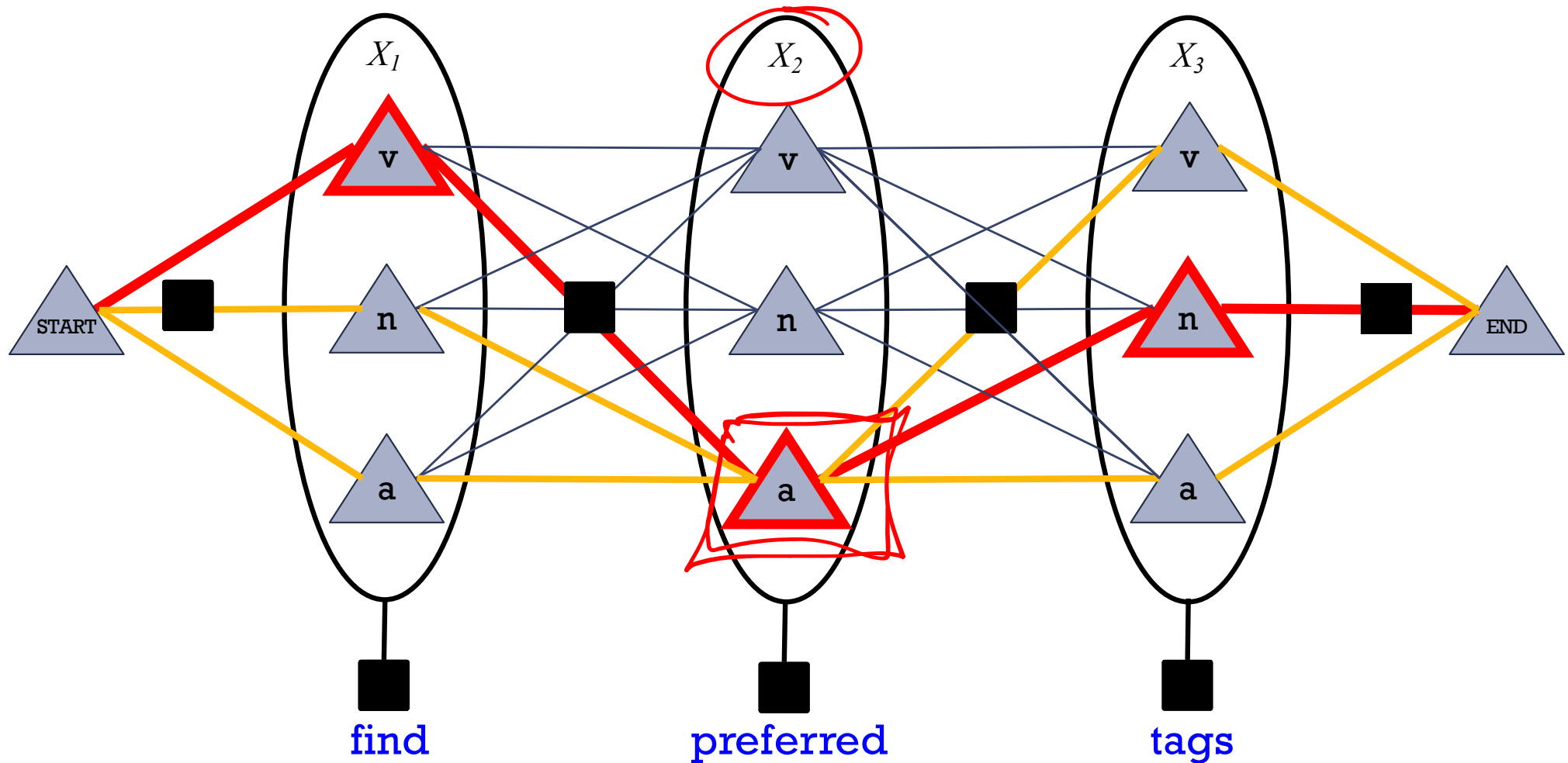
- So  $p(v \ a \ n) = (1/Z) * \text{product of 7 numbers}$
- Numbers associated with edges and nodes of path
- Most probable assignment = **path with highest product**

# Viterbi Algorithm: Most Probable Assignment



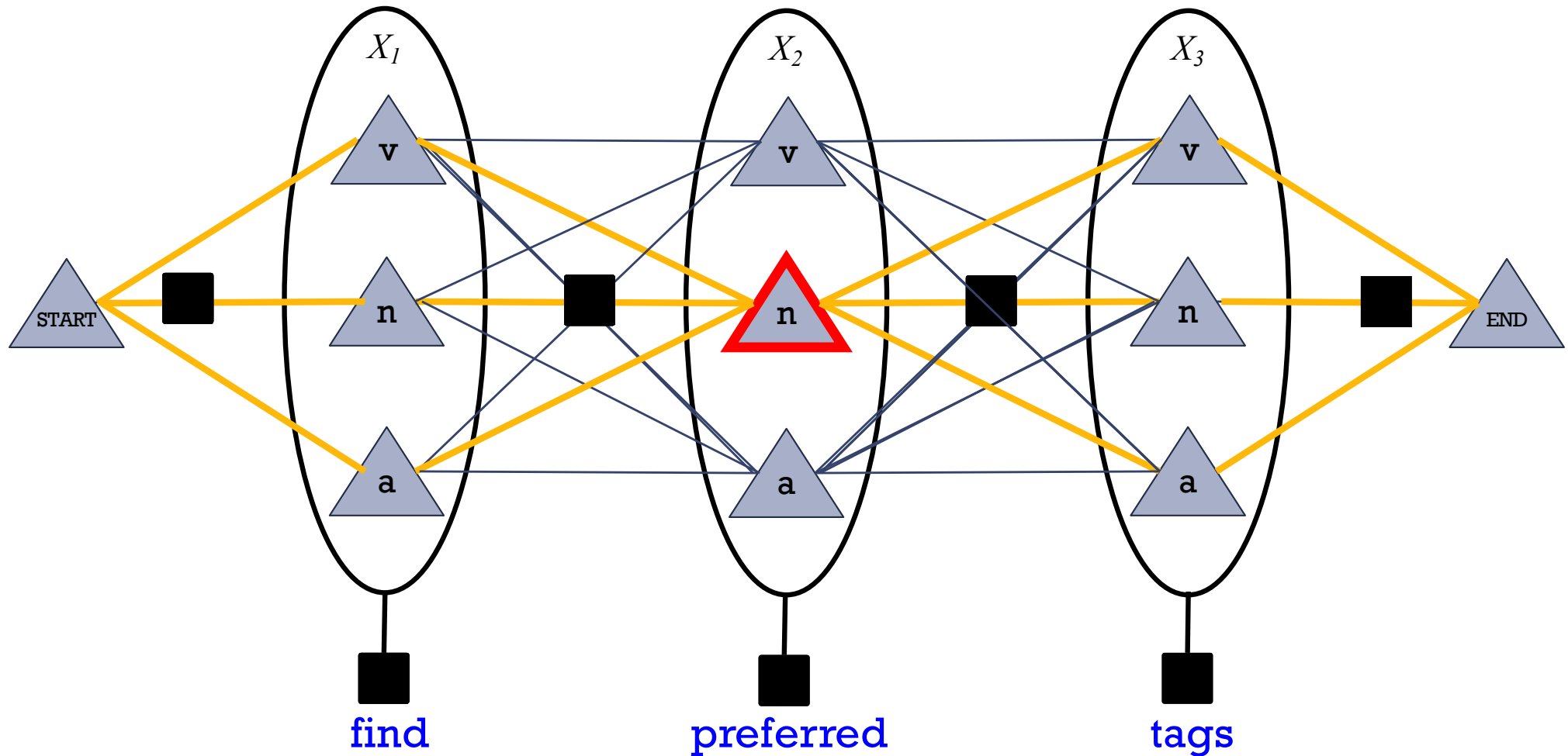
- So  $p(v \ a \ n)$  =  $(1/Z)$  \* product weight of one path

# Forward-Backward Algorithm: Finds Marginals

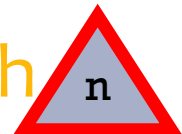


- So  $p(v \ a \ n) = (1/Z) * \text{product weight of one path}$
- Marginal probability  $p(X_2 = a)$   
 $= (1/Z) * \text{total weight of all paths through } a$

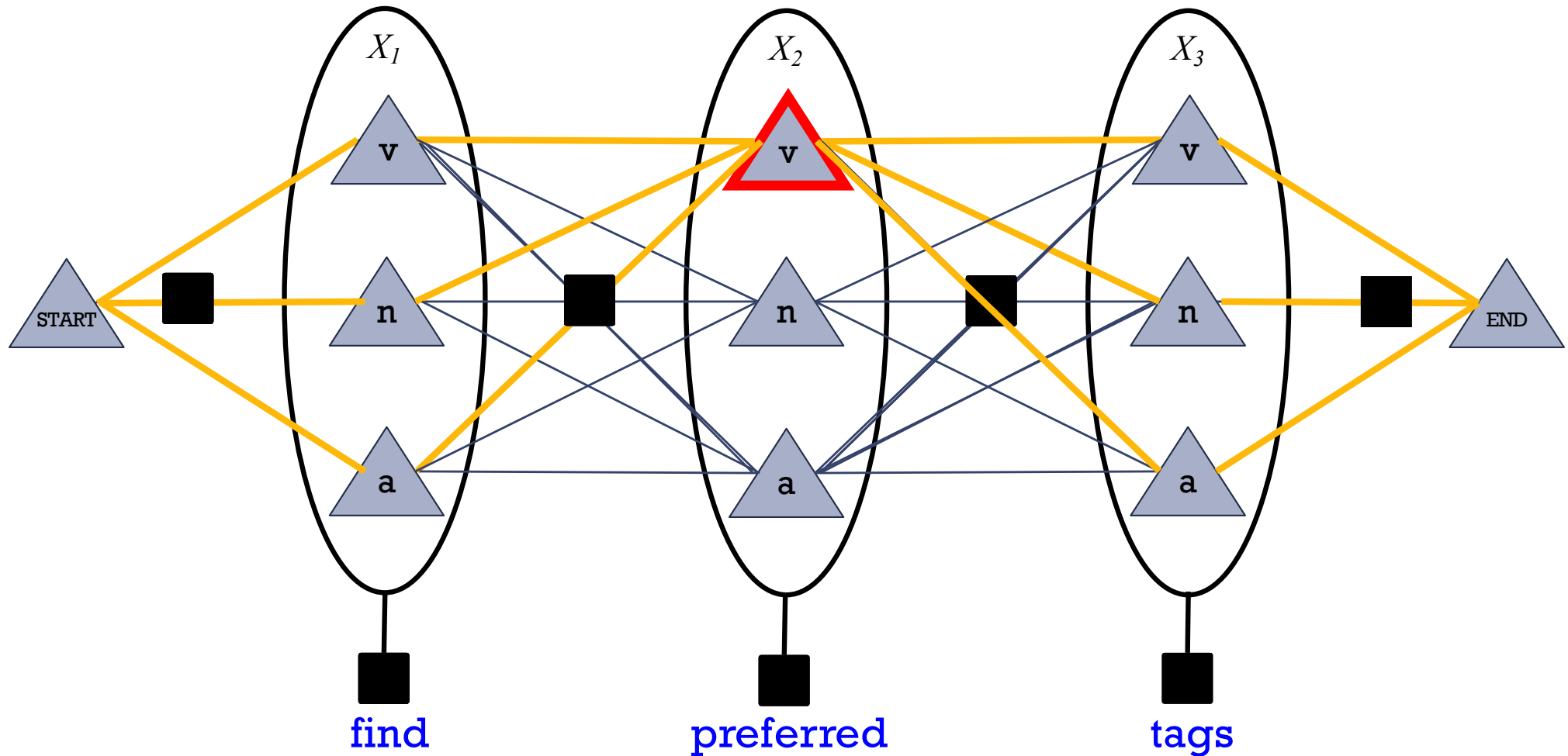
# Forward-Backward Algorithm: Finds Marginals



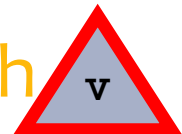
- So  $p(v \ a \ n) = (1/Z) * \text{product weight of one path}$
- Marginal probability  $p(X_2 = a)$   
 $= (1/Z) * \text{total weight of all paths through}$



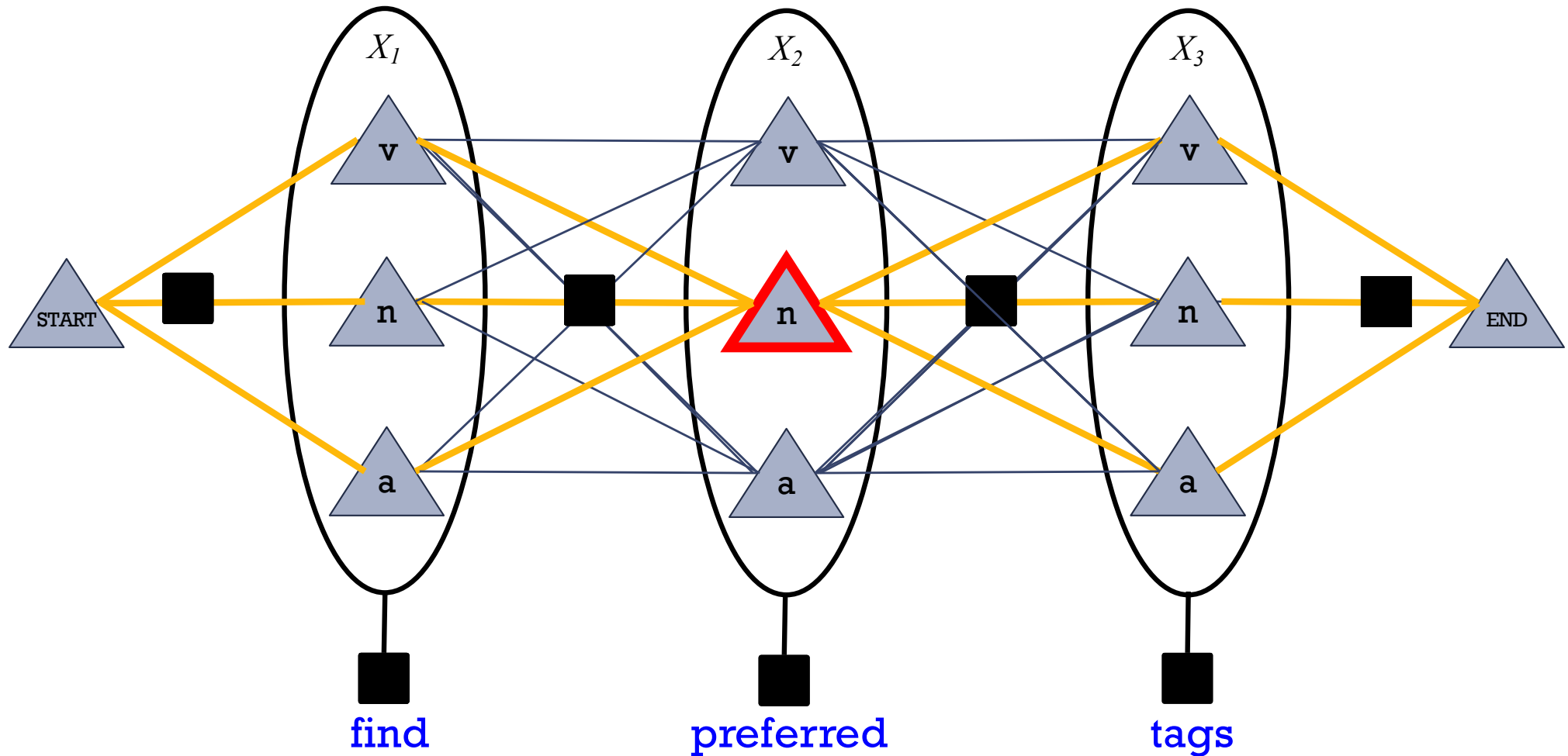
# Forward-Backward Algorithm: Finds Marginals



- So  $p(\mathbf{v} \mathbf{a} \mathbf{n}) = (1/Z) * \text{product weight of one path}$
- Marginal probability  $p(X_2 = a)$   
 $= (1/Z) * \text{total weight of all paths through}$



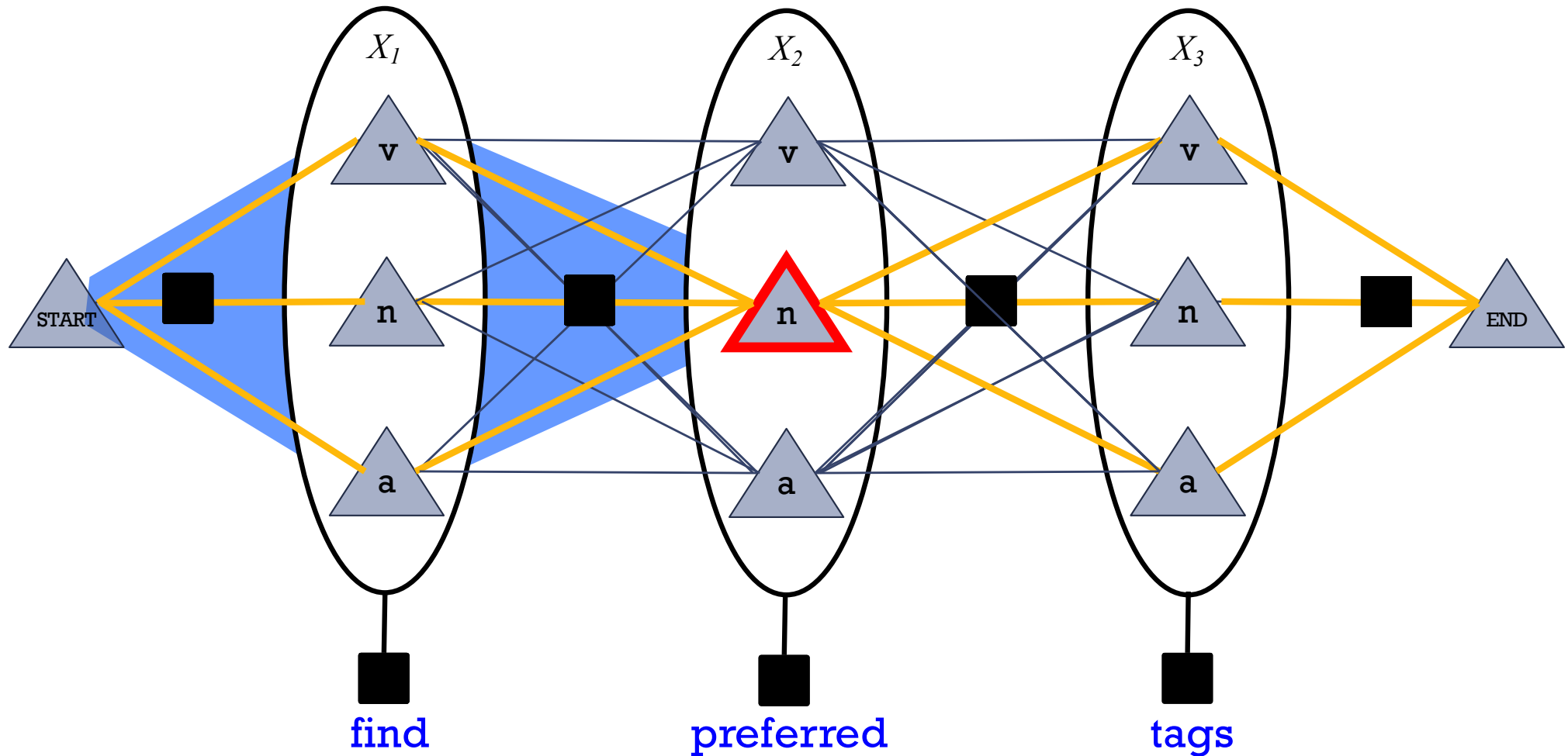
# Forward-Backward Algorithm: Finds Marginals



- So  $p(\mathbf{v} \mathbf{a} \mathbf{n}) = (1/Z) * \text{product weight of one path}$
- Marginal probability  $p(X_2 = a)$   
 $= (1/Z) * \text{total weight of all paths through } \mathbf{n}$



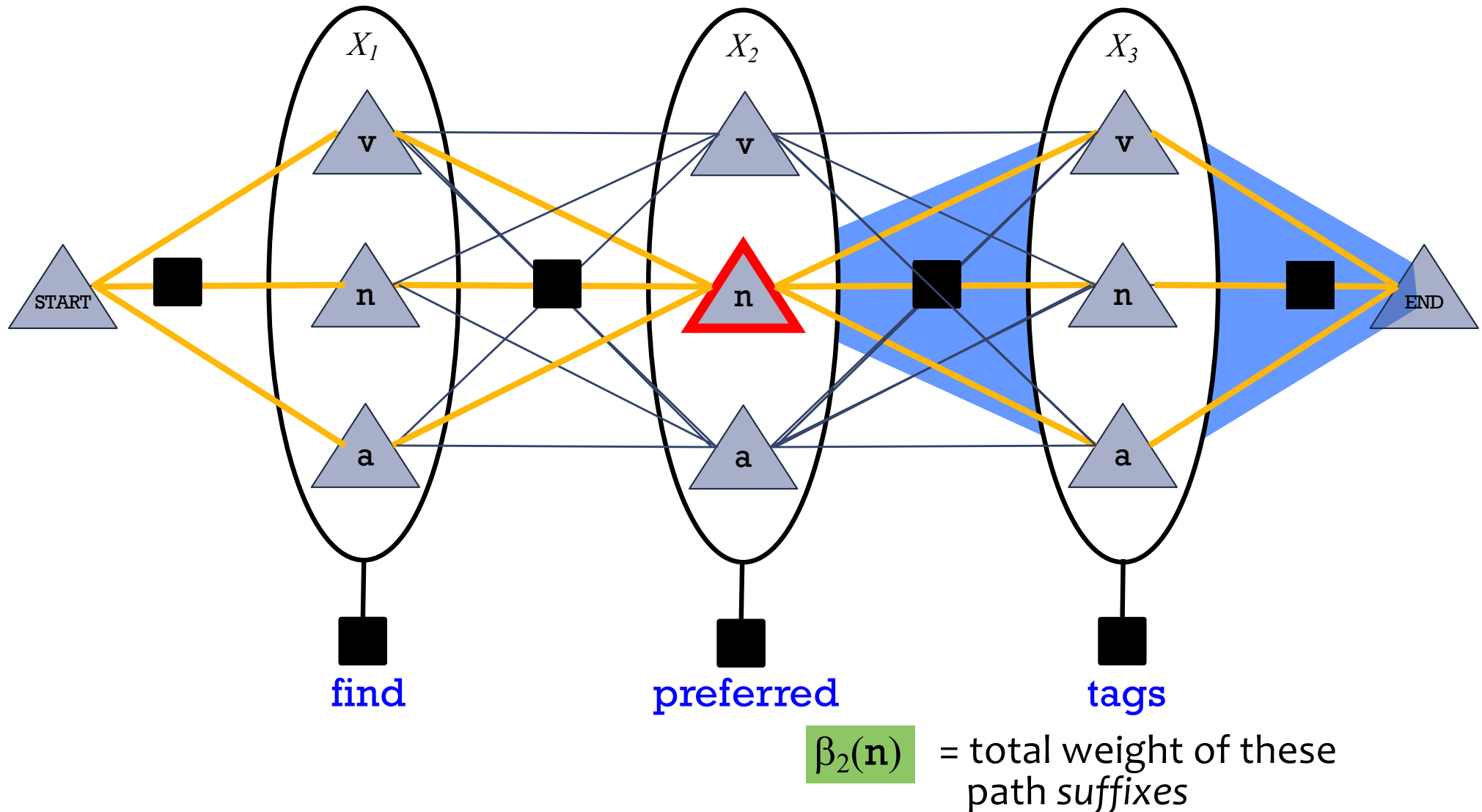
# Forward-Backward Algorithm: Finds Marginals



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

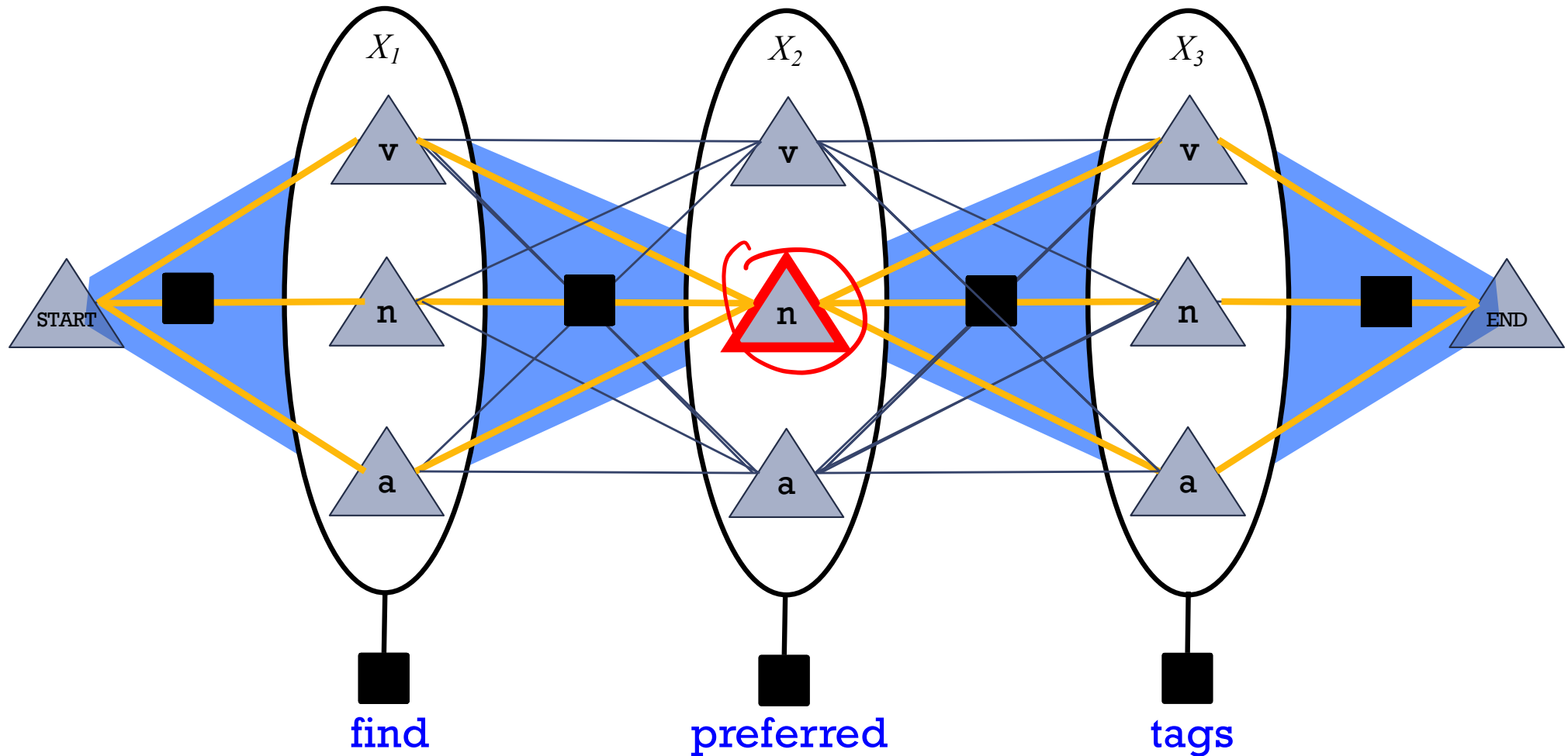
(found by dynamic programming: matrix-vector products)

# Forward-Backward Algorithm: Finds Marginals



(found by dynamic programming: matrix-vector products)

# Forward-Backward Algorithm: Finds Marginals



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

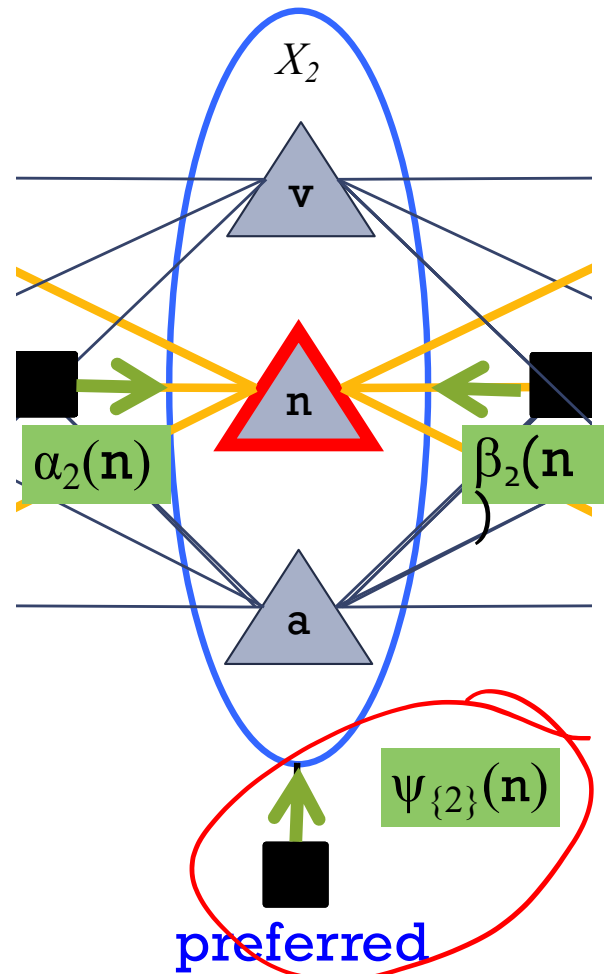
$\beta_2(n)$  = total weight of these path suffixes (x + y + z)

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  $\alpha(\mathbf{n}) \cdot \beta(\mathbf{n})$  isn't enough.

The extra weight is the opinion of the unigram factor at this variable.

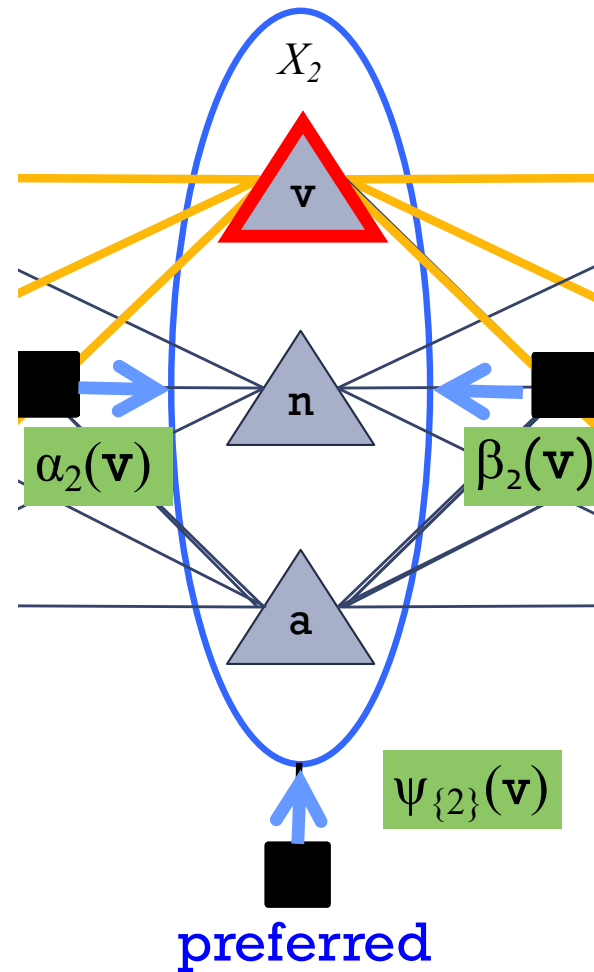


“belief that  $X_2 = \mathbf{n}$ ”

total weight of *all* paths through  $\triangle n$

$$= \underbrace{\alpha_2(\mathbf{n}) \quad \psi_{\{2\}}(\mathbf{n}) \quad \beta_2(\mathbf{n})}_{\text{total weight}}$$

# Forward-Backward Algorithm: Finds Marginals



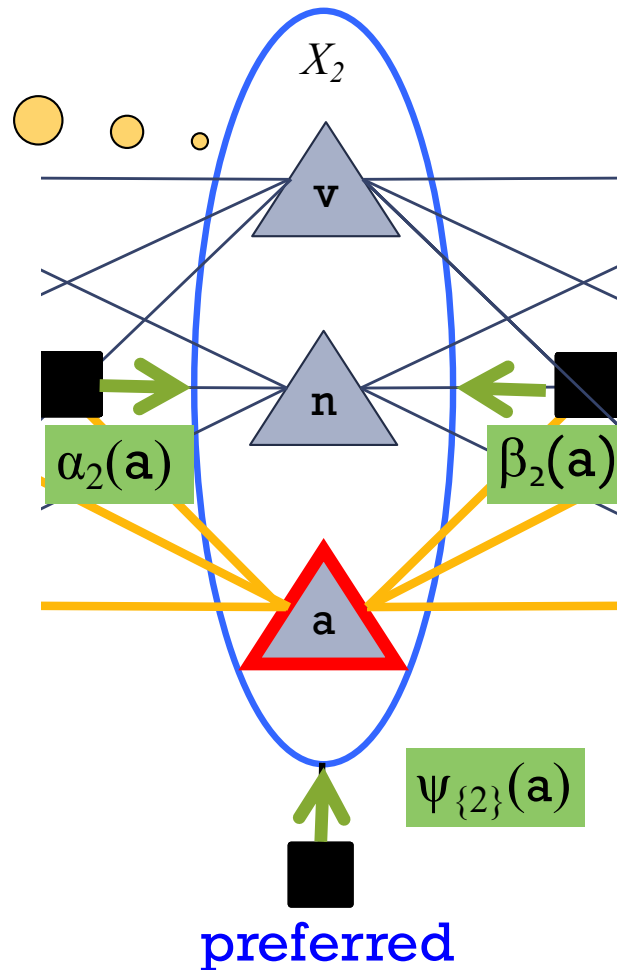
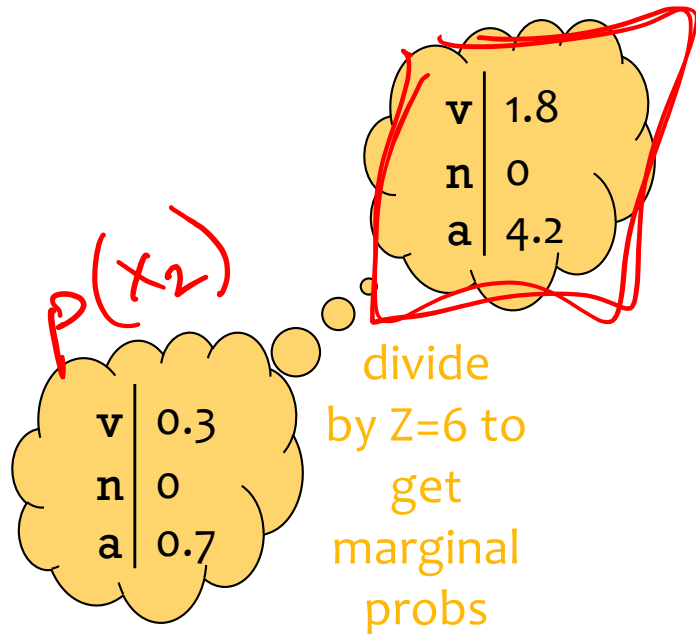
“belief that  $X_2 = \mathbf{v}$ ”

“belief that  $X_2 = \mathbf{n}$ ”

total weight of *all* paths through 

$$= \alpha_2(\mathbf{v}) \psi_{\{2\}}(\mathbf{v}) \beta_2(\mathbf{v})$$

# Forward-Backward Algorithm: Finds Marginals



“belief that  $X_2 = \mathbf{v}$ ”

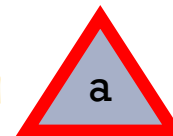
“belief that  $X_2 = \mathbf{n}$ ”

“belief that  $X_2 = \mathbf{a}$ ”

sum =  $Z$   
(total probability of all paths)

total weight of all paths through

$$= \alpha_2(\mathbf{a}) \psi_{\{2\}}(\mathbf{a}) \beta_2(\mathbf{a})$$



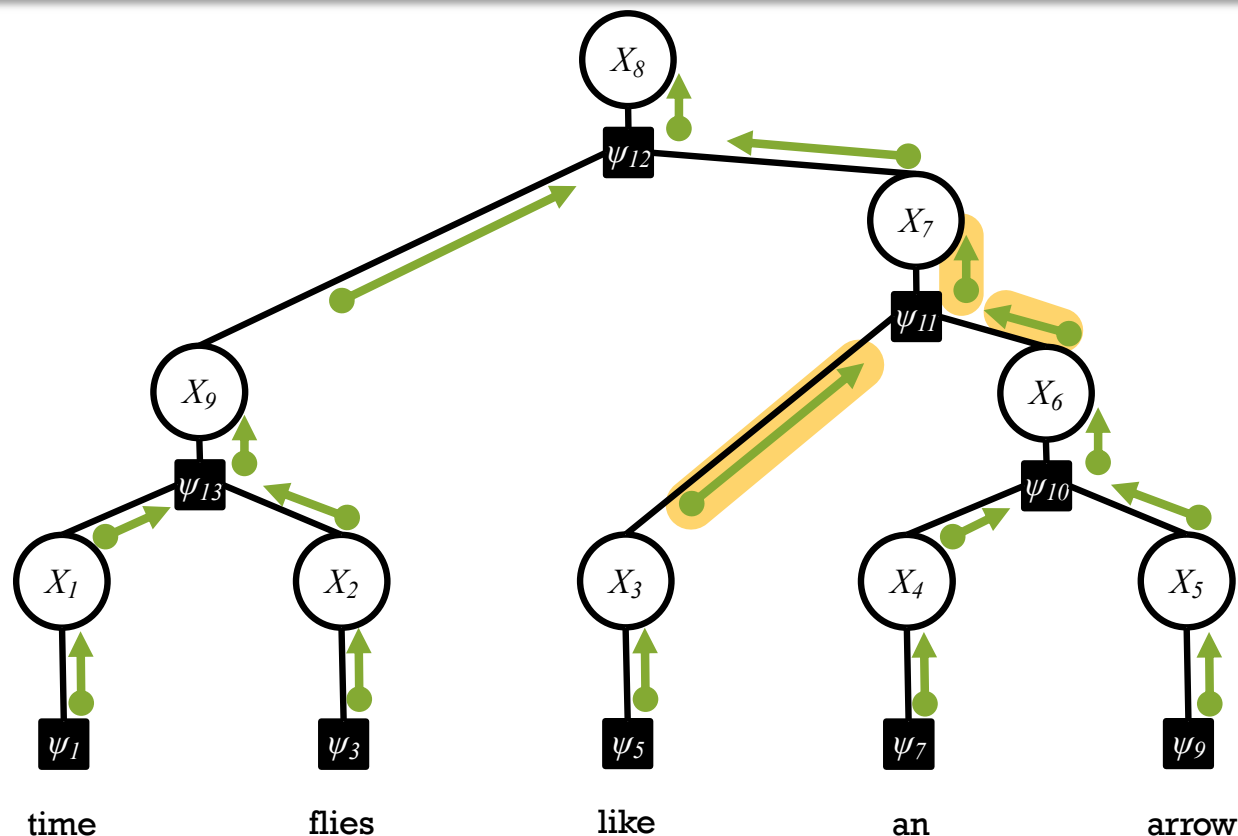
# **BP AS DYNAMIC PROGRAMMING**

# (Acyclic) Belief Propagation

In a factor graph with no cycles:

1. Pick any node to serve as the root.
2. Send messages from the **leaves** to the **root**.
3. Send messages from the **root** to the **leaves**.

A node computes an outgoing message along an edge only after it has received incoming messages along all its other edges.



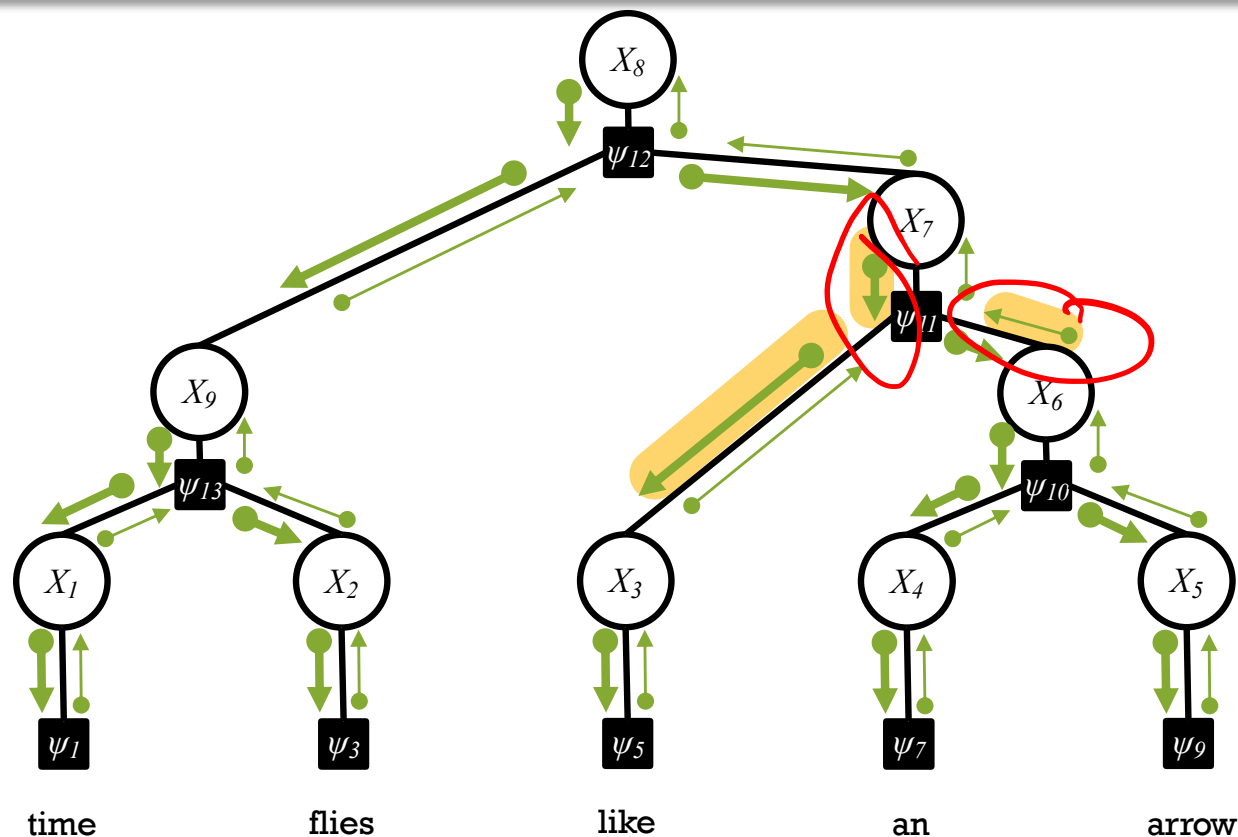


# (Acyclic) Belief Propagation

In a factor graph with no cycles:

1. Pick any node to serve as the root.
2. Send messages from the **leaves** to the **root**.
3. Send messages from the **root** to the **leaves**.

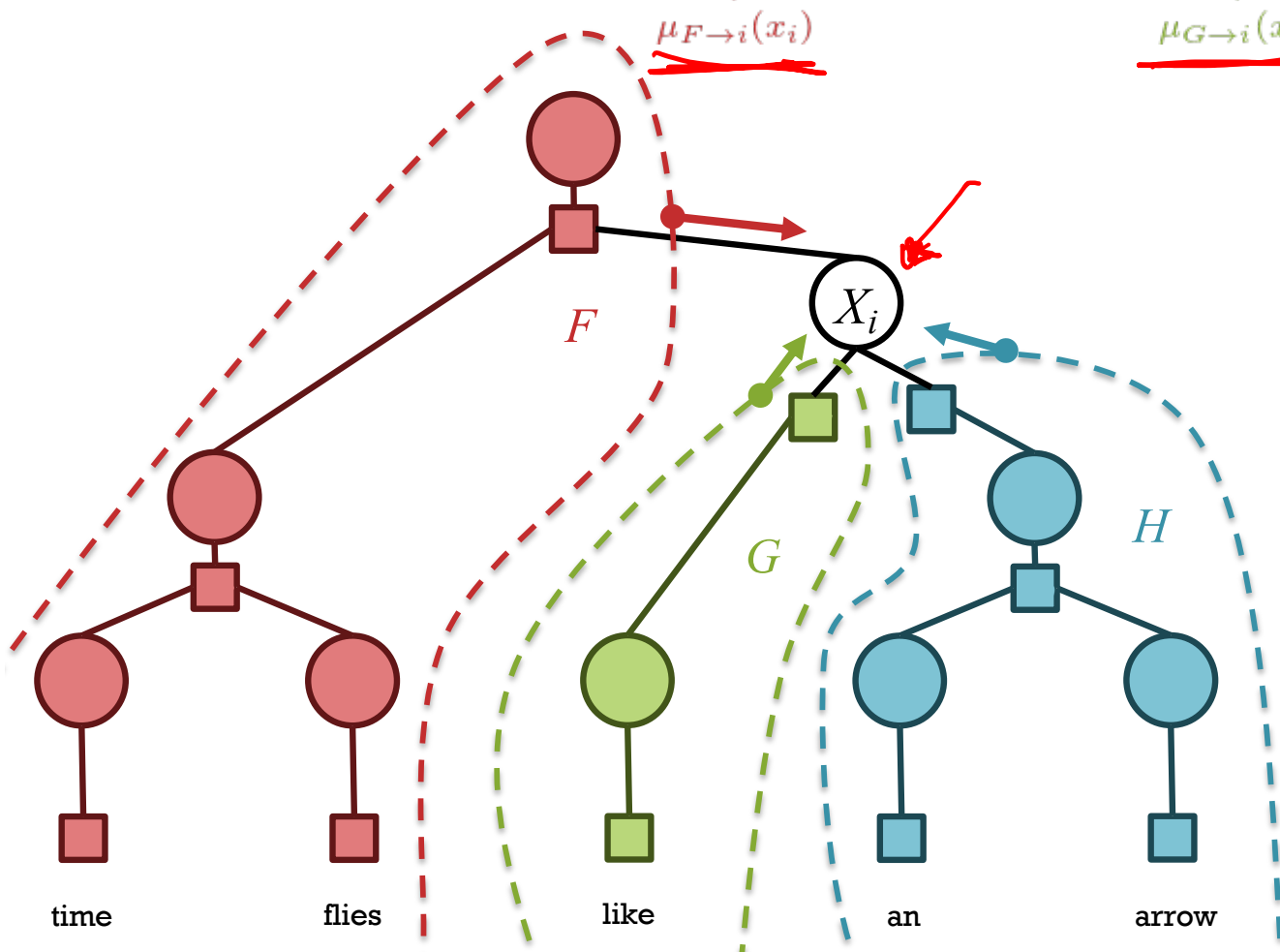
A node computes an outgoing message along an edge only after it has received incoming messages along all its other edges.



# Acyclic BP as Dynamic Programming

$$\underline{p(X_i = x_i)} \propto \underline{b_i(x_i)} = \frac{1}{c} \sum_{\mathbf{x}: \mathbf{x}[i] = x_i} \prod_{\alpha} \psi_{\alpha}(\mathbf{x}_{\alpha})$$

$$= \underbrace{\left( \sum_{\mathbf{x}: \mathbf{x}[i] = x_i} \prod_{\alpha \subseteq F} \psi_{\alpha}(\mathbf{x}_{\alpha}) \right)}_{\mu_{F \rightarrow i}(x_i)} \underbrace{\left( \sum_{\mathbf{x}: \mathbf{x}[i] = x_i} \prod_{\alpha \subseteq G} \psi_{\alpha}(\mathbf{x}_{\alpha}) \right)}_{\mu_{G \rightarrow i}(x_i)} \underbrace{\left( \sum_{\mathbf{x}: \mathbf{x}[i] = x_i} \prod_{\alpha \subseteq H} \psi_{\alpha}(\mathbf{x}_{\alpha}) \right)}_{\mu_{H \rightarrow i}(x_i)}$$



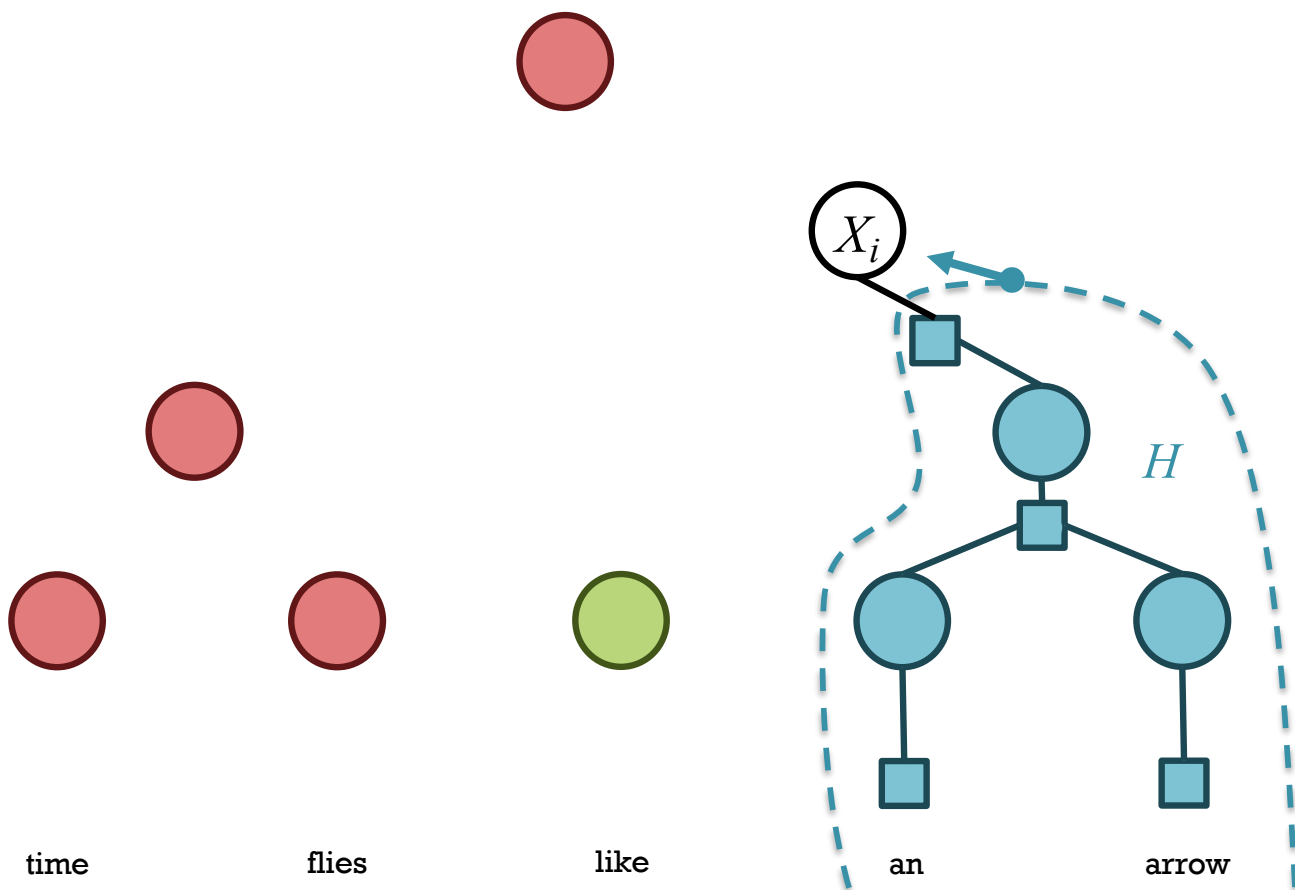
## Subproblem:

Inference using just the factors in subgraph  $H$

Figure adapted from Burkett & Klein (2012)<sup>40</sup>

# Acyclic BP as Dynamic Programming

$$\begin{aligned}
 p(X_i = x_i) &\propto b_i(x_i) = \sum_{\mathbf{x}: \mathbf{x}[i]=x_i} \prod_{\alpha} \psi_{\alpha}(\mathbf{x}_{\alpha}) \\
 &= \underbrace{\left( \sum_{\mathbf{x}: \mathbf{x}[i]=x_i} \prod_{\alpha \subseteq F} \psi_{\alpha}(\mathbf{x}_{\alpha}) \right)}_{\mu_{F \rightarrow i}(x_i)} \underbrace{\left( \sum_{\mathbf{x}: \mathbf{x}[i]=x_i} \prod_{\alpha \subseteq G} \psi_{\alpha}(\mathbf{x}_{\alpha}) \right)}_{\mu_{G \rightarrow i}(x_i)} \underbrace{\left( \sum_{\mathbf{x}: \mathbf{x}[i]=x_i} \prod_{\alpha \subseteq H} \psi_{\alpha}(\mathbf{x}_{\alpha}) \right)}_{\mu_{H \rightarrow i}(x_i)}
 \end{aligned}$$



## Subproblem:

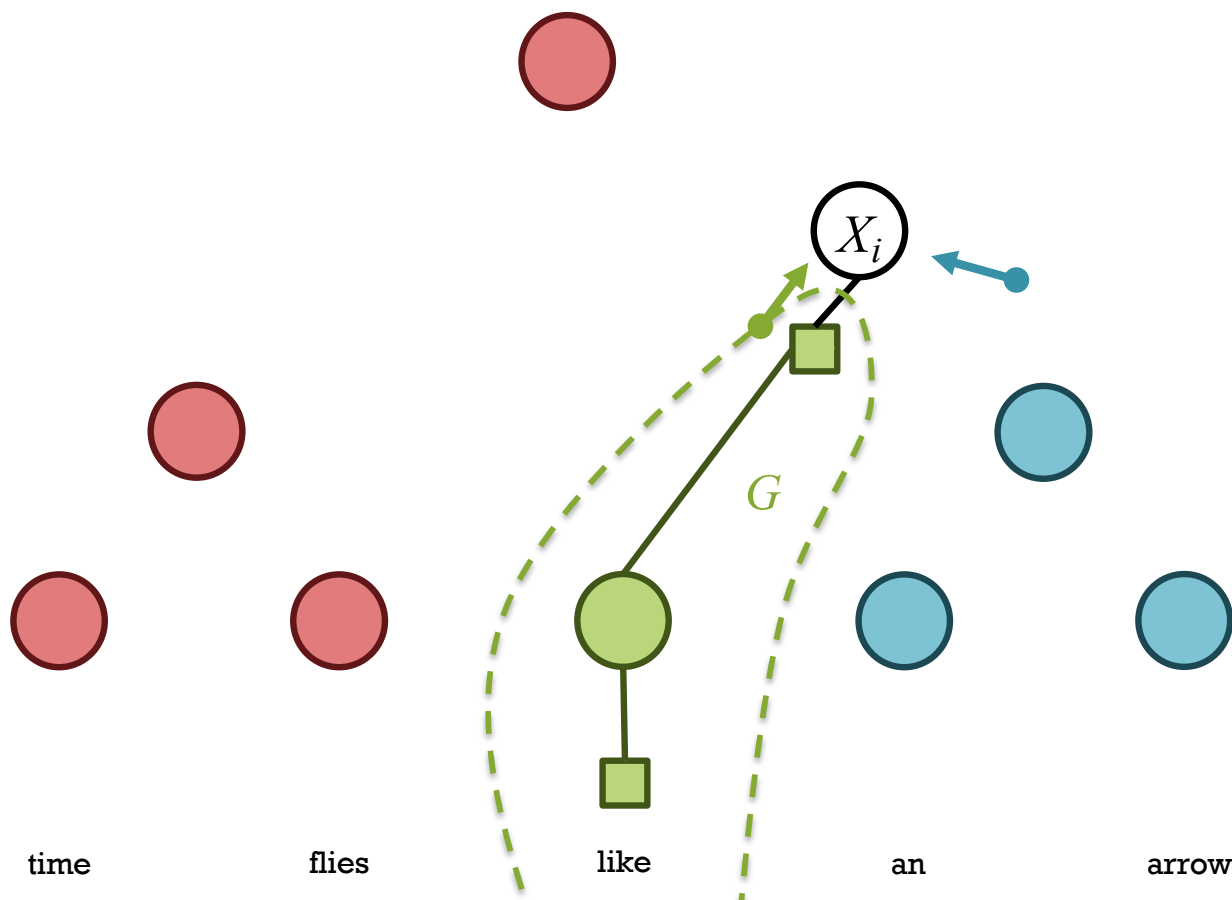
Inference using just the factors in subgraph  $H$

The marginal of  $X_i$  in that smaller model is the message sent to  $X_i$  from subgraph  $H$

*Message to a variable*

# Acyclic BP as Dynamic Programming

$$\begin{aligned}
 p(X_i = x_i) \propto b_i(x_i) &= \sum_{\mathbf{x}: \mathbf{x}[i] = x_i} \prod_{\alpha} \psi_{\alpha}(\mathbf{x}_{\alpha}) \\
 &= \underbrace{\left( \sum_{\mathbf{x}: \mathbf{x}[i] = x_i} \prod_{\alpha \subseteq F} \psi_{\alpha}(\mathbf{x}_{\alpha}) \right)}_{\mu_{F \rightarrow i}(x_i)} \underbrace{\left( \sum_{\mathbf{x}: \mathbf{x}[i] = x_i} \prod_{\alpha \subseteq G} \psi_{\alpha}(\mathbf{x}_{\alpha}) \right)}_{\mu_{G \rightarrow i}(x_i)} \underbrace{\left( \sum_{\mathbf{x}: \mathbf{x}[i] = x_i} \prod_{\alpha \subseteq H} \psi_{\alpha}(\mathbf{x}_{\alpha}) \right)}_{\mu_{H \rightarrow i}(x_i)}
 \end{aligned}$$



## Subproblem:

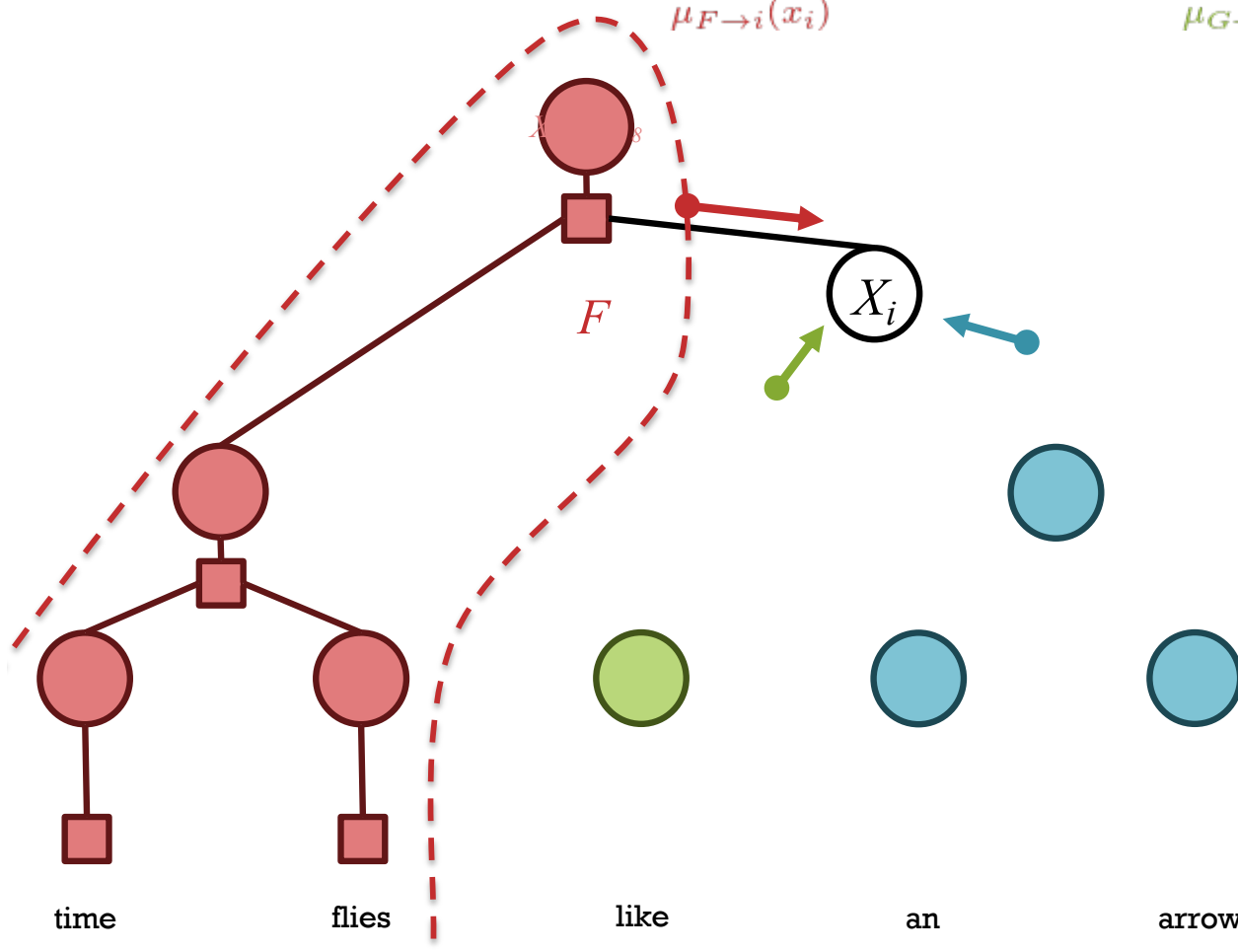
Inference using just the factors in subgraph  $H$

The marginal of  $X_i$  in that smaller model is the message sent to  $X_i$  from subgraph  $H$

*Message to a variable*

# Acyclic BP as Dynamic Programming

$$\begin{aligned}
 p(X_i = x_i) \propto b_i(x_i) &= \sum_{\mathbf{x}: \mathbf{x}[i] = x_i} \prod_{\alpha} \psi_{\alpha}(\mathbf{x}_{\alpha}) \\
 &= \underbrace{\left( \sum_{\mathbf{x}: \mathbf{x}[i] = x_i} \prod_{\alpha \subseteq F} \psi_{\alpha}(\mathbf{x}_{\alpha}) \right)}_{\mu_{F \rightarrow i}(x_i)} \underbrace{\left( \sum_{\mathbf{x}: \mathbf{x}[i] = x_i} \prod_{\alpha \subseteq G} \psi_{\alpha}(\mathbf{x}_{\alpha}) \right)}_{\mu_{G \rightarrow i}(x_i)} \underbrace{\left( \sum_{\mathbf{x}: \mathbf{x}[i] = x_i} \prod_{\alpha \subseteq H} \psi_{\alpha}(\mathbf{x}_{\alpha}) \right)}_{\mu_{H \rightarrow i}(x_i)}
 \end{aligned}$$



## Subproblem:

Inference using just the factors in subgraph  $H$

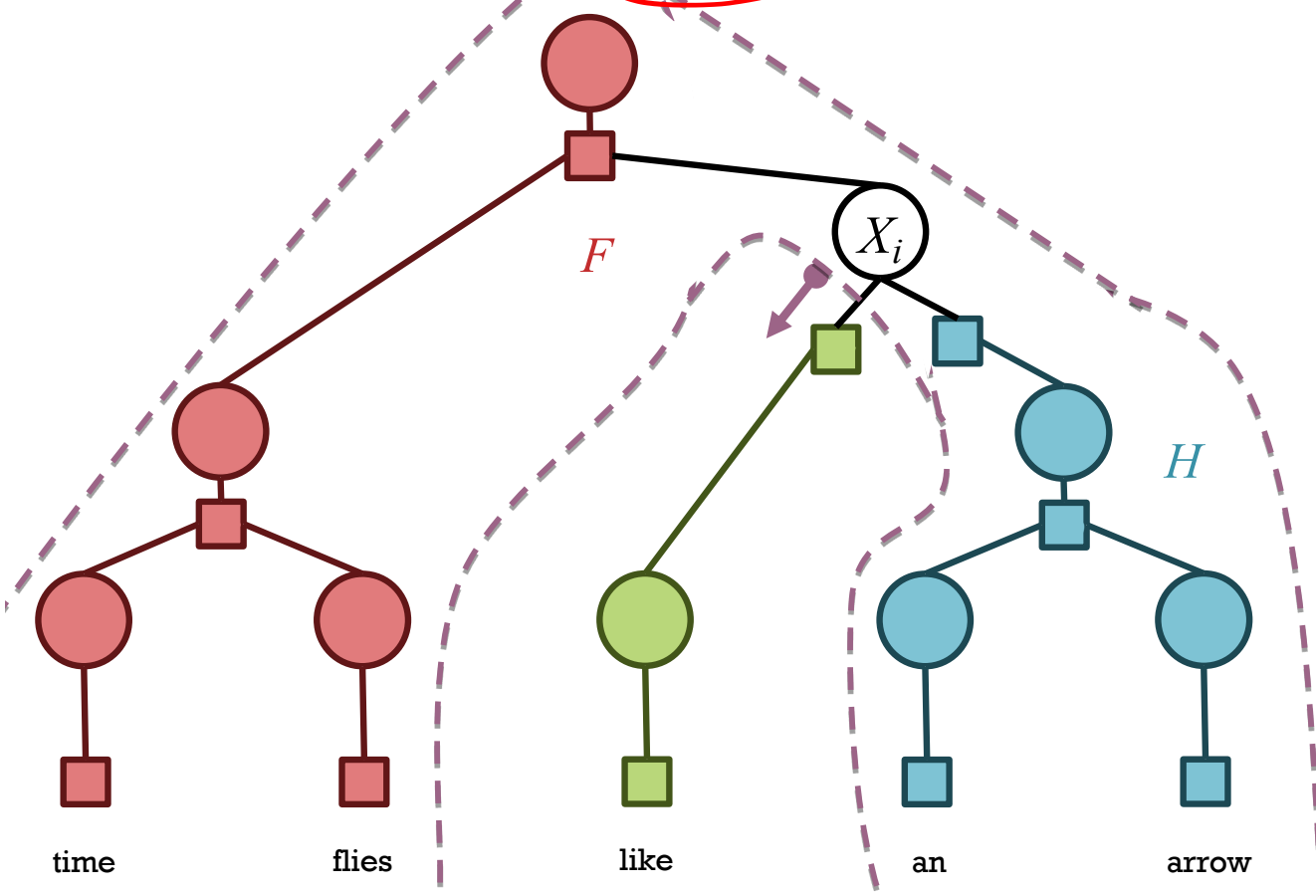
The marginal of  $X_i$  in that smaller model is the message sent to  $X_i$  from subgraph  $H$

*Message to a variable*

# Acyclic BP as Dynamic Programming

$$p(X_i = x_i) \propto b_i(x_i) = \sum_{\mathbf{x}: \mathbf{x}[i]=x_i} \prod_{\alpha} \psi_{\alpha}(\mathbf{x}_{\alpha})$$

$$= \underbrace{\left( \sum_{\mathbf{x}: \mathbf{x}[i]=x_i} \prod_{\alpha \subseteq F} \psi_{\alpha}(\mathbf{x}_{\alpha}) \right)}_{\mu_{F \rightarrow i}(x_i)} \underbrace{\left( \sum_{\mathbf{x}: \mathbf{x}[i]=x_i} \prod_{\alpha \subseteq G} \psi_{\alpha}(\mathbf{x}_{\alpha}) \right)}_{\mu_{G \rightarrow i}(x_i)} \underbrace{\left( \sum_{\mathbf{x}: \mathbf{x}[i]=x_i} \prod_{\alpha \subseteq H} \psi_{\alpha}(\mathbf{x}_{\alpha}) \right)}_{\mu_{H \rightarrow i}(x_i)}$$



## Subproblem:

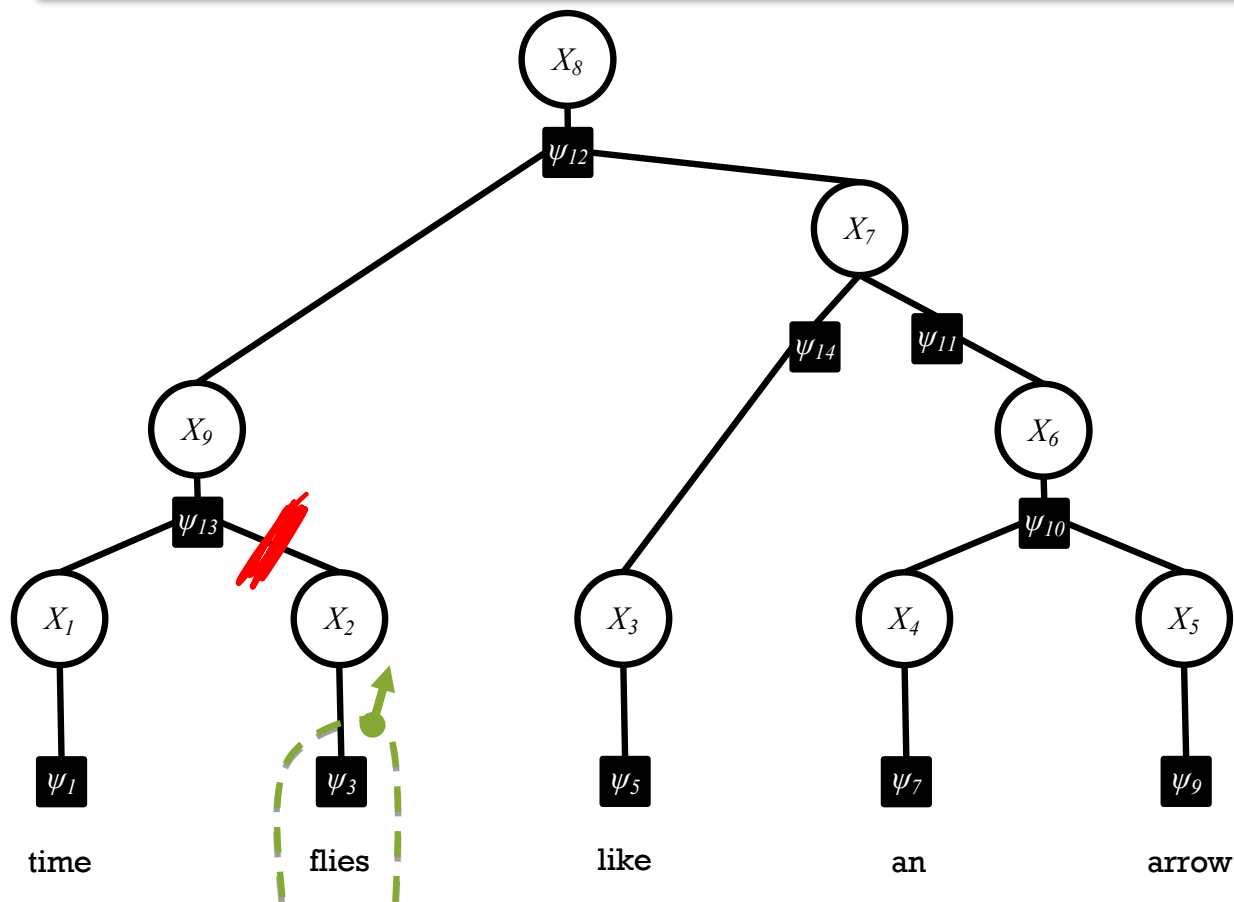
Inference using just the factors in subgraph  $F \cup H$

The marginal of  $X_i$  in that smaller model is the message sent by  $X_i$  out of subgraph  $F \cup H$

*Message from a variable*

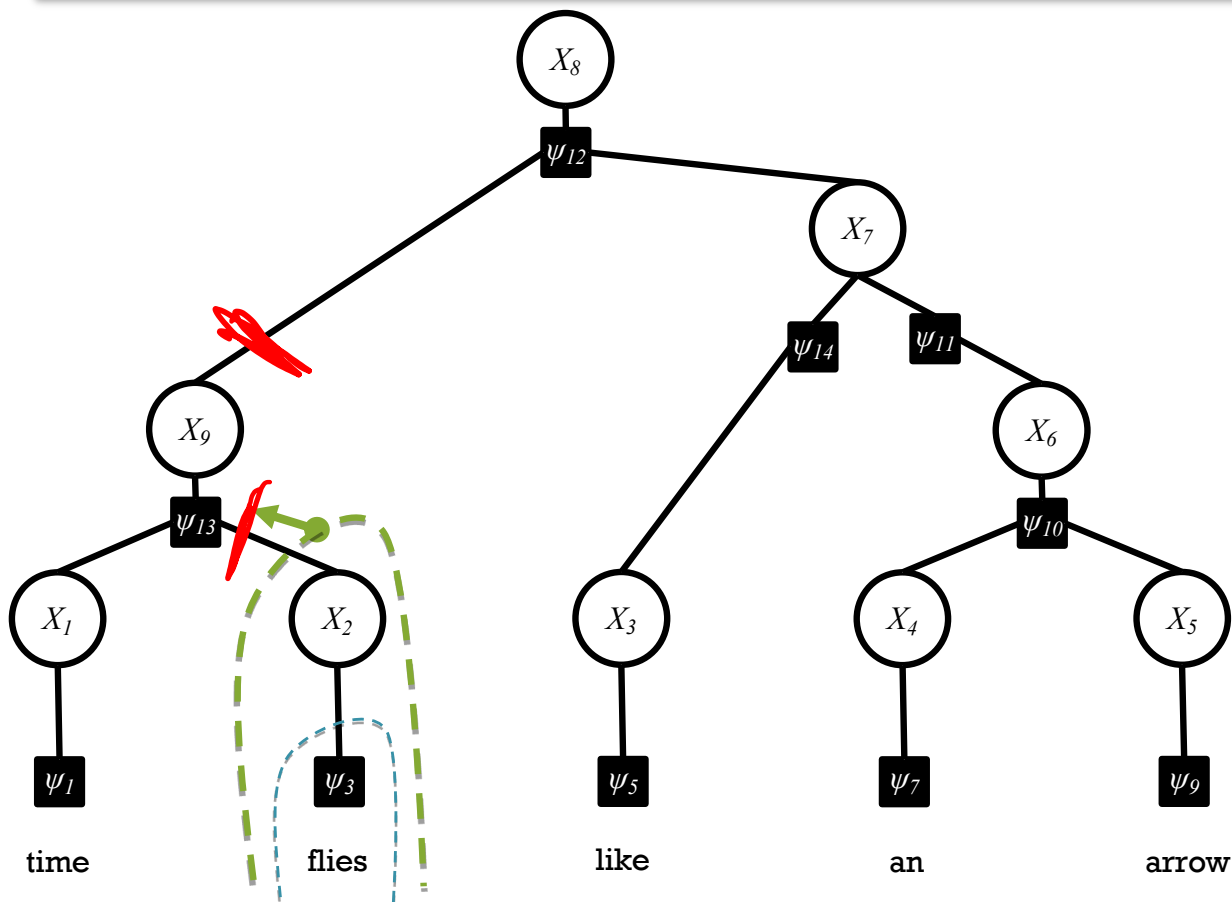
# Acyclic BP as Dynamic Programming

- If you want the **marginal**  $p_i(x_i)$  where  $X_i$  has degree  $k$ , you can think of that summation as a **product of  $k$  marginals** computed on smaller subgraphs.
- Each subgraph is obtained by **cutting** some edge of the tree.
- The message-passing algorithm uses **dynamic programming** to compute the marginals on all such subgraphs, working from **smaller to bigger**. So you can compute all the marginals.



# Acyclic BP as Dynamic Programming

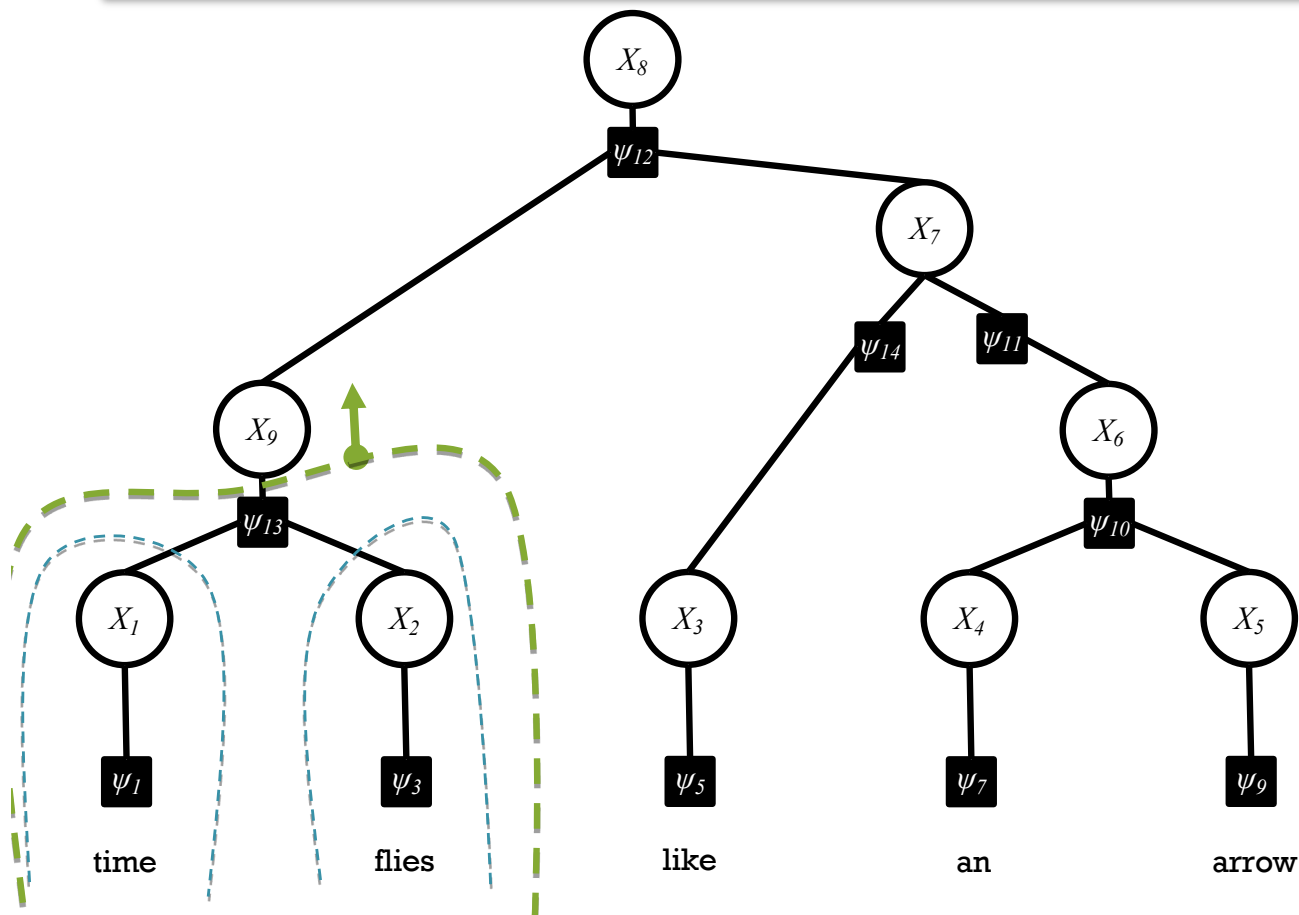
- If you want the **marginal**  $p_i(x_i)$  where  $X_i$  has degree  $k$ , you can think of that summation as a **product of  $k$  marginals** computed on smaller subgraphs.
- Each subgraph is obtained by **cutting** some edge of the tree.
- The message-passing algorithm uses **dynamic programming** to compute the marginals on all such subgraphs, working from **smaller to bigger**. So you can compute all the marginals.





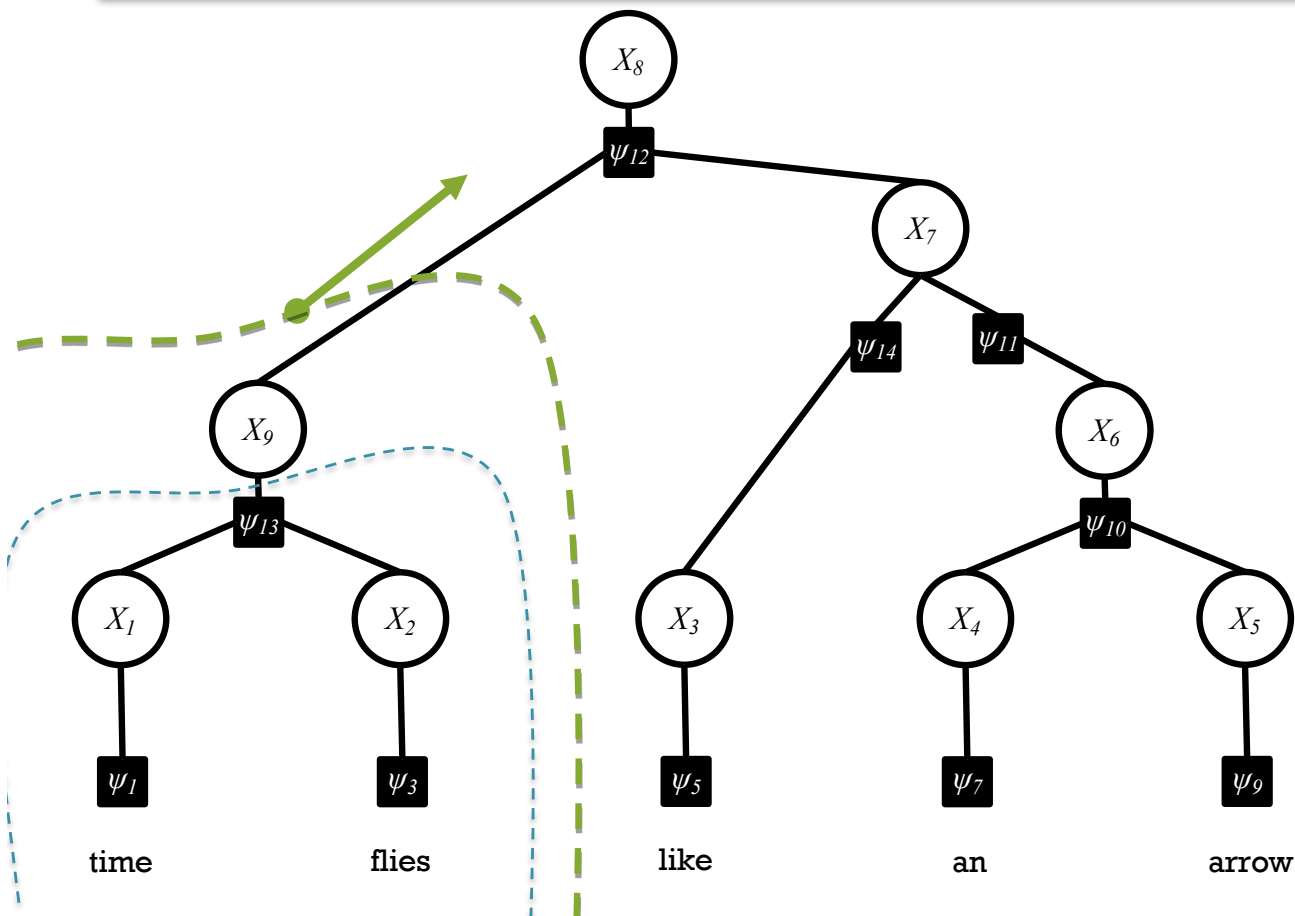
# Acyclic BP as Dynamic Programming

- If you want the **marginal**  $p_i(x_i)$  where  $X_i$  has degree  $k$ , you can think of that summation as a **product of  $k$  marginals** computed on smaller subgraphs.
- Each subgraph is obtained by **cutting** some edge of the tree.
- The message-passing algorithm uses **dynamic programming** to compute the marginals on all such subgraphs, working from **smaller to bigger**. So you can compute all the marginals.



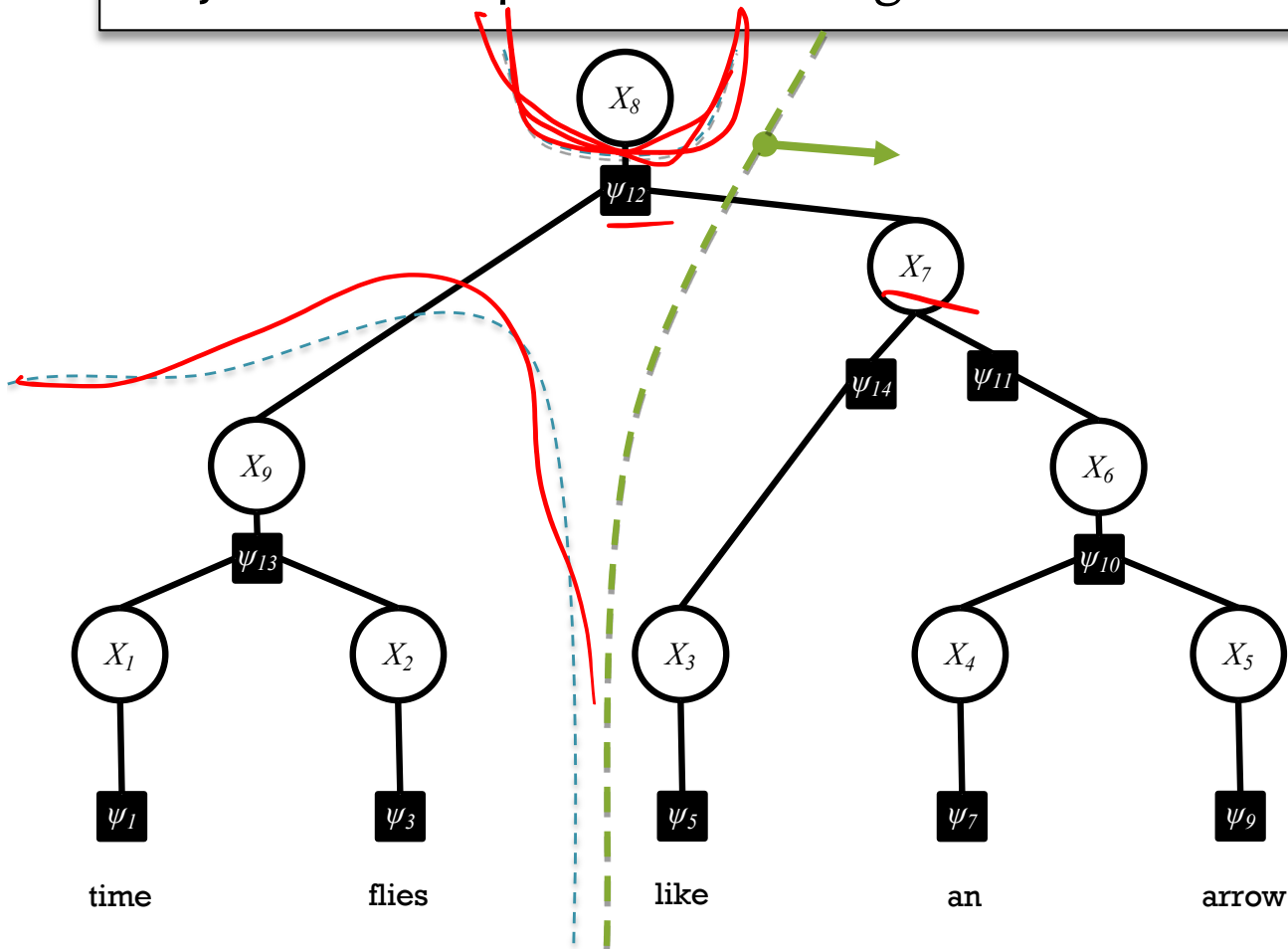
# Acyclic BP as Dynamic Programming

- If you want the **marginal**  $p_i(x_i)$  where  $X_i$  has degree  $k$ , you can think of that summation as a **product of  $k$  marginals** computed on smaller subgraphs.
- Each subgraph is obtained by **cutting** some edge of the tree.
- The message-passing algorithm uses **dynamic programming** to compute the marginals on all such subgraphs, working from **smaller to bigger**. So you can compute all the marginals.



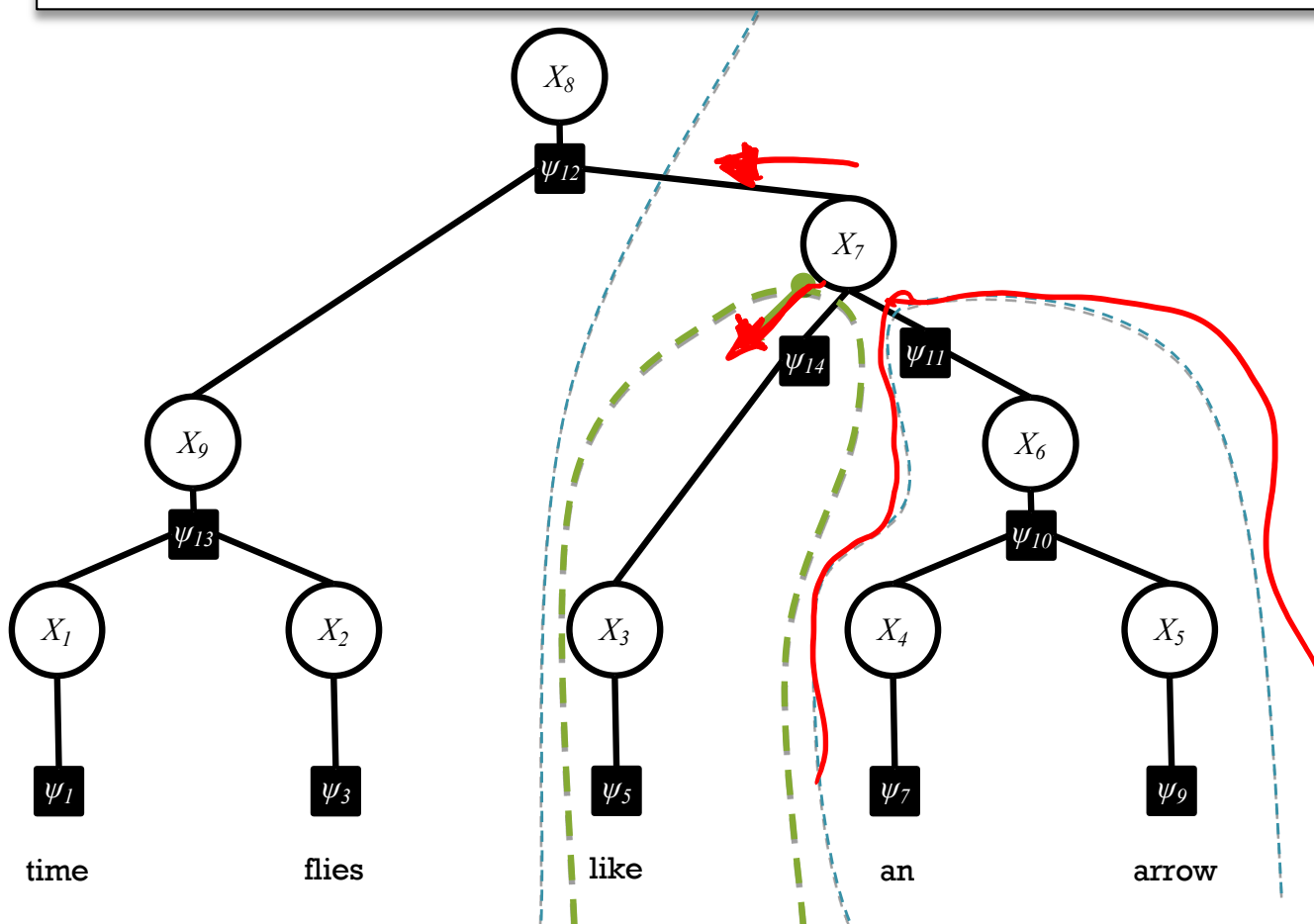
# Acyclic BP as Dynamic Programming

- If you want the **marginal**  $p_i(x_i)$  where  $X_i$  has degree  $k$ , you can think of that summation as a **product of  $k$  marginals** computed on smaller subgraphs.
- Each subgraph is obtained by **cutting** some edge of the tree.
- The message-passing algorithm uses **dynamic programming** to compute the marginals on all such subgraphs, working from **smaller to bigger**. So you can compute all the marginals.



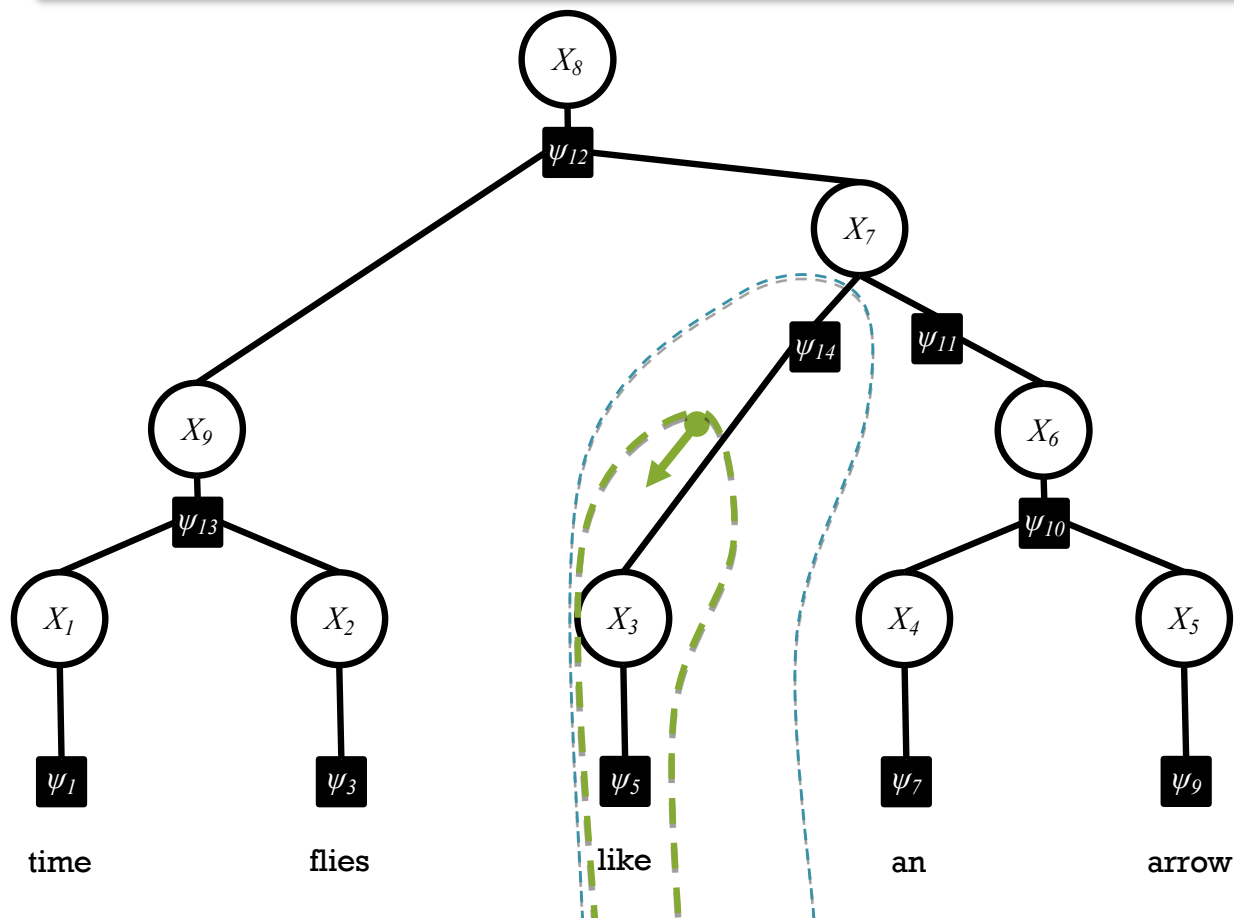
# Acyclic BP as Dynamic Programming

- If you want the **marginal**  $p_i(x_i)$  where  $X_i$  has degree  $k$ , you can think of that summation as a **product of  $k$  marginals** computed on smaller subgraphs.
- Each subgraph is obtained by **cutting** some edge of the tree.
- The message-passing algorithm uses **dynamic programming** to compute the marginals on all such subgraphs, working from **smaller to bigger**. So you can compute all the marginals.



# Acyclic BP as Dynamic Programming

- If you want the **marginal**  $p_i(x_i)$  where  $X_i$  has degree  $k$ , you can think of that summation as a **product of  $k$  marginals** computed on smaller subgraphs.
- Each subgraph is obtained by **cutting** some edge of the tree.
- The message-passing algorithm uses **dynamic programming** to compute the marginals on all such subgraphs, working from **smaller to bigger**. So you can compute all the marginals.



Exact MAP inference for factor trees

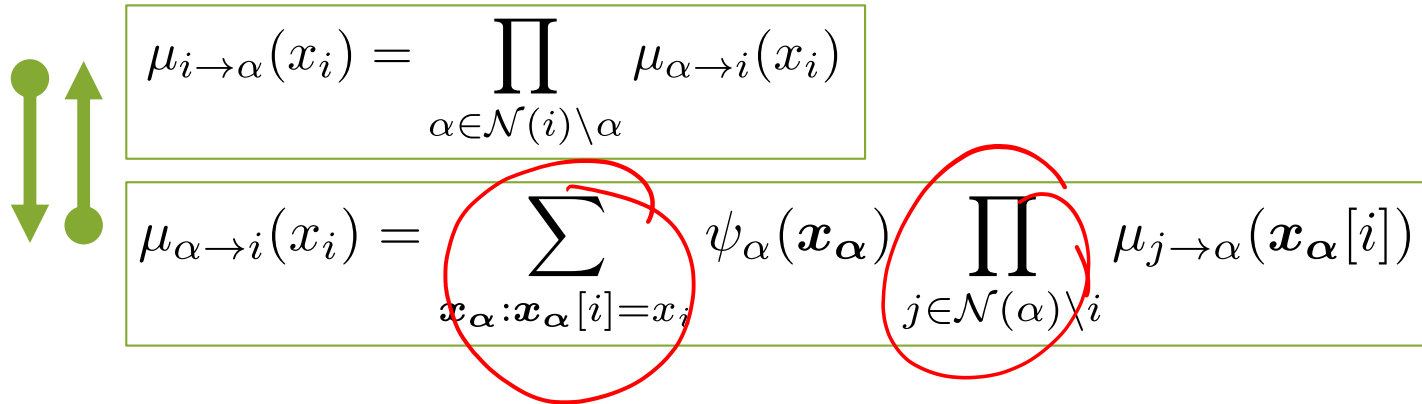
# **MAX-PRODUCT BELIEF PROPAGATION**

# Max-product Belief Propagation

- **Sum-product BP** can be used to  
compute the marginals,  $p_i(X_i)$   
compute the partition function,  $Z$
- **Max-product BP** can be used to  
compute the most likely assignment,  
 $X^* = \operatorname{argmax}_X p(X)$

# Max-product Belief Propagation

- Change the sum to a max:


$$\mu_{i \rightarrow \alpha}(x_i) = \prod_{\alpha \in \mathcal{N}(i) \setminus \alpha} \mu_{\alpha \rightarrow i}(x_i)$$
$$\mu_{\alpha \rightarrow i}(x_i) = \sum_{\mathbf{x}_\alpha: \mathbf{x}_\alpha[i] = x_i} \psi_\alpha(\mathbf{x}_\alpha) \prod_{j \in \mathcal{N}(\alpha) \setminus i} \mu_{j \rightarrow \alpha}(\mathbf{x}_\alpha[j])$$


- **Max-product BP** computes **max-marginals**
  - The max-marginal  $b_i(x_i)$  is the (unnormalized) probability of the MAP assignment under the constraint  $X_i = x_i$ .
  - For an acyclic graph, the MAP assignment (assuming there are no ties) is given by:

$$x_i^* = \arg \max_{x_i} b_i(x_i)$$



# Max-product Belief Propagation

- Change the sum to a max:



$$\mu_{i \rightarrow \alpha}(x_i) = \prod_{\alpha \in \mathcal{N}(i) \setminus \alpha} \mu_{\alpha \rightarrow i}(x_i)$$

$$\mu_{\alpha \rightarrow i}(x_i) = \max_{\mathbf{x}_\alpha : \mathbf{x}_\alpha[i] = x_i} \psi_\alpha(\mathbf{x}_\alpha) \prod_{j \in \mathcal{N}(\alpha) \setminus i} \mu_{j \rightarrow \alpha}(\mathbf{x}_\alpha[j])$$

- **Max-product BP** computes **max-marginals**
  - The max-marginal  $b_i(x_i)$  is the (unnormalized) probability of the MAP assignment under the constraint  $X_i = x_i$ .
  - For an acyclic graph, the MAP assignment (assuming there are no ties) is given by:

$$\underline{x_i^*} = \arg \max_{\underline{x_i}} \underline{b_i(x_i)}$$

# Deterministic Annealing


**Motivation:** Smoothly transition from sum-product to max-product

1. Incorporate inverse temperature parameter into each factor:

**Annealed Joint Distribution**

$$p(\mathbf{x}) = \frac{1}{Z} \prod_{\alpha} \left[ \psi_{\alpha}(\mathbf{x}_{\alpha})^{\frac{1}{T}} \right]$$

1. Send messages as usual for sum-product BP
2. Anneal  $T$  from  $1$  to  $0$ :

|                   |                                                                                                   |
|-------------------|---------------------------------------------------------------------------------------------------|
| $T = 1$           | Sum-product  |
| $T \rightarrow 0$ | Max-product                                                                                       |

3. Take resulting beliefs to power  $T$

# Semirings

- Sum-product  $+/*$  and max-product  $\max/*$  are commutative semirings
- We can run BP with any such commutative semiring

$$\mu_{i \rightarrow \alpha}(x_i) = \prod_{\alpha \in \mathcal{N}(i) \setminus \alpha} \mu_{\alpha \rightarrow i}(x_i)$$

$$\mu_{\alpha \rightarrow i}(x_i) = \sum_{\mathbf{x}_\alpha: \mathbf{x}_\alpha[i] = x_i} \psi_\alpha(\mathbf{x}_\alpha) \prod_{j \in \mathcal{N}(\alpha) \setminus i} \mu_{j \rightarrow \alpha}(\mathbf{x}_\alpha[j])$$

- In practice, multiplying many small numbers together can yield underflow
  - instead of using  $+/*$ , we use log-add/ $+$
  - Instead of using  $\max/*$ , we use  $\max/+$

$$\otimes = +$$

$$\oplus = \log\text{-add}(a, b) = \log(\exp(a) + \exp(b))$$

Exact inference for linear chain models

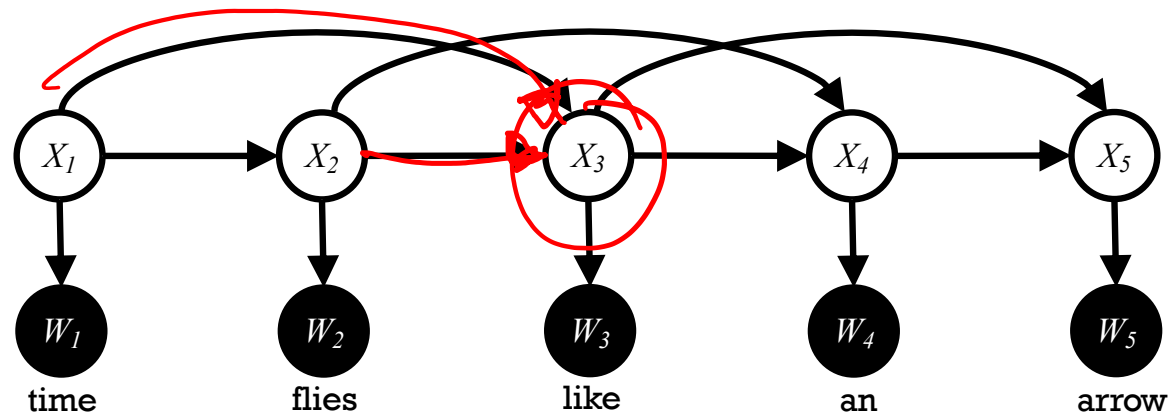
# **FORWARD-BACKWARD AND VITERBI ALGORITHMS**

# Forward-Backward Algorithm

- Sum-product BP on an HMM is called the **forward-backward algorithm**
- Max-product BP on an HMM is called the **Viterbi algorithm**

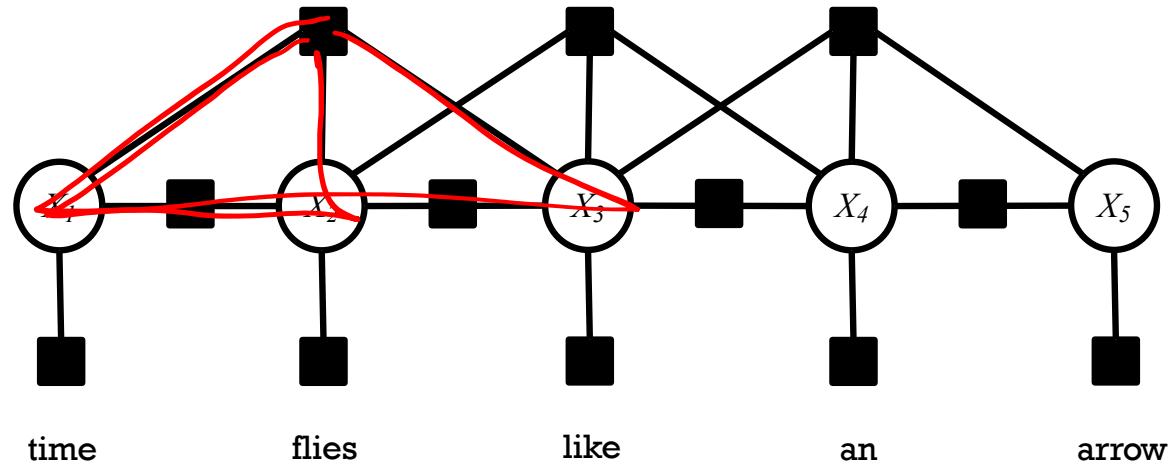
# Forward-Backward Algorithm

Trigram HMM is not a tree, even when converted to a factor graph



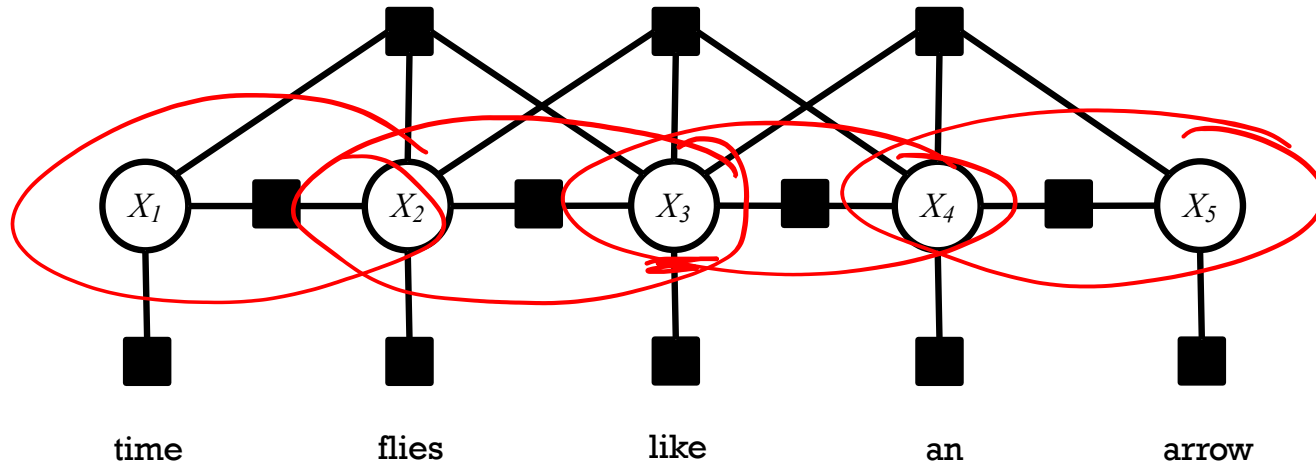
# Forward-Backward Algorithm

Trigram HMM is not a tree, even when converted to a factor graph



# Forward-Backward Algorithm

Trigram HMM is not a tree, even when converted to a factor graph



**Trick: (See also Sha & Pereira (2003))**

- Replace each variable domain with its cross product  
e.g.  $\{B, I, O\} \rightarrow \{BB, BI, BO, IB, II, IO, OB, OI, OO\}$
- Replace each pair of variables with a single one. For all  $i$ ,  $y_{i,i+1} = (x_i, x_{i+1})$
- Add features with weight  $-\infty$  that disallow illegal configurations between pairs of the new variables  
e.g. **legal** = BI and IO **illegal** = II and OO
- This is effectively a special case of the junction tree algorithm



# Summary

## 1. Factor Graphs

- Alternative representation of directed / undirected graphical models
- Make the cliques of an undirected GM explicit

## 2. Variable Elimination

- Simple and general approach to exact inference
- Just a matter of being clever when computing sum-products

## 3. Sum-product Belief Propagation

- Computes all the marginals and the partition function in only twice the work of Variable Elimination

## 4. Max-product Belief Propagation

- Identical to sum-product BP, but changes the semiring
- Computes: max-marginals, probability of MAP assignment, and (with backpointers) the MAP assignment itself.

# LEARNING FOR MRFS + CRFs

# Machine Learning

The **data** inspires  
the structures  
we want to  
predict



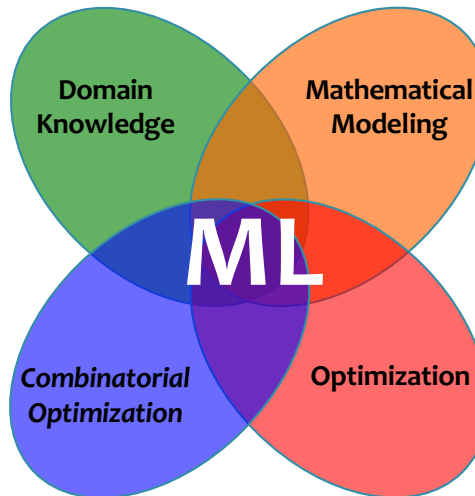
Our **model**  
defines a score  
for each structure

It also tells us  
what to optimize



**Inference** finds  
{best structure, marginals,  
partition function} for a  
new observation

(**Inference** is usually  
called as a subroutine  
in learning)



**Learning** tunes the  
parameters of the  
model

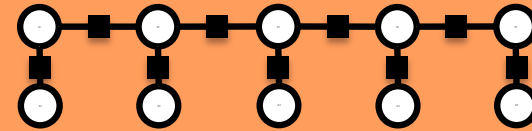
## 1. Data

$$\mathcal{D} = \{x^{(n)}\}_{n=1}^N$$

|           |            |            |           |            |            |
|-----------|------------|------------|-----------|------------|------------|
| Sample 1: | n<br>time  | v<br>flies | p<br>like | d<br>an    | n<br>from  |
| Sample 2: | n<br>time  | n<br>flies | v<br>like | d<br>an    | n<br>from  |
| Sample 3: | n<br>flies | v<br>fly   | p<br>with | n<br>their | n<br>bring |
| Sample 4: | p<br>with  | n<br>time  | n<br>you  | v<br>will  | v<br>see   |

## 2. Model

$$p(x | \theta) = \frac{1}{Z(\theta)} \prod_{C \in \mathcal{C}} \psi_C(x_C)$$



## 3. Objective

$$\ell(\theta; \mathcal{D}) = \sum_{n=1}^N \log p(x^{(n)} | \theta)$$

## 5. Inference

### 1. Marginal Inference

$$p(x_C) = \sum_{x': x'_C = x_C} p(x' | \theta)$$

### 2. Partition Function

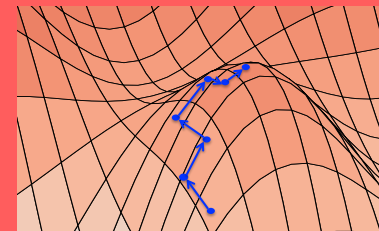
$$Z(\theta) = \sum_x \prod_{C \in \mathcal{C}} \psi_C(x_C)$$

### 3. MAP Inference

$$\hat{x} = \operatorname{argmax}_x p(x | \theta)$$

## 4. Learning

$$\theta^* = \operatorname{argmax}_{\theta} \ell(\theta; \mathcal{D})$$

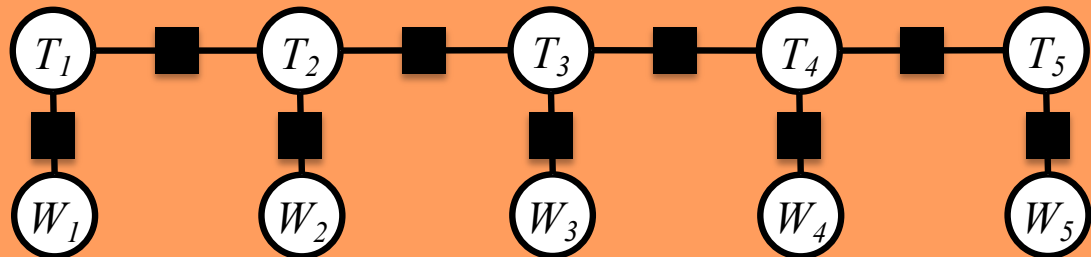


# 1. Data

Given training examples:  $\mathcal{D} = \{x^{(n)}\}_{n=1}^N$

|           |                               |                               |                              |                               |                               |
|-----------|-------------------------------|-------------------------------|------------------------------|-------------------------------|-------------------------------|
| Sample 1: | <div>n</div> <div>time</div>  | <div>v</div> <div>flies</div> | <div>p</div> <div>like</div> | <div>d</div> <div>an</div>    | <div>n</div> <div>arrow</div> |
| Sample 2: | <div>n</div> <div>time</div>  | <div>n</div> <div>flies</div> | <div>v</div> <div>like</div> | <div>d</div> <div>an</div>    | <div>n</div> <div>arrow</div> |
| Sample 3: | <div>n</div> <div>flies</div> | <div>v</div> <div>fly</div>   | <div>p</div> <div>with</div> | <div>n</div> <div>their</div> | <div>n</div> <div>wings</div> |
| Sample 4: | <div>p</div> <div>with</div>  | <div>n</div> <div>time</div>  | <div>n</div> <div>you</div>  | <div>v</div> <div>will</div>  | <div>v</div> <div>see</div>   |

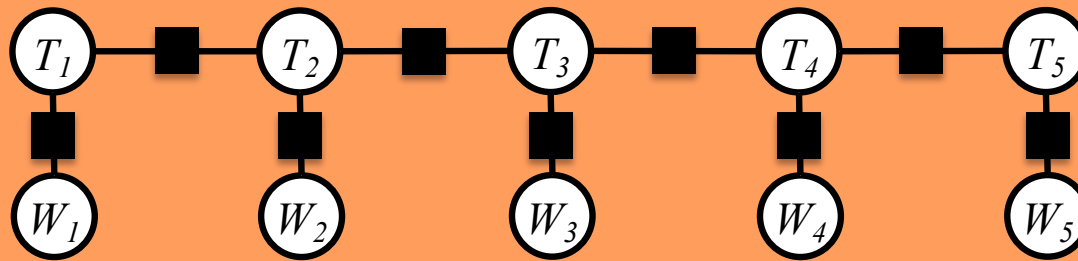
## 2. Model



## 2. Model

Define the model to be an MRF:

$$p(\mathbf{x} \mid \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \prod_{C \in \mathcal{C}} \psi_C(\mathbf{x}_C)$$



## 3. Objective

Choose the objective to be log-likelihood:

(Assign high probability to the things we observe and low probability to everything else)

$$\ell(\theta; \mathcal{D}) = \sum_{n=1}^N \log p(\mathbf{x}^{(n)} \mid \boldsymbol{\theta})$$

### 3. Objective

Choose the objective to be log-likelihood:

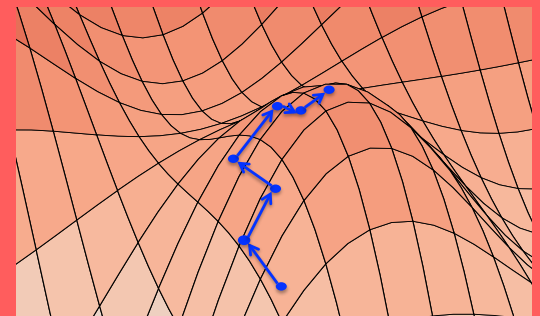
(Assign high probability to the things we observe and low probability to everything else)

$$\ell(\theta; \mathcal{D}) = \sum_{n=1}^N \log p(\mathbf{x}^{(n)} \mid \theta)$$

### 4. Learning

Tune the parameters to maximize the objective function

$$\theta^* = \operatorname{argmax}_{\theta} \ell(\theta; \mathcal{D})$$



# 3. Objective

Choose the objective to be log-likelihood:

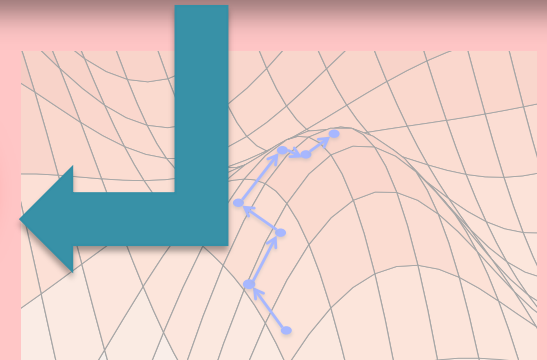
(Assign high probability to the things we observe and low probability to everything else)

## Goals for Today's Lecture

1. Consider different parameterizations
2. Optimize this objective function

Tune the parameter function

$$\theta^* = \operatorname{argmax}_{\theta} \ell(\theta; \mathcal{D})$$





# 5. Inference

Three Tasks:

## 1. Marginal Inference

Compute marginals of variables and cliques

$$p(x_i) = \sum_{\mathbf{x}' : x'_i = x_i} p(\mathbf{x}' \mid \boldsymbol{\theta}) \quad \Bigg| \quad p(\mathbf{x}_C) = \sum_{\mathbf{x}' : \mathbf{x}'_C = \mathbf{x}_C} p(\mathbf{x}' \mid \boldsymbol{\theta})$$

## 2. Partition Function

Compute the normalization constant

$$Z(\boldsymbol{\theta}) = \sum_{\mathbf{x}} \prod_{C \in \mathcal{C}} \psi_C(\mathbf{x}_C)$$

## 3. MAP Inference

Compute variable assignment with highest probability

$$\hat{\mathbf{x}} = \operatorname{argmax}_{\mathbf{x}} p(\mathbf{x} \mid \boldsymbol{\theta})$$

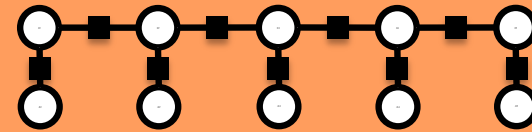
## 1. Data

$$\mathcal{D} = \{\mathbf{x}^{(n)}\}_{n=1}^N$$

|           |                                                                                   |                                                                                   |                                                                                   |                                                                                   |                                                                                   |
|-----------|-----------------------------------------------------------------------------------|-----------------------------------------------------------------------------------|-----------------------------------------------------------------------------------|-----------------------------------------------------------------------------------|-----------------------------------------------------------------------------------|
| Sample 1: |  |  |  |  |  |
|           | time                                                                              | flies                                                                             | like                                                                              | an                                                                                | arrow                                                                             |
| Sample 2: |  |  |  |  |  |
|           | time                                                                              | flies                                                                             | like                                                                              | an                                                                                | arrow                                                                             |
| Sample 3: |  |  |  |  |  |
|           | flies                                                                             | fly                                                                               | with                                                                              | their                                                                             | ring                                                                              |
| Sample 4: |  |  |  |  |  |
|           | with                                                                              | time                                                                              | you                                                                               | will                                                                              | see                                                                               |

## 2. Model

$$p(\mathbf{x} \mid \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \prod_{C \in \mathcal{C}} \psi_C(\mathbf{x}_C)$$



## 3. Objective

$$\ell(\boldsymbol{\theta}; \mathcal{D}) = \sum_{n=1}^N \log p(\mathbf{x}^{(n)} \mid \boldsymbol{\theta})$$

## 5. Inference

### 1. Marginal Inference

$$p(\mathbf{x}_C) = \sum_{\mathbf{x}': \mathbf{x}'_C = \mathbf{x}_C} p(\mathbf{x}' \mid \boldsymbol{\theta})$$

### 2. Partition Function

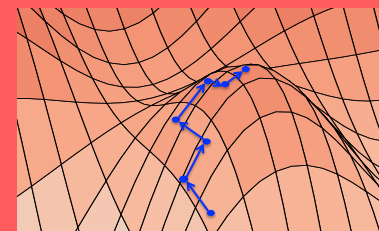
$$Z(\boldsymbol{\theta}) = \sum_{\mathbf{x}} \prod_{C \in \mathcal{C}} \psi_C(\mathbf{x}_C)$$

### 3. MAP Inference

$$\hat{\mathbf{x}} = \operatorname{argmax}_{\mathbf{x}} p(\mathbf{x} \mid \boldsymbol{\theta})$$

## 4. Learning

$$\boldsymbol{\theta}^* = \operatorname{argmax}_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; \mathcal{D})$$



# MLE for Undirected GMs

- Today's parameter estimation assumptions:
  1. The graphical model structure is given
  2. Every variable appears in the training examples

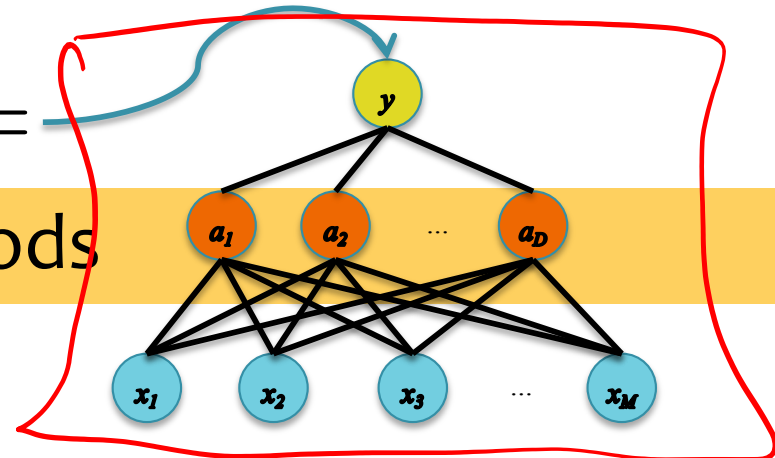
# Questions

1. What does the **likelihood objective** accomplish?
2. Is likelihood the *right objective* function?
3. **How do we optimize** the objective function (i.e. learn)?
4. What **guarantees** does the optimizer provide?
5. (What is the **mapping from data → model**? In what ways can we incorporate our domain knowledge? How does this impact learning?)

# Options for MLE of MRFs

- **Setting I:**  $\psi_C(\mathbf{x}_C) = \theta_{C, \mathbf{x}_C}$ 
  - A. MLE by inspection (Decomposable Models)
  - B. Iterative Proportional Fitting (IPF)
- **Setting II:**  $\psi_C(\mathbf{x}_C) = \exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}_C))$ 
  - C. Generalized Iterative Scaling
  - D. Gradient-based Methods

- **Setting III:**  $\psi_C(\mathbf{x}_C) =$ 
  - E. Gradient-based Methods

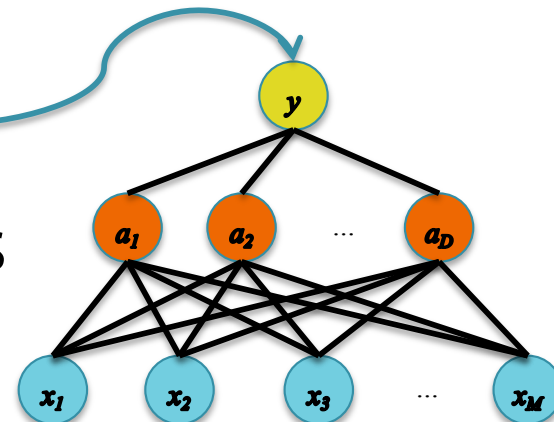


# **LOG-LINEAR PARAMETERIZATION OF CONDITIONAL RANDOM FIELD**

# Options for MLE of MRFs

- **Setting I:**  $\psi_C(\mathbf{x}_C) = \theta_{C, \mathbf{x}_C}$ 
  - A. MLE by inspection (Decomposable Models)
  - B. Iterative Proportional Fitting (IPF)
- **Setting II:**  $\psi_C(\mathbf{x}_C) = \exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}_C))$ 
  - C. Generalized Iterative Scaling
  - D. Gradient-based Methods

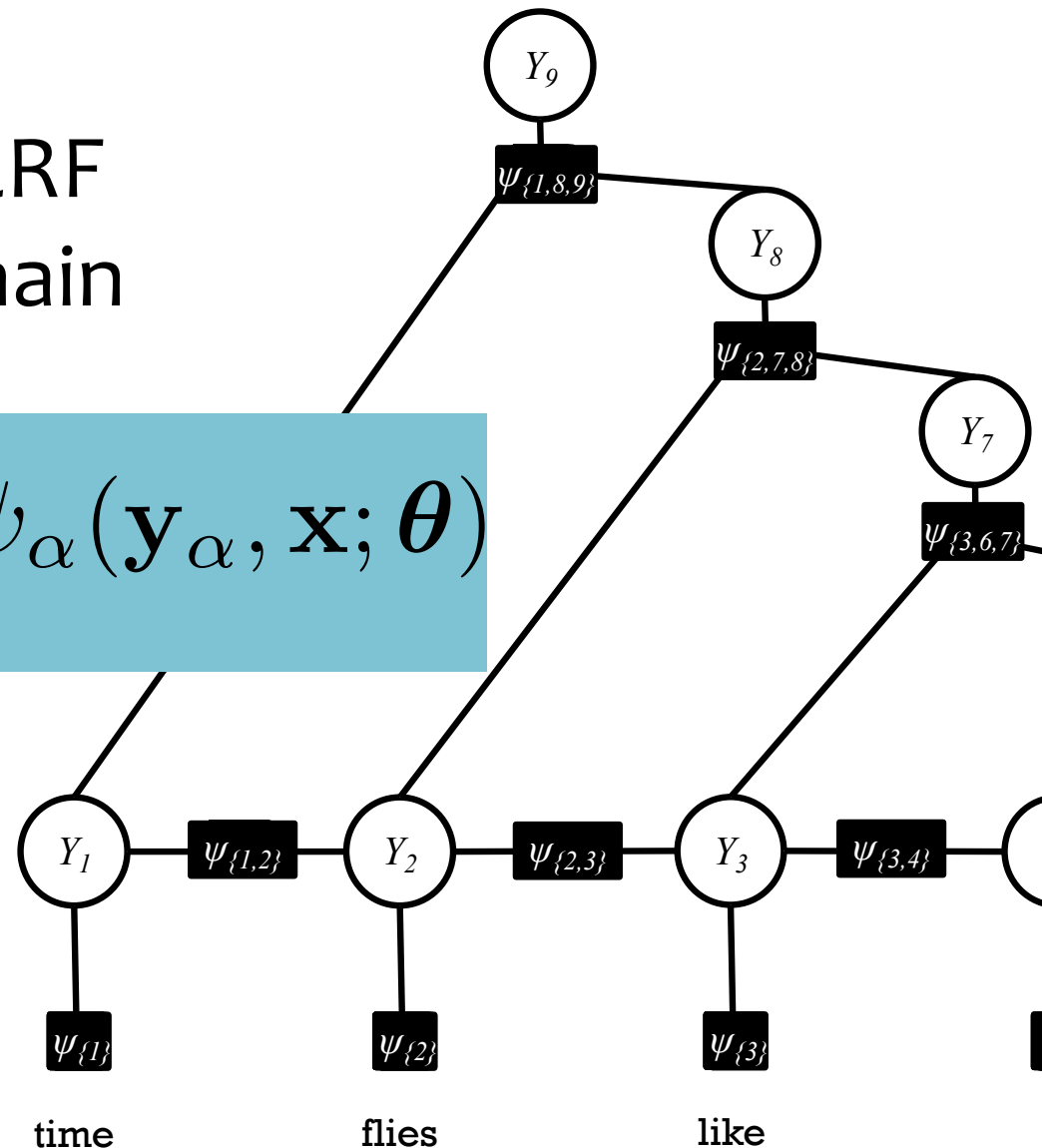
- **Setting III:**  $\psi_C(\mathbf{x}_C) =$ 
  - E. Gradient-based Methods



# General CRF

The topology of the graphical model for a CRF doesn't have to be a chain

$$p_{\theta}(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{\alpha} \psi_{\alpha}(\mathbf{y}_{\alpha}, \mathbf{x}; \theta)$$





# Log-linear CRF Parameterization

$$p_{\theta}(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{\alpha} \psi_{\alpha}(\mathbf{y}_{\alpha}, \mathbf{x}; \theta)$$

Define each potential function in terms of a fixed set of feature functions:

$$\psi_{\alpha}(\mathbf{y}_{\alpha}, \mathbf{x}; \theta) = \exp(\theta \cdot \mathbf{f}_{\alpha}(\mathbf{y}_{\alpha}, \mathbf{x}))$$

Predicted  
variables

Observed  
variables

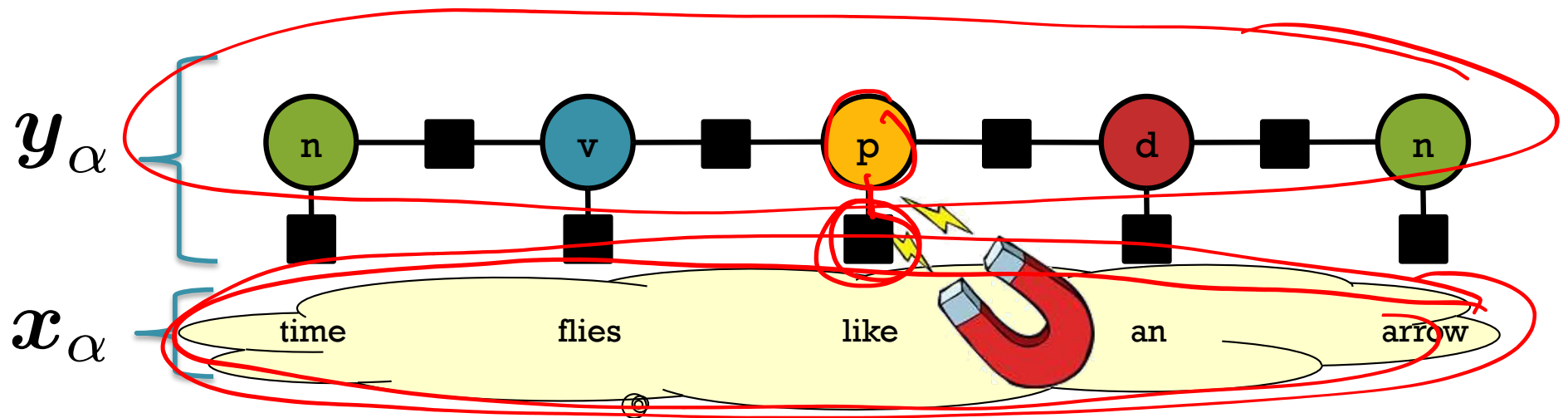
positive  
or  
non-negative

$$|\theta| = |f_{\alpha}(\cdot)|$$

# Log-linear CRF Parameterization

Define each potential function in terms of a fixed set of feature functions:

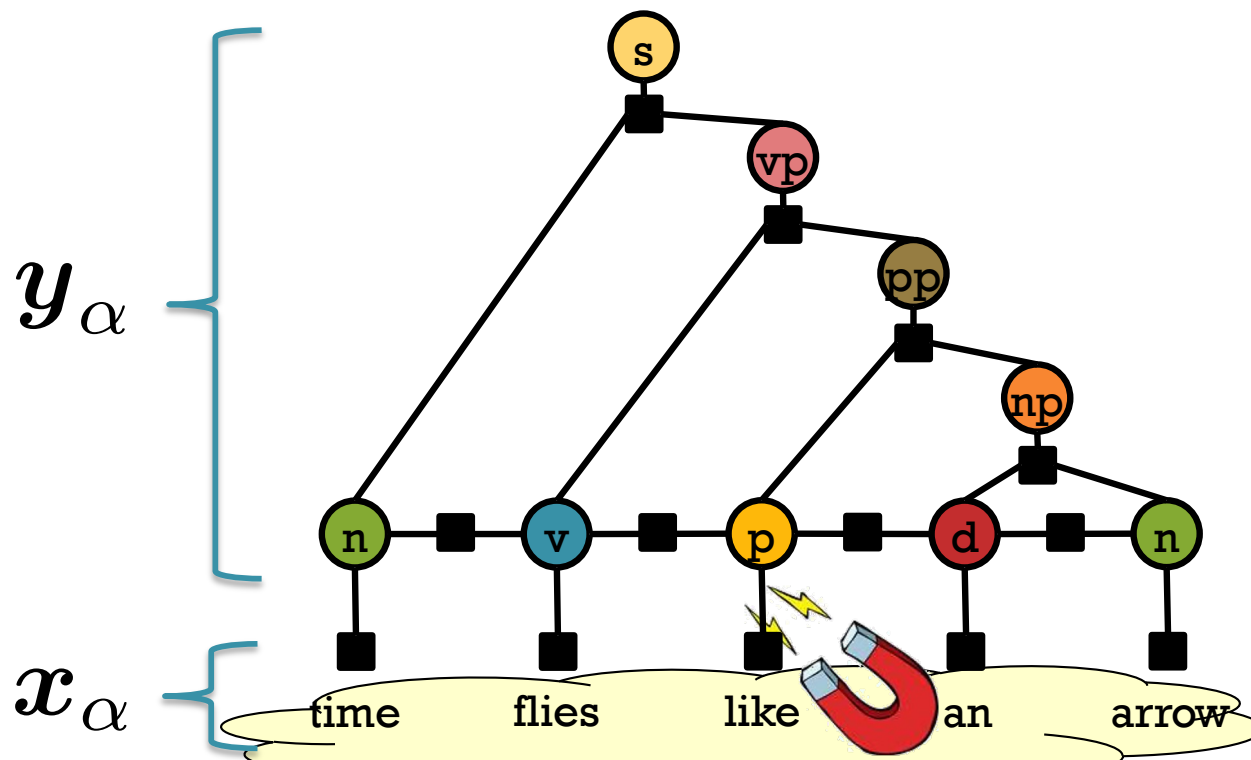
$$\psi_{\alpha}(\mathbf{y}_{\alpha}, \mathbf{x}; \boldsymbol{\theta}) = \exp(\boldsymbol{\theta} \cdot \mathbf{f}_{\alpha}(\mathbf{y}_{\alpha}, \mathbf{x}))$$



# Log-linear CRF Parameterization

Define each potential function in terms of a fixed set of feature functions:

$$\psi_{\alpha}(\mathbf{y}_{\alpha}, \mathbf{x}; \boldsymbol{\theta}) = \exp(\boldsymbol{\theta} \cdot \mathbf{f}_{\alpha}(\mathbf{y}_{\alpha}, \mathbf{x}))$$



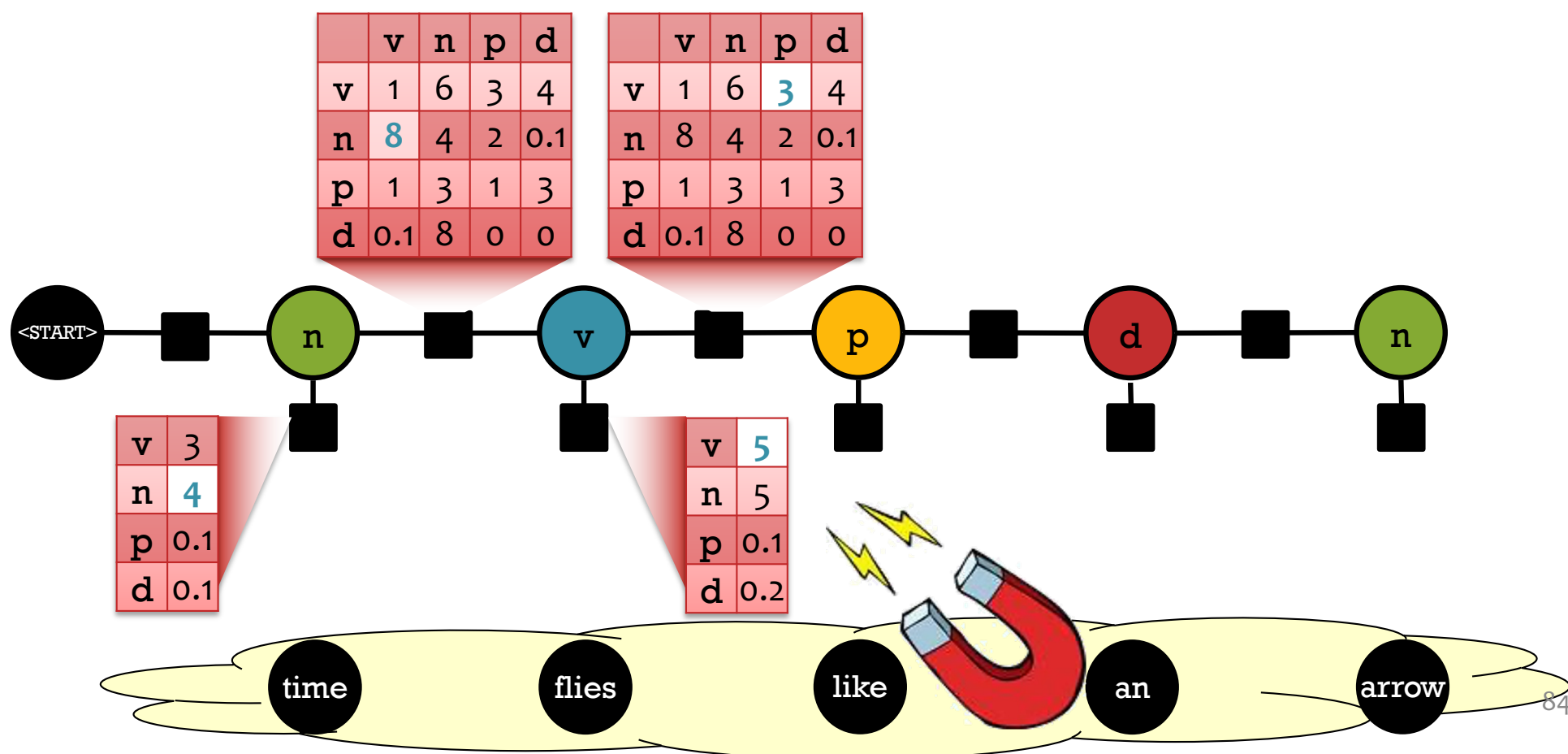
Conditional Random Fields (CRFs) for time series data

# **LINEAR-CHAIN CRFS** **(LOG-LINEAR PARAMETERIZATION)**

# Conditional Random Field (CRF)

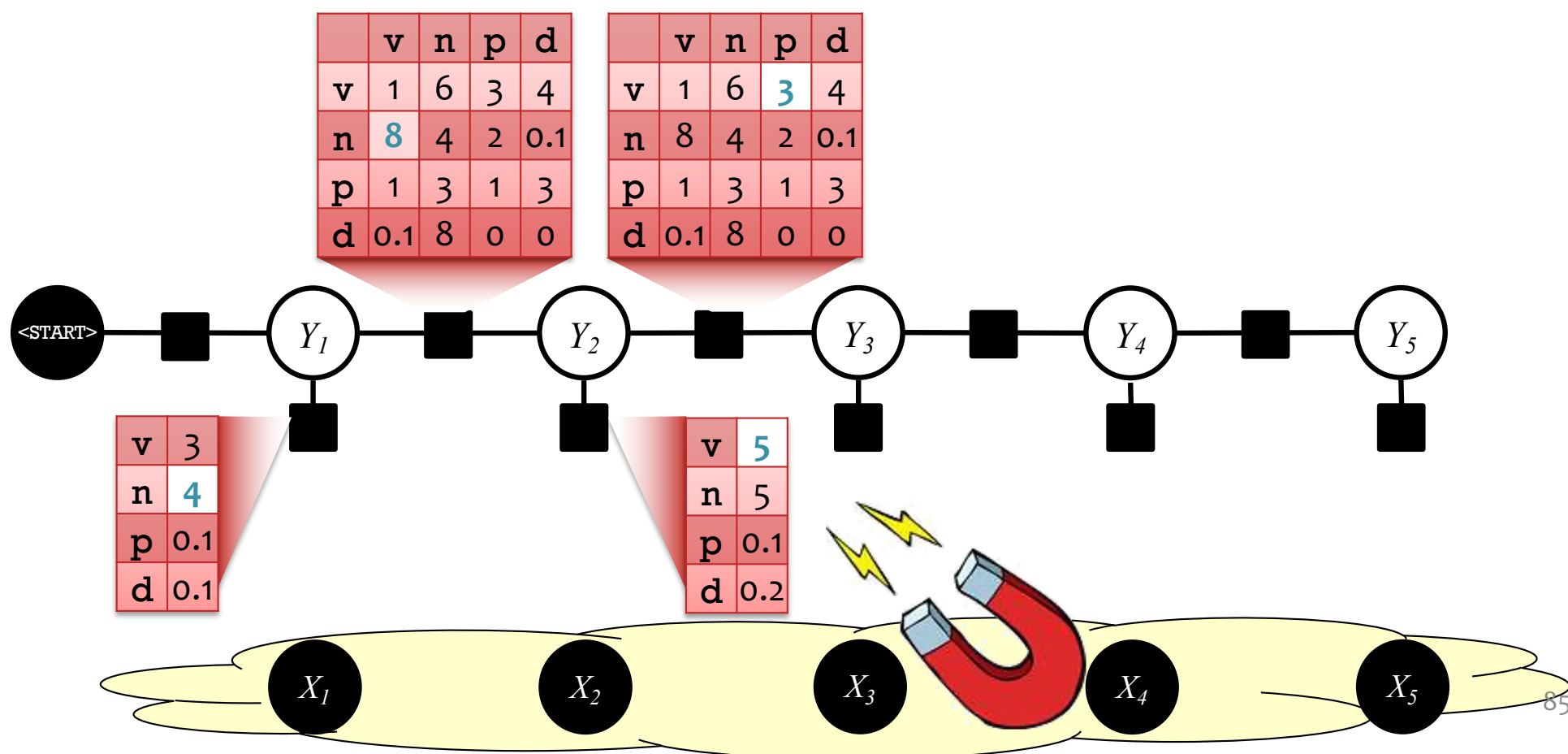
Conditional distribution over tags  $X_i$  given words  $w_i$ .  
The factors and  $Z$  are now specific to the sentence  $w$ .

$$p(n, v, p, d, n \mid \text{time, flies, like, an, arrow}) = \frac{1}{Z} (4 * 8 * 5 * 3 * \dots)$$



# Conditional Random Field (CRF)

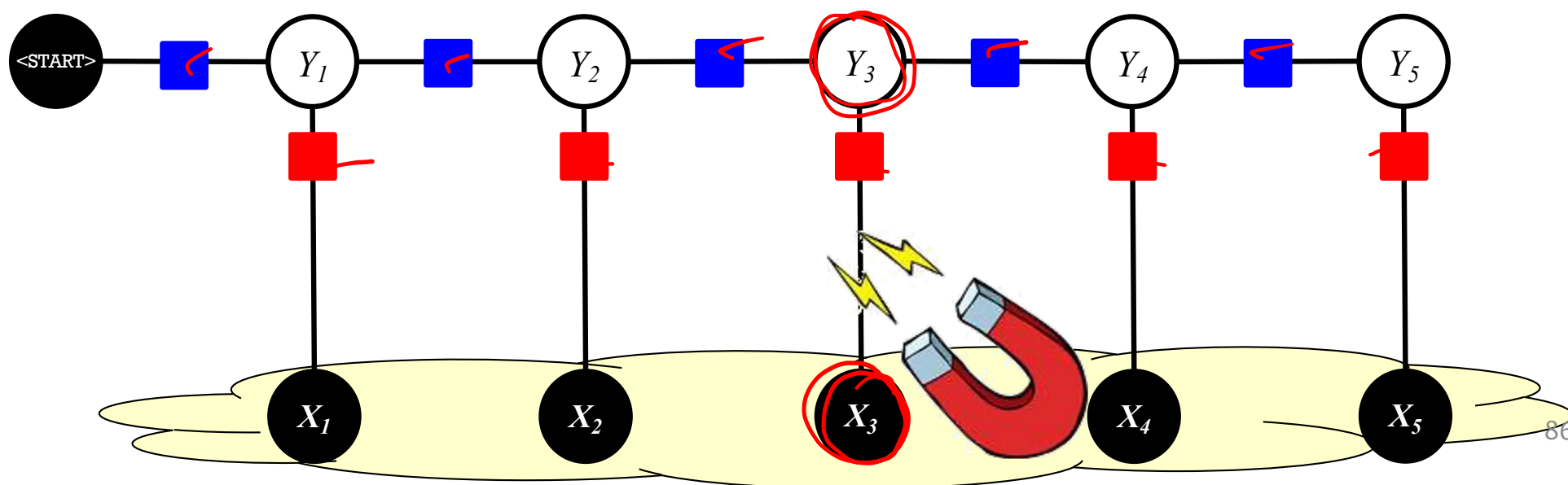
Recall: Shaded nodes in a graphical model are **observed**



# Conditional Random Field (CRF)

This **linear-chain CRF** is just **like an HMM**, except that its factors are **not** necessarily probability distributions

$$\begin{aligned}
 p(\mathbf{y}|\mathbf{x}) &= \frac{1}{Z(\mathbf{x})} \prod_{k=1}^K \underbrace{\psi_{\text{em}}(y_k, x_k)}_{\text{red}} \underbrace{\psi_{\text{tr}}(y_k, y_{k-1})}_{\text{blue}} \\
 &= \frac{1}{Z(\mathbf{x})} \prod_{k=1}^K \underbrace{\exp(\theta \cdot \mathbf{f}_{\text{em}}(y_k, x_k))}_{\text{red}} \underbrace{\exp(\theta \cdot \mathbf{f}_{\text{tr}}(y_k, y_{k-1}))}_{\text{blue}}
 \end{aligned}$$



# Exercise

Handwritten notes:  $\vec{y} \in \mathcal{Y}_{\vec{x}}$ ,  $y \in \mathcal{Y}_{\vec{x}}$ ,  $|\mathcal{Y}_{\vec{x}}| = T^k$ ,  $|\mathcal{Y}'_{\vec{x}}| = T^k$

$$p(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{k=1}^K \psi_{\text{em}}(y_k, x_k) \psi_{\text{tr}}(y_k, y_{k-1})$$

Handwritten note:  $\vec{y} = \text{seq}(\vec{y})$

$$= \frac{1}{Z(\mathbf{x})} \prod_{k=1}^K \exp(\boldsymbol{\theta} \cdot \mathbf{f}_{\text{em}}(y_k, x_k)) \exp(\boldsymbol{\theta} \cdot \mathbf{f}_{\text{tr}}(y_k, y_{k-1}))$$

Handwritten note:  $= \exp(\boldsymbol{\theta} [\mathbf{f}_{\text{em}}(\cdot) + \mathbf{f}_{\text{tr}}(\cdot)])$

Q1

**Select All That Apply:** Which model does the above distribution share the most in common with?

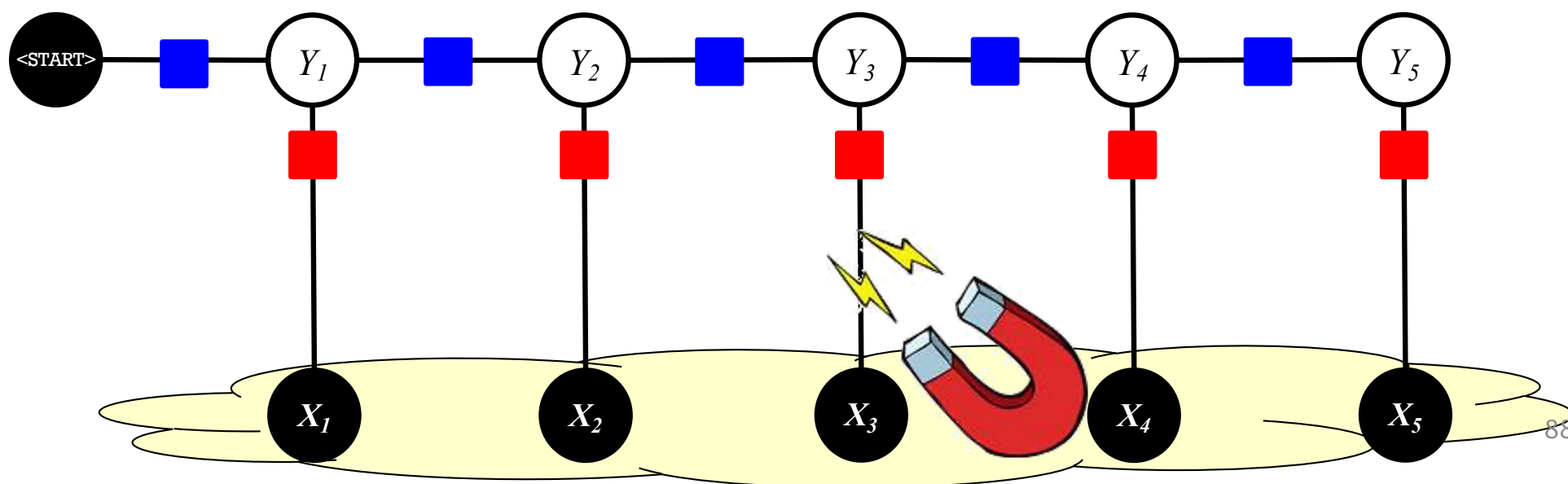
- A. Hidden Markov Model 87% ☒
  - B. Bernoulli Naïve Bayes 32% ☐
  - C. Gaussian Naïve Bayes 53% ☐
  - D. Logistic Regression 31% ☐
- Handwritten note: *Something*



# Conditional Random Field (CRF)

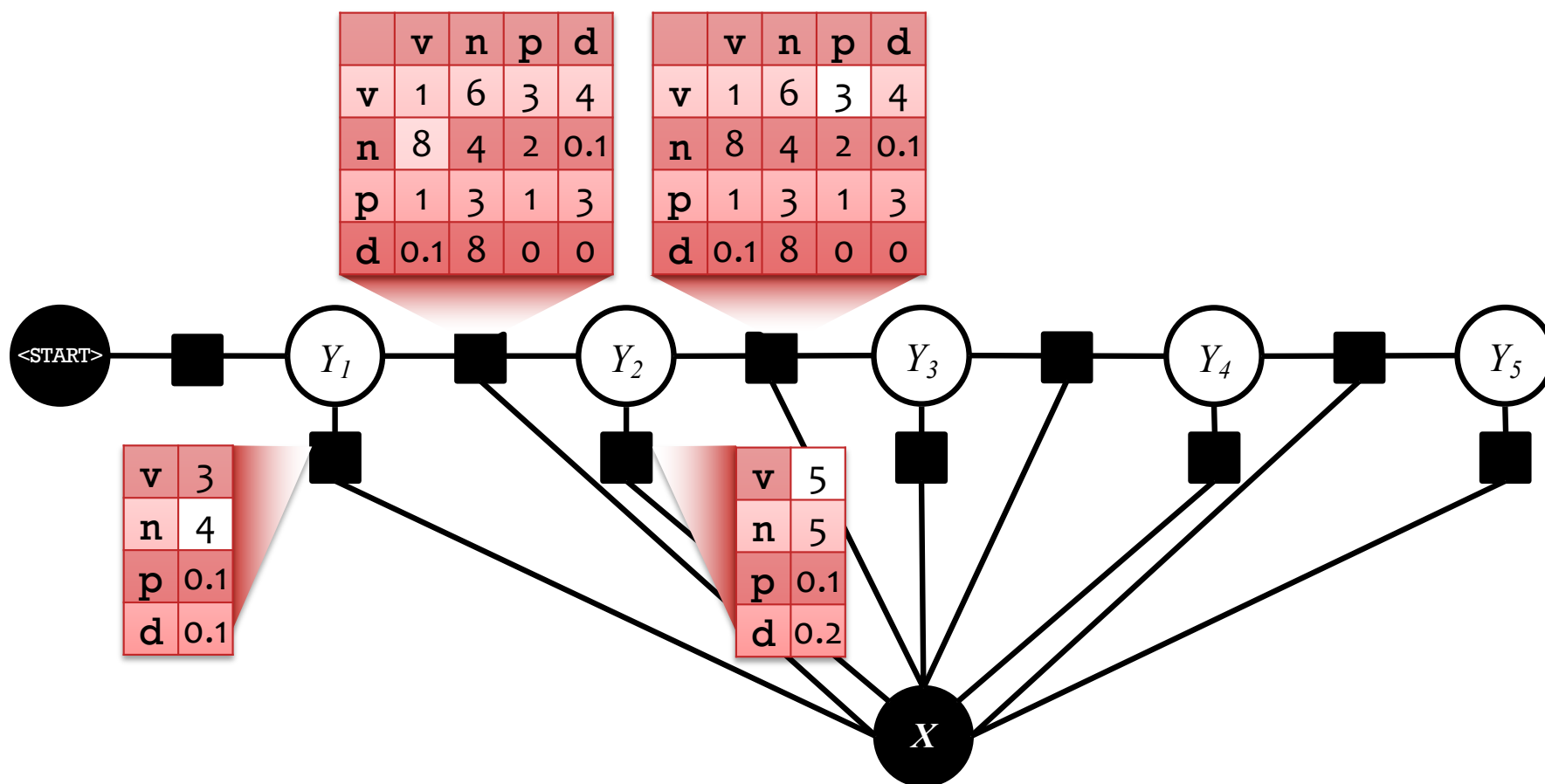
This **linear-chain CRF** is just **like an HMM**, except that its factors are **not** necessarily probability distributions

$$\begin{aligned} p(\mathbf{y}|\mathbf{x}) &= \frac{1}{Z(\mathbf{x})} \prod_{k=1}^K \psi_{\text{em}}(y_k, x_k) \psi_{\text{tr}}(y_k, y_{k-1}) \\ &= \frac{1}{Z(\mathbf{x})} \prod_{k=1}^K \exp(\boldsymbol{\theta} \cdot \mathbf{f}_{\text{em}}(y_k, x_k)) \exp(\boldsymbol{\theta} \cdot \mathbf{f}_{\text{tr}}(y_k, y_{k-1})) \end{aligned}$$



# Conditional Random Field (CRF)

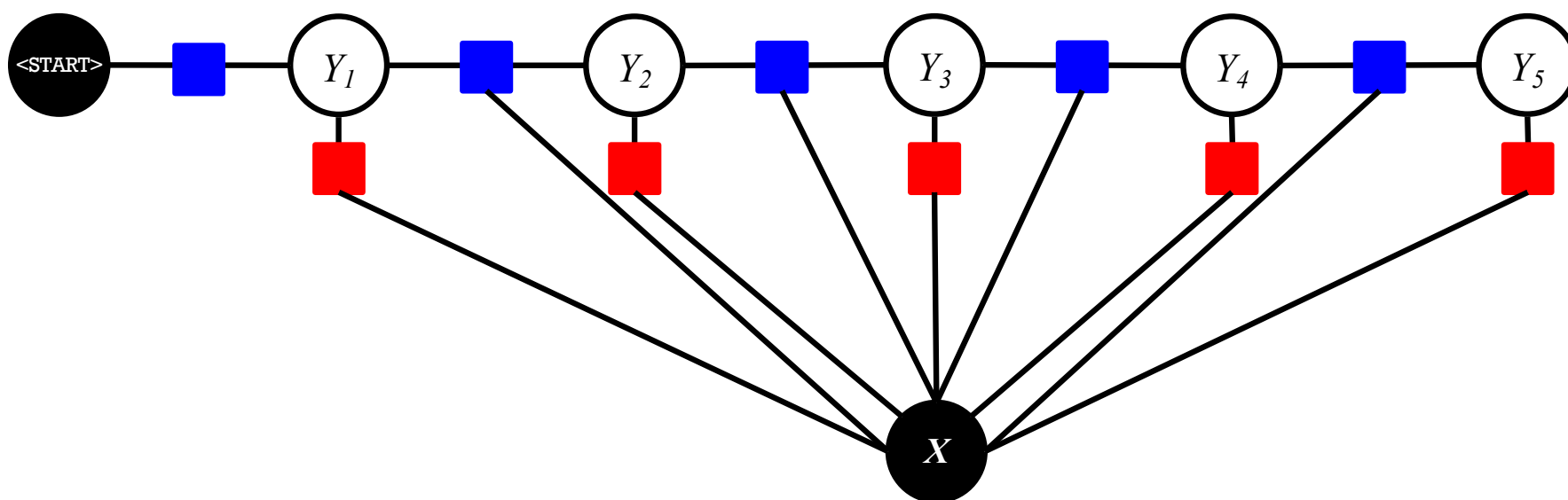
- That is the **vector**  $X$
- Because it's observed, we can condition on it for free
- Conditioning is how we converted from the MRF to the CRF (i.e. when taking a slice of the emission factors)



# Conditional Random Field (CRF)

- This is the **standard** linear-chain CRF definition
- It permits rich, overlapping features of the vector  $\mathbf{x}$

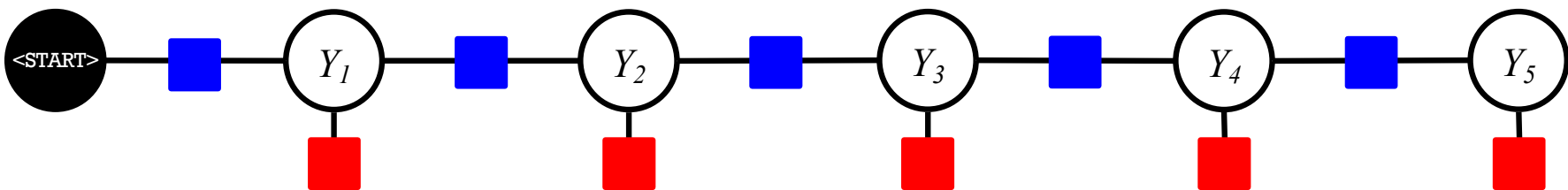
$$p(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{k=1}^K \psi_{\text{em}}(y_k, \mathbf{x}) \psi_{\text{tr}}(y_k, y_{k-1}, \mathbf{x})$$
$$= \frac{1}{Z(\mathbf{x})} \prod_{k=1}^K \exp(\theta \cdot \mathbf{f}_{\text{em}}(y_k, \mathbf{x})) \exp(\theta \cdot \mathbf{f}_{\text{tr}}(y_k, y_{k-1}, \mathbf{x}))$$



# Conditional Random Field (CRF)

- This is the **standard** linear-chain CRF definition
- It permits rich, overlapping features of the vector  $\mathbf{X}$

$$p(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{k=1}^K \psi_{\text{em}}(y_k, \mathbf{x}) \psi_{\text{tr}}(y_k, y_{k-1}, \mathbf{x})$$
$$= \frac{1}{Z(\mathbf{x})} \prod_{k=1}^K \exp(\boldsymbol{\theta} \cdot \mathbf{f}_{\text{em}}(y_k, \mathbf{x})) \exp(\boldsymbol{\theta} \cdot \mathbf{f}_{\text{tr}}(y_k, y_{k-1}, \mathbf{x}))$$



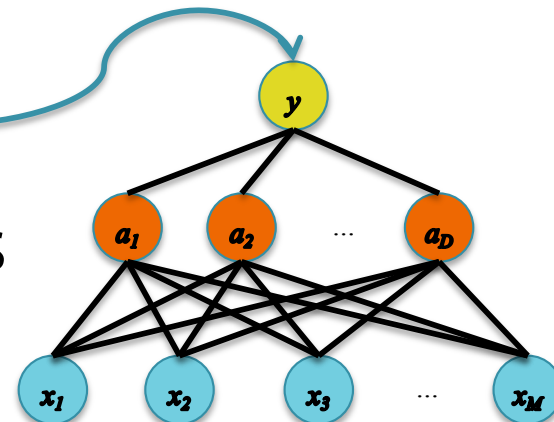
**Visual Notation:** Usually we draw a CRF **without** showing the variable corresponding to  $\mathbf{X}$

# **MRF AND CRF LEARNING (LOG-LINEAR PARAMETERIZATION)**

# Options for MLE of MRFs

- **Setting I:**  $\psi_C(\mathbf{x}_C) = \theta_{C, \mathbf{x}_C}$ 
  - A. MLE by inspection (Decomposable Models)
  - B. Iterative Proportional Fitting (IPF)
- **Setting II:**  $\psi_C(\mathbf{x}_C) = \exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}_C))$ 
  - C. Generalized Iterative Scaling
  - D. Gradient-based Methods

- **Setting III:**  $\psi_C(\mathbf{x}_C) =$ 
  - E. Gradient-based Methods



# MRF and CRF Learning

## ***Whiteboard***

- log-linear MRF model (i.e. with feature based potentials)
- log-linear MRF derivatives
- log-linear MRF training with SGD
- log-linear CRF model (i.e. with feature based potentials)
- log-linear CRF derivatives
- log-linear CRF training with SGD