



10-708 Probabilistic Graphical Models

Machine Learning Department
School of Computer Science
Carnegie Mellon University



Variable Elimination + Belief Propagation

Matt Gormley
Lecture 5
Feb., 15 2021

Q&A

Q: Is Homework 1 representative of future assignments?

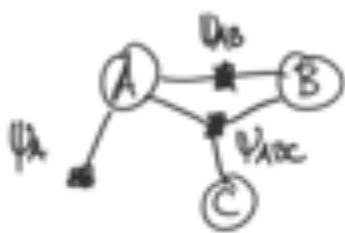
A: Not really...

Recall that the remaining assignments will involve a written *and* programming component, whereas HW1 has just a written section.

Reminders

- **Homework 1: PGM Representation**
 - Out: Mon, Feb. 15
 - Due: Mon, Feb. 22 at 11:59pm
- **Homework 2: Exact inference and supervised learning (CRF+RNN)**
 - Out: Mon, Feb. 22
 - Due: Mon, Mar. 08 at 11:59pm

Ex: Factor Graph over Binary Variables



a	$\psi_A(a)$
0	2
1	7

a	b	$\psi_{AB}(a,b)$
0	0	4
0	1	3
1	0	6
1	1	2

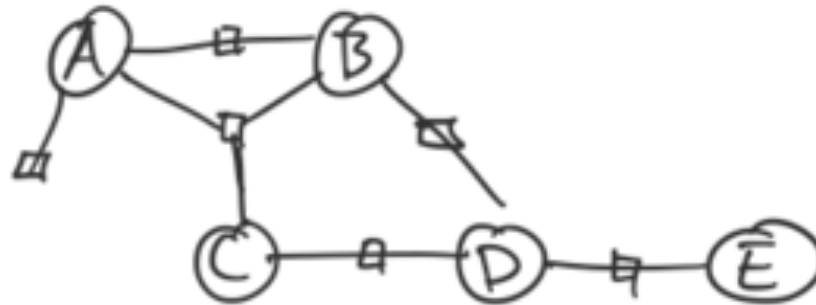
a	b	c	$\psi_{ABC}(a,b,c)$
0	0	0	6
0	0	1	2
0	1	0	1
...
1	1	1	5

$$P(A=a, B=b, C=c) = p(a,b,c) = \frac{1}{Z} \underbrace{\psi_A(a) \psi_{AB}(a,b) \psi_{ABC}(a,b,c)}_{s(a,b,c)} \Rightarrow Z = \sum_a \sum_b \sum_c s(a,b,c)$$

a	b	c	ψ_A	ψ_{AB}	ψ_{ABC}	$s(\cdot)$	$p(\cdot)$
0	0	0	2	4	6	48	48/Z
0	0	1	2	3	2	16	16/Z
0	1	0	2	3	1	6	6/Z
...
1	1	1	7	2	5	+ 70	70/Z
							Z

Ex: Marginal Inference

Ex:



Marginal Probability:

$$P(A=a) = p(a) = \sum_b \sum_c \sum_d \sum_e p(a, b, c, d, e)$$

$$P(B=b) = p(b) = \sum_a \sum_c \sum_d \sum_e p(a, b, c, d, e)$$

$$p(a, b, c) = \sum_d \sum_e p(a, b, c, d, e)$$

"marginalized out d and e"

BRUTE FORCE INFERENCE

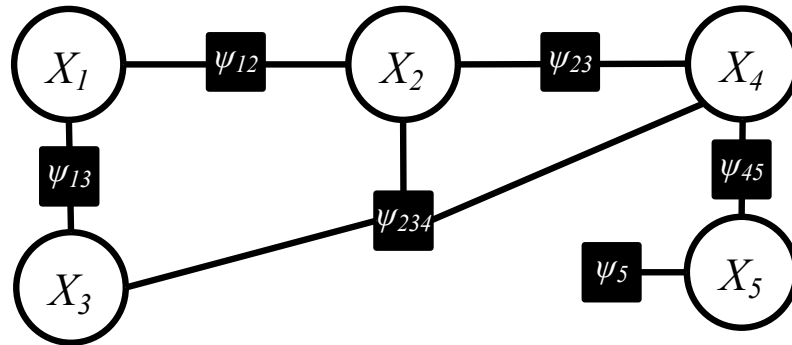
Brute Force (Naïve) Inference

For all i , suppose the **range** of X_i is $\{0, 1, 2\}$.

Let $k=3$ denote the **size of the range**.

The distribution **factorizes** as:

$$S(\mathbf{x}) = \psi_{12}(x_1, x_2) \psi_{13}(x_1, x_3) \psi_{24}(x_2, x_4) \\ \psi_{234}(x_2, x_3, x_4) \psi_{45}(x_4, x_5) \psi_5(x_5)$$



Naively, we compute the **partition function** as:

$$Z = \sum_{x_1} \sum_{x_2} \sum_{x_3} \sum_{x_4} \sum_{x_5} S(\mathbf{x})$$

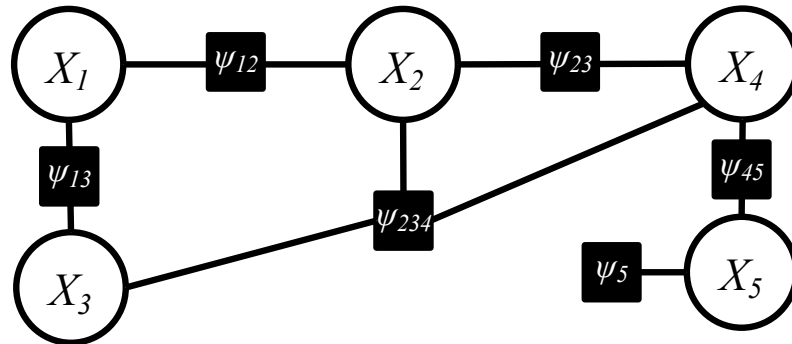
Brute Force (Naïve) Inference

For all i , suppose the **range** of X_i is $\{0, 1, 2\}$.

Let $k=3$ denote the **size of the range**.

The distribution **factorizes** as:

$$s(\mathbf{x}) = \psi_{12}(x_1, x_2) \psi_{13}(x_1, x_3) \psi_{23}(x_2, x_3) \psi_{24}(x_2, x_4) \psi_{234}(x_2, x_3, x_4) \psi_{45}(x_4, x_5) \psi_5(x_5)$$



Naively, we compute the **partition function** as:

$$Z = \sum_{x_1} \sum_{x_2} \sum_{x_3} \sum_{x_4} \sum_{x_5} s(\mathbf{x})$$

$s(\mathbf{x})$ can be represented as a joint probability table with 3^5 entries:

x_1	x_2	x_3	x_4	x_5	$s(\mathbf{x})$
0	0	0	0	0	0.019517693
0	0	0	0	1	0.017090249
0	0	0	0	2	0.014885825
0	0	0	1	0	0.024117638
0	0	0	1	1	0.000925849
0	0	0	1	2	0.028112576
0	0	0	2	0	0.028050205
0	0	0	2	1	0.004812689
0	0	0	2	2	0.007987737
0	0	1	0	0	0.028433687
0	0	1	0	1	0.037073469
0	0	1	0	2	0.013558227
0	0	1	1	0	0.019479016
0	0	1	1	1	0.012312901
0	0	1	1	2	0.023439775
0	0	1	2	0	0.038206131
0	0	1	2	1	0.038996005
0	0	1	2	2	0.041458783
0	0	2	0	0	0.044616806
0	0	2	0	1	0.020846989
0	0	2	0	2	0.03006475
0	0	2	1	0	0.048436964
0	0	2	1	1	0.02854376
0	0	2	1	2	0.029191506
0	0	2	2	0	0.031531118
0	0	2	2	1	0.005132392
0	0	2	2	2	0.032027091
...

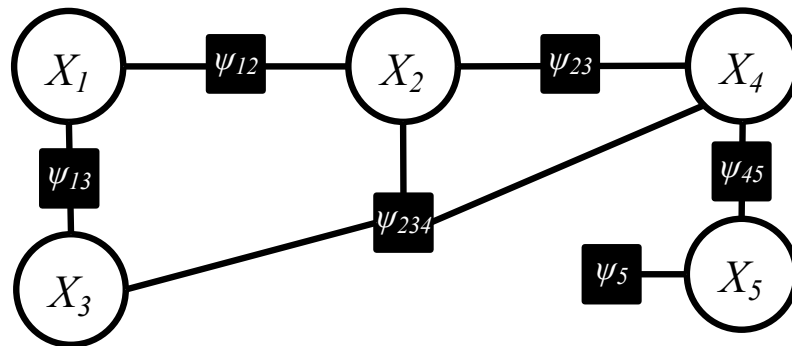
Brute Force (Naïve) Inference

For all i , suppose the **range** of X_i is $\{0, 1, 2\}$.

Let $k=3$ denote the **size of the range**.

The distribution **factorizes** as:

$$S(\mathbf{x}) = \psi_{12}(x_1, x_2) \psi_{13}(x_1, x_3) \psi_{23}(x_2, x_3) \psi_{24}(x_2, x_4) \psi_{234}(x_2, x_3, x_4) \psi_{45}(x_4, x_5) \psi_5(x_5)$$



Naively, we compute the **partition function** as:

$$Z = \sum_{x_1} \sum_{x_2} \sum_{x_3} \sum_{x_4} \sum_{x_5} S(\mathbf{x})$$

$s(\mathbf{x})$ can be represented as a joint probability table with 3^5 entries:

x_1	x_2	x_3	x_4	x_5	$s(\mathbf{x})$
0	0	0	0	0	0.019517693
0	0	0	0	1	0.017090249
0	0	0	0	2	0.014885825
0	0	0	1	0	0.024117638
0	0	0	1	1	0.000925849
0	0	0	1	2	0.028112576
0	0	0	2	0	0.028050205
0	0	0	2	1	0.004812689
0	0	0	2	2	0.007987737
0	0	1	0	0	0.028433687
0	0	1	0	1	0.037073469
0	0	1	0	2	0.013558227
0	0	1	1	0	0.019479016
0	0	1	1	1	0.012312901
0	0	1	1	2	0.023439775
0	0	1	2	0	0.038206131
0	0	1	2	1	0.038996005
0	0	1	2	2	0.041458783
0	0	2	0	0	0.044616806
0	0	2	0	1	0.020846989
0	0	2	0	2	0.03006475
0	0	2	1	0	0.048436964
0	0	2	1	1	0.02854376

Naïve computation of Z requires 3^5 additions.

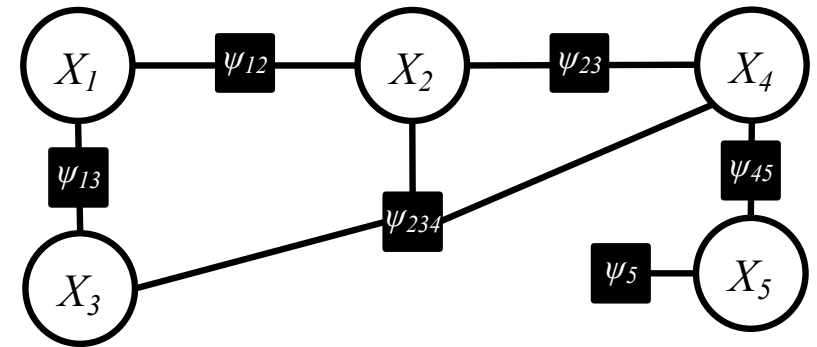
Can we do better?

Simple and general exact inference for graphical models

VARIABLE ELIMINATION

The Variable Elimination Algorithm

Instead, capitalize on the factorization of $s(\mathbf{x})$.



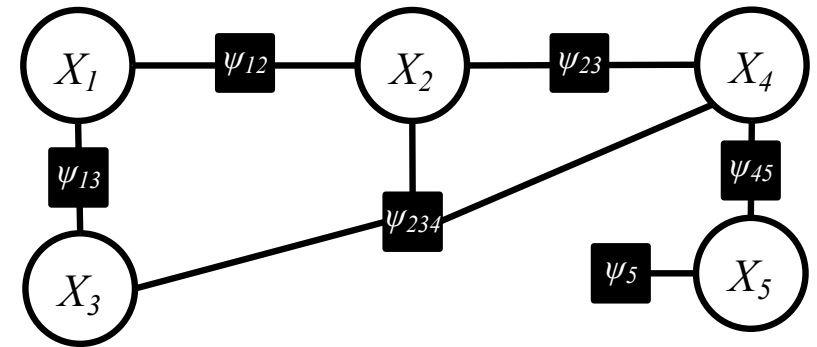
$$\begin{aligned}
 Z &= \sum_{x_1} \sum_{x_2} \sum_{x_3} \sum_{x_4} \sum_{x_5} \psi_{12}(x_1, x_2) \psi_{13}(x_1, x_3) \psi_{24}(x_2, x_4) \psi_{234}(x_2, x_3, x_4) \psi_{45}(x_4, x_5) \psi_5(x_5) \\
 &= \sum_{x_1} \sum_{x_2} \sum_{x_3} \sum_{x_4} \psi_{12}(x_1, x_2) \psi_{13}(x_1, x_3) \psi_{24}(x_2, x_4) \psi_{234}(x_2, x_3, x_4) \underbrace{\sum_{x_5} \psi_{45}(x_4, x_5) \psi_5(x_5)}_{\text{factor}}
 \end{aligned}$$

This “factor” is a much smaller table with 3^2 entries:

x_4	x_5	$s(x_4, x_5)$
0	0	0.019517693
0	1	0.017090249
0	2	0.014885825
1	0	0.024117638
1	1	0.000925849
1	2	0.028112576
2	0	0.028050205
2	1	0.004812689
2	2	0.007987737

The Variable Elimination Algorithm

Instead, capitalize on the factorization of $s(\mathbf{x})$.



$$\begin{aligned}
 Z &= \sum_{x_1} \sum_{x_2} \sum_{x_3} \sum_{x_4} \sum_{x_5} \psi_{12}(x_1, x_2) \psi_{13}(x_1, x_3) \psi_{24}(x_2, x_4) \psi_{234}(x_2, x_3, x_4) \psi_{45}(x_4, x_5) \psi_5(x_5) \\
 &= \sum_{x_1} \sum_{x_2} \sum_{x_3} \sum_{x_4} \psi_{12}(x_1, x_2) \psi_{13}(x_1, x_3) \psi_{24}(x_2, x_4) \psi_{234}(x_2, x_3, x_4) \underbrace{\sum_{x_5} \psi_{45}(x_4, x_5) \psi_5(x_5)}_{\text{factor}}
 \end{aligned}$$

Only 3^2 additions are needed to marginalize out x_5 .

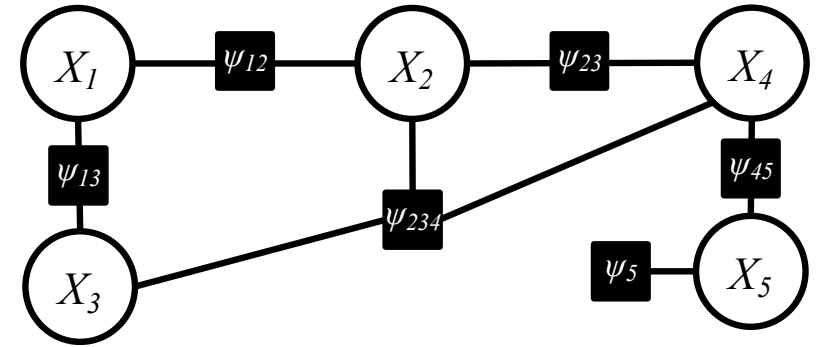
We denote the **marginal's table** by $m_5(x_4)$.

This “factor” is a much smaller table with 3 entries:

x_4	$m_5(x_4)$
0	0.019517693
1	0.017090249
2	0.014885825

The Variable Elimination Algorithm

Instead, capitalize on the factorization of $s(\mathbf{x})$.

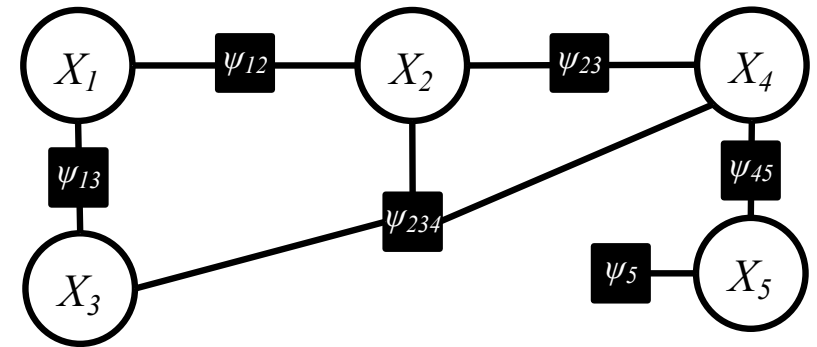


$$\begin{aligned}
 Z &= \sum_{x_1} \sum_{x_2} \sum_{x_3} \sum_{x_4} \sum_{x_5} \psi_{12}(x_1, x_2) \psi_{13}(x_1, x_3) \psi_{24}(x_2, x_4) \psi_{234}(x_2, x_3, x_4) \psi_{45}(x_4, x_5) \psi_5(x_5) \\
 &= \sum_{x_1} \sum_{x_2} \sum_{x_3} \sum_{x_4} \psi_{12}(x_1, x_2) \psi_{13}(x_1, x_3) \psi_{24}(x_2, x_4) \psi_{234}(x_2, x_3, x_4) \sum_{x_5} \psi_{45}(x_4, x_5) \psi_5(x_5) \\
 &= \sum_{x_1} \sum_{x_2} \sum_{x_3} \sum_{x_4} \psi_{12}(x_1, x_2) \psi_{13}(x_1, x_3) \psi_{24}(x_2, x_4) \psi_{234}(x_2, x_3, x_4) m_5(x_4)
 \end{aligned}$$

$$m_5(x_4) \triangleq \sum_{x_5} \psi_{45}(x_4, x_5) \psi_5(x_5)$$

The Variable Elimination Algorithm

Instead, capitalize on the factorization of $s(\mathbf{x})$.



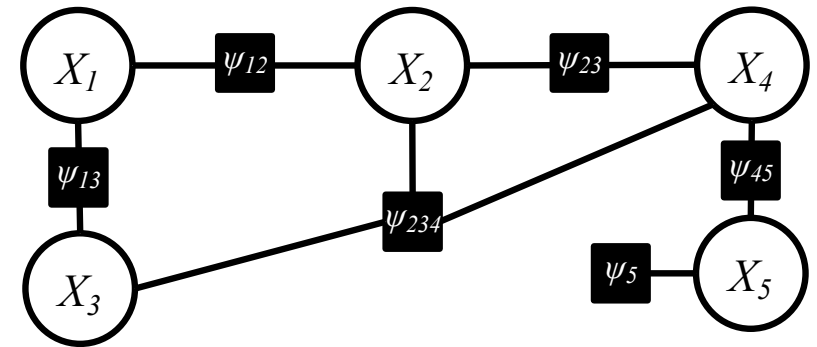
$$\begin{aligned}
 Z &= \sum_{x_1} \sum_{x_2} \sum_{x_3} \sum_{x_4} \sum_{x_5} \psi_{12}(x_1, x_2) \psi_{13}(x_1, x_3) \psi_{24}(x_2, x_4) \psi_{234}(x_2, x_3, x_4) \psi_{45}(x_4, x_5) \psi_5(x_5) \\
 &= \sum_{x_1} \sum_{x_2} \sum_{x_3} \sum_{x_4} \psi_{12}(x_1, x_2) \psi_{13}(x_1, x_3) \psi_{24}(x_2, x_4) \psi_{234}(x_2, x_3, x_4) \sum_{x_5} \psi_{45}(x_4, x_5) \psi_5(x_5) \\
 &= \sum_{x_1} \sum_{x_2} \sum_{x_3} \sum_{x_4} \underbrace{\psi_{12}(x_1, x_2) \psi_{13}(x_1, x_3) \psi_{24}(x_2, x_4) \psi_{234}(x_2, x_3, x_4)}_{m_5(x_4)} m_5(x_4)
 \end{aligned}$$

This “factor” is still a 3^4 table so apply the same trick again.

$$m_5(x_4) \triangleq \sum_{x_5} \psi_{45}(x_4, x_5) \psi_5(x_5)$$

The Variable Elimination Algorithm

Instead, capitalize on the factorization of $s(\mathbf{x})$.



$$\begin{aligned}
 Z &= \sum_{x_1} \sum_{x_2} \psi_{12}(x_1, x_2) \sum_{x_3} \psi_{13}(x_1, x_3) \sum_{x_4} \psi_{24}(x_2, x_4) \psi_{234}(x_2, x_3, x_4) \sum_{x_5} \psi_{45}(x_4, x_5) \psi_5(x_5) \\
 &= \sum_{x_1} \sum_{x_2} \psi_{12}(x_1, x_2) \sum_{x_3} \psi_{13}(x_1, x_3) \sum_{x_4} \psi_{24}(x_2, x_4) \psi_{234}(x_2, x_3, x_4) m_5(x_4) \\
 &= \sum_{x_1} \sum_{x_2} \psi_{12}(x_1, x_2) \sum_{x_3} \psi_{13}(x_1, x_3) m_4(x_2, x_3) \\
 &= \sum_{x_1} \sum_{x_2} \psi_{12}(x_1, x_2) m_3(x_1, x_2) \\
 &= \sum_{x_1} m_2(x_1)
 \end{aligned}$$

3^2 additions

3^3 additions

3^3 additions

3^2 additions

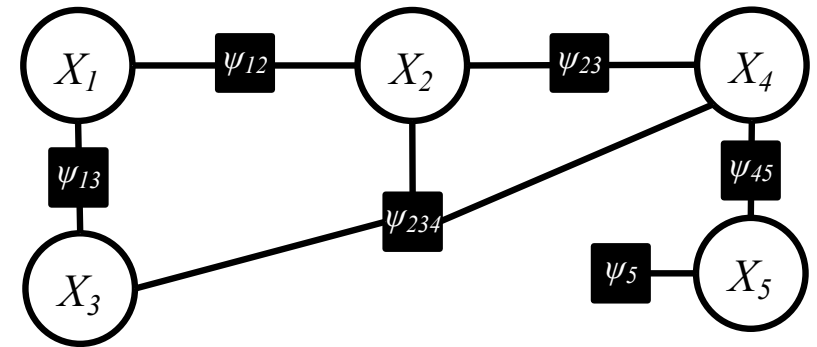
3 additions

Naïve solution requires $3^5 = 243$ additions.

Variable elimination only requires $3 + 3^2 + 3^3 + 3^3 + 3^2 = 75$ additions.

The Variable Elimination Algorithm

The same trick can be used to compute **marginal probabilities**. Just choose the variable elimination order such that the query variables are last.



$$\begin{aligned}
 p(x_1) &= \frac{1}{Z} \sum_{x_2} \psi_{12}(x_1, x_2) \sum_{x_3} \psi_{13}(x_1, x_3) \sum_{x_4} \psi_{24}(x_2, x_4) \psi_{234}(x_2, x_3, x_4) \sum_{x_5} \psi_{45}(x_4, x_5) \psi_5(x_5) \\
 &= \frac{1}{Z} \sum_{x_2} \psi_{12}(x_1, x_2) \sum_{x_3} \psi_{13}(x_1, x_3) \sum_{x_4} \psi_{24}(x_2, x_4) \psi_{234}(x_2, x_3, x_4) m_5(x_4) \\
 &= \frac{1}{Z} \sum_{x_2} \psi_{12}(x_1, x_2) \sum_{x_3} \psi_{13}(x_1, x_3) m_4(x_2, x_3) \\
 &= \frac{1}{Z} \sum_{x_2} \psi_{12}(x_1, x_2) m_3(x_1, x_2) \\
 &= \frac{1}{Z} m_2(x_1)
 \end{aligned}$$

3^2 additions

3^3 additions

3^3 additions

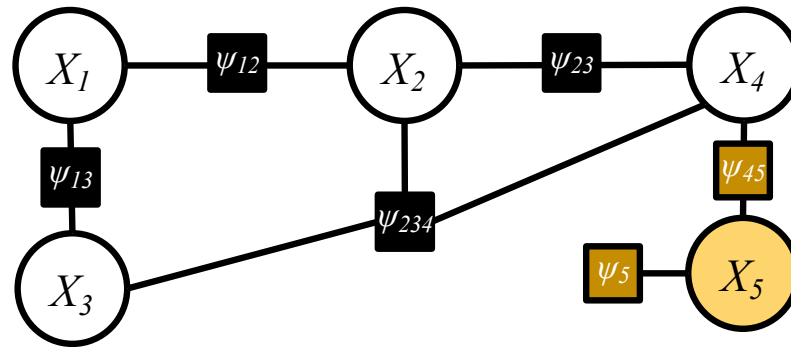
3^2 additions

3 different values on LHS

For directed graphs, $Z = 1$.

For undirected graphs, if we compute each (unnormalized) value on the LHS, we can sum them to get Z .

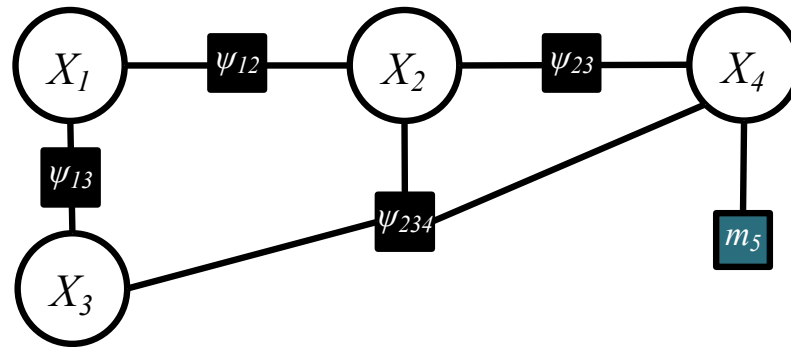
The Variable Elimination Algorithm



$$\begin{aligned}
 Z &= \sum_{x_1} \sum_{x_2} \psi_{12}(x_1, x_2) \sum_{x_3} \psi_{13}(x_1, x_3) \sum_{x_4} \psi_{24}(x_2, x_4) \psi_{234}(x_2, x_3, x_4) \sum_{x_5} \psi_{45}(x_4, x_5) \psi_5(x_5) \\
 &= \sum_{x_1} \sum_{x_2} \psi_{12}(x_1, x_2) \sum_{x_3} \psi_{13}(x_1, x_3) \sum_{x_4} \psi_{24}(x_2, x_4) \psi_{234}(x_2, x_3, x_4) m_5(x_4) \\
 &= \sum_{x_1} \sum_{x_2} \psi_{12}(x_1, x_2) \sum_{x_3} \psi_{13}(x_1, x_3) m_4(x_2, x_3)
 \end{aligned}$$

In a factor graph, variable **elimination** corresponds to replacement of a subgraph with a factor.

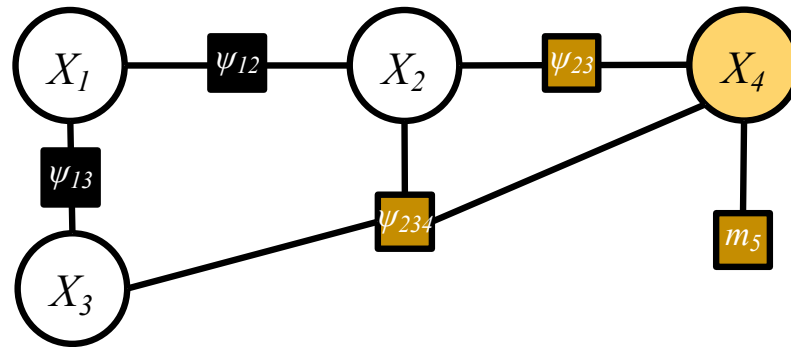
The Variable Elimination Algorithm



$$\begin{aligned}
 Z &= \sum_{x_1} \sum_{x_2} \psi_{12}(x_1, x_2) \sum_{x_3} \psi_{13}(x_1, x_3) \sum_{x_4} \psi_{24}(x_2, x_4) \psi_{234}(x_2, x_3, x_4) \sum_{x_5} \psi_{45}(x_4, x_5) \psi_5(x_5) \\
 &= \sum_{x_1} \sum_{x_2} \psi_{12}(x_1, x_2) \sum_{x_3} \psi_{13}(x_1, x_3) \sum_{x_4} \psi_{24}(x_2, x_4) \psi_{234}(x_2, x_3, x_4) m_5(x_4) \\
 &= \sum_{x_1} \sum_{x_2} \psi_{12}(x_1, x_2) \sum_{x_3} \psi_{13}(x_1, x_3) m_4(x_2, x_3)
 \end{aligned}$$

In a factor graph, variable **elimination** corresponds to replacement of a subgraph with a factor.

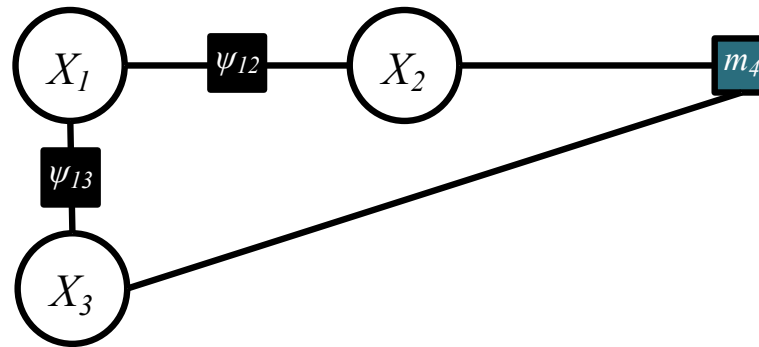
The Variable Elimination Algorithm



$$\begin{aligned}
 Z &= \sum_{x_1} \sum_{x_2} \psi_{12}(x_1, x_2) \sum_{x_3} \psi_{13}(x_1, x_3) \sum_{x_4} \psi_{24}(x_2, x_4) \psi_{234}(x_2, x_3, x_4) \sum_{x_5} \psi_{45}(x_4, x_5) \psi_5(x_5) \\
 &= \sum_{x_1} \sum_{x_2} \psi_{12}(x_1, x_2) \sum_{x_3} \psi_{13}(x_1, x_3) \sum_{x_4} \psi_{24}(x_2, x_4) \psi_{234}(x_2, x_3, x_4) m_5(x_4) \\
 &= \sum_{x_1} \sum_{x_2} \psi_{12}(x_1, x_2) \sum_{x_3} \psi_{13}(x_1, x_3) m_4(x_2, x_3)
 \end{aligned}$$

In a factor graph, variable **elimination** corresponds to replacement of a subgraph with a factor.

The Variable Elimination Algorithm



$$\begin{aligned}
 Z &= \sum_{x_1} \sum_{x_2} \psi_{12}(x_1, x_2) \sum_{x_3} \psi_{13}(x_1, x_3) \sum_{x_4} \psi_{24}(x_2, x_4) \psi_{234}(x_2, x_3, x_4) \sum_{x_5} \psi_{45}(x_4, x_5) \psi_5(x_5) \\
 &= \sum_{x_1} \sum_{x_2} \psi_{12}(x_1, x_2) \sum_{x_3} \psi_{13}(x_1, x_3) \sum_{x_4} \psi_{24}(x_2, x_4) \psi_{234}(x_2, x_3, x_4) m_5(x_4) \\
 &= \sum_{x_1} \sum_{x_2} \psi_{12}(x_1, x_2) \sum_{x_3} \psi_{13}(x_1, x_3) m_4(x_2, x_3)
 \end{aligned}$$

In a factor graph, variable **elimination** corresponds to replacement of a subgraph with a factor.

Variable Elimination for Marginal Inference

Algorithm 1a: Variable Elimination for Marginal Inference

Input: the factor graph and the query variable

Output: the marginal distribution for the query variable

- a. Run a breadth-first-search starting at the query variable to obtain an ordering of the variable nodes
- b. Reverse that ordering
- c. Eliminate each variable in the reversed ordering using Algorithm 2

Algorithm 2: Eliminate One Variable

Input: the variable to be eliminated

Output: new factor graph with the variable marginalized out

- a. Find the input variable and its neighboring factors -- call this set the eliminated set
- b. Replace the eliminated set with a new factor
 - a. The neighbors of the new factor should be all the neighbors of all the factors in the eliminated set
 - b. The new factor should assign a score to each possible assignment of its neighboring variables
 - c. Said score should be identical to the product of the factors it is replacing, summing over the eliminated variable

Variable Elimination for Marginal Inference

Algorithm 1b: Variable Elimination for the Partition Function

Input: the factor graph

Output: the partition function

- a. Run a breadth-first-search starting at an arbitrary variable to obtain an ordering of the variable nodes
- b. Eliminate each variable in the ordering using Algorithm 2

Algorithm 2: Eliminate One Variable

Input: the variable to be eliminated

Output: new factor graph with the variable marginalized out

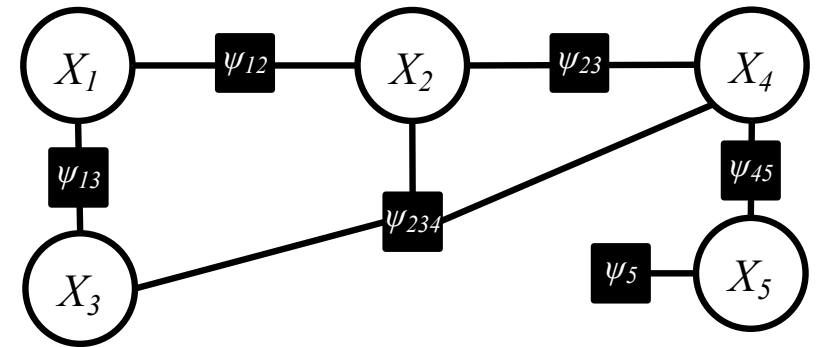
- a. Find the input variable and its neighboring factors -- call this set the eliminated set
- b. Replace the eliminated set with a new factor
 - a. The neighbors of the new factor should be all the neighbors of all the factors in the eliminated set
 - b. The new factor should assign a score to each possible assignment of its neighboring variables
 - c. Said score should be identical to the product of the factors it is replacing, summing over the eliminated variable

Variable Elimination

Whiteboard:

- Ex: Variable Elimination as factor replacement

Variable Elimination Complexity



In-Class Exercise: *Fill in the blank*

Brute force, naïve,
inference is $O(\underline{\hspace{1cm}})$

Variable elimination
is $O(\underline{\hspace{1cm}})$

where $n = \#$ of variables
 $k = \max \#$ values a variable can take
 $r = \#$ variables participating in
largest “intermediate” table

Exact Inference

Variable Elimination

- *Uses*
 - Computes the **partition function** of **any** factor graph
 - Computes the **marginal probability** of a **query variable** in **any** factor graph
- *Limitations*
 - Only computes the marginal for **one variable at a time** (i.e. need to re-run variable elimination for each variable if you need them all)
 - **Elimination order** affects runtime

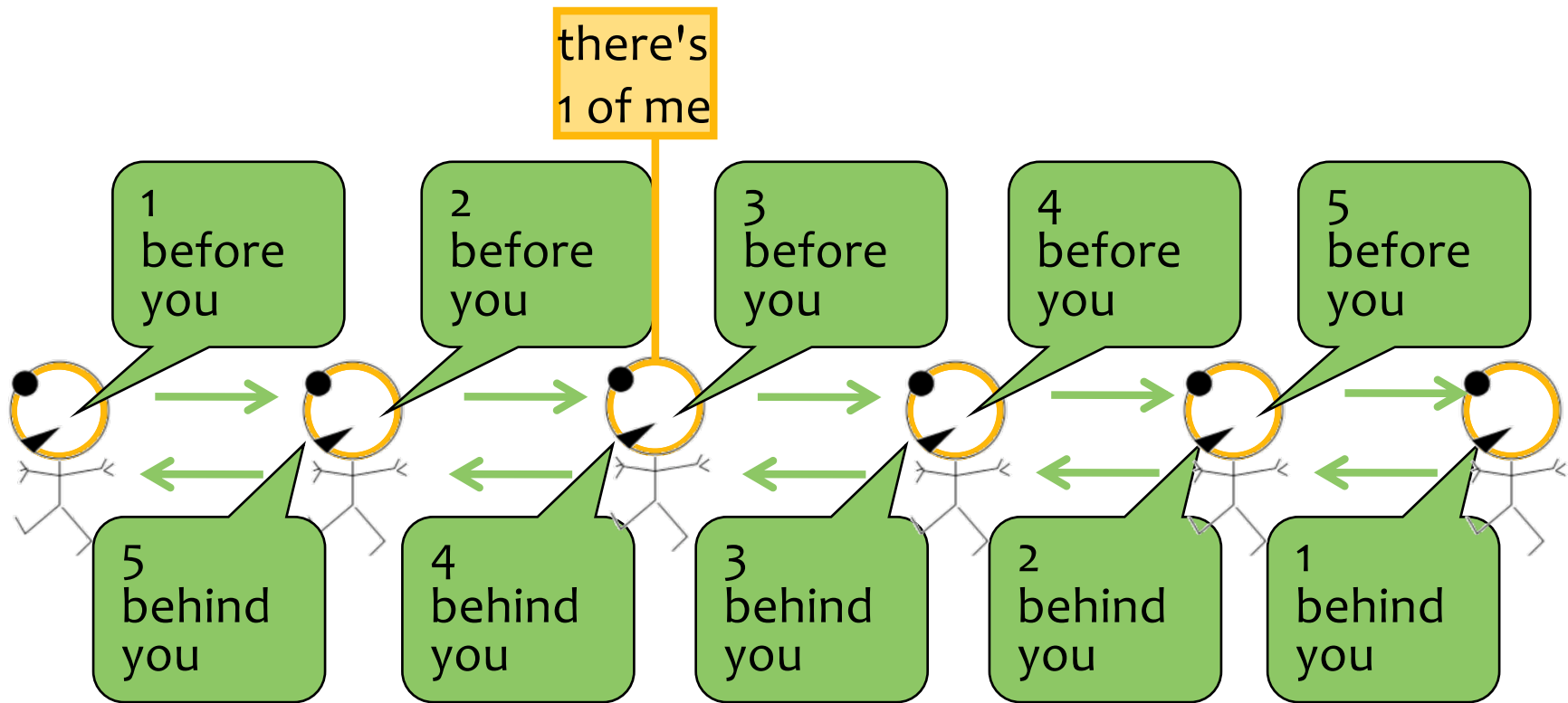
Belief Propagation

- *Uses*
 - Computes the **partition function** of **any acyclic** factor graph
 - Computes **all marginal probabilities** of factors and variables at once, for **any acyclic** factor graph
- *Limitations*
 - Only **exact** on acyclic factor graphs (though we'll consider its “loopy” variant later)
 - **Message passing order** affects runtime (but the obvious topological ordering always works best)

MESSAGE PASSING

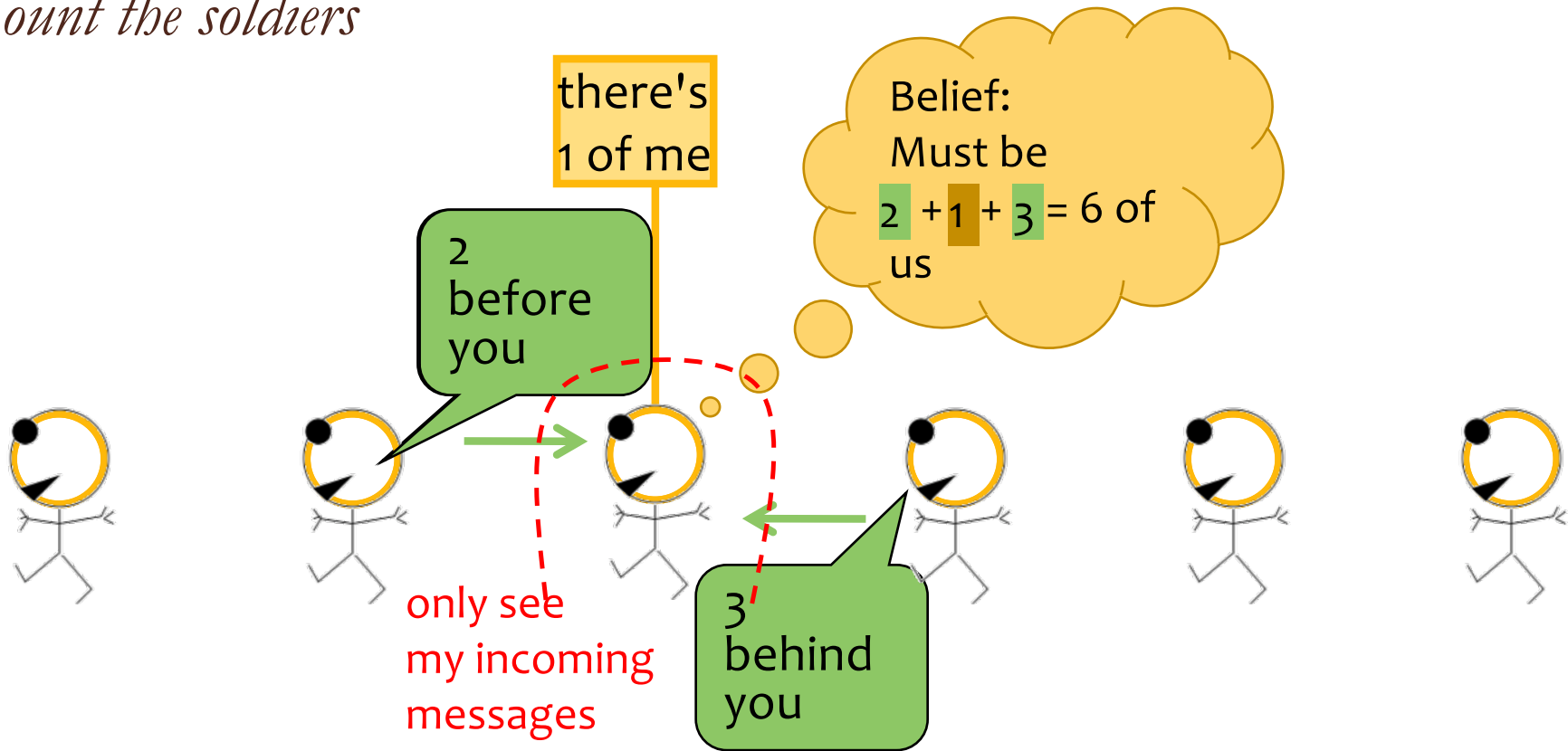
Great Ideas in ML: Message Passing

Count the soldiers



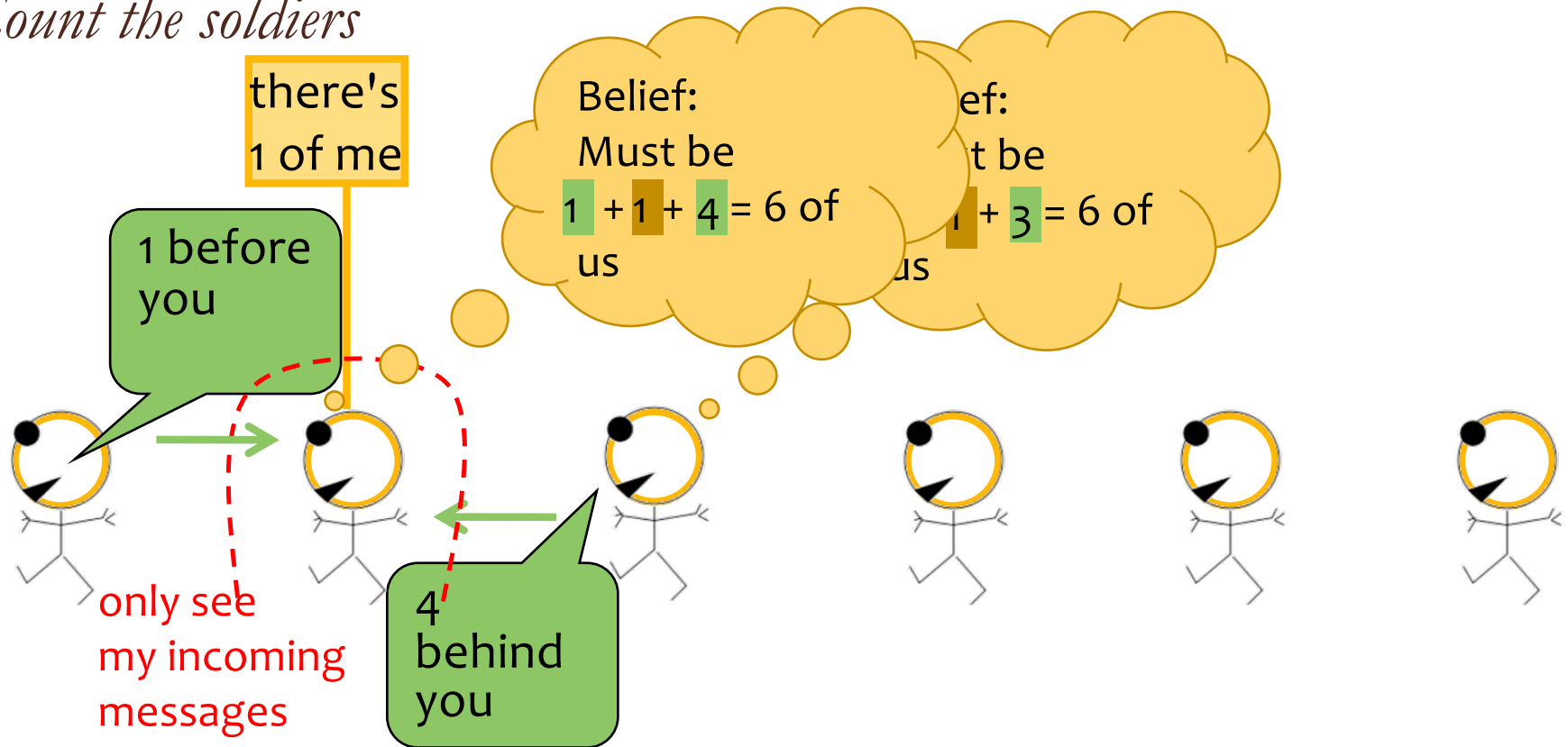
Great Ideas in ML: Message Passing

Count the soldiers



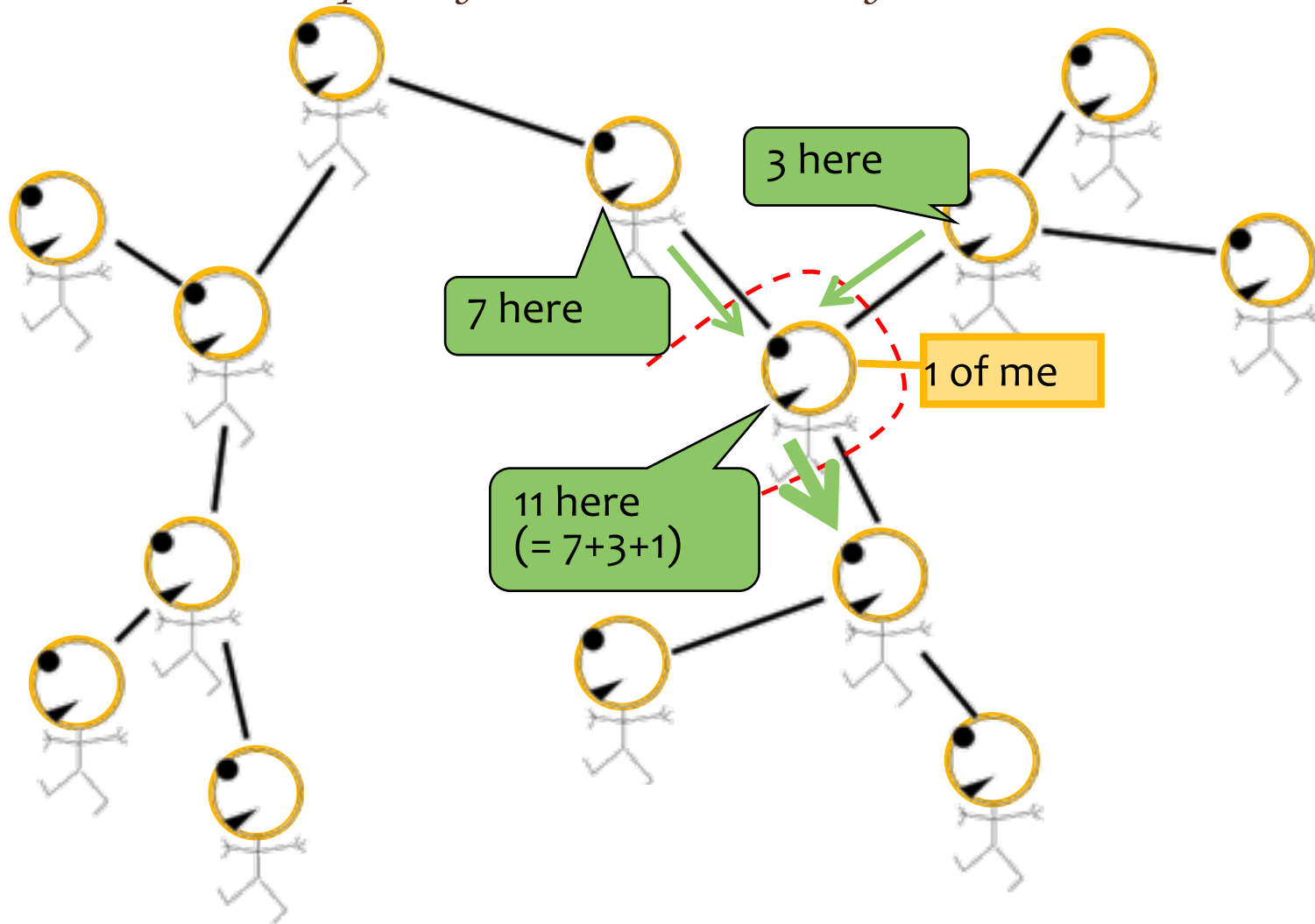
Great Ideas in ML: Message Passing

Count the soldiers



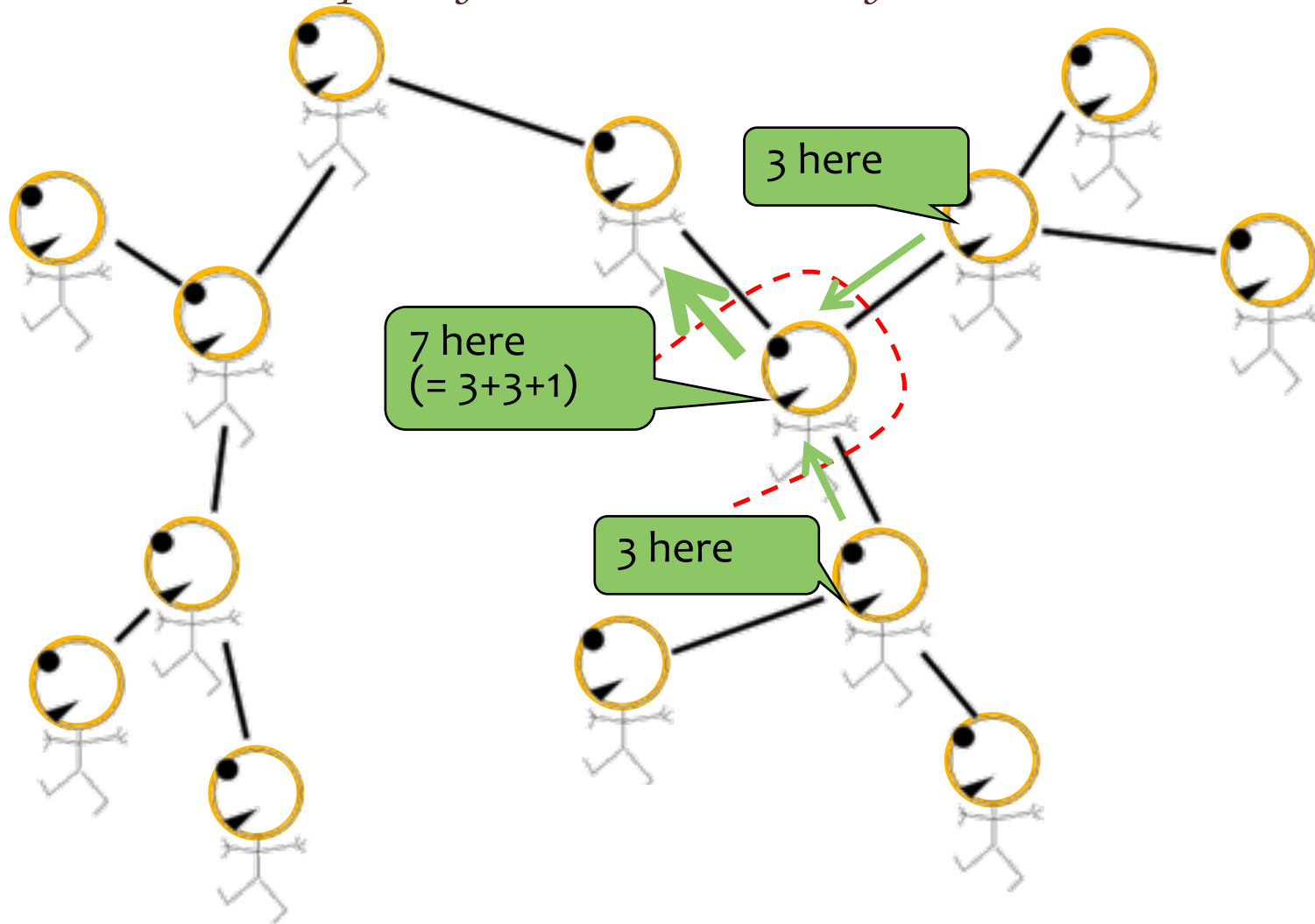
Great Ideas in ML: Message Passing

Each soldier receives reports from all branches of tree



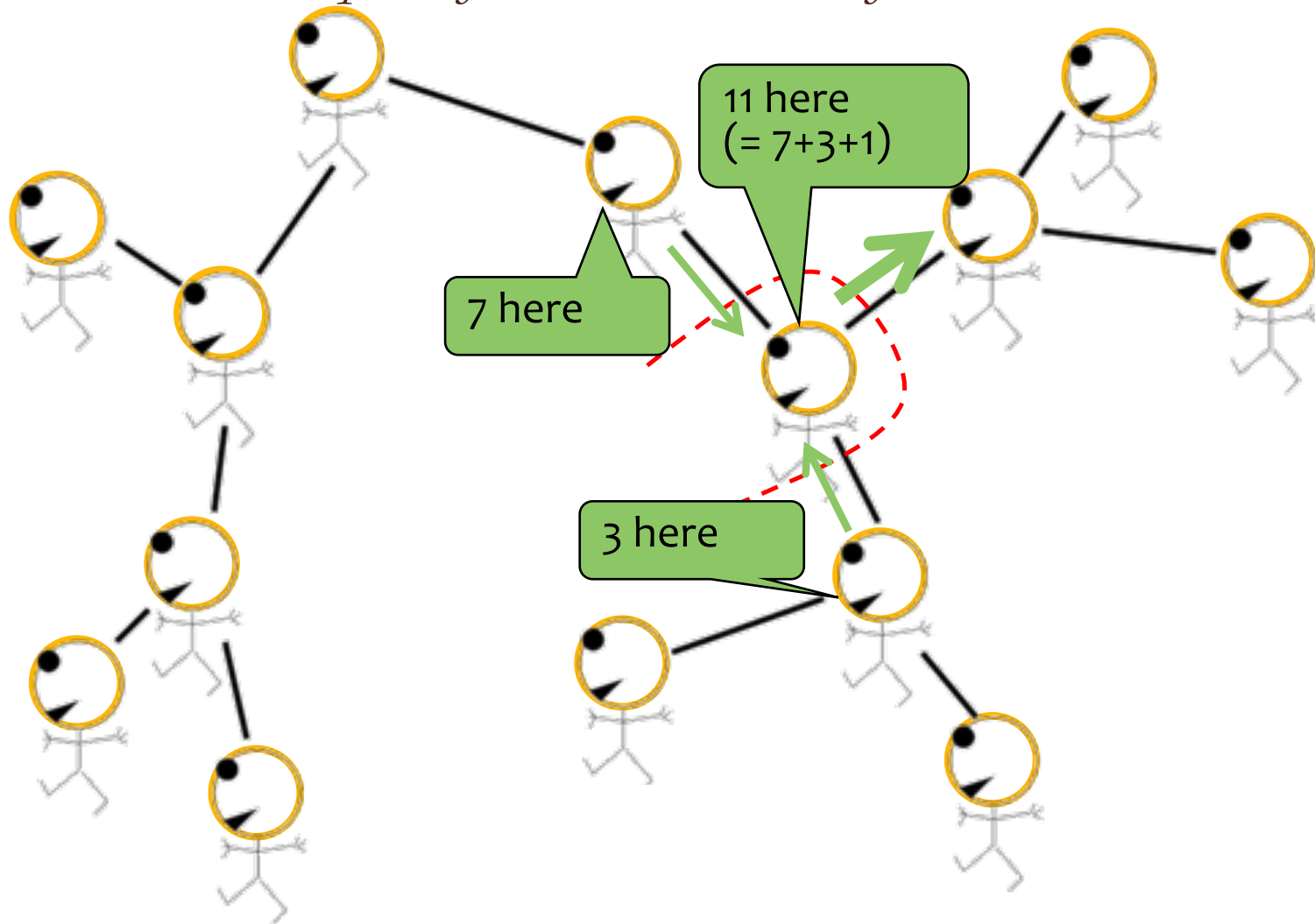
Great Ideas in ML: Message Passing

Each soldier receives reports from all branches of tree



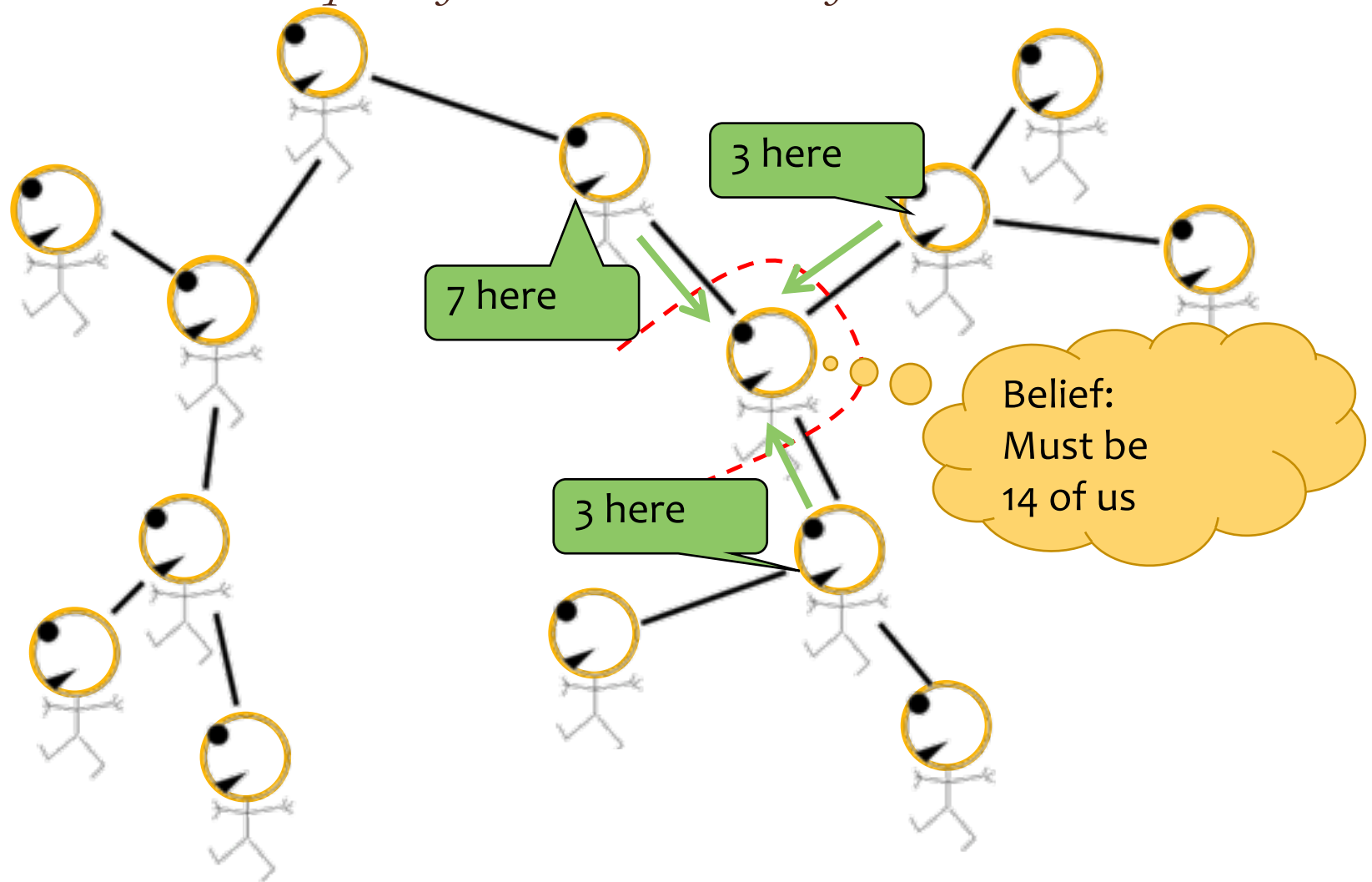
Great Ideas in ML: Message Passing

Each soldier receives reports from all branches of tree



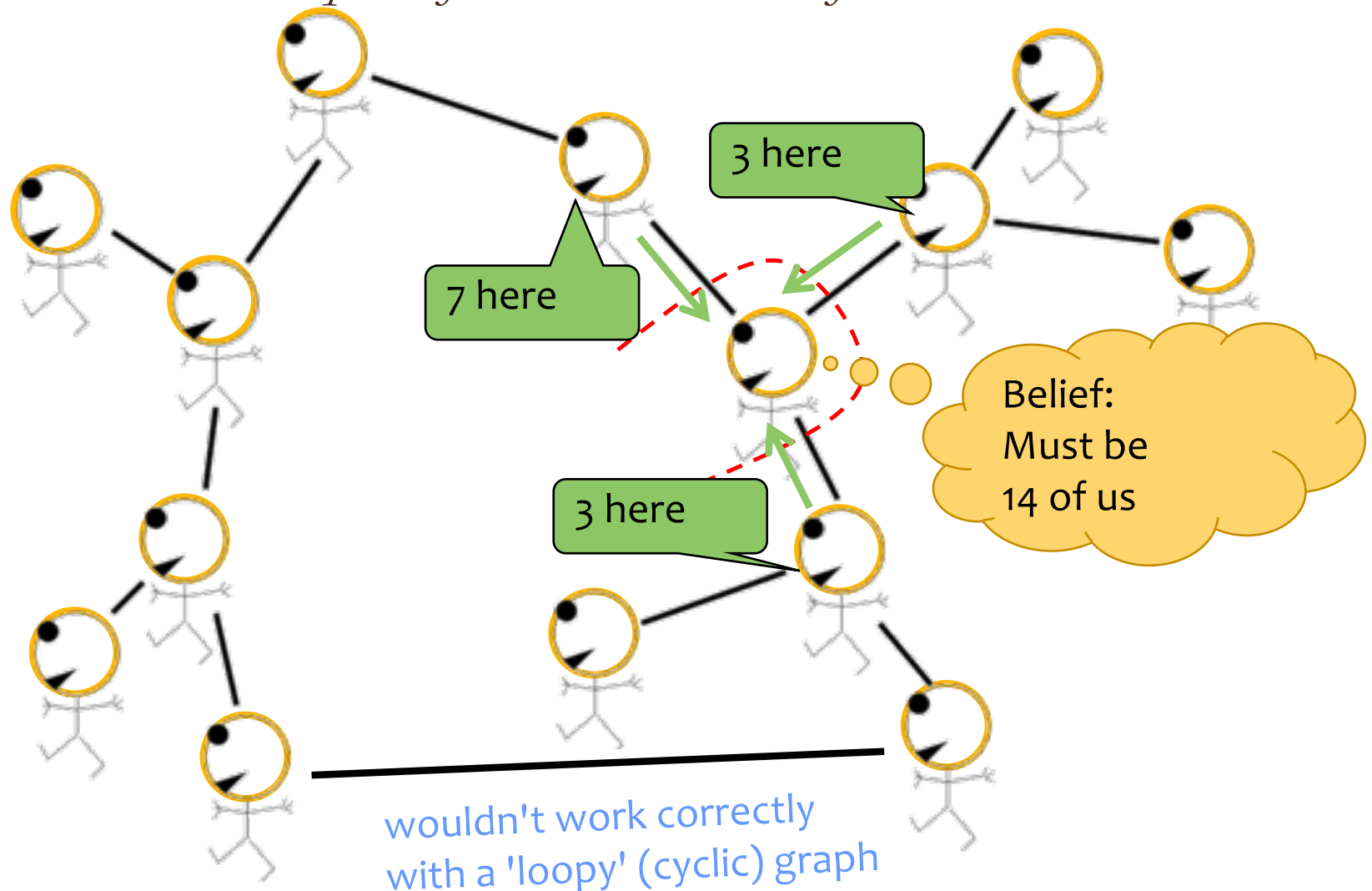
Great Ideas in ML: Message Passing

Each soldier receives reports from all branches of tree



Great Ideas in ML: Message Passing

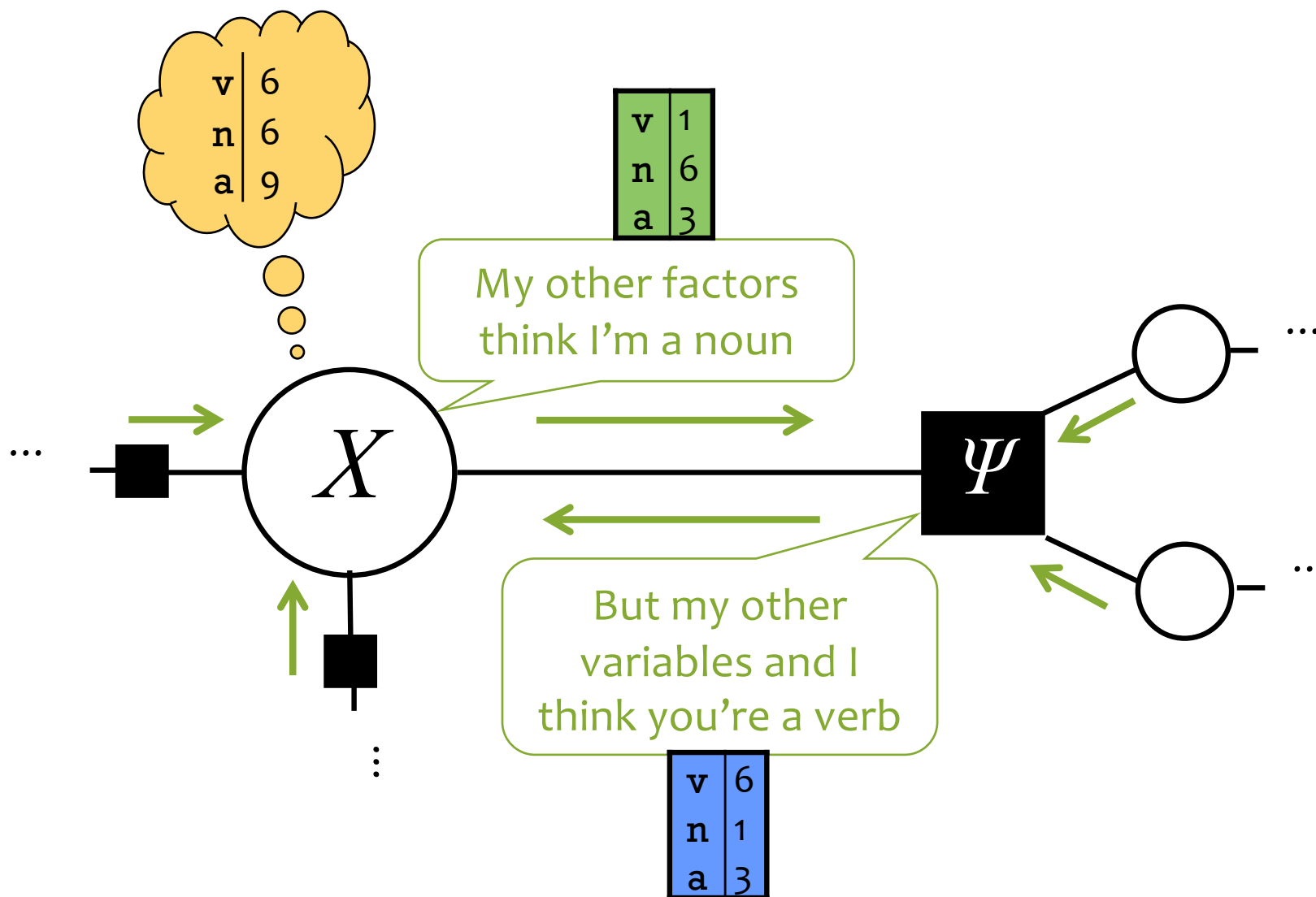
Each soldier receives reports from all branches of tree



Exact marginal inference for factor trees

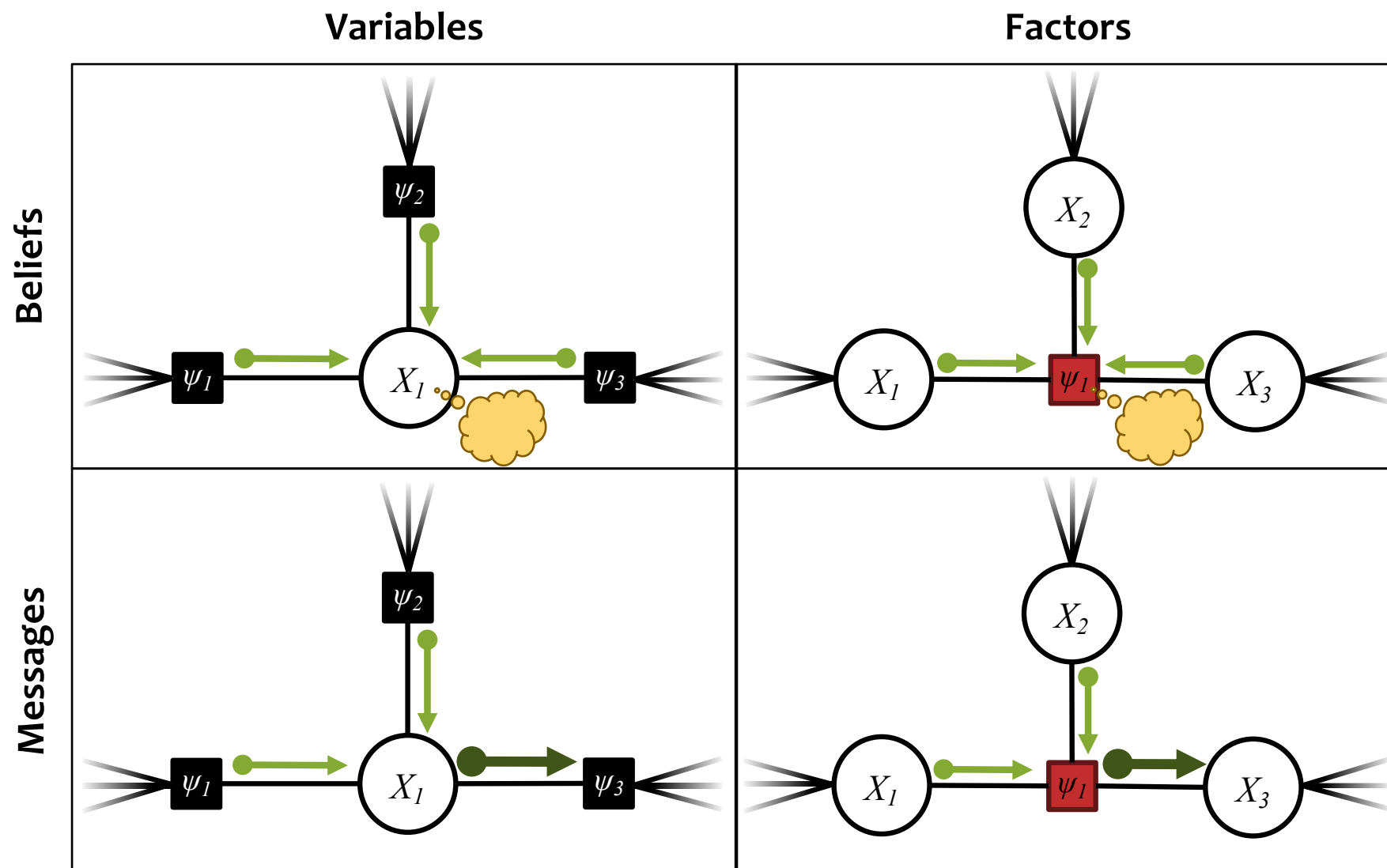
SUM-PRODUCT BELIEF PROPAGATION

Message Passing in Belief Propagation



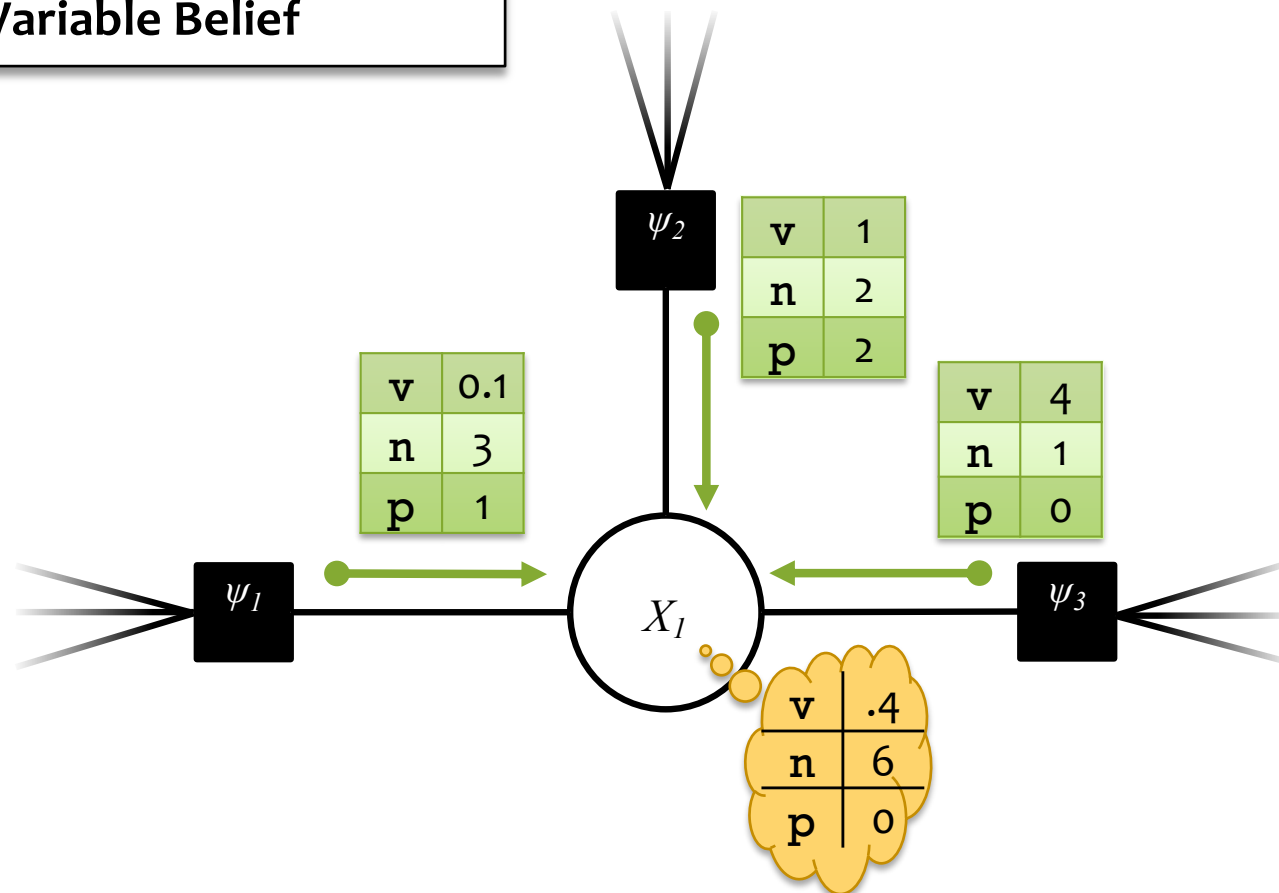
Both of these messages judge the possible values of variable X .
Their product = belief at X = product of all 3 messages to X .

Sum-Product Belief Propagation



Sum-Product Belief Propagation

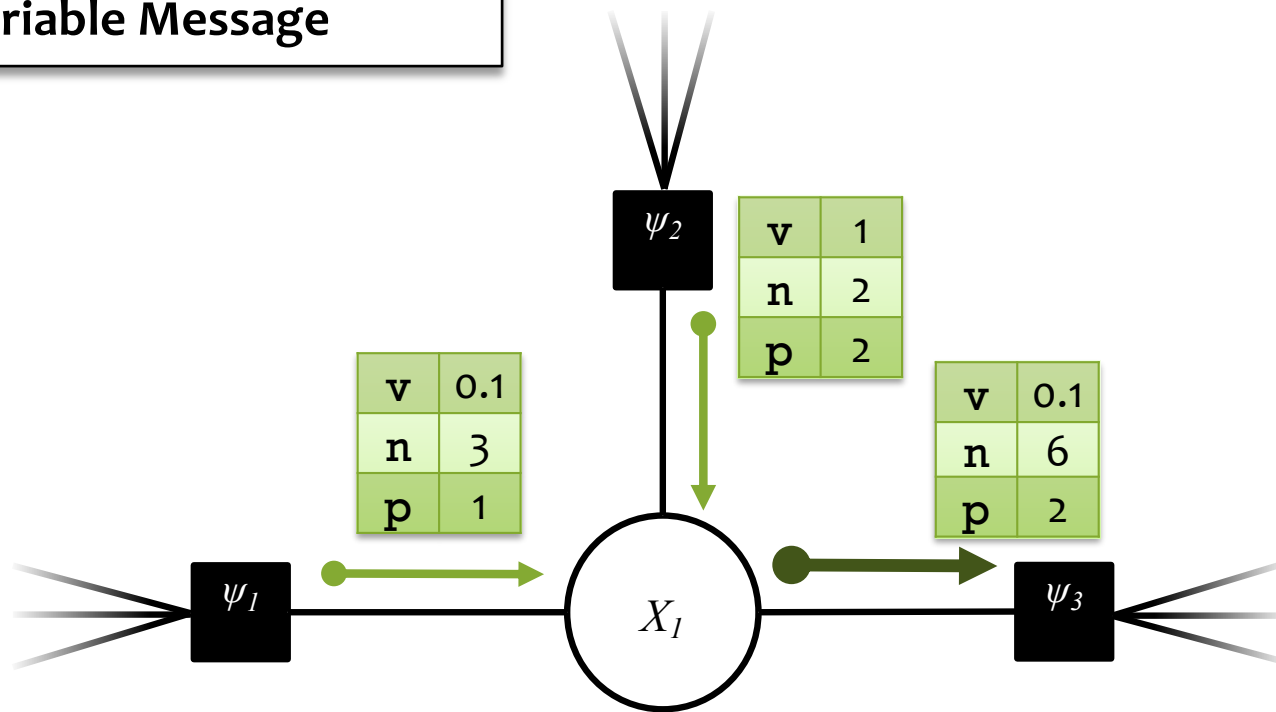
Variable Belief



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

Sum-Product Belief Propagation

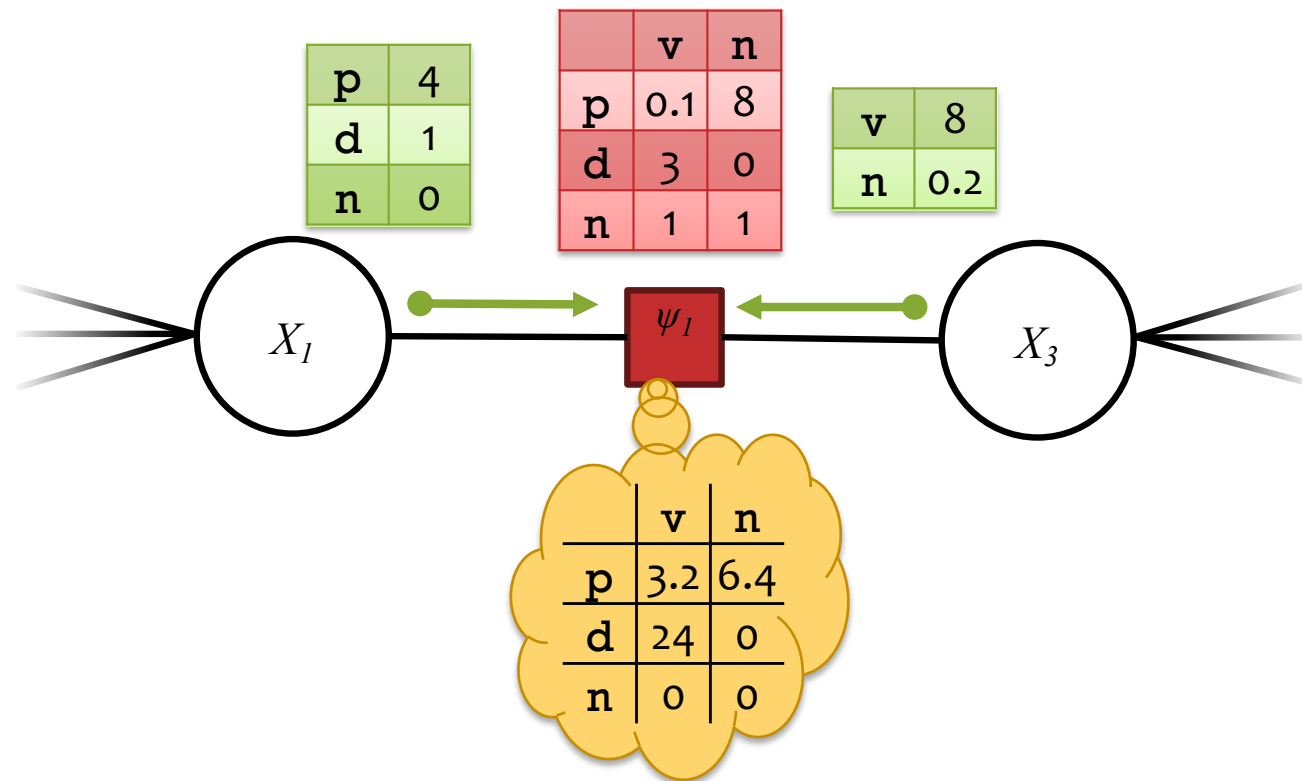
Variable Message



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

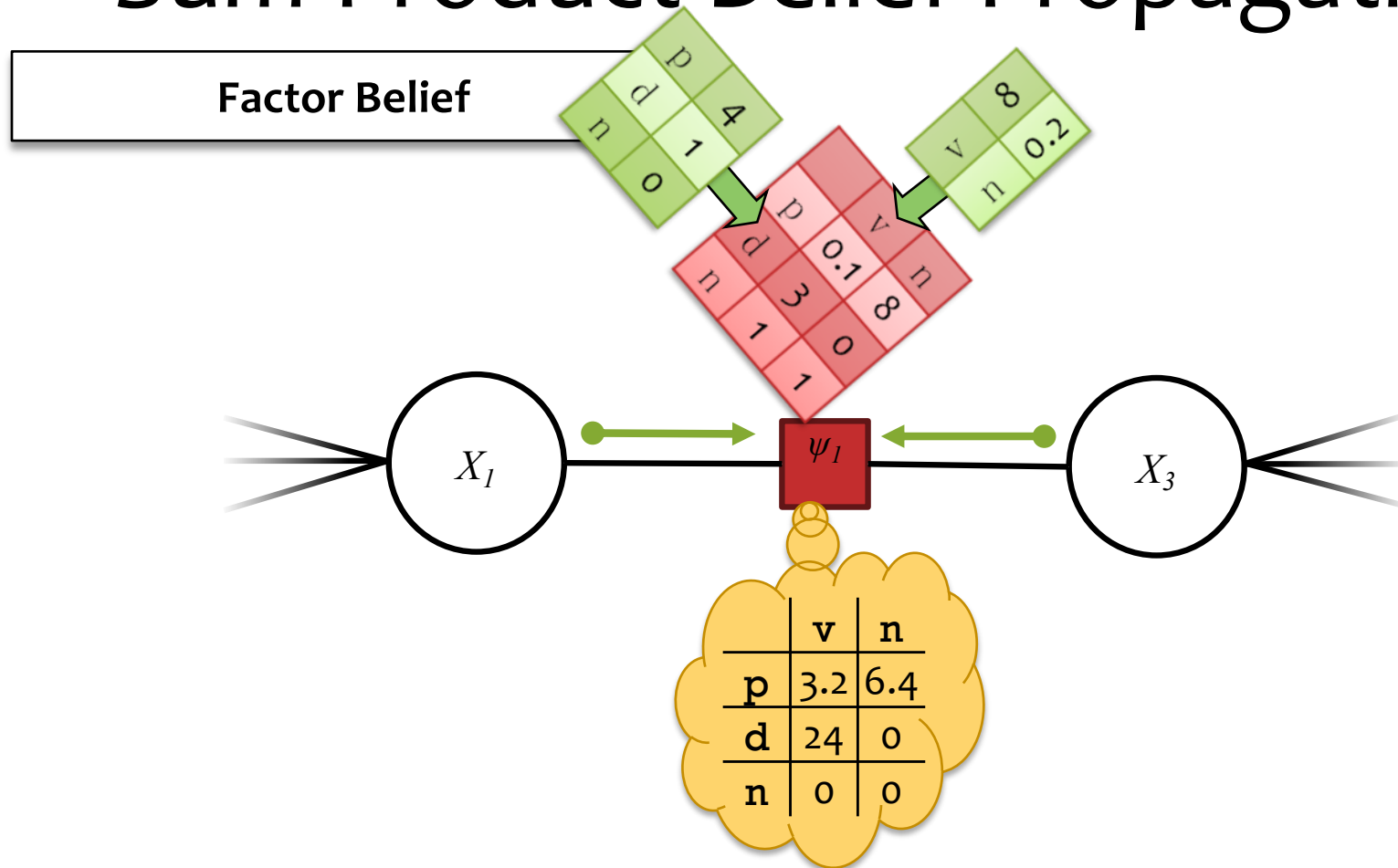
Sum-Product Belief Propagation

Factor Belief



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

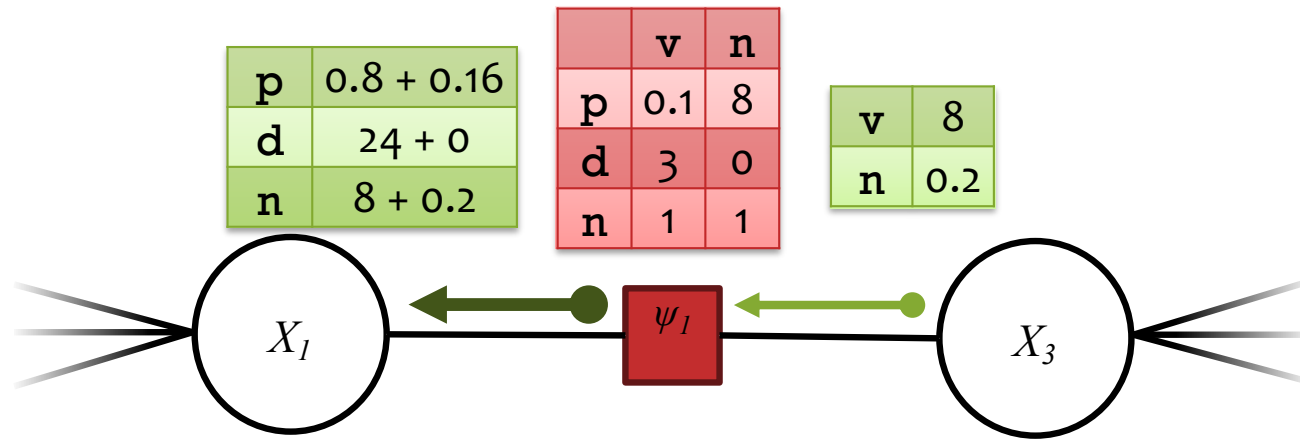
Sum-Product Belief Propagation



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

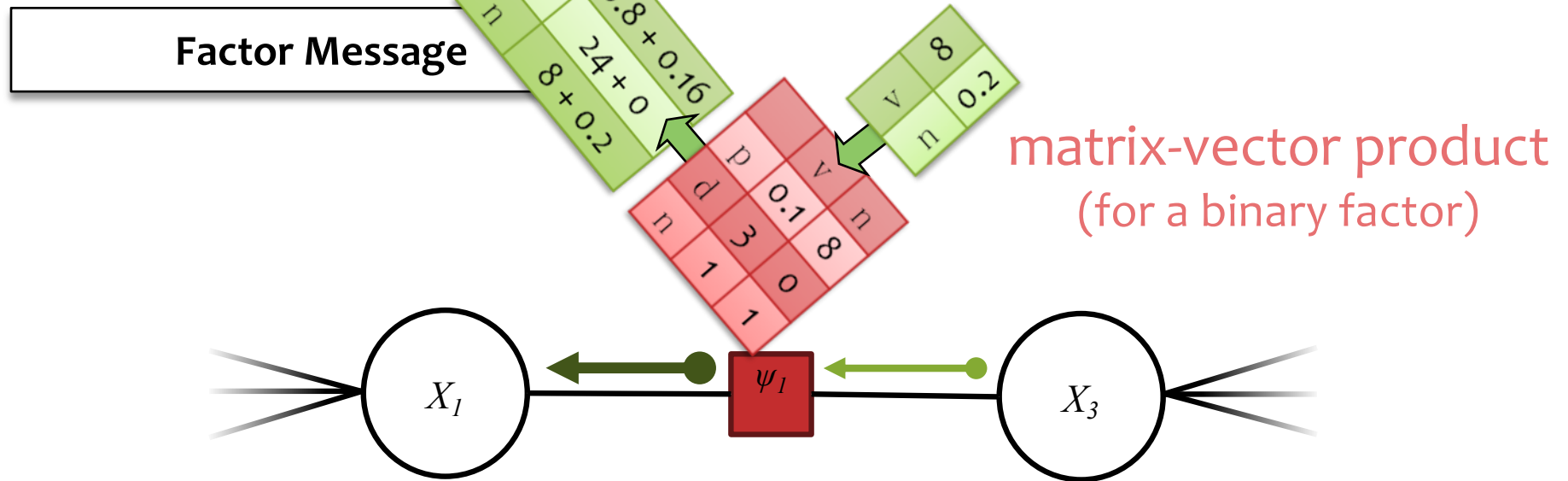
Sum-Product Belief Propagation

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**.

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

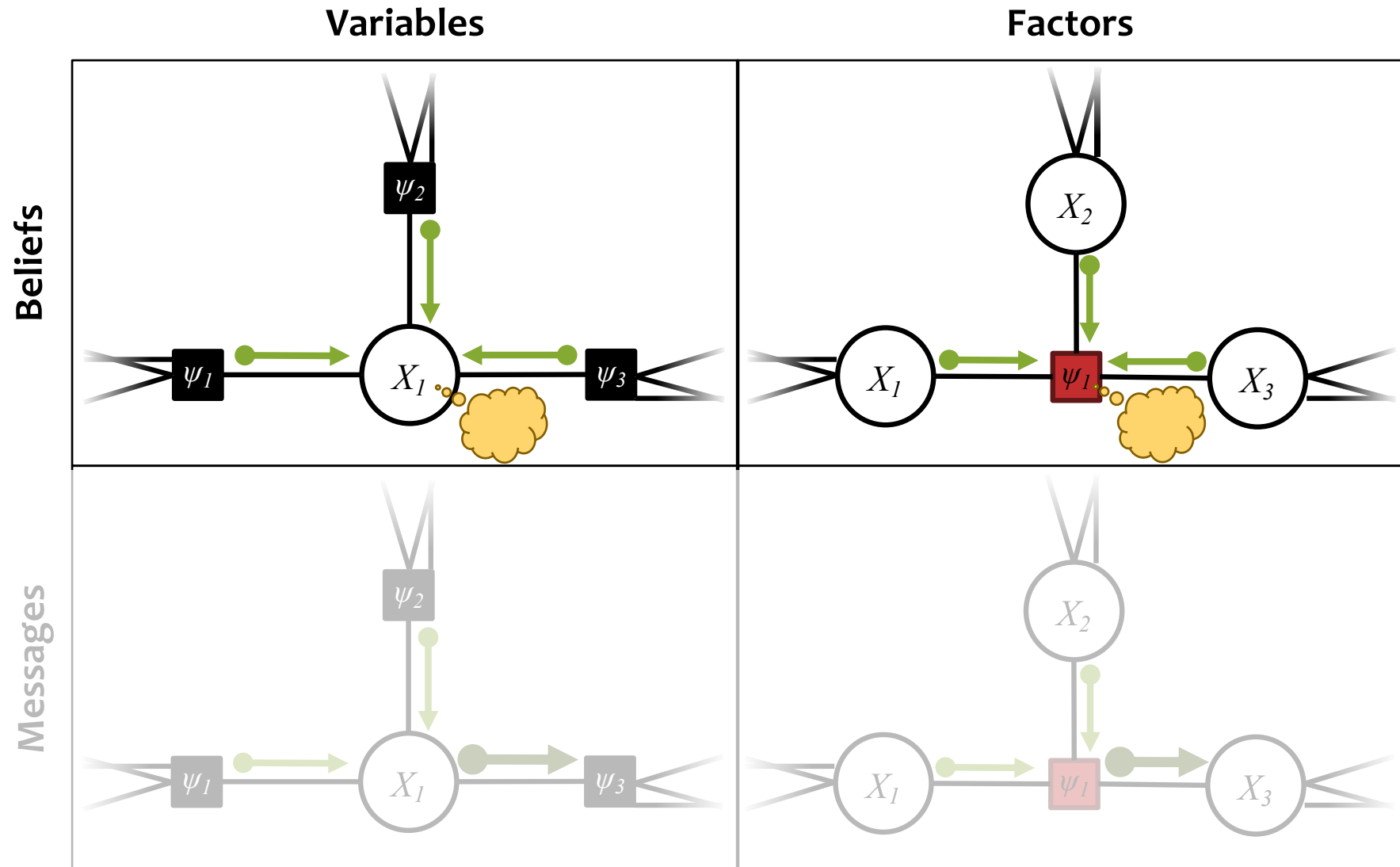
1. Compute the beliefs (unnormalized marginals).

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

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

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

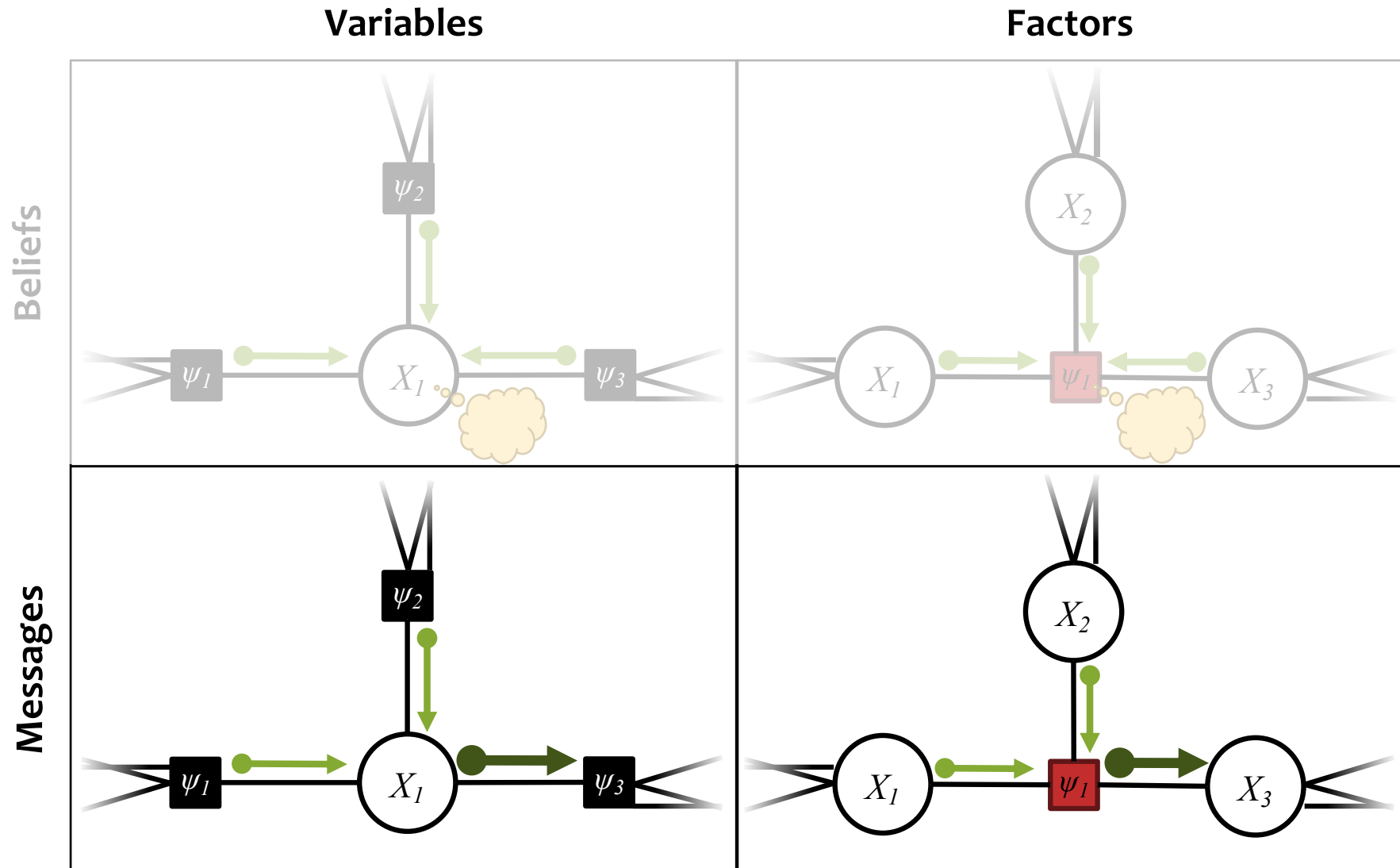
Sum-Product Belief Propagation



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

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

Sum-Product Belief Propagation

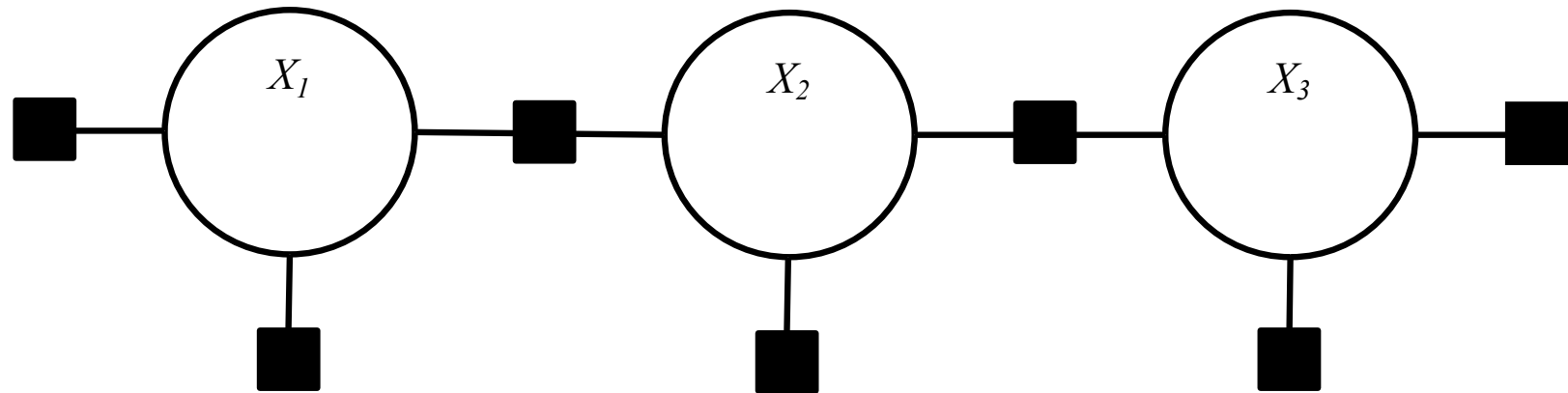


$$\mu_{i \rightarrow \alpha}(x_i) = \prod_{\alpha \in \mathcal{N}(i) \setminus \alpha} \mu_{\alpha \rightarrow i}(x_i)$$

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

**FORWARD BACKWARD AS
SUM-PRODUCT BP**

CRF Tagging Model



find

preferred

tags

Could be verb or noun

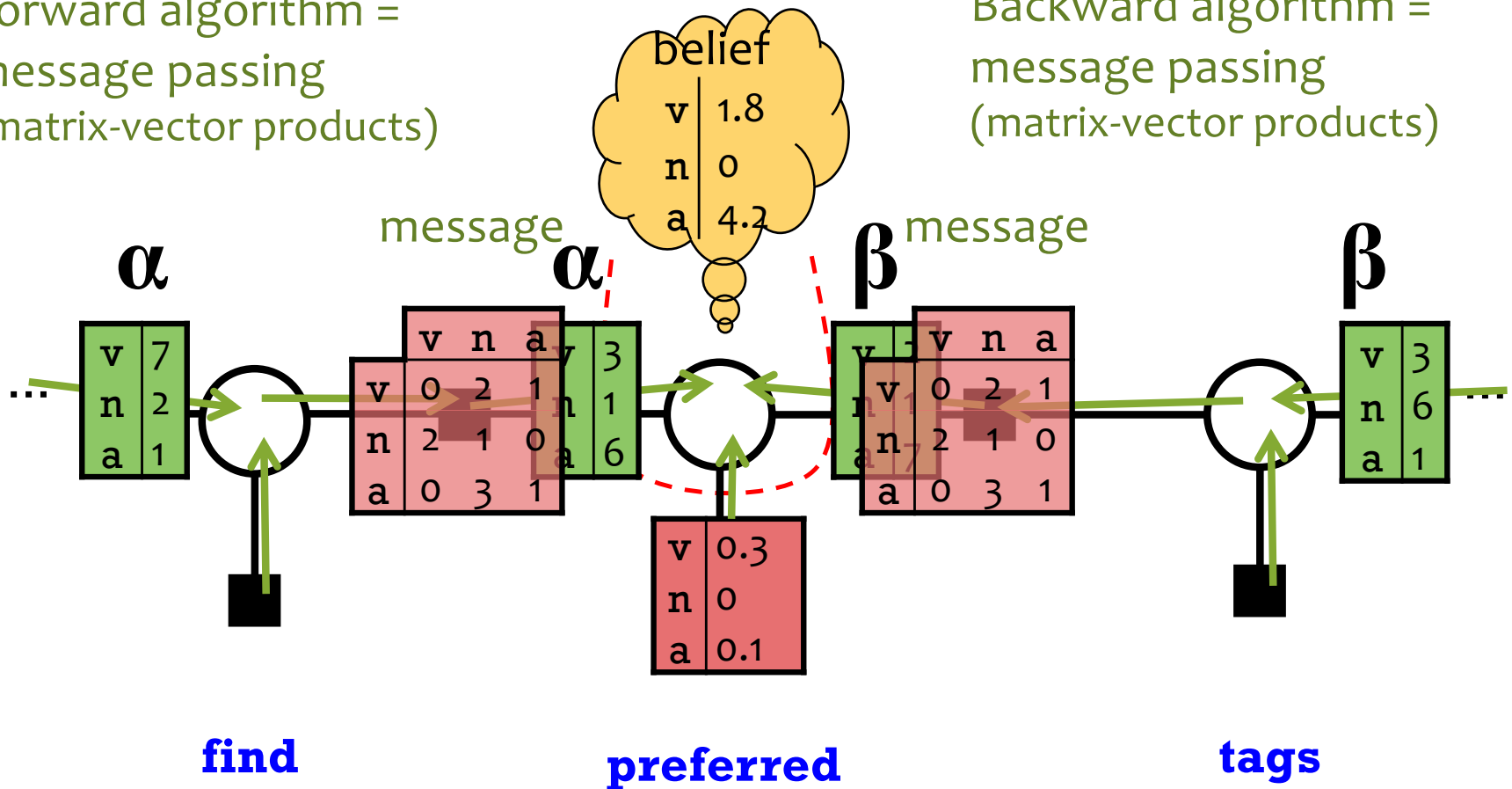
Could be adjective or verb

Could be noun or verb

CRF Tagging by Belief Propagation

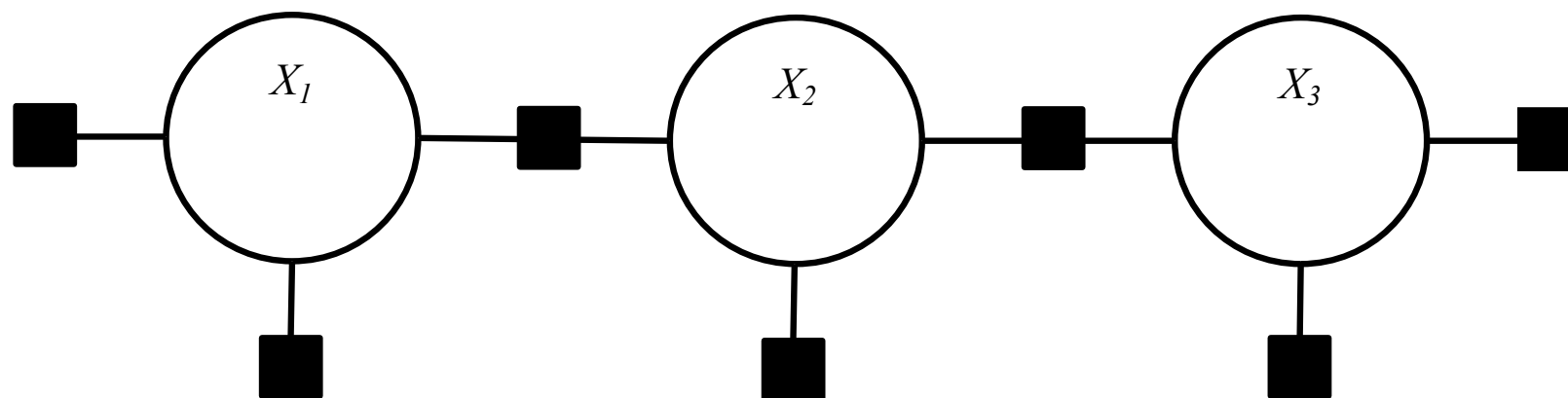
Forward algorithm =
message passing
(matrix-vector products)

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

Could be verb or noun

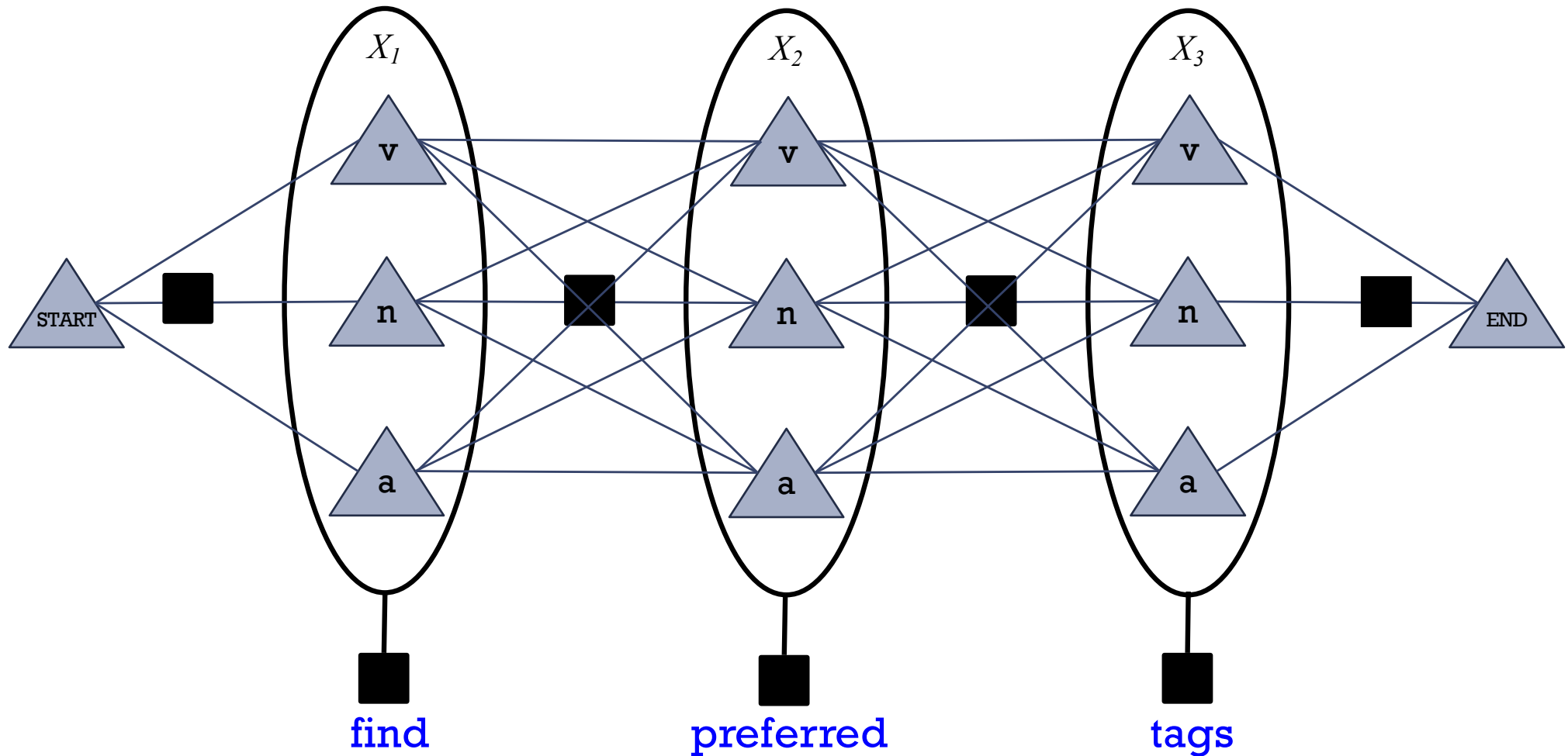
preferred

Could be adjective or verb

tags

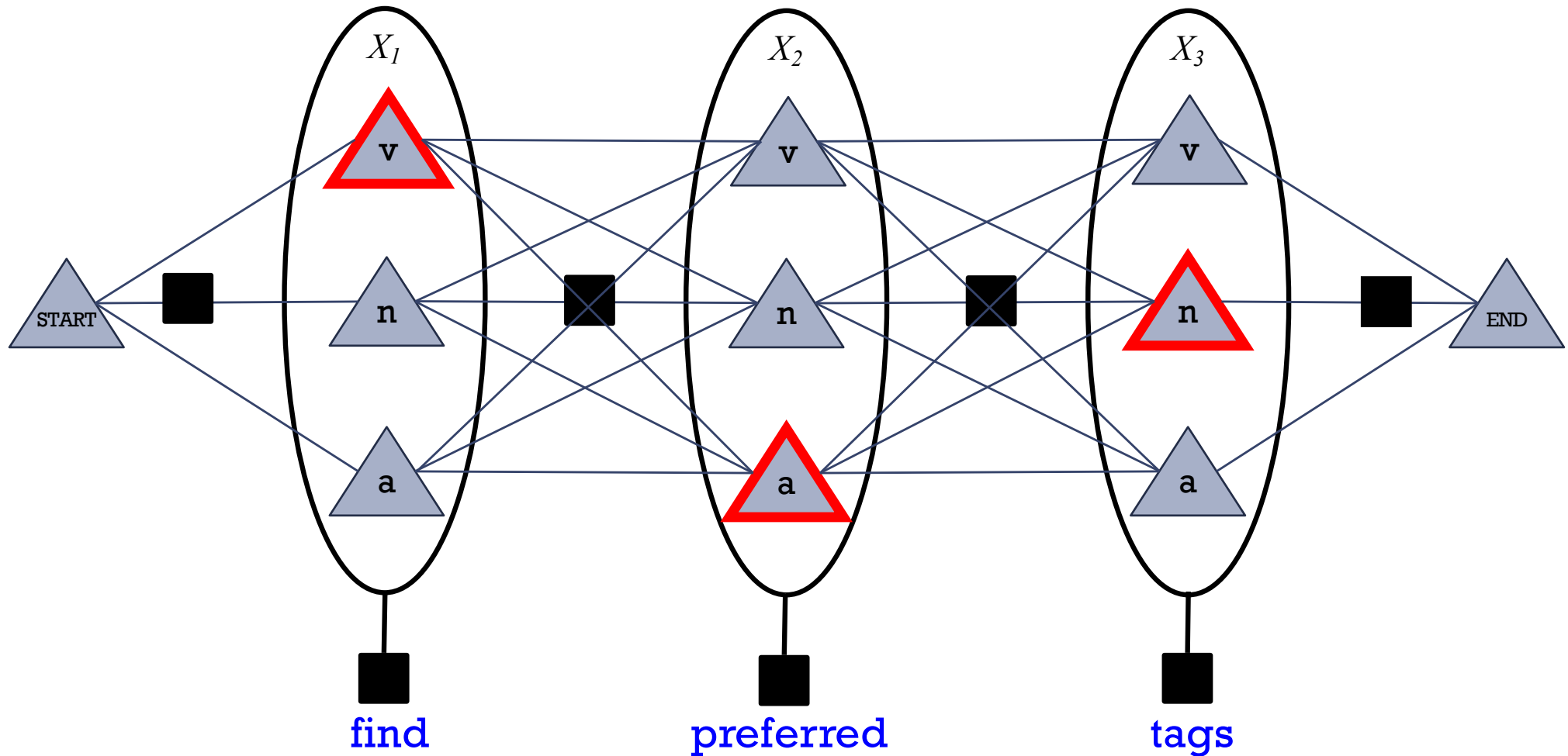
Could be noun or verb

So Let's Review Forward-Backward ...



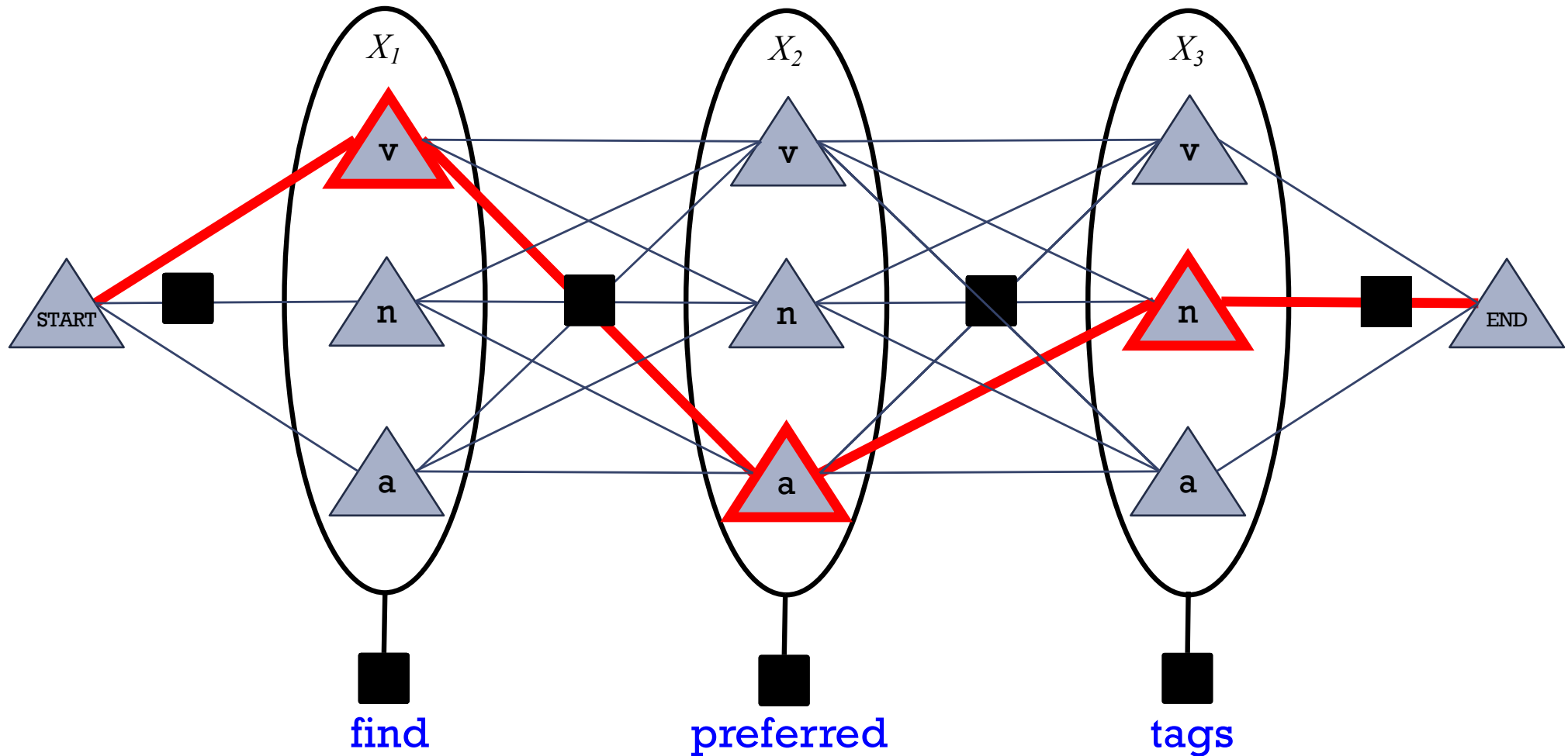
- Show the possible *values* for each variable

So Let's Review Forward-Backward ...



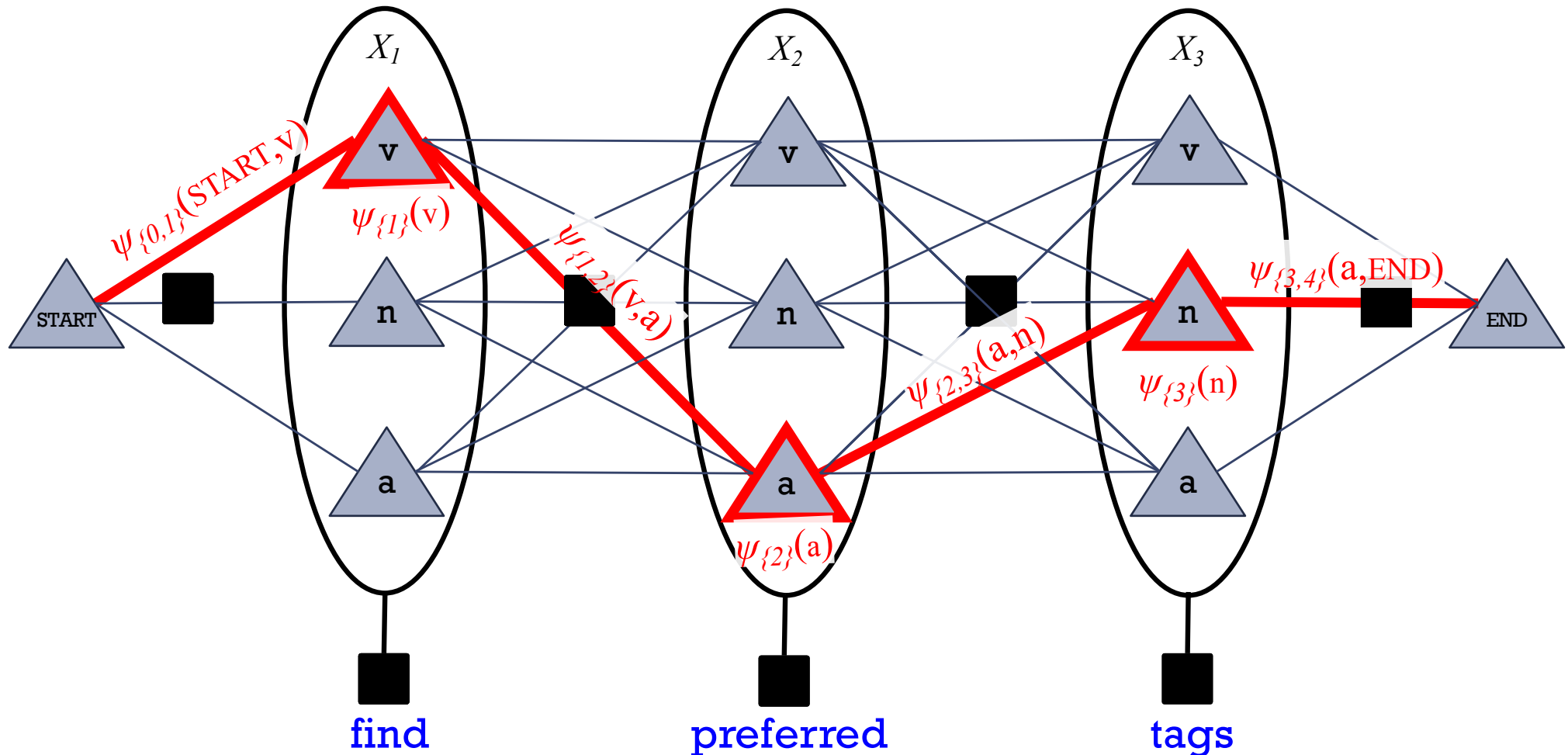
- Let's show the possible *values* for each variable
- One possible assignment

So Let's Review Forward-Backward ...



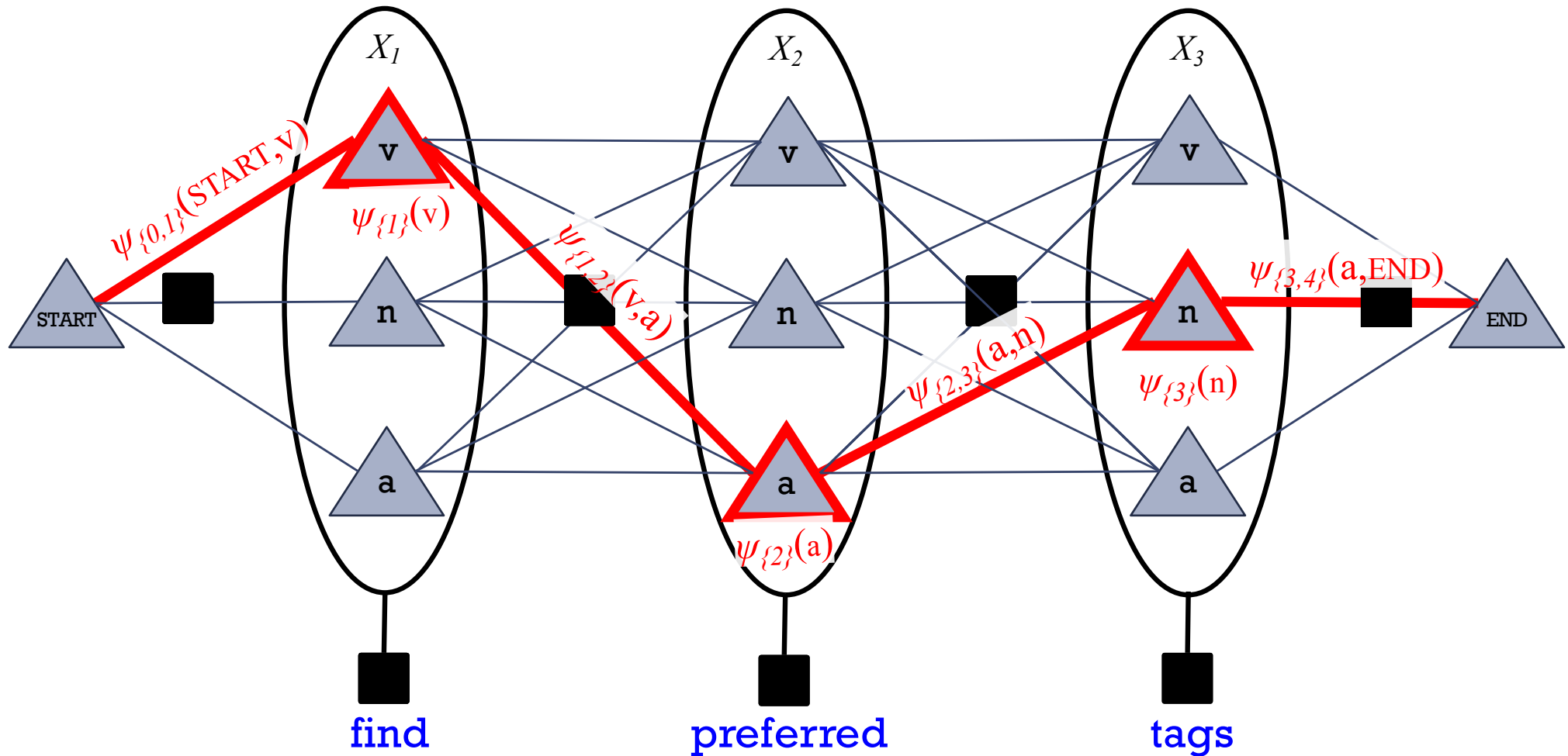
- Let's show the possible values for each variable
- One possible assignment
- And what the 7 factors **think of it** ...

Viterbi Algorithm: Most Probable Assignment



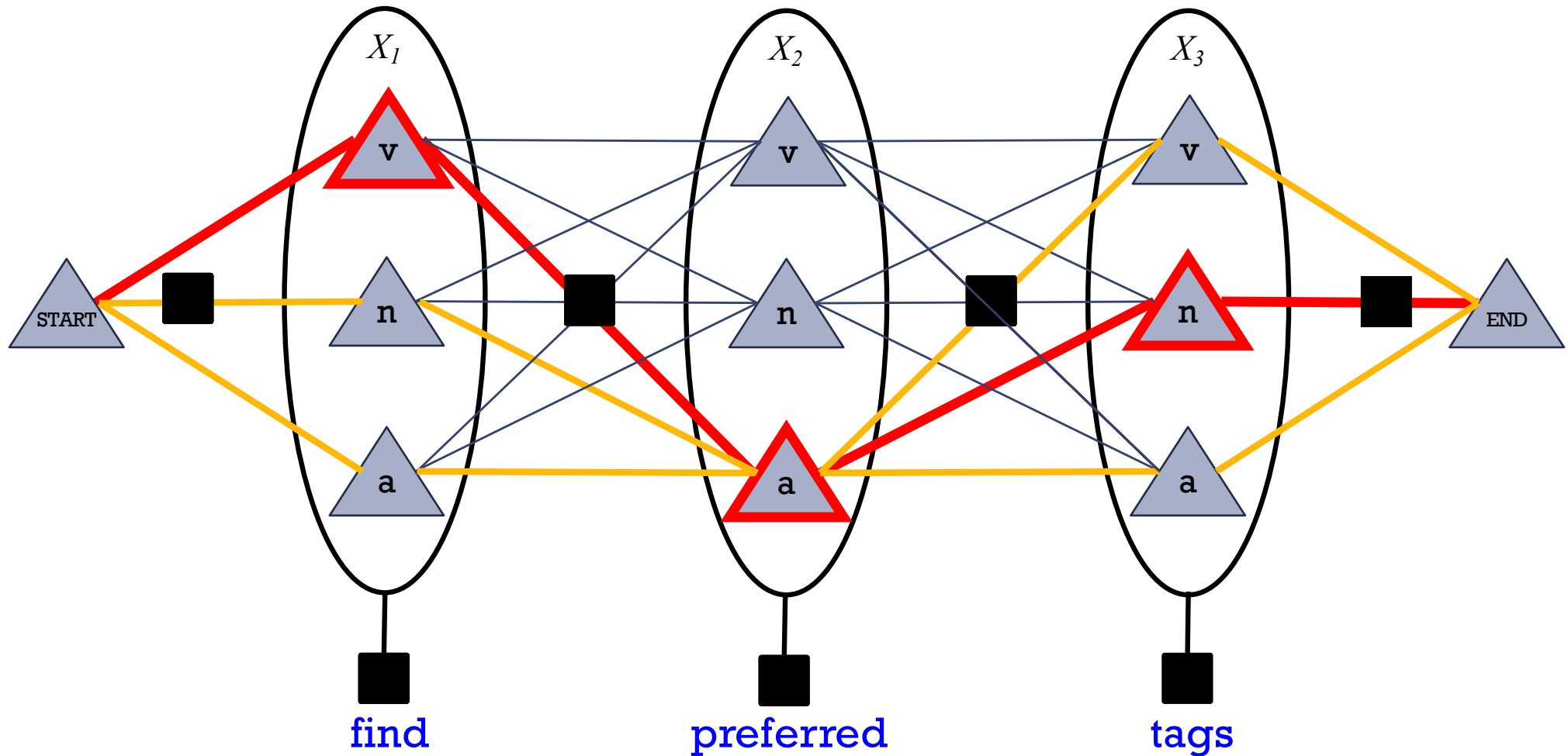
- So $p(v \ a \ n) = (1/Z) * \text{product of 7 numbers}$
- Numbers associated with edges and nodes of path
- Most probable assignment = **path with highest product**

Viterbi Algorithm: Most Probable Assignment

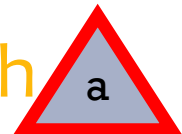


- So $p(v \ a \ n) = (1/Z) * \text{product weight of one path}$

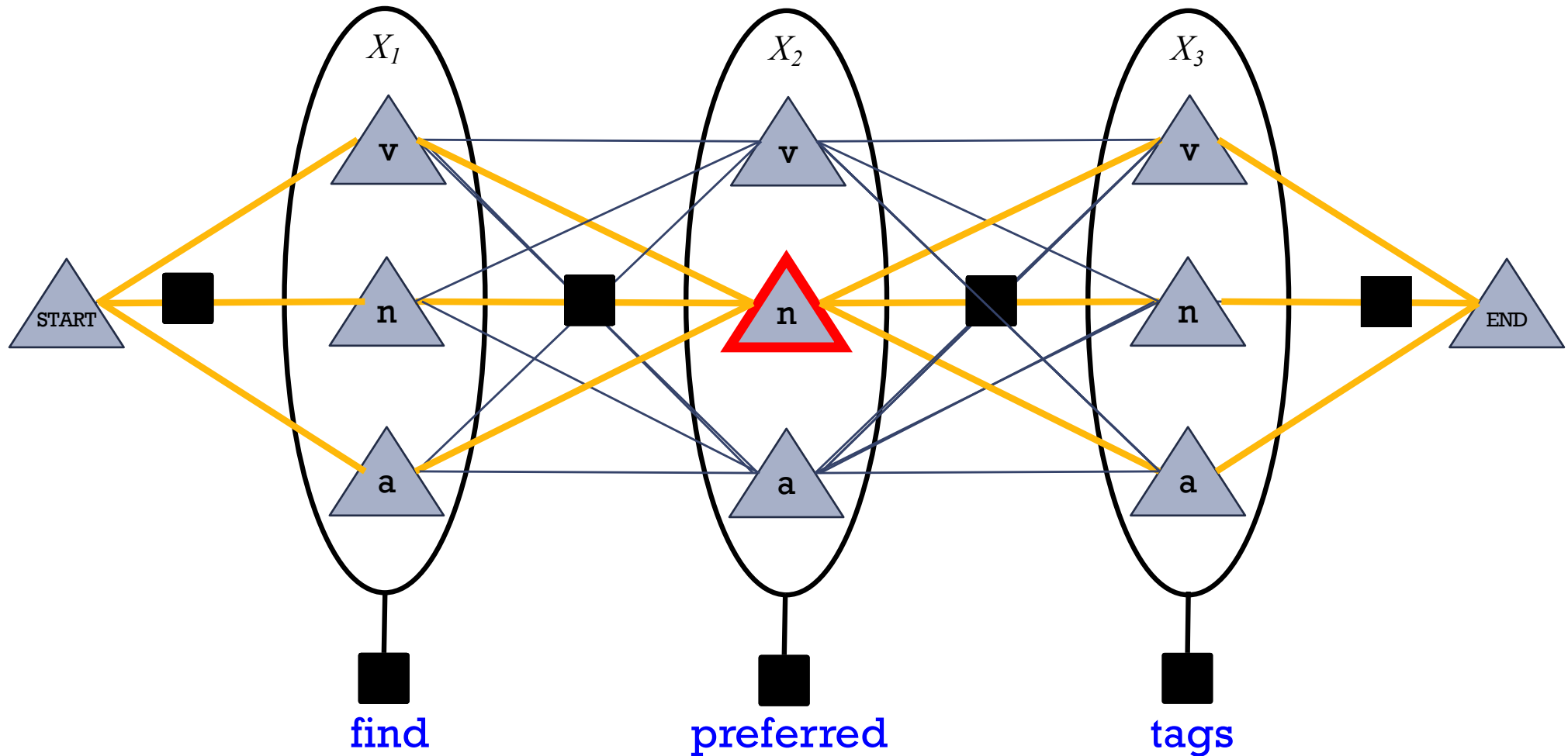
Forward-Backward Algorithm: Finds Marginals



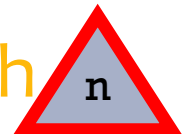
- So $p(v \ a \ n) = (1/Z) * \text{product weight of one path}$
- Marginal probability $p(X_2 = a)$
 $= (1/Z) * \text{total weight of all paths through}$



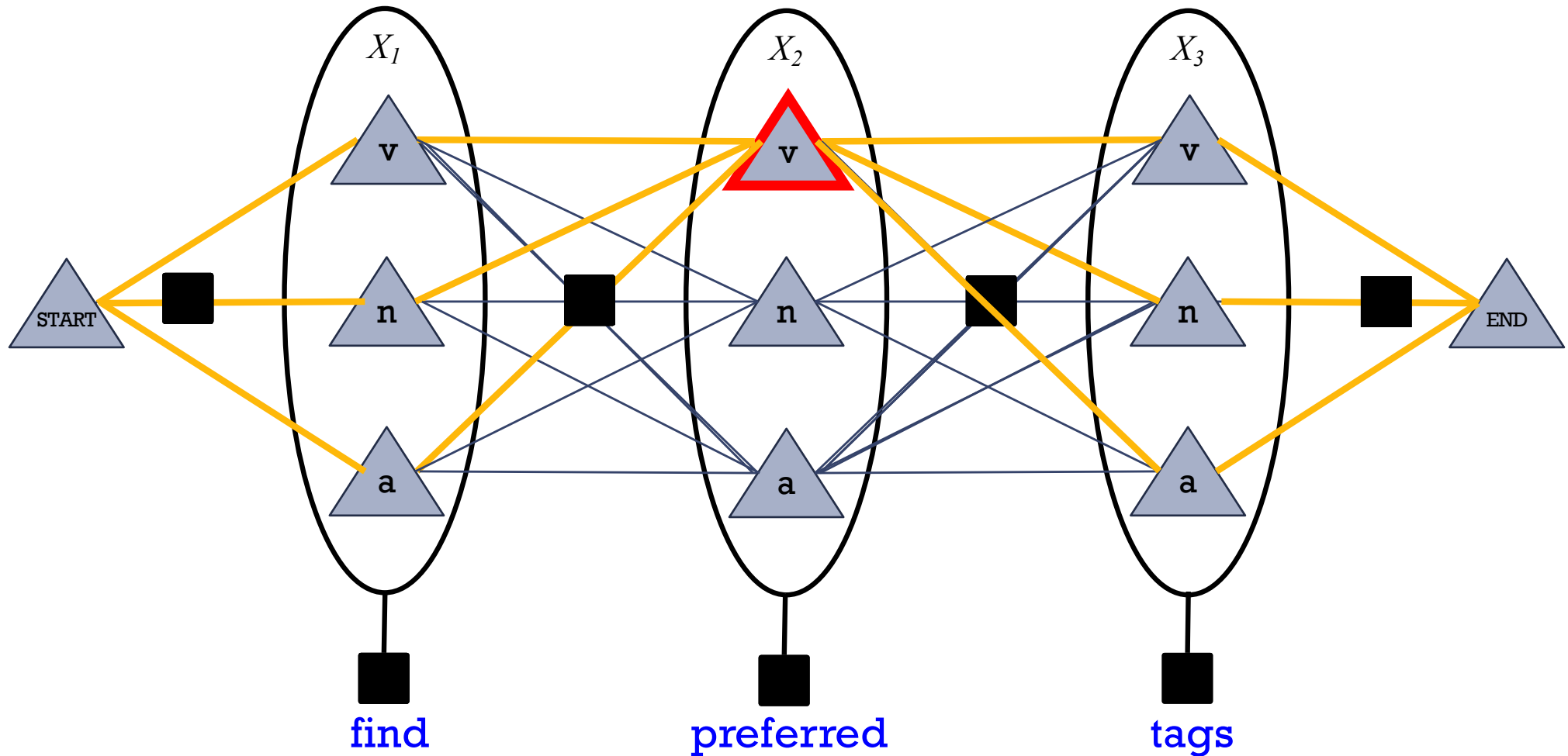
Forward-Backward Algorithm: Finds Marginals



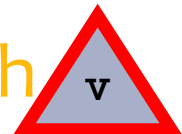
- So $p(v \ a \ n) = (1/Z) * \text{product weight of one path}$
- Marginal probability $p(X_2 = a)$
 $= (1/Z) * \text{total weight of all paths through}$



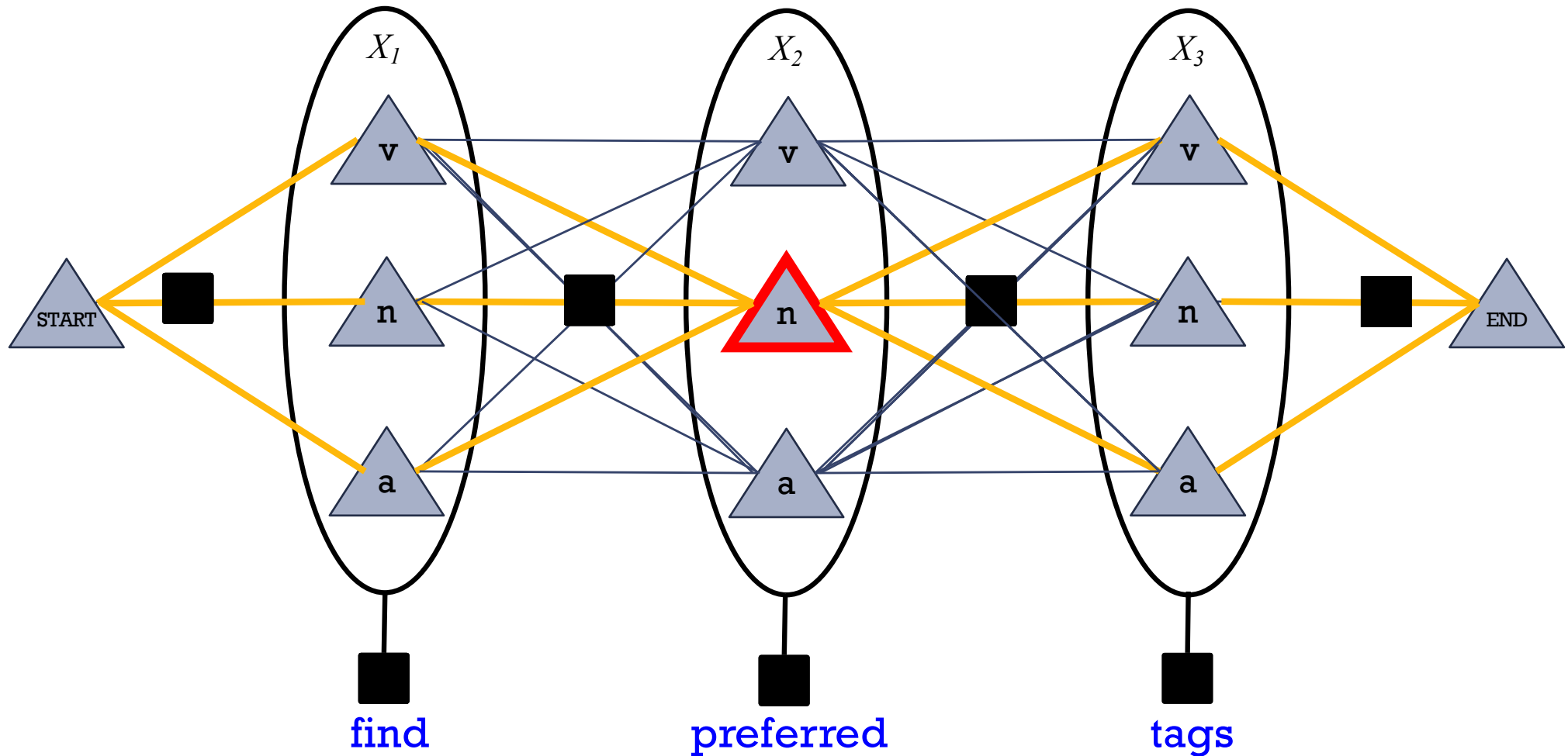
Forward-Backward Algorithm: Finds Marginals



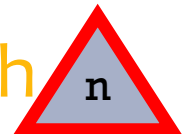
- So $p(\mathbf{v} \mathbf{a} \mathbf{n}) = (1/Z) * \text{product weight of one path}$
- Marginal probability $p(X_2 = a)$
 $= (1/Z) * \text{total weight of all paths through}$



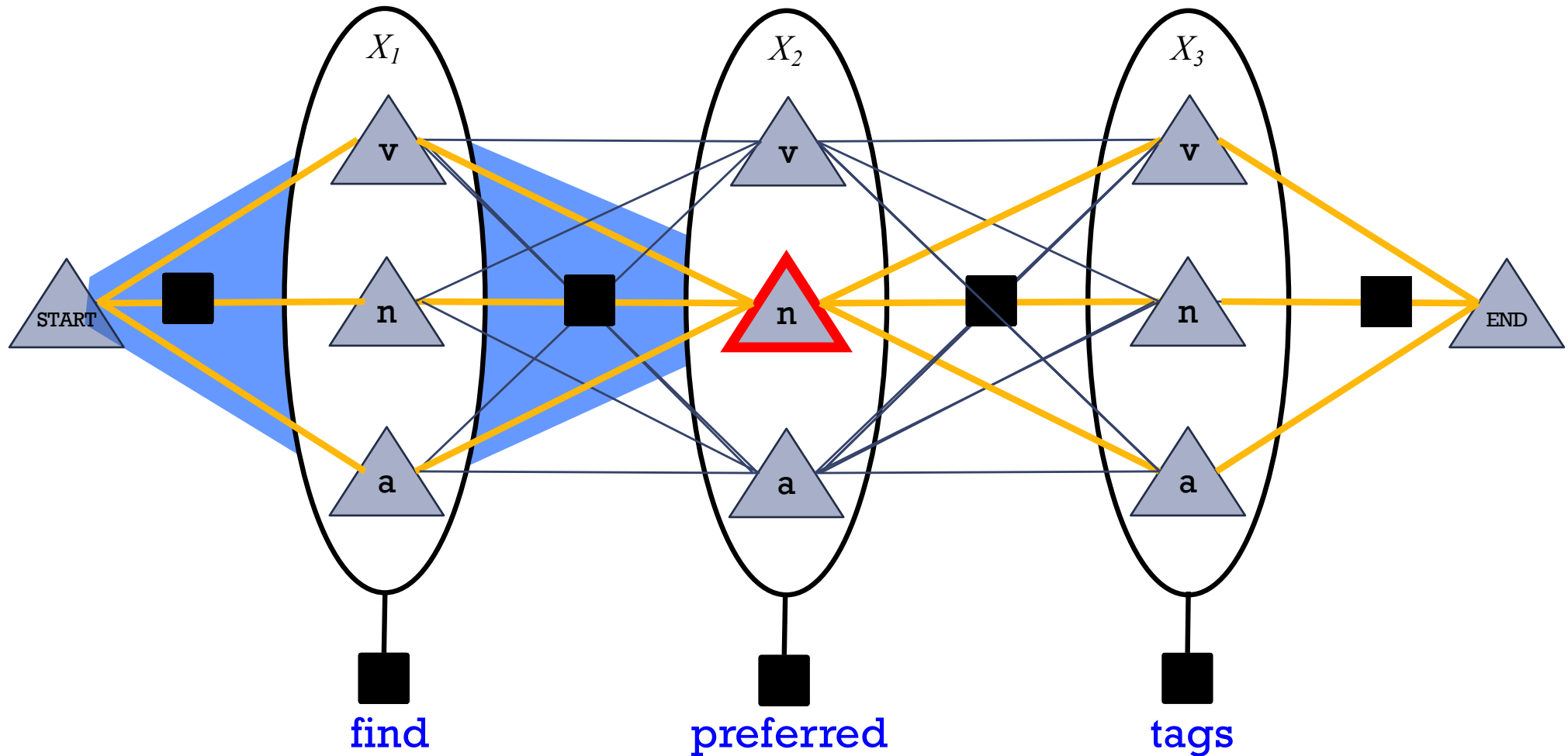
Forward-Backward Algorithm: Finds Marginals



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



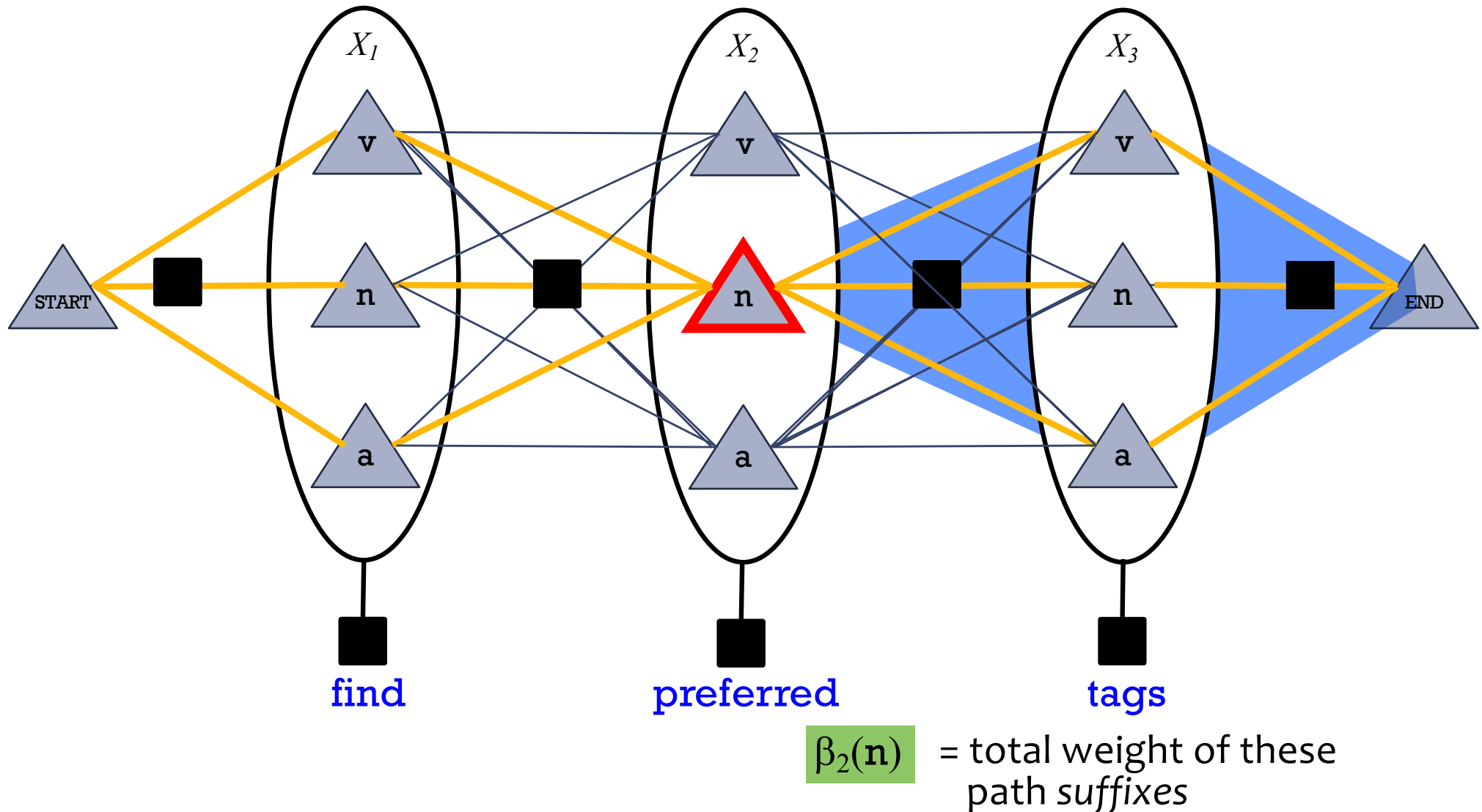
Forward-Backward Algorithm: Finds Marginals



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

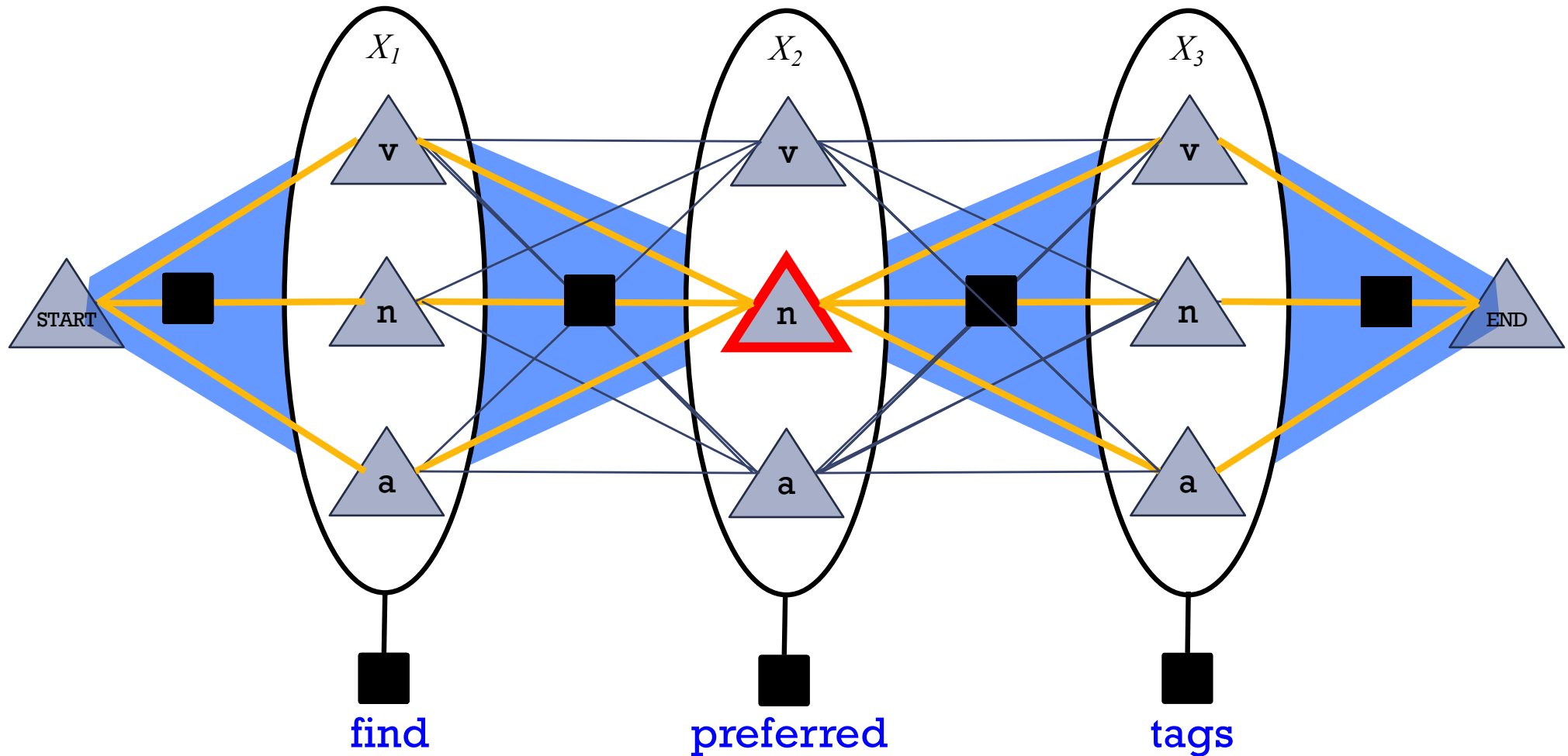
(found by dynamic programming: matrix-vector products)

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 ($a + b + c$)

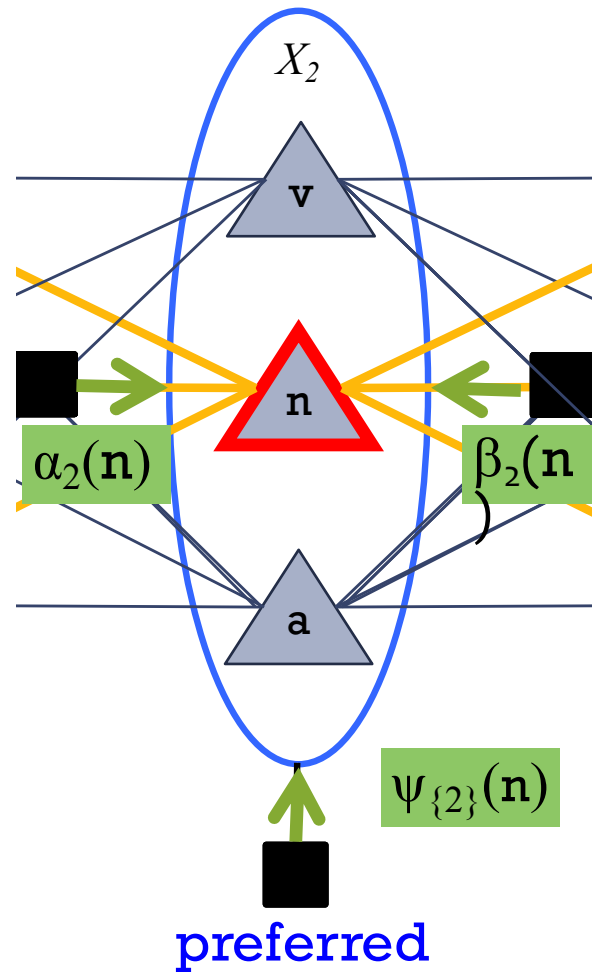
$\beta_2(n)$ = total weight of these path suffixes ($x + y + z$)

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.

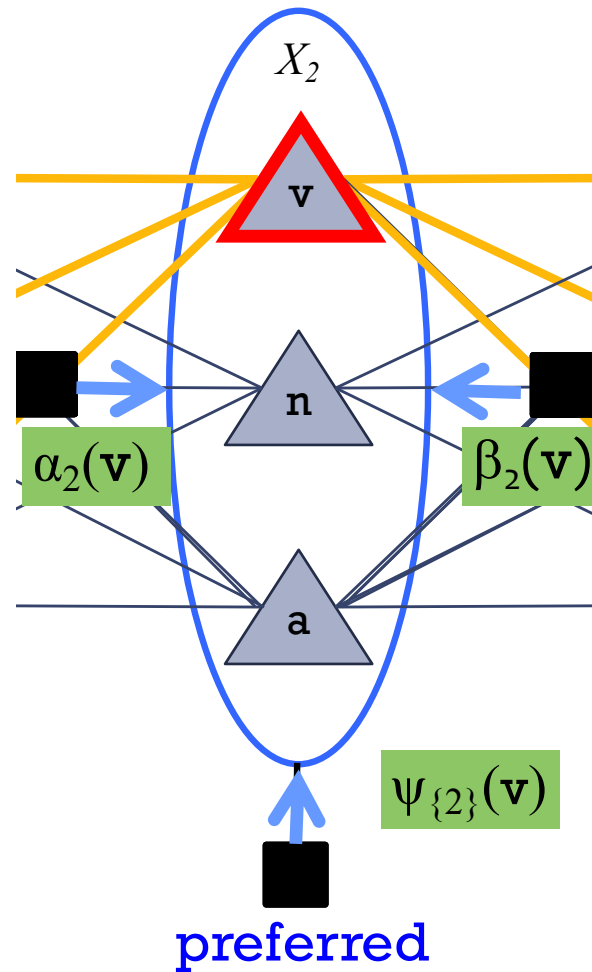


“belief that $X_2 = \mathbf{n}$ ”

total weight of *all* paths through 

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

Forward-Backward Algorithm: Finds Marginals

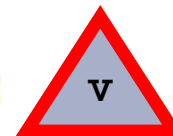


“belief that $X_2 = \mathbf{v}$ ”

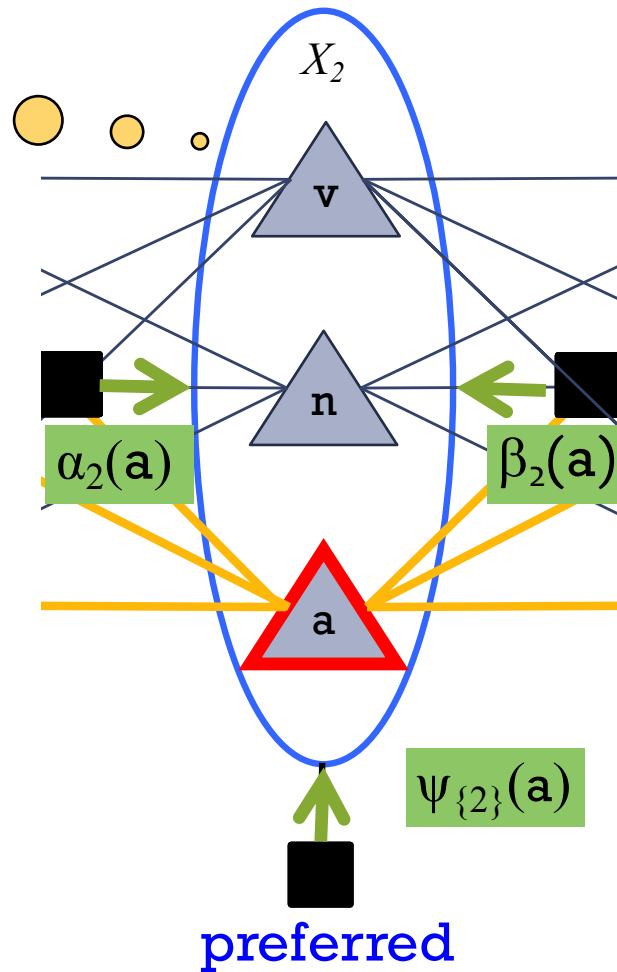
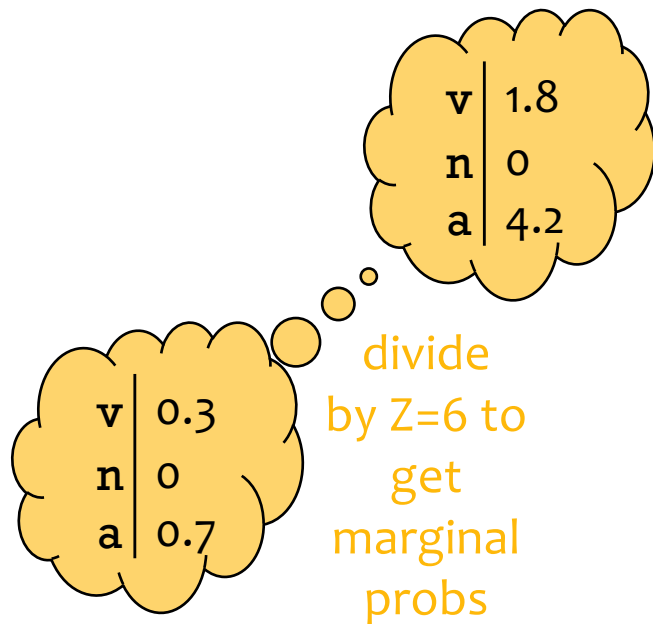
“belief that $X_2 = \mathbf{n}$ ”

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 

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

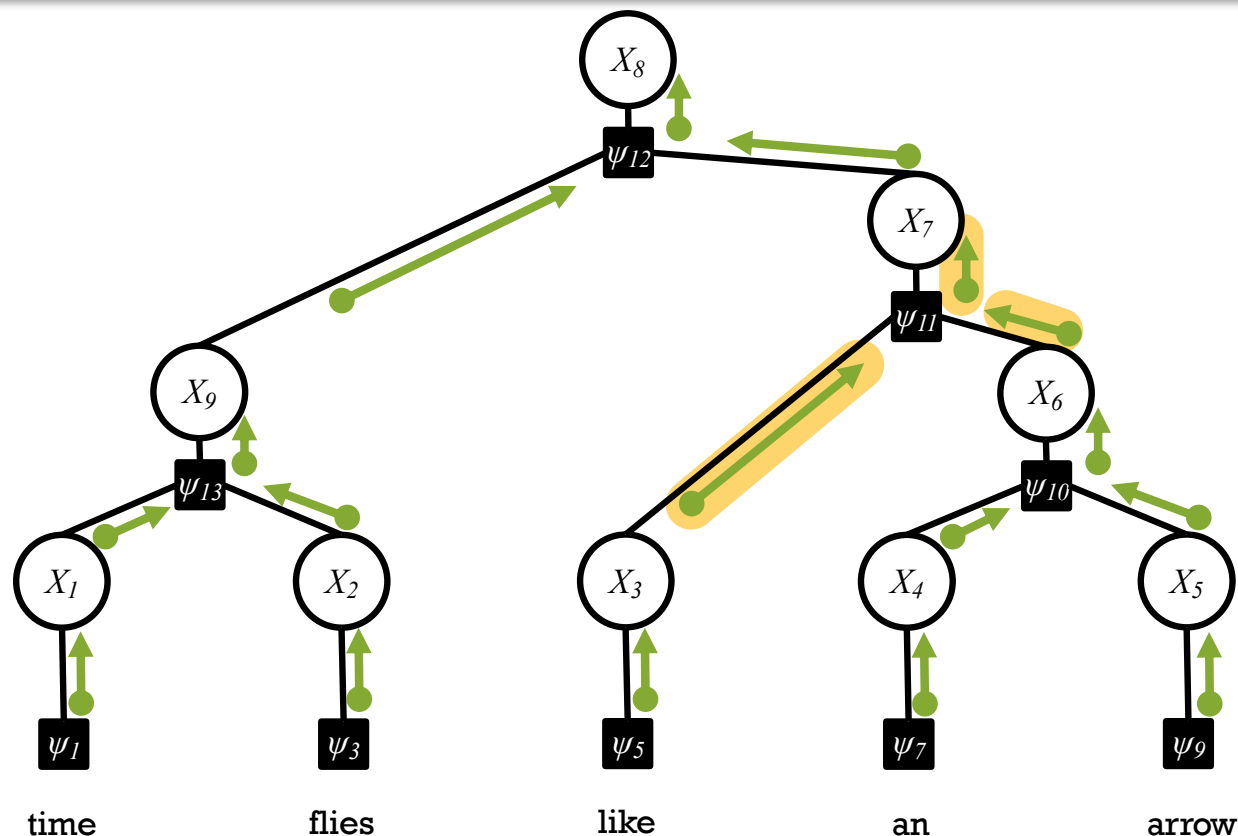
BP AS DYNAMIC PROGRAMMING

(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.

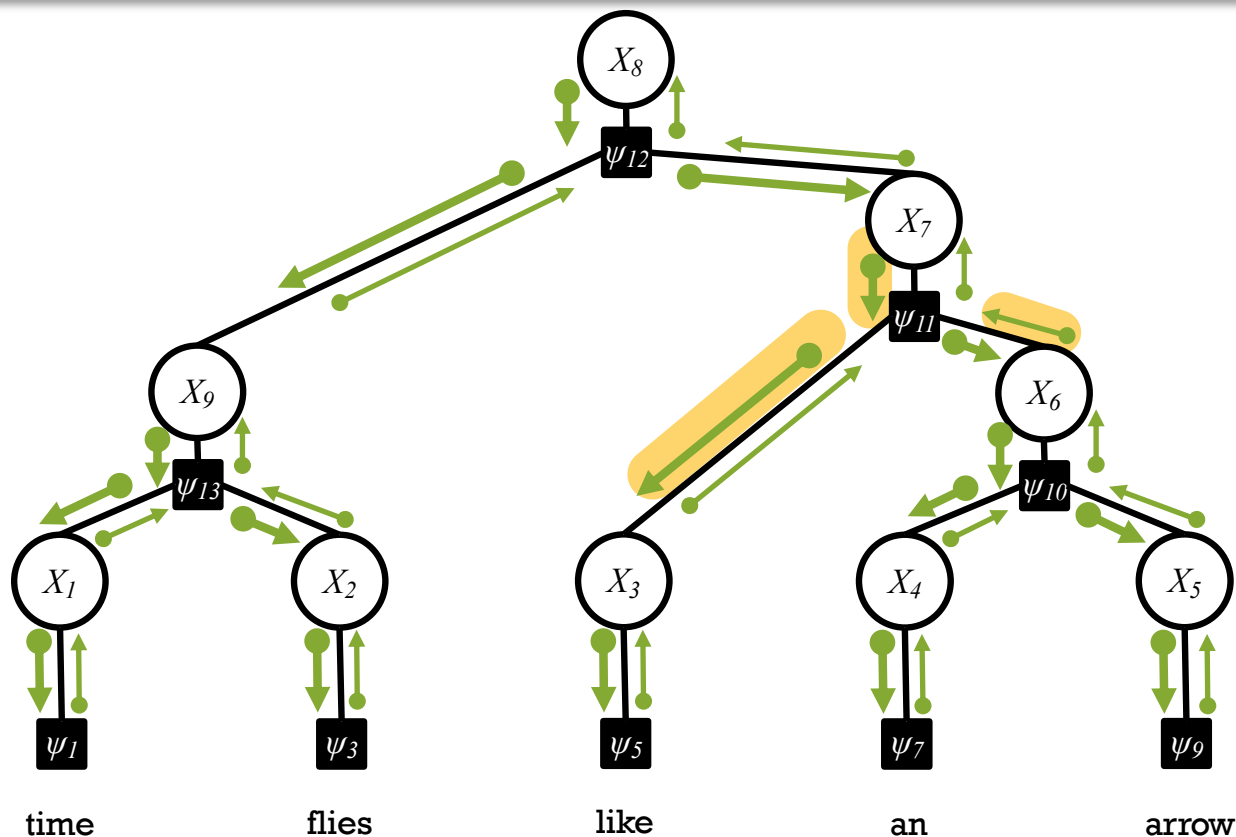


(Acyclic) Belief Propagation

In a factor graph with no cycles:

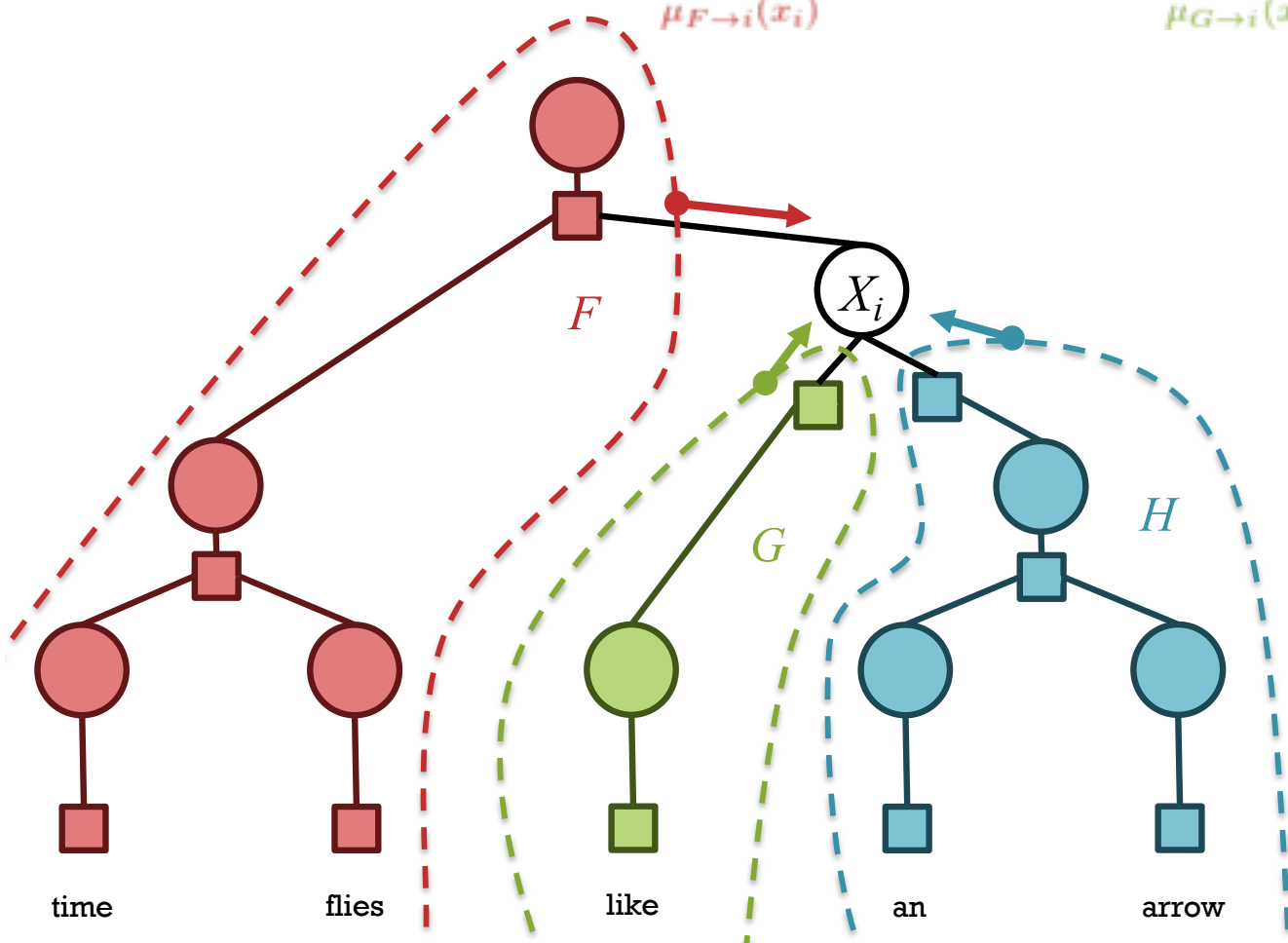
1. Pick any node to serve as the root.
2. Send messages from the **leaves** to the **root**.
3. Send messages from the **root** to the **leaves**.

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

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



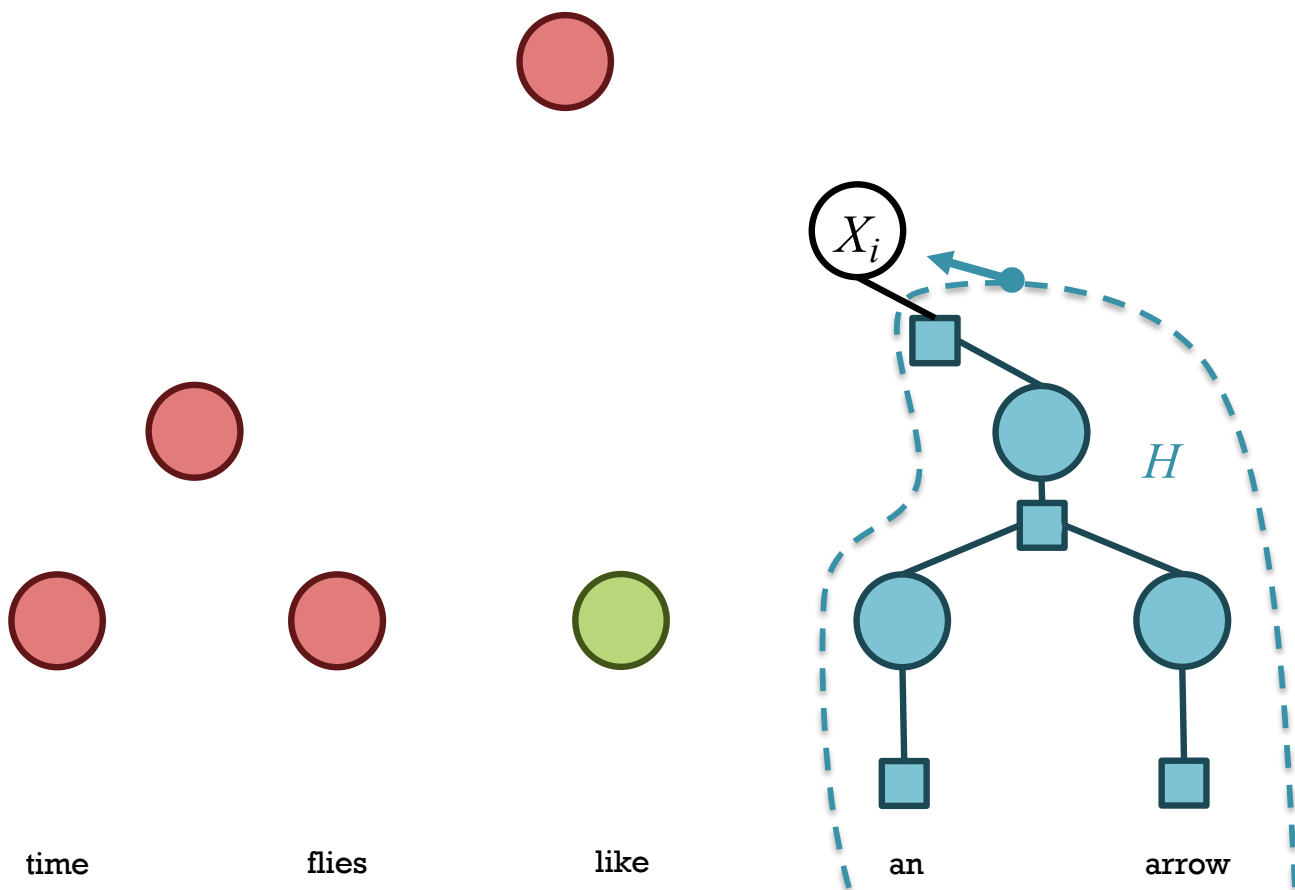
Subproblem:

Inference using just the factors in subgraph H

Figure adapted from
Burkett & Klein (2012)⁷¹

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:

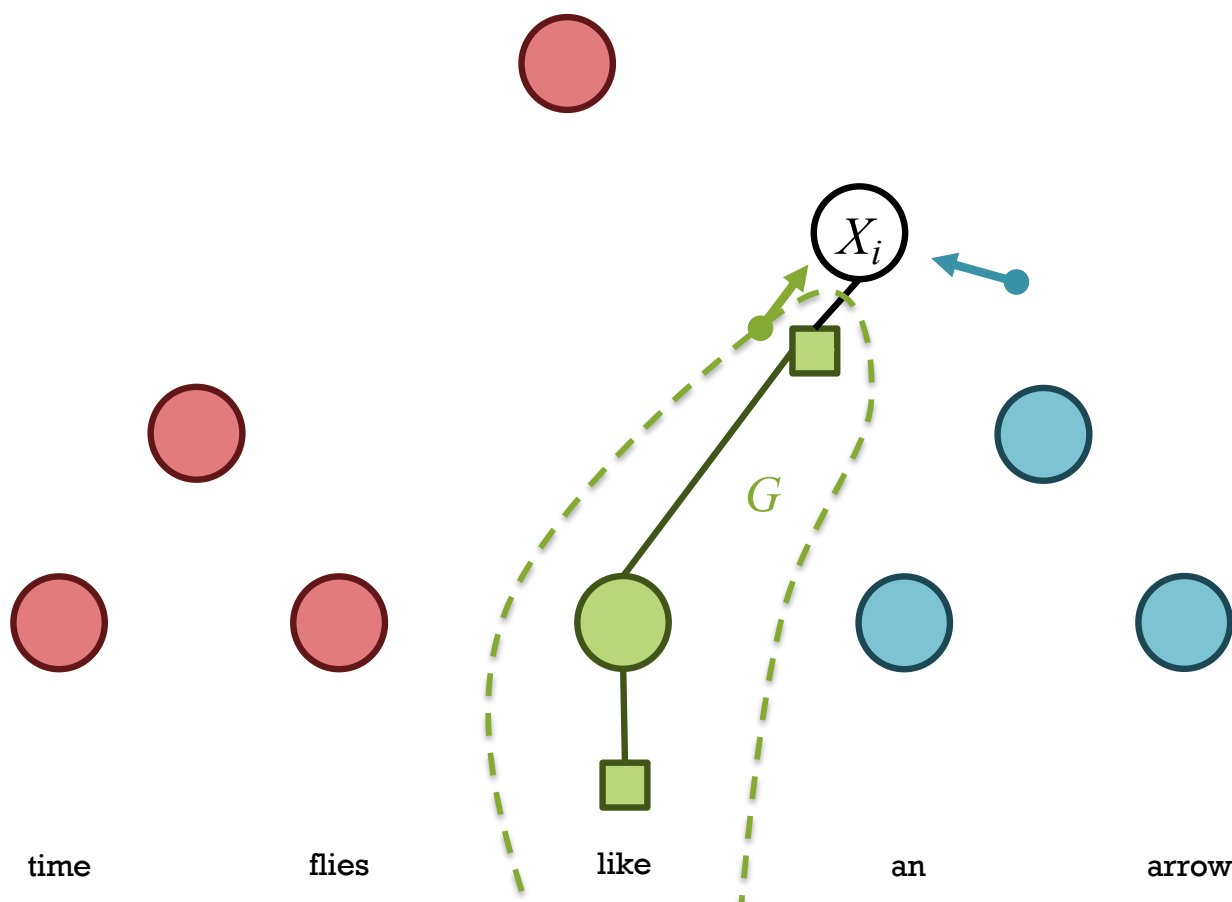
Inference using just the factors in subgraph H

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:

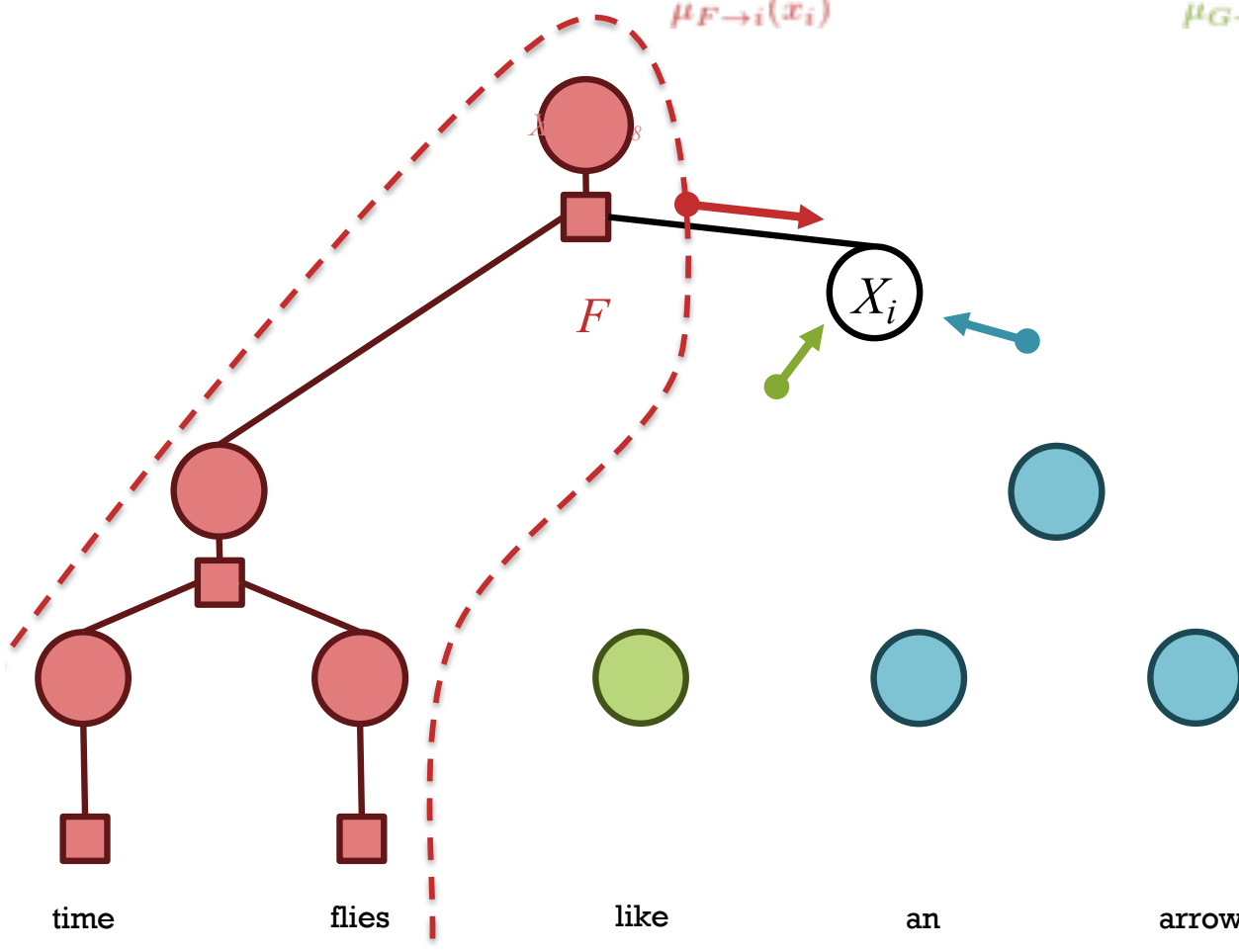
Inference using just the factors in subgraph H

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

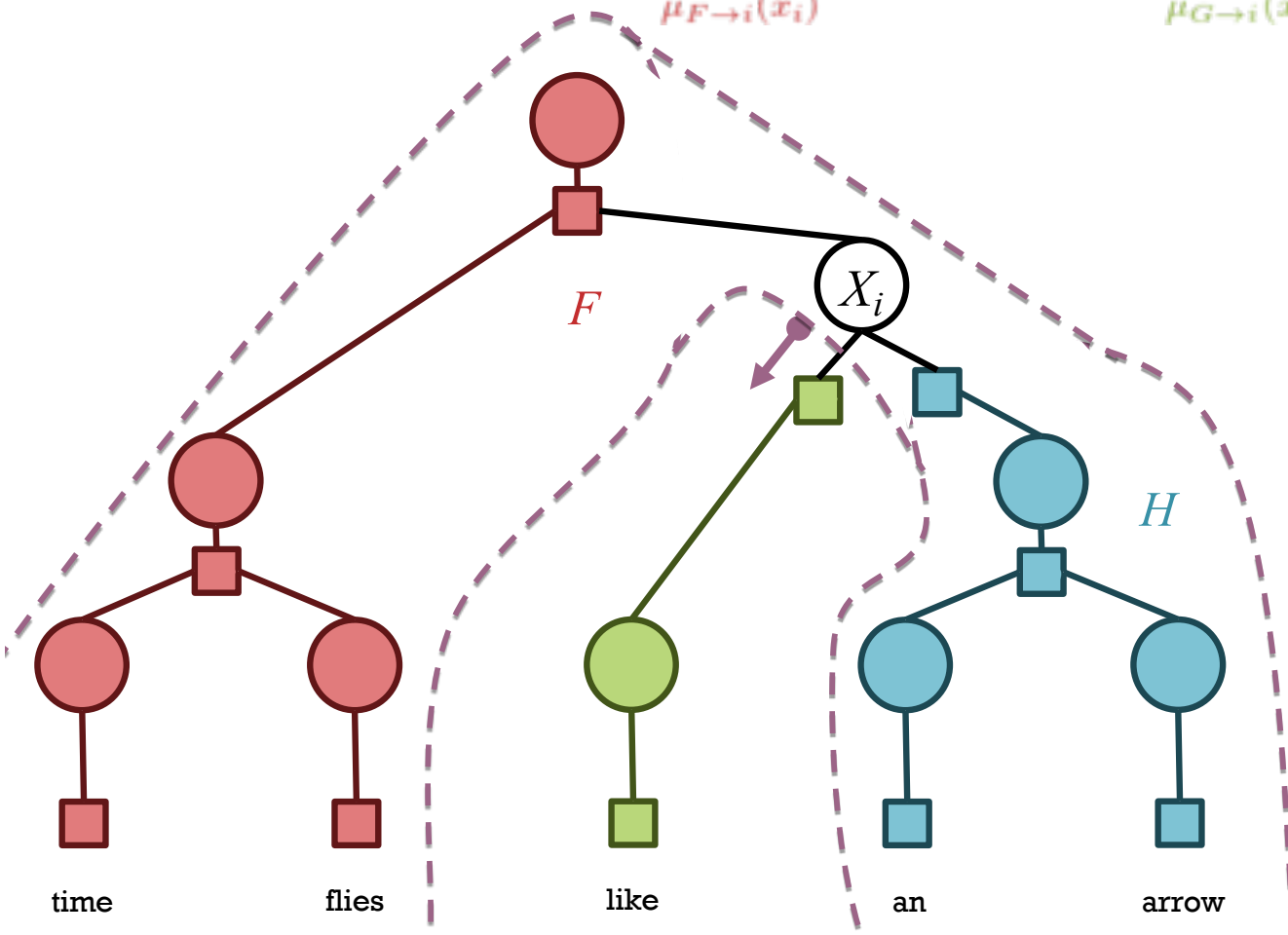
The marginal of X_i in that smaller model is the message sent to X_i from subgraph H

*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:

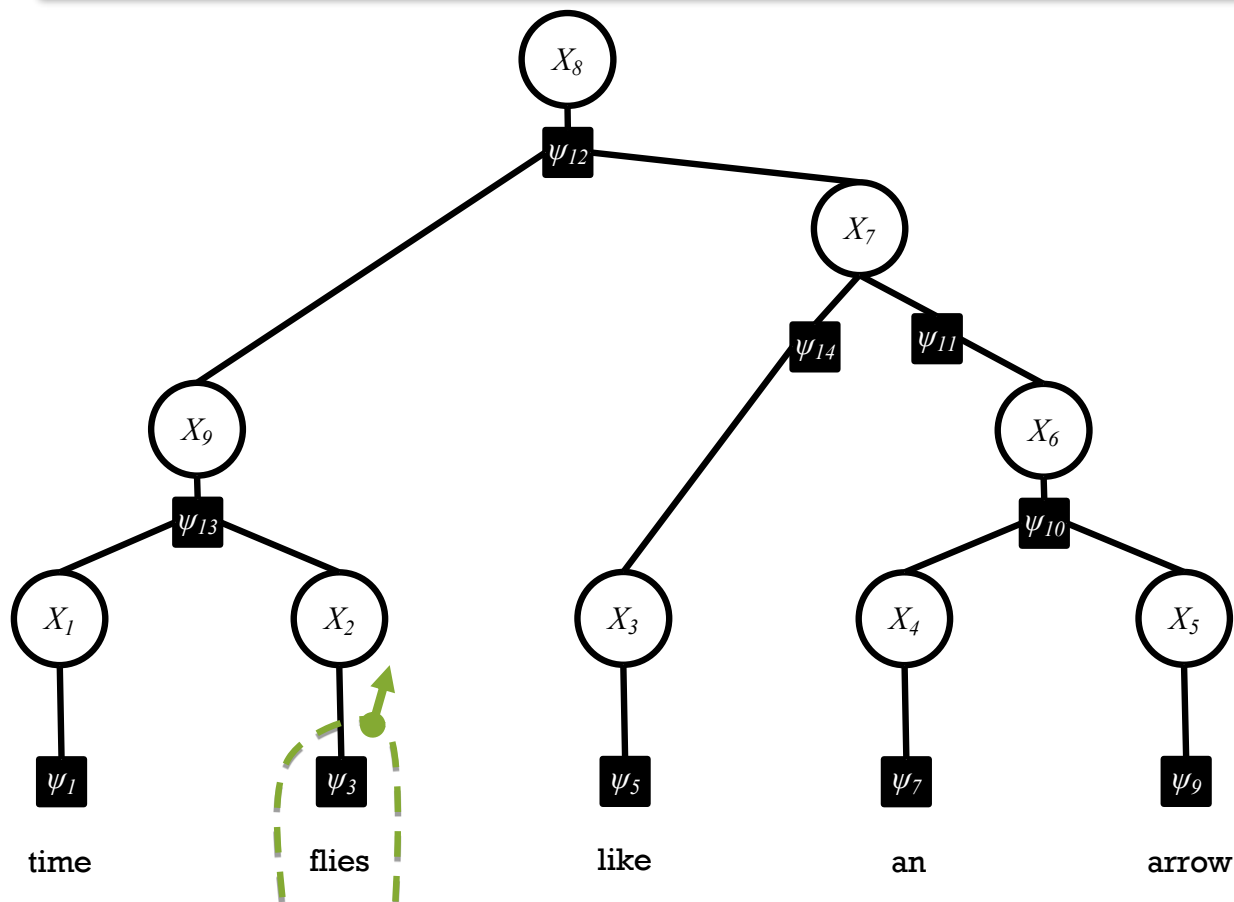
Inference using just the factors in subgraph $F \cup H$

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

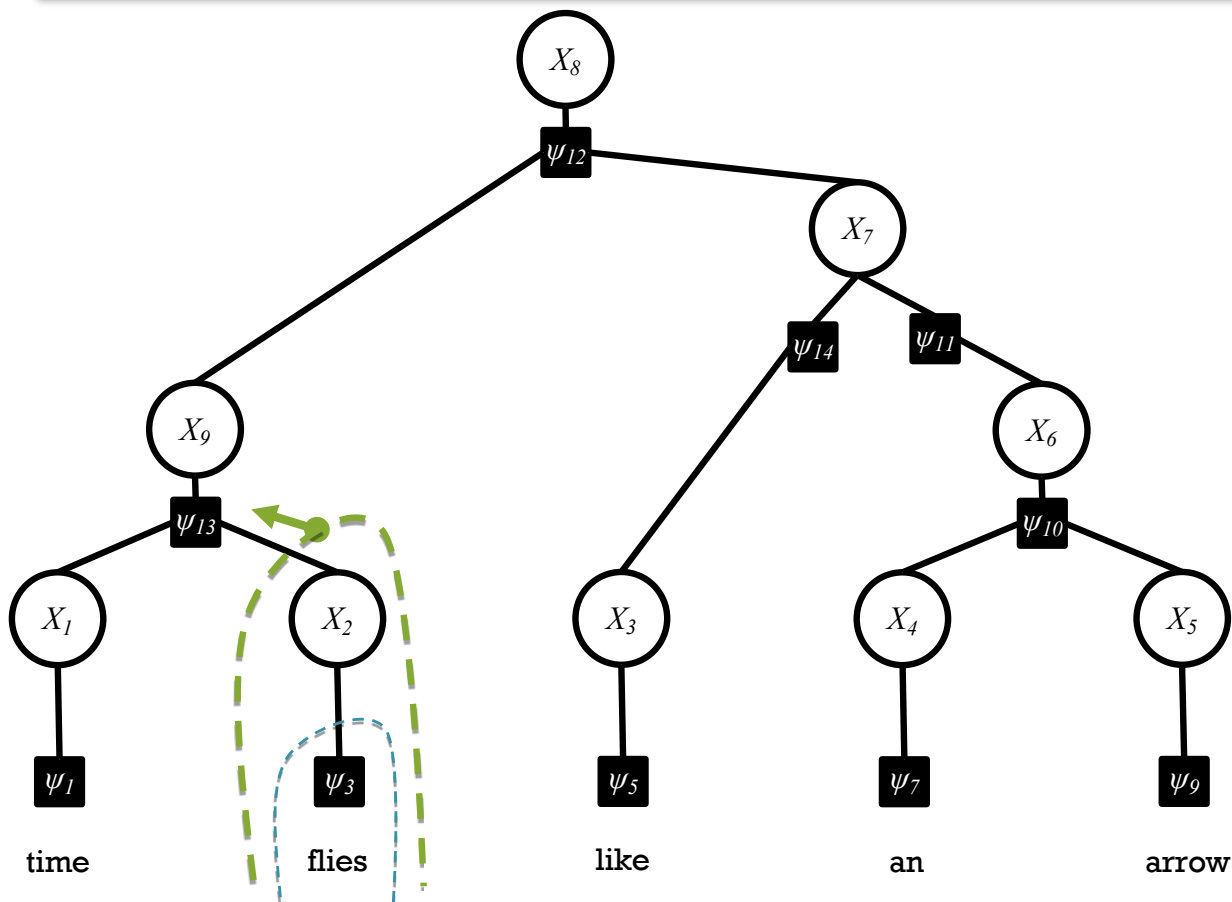
Acyclic BP as Dynamic Programming

- 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.



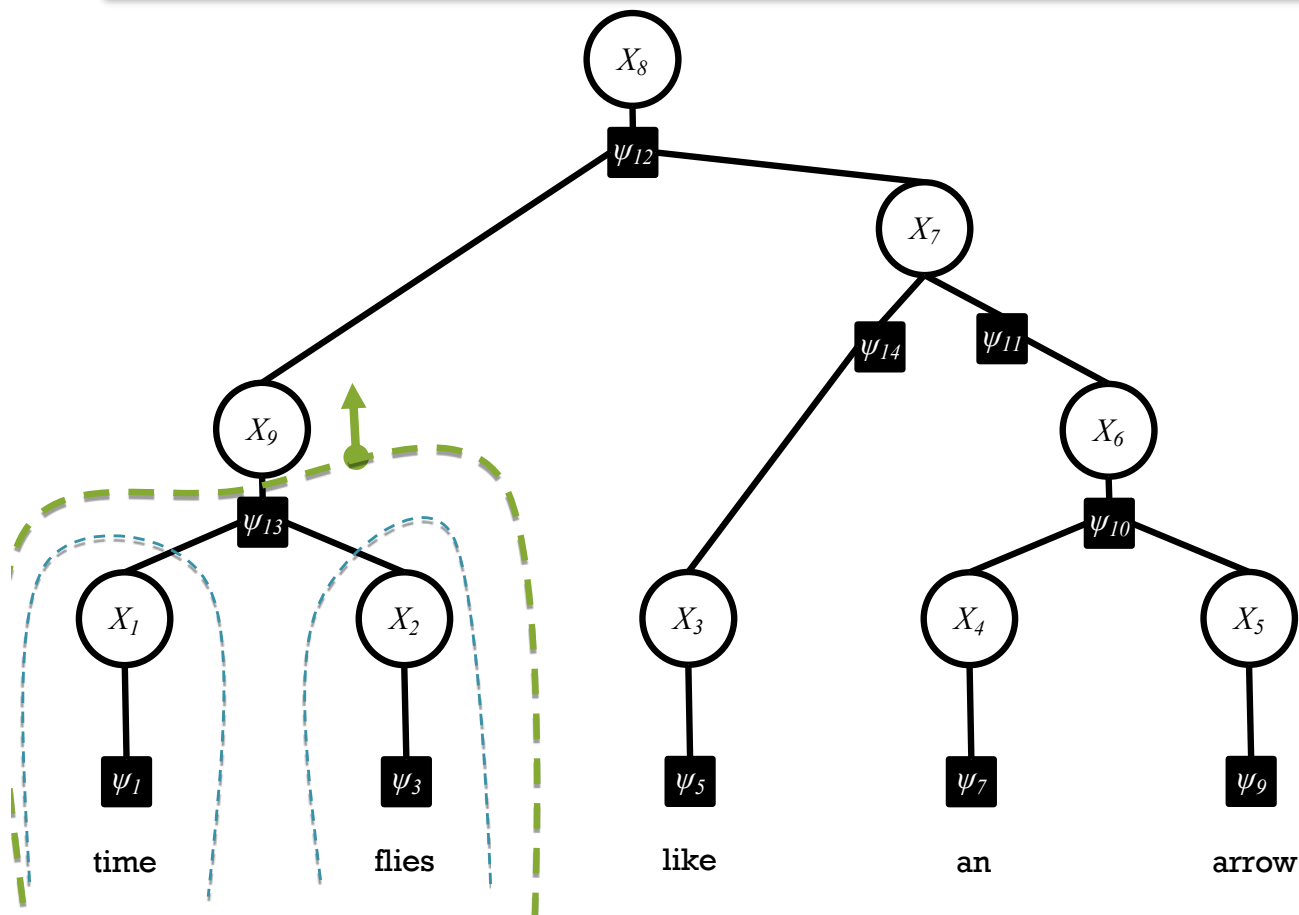
Acyclic BP as Dynamic Programming

- 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.



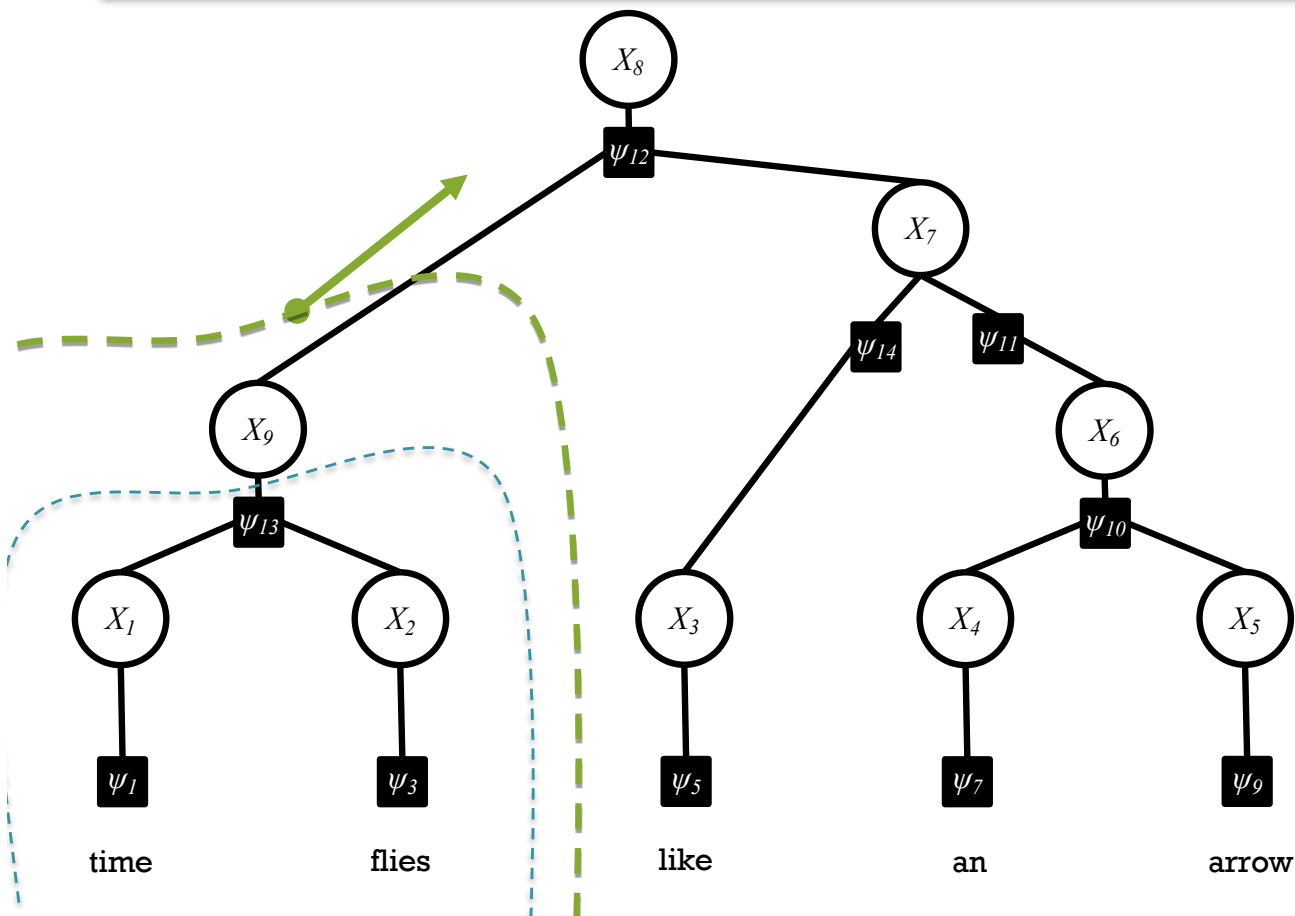
Acyclic BP as Dynamic Programming

- 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.



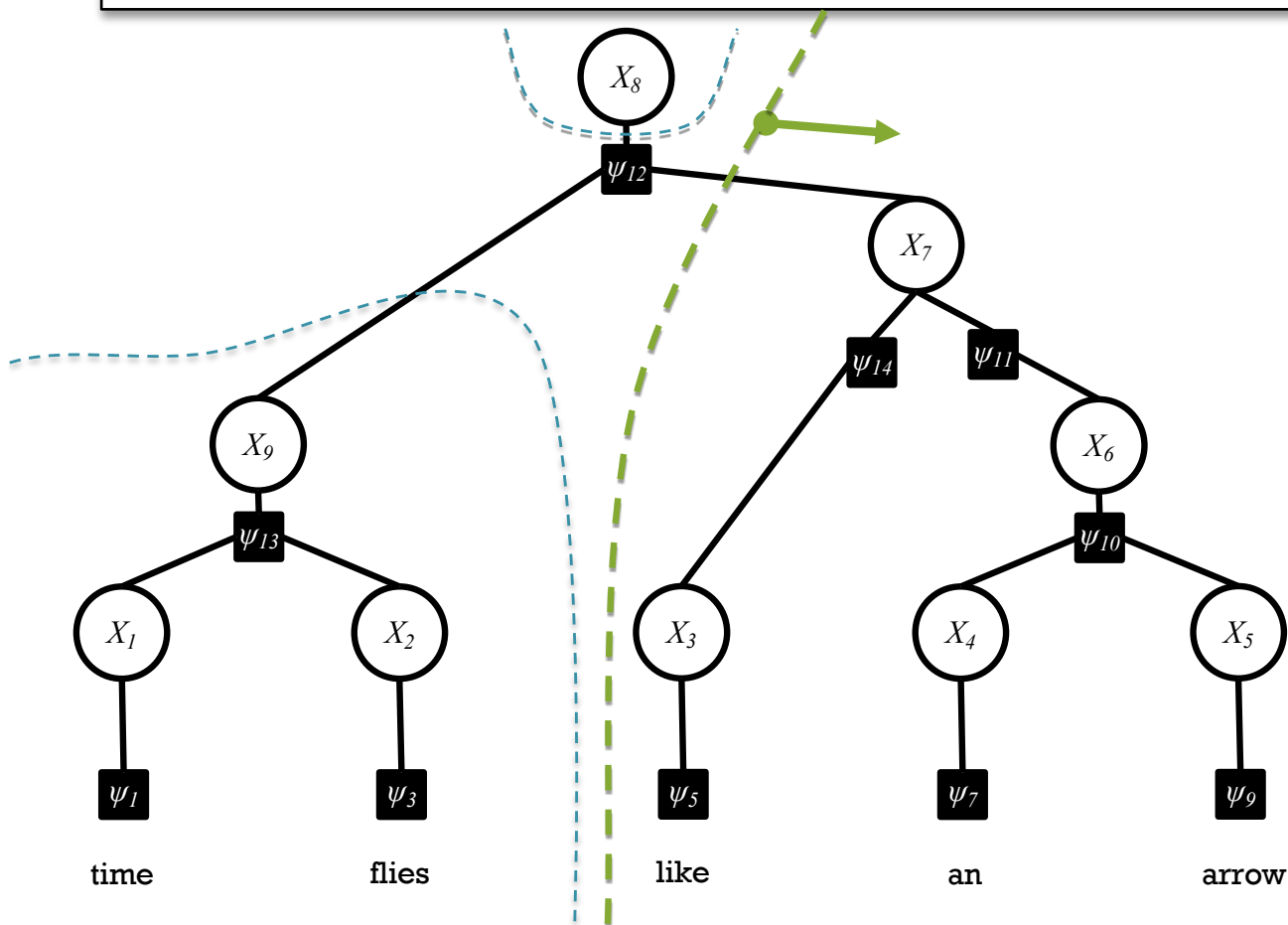
Acyclic BP as Dynamic Programming

- 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.



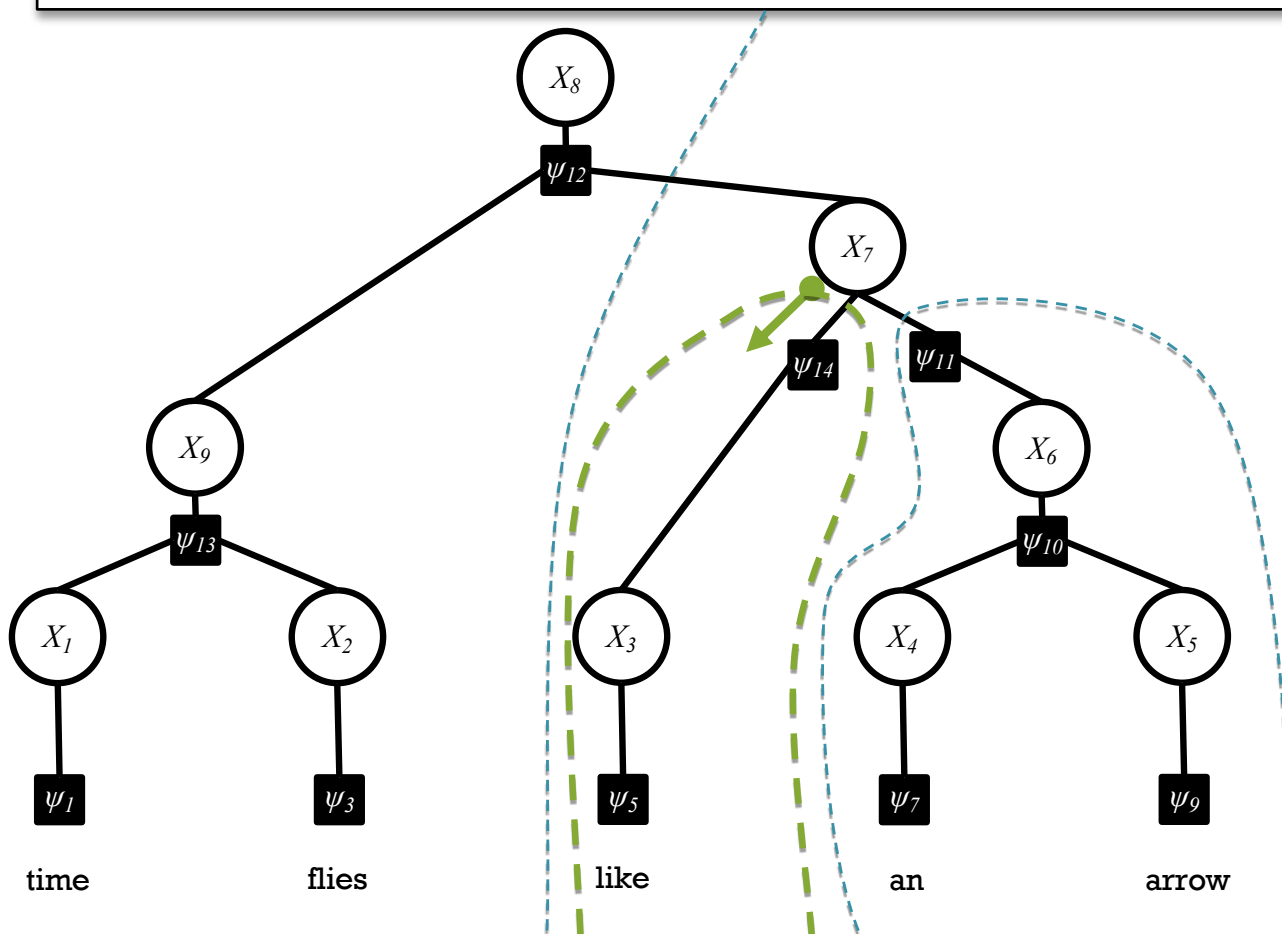
Acyclic BP as Dynamic Programming

- 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.



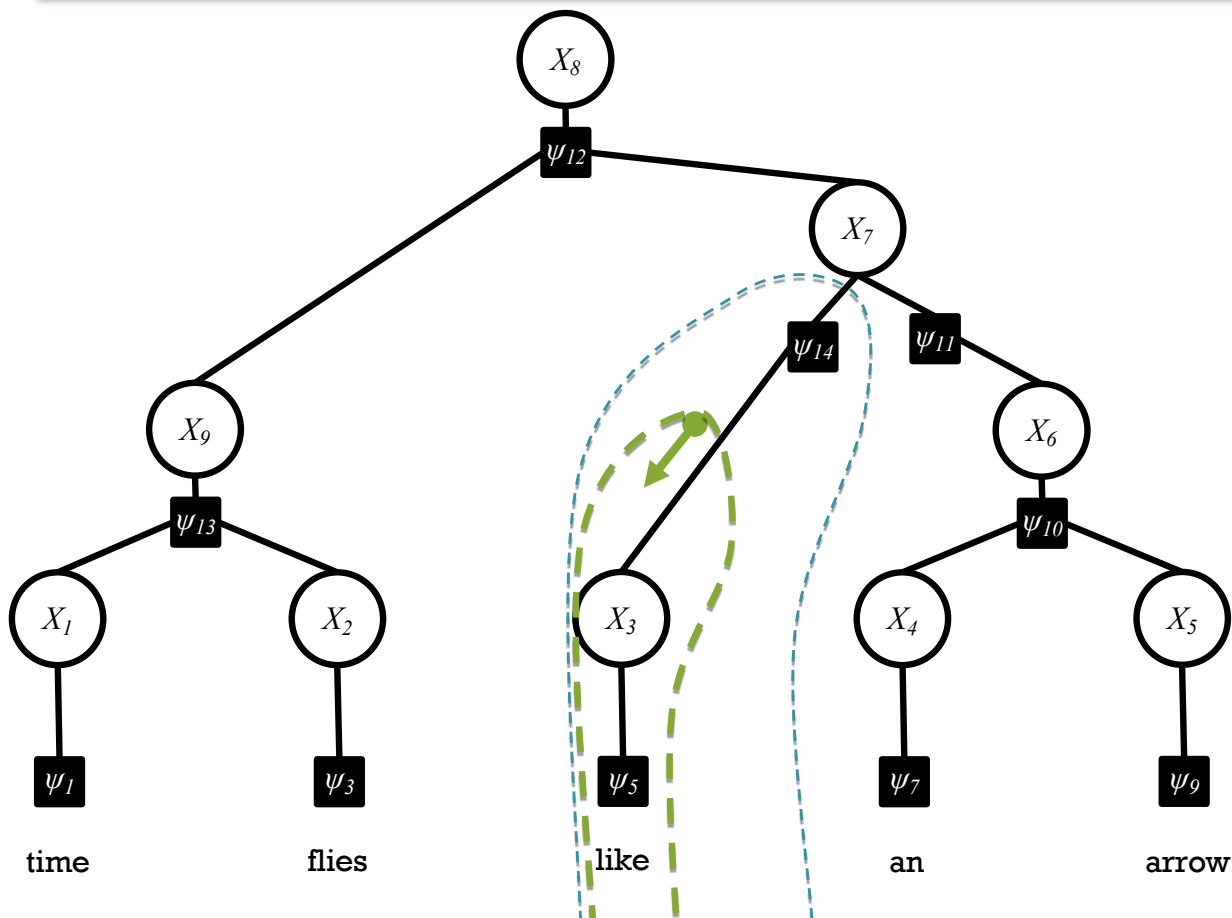
Acyclic BP as Dynamic Programming

- 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:

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

Deterministic Annealing

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

1. Incorporate inverse temperature parameter into each factor:

Annealed Joint Distribution

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

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 $\max/+$

Exact inference for linear chain models

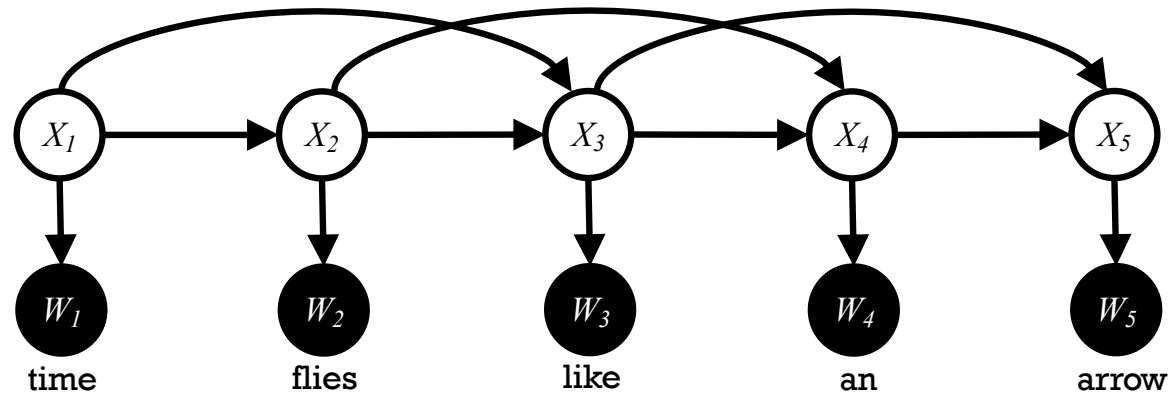
FORWARD-BACKWARD AND VITERBI ALGORITHMS

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**

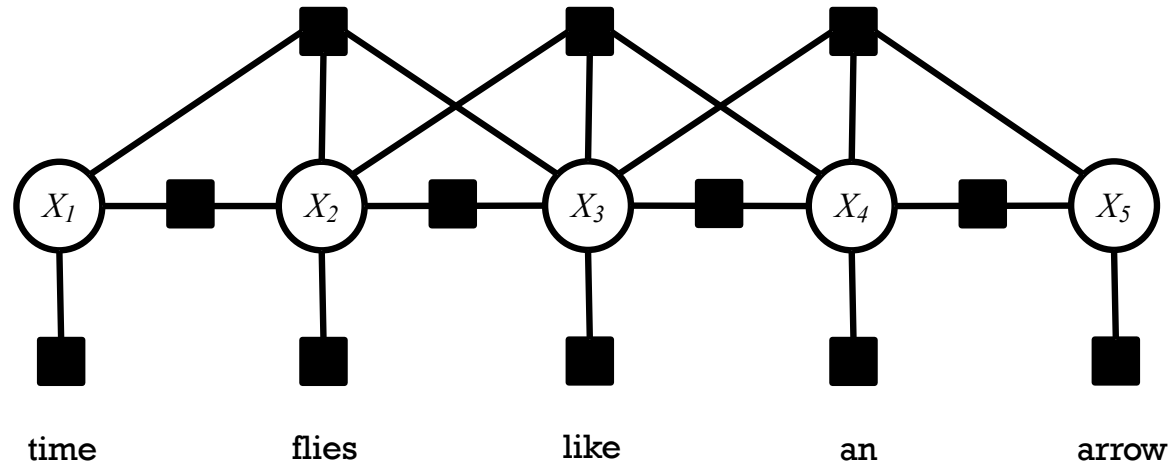
Forward-Backward Algorithm

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



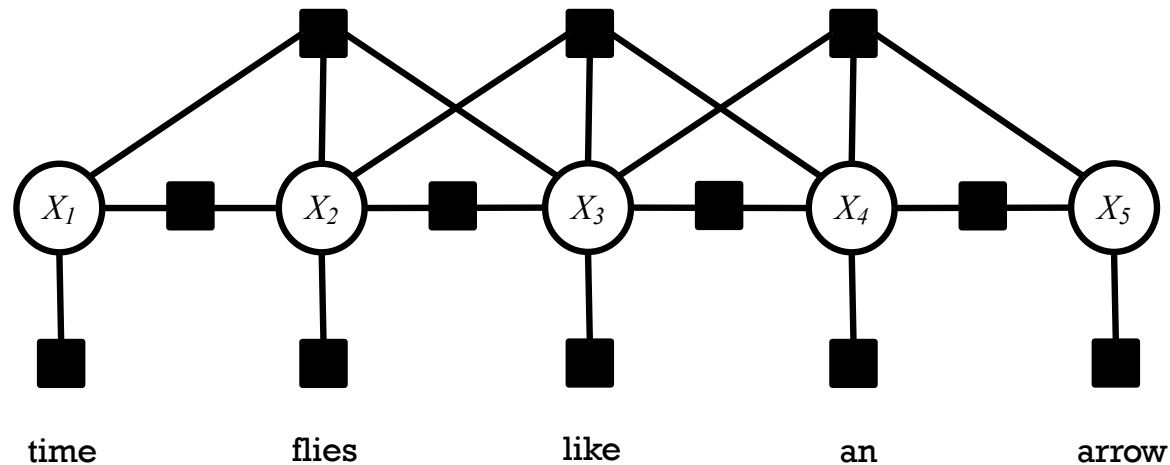
Forward-Backward Algorithm

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.

An example of why we need approximate inference

EXACT INFERENCE ON GRID CRF

Application: Pose Estimation

$\phi_i(y_i, x) \in \mathbb{R}^{\approx 1000}$: local image representation, e.g. HoG

$\rightarrow \langle w_i, \phi_i(y_i, x) \rangle$: local confidence map

$\phi_{i,j}(y_i, y_j) = \text{good_fit}(y_i, y_j) \in \mathbb{R}^1$: test for geometric fit

$\rightarrow \langle w_{ij}, \phi_{ij}(y_i, y_j) \rangle$: penalizer for unrealistic poses

together: $\operatorname{argmax}_y p(y|x)$ is sanitized version of local cues



original



local classification



local + geometry

Feature Functions for CRF in Vision

$\phi_i(y_i, x)$: local representation, high-dimensional

→ $\langle w_i, \phi_i(y_i, x) \rangle$: local classifier

$\phi_{i,j}(y_i, y_j)$: prior knowledge, low-dimensional

→ $\langle w_{ij}, \phi_{ij}(y_i, y_j) \rangle$: penalize outliers

learning adjusts parameters:

- ▶ unary w_i : learn local classifiers and their importance
- ▶ binary w_{ij} : learn importance of smoothing/penalization

$\operatorname{argmax}_y p(y|x)$ is cleaned up version of local prediction

Case Study: Image Segmentation

- Image segmentation (FG/BG) by modeling of interactions btw RVs
 - Images are noisy.
 - Objects occupy continuous regions in an image.

[Nowozin, Lampert 2012]



Input image



Pixel-wise separate
optimal labeling



Locally-consistent
joint optimal labeling

$$Y^* = \arg \max_{y \in \{0,1\}^n} \left[\overbrace{\sum_{i \in S} V_i(y_i, X)}^{\text{Unary Term}} + \overbrace{\sum_{i \in S} \sum_{j \in N_i} V_{i,j}(y_i, y_j)}^{\text{Pairwise Term}} \right].$$

© Eric Xing @ CMU, 2005-2015

Y : labels

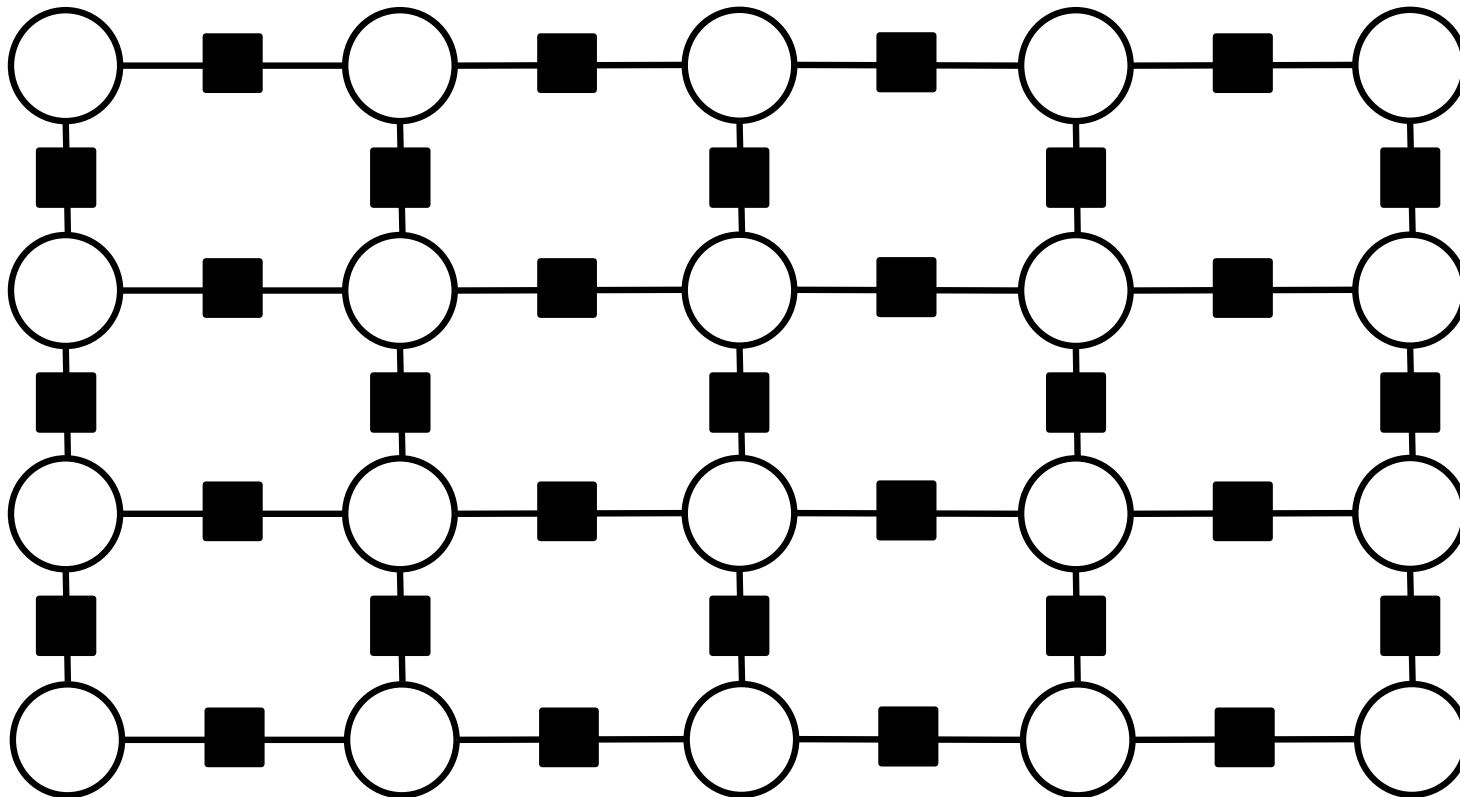
X : data (features)

S : pixels

N_i : neighbors of pixel i

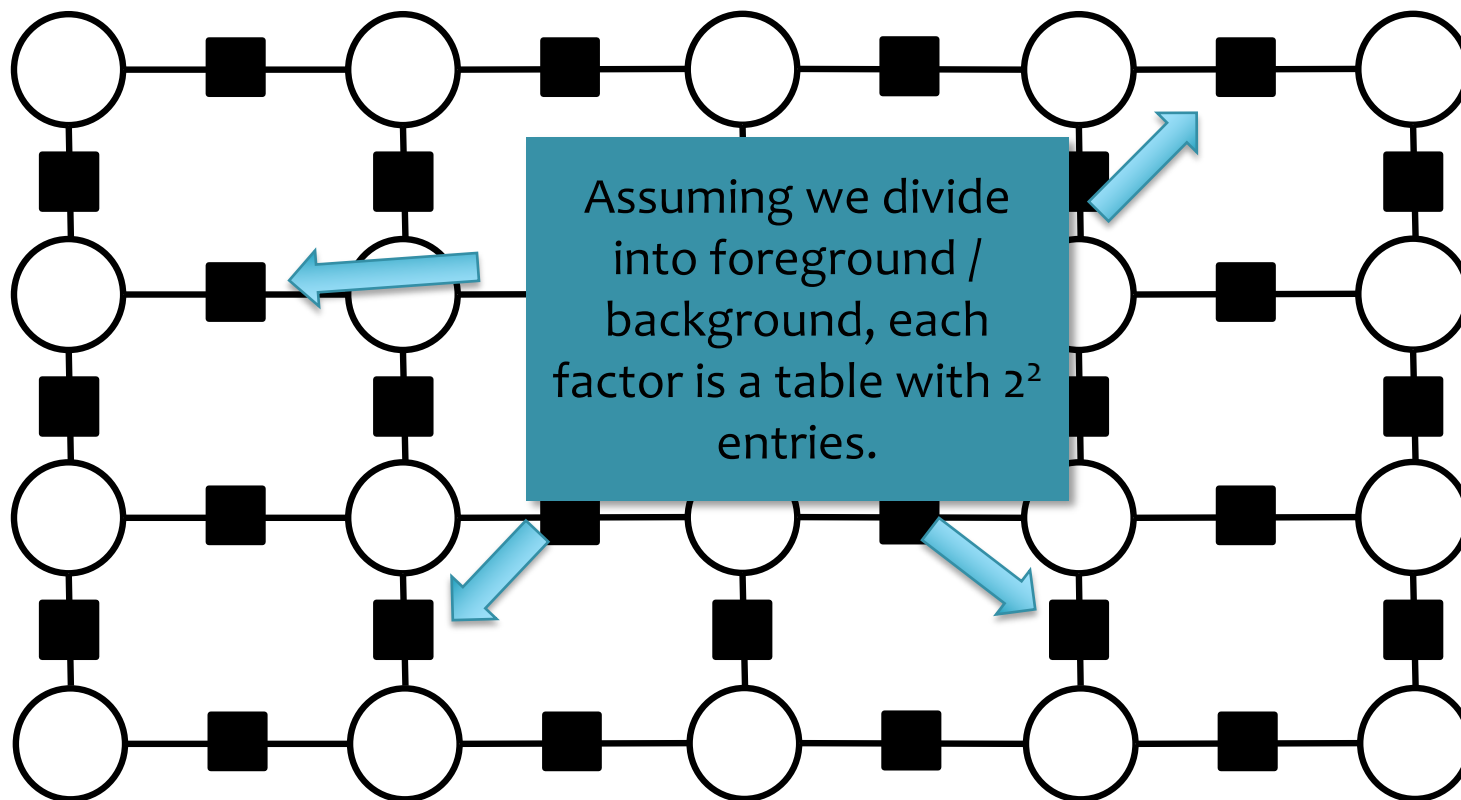
Grid CRF

- Suppose we want to image segmentation using a grid model



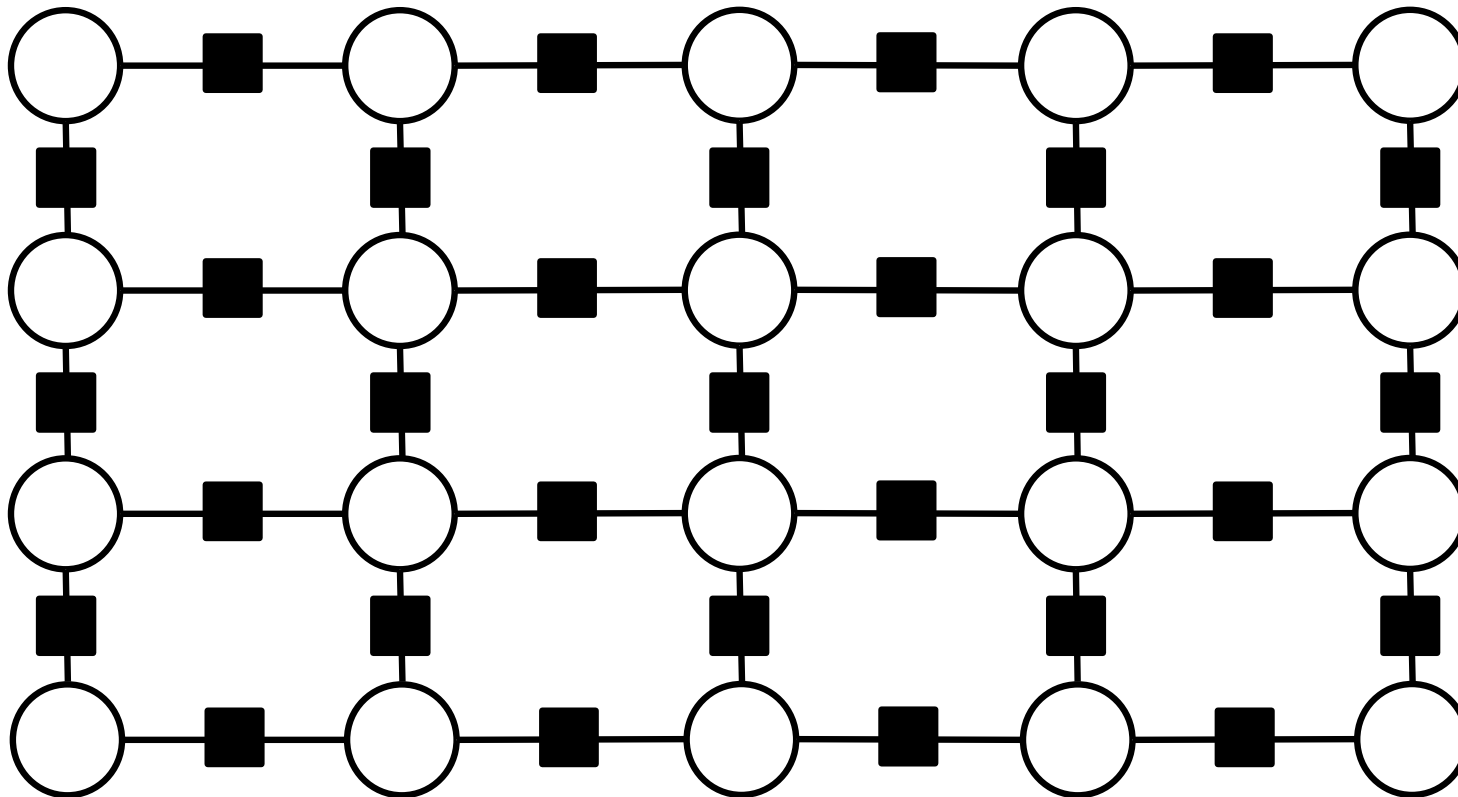
Grid CRF

- Suppose we want to image segmentation using a grid model



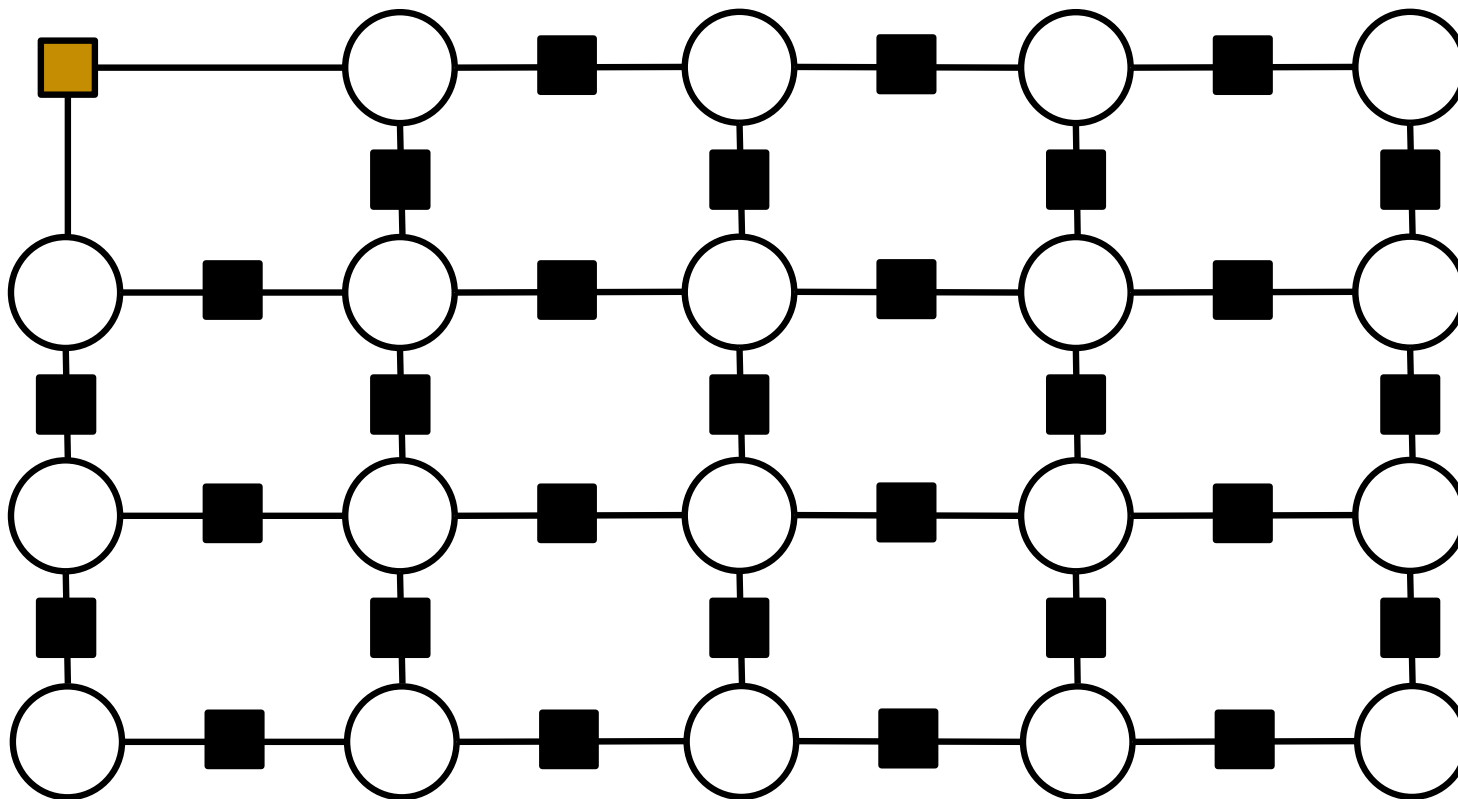
Grid CRF

- Suppose we want to image segmentation using a grid model
- What happens when we run variable elimination?



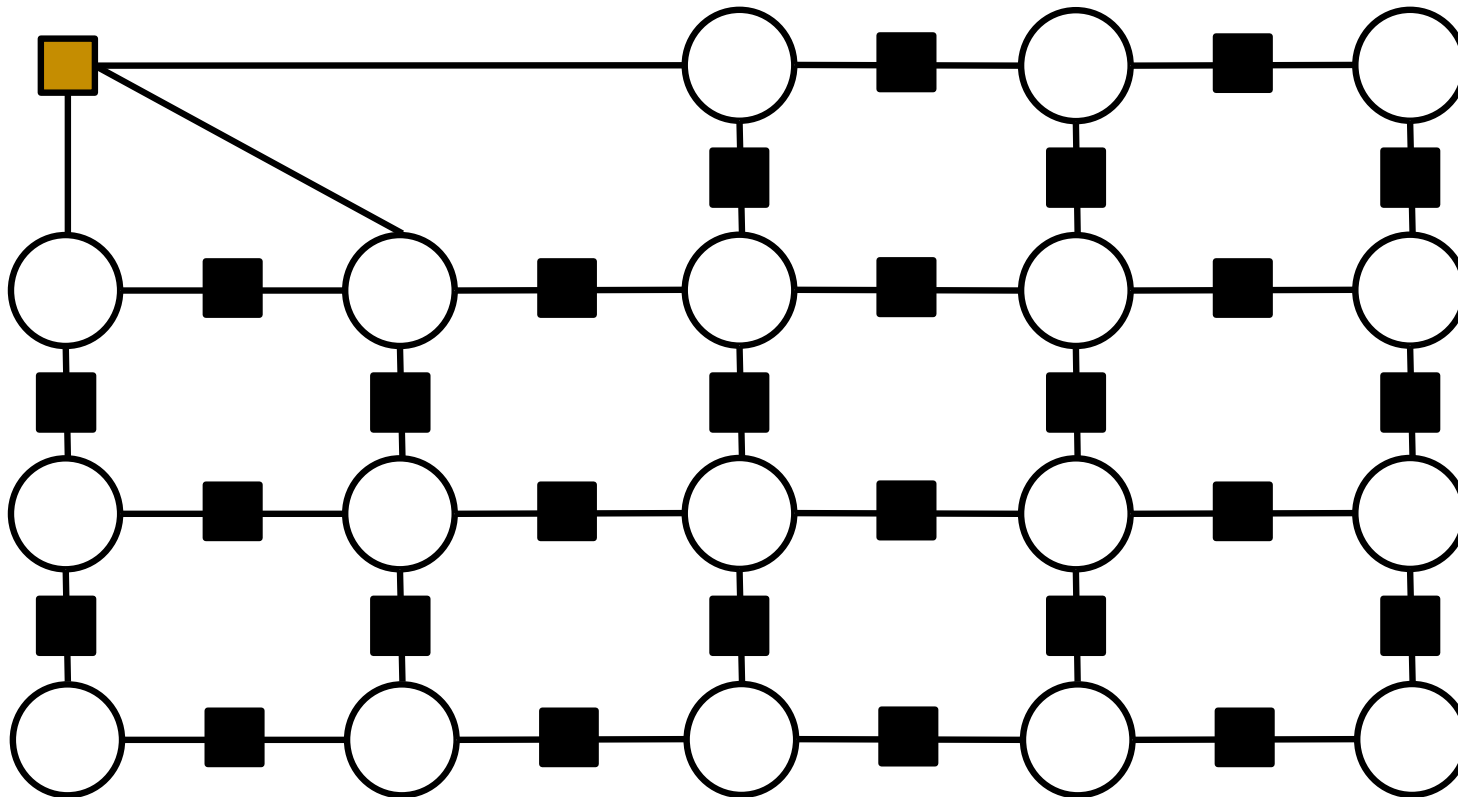
Grid CRF

- Suppose we want to image segmentation using a grid model
- What happens when we run variable elimination?



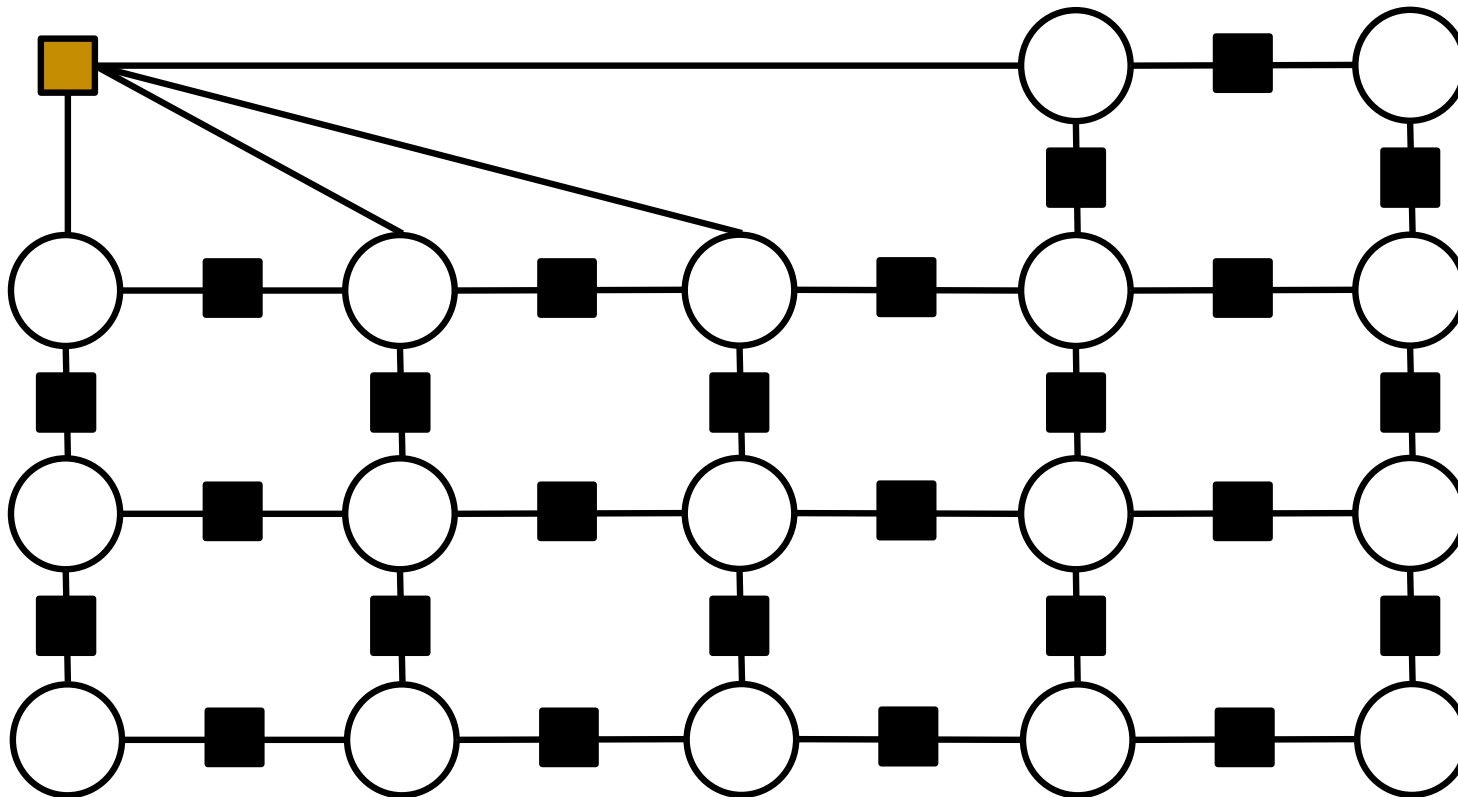
Grid CRF

- Suppose we want to image segmentation using a grid model
- What happens when we run variable elimination?



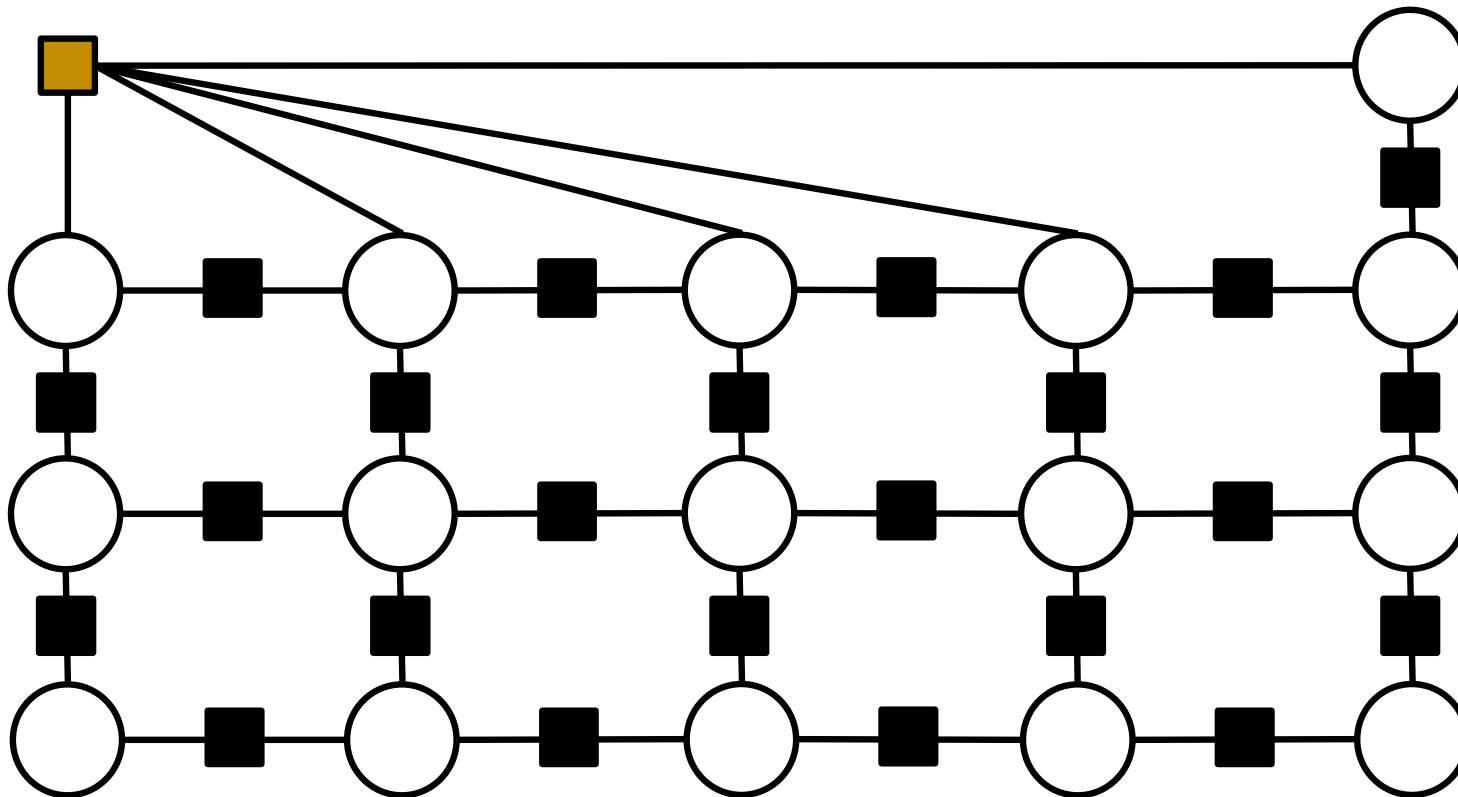
Grid CRF

- Suppose we want to image segmentation using a grid model
- What happens when we run variable elimination?



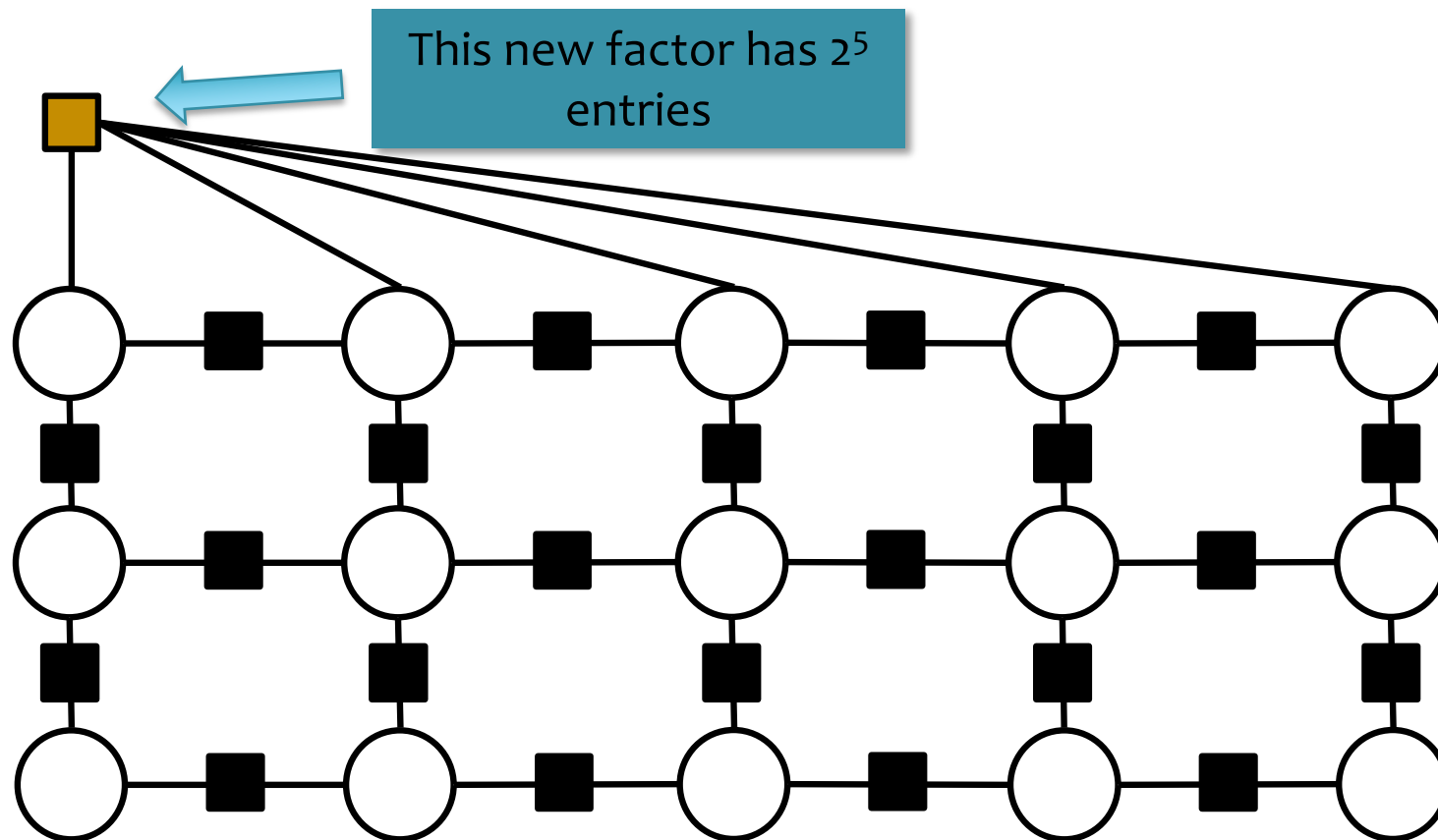
Grid CRF

- Suppose we want to image segmentation using a grid model
- What happens when we run variable elimination?



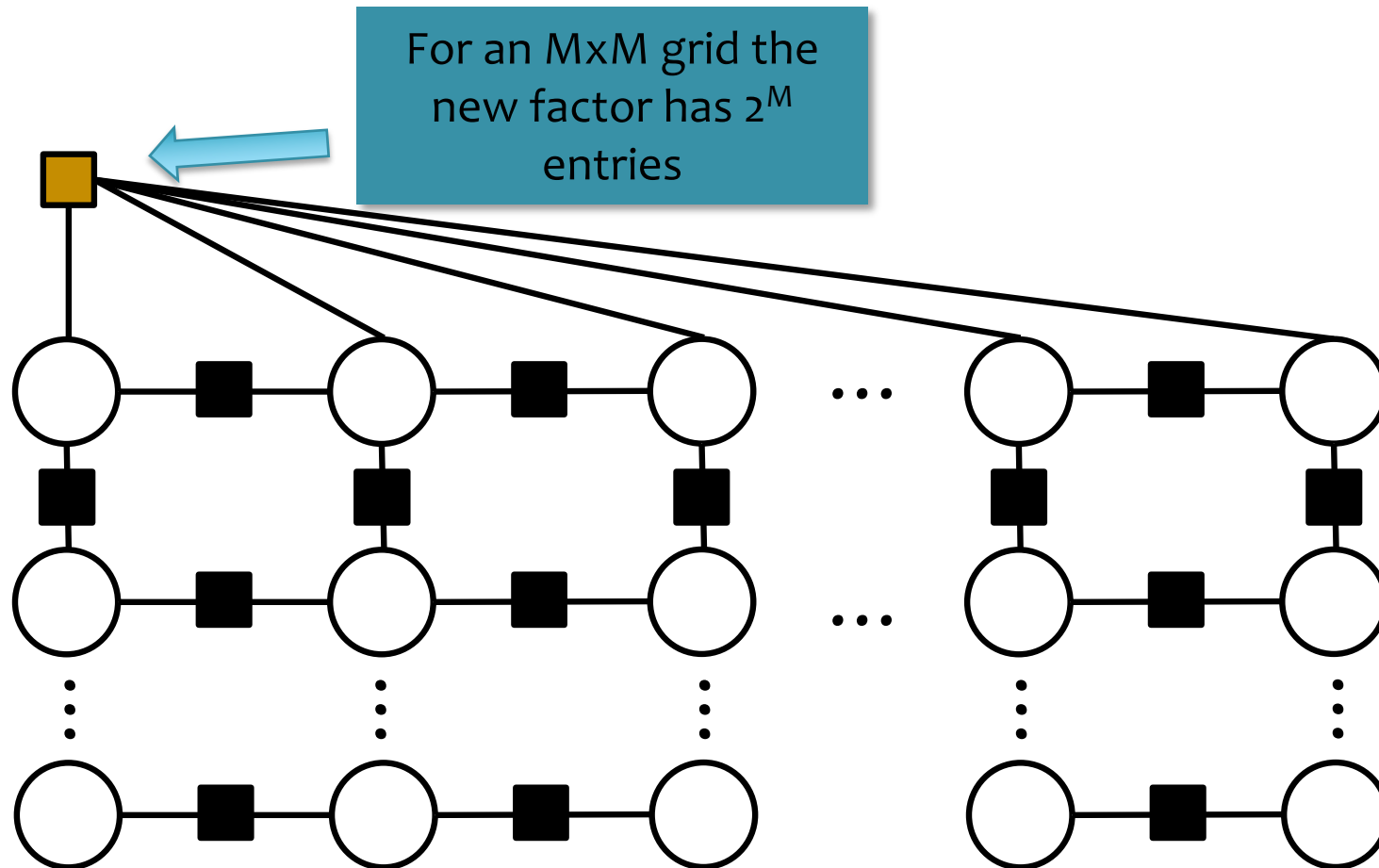
Grid CRF

- Suppose we want to image segmentation using a grid model
- What happens when we run variable elimination?



Grid CRF

- Suppose we want to image segmentation using a grid model
- What happens when we run variable elimination?



Grid CRF

- Suppose we want to image segmentation using a grid model
- What happens when we run variable elimination?

