# Exact Inference

Kayhan Batmanghelich

# Probabilistic Inference and Learning

- We now have compact representations of probability distributions:  **Graphical Models**

- A GM $M$ describes a unique probability distribution $P$

- Typical tasks:

  - Task 1: How do we answer **queries** about $P_M$, e.g., $P_M(X|Y)$ ?

    - We use **inference** as a name for the process of computing answers to such queries

  - Task 2: How do we estimate a **plausible model** $M$ from data $D$?

    By this  I mean the graph structure and/or the parameters

    i. We use **learning** as a name for the process of obtaining point estimate of $M$.

    ii. But for *Bayesian*, they seek $p(M|D)$, which is actually an **inference** problem.

    iii. When not all variables are observable, even computing point estimate of $M$ need to do **inference** to impute the *missing data*.

# Query 1: Likelihood

- Most of the queries one may ask involve **evidence**
  - Evidence **e** is an assignment of values to a set **E** variables in the domain
  - Without loss of generality $\mathbf{E} = \{ X_{k+1}, ..., X_n \}$

- Simplest query: compute probability of evidence

$$P(\mathbf{e}) = \sum_{x_1} \cdots \sum_{x_k} P(x_1, ..., x_k, \mathbf{e})$$

  - this is often referred to as computing the **likelihood** of **e**

# Query 2: Conditional Probability

- Often we are interested in the **conditional probability distribution** of a variable given the evidence

$$P(X \mid \mathbf{e}) = \frac{P(X, \mathbf{e})}{P(\mathbf{e})} = \frac{P(X, \mathbf{e})}{\sum_x P(X = x, \mathbf{e})}$$

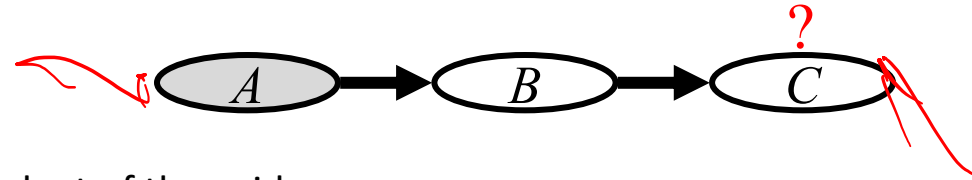  - this is the **_a posteriori_ belief** in $X$, given evidence $\mathbf{e}$

- **Marginalization:** the process of summing out the "unobserved" or "don't care" variables $z$ is called.

$$P(\mathbf{Y} \mid e) = \sum_z P(\mathbf{Y}, \mathbf{Z} = \mathbf{z} \mid \mathbf{e})$$

# Applications of a posteriori Belief

- **Prediction** : what is the probability of an outcome given the starting condition (eg $P(C|A)$ )



  - the query node is a descendent of the evidence

- **Diagnosis**: what is the probability of disease/fault given symptoms (eg $P(A|C)$ )



  - the query node an ancestor of the evidence

- ***Note***: The directionality of information flow between variables is not restricted by the directionality of the edges in a GM.

- You will see more application during **Learning** under partial observation

# Query 3: Most Probable Assignment

$$\max_{x_1 \cdots x_4} P(x_1, \cdots x_4 | e)$$

- In this query we want to find the most probable joint assignment (MPA) for *some* variables of interest

$$e = \text{Cat sit on the Chair}$$

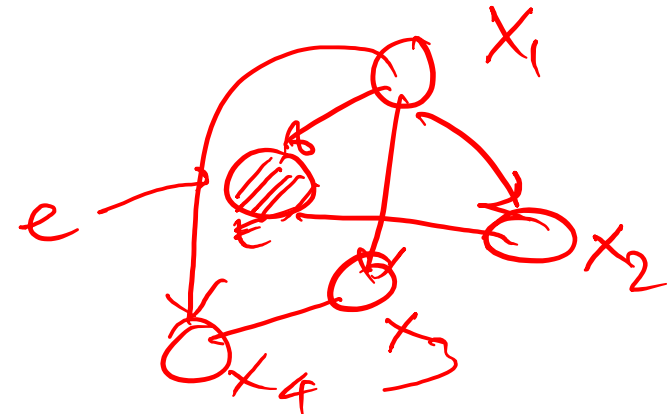- Such reasoning is usually performed under some given evidence e, and ignoring (the values of) other variables $z$ :

$$\mathrm{MPA}(\mathbf{Y} \,|\, \mathbf{e}) = \arg\max_{\mathbf{y} \in \mathcal{Y}} P(\mathbf{y} \,|\, \mathbf{e}) = \arg\max_{\mathbf{y} \in \mathcal{Y}} \sum_{\mathbf{z}} P(\mathbf{y}, \mathbf{z} \,|\, \mathbf{e})$$

  - this is the maximum *a posteriori* configuration of y.

$$P(X_1, X_2, X_3, X_4 | E)$$

# Applications of MPA

- ## Classification
  - find most likely label, given the evidence

- ## Explanation
  - what is the most likely scenario, given the evidence

Cautionary note:

- The MPA of a variable depends on its "context"---the set of variables been jointly queried
- Example:
  - MPA of $Y_1$?
  - MPA of $(Y_1, Y_2)$?

$\max \quad P(Y_1)$

$\sum_{Y_2} P(Y_1, Y_2)$

| $Y_1$ | $Y_2$ | $P(Y_1, Y_2)$ |
|-------|-------|---------------|
| 0 | 0 | 0.35 |
| 0 | 1 | 0.05 |
| 1 | 0 | 0.3 |
| 1 | 1 | 0.3 |

0.4

0.6

# Complexity of Inference

**Thm**:

Computing $P(X = \mathrm{x} \mid e)$ in a GM is NP-hard

- Hardness does not mean we cannot solve inference. It simply says that there exist difficult inference problems. It depends on the structure.

  - It implies that we cannot find a general procedure that works efficiently for arbitrary GMs
  - For particular families of GMs, we can have provably efficient procedures

# Landscape of inference algorithms

- **Exact inference algorithms**
  - The elimination algorithm
  - Message-passing algorithm (sum-product, belief propagation)
  - The junction tree algorithms

*Classical*

- **Approximate inference techniques**
  - Stochastic simulation / sampling methods
  - Markov chain Monte Carlo methods
  - Variational algorithms

# Variable Elimination

# Marginalization and Elimination

$E = \{e_1, \ldots, e_k\}$

$E = e_1$

- A signal transduction pathway:



$A \rightarrow B \rightarrow C \rightarrow D \rightarrow E$

What is the likelihood that protein E is active?

- Query: $P(e)$

$$P(e) = \sum_d \sum_c \sum_b \sum_a P(a,b,c,d,e)$$

a naïve summation needs to enumerate over an exponential number of terms

- By chain decomposition, we get

$$= \sum_d \sum_c \sum_b \sum_a P(a)P(b|a)P(c|b)P(d|c)P(e|d)$$

# Elimination on Chains



- Rearranging terms (First sum $a$, then $b$, ...)

$$P(e) = \sum_d \sum_c \sum_b \boxed{\sum_a P(a)P(b \mid a)} P(c \mid b)P(d \mid c)P(e \mid d)$$

$$= \sum_d \sum_c \sum_b P(c \mid b)P(d \mid c)P(e \mid d)\sum_a P(a)P(b \mid a)$$

# Elimination on Chains



- Now we can perform innermost summation

$$P(e) = \sum_d \sum_c \sum_b P(c \mid b) P(d \mid c) P(e \mid d) \left( \sum_a P(a) P(b \mid a) \right)$$

$$= \sum_d \sum_c \sum_b P(c \mid b) P(d \mid c) P(e \mid d) \, p(b)$$

- This summation "eliminates" one variable from our summation argument at a "local cost".

# Elimination in Chains



- Rearranging and then summing again, we get

$$P(e) = \sum_d \sum_c \sum_b P(c \mid b) P(d \mid c) P(e \mid d) p(b)$$

$$= \sum_d \sum_c P(d \mid c) P(e \mid d) \sum_b P(c \mid b) p(b)$$

$$= \sum_d \sum_c P(d \mid c) P(e \mid d) p(c)$$

# Elimination in Chains

$A \left( M_A \left( M_B \left( M_e \ M_D \right) \right) \right) \cdots$



- Eliminate nodes one by one all the way to the end, we get

$$P(e) = \sum_d P(e \mid d) p(d)$$

- Complexity:
  - A clever elimination takes $O(nk^2)$ *versus* $O(k^n)$ for the naïve approach.

# Associativity of matrix multiplication



$$p(x_{t+1} = i | x_t = j) = M_{ij}$$

$$\begin{bmatrix} 0.7 & 0.5 & 0 \\ 0.3 & 0.3 & 0.5 \\ 0 & 0.2 & 0.5 \end{bmatrix}$$

Room 1    Room 2    Room 3

$$p(x_1, \cdots, x_T) = p(x_1) \prod_{t=1}^{T-1} p(x_{t+1} | x_t)$$
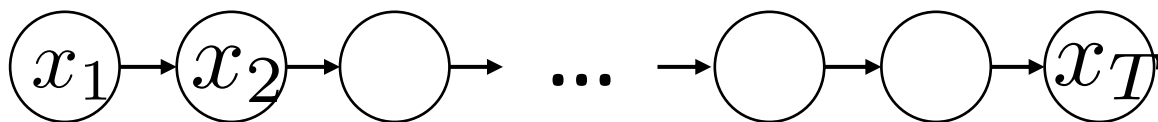
# Associativity of matrix multiplication

$P = P(x_2|x_1)P(x_1)$

$P(x) = v$   $v \in \mathbb{R}^3$

$Mv = \lambda v$

$MP = P$

$$p(x_{t+1} = i | x_t = j) = M_{ij}$$

$$\begin{bmatrix} 0.7 & 0.5 & 0 \\ 0.3 & 0.3 & 0.5 \\ 0 & 0.2 & 0.5 \end{bmatrix}$$

Room 1    Room 2    Room 3
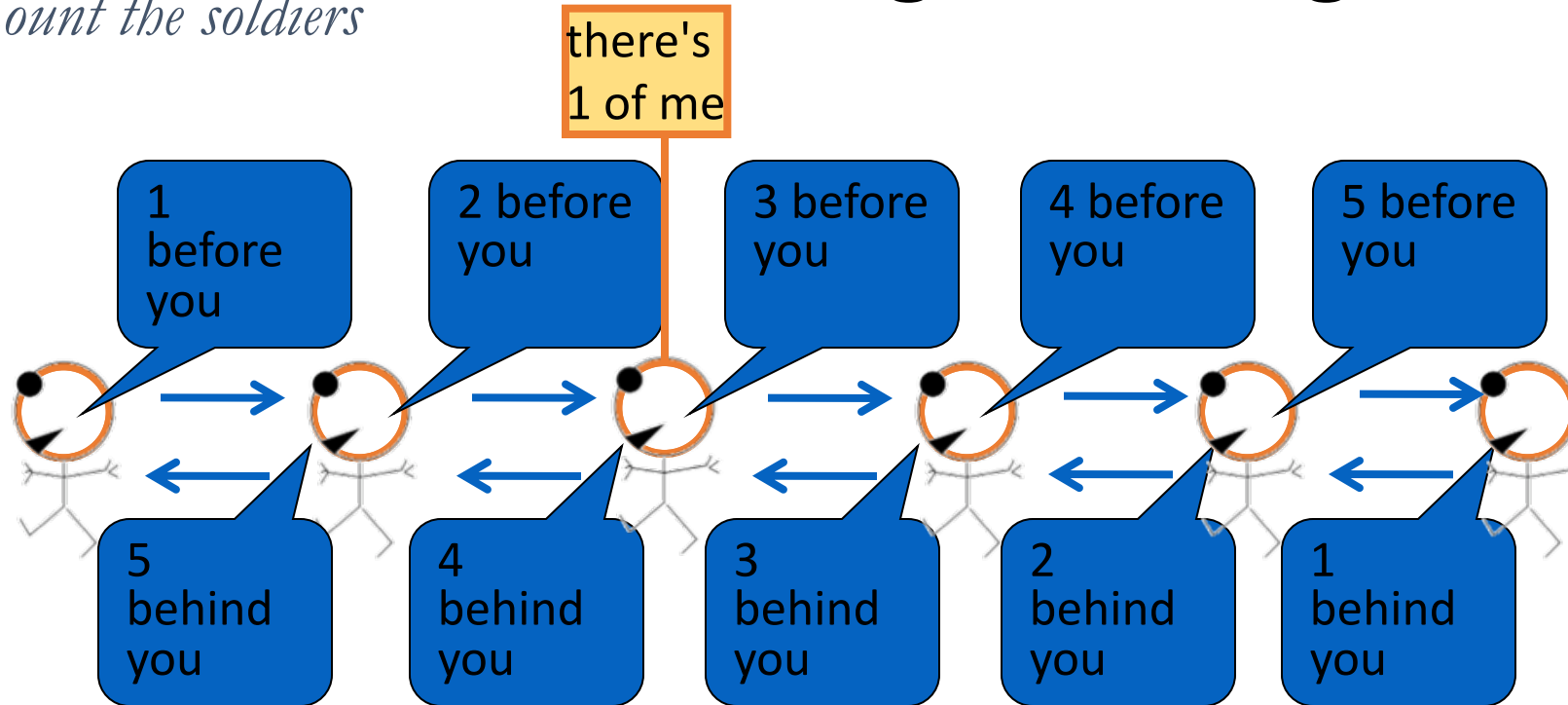
$$x_1 \rightarrow x_2 \rightarrow \bigcirc \rightarrow \dots \rightarrow \bigcirc \rightarrow \bigcirc \rightarrow x_T$$

$$p(x_5 = i | x_1 = 1) = \sum_{x_4, x_3, x_2} p(x_5|x_4)p(x_4|x_3)p(x_3|x_2)p(x_2|x_1 = 1)$$

$$= [M^4 v]_i$$

$M_{\infty}P_{\infty} = P_{\infty}$

$T \rightarrow \infty$

# Great Ideas in ML: Message Passing

*Count the soldiers*

# Great Ideas in ML: Message Passing

*Count the soldiers*

there's 1 of me

Belief:
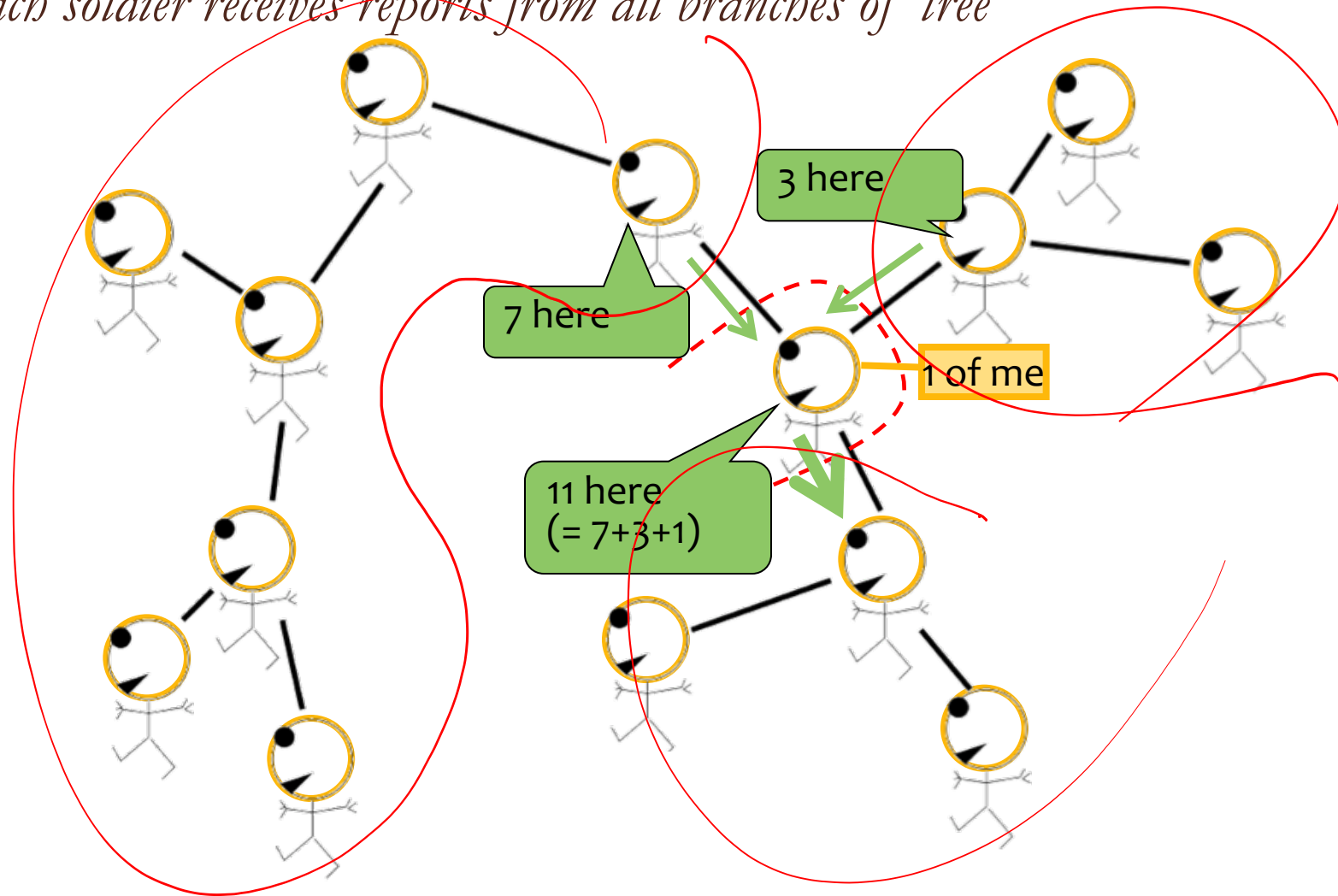Must be
2 +1 +3 = 6
of us

2 before you

only see my incoming messages

3 behind you

# Great Ideas in ML: Message Passing

*Count the soldiers*

there's 1 of me

Belief: Must be **1** + **1** + **4** = 6 of us

ef: t be **1** + **3** = 6 of us

1 before you

only see my incoming messages

4 behind you

# What about a general DAG?

# Great Ideas in ML: Message Passing

*Each soldier receives reports from all branches of tree*

3 here

7 here

1 of me

11 here
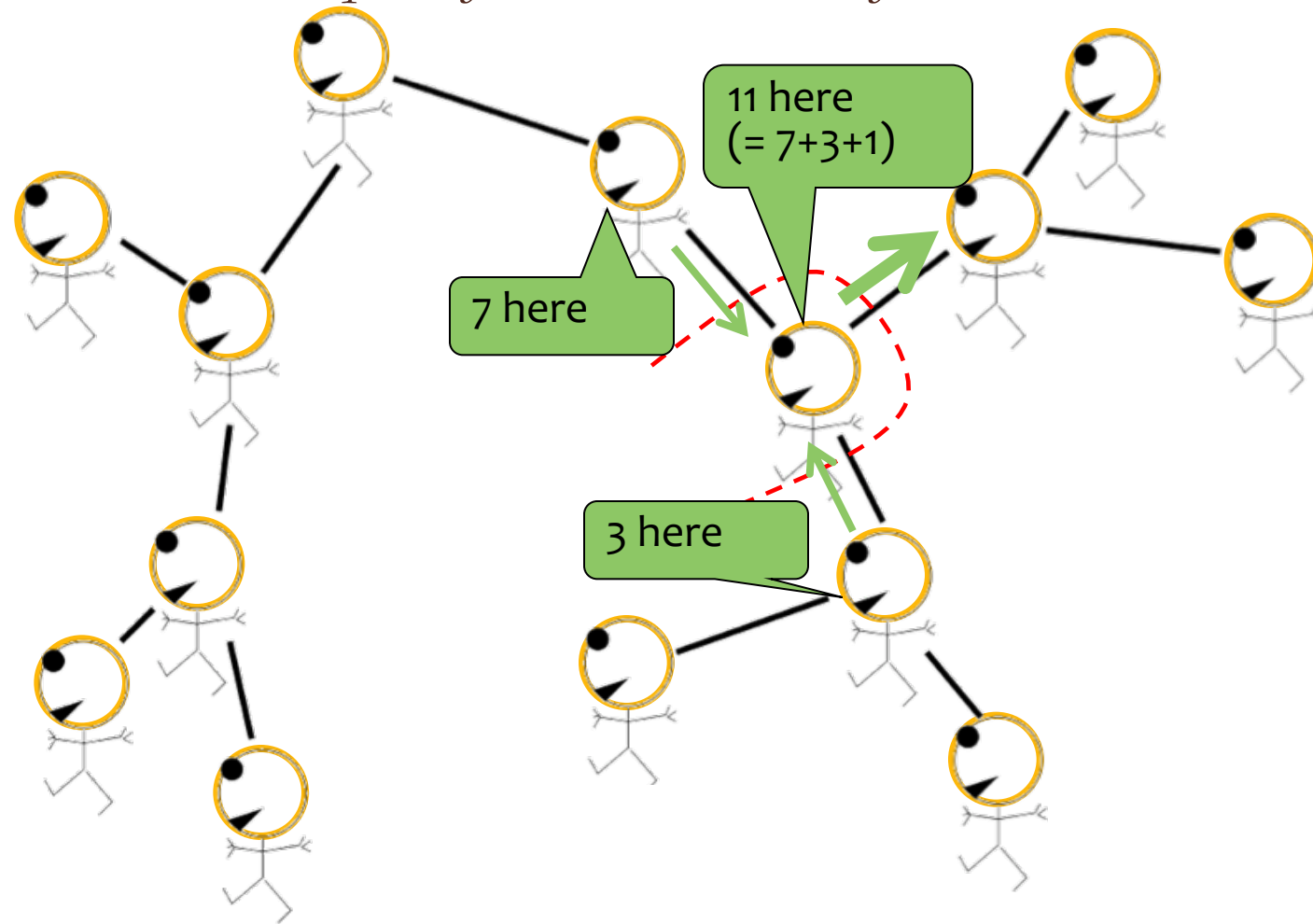(= 7+3+1)

adapted from MacKay (2003) textbook

# Great Ideas in ML: Message Passing

*Each soldier receives reports from all branches of tree*

# Great Ideas in ML: Message Passing

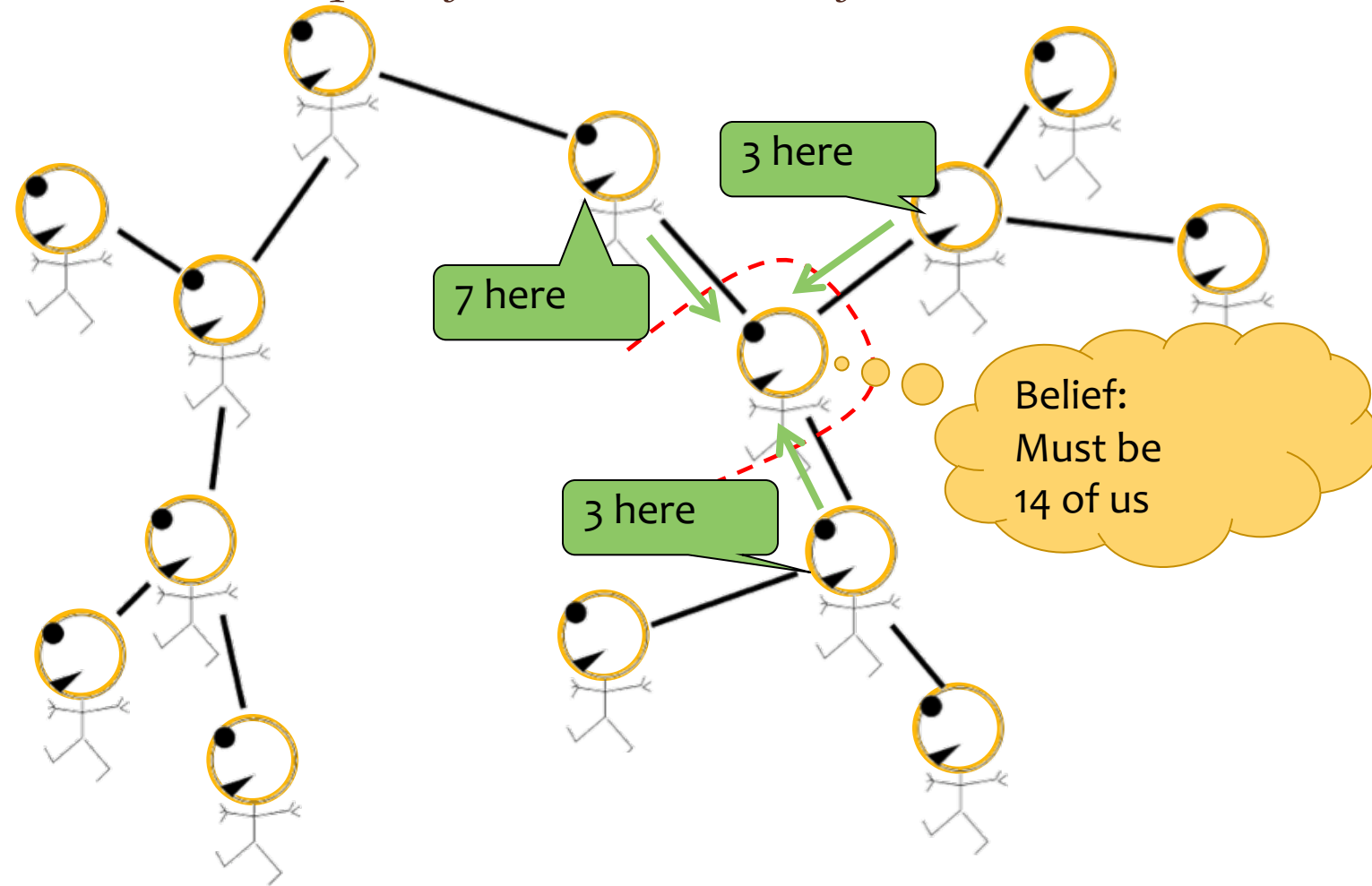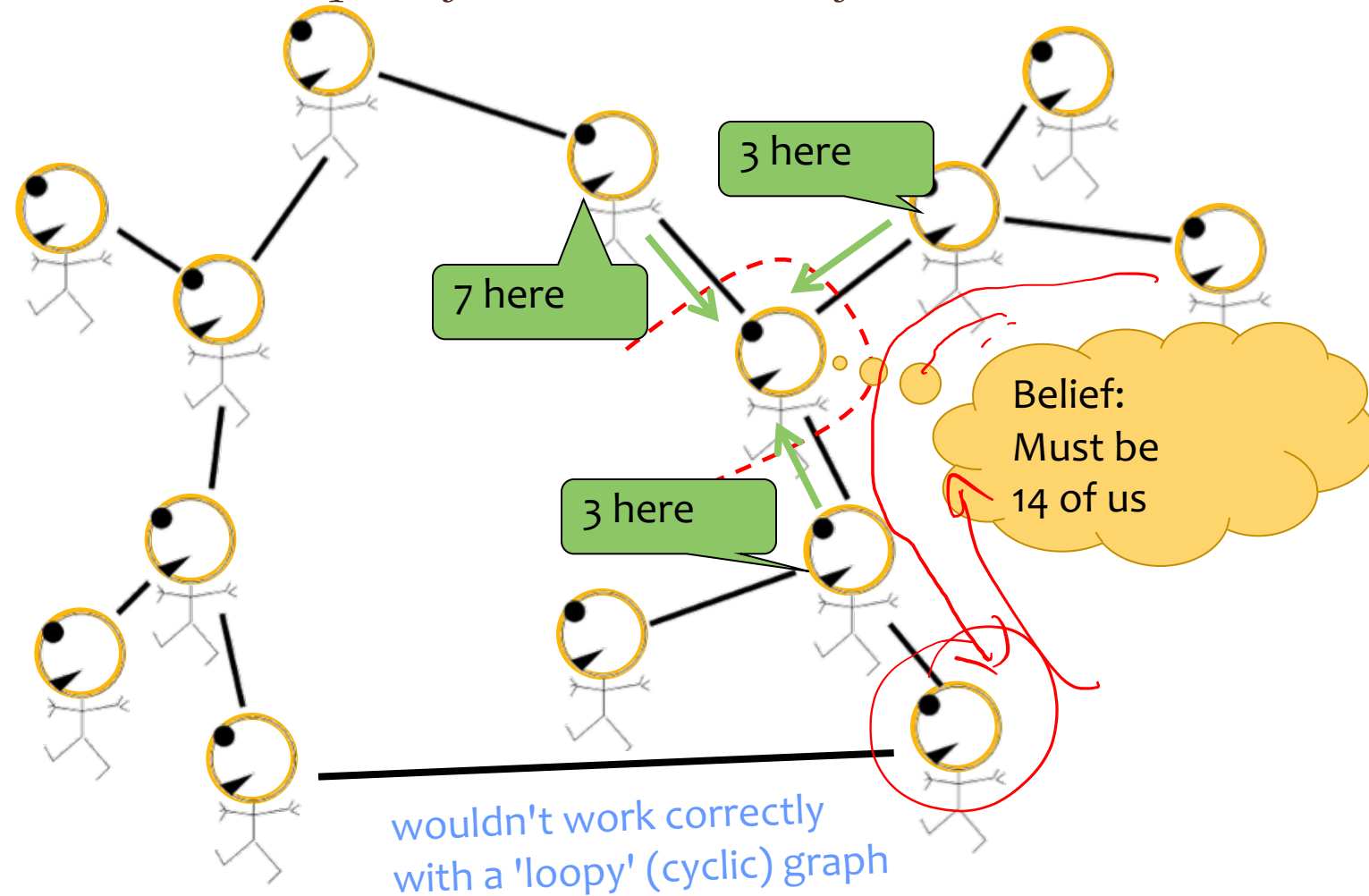*Each soldier receives reports from all branches of tree*

# Great Ideas in ML: Message Passing

*Each soldier receives reports from all branches of tree*



adapted from MacKay (2003) textbook

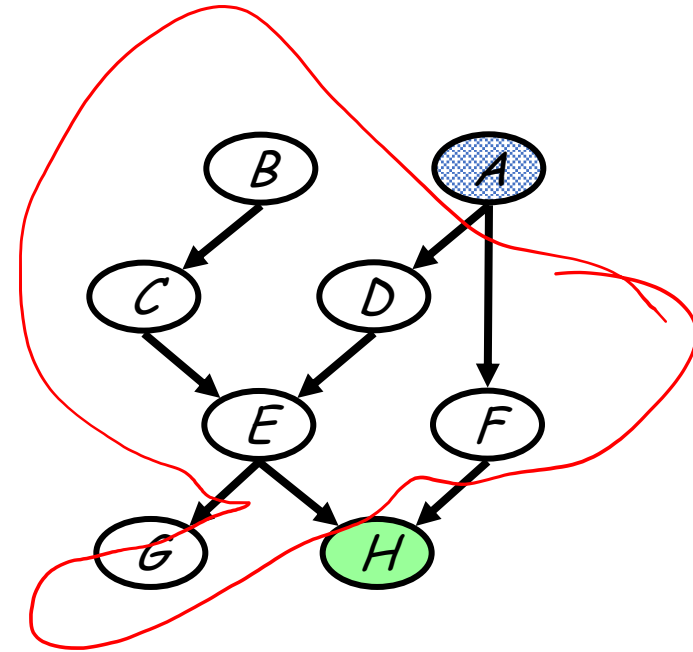# Great Ideas in ML: Message Passing

*Each soldier receives reports from all branches of tree*



3 here

7 here

3 here

Belief:
Must be
14 of us

wouldn't work correctly
with a 'loopy' (cyclic) graph

adapted from MacKay (2003) textbook

# A more complex network

A food web

$P(A \mid H)$



What is the probability that hawks are leaving given that the grass condition is poor?

# Example: Variable Elimination

- Query: $P(A \mid h)$
  - Need to eliminate: $B,C,D,E,F,G,H$

- Initial factors:

$$P(a)P(b)P(c \mid b)P(d \mid a)P(e \mid c,d)P(f \mid a)P(g \mid e)P(h \mid e,f)$$

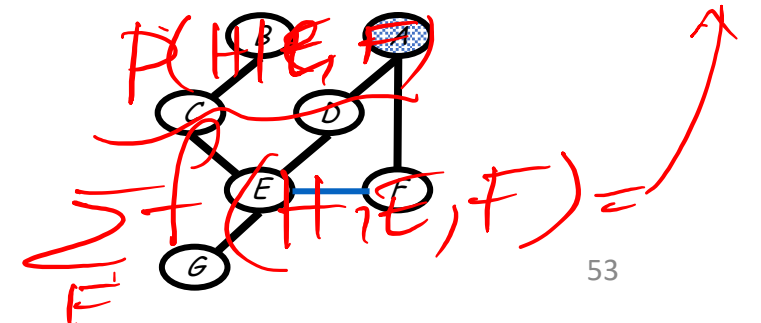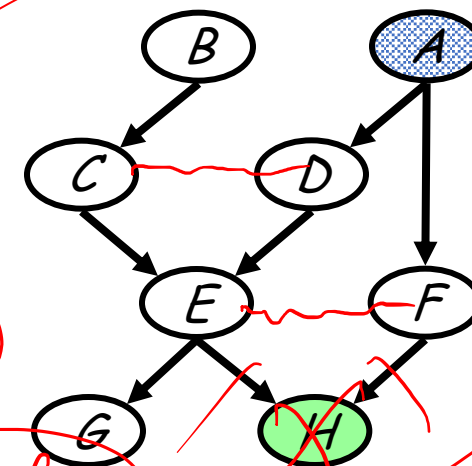- Choose an elimination order: $H,G,F,E,D,C,B$

- Step 1:
  - **Conditioning** (fix the evidence node (i.e., $h$) on its observed value (i.e., $\tilde{h}$)):

$$m_h(e,f) = p(h = \tilde{h} \mid e,f)$$

  - This step is isomorphic to a marginalization step:

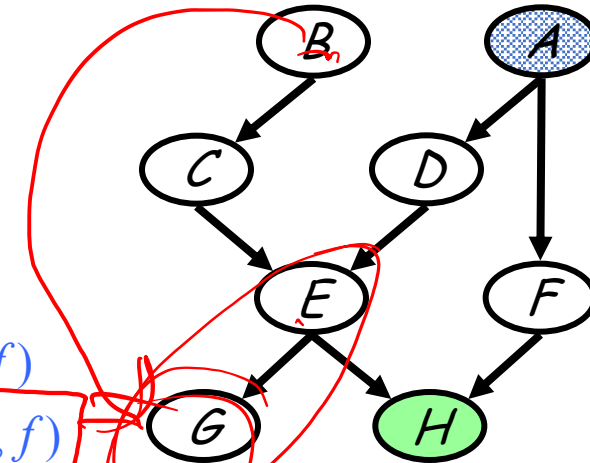$$m_h(e,f) = \sum_h p(h \mid e,f)\delta(h = \tilde{h})$$

# Example: Variable Elimination

- Query: $P(A \mid h)$
  - Need to eliminate: $B, C, D, E, F, G$

- Initial factors:

$$P(a)P(b)P(c \mid b)P(d \mid a)P(e \mid c,d)P(f \mid a)P(g \mid e)P(h \mid e,f)$$
$$\Rightarrow P(a)P(b)P(c \mid b)P(d \mid a)P(e \mid c,d)P(f \mid a)P(g \mid e)m_h(e,f)$$

- Step 2: Eliminate $G$

$$m_g(e) = \sum_g p(g \mid e) = 1$$

# Example: Variable Elimination

- Query: $P(A \mid h)$
  - Need to eliminate: $B, C, D, E, F$

- Initial factors:

$$P(a)P(b)P(c \mid b)P(d \mid a)P(e \mid c, d)P(f \mid a)P(g \mid e)P(h \mid e, f)$$
$$\Rightarrow P(a)P(b)P(c \mid b)P(d \mid a)P(e \mid c, d)P(f \mid a)P(g \mid e)m_h(e, f)$$
$$\Rightarrow P(a)P(b)P(c \mid b)P(d \mid a)P(e \mid c, d)P(f \mid a)m_h(e, f)$$

$$m_f(e, a)$$

- Step 3: Eliminate $F$

$$m_f(e, a) = \sum_f p(f \mid a)m_h(e, f)$$

# Example: Variable Elimination

- Query: $P(A \mid h)$
  - Need to eliminate: $B, C, D, E$

- Initial factors:

$$P(a)P(b)P(c \mid b)P(d \mid a)P(e \mid c,d)P(f \mid a)P(g \mid e)P(h \mid e,f)$$
$$\Rightarrow P(a)P(b)P(c \mid b)P(d \mid a)P(e \mid c,d)P(f \mid a)P(g \mid e)m_h(e,f)$$
$$\Rightarrow P(a)P(b)P(c \mid b)P(d \mid a)P(e \mid c,d)P(f \mid a)m_h(e,f)$$
$$\Rightarrow P(a)P(b)P(c \mid b)P(d \mid a)P(e \mid c,d)m_f(a,e)$$

$m(a,c,d)$
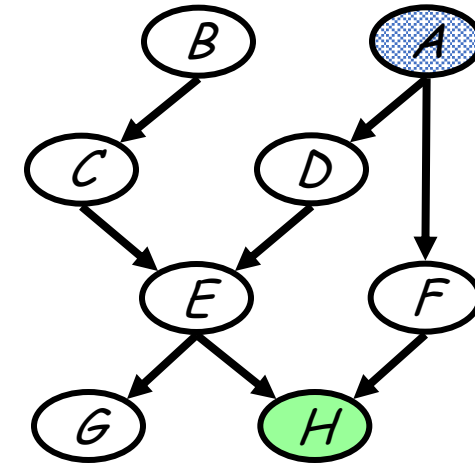
- Step 4: Eliminate $E$

$$m_e(a,c,d) = \sum_e p(e \mid c,d)m_f(a,e)$$

# Example: Variable Elimination

- Query: $P(A \mid h)$
  - Need to eliminate: $B, C, D$

- Initial factors:

$$P(a)P(b)P(c \mid b)P(d \mid a)P(e \mid c,d)P(f \mid a)P(g \mid e)P(h \mid e,f)$$

$$\Rightarrow P(a)P(b)P(c \mid b)P(d \mid a)P(e \mid c,d)P(f \mid a)P(g \mid e)m_h(e,f)$$

$$\Rightarrow P(a)P(b)P(c \mid b)P(d \mid a)P(e \mid c,d)P(f \mid a)m_h(e,f)$$

$$\Rightarrow P(a)P(b)P(c \mid b)P(d \mid a)P(e \mid c,d)m_f(a,e)$$

$$\Rightarrow P(a)P(b)P(c \mid b)P(d \mid a)m_e(a,c,d)$$

$$m_d(a,c)$$

- Step 5: Eliminate $D$

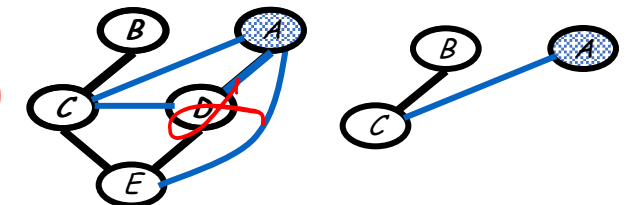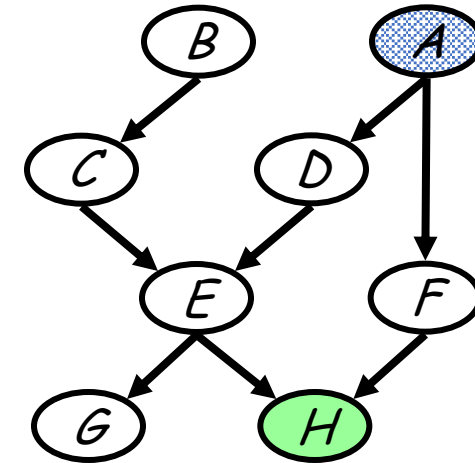$$m_d(a,c) = \sum_d p(d \mid a)m_e(a,c,d)$$

# Example: Variable Elimination

- Query: *P(A |h)*
  - Need to eliminate: *B,C*

- Initial factors:

$$P(a)P(b)P(c|d)P(d|a)P(e|c,d)P(f|a)P(g|e)P(h|e,f)$$
$$\Rightarrow P(a)P(b)P(c|d)P(d|a)P(e|c,d)P(f|a)P(g|e)m_h(e,f)$$
$$\Rightarrow P(a)P(b)P(c|d)P(d|a)P(e|c,d)P(f|a)m_h(e,f)$$
$$\Rightarrow P(a)P(b)P(c|d)P(d|a)P(e|c,d)m_f(a,e)$$
$$\Rightarrow P(a)P(b)P(c|d)P(d|a)m_e(a,c,d)$$
$$\Rightarrow P(a)P(b)P(c|d)m_d(a,c)$$

$$m_c(a,b)$$

- Step 6: Eliminate *C*

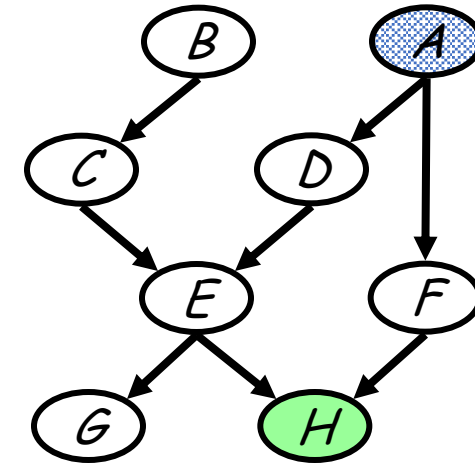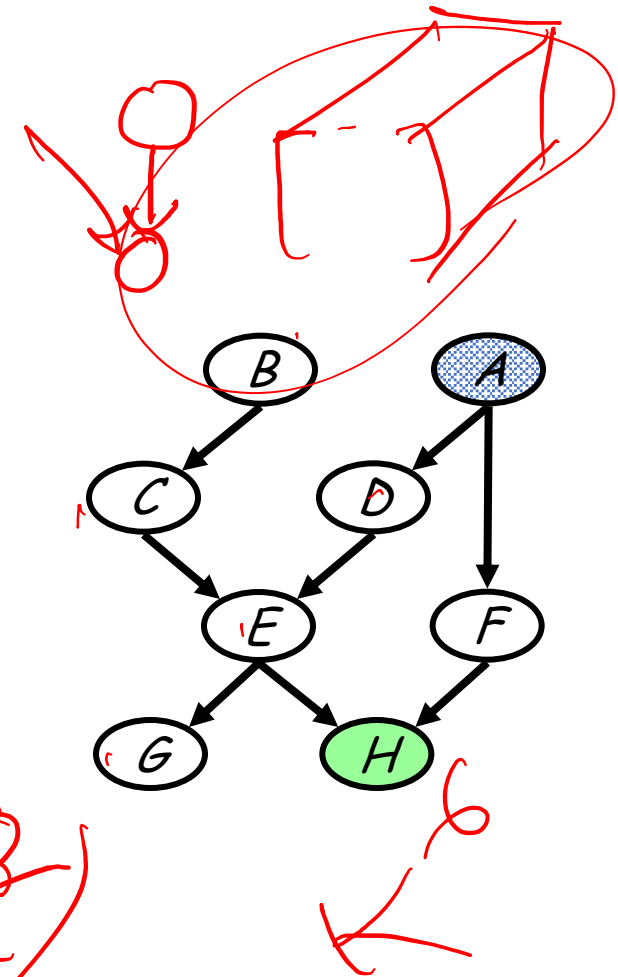$$m_c(a,b) = \sum_c p(c|b)m_d(a,c)$$

# Example: Variable Elimination

- Query: $P(A \mid h)$
  - Need to eliminate: $B$

- Initial factors:

$$P(a)P(b)P(c \mid d)P(d \mid a)P(e \mid c,d)P(f \mid a)P(g \mid e)P(h \mid e,f)$$
$$\Rightarrow P(a)P(b)P(c \mid d)P(d \mid a)P(e \mid c,d)P(f \mid a)P(g \mid e)m_h(e,f)$$
$$\Rightarrow P(a)P(b)P(c \mid d)P(d \mid a)P(e \mid c,d)P(f \mid a)m_h(e,f)$$
$$\Rightarrow P(a)P(b)P(c \mid d)P(d \mid a)P(e \mid c,d)m_f(a,e)$$
$$\Rightarrow P(a)P(b)P(c \mid d)P(d \mid a)m_e(a,c,d)$$
$$\Rightarrow P(a)P(b)P(c \mid d)m_d(a,c)$$
$$\Rightarrow P(a)P(b)m_c(a,b)$$
$$m_b(a)$$

- Step 7: Eliminate $B$

$$m_b(a) = \sum_b p(b)m_c(a,b)$$

# Example: Variable Elimination

- Query: $P(A \mid h)$
  - Need to eliminate: $B$

- Initial factors:

$$P(a)P(b)P(c \mid d)P(d \mid a)P(e \mid c,d)P(f \mid a)P(g \mid e)P(h \mid e,f)$$
$$\Rightarrow P(a)P(b)P(c \mid d)P(d \mid a)P(e \mid c,d)P(f \mid a)P(g \mid e)m_h(e,f)$$
$$\Rightarrow P(a)P(b)P(c \mid d)P(d \mid a)P(e \mid c,d)P(f \mid a)m_h(e,f)$$
$$\Rightarrow P(a)P(b)P(c \mid d)P(d \mid a)P(e \mid c,d)m_f(a,e)$$
$$\Rightarrow P(a)P(b)P(c \mid d)P(d \mid a)m_e(a,c,d)$$
$$\Rightarrow P(a)P(b)P(c \mid d)m_d(a,c)$$
$$\Rightarrow P(a)P(b)m_c(a,b)$$
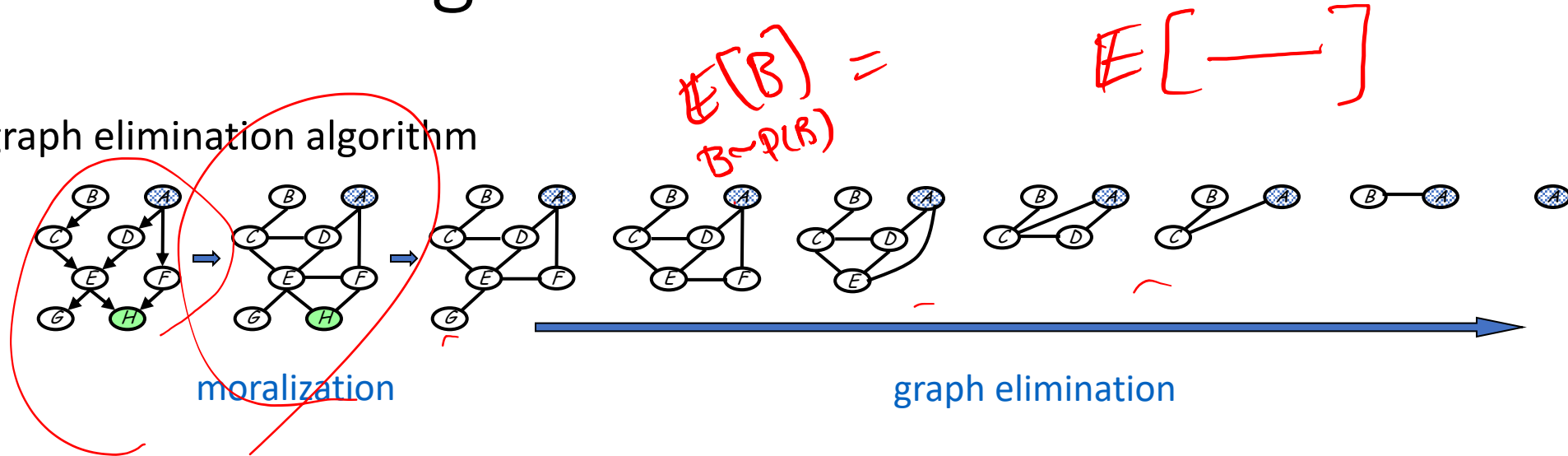$$\Rightarrow P(a)m_b(a)$$

- Step 8: Wrap-up

$$p(a,\widetilde{h}) = p(a)m_b(a), \quad p(\widetilde{h}) = \sum_a p(a)m_b(a)$$
$$\Rightarrow P(a \mid \widetilde{h}) = \frac{p(a)m_b(a)}{\sum_a p(a)m_b(a)}$$

# Understanding Variable Elimination

- A graph elimination algorithm



moralization                    graph elimination

$$\mathbb{E}[B] =$$
$$B \sim P(B)$$
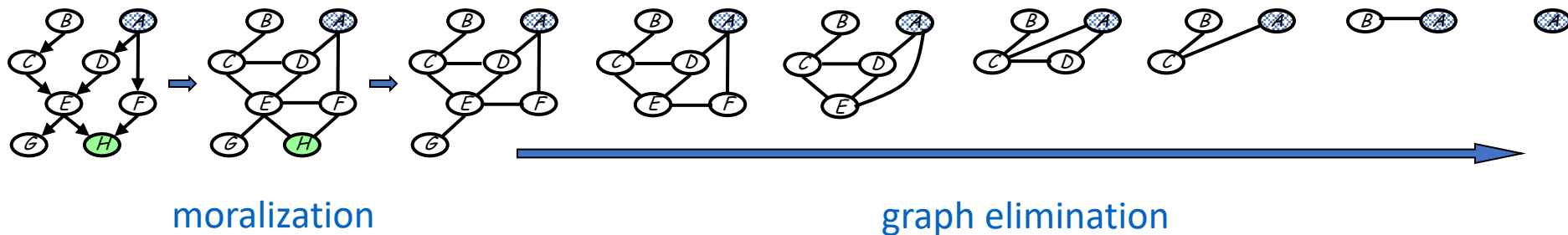
$$\mathbb{E}[\quad\text{---}\quad]$$
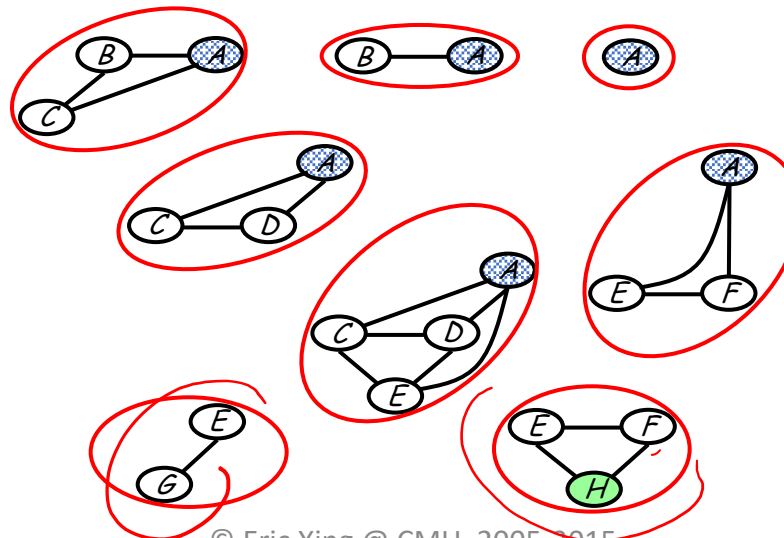
# Graph elimination

- Begin with the undirected GM or moralized BN

- Graph $G(V, E)$ and elimination ordering $I$

- Eliminate next node in the ordering $I$
  - Removing the node from the graph
  - Connecting the remaining neighbors of the nodes

- The reconstituted graph $G'(V, E')$
  - Retain the edges that were created during the elimination procedure
  - The graph-theoretic property: the factors resulted during variable elimination are captured by recording the elimination clique
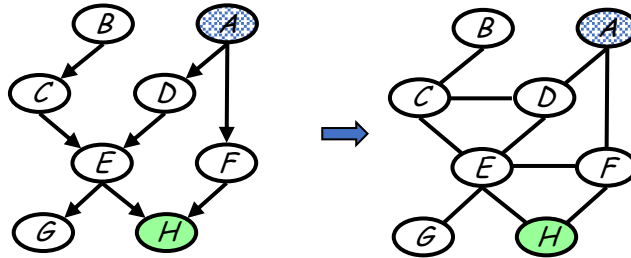
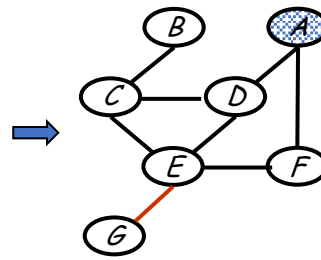# Understanding Variable Elimination

- A graph elimination algorithm



moralization          graph elimination

- Intermediate terms correspond to the cliques resulted from elimination

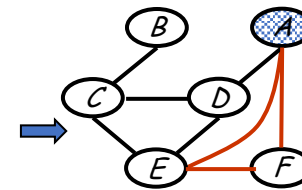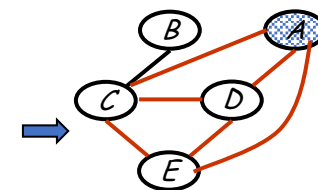# Elimination Cliques
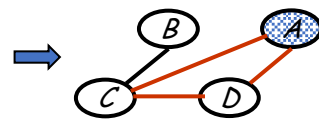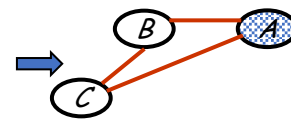


Messages sent: $m_h(e,f)$     $m_g(e)$     $m_f(e,a)$     $m_e(a,c,d)$

Messages sent: $m_d(a,c)$     $m_c(a,b)$     $m_b(a)$
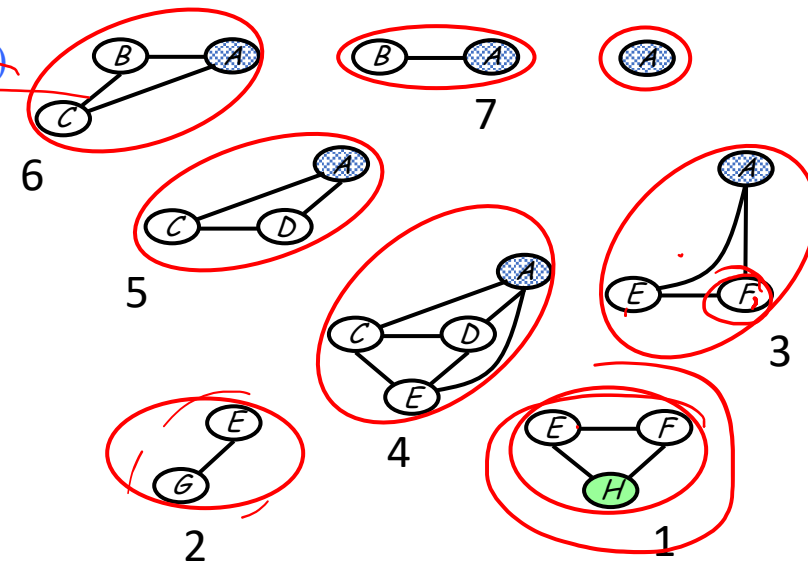
# Graph elimination and marginalization

- Induced dependency during marginalization vs. elimination clique
  - Summation <-> elimination
  - Intermediate term <-> elimination clique

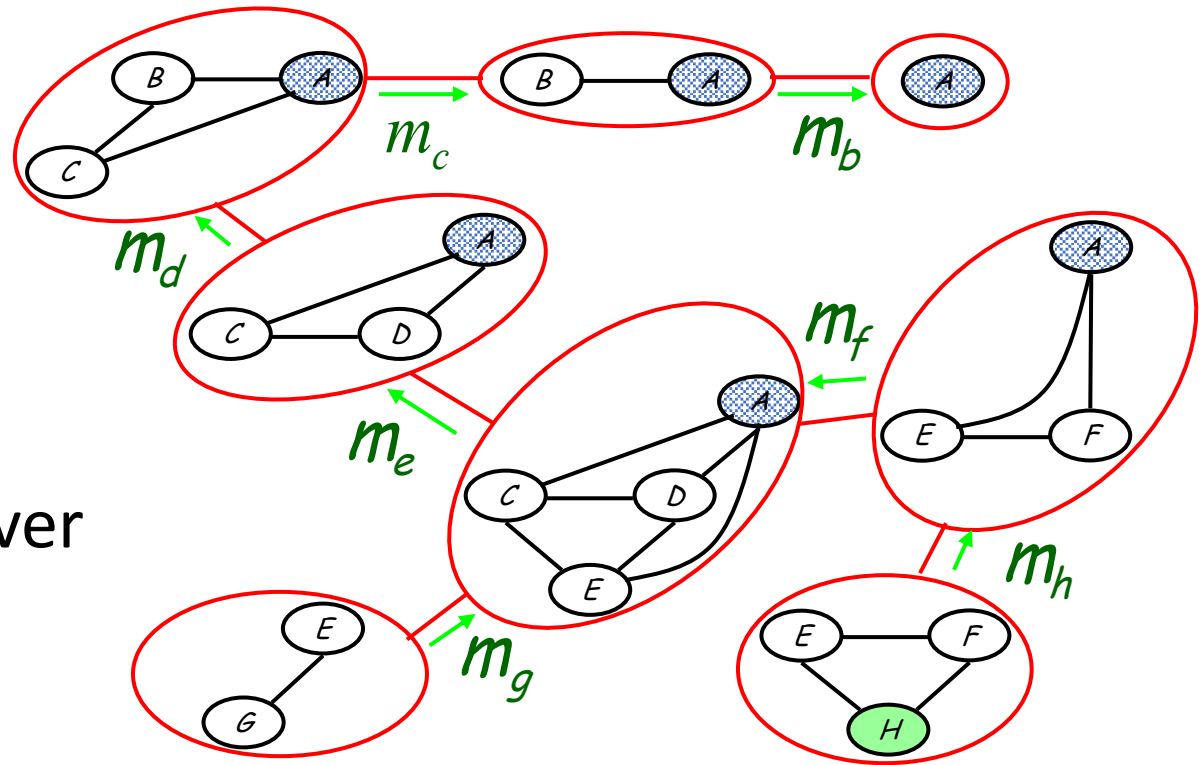$$P(a)P(b)P(c|d)P(d|a)P(e|c,d)P(f|a)P(g|e)P(h|e,f)$$

1. $\Rightarrow P(a)P(b)P(c|d)P(d|a)P(e|c,d)P(f|a)P(g|e)m_h(e,f)$
2. $\Rightarrow P(a)P(b)P(c|d)P(d|a)P(e|c,d)P(f|a)m_h(e,f)$
3. $\Rightarrow P(a)P(b)P(c|d)P(d|a)P(e|c,d)m_f(a,e)$
4. $\Rightarrow P(a)P(b)P(c|d)P(d|a)m_e(a,c,d)$
5. $\Rightarrow P(a)P(b)P(c|d)m_d(a,c)$
6. $\Rightarrow P(a)P(b)m_c(a,b)$
7. $\Rightarrow P(a)m_b(a)$

# A clique tree

**Message** from one C1 to C2:
**Multiply** all incoming messages
with the local factor and **sum** over
variables that are not shared

$$m_e(a,c,d)$$
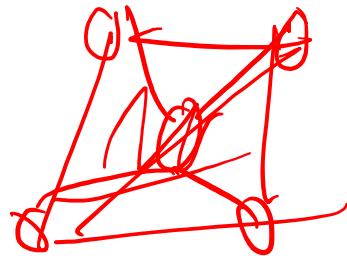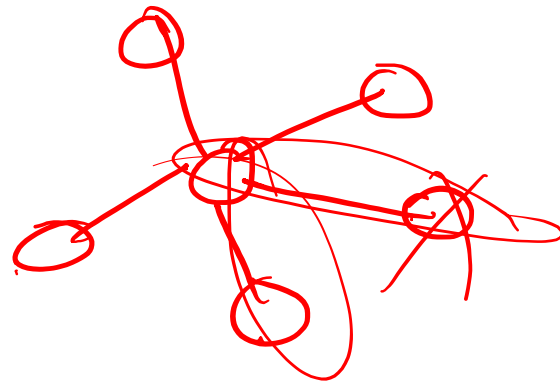$$= \sum_e p(e \mid c,d) m_g(e) m_f(a,e)$$

# Complexity

- The overall complexity is determined by the number of the largest elimination clique
  - What is the largest elimination clique? – a pure graph theoretic question

  - **Tree-width** $k$: one less than the smallest achievable value of the cardinality of the largest elimination clique, ranging over all possible elimination ordering

  - "good" elimination orderings lead to **small cliques** and hence reduce complexity (what will happen if we eliminate "e" first in the above graph?)

  - Find the best elimination ordering of a graph --- NP-hard
  - → Inference is NP-hard

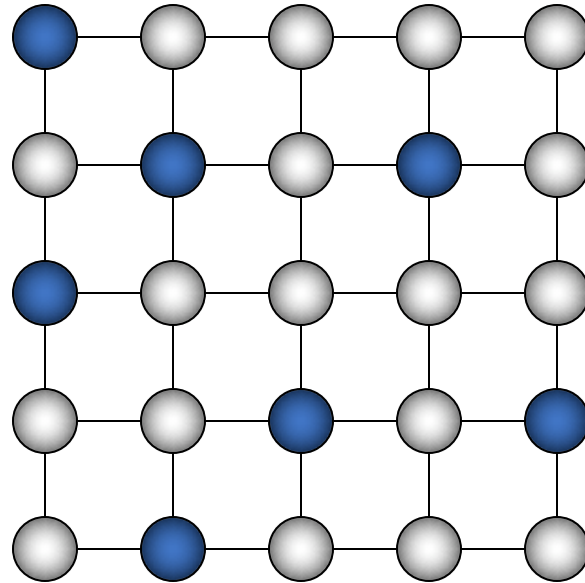  - But there often exist "obvious" optimal or near-opt elimination ordering
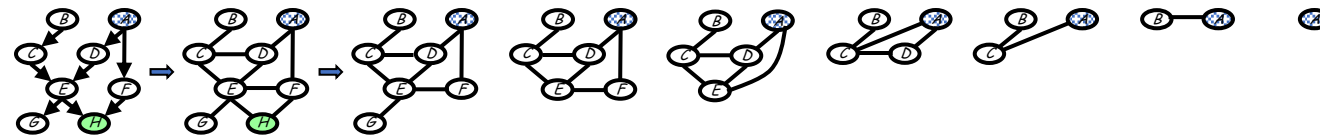
# Examples
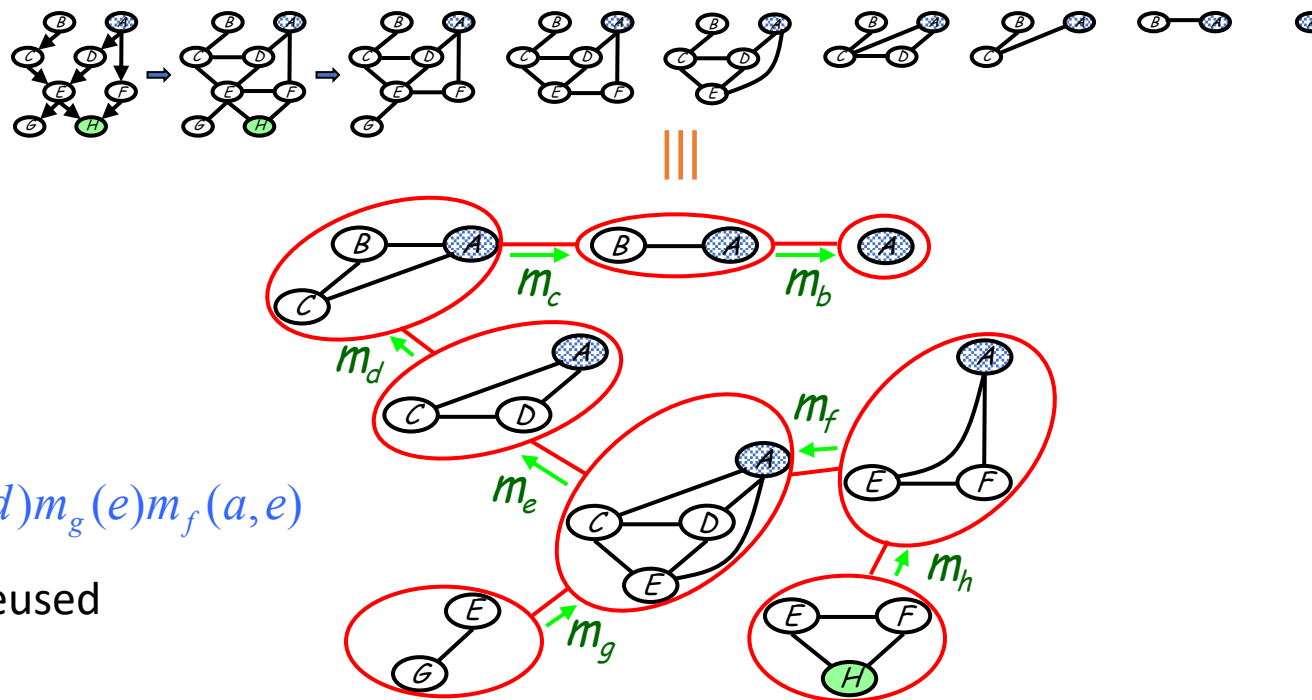
- Star

- Tree

# More example: Ising model

# Limitation of Procedure Elimination

- Limitation

# From Elimination to Message Passing

- Our algorithm so far answers only one query (e.g., on one node), do we need to do a complete elimination for every such query?

- Elimination ≡ message passing on a **clique tree**



$$m_e(a,c,d)$$
$$= \sum_e p(e|c,d)m_g(e)m_f(a,e)$$

- Messages can be reused

# From Elimination to Message Passing

- Our algorithm so far answers only one query (e.g., on one node), do we need to do a complete elimination for every such query?
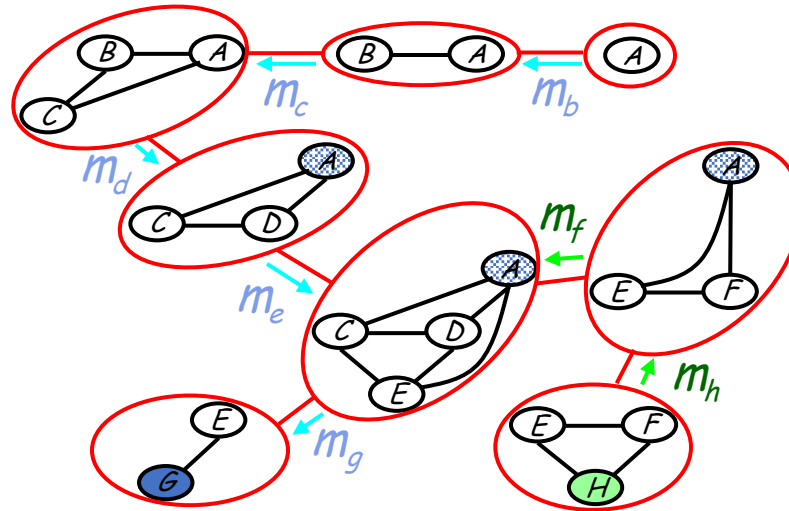
- Elimination ≡ message passing on a **clique tree**
  - **Another query …**



- Messages $m_f$ and $m_h$ are reused, others need to be recomputed

# Summary

- The simple Eliminate algorithm captures the key algorithmic Operation underlying probabilistic inference:

    --- That of taking a sum over product of potential functions

- What can we say about the overall computational complexity of the algorithm? In particular, how can we control the "size" of the summands that appear in the sequence of summation operation.

- The computational complexity of the Eliminate algorithm can be reduced to purely graph-theoretic considerations.

- This graph interpretation will also provide hints about how to design improved inference algorithm that overcome the limitation of Eliminate.