

# CS 581 – ADVANCED ARTIFICIAL INTELLIGENCE

## TOPIC: BAYESIAN NETWORKS



**Mustafa Bilgic**



<http://www.cs.iit.edu/~mbilgic>



<https://twitter.com/bilgicm>

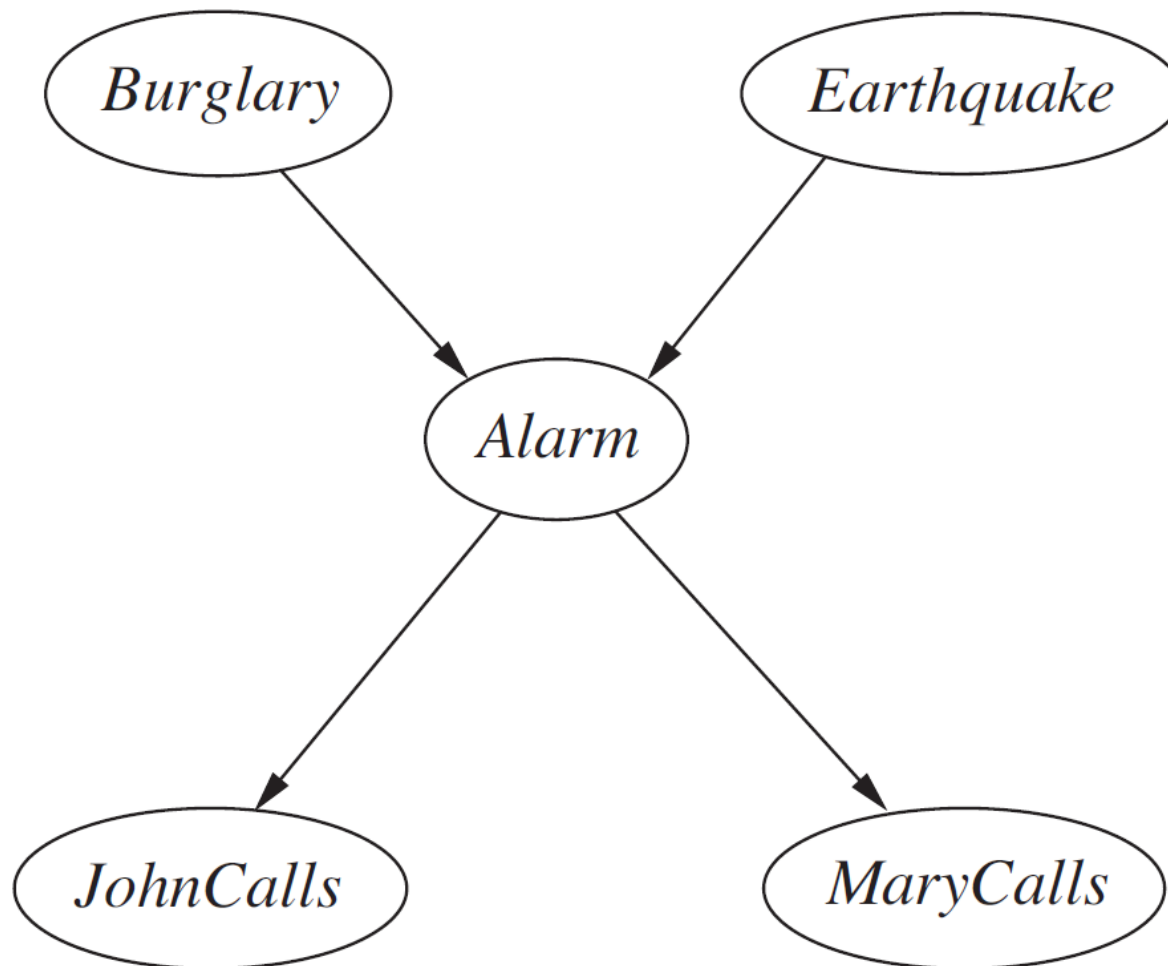
# MOTIVATION

- Efficient, intuitive, and modular representation of probability distributions
  - Represent joint and conditional distributions
- Structured and efficient inference
  - Answer probability and MAP queries
- BN structure represents correlation but can be used to answer causality questions under certain conditions

# AN EXAMPLE

- Five binary variables
  - Earthquake, Burglary, Alarm, MaryCalls, JohnCalls
- Assume the following
  - E and B are uncorrelated
  - E and M are related only through A; similarly, E, J, and A
  - B and M are related only through A; similarly, B, J, and A
  - M and J are directly related through A; undirectly related through E and B; otherwise, M and J are unrelated
- One approach
  - Represent and estimate the full joint  $P(E, B, A, M, J)$ 
    - How many independent parameters?
    - What can you tell about the relationships between the variables?
- Alternative approach
  - Bayesian network (next slide)

# BURGLARY EXAMPLE



# POSSIBLE QUERIES

- $P(B \mid J = \text{true})$
- $P(B \mid M = \text{true}, J = \text{true})$
- $P(M \mid B = \text{true})$
- $P(M \mid B = \text{false})$
- $P(M, J \mid B = \text{true})$
- $P(M \mid J = \text{true})$
- ...

# WE'LL COVER

- Bayesian networks (in detail)
  - [https://en.wikipedia.org/wiki/Bayesian\\_network](https://en.wikipedia.org/wiki/Bayesian_network)
- Hidden Markov Models (in detail)
  - [https://en.wikipedia.org/wiki/Hidden\\_Markov\\_model](https://en.wikipedia.org/wiki/Hidden_Markov_model)
- Dynamic Bayesian networks (brief)
  - [https://en.wikipedia.org/wiki/Dynamic\\_Bayesian\\_network](https://en.wikipedia.org/wiki/Dynamic_Bayesian_network)
- Influence diagrams (in detail)
  - [https://en.wikipedia.org/wiki/Influence\\_diagram](https://en.wikipedia.org/wiki/Influence_diagram)
- Causal networks (brief)

# BAYESIAN NETWORKS

- Random variables = nodes
- Direct relationships = directed edges
- BNs capture independencies
  - More compact than full joint representation
- Graphs provide
  - Graph theory / efficient reasoning
  - Intuition

# DIRECTED GRAPHS

- A **graph** consists of **nodes** and **edges**
- **Nodes:**  $\mathcal{X} = \{X_1, X_2, \dots, X_n\}$
- **Undirected Edge:**  $X_i - X_j$
- **Directed Edge:**  $X_i \rightarrow X_j$
- A graph is **directed** if its *all* edges are directed



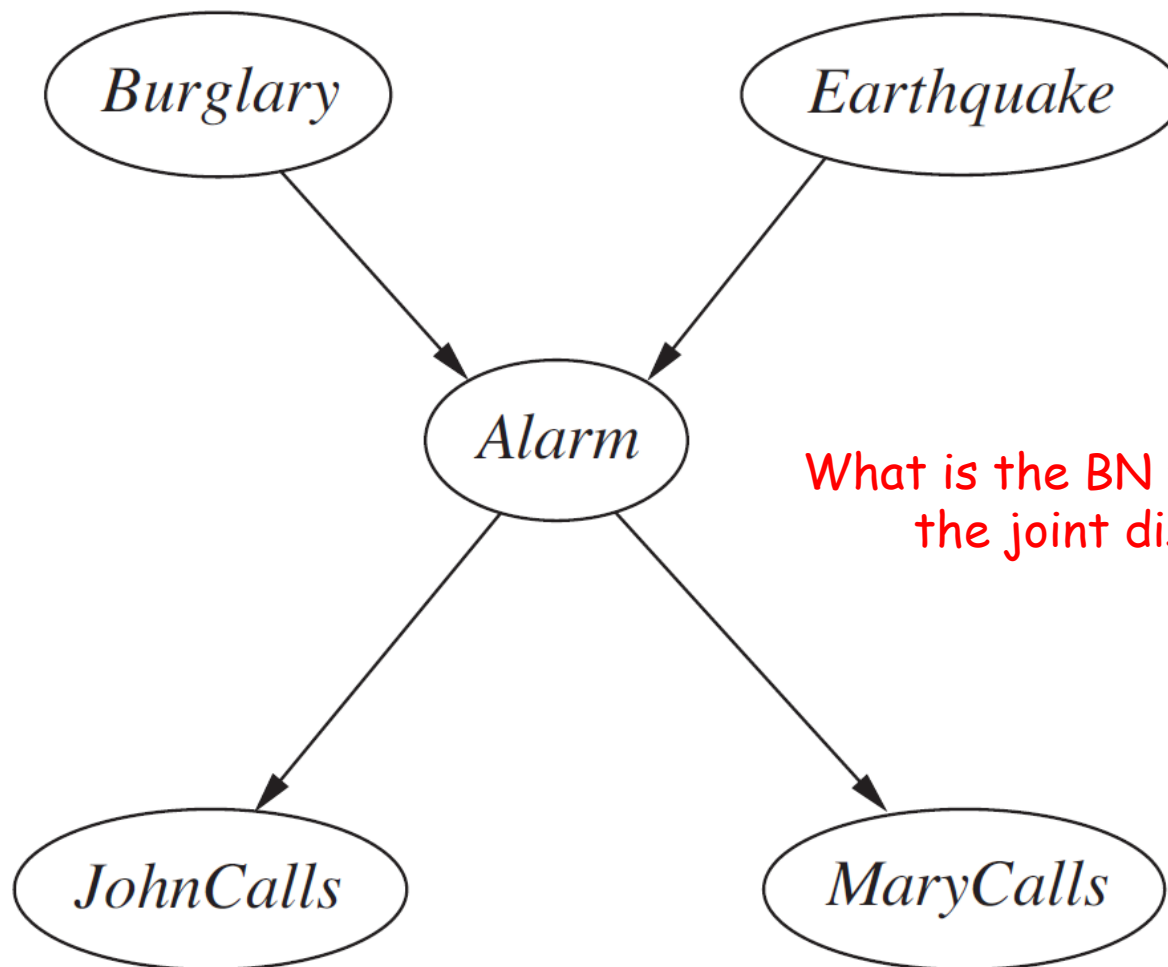
# RELATIONSHIPS

- $X_i \rightarrow X_j$ 
  - $X_i$  is the **parent**
  - $X_j$  is the **child**
- $X_i$  is an **ancestor** of  $X_j$  if there is a directed path from  $X_i$  to  $X_j$
- $X_i$  is a **descendant** of  $X_j$  if there is a directed path from  $X_j$  to  $X_i$
- **Nondescendants**( $X_i$ )  $\equiv \mathcal{X} \setminus \text{Descendants}(X_i)$

# BAYESIAN NETWORK FACTORIZATION

$$P(X_1, \dots, X_n) = \prod_i P(X_i \mid Pa(X_i))$$

BURGLARY EXAMPLE  $P(B, E, A, J, M)$

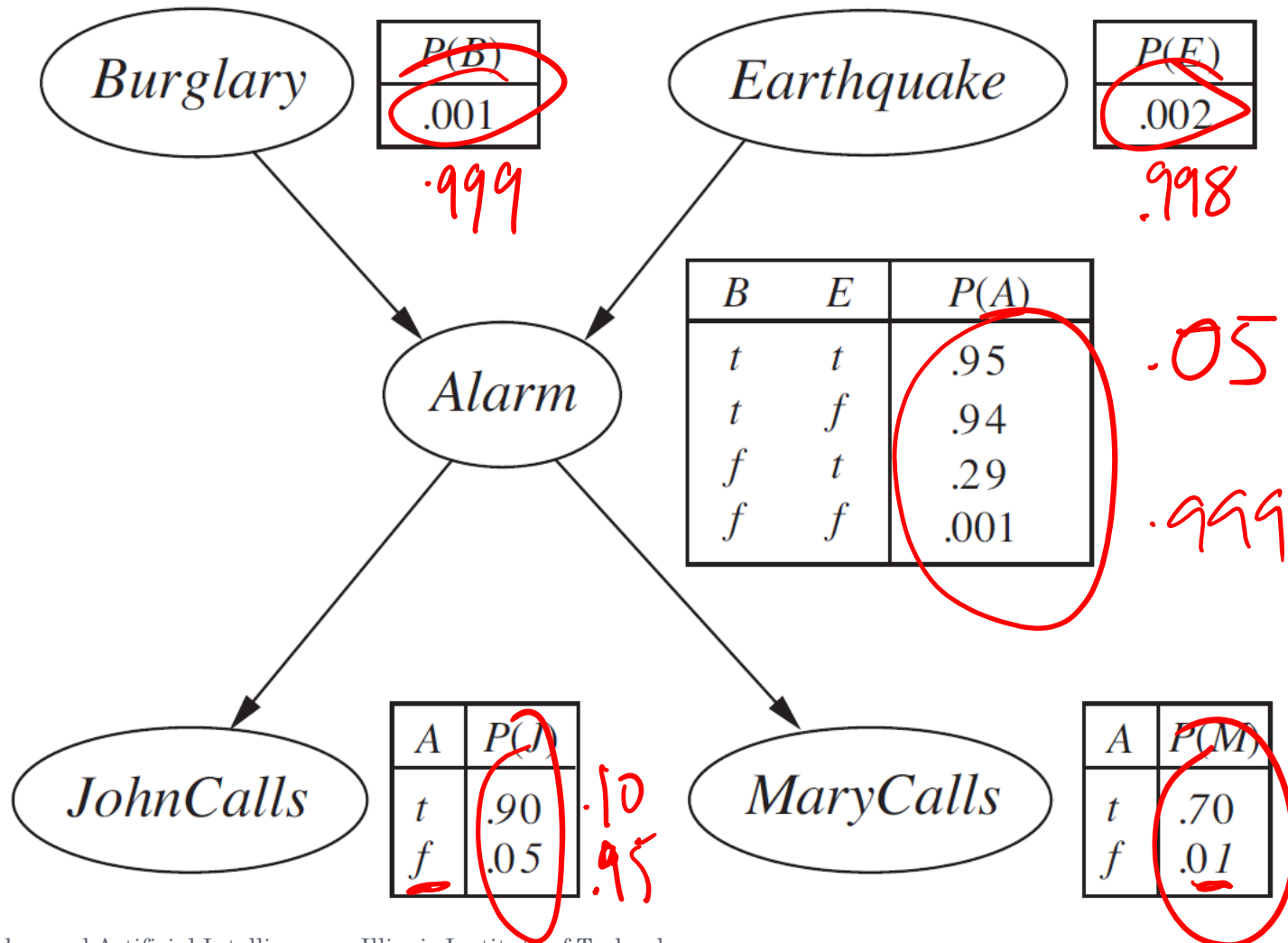


What is the BN factorization of the joint distribution?

$$= P(B) \cdot P(E) \cdot P(A|B, E) \cdot P(M|A) \cdot P(J|A)$$

# BURGLARY EXAMPLE

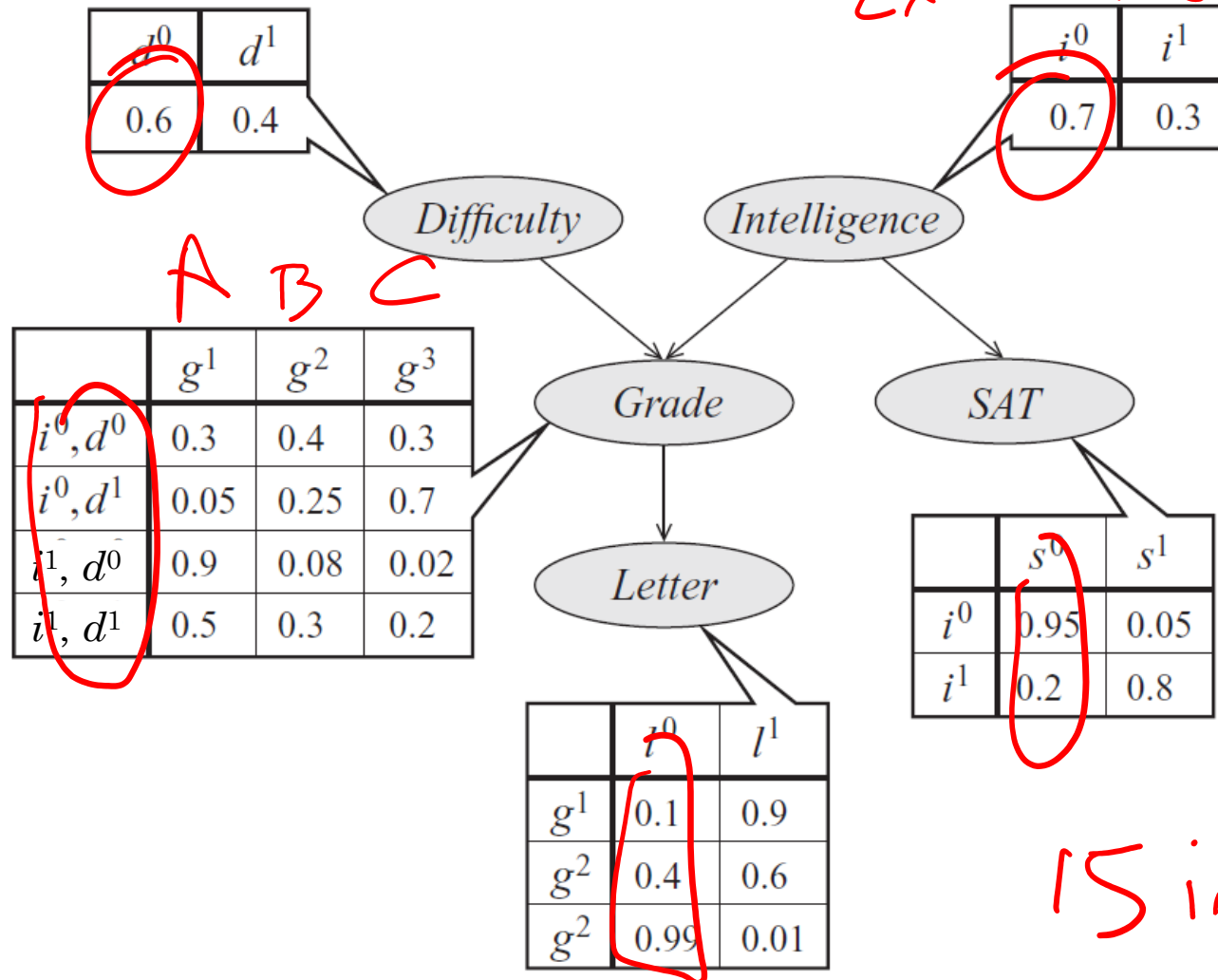
10 ind vs 31



# STUDENT EXAMPLE

$$P(D, I, G, S, L)$$

$$2 \times 2 \times 3 \times 2 \times 2 - 1 = 47$$



# INDEPENDENCIES

- X is independent of its non-descendants given its parents
  - $X \perp \text{Non-descendants}(X) \mid \text{Parents}(X)$
- D-separation

# INDEPENDENCIES – D-SEPARATION

- Definition: Observed  $\equiv$  Its value is known
- Causal trail
  - $X \rightarrow Y \rightarrow Z$ ; E.g., Burglary  $\rightarrow$  Alarm  $\rightarrow$  MaryCalls
  - X and Z are independent if Y is observed
- Evidential trail
  - $X \leftarrow Y \leftarrow Z$ ; E.g., MaryCalls  $\leftarrow$  Alarm  $\leftarrow$  Burglary
  - X and Z are independent if Y is observed
- Common cause
  - $X \leftarrow Y \rightarrow Z$ ; E.g., JohnCalls  $\leftarrow$  Alarm  $\rightarrow$  MaryCalls
  - X and Z are independent if Y is observed
- Common effect
  - $X \rightarrow Y \leftarrow Z$ ; E.g., Burglary  $\rightarrow$  Alarm  $\leftarrow$  Earthquake
  - X and Z are marginally independent, but they become dependent if Y or any of Y's descendants are observed

# EXAMPLES

- X causes Y and Y causes Z; no direct relationship between X and Z
  - $X \rightarrow Y \rightarrow Z$
  - Nothing is marginally independent of each other
  - $Z \perp X \mid Y$
- Y causes both X and Z; no direct relationship between X and Z
  - $X \leftarrow Y \rightarrow Z$
  - Nothing is marginally independent of each other
  - $Z \perp X \mid Y$
- Both X and Z cause Y; no direct relationship between X and Z
  - $X \rightarrow Y \leftarrow Z$
  - X and Z are marginally independent
  - X and Z become dependent when the value of Y is known



# INDEPENDENCE $\Leftrightarrow$ FACTORIZATION

- Independence  $\Rightarrow$  Factorization
- Factorization  $\Rightarrow$  Independence

# REASONING PATTERNS

## ○ Causal reasoning

- From causes to effects
  - E.g., Burglary to Alarm to MaryCalls
  - E.g., Intelligence to Grade to Letter

## ○ Evidential reasoning

- From effects to the causes
  - E.g., JohnCalls to Alarm to Earthquake
  - E.g., Letter to Grade to Difficulty

## ○ Explaining away/inter-causal reasoning

- Causes of a common effect interact
  - E.g., Earthquake, Burglary, and Alarm (and Alarm's descendants)
  - E.g., Difficulty, Intelligence, and Grade (and Grade's descendants)

# INFERENCE IN BAYESIAN NETWORKS

- There are several methods, some are exact and some are approximate
- We will study two in this class
  - *Variable elimination*
  - *Message passing*

# VARIABLE ELIMINATION

- Let
  - $\mathbf{V}$  be the set of all variables,  $\mathbf{Q}$  be the set of query variables,  $\mathbf{E}$  be the set of evidence variables
  - $P(\mathbf{Q} | \mathbf{E})$  be the query
- 1. Write down the joint dist. using the Bayesian network structure
- 2. Set the variables in  $\mathbf{E}$  to their respective values
- 3. Sum over all variables in  $\mathbf{V} \setminus (\mathbf{Q} \cup \mathbf{E})$ 
  - a) Pick an order for variables in  $\mathbf{V} \setminus (\mathbf{Q} \cup \mathbf{E})$
  - b) For each variable  $V_i$  in  $\mathbf{V} \setminus (\mathbf{Q} \cup \mathbf{E})$ , create a new factor by
    - Multiplying all the factors that contains  $V_i$ , and
    - Summing over possible values of  $V_i$
- 4. Normalize the last remaining factor (this step is unnecessary if  $\mathbf{E}$  is empty)

# IRRELEVANT

- Let
  - $\mathbf{V}$  be the set of all variables,  $\mathbf{Q}$  be the set of query variables,  $\mathbf{E}$  be the set of evidence variables
  - $P(\mathbf{Q} | \mathbf{E})$  be the query
- $Y \in \mathbf{V} \setminus \{\mathbf{Q} \cup \mathbf{E}\}$  is irrelevant iff
  - $Y \notin \text{Ancestors of } \{\mathbf{Q} \cup \mathbf{E}\}$ 
    - or
  - $Y \perp \mathbf{Q} | \mathbf{E}$
- Examples

# VARIABLE ELIMINATION EXAMPLES

- See OneNote

# MESSAGE PASSING

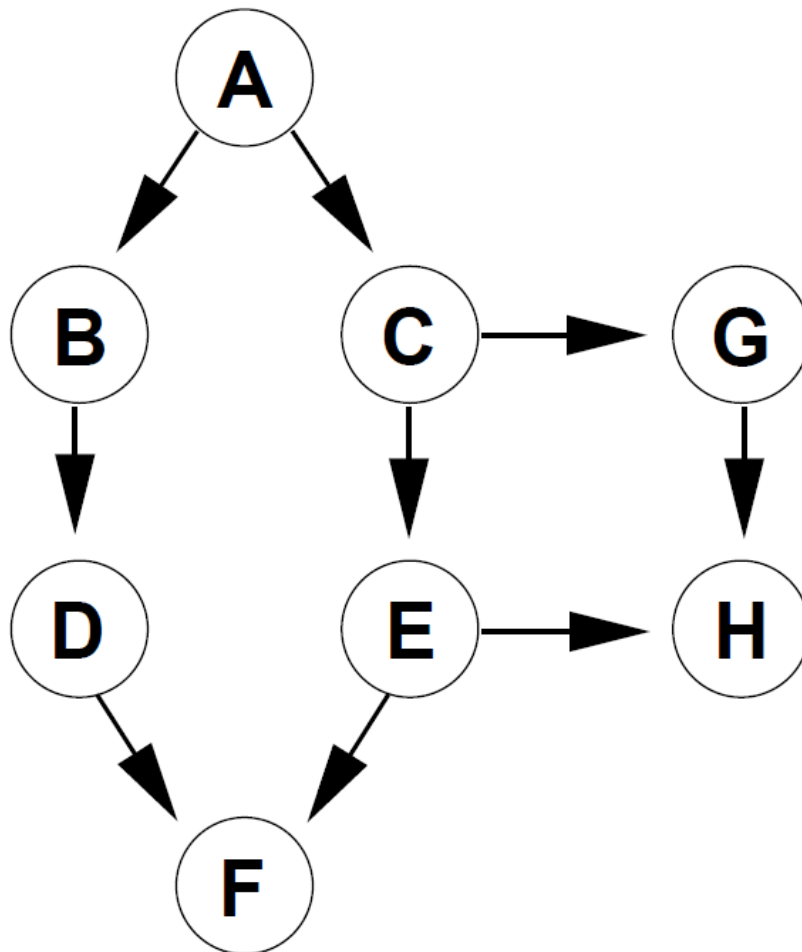
- Junction tree algorithm
- See OneNote for an example

# MESSAGE PASSING - MOTIVATION

- We are interested in multiple marginal/conditional probabilities
- In variable elimination, we define our target upfront and then eliminate the others
- If we need probabilities for other variables, there is no apparent way of reusing shared computations
- In the student example, assume that I'm interested in  $P(G)$  and  $P(L)$ . What are some of the shared computations?



# EXAMPLE



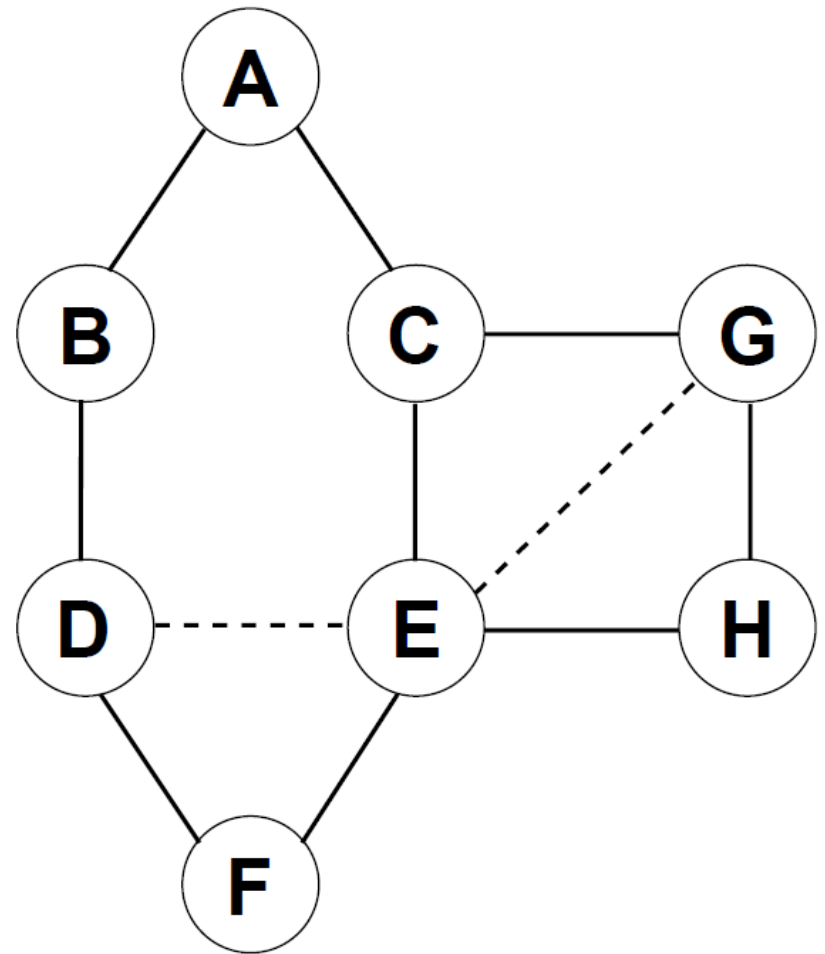
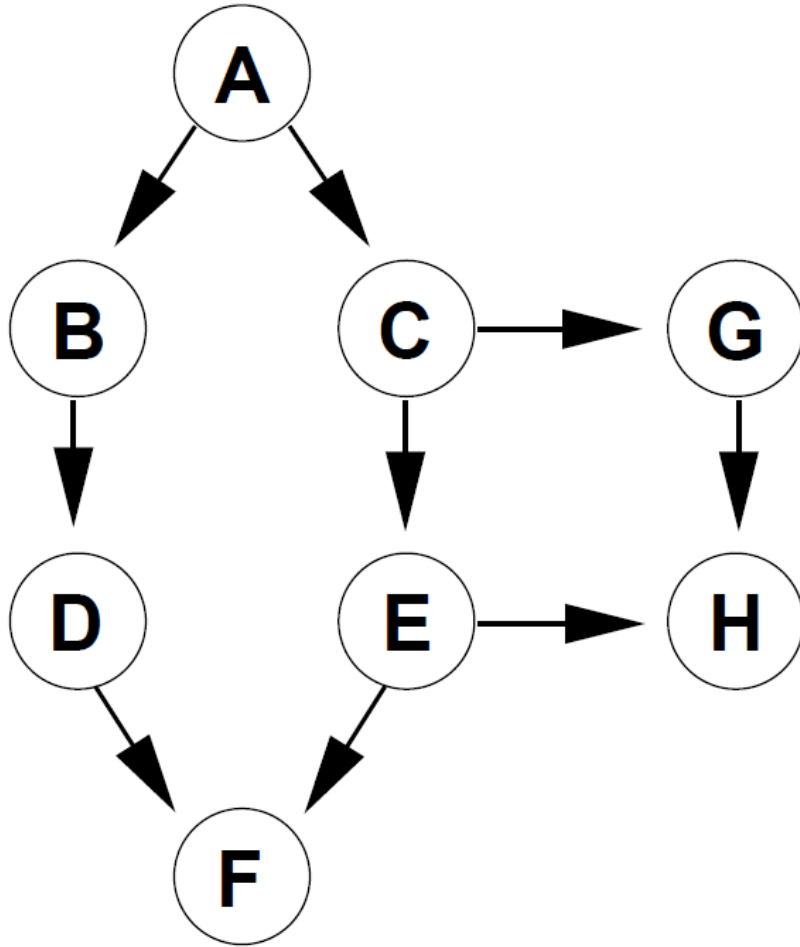
Calculate  $P(H)$  using variable elimination

Now, calculate  $P(G)$  using variable elimination

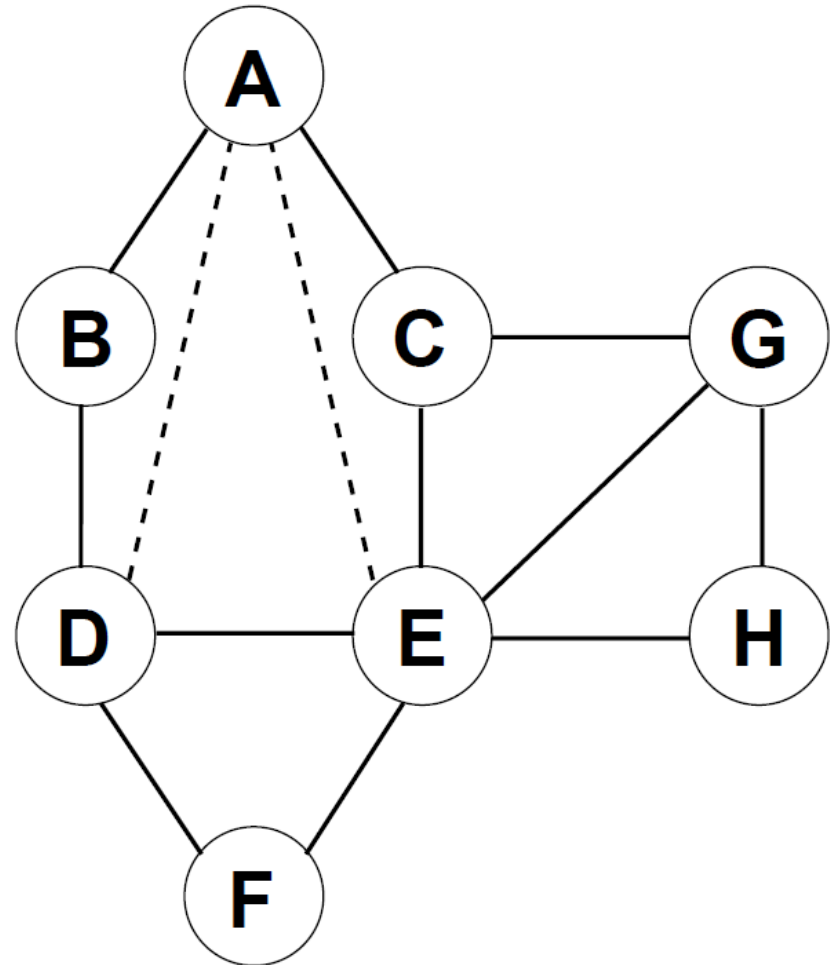
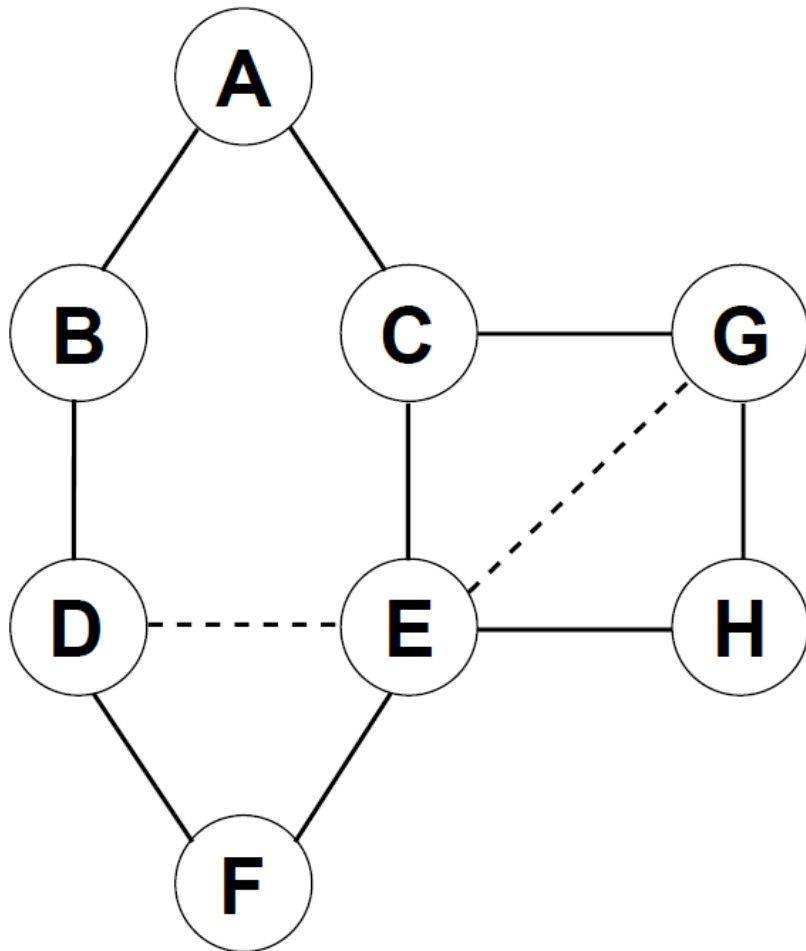
# VARIABLE ELIMINATION AS GRAPH TRANSFORMATION

- First, construct the moral graph
- Then, eliminate variables so that each elimination introduces the fewest number of edges
- Take note of the factors

## EXAMPLE - MORALIZE



ELIMINATION ORDER: H, G, F, C, B, D, E, A



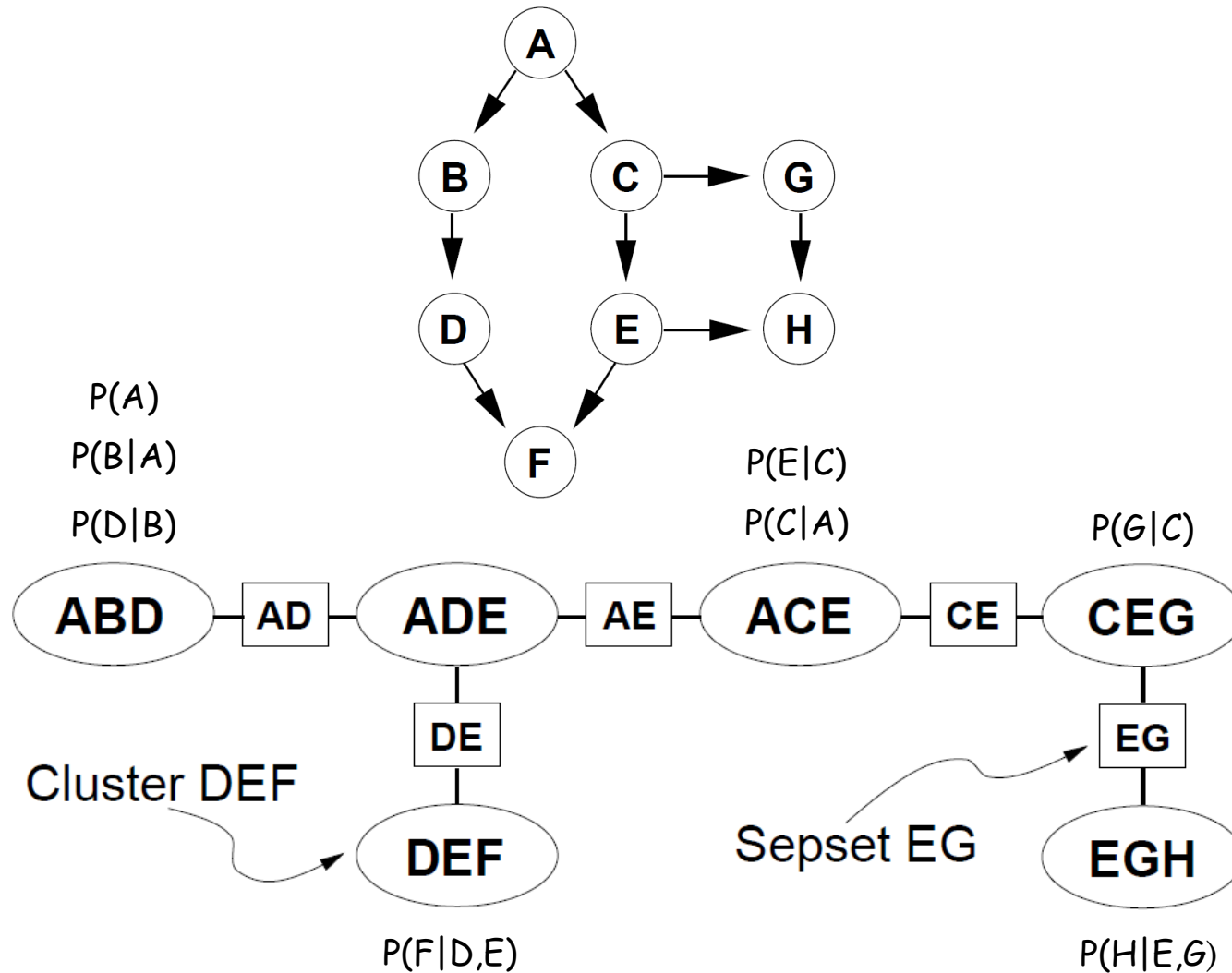
# CLUSTER GRAPH

- A *cluster graph*  $\mathcal{U}$  for a set of factors  $\Phi$  over  $\mathcal{X}$  is an undirected graph, each of whose nodes are associated with a cluster  $C_i \subseteq \mathcal{X}$ .
- A cluster graph must be *family preserving* – each factor  $\phi \in \Phi$  must be associated with a cluster  $C_i$ , denoted as  $\alpha(\phi)$ , such that  $\text{Scope}[\phi] \subseteq C_i$ .
- Each edge between a pair of clusters  $C_i$  and  $C_j$  is associated with a *sepset*  $S_{ij} \subseteq C_i \cap C_j$

# RUNNING INTERSECTION PROPERTY

- Let  $\mathcal{T}$  be a cluster tree.  $\mathcal{T}$  has *running intersection property* if, whenever there is a variable  $X$  such that  $X \in C_i$  and  $X \in C_j$ , then  $X$  is also in every cluster in the unique path in  $\mathcal{T}$  between  $C_i$  and  $C_j$ .
- A cluster tree that satisfies the running intersection property is also called the *join/clique/junction tree*.
- **Theorem:** A cluster tree obtained through a run of variable elimination satisfies the running intersection property; that is, it is a clique tree.

# EXAMPLE CLIQUE TREE

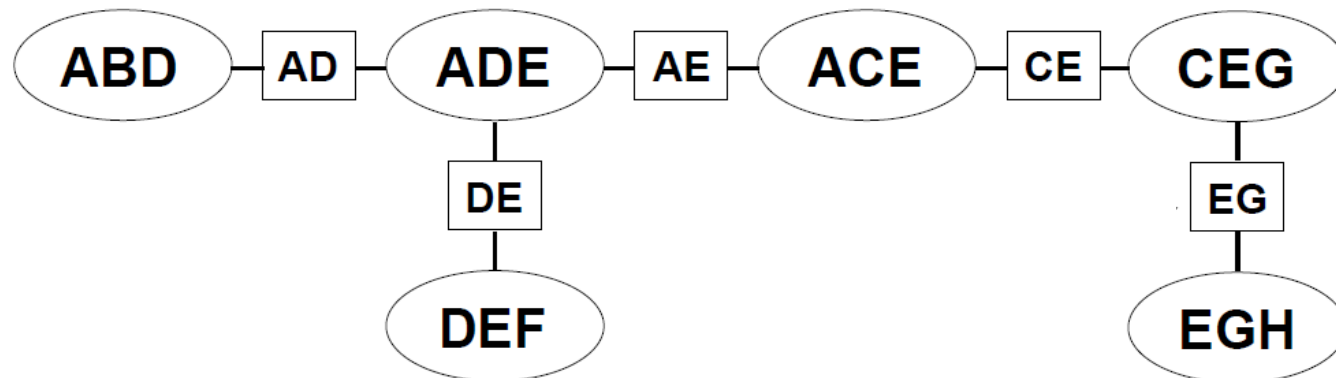
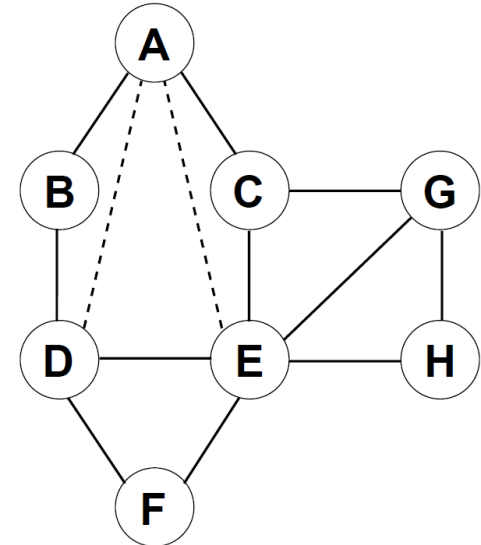
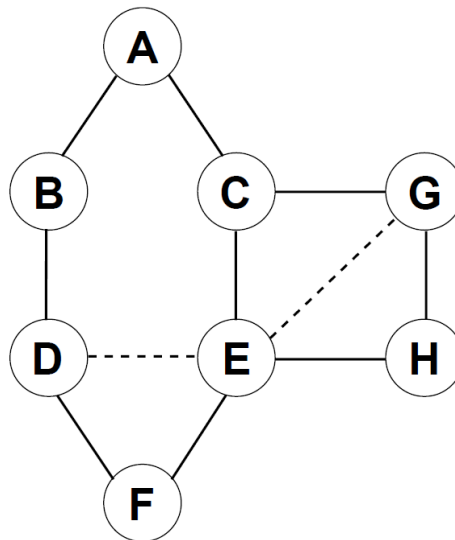
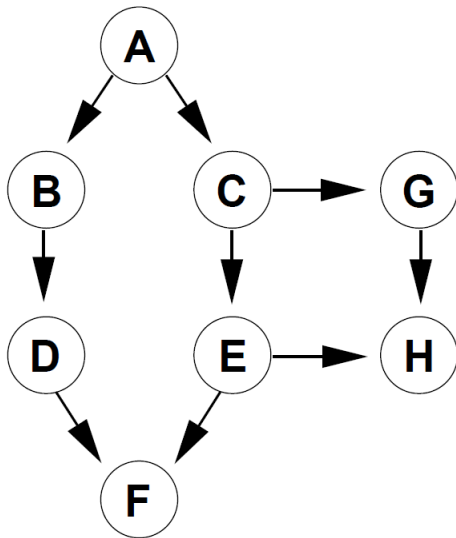


# CONSTRUCT A CLIQUE TREE

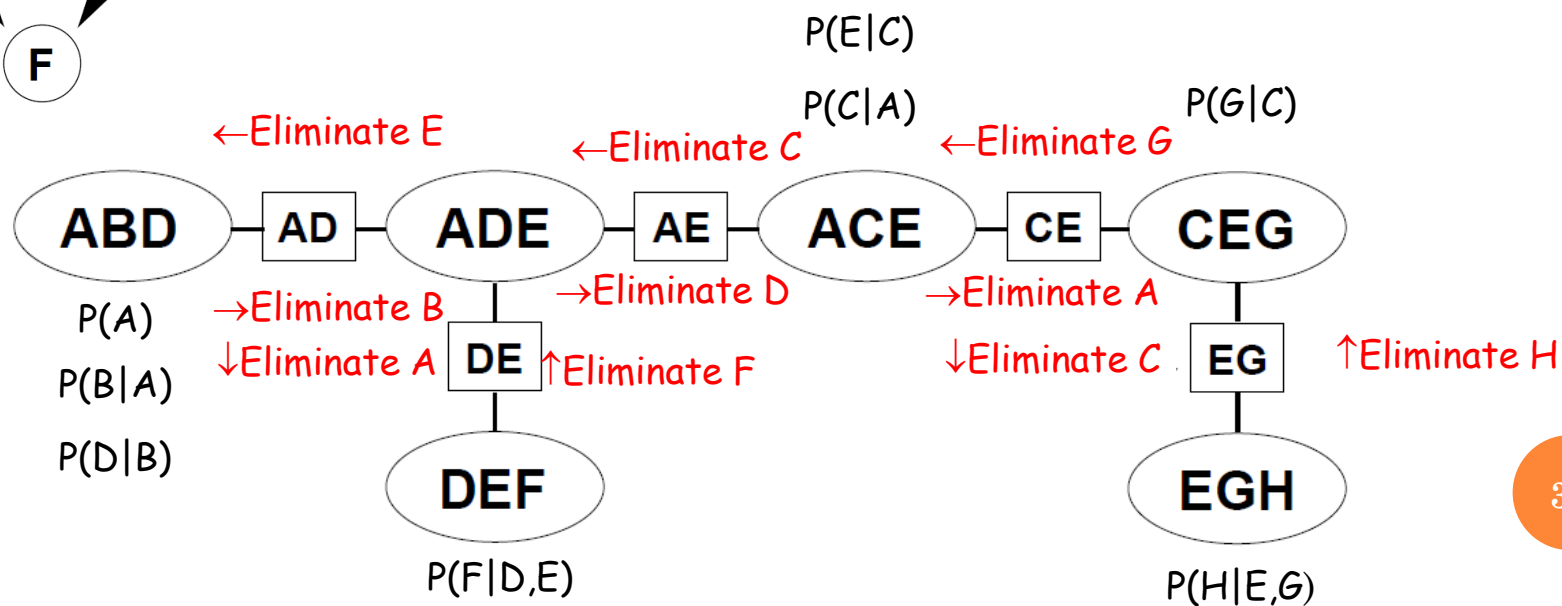
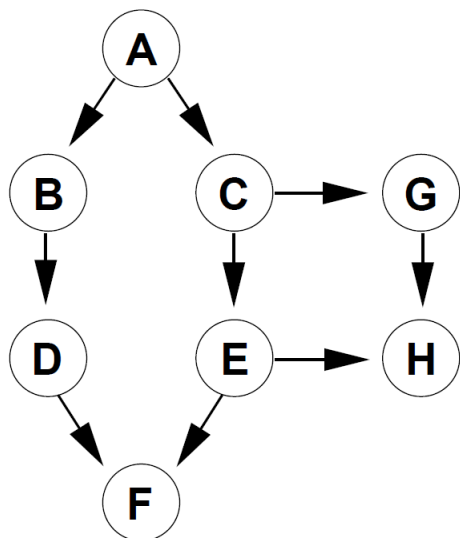
1. Moralize the graph
2. Pick a variable elimination order
3. Eliminate the variables, noting the maximal cliques
4. The cliques are the nodes of the tree
5. Until a tree is formed (i.e.,  $n-1$  edges are added)
  - a) Connect two disconnected components by a maximal size sepset



ELIMINATION ORDER: H, G, F, C, B, D, E, A



# VARIABLE ELIMINATION ON JUNCTION TREE



# MESSAGE PASSING ON JUNCTION TREE

- Clusters receive from and send messages to its neighbors
- Each message pass consists of elimination of one or more variables
- A cluster  $C_i$  is ready to send a message to its neighbor  $C_j$ , when it receives messages from its *all other* neighbors
- A message from  $C_i$  to  $C_j$  is computed as follows
  - $C_i$  multiplies all the factors assigned to it, and all the messages it received from its *other* neighbors
  - It sums out  $C_i \setminus S_{ij}$

## A MESSAGE

$$\delta_{i \rightarrow j} = \sum_{C_i \setminus S_{ij}} \left( \left( \prod_{\phi: \alpha(\phi)=i} \phi \right) \times \left( \prod_{k \in (Nb_i - \{j\})} \delta_{k \rightarrow i} \right) \right)$$

# BELIEF

$$\beta_i = \left( \prod_{\phi: \alpha(\phi)=i} \phi \right) \times \left( \prod_{k \in Nb_i} \delta_{k \rightarrow i} \right)$$

# LINEAR GAUSSIAN EXERCISE

## ○ Given

- $p(X) \sim N(\mu_X; \sigma_X^2)$
- $p(Y | X) \sim N(\beta_0 + \beta_1 \mu_X; \sigma_Y^2)$

## ○ Calculate

- $p(Y)$
- $p(Y | X = 5)$
- $p(X | Y = 5)$
- $p(X, Y)$