

# CS 583: PROBABILISTIC GRAPHICAL MODELS

## TOPIC: VARIABLE ELIMINATION



**Mustafa Bilgic**

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

 <https://twitter.com/bilgicm>

# TASK

- Given a graphical model over  $\mathcal{X}$  (structure and parameters)
- Compute  $P(\mathbf{Y} \mid \mathbf{e})$ , where  $\mathbf{Y} \subseteq \mathcal{X}$  and  $\mathbf{E} \subseteq \mathcal{X}$
- There are several approaches
  - Exact inference
    - Variable elimination
    - Belief propagation
  - Approximate inference
    - Sampling
- Today, we'll cover variable elimination

# VARIABLE ELIMINATION

- $P(\mathbf{Y} \mid \mathbf{e}) = P(\mathbf{Y}, \mathbf{e}) / P(\mathbf{e})$
- $\mathbf{W} = \mathcal{X} - \mathbf{Y} - \mathbf{E}$
- $P(\mathbf{y}, \mathbf{e}) = \sum_{\mathbf{w}} P(\mathbf{y}, \mathbf{e}, \mathbf{w})$
- $P(\mathbf{e}) = \sum_{\mathbf{y}, \mathbf{w}} P(\mathbf{y}, \mathbf{e}, \mathbf{w})$
- Or, better yet:  $P(\mathbf{e}) = \sum_{\mathbf{y}} P(\mathbf{y}, \mathbf{e})$

$$P(Y, E) = \sum_w P(Y, E, W)$$

- $P(Y, E, W)$  can be represented as
  - $\prod P(X_i | \text{Pa}(X_i))$
  - $1/Z \prod \phi(D_i)$
- The problem with  $P(\mathbf{y}, \mathbf{e}) = \sum_w P(\mathbf{y}, \mathbf{e}, \mathbf{w})$  is that the joint representation is exponential
  - The very first problem we were trying to avoid

# COMPLEXITY

- Unfortunately, exact inference is  $\mathcal{NP}$ -hard in worst case
  - Proof: pages 288 and 289. Reduction from 3-SAT
- Approximate inference is also  $\mathcal{NP}$ -hard in worst case
  - Proof: pages 291 and 292.
- Good news:
  - In general, we care about the cases we encounter in practice, not the worst-case scenario

# KEY IDEA

- Summation can be moved inside
- $\sum_x \sum_y x * y = \sum_x x * (\sum_y y)$
- If  $x$  has  $n$  and  $y$  has  $m$  possible values, how many operations are needed, if we use
  - $\sum_x \sum_y x * y$  ?
  - $\sum_x x * (\sum_y y)$  ?

# OUTLINE

- First, focus on Bayesian networks
  - Simple linear chains
  - More complex structures
- Two cases
  - Marginal queries:  $\mathbf{E} = \emptyset$
  - Conditional queries:  $\mathbf{E} \neq \emptyset$

# VARIABLE ELIMINATION

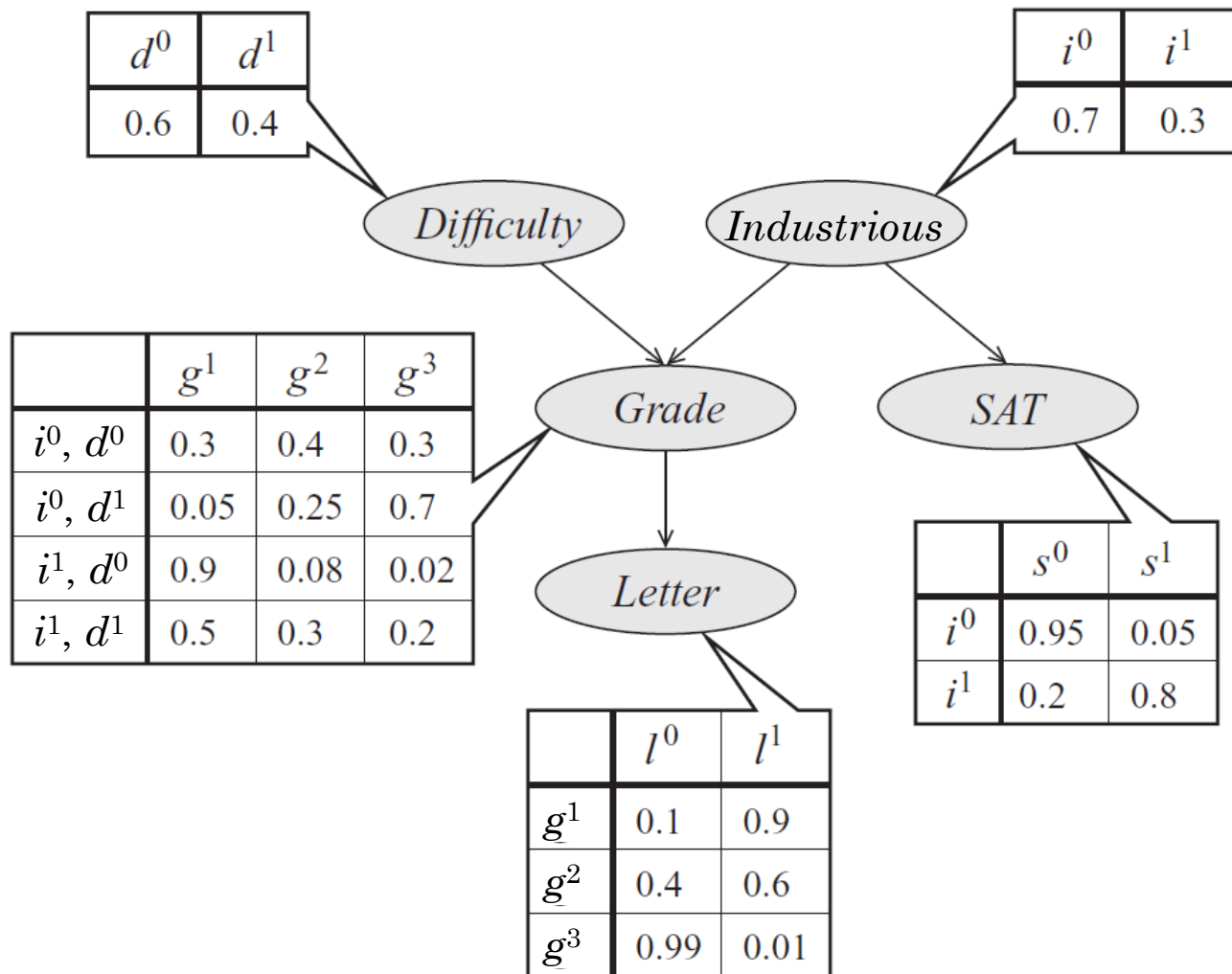
- $\mathcal{X}$ : all variables,  $\mathbf{Y}$ : query variables,  $\mathbf{E}$ : evidence variables,  $\mathbf{W} = \mathcal{X} - \mathbf{Y} - \mathbf{E}$ : remaining variables
- 1. Write down the joint for  $P(\mathcal{X})$
- 2. Set  $X_i \in \mathbf{E}$  to their values
- 3. Pick an order for  $X_j \in \mathbf{W}$
- 4. Sum out each  $X_j$  from the joint
  - a) Multiply the factors  $\phi(X_j, Z_1), \dots, \phi(X_j, Z_k)$  to create  $\psi(X_j, Z_1, \dots, Z_k)$
  - b) Sum out  $X_j$  from  $\psi(X_j, Z_1, \dots, Z_k)$  to create  $\tau(Z_1, \dots, Z_k)$
- 5. What remains is  $\tau(\mathbf{Y}, \mathbf{e})$ . Normalize it to get  $P(\mathbf{Y} \mid \mathbf{e})$ .



# EXAMPLES – LINEAR CHAIN BNS

- $A \rightarrow B$ 
  - $P(B) = ?$
  - $P(A) = ?$
- $X_1 \rightarrow X_2 \rightarrow \dots \rightarrow X_n$ 
  - $P(X_i)$  where  $1 \leq i \leq n$
- How many operations are needed if we compute the full joint distribution vs. if we use variable elimination?

# STUDENT NETWORK EXAMPLE

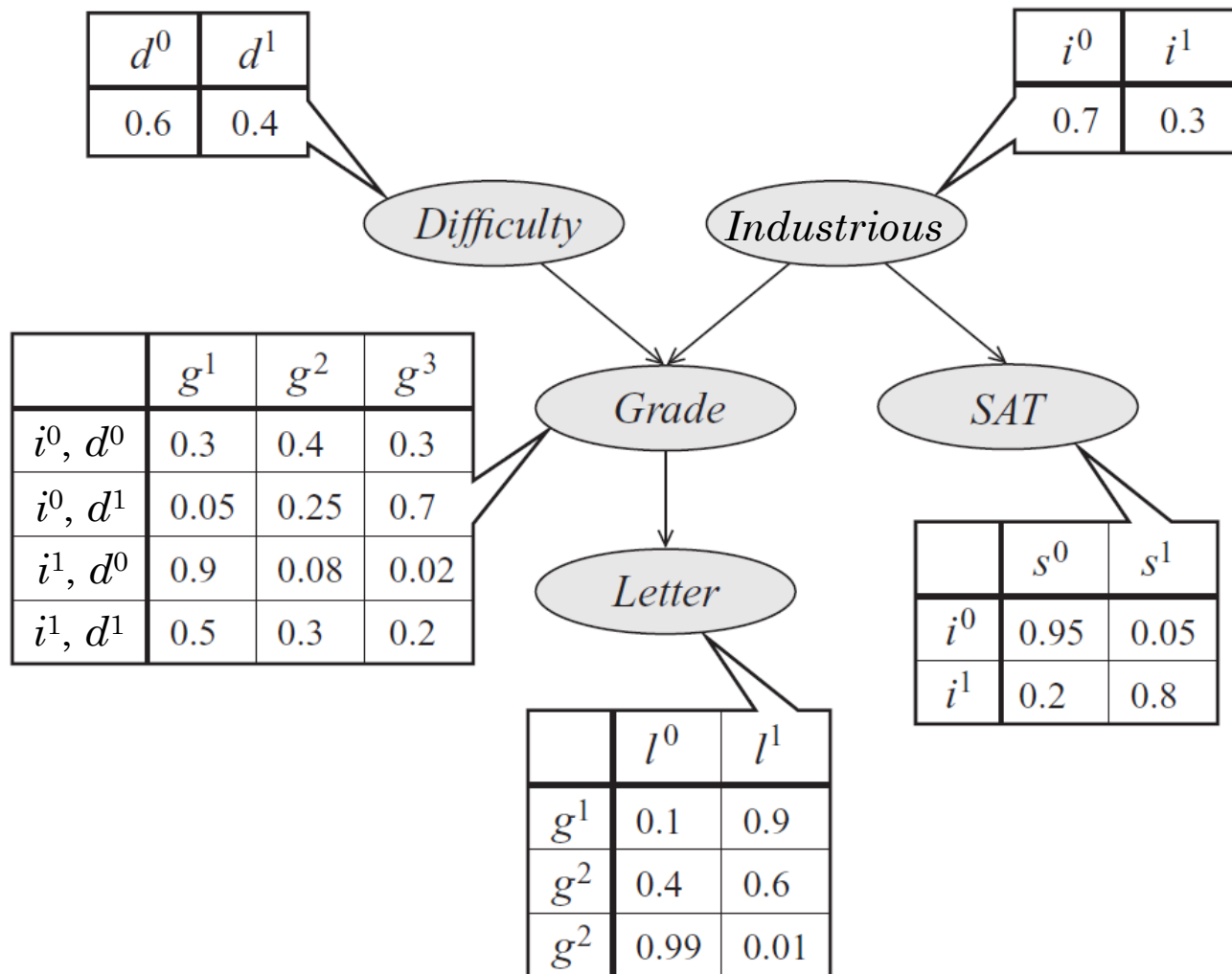


$P(D)$   
 $P(I)$   
 $P(S)$   
 $P(G)$   
 $P(L)$

# $P(Y | E)$ EXAMPLES

- $A \rightarrow B$ 
  - $P(B | A=t)$
  - $P(A=t | B=t)$
- $X_1 \rightarrow X_2 \rightarrow \dots \rightarrow X_n$ 
  - $P(X_i | X_j=x_j)$
  - $P(X_i | \mathbf{X}_j=\mathbf{x}_j)$

# STUDENT NETWORK EXAMPLE



$$\begin{aligned}
 &P(S \mid I=i^0) \\
 &P(D \mid I=i^0) \\
 &P(I \mid G=g^1) \\
 &P(D \mid G=g^1) \\
 &P(D \mid I=i^0, G=g^1)
 \end{aligned}$$

# P(L) – ORDER: I, S, D, G

Variable	All Factors	Participates	New Factor After *	# *s	New Factor After +	# +s	# Ops
I	P(I), P(D), P(S   I), P(G   D, I), P(L   G)	P(I), P(S   I), P(G   D, I)	$\psi_1(G, D, S, I)$	$2*3*2*2*2=48$	$\tau_1(G, D, S)$	$1*3*2*2=12$	60
S	P(D), P(L   G), $\tau_1(G, D, S)$	$\tau_1(G, D, S)$	$\psi_2(G, D, S)$	0	$\tau_2(G, D)$	$1*3*2=6$	6
D	P(D), P(L   G), $\tau_2(G, D)$	P(D), $\tau_2(G, D)$	$\psi_3(G, D)$	$1*3*2=6$	$\tau_3(G)$	$1*3$	9
G	P(L   G), $\tau_3(G)$	P(L   G), $\tau_3(G)$	$\psi_4(L, G)$	$1*2*3=6$	$\tau_4(L)$	$2*2=4$	10
Normalize	$\tau_4(L)$					1	3 (2 divs)
Total							88

# P(L) – ORDER: S, I, D, G

Variable	All Factors	Participates	New Factor After *	# *s	New Factor After +	# +s	# Ops
S	P(I), P(D), P(S   I), P(G   D, I), P(L   G)	P(S   I)	$\psi_1(I, S)$	0	$\tau_1(I)$	$1 * 2 = 2$	2
I	P(I), P(D), P(G   D, I), P(L   G) $\tau_1(I)$	P(I), P(G   D, I), $\tau_1(I)$	$\psi_2(G, D, I)$	$2 * 3 * 2 * 2 = 24$	$\tau_2(G, D)$	$1 * 3 * 2 = 6$	30
D	P(D), P(L   G), $\tau_2(G, D)$	P(D), $\tau_2(G, D)$	$\psi_3(G, D)$	$1 * 3 * 2 = 6$	$\tau_3(G)$	$1 * 3$	9
G	P(L   G), $\tau_3(G)$	P(L   G), $\tau_3(G)$	$\psi_4(L, G)$	$1 * 2 * 3 = 6$	$\tau_4(L)$	$2 * 2 = 4$	10
Normalize	$\tau_4(L)$					1	3 (2 divs)
Total							54

# MARKOV NETWORK EXAMPLE

A	B	$\phi(A,B)$	B	C	$\phi(B,C)$	A	B	C	$\phi(A,B)*\phi(B,C)$	P(A,B,C)
T	T	5	T	T	1	T	T	T	5	0.11
T	F	1	T	F	2	T	T	F	10	0.22
F	T	1	F	T	6	T	F	T	6	0.13
F	F	3	F	F	1	T	F	F	1	0.02
						F	T	T	1	0.02
						F	T	F	2	0.04
						F	F	T	18	0.39
						F	F	F	3	0.07
						Z			46	1.00

A	B	P(A,B)	B	C	P(B,C)
T	T	0.33	T	T	0.13
T	F	0.15	T	F	0.26
F	T	0.07	F	T	0.52
F	F	0.46	F	F	0.09

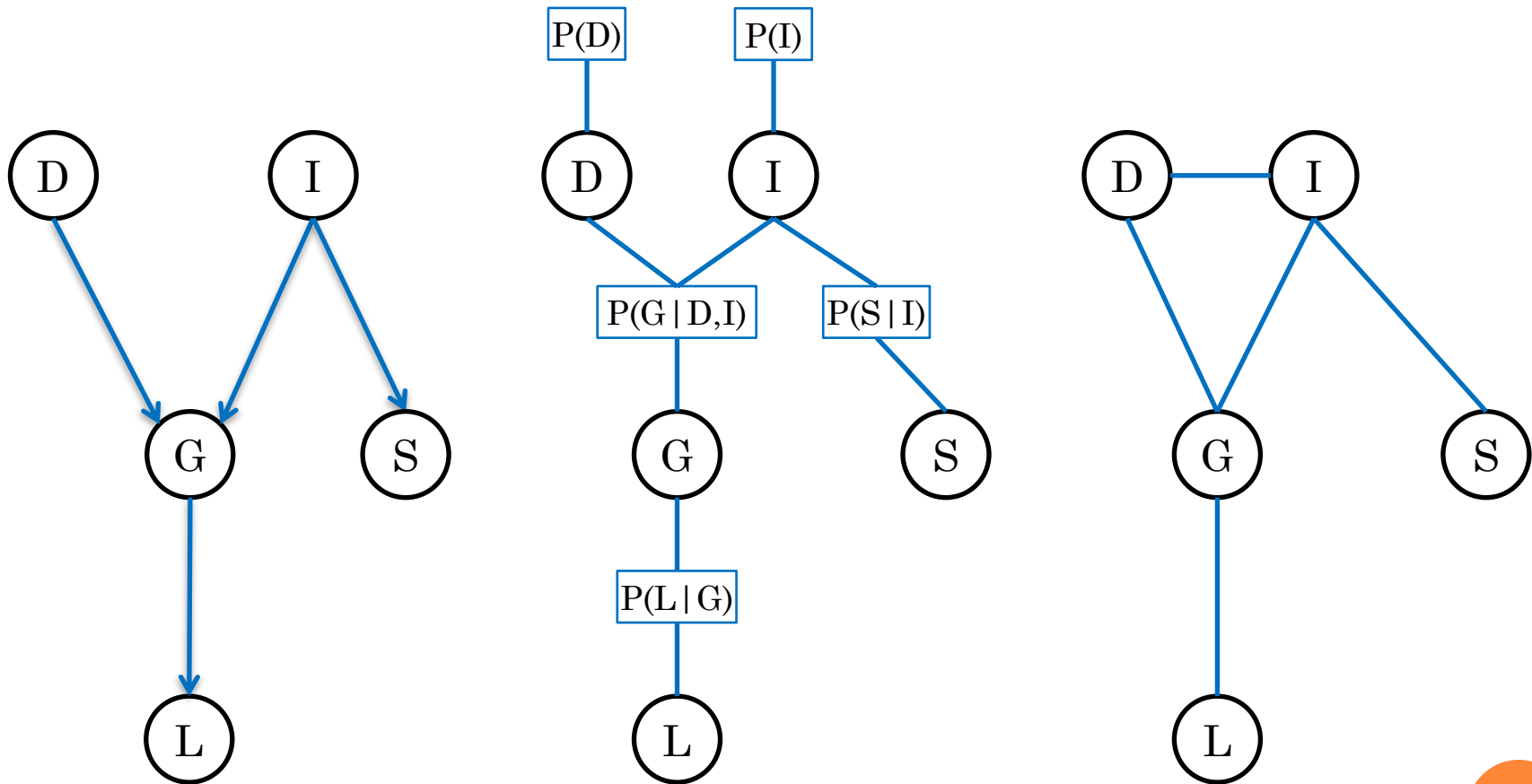
A	P(A)	B	P(B)	C	P(C)
T	0.48	T	0.39	T	0.65
F	0.52	F	0.61	F	0.35

# ELIMINATION AS GRAPH TRANSFORMATION

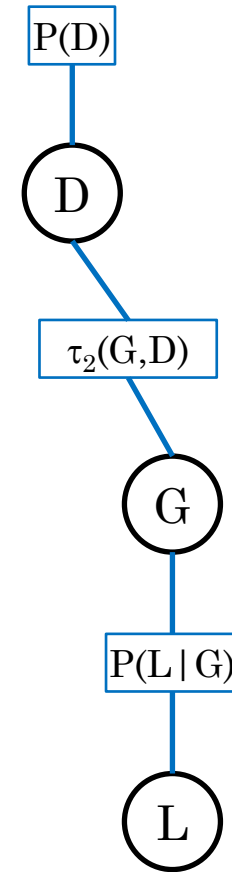
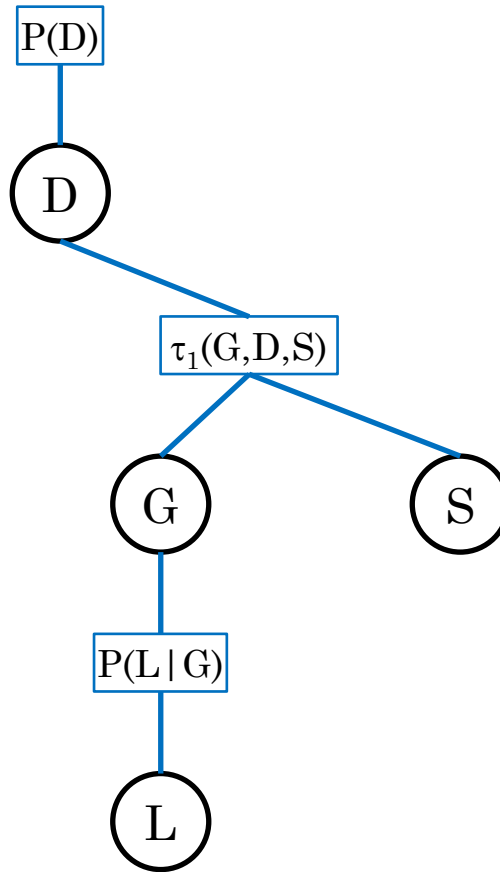
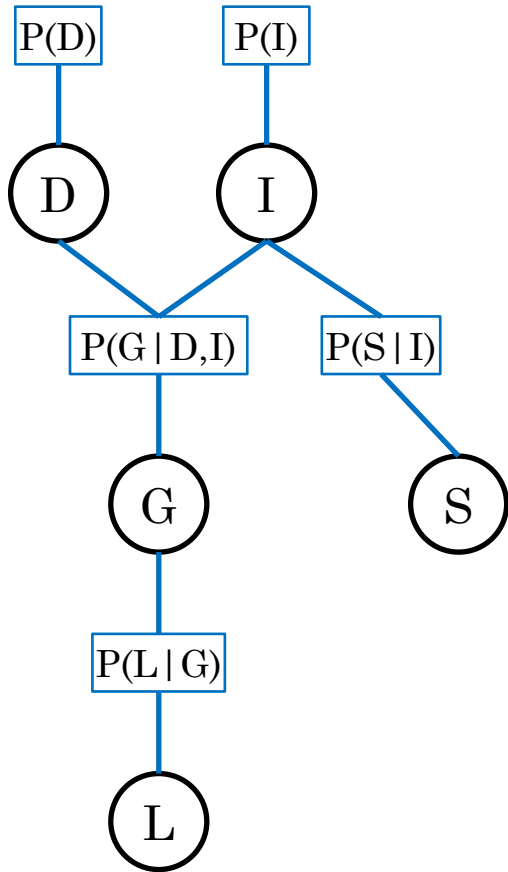
- Eliminating  $X$ 
  - Multiply all the factors  $X$  participates in
  - Sum out  $X$
- Graph transformation (need to be moralized first)
  - Connect all of  $X$ 's neighbors
  - Remove  $X$



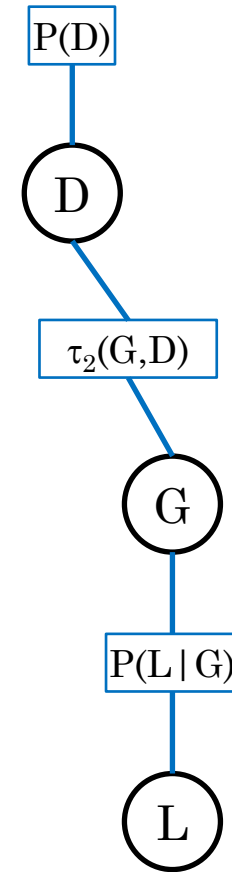
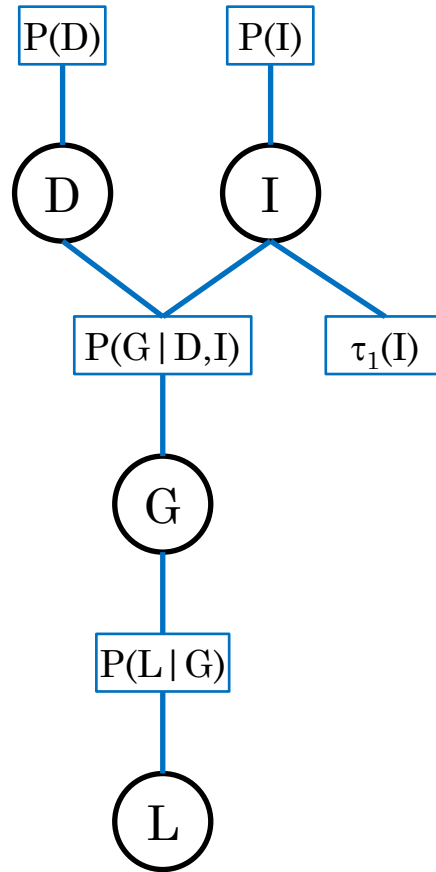
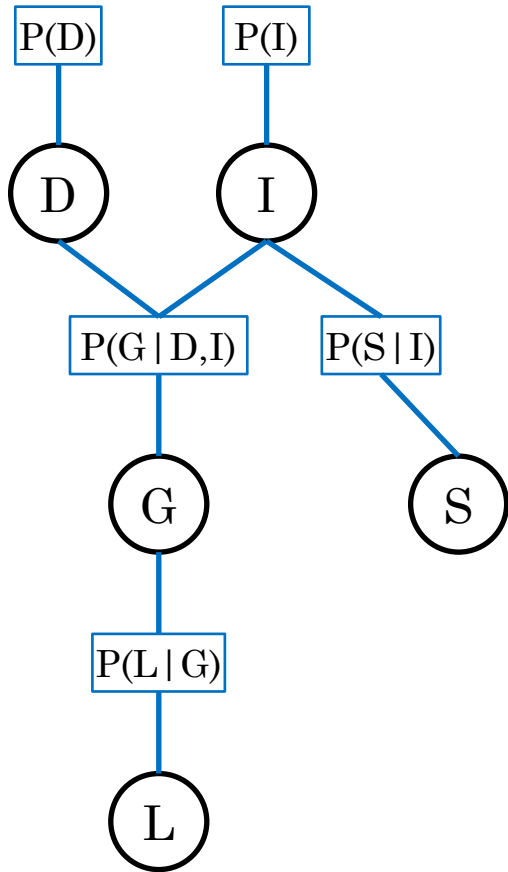
# REPRESENTATION



IF WE FIRST ELIMINATE  $I$  THEN  $S$



IF WE FIRST ELIMINATE  $S$  THEN  $I$



# FINDING GOOD ELIMINATION ORDERINGS

- Finding the best order is NP-hard
  - Best = optimal time and space complexity
- Heuristics
  - Min-neighbors
  - Min-fill
  - Weighted versions of min-neighbors and min-fill

# IRRELEVANT NODES IN BNS

- $\mathcal{X}$ : all variables,  $\mathbf{Y}$ : query variables,  $\mathbf{E}$ : evidence variables,  $\mathbf{W} = \mathcal{X} - \mathbf{Y} - \mathbf{E}$ : remaining variables
- A node  $X_i \in \mathbf{W}$  is irrelevant for the query  $P(\mathbf{Y} | \mathbf{e})$  if it can be removed from the network without effecting the value of  $P(\mathbf{Y} | \mathbf{e})$
- Obvious:
  - If  $\mathbf{Z} \subseteq \mathbf{W}$  is d-separated from  $\mathbf{Y}$  given  $\mathbf{E}$ , then  $\mathbf{Z}$  is irrelevant
- Perhaps less obvious:
  - Let  $\mathbf{Z}$  be ancestors of  $\mathbf{Y} \cup \mathbf{E}$ . Then  $\mathbf{W} \setminus \mathbf{Z}$  is irrelevant