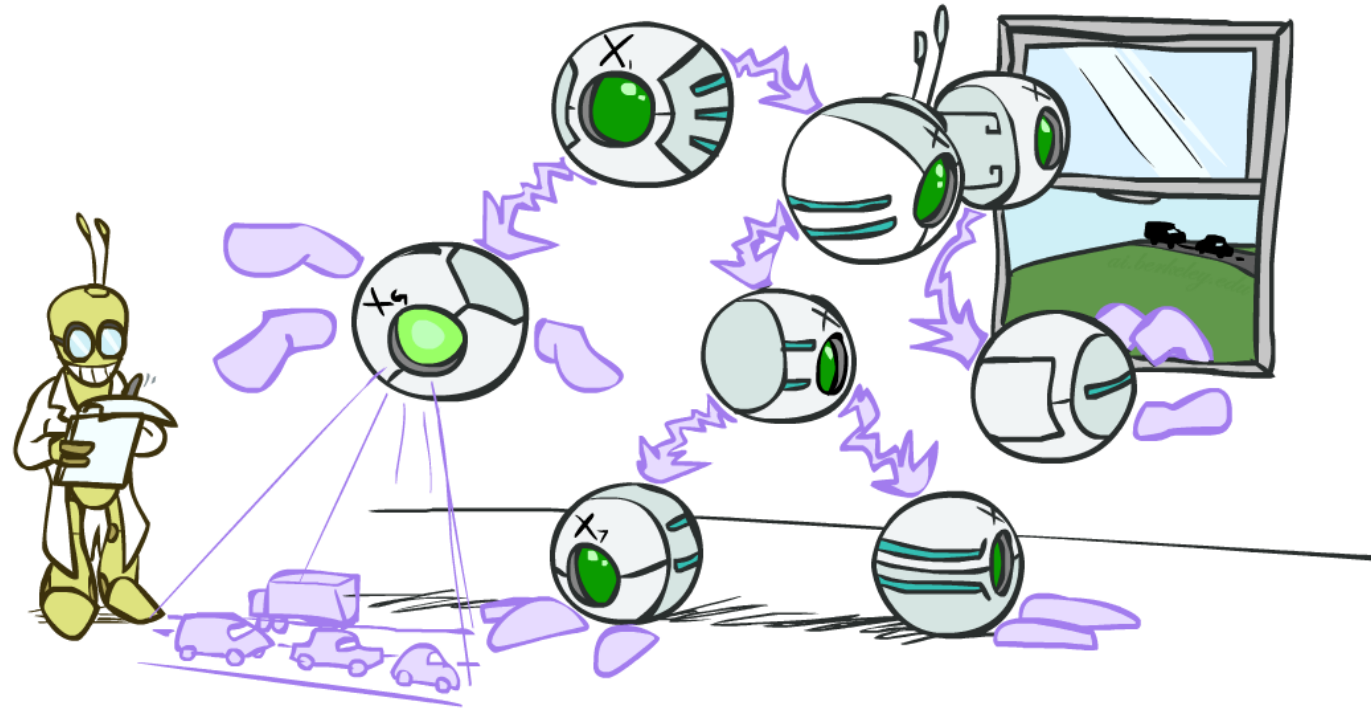# Announcement

- Programming 2 due 11:59pm today!

- Programming 3
  - Due: April 12, 11:59pm
- Homework 3
  - Due: April 7, 11:59pm

# Bayes Nets: Exact Inference
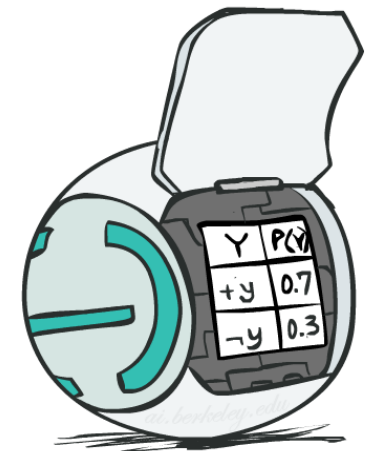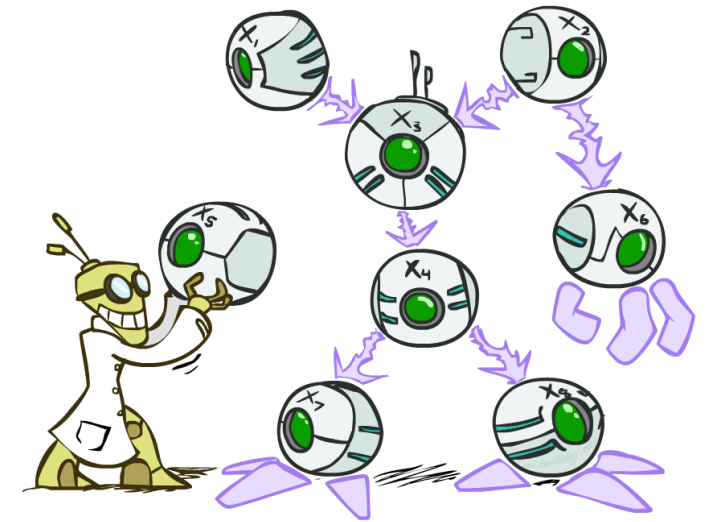


AIMA Chapter 14.4, PRML Chapter 8.4

# Bayes' Net Representation

- A directed, acyclic graph, one node per random variable

- A conditional probability table (CPT) for each node

  - A collection of distributions over X, one for each combination of parents' values
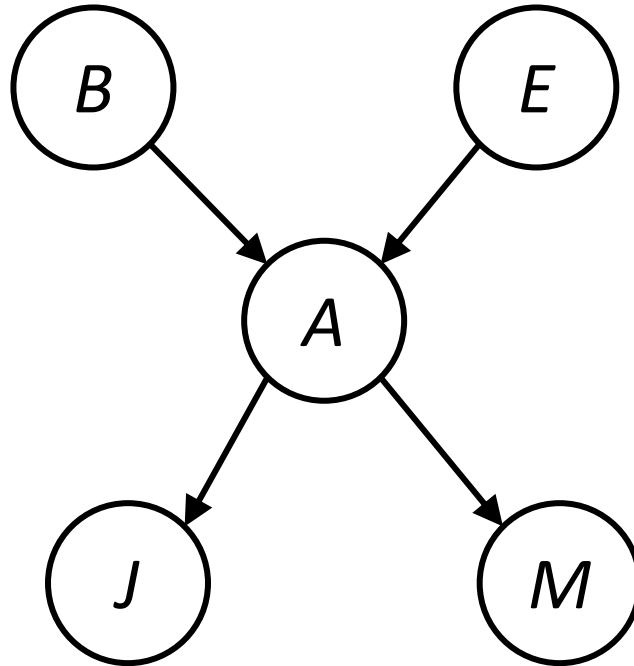
  $$P(X|a_1 \ldots a_n)$$

- Bayes' nets implicitly encode joint distributions

  - As a product of local conditional distributions

  - To see what probability a BN gives to a full assignment, multiply all the relevant conditionals together:

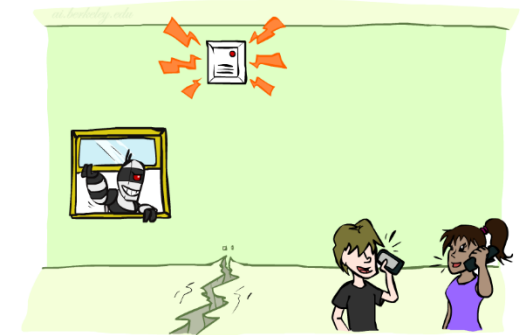  $$P(x_1, x_2, \ldots x_n) = \prod_{i=1}^{n} P(x_i | parents(X_i))$$

# Example: Alarm Network

| B | P(B) |
|---|------|
| +b | 0.001 |
| -b | 0.999 |

| E | P(E) |
|---|------|
| +e | 0.002 |
| -e | 0.998 |



| A | J | P(J\|A) |
|---|---|--------|
| +a | +j | 0.9 |
| +a | -j | 0.1 |
| -a | +j | 0.05 |
| -a | -j | 0.95 |

| A | M | P(M\|A) |
|---|---|--------|
| +a | +m | 0.7 |
| +a | -m | 0.3 |
| -a | +m | 0.01 |
| -a | -m | 0.99 |

| B | E | A | P(A\|B,E) |
|---|---|---|----------|
| +b | +e | +a | 0.95 |
| +b | +e | -a | 0.05 |
| +b | -e | +a | 0.94 |
| +b | -e | -a | 0.06 |
| -b | +e | +a | 0.29 |
| -b | +e | -a | 0.71 |
| -b | -e | +a | 0.001 |
| -b | -e | -a | 0.999 |

$$P(+b, -e, +a, -j, +m) =$$
$$P(+b)P(-e)P(+a|+b, -e)P(-j|+a)P(+m|+a) =$$

# Example: Alarm Network

| B | P(B) |
|---|---|
| +b | 0.001 |
| -b | 0.999 |

| E | P(E) |
|---|---|
| +e | 0.002 |
| -e | 0.998 |

| A | J | P(J|A) |
|---|---|---|
| +a | +j | 0.9 |
| +a | -j | 0.1 |
| -a | +j | 0.05 |
| -a | -j | 0.95 |

| A | M | P(M|A) |
|---|---|---|
| +a | +m | 0.7 |
| +a | -m | 0.3 |
| -a | +m | 0.01 |
| -a | -m | 0.99 |

| B | E | A | P(A|B,E) |
|---|---|---|---|
| +b | +e | +a | 0.95 |
| +b | +e | -a | 0.05 |
| +b | -e | +a | 0.94 |
| +b | -e | -a | 0.06 |
| -b | +e | +a | 0.29 |
| -b | +e | -a | 0.71 |
| -b | -e | +a | 0.001 |
| -b | -e | -a | 0.999 |



$$P(+b, -e, +a, -j, +m) =$$

$$P(+b)P(-e)P(+a|+b,-e)P(-j|+a)P(+m|+a) =$$
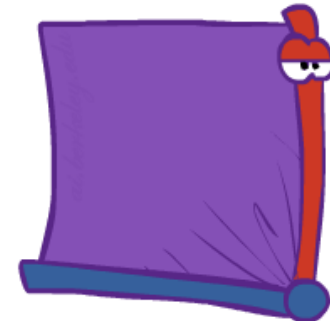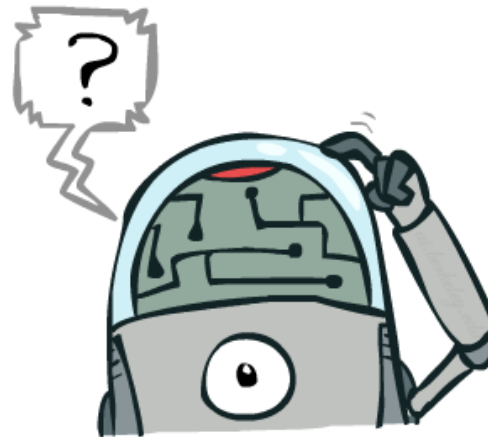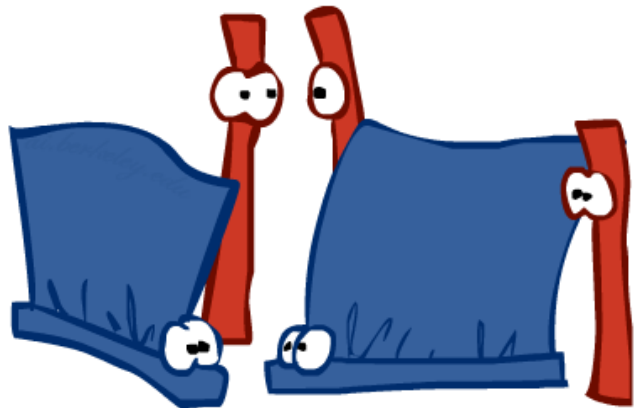
$$0.001 \times 0.998 \times 0.94 \times 0.1 \times 0.7$$

# Probabilistic Inference

- Enumeration (exact, exponential complexity)

- Variable elimination (exact, worst-case exponential complexity, often better)

- Inference is NP-complete

- Sampling (approximate)

# Inference

- Inference: calculating some useful quantity from a probability model (joint probability distribution)

- Examples:
  - Posterior marginal probability
    - $P(Q|e_1,..,e_k)$
    - E.g., what disease might I have?
  - Most likely explanation:
    - $\text{argmax}_q \, P(Q=q|e_1,..,e_k)$
    - E.g., what did he say?

# Inference by Enumeration
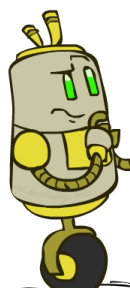
- General case:
  - Evidence variables: $E_1 \ldots E_k = e_1 \ldots e_k$
  - Query variable: $Q$ $\left.\begin{array}{l} \\ \\ \end{array}\right\}$ $X_1, X_2, \ldots X_n$
  - Hidden variables: $H_1 \ldots H_r$ *All variables*

- We want:

$$P(Q|e_1 \ldots e_k)$$

- Step 1: Select the entries consistent with the evidence

- Step 2: Sum out H to get joint of Query and evidence

- Step 3: Normalize



$$P(Q, e_1 \ldots e_k) = \sum_{h_1 \ldots h_r} P(\underbrace{Q, h_1 \ldots h_r, e_1 \ldots e_k}_{X_1, X_2, \ldots X_n})$$

$$\times \frac{1}{Z}$$

$$Z = \sum_q P(Q, e_1 \cdots e_k)$$

$$P(Q|e_1 \cdots e_k) = \frac{1}{Z} P(Q, e_1 \cdots e_k)$$
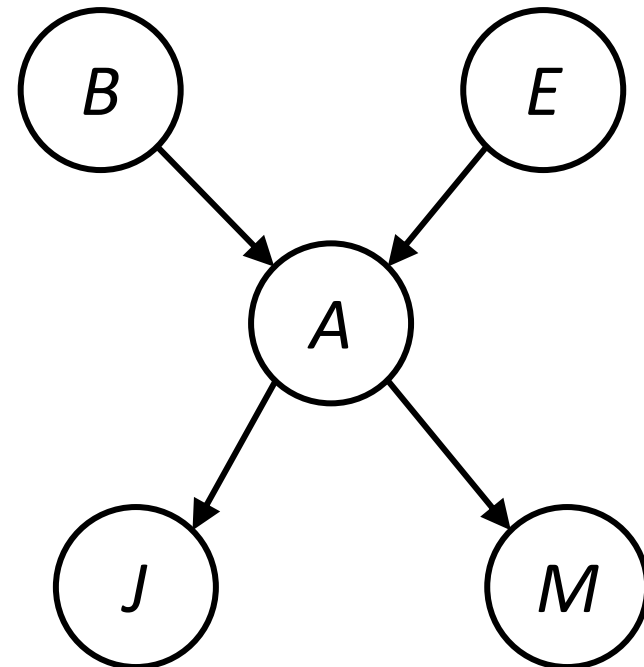
# Inference by Enumeration in Bayes Net

- The joint distribution can be computed from a BN by multiplying the conditional distributions

- Then we can do inference by enumeration
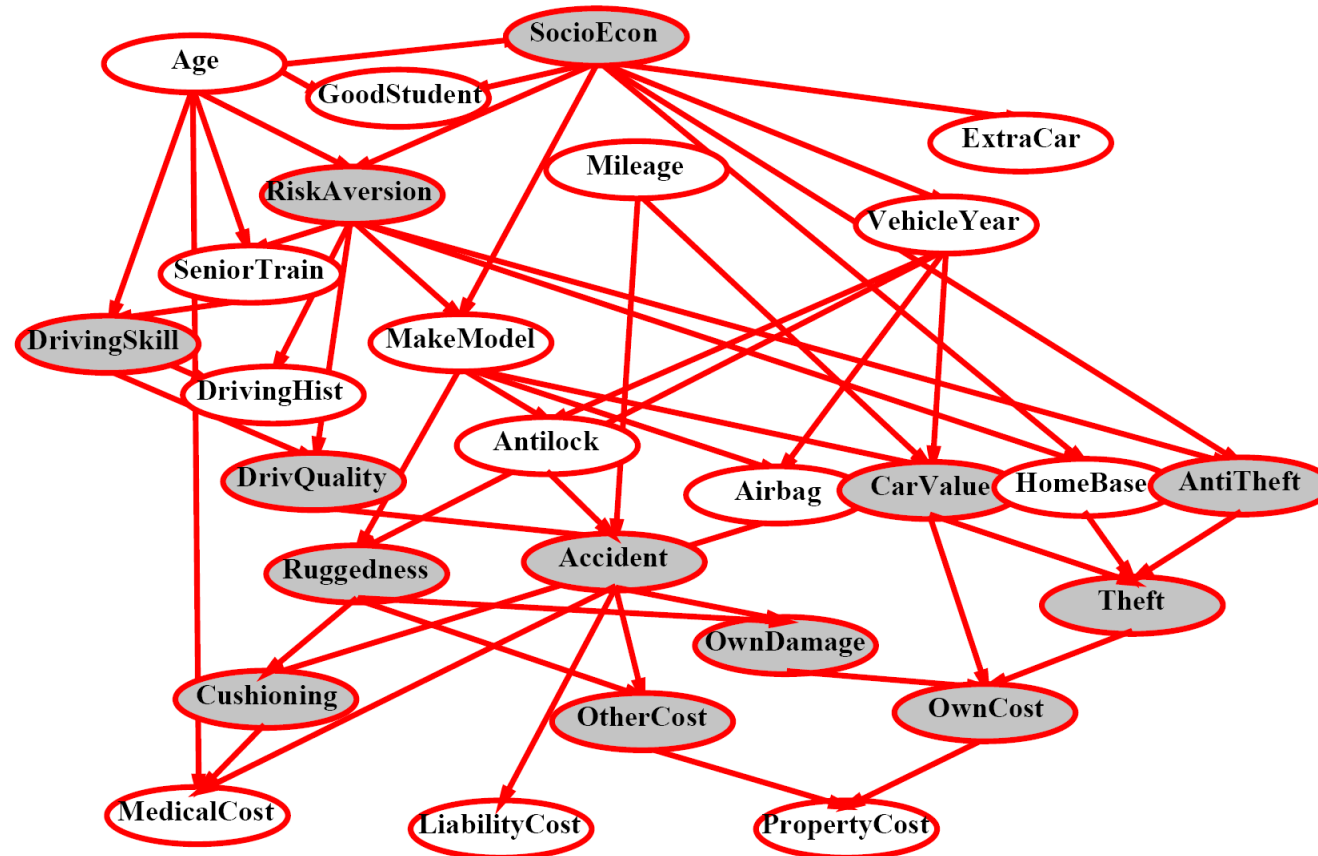
$$P(B \mid +j, +m) \quad \propto_B \quad P(B, +j, +m)$$

$$= \sum_{e,a} P(B, e, a, +j, +m)$$

$$= \sum_{e,a} P(B)P(e)P(a|B,e)P(+j|a)P(+m|a)$$
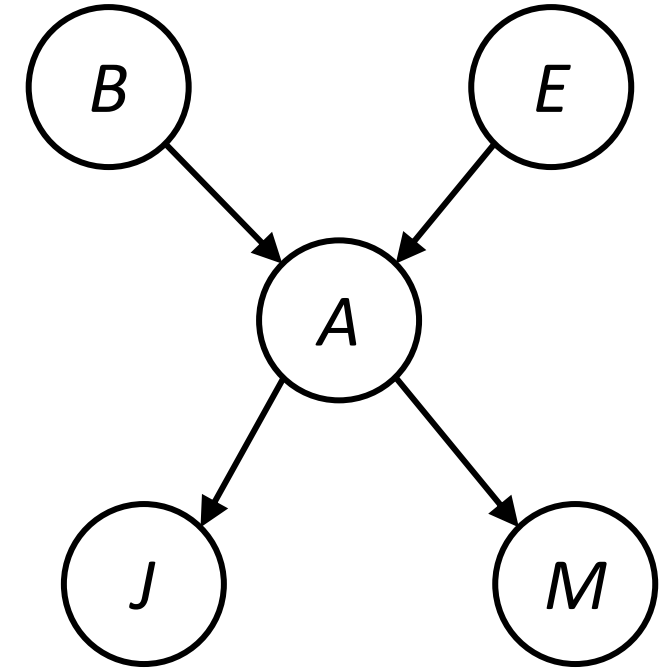
- Problem: sums of **exponentially many** products!

$$P(Antilock|observed\ variables) = ?$$

# Inference by Enumeration in Bayes Net

$$P(B \mid +j, +m) \propto_B P(B, +j, +m)$$

$$= \sum_{e,a} P(B, e, a, +j, +m)$$

$$= \sum_{e,a} P(B)P(e)P(a|B,e)P(+j|a)P(+m|a)$$

$$= P(B)P(+e)P(+a|B,+e)P(+j|+a)P(+m|+a) + P(B)P(+e)P(-a|B,+e)P(+j|-a)P(+m|-a)$$
$$P(B)P(-e)P(+a|B,-e)P(+j|+a)P(+m|+a) + P(B)P(-e)P(-a|B,-e)P(+j|-a)P(+m|-a)$$

Lots of repeated subexpressions!

# Can we do better?

- Consider **uwy + uwz + uxy + uxz + vwy + vwz + vxy +vxz**
  - 16 multiplies, 7 adds
  - Lots of repeated subexpressions!
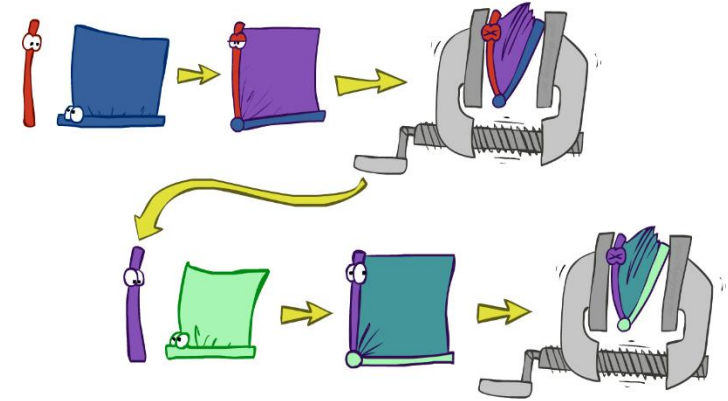- Rewrite as **(u+v)(w+x)(y+z)**
  - 2 multiplies, 3 adds

# Can we do better?

- $\sum_{e,a} P(B)\ P(e)\ P(a|B,e)\ P(j|a)\ P(m|a)$

- $= P(B)P(e)P(a|B,e)P(j|a)P(m|a) + P(B)P(\neg e)P(a|B,\neg e)P(j|a)P(m|a)$

  $+ P(B)P(e)P(\neg a|B,e)P(j|\neg a)P(m|\neg a) + P(B)P(\neg e)P(\neg a|B,\neg e)P(j|\neg a)P(m|\neg a)$
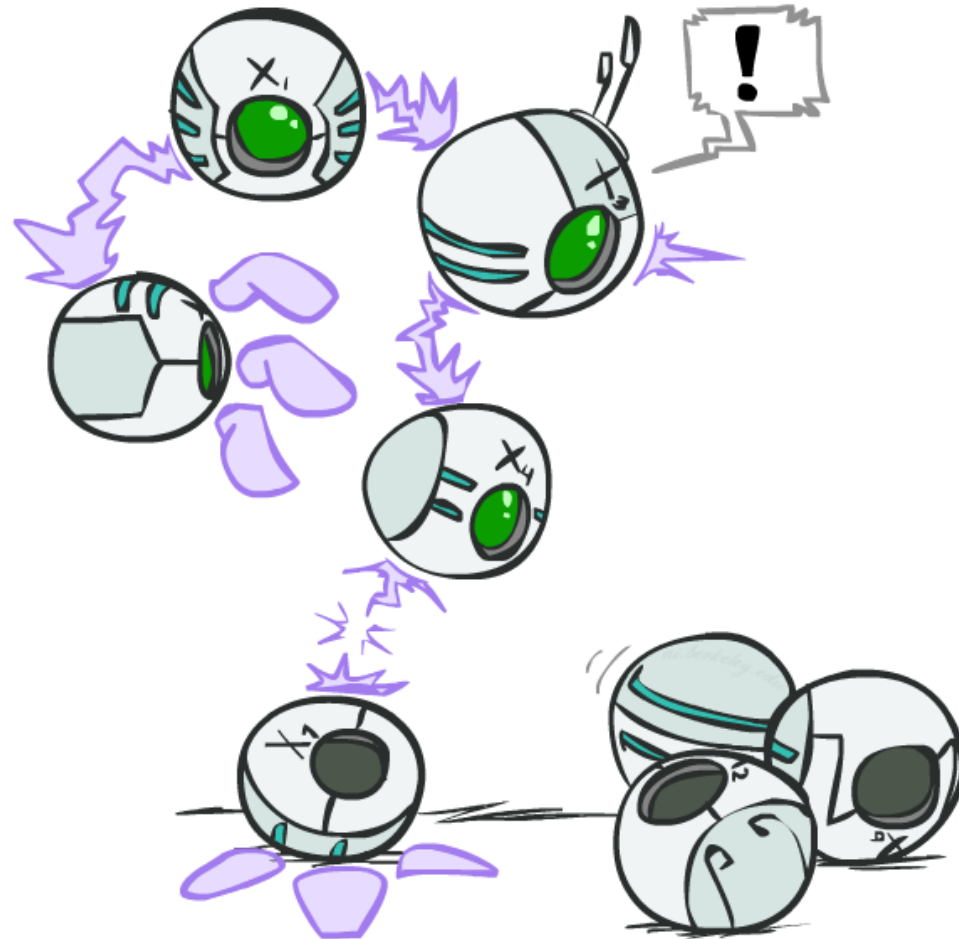
  Lots of repeated subexpressions!

# Variable elimination: The basic ideas

- Move summations inwards as far as possible

    - $P(B \mid j, m) = \alpha \sum_{e,a} P(B)\, P(e)\, P(a|B,e)\, P(j|a)\, P(m|a)$

    - $= \alpha\, P(B) \sum_{e} P(e) \sum_{a} P(a|B,e)\, P(j|a)\, P(m|a)$

- Do the calculation from the inside out

    - I.e., sum over *a* first, the sum over *e*

    - Problem: *P(a|B,e)* isn't a single number, it's a bunch of different numbers depending on the values of *B* and *e*

    - Solution: use arrays of numbers (of various dimensions) with appropriate operations on them; these are called ***factors***

# Operations on Factors

# Factors

- A factor is a multi-dimensional array to represent $P(Y_1 \dots Y_N \mid X_1 \dots X_M)$
  - If a variable is assigned (represented with lower-case), its dimension is missing (selected) from the array

  - Joint distribution: P(X,Y)
    - Entries P(x,y) for all x, y
    - Sums to 1

  - Selected joint: P(x,Y)
    - A slice of the joint distribution
    - Entries P(x,y) for fixed x, all y
    - Sums to P(x)

$P(T, W)$

| T | W | P |
|---|---|---|
| hot | sun | 0.4 |
| hot | rain | 0.1 |
| cold | sun | 0.2 |
| cold | rain | 0.3 |

$P(cold, W)$

| T | W | P |
|---|---|---|
| cold | sun | 0.2 |
| cold | rain | 0.3 |

# Factors

- A factor is a multi-dimensional array to represent P(Y$_1$ ... Y$_N$ | X$_1$ ... X$_M$)
  - If a variable is assigned (represented with lower-case), its dimension is missing (selected) from the array

  - Single conditional: P(Y | x)
    - Entries P(y | x) for fixed x, all y
    - Sums to 1

  - Family of conditionals:
  
  P(X |Y)
    - Multiple conditionals
    - Entries P(x | y) for all x, y
    - Sums to |Y|

$P(W|cold)$

| T | W | P |
|------|------|-----|
| cold | sun | 0.4 |
| cold | rain | 0.6 |

$P(W|T)$

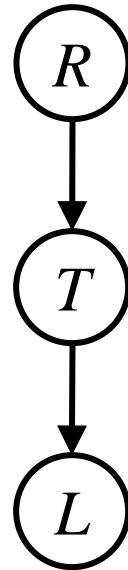| T | W | P |
|------|------|-----|
| hot | sun | 0.8 |
| hot | rain | 0.2 |
| cold | sun | 0.4 |
| cold | rain | 0.6 |

$P(W|hot)$

$P(W|cold)$

# Factors

- A **factor** is a multi-dimensional array to represent $P(Y_1 \ldots Y_N \mid X_1 \ldots X_M)$
  - If a variable is assigned (represented with lower-case), its dimension is missing (selected) from the array

  - Specified family: P( y | X )
    - Entries P(y | x) for fixed y, but for all x
    - Sums to … who knows!

$$P(rain|T)$$

| T | W | P |
|------|------|-----|
| hot | rain | 0.2 |
| cold | rain | 0.6 |

$P(rain|hot)$

$P(rain|cold)$

# Example: Traffic Domain

- **Random Variables**
  - R: Raining
  - T: Traffic
  - L: Late for class!

$P(R)$

| +r | 0.1 |
|----|-----|
| -r | 0.9 |

$P(T|R)$

| +r | +t | 0.8 |
|----|----|-----|
| +r | -t | 0.2 |
| -r | +t | 0.1 |
| -r | -t | 0.9 |

$P(L|T)$

| +t | +l | 0.3 |
|----|----|-----|
| +t | -l | 0.7 |
| -t | +l | 0.1 |
| -t | -l | 0.9 |

# Running Example: Traffic Domain

- Initial factors are local CPTs (one per node)

$P(R)$

| +r | 0.1 |
|----|-----|
| -r | 0.9 |

$P(T|R)$

| +r | +t | 0.8 |
|----|----|-----|
| +r | -t | 0.2 |
| -r | +t | 0.1 |
| -r | -t | 0.9 |

$P(L|T)$

| +t | +l | 0.3 |
|----|----|-----|
| +t | -l | 0.7 |
| -t | +l | 0.1 |
| -t | -l | 0.9 |

- Any known values are selected
  - E.g. if we know $L = +\ell$, the initial factors are

$P(R)$

| +r | 0.1 |
|----|-----|
| -r | 0.9 |

$P(T|R)$

| +r | +t | 0.8 |
|----|----|-----|
| +r | -t | 0.2 |
| -r | +t | 0.1 |
| -r | -t | 0.9 |

$P(+\ell|T)$

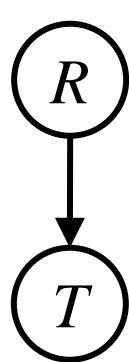| +t | +l | 0.3 |
|----|----|-----|
| -t | +l | 0.1 |

# Operation 1: Join Factors

- First basic operation: joining factors
    - Just like a database join
    - Given multiple factors, build a new factor over the union of the variables involved
    - Each entry is computed by pointwise products

- Example:

$$P(R) \quad \times \quad P(T|R) \quad \Longrightarrow \quad P(R,T)$$

$R$

| +r | 0.1 |
|----|-----|
| -r | 0.9 |

| +r | +t | 0.8 |
|----|----|-----|
| +r | -t | 0.2 |
| -r | +t | 0.1 |
| -r | -t | 0.9 |

$T$

| +r | +t | 0.08 |
|----|----|------|
| +r | -t | 0.02 |
| -r | +t | 0.09 |
| -r | -t | 0.81 |

$R,T$

$$\forall r,t: \quad P(r,t) = P(r) \cdot P(t|r)$$

# Operation 2: Eliminate

- Second basic operation: eliminating a variable

  - Take a factor and sum out (marginalize) a variable

- Example:

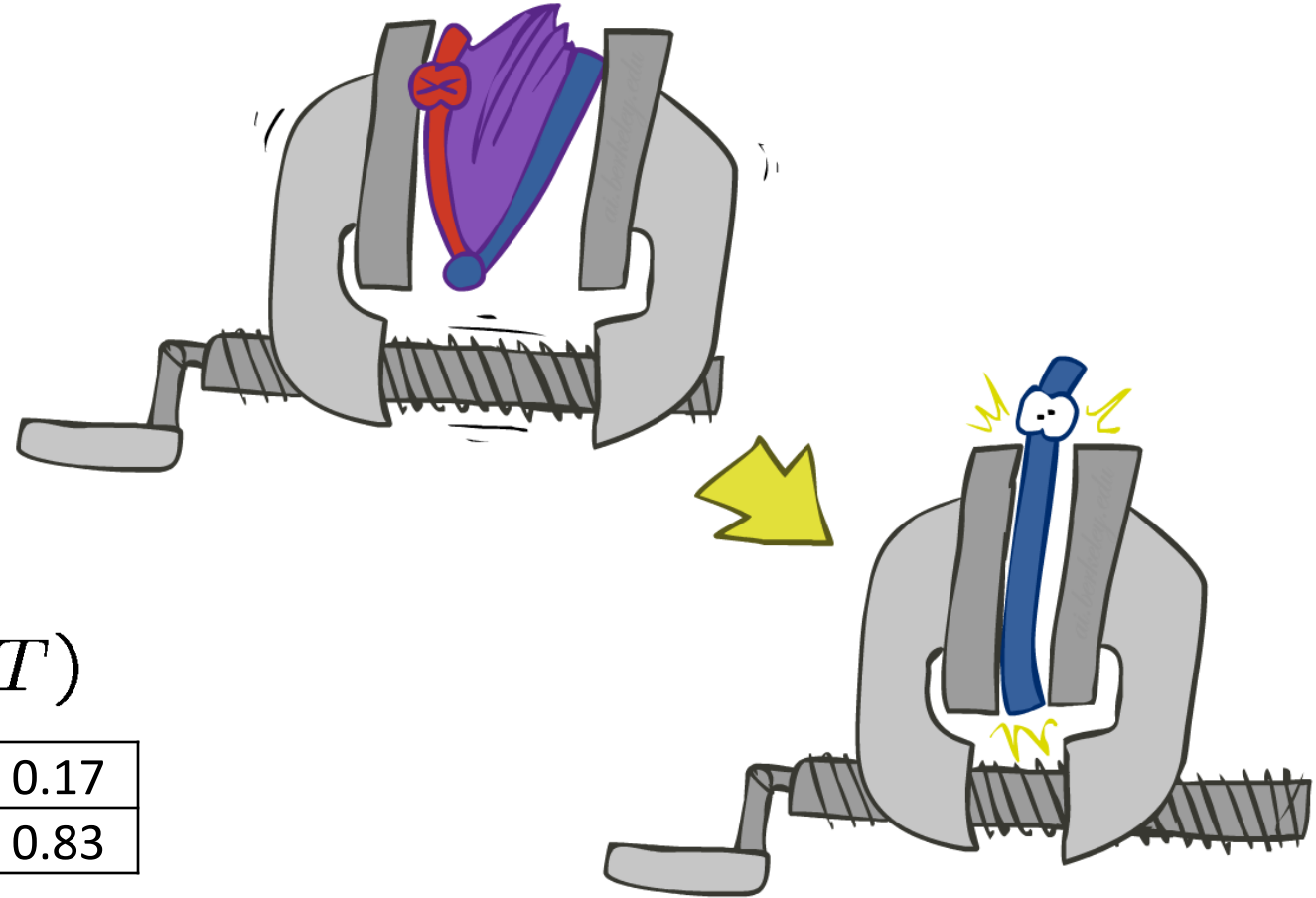$$P(R, T)$$

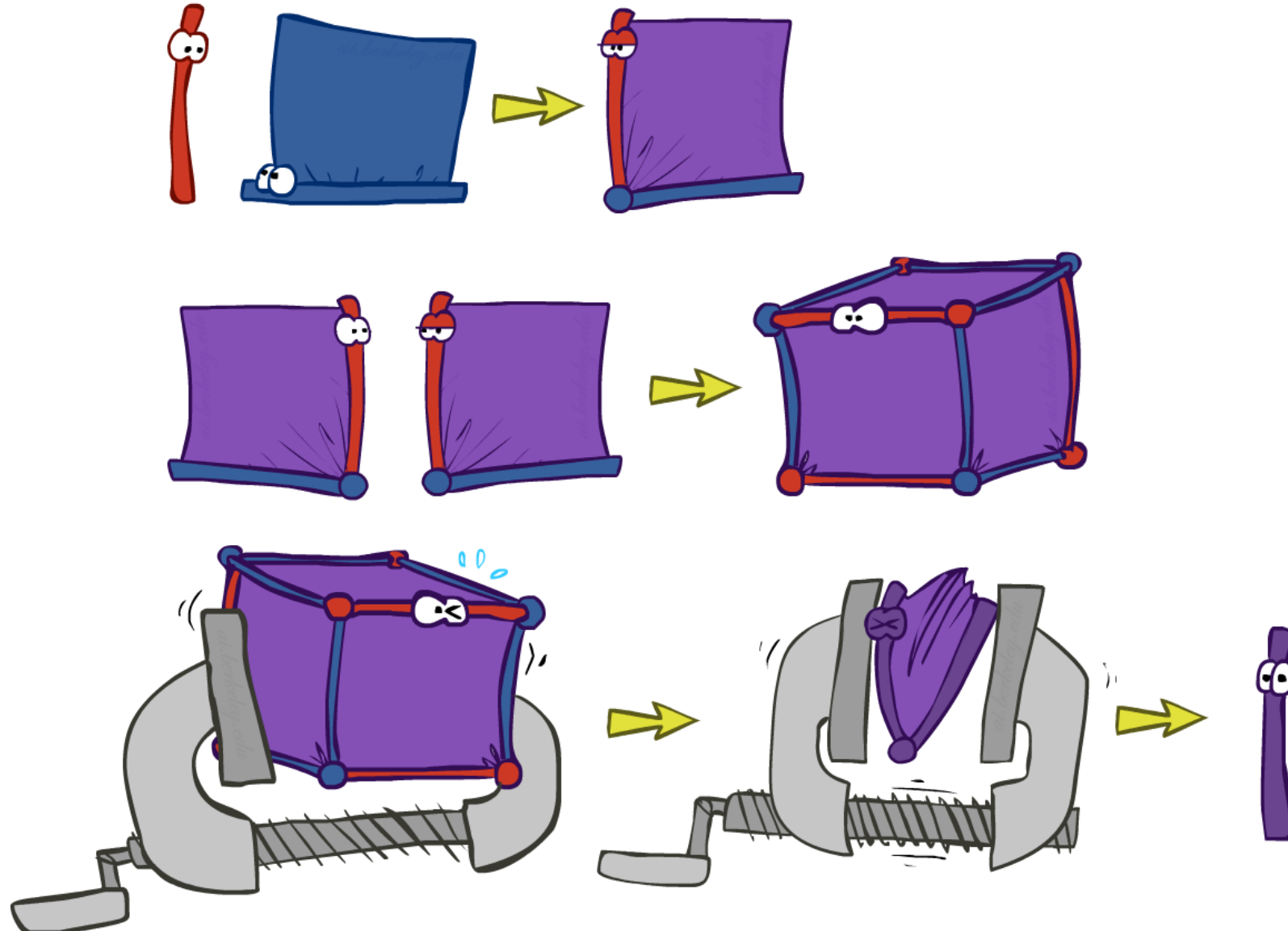| | | |
|---|---|---|
| +r | +t | 0.08 |
| +r | -t | 0.02 |
| -r | +t | 0.09 |
| -r | -t | 0.81 |

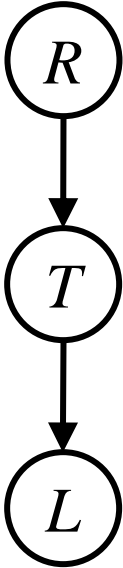sum $R$ $\longrightarrow$

$$P(T)$$

| | |
|---|---|
| +t | 0.17 |
| -t | 0.83 |

# Inference by Enumeration in BN = Multiple Join + Multiple Eliminate

# Computing *P*(*L*): Multiple Joins



$P(R)$

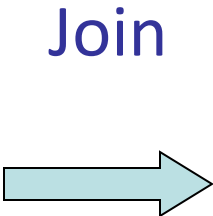| +r | 0.1 |
|---|---|
| -r | 0.9 |

$P(T|R)$

| +r | +t | 0.8 |
|---|---|---|
| +r | -t | 0.2 |
| -r | +t | 0.1 |
| -r | -t | 0.9 |

$P(L|T)$

| +t | +l | 0.3 |
|---|---|---|
| +t | -l | 0.7 |
| -t | +l | 0.1 |
| -t | -l | 0.9 |

Join →

$P(R,T)$

| +r | +t | 0.08 |
|---|---|---|
| +r | -t | 0.02 |
| -r | +t | 0.09 |
| -r | -t | 0.81 |

$P(L|T)$

| +t | +l | 0.3 |
|---|---|---|
| +t | -l | 0.7 |
| -t | +l | 0.1 |
| -t | -l | 0.9 |

Join →

$R, T, L$

$P(R,T,L)$

| +r | +t | +l | 0.024 |
|---|---|---|---|
| +r | +t | -l | 0.056 |
| +r | -t | +l | 0.002 |
| +r | -t | -l | 0.018 |
| -r | +t | +l | 0.027 |
| -r | +t | -l | 0.063 |
| -r | -t | +l | 0.081 |
| -r | -t | -l | 0.729 |

# Computing $P(L)$: Multiple Elimination

$R, T, L$      $T, L$      $L$

$P(R, T, L)$

| | | | |
|---|---|---|---|
| +r | +t | +l | 0.024 |
| +r | +t | -l | 0.056 |
| +r | -t | +l | 0.002 |
| +r | -t | -l | 0.018 |
| -r | +t | +l | 0.027 |
| -r | +t | -l | 0.063 |
| -r | -t | +l | 0.081 |
| -r | -t | -l | 0.729 |

**Sum out R**

$P(T, L)$

| | | |
|---|---|---|
| +t | +l | 0.051 |
| +t | -l | 0.119 |
| -t | +l | 0.083 |
| -t | -l | 0.747 |

**Sum out T**

$P(L)$

| | |
|---|---|
| +l | 0.134 |
| -l | 0.866 |

*A factor of exponential size!*

# Variable Elimination

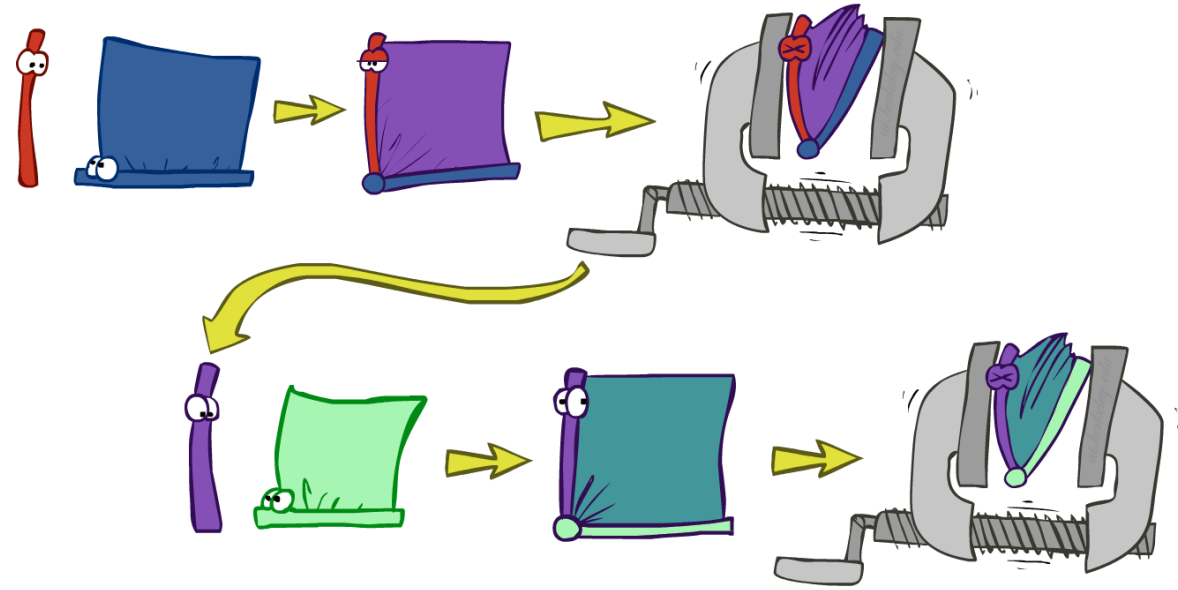# Inference by Enumeration vs. Variable Elimination

- **Why is inference by enumeration so slow?**
  - You join up the whole joint distribution before you sum out the hidden variables

- **Idea: interleave joining and elimination!**
  - Called "Variable Elimination"
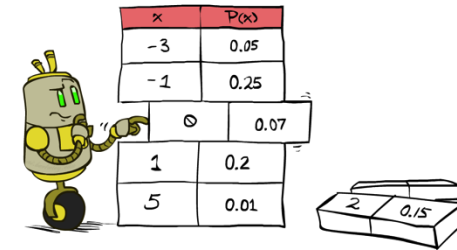  - Still NP-hard, but usually much faster than inference by enumeration

# Variable Elimination

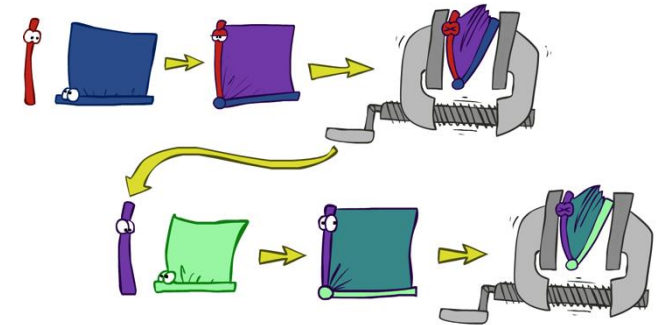- Query: $P(Q|E_1 = e_1, \ldots E_k = e_k)$

- Start with initial factors:
  - Local CPTs (but instantiated by evidence)
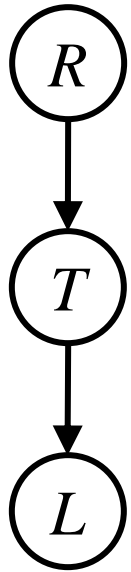
- While there are still hidden variables (not Q or evidence):
  - Pick a hidden variable H
  - Join all factors mentioning H
  - Eliminate (sum out) H

- Join all remaining factors and normalize

# Traffic Domain

$$P(L) = ?$$

$R \rightarrow T \rightarrow L$

- **Inference by Enumeration**

$$= \sum_t \sum_r P(L|t)P(r)P(t|r)$$

Join on r

Join on t

Eliminate r

Eliminate t

- **Variable Elimination**

$$= \sum_t P(L|t) \sum_r P(r)P(t|r)$$

Join on r

Eliminate r

Join on t

Eliminate t

# Variable Elimination

$P(R)$

| +r | 0.1 |
|----|-----|
| -r | 0.9 |

$P(R, T)$

| +r | +t | 0.08 |
|----|----|------|
| +r | -t | 0.02 |
| -r | +t | 0.09 |
| -r | -t | 0.81 |

$P(T)$

| +t | 0.17 |
|----|------|
| -t | 0.83 |

$P(T|R)$

| +r | +t | 0.8 |
|----|----|-----|
| +r | -t | 0.2 |
| -r | +t | 0.1 |
| -r | -t | 0.9 |

$R, T$

$T$

$T, L$

$L$

$L$

$L$

$P(L|T)$

| +t | +l | 0.3 |
|----|----|-----|
| +t | -l | 0.7 |
| -t | +l | 0.1 |
| -t | -l | 0.9 |

$P(L|T)$

| +t | +l | 0.3 |
|----|----|-----|
| +t | -l | 0.7 |
| -t | +l | 0.1 |
| -t | -l | 0.9 |

$P(L|T)$

| +t | +l | 0.3 |
|----|----|-----|
| +t | -l | 0.7 |
| -t | +l | 0.1 |
| -t | -l | 0.9 |

$P(T, L)$

| +t | +l | 0.051 |
|----|----|-------|
| +t | -l | 0.119 |
| -t | +l | 0.083 |
| -t | -l | 0.747 |

$P(L)$

| +l | 0.134 |
|----|-------|
| -l | 0.866 |

# Evidence

- If evidence, start with factors that select that evidence
  - No evidence uses these initial factors:

$P(R)$

| +r | 0.1 |
|----|-----|
| -r | 0.9 |

$P(T|R)$

| +r | +t | 0.8 |
|----|----|-----|
| +r | -t | 0.2 |
| -r | +t | 0.1 |
| -r | -t | 0.9 |

$P(L|T)$

| +t | +l | 0.3 |
|----|----|-----|
| +t | -l | 0.7 |
| -t | +l | 0.1 |
| -t | -l | 0.9 |

  - Computing $P(L| + r)$ the initial factors become:

$P(+r)$

| +r | 0.1 |
|----|-----|

$P(T| + r)$

| +r | +t | 0.8 |
|----|----|-----|
| +r | -t | 0.2 |

$P(L|T)$

| +t | +l | 0.3 |
|----|----|-----|
| +t | -l | 0.7 |
| -t | +l | 0.1 |
| -t | -l | 0.9 |

- We eliminate all vars other than query + evidence

# Evidence II

- Result will be a selected joint of query and evidence
  - E.g. for P(L | +r), we would end up with:

$$P(+r, L)$$

| +r | +l | 0.026 |
|----|----|-------|
| +r | -l | 0.074 |

Normalize

$$P(L| + r)$$

| +l | 0.26 |
|----|------|
| -l | 0.74 |

- To get our answer, just normalize this!

- That's it!

# Example

$$P(B|j,m) \propto P(B,j,m)$$

$$\boxed{P(B) \qquad P(E) \qquad P(A|B,E) \qquad P(j|A) \qquad P(m|A)}$$

Choose A

$$P(A|B,E)$$
$$P(j|A) \qquad \times \qquad \Sigma \qquad P(j,m|B,E)$$
$$P(m|A)$$

$$\boxed{P(B) \qquad P(E) \qquad P(j,m|B,E)}$$

# Example

$$P(B) \qquad P(E) \qquad P(j,m|B,E)$$

Choose E

$$P(E)$$
$$P(j,m|B,E)$$

$\times$ → $\Sigma$ → $P(j,m|B)$

$$P(B) \qquad\qquad P(j,m|B)$$

Finish with B

$$P(B)$$
$$P(j,m|B)$$

$\times$ → $P(j,m,B)$ → Normalize → $P(B|j,m)$

# Same Example in Equations

$$P(B|j,m) \propto P(B,j,m)$$

$$\boxed{P(B) \qquad P(E) \qquad P(A|B,E) \qquad P(j|A) \qquad P(m|A)}$$

$$
\begin{aligned}
P(B|j,m) \;\propto\; & P(B,j,m) \\
= & \sum_{e,a} P(B,j,m,e,a) \\
= & \sum_{e,a} P(B)P(e)P(a|B,e)P(j|a)P(m|a) \\
= & \sum_{e} P(B)P(e) \sum_{a} P(a|B,e)P(j|a)P(m|a) \\
= & \sum_{e} P(B)P(e)f_1(B,e,j,m) \\
= & P(B) \sum_{e} P(e)f_1(B,e,j,m) \\
= & P(B)f_2(B,j,m)
\end{aligned}
$$

marginal can be obtained from joint by summing out

use Bayes' net joint distribution expression

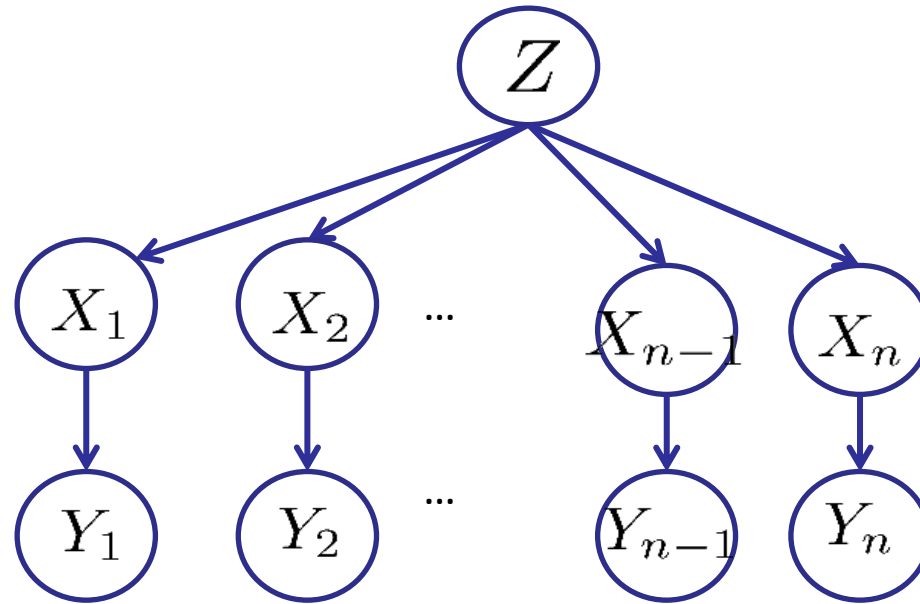use x*(y+z) = xy + xz

joining on a, and then summing out gives $f_1$

use x*(y+z) = xy + xz

joining on e, and then summing out gives $f_2$

**All we are doing is exploiting uwy + uwz + uxy + uxz + vwy + vwz + vxy +vxz = (u+v)(w+x)(y+z) to improve computational efficiency!**
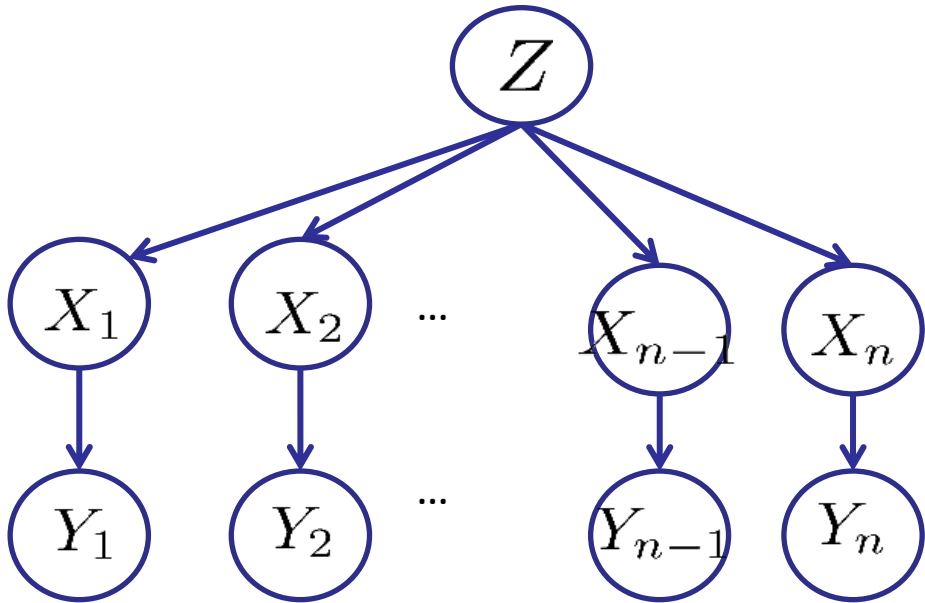
# Variable Elimination Ordering

- Query: $P(X_n | y_1, ..., y_n)$
- Two different orderings: $Z, X_1, ..., X_{n-1}$ and $X_1, ..., X_{n-1}, Z$.
- What is the size of the maximum factor generated for each of the orderings?

# Variable Elimination Ordering

- Z, X$_1$, …, X$_{n-1}$

$$P(Z)P(X_1|Z)P(X_2|Z),\ldots,P(X_n|Z)$$

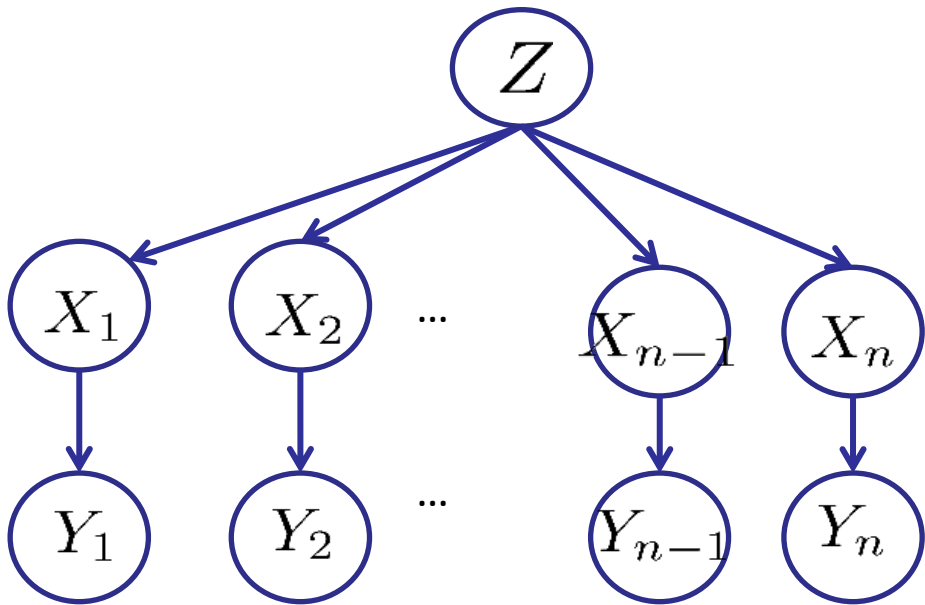$$f_1(X_1, X_2, \ldots, X_n)$$

$$f_2(y_1, X_2, \ldots, X_n)$$

$$f_3(y_1, y_2, \ldots, X_n)$$

$$2^{n+1}$$

# Variable Elimination Ordering

- $X_1, ..., X_{n-1}, Z$



$$P(X_1|Z)P(y_1|X_1)$$

$$\Downarrow$$

$$f_1(Z, y_1)$$

$$\vdots$$

$$f_1(Z, y_1), f_2(Z, y_2), \ldots, f_{n-1}(Z, y_{n-1}), P(Z), P(X_n|Z)$$

$$\Downarrow$$

$$f_n(X_n, y_1, ..., y_{n-1})$$

$$\vdots$$

$$2^2$$

# VE: Computational Complexity

- The size of the largest factor determines the time and space complexity of VE

- The elimination ordering can greatly affect the size of the largest factor.
  - E.g., previous slide's example $2^{n+1}$ vs. $2^2$


- Does there always exist an ordering that only results in small factors?
  - No!

# Reduction from 3SAT

$$(x_1 \lor x_2 \lor \neg x_3) \land (\neg x_1 \lor x_3 \lor \neg x_4) \land (x_2 \lor \neg x_2 \lor x_4) \land (\neg x_3 \lor \neg x_4 \lor \neg x_5) \land (x_2 \lor x_5 \lor x_7) \land (x_4 \lor x_5 \lor x_6) \land (\neg x_5 \lor x_6 \lor \neg x_7) \land (\neg x_5 \lor \neg x_6 \lor x_7)$$

$P(X_i = 0) = P(X_i = 1) = 0.5$

$Y_1 = X_1 \lor X_2 \lor \neg X_3$

...

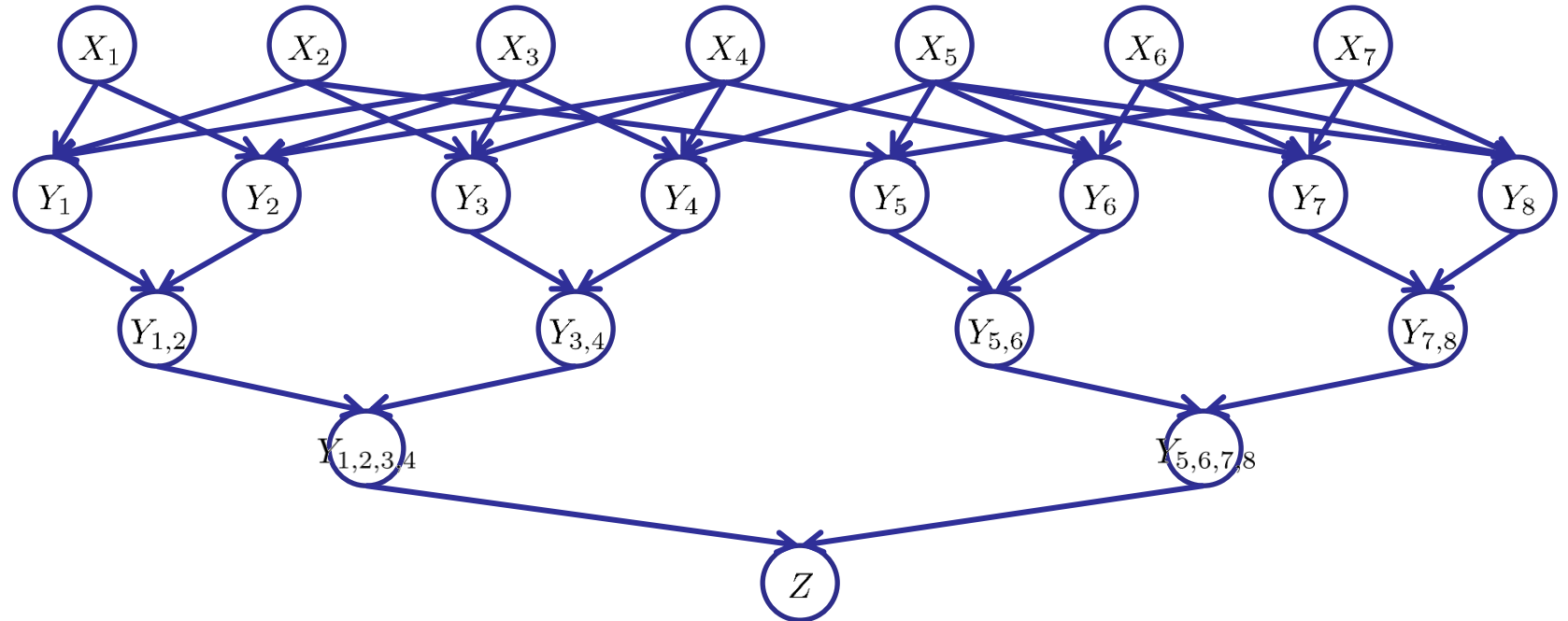$Y_8 = \neg X_5 \lor X_6 \lor X_7$

$Y_{1,2} = Y_1 \land Y_2$

...

$Y_{7,8} = Y_7 \land Y_8$

$Y_{1,2,3,4} = Y_{1,2} \land Y_{3,4}$

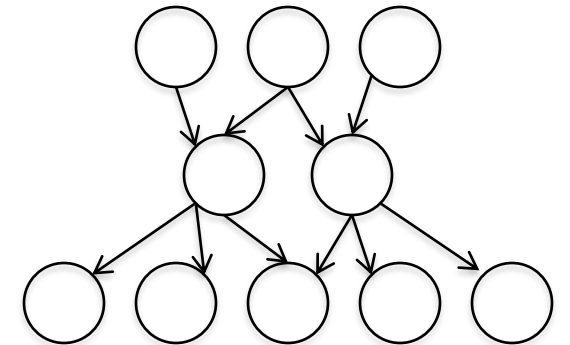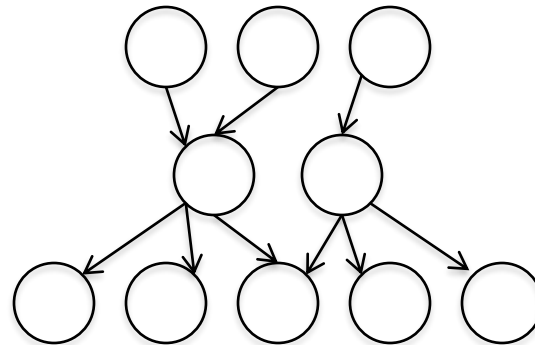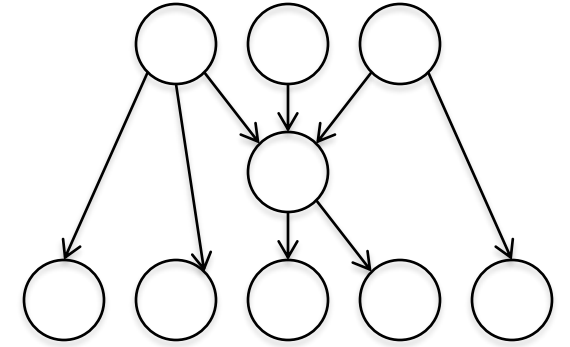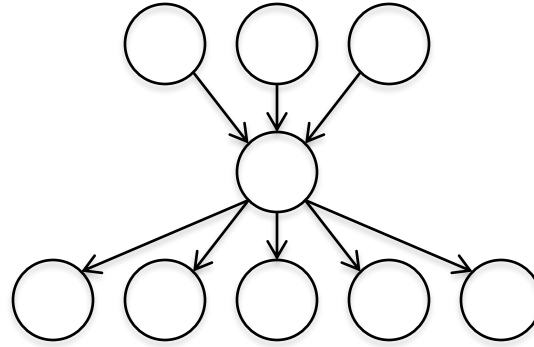$Y_{5,6,7,8} = Y_{5,6} \land Y_{7,8}$

$Z = Y_{1,2,3,4} \land Y_{5,6,7,8}$



- If we can answer P(z) equal to zero or not, we answered whether the 3-SAT problem has a solution.

- Hence inference in Bayes' nets is NP-hard. No known efficient probabilistic inference in general.
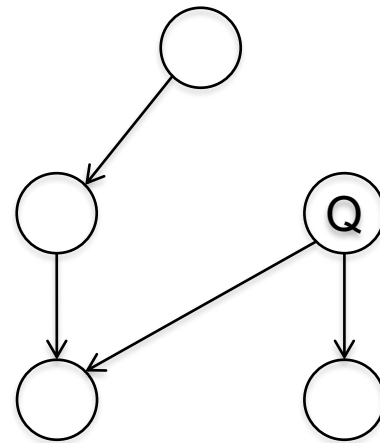
# Polytrees

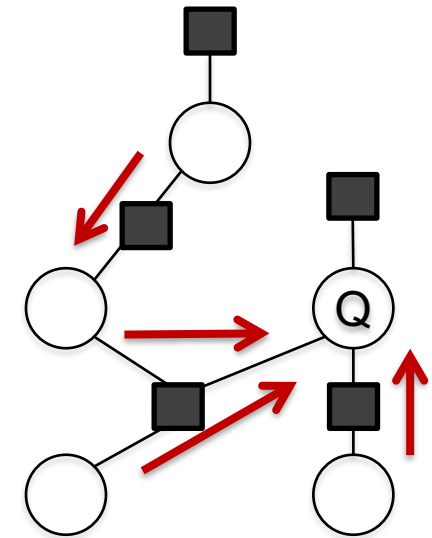- A polytree is a directed graph with no undirected cycles

# Variable Elimination on Polytrees

- For poly-tree BNs, the complexity of VE is ***linear in the BN size*** (number of CPT entries) with the following elimination ordering:
  - Convert to a factor graph
  - Take Q as the root
  - Eliminate from the leaves towards the root
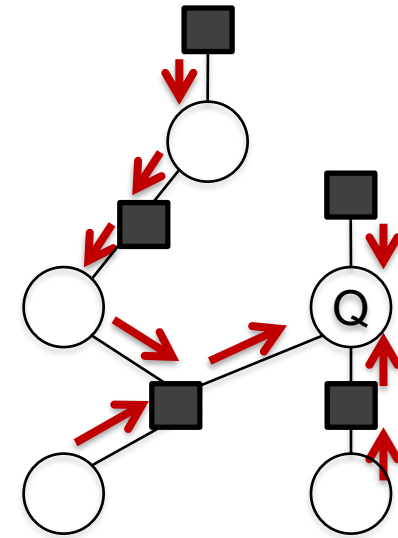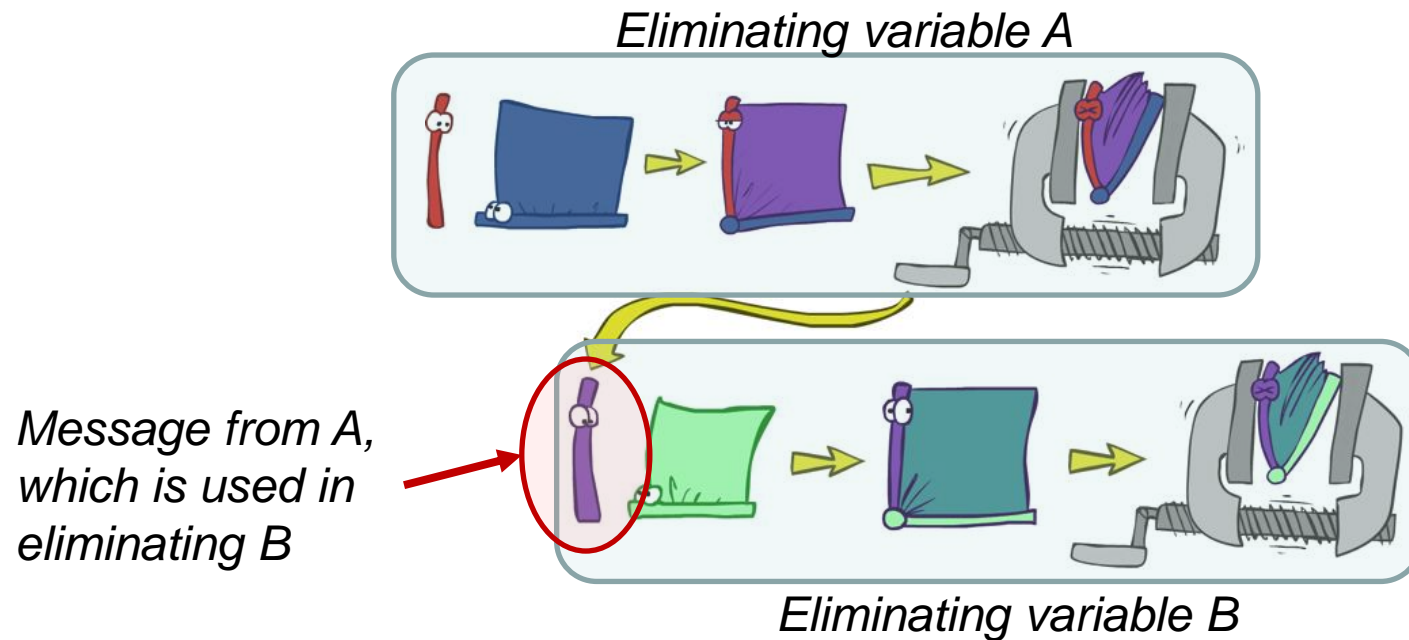


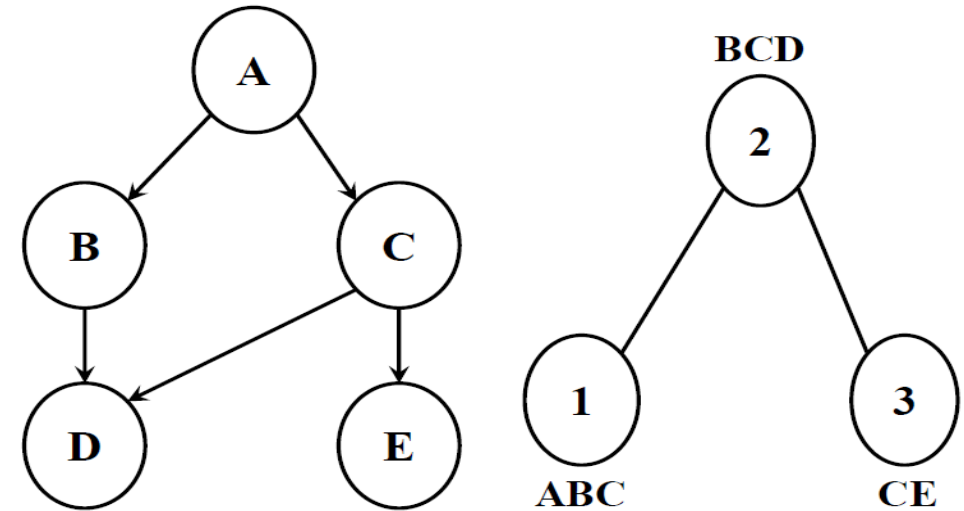Bayesian Network                    Factor Graph

# Variable Elimination on Polytrees

- VE for poly-tree BNs is equivalent to
  - Sum-product message passing algorithm or belief propagation algorithm
  (i.e., passing messages/beliefs from leaf nodes to the root node)

*Eliminating variable A*



*Message from A, which is used in eliminating B*
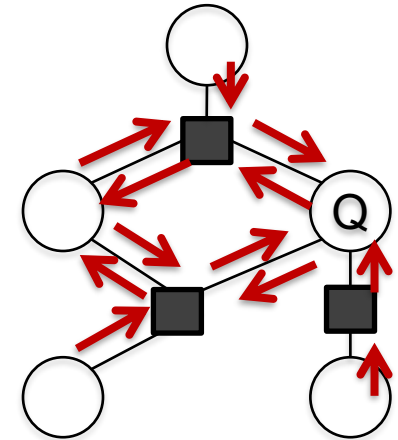
*Eliminating variable B*

# Message Passing on General Graphs

- Exact inference: Junction Tree Algorithm

  - Group individual nodes to form cluster nodes in such a way that the resulting network is a polytree (called a junction tree or join tree)

  - Run a sum-product-like algorithm on the junction tree.

  - *Intractable* on graphs with large cliques (i.e., large tree-width).

# Message Passing on General Graphs

- Approximate inference: Loopy Belief Propagation
  - Simply pass the messages on the general graph
    - Will not terminate with loops
    - Run until convergence (not guaranteed!)
  - *Approximate* but *tractable* for large graphs.
  - Sometime works well, sometimes not at all.

# Summary

- Exact inference of Bayesian networks
  - Enumeration
    - exponential complexity
  - Variable Eliminating
    - worst-case exponential complexity, often better
  - VE on polytrees
    - linear complexity
    - = message passing
  - Message passing on general graphs
    - Junction Tree Algorithm
    - Loopy Belief Propagation: no longer exact