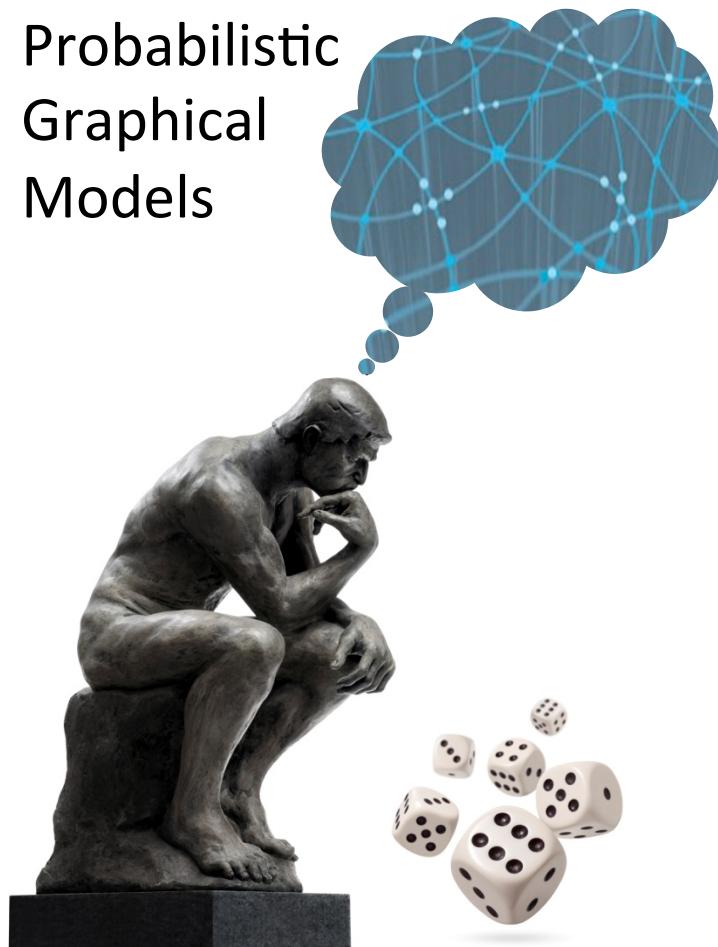


Probabilistic  
Graphical  
Models



Inference

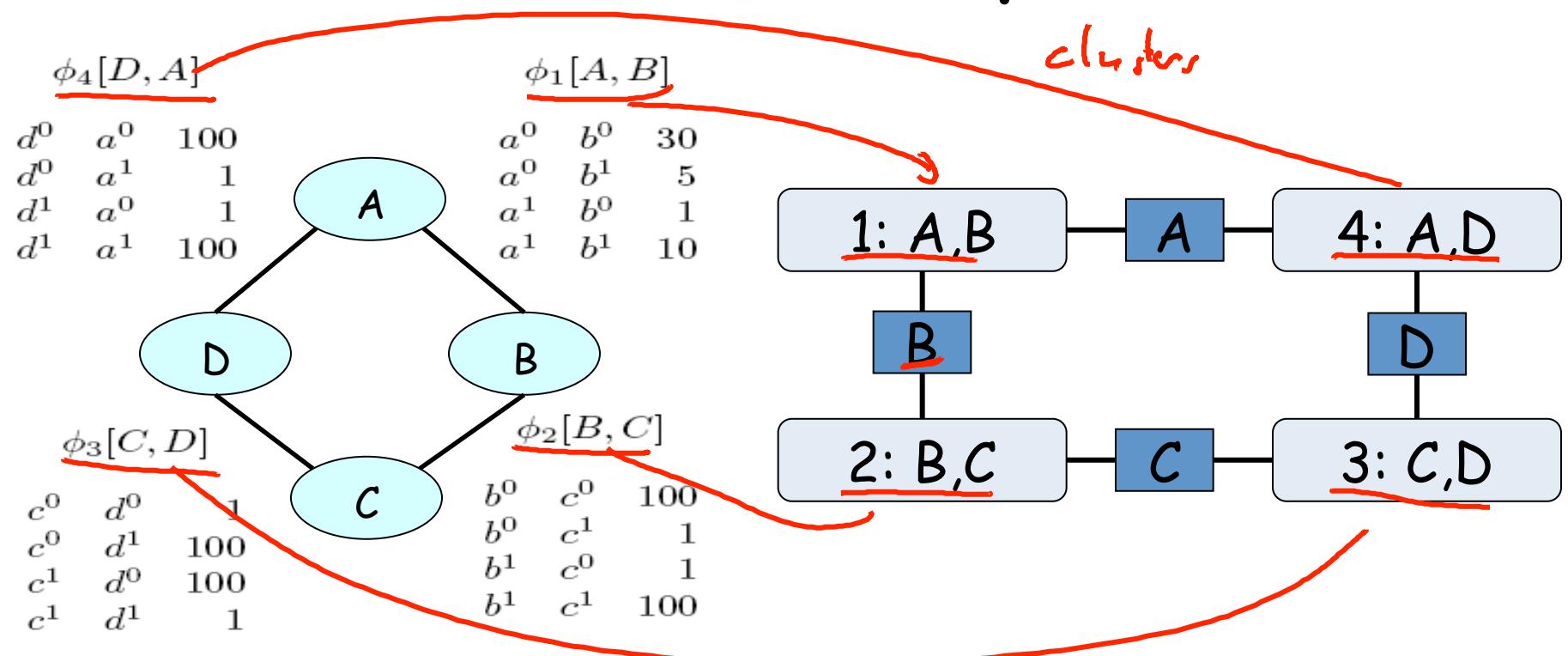
---

Message Passing

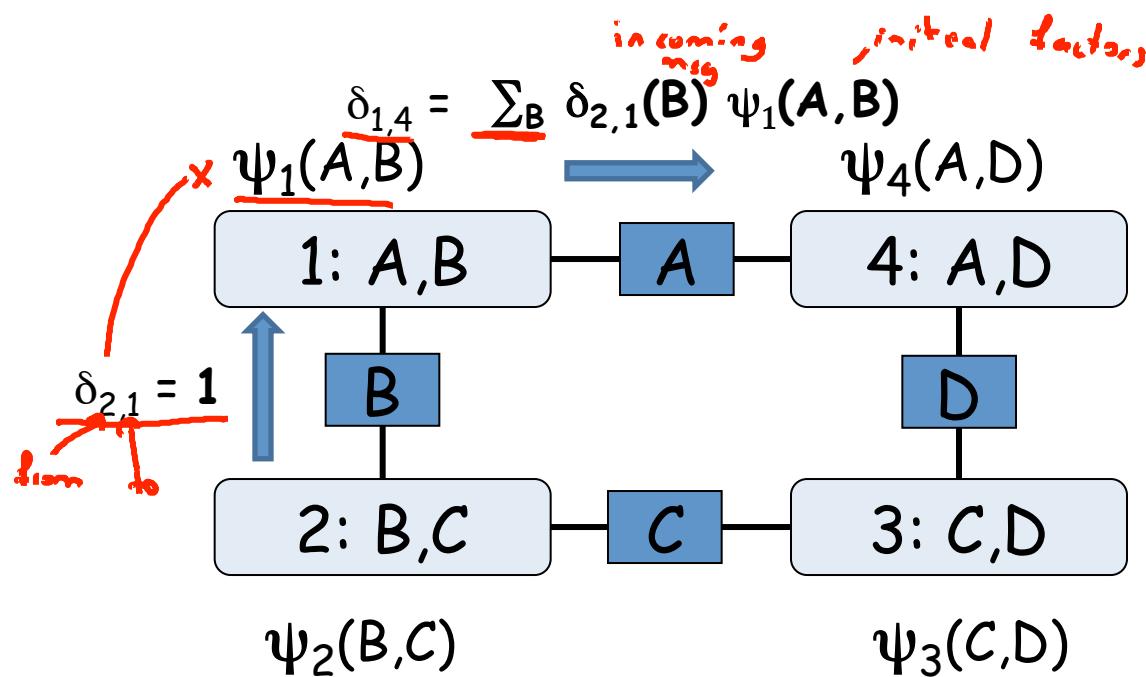
---

Belief  
Propagation  
Algorithm

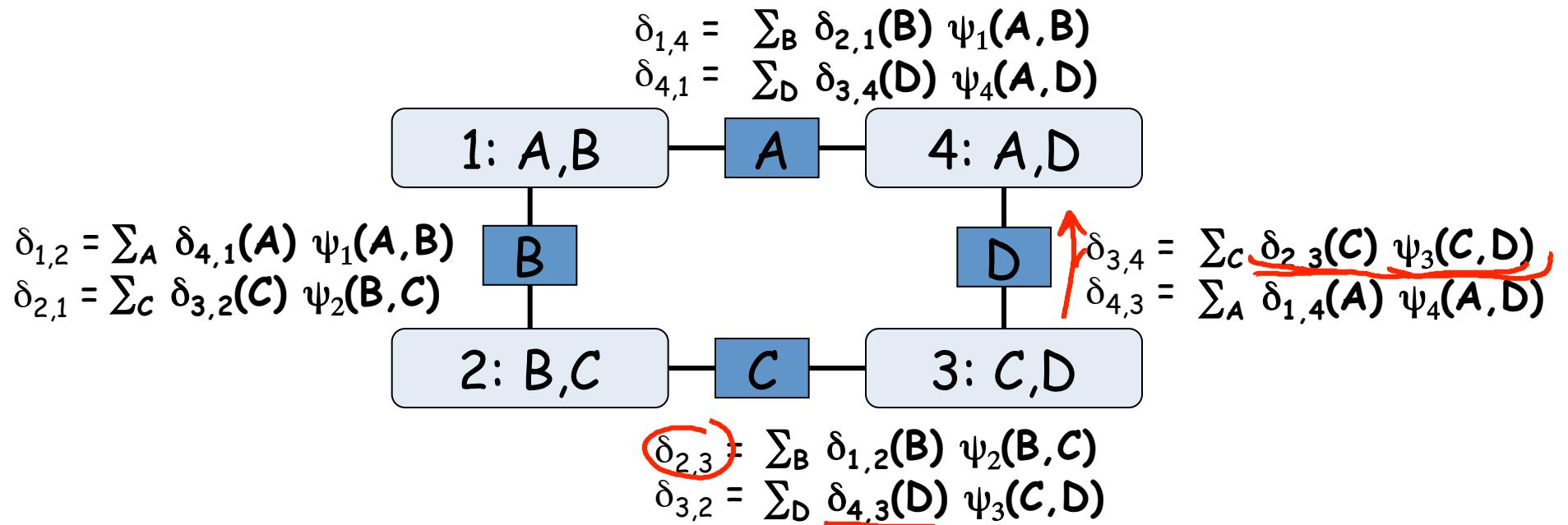
# Cluster Graph



# Passing Messages



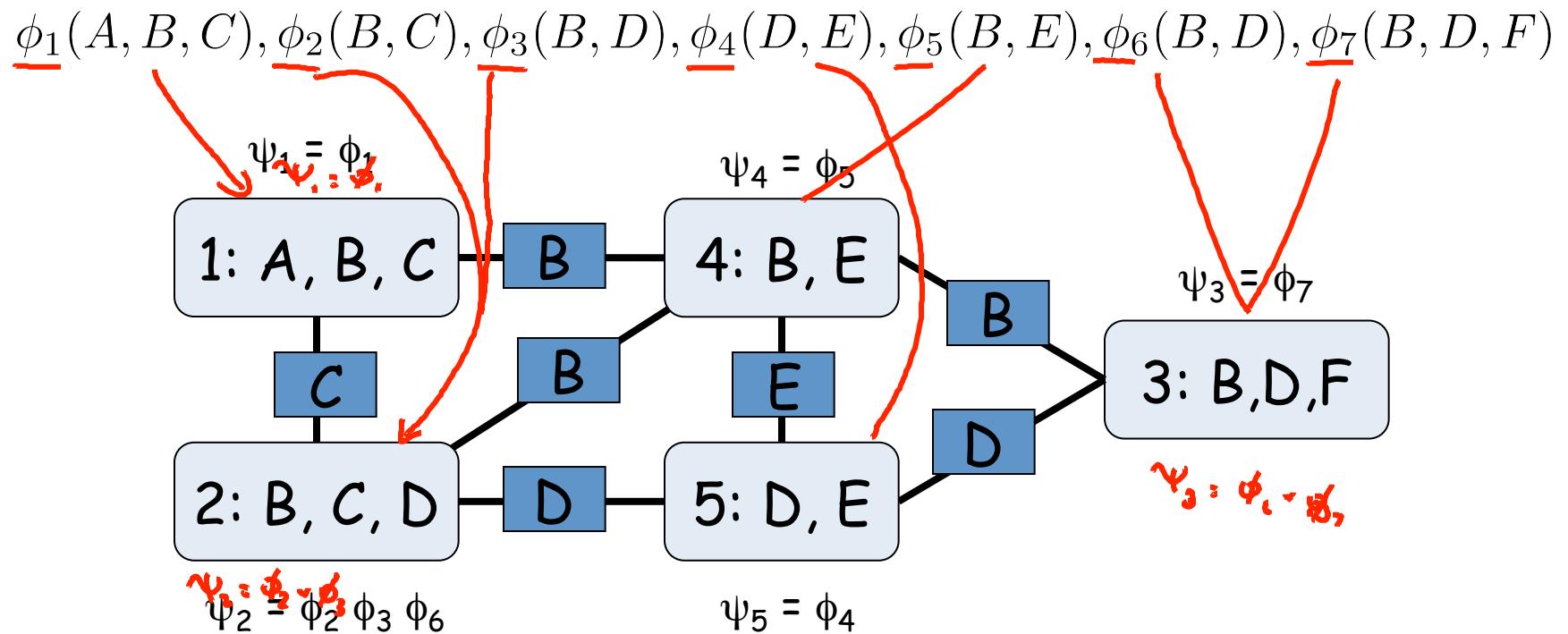
# Passing Messages



# Cluster Graphs

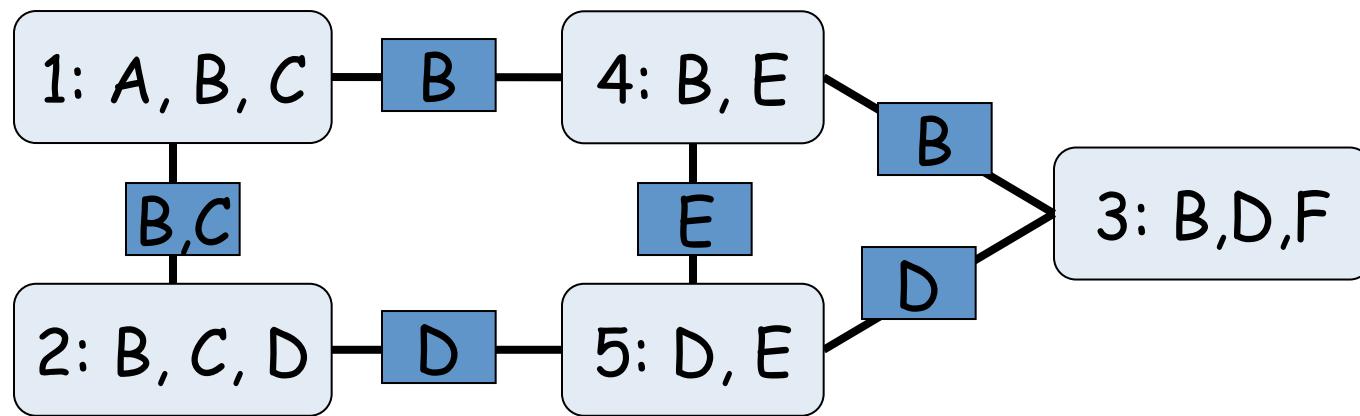
- Undirected graph such that:
  - nodes are clusters  $C_i \subseteq \{X_1, \dots, X_n\}$  *Subsets of variables*
  - edge between  $C_i$  and  $C_j$  associated with sepset  $S_{i,j} \subseteq C_i \cap C_j$  *Variables that they talk about*
- Given set of factors  $\Phi$ , we assign each  $\phi_k$  to a cluster  $C_{\alpha(k)}$  s.t. Scope [ $\phi_k$ ]  $\subseteq C_{\alpha(k)}$
- Define  $\psi_i(C_i) = \prod_{k: \alpha(k)=i} \phi_k$  *all factors assigned to it*

# Example Cluster Graph



# Different Cluster Graph

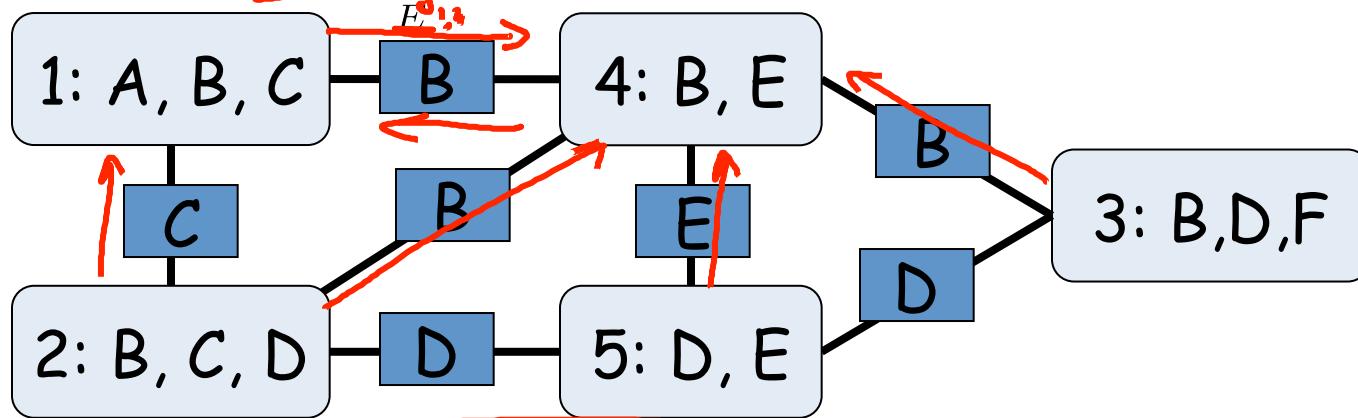
$\phi_1(A, B, C), \phi_2(B, C), \phi_3(B, D), \phi_4(D, E), \phi_5(B, E), \phi_6(B, D), \phi_7(B, D, F)$



# Message Passing

$$\delta_{1 \rightarrow 4}(B) = \sum_{A,C} \psi_1(A, B, C) \delta_{2 \rightarrow 1}(C)$$

$$\delta_{4 \rightarrow 1}(B) = \sum_{E} \psi_4(B, E) \times \delta_{2 \rightarrow 4}(B) \times \delta_{5 \rightarrow 4}(E) \times \delta_{3 \rightarrow 4}(B)$$



$$\delta_{i \rightarrow j}(S_{i,j}) = \sum_{C_i - S_{i,j}} \psi_i \times \prod_{k \in (\mathcal{N}_i - \{j\})} \delta_{k \rightarrow i}$$

incoming msgs  
only from  $j$ ;

Daphne Koller

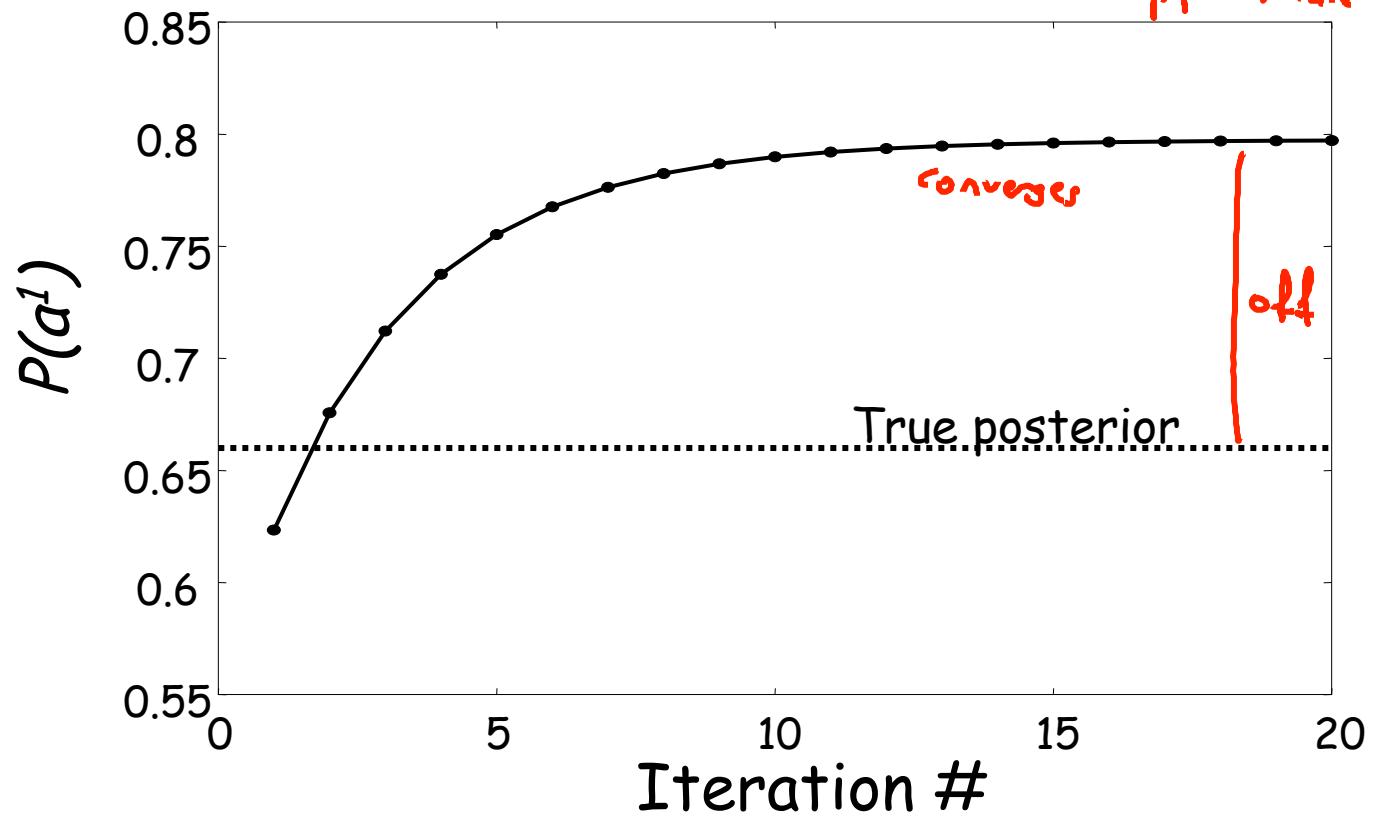
# Belief Propagation Algorithm

- Assign each factor  $\phi_k \in \Phi$  to a cluster  $C_{\alpha(k)}$
- Construct initial potentials  $\psi_i(C_i) = \prod_{k:\alpha(k)=i} \phi_k$
- Initialize all messages to be 1
- Repeat until when?
  - Select edge  $(i,j)$  and pass message

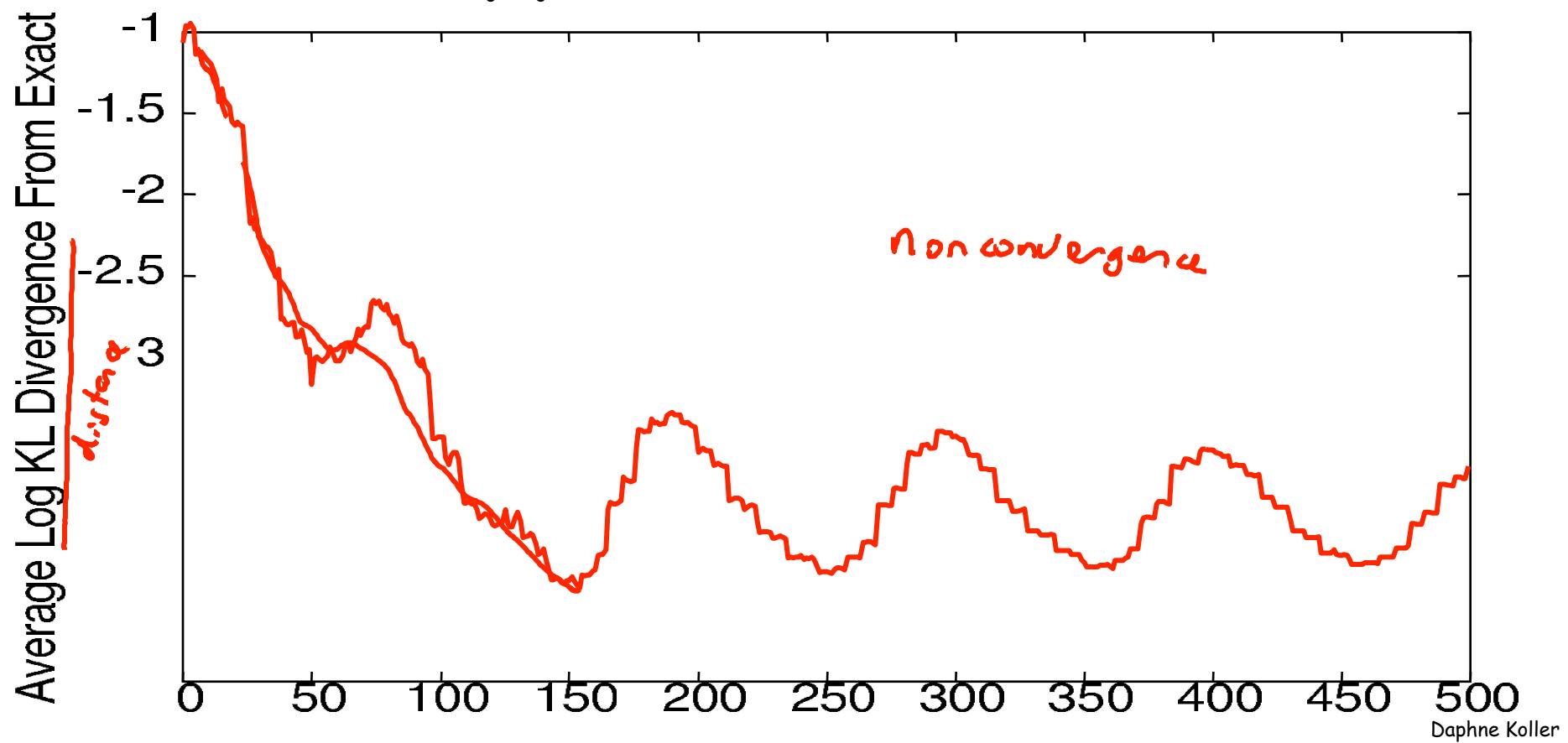
$$\delta_{i \rightarrow j}(S_{i,j}) = \sum_{C_i - S_{i,j}} \psi_i \times \prod_{k \in (\mathcal{N}_i - \{j\})} \delta_{k \rightarrow i}$$

- Compute  $\beta_i(C_i) = \psi_i \times \prod_{k \in \mathcal{N}_i} \delta_{k \rightarrow i}$  — all neighbors

# Belief Propagation Run



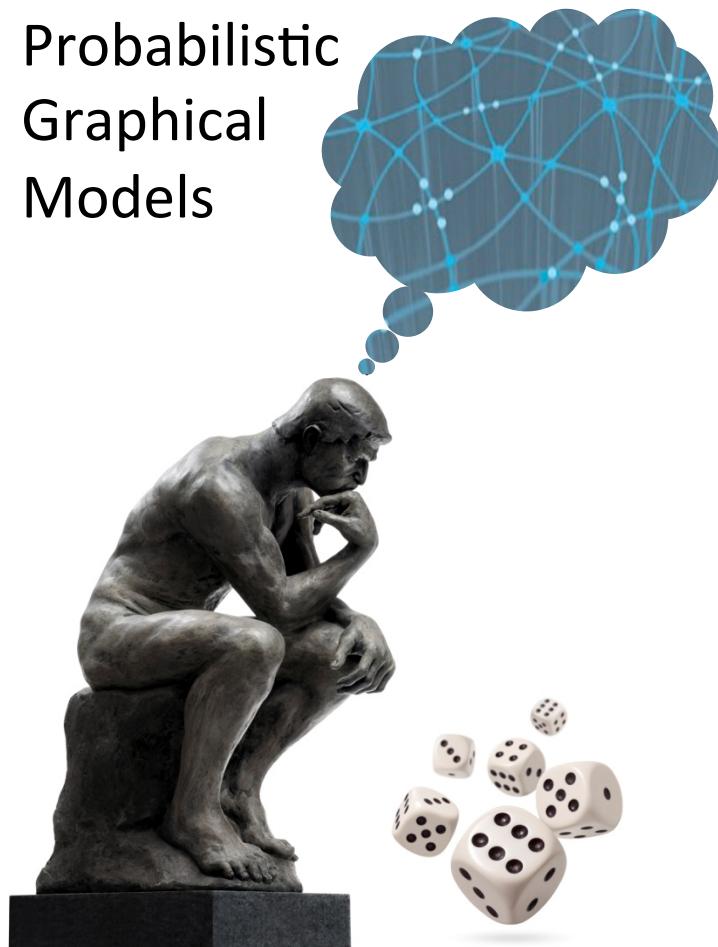
# Different BP Run



# Summary

- Graph of clusters connected by sepsets
- Adjacent clusters pass information to each other about variables in sepset
  - Message from i to j summarizes everything i knows, except information obtained from j
- Algorithm may not converge *not marginals +  $P_j$*
- The resulting beliefs are pseudo-marginals
- Nevertheless, very useful in practice

Probabilistic  
Graphical  
Models



Inference

---

Message Passing

---

# Cluster Graph Properties

# Cluster Graphs

- Undirected graph such that:
  - nodes are clusters  $C_i \subseteq \{X_1, \dots, X_n\}$
  - edge between  $C_i$  and  $C_j$  associated with sepset  $S_{i,j} \subseteq C_i \cap C_j$

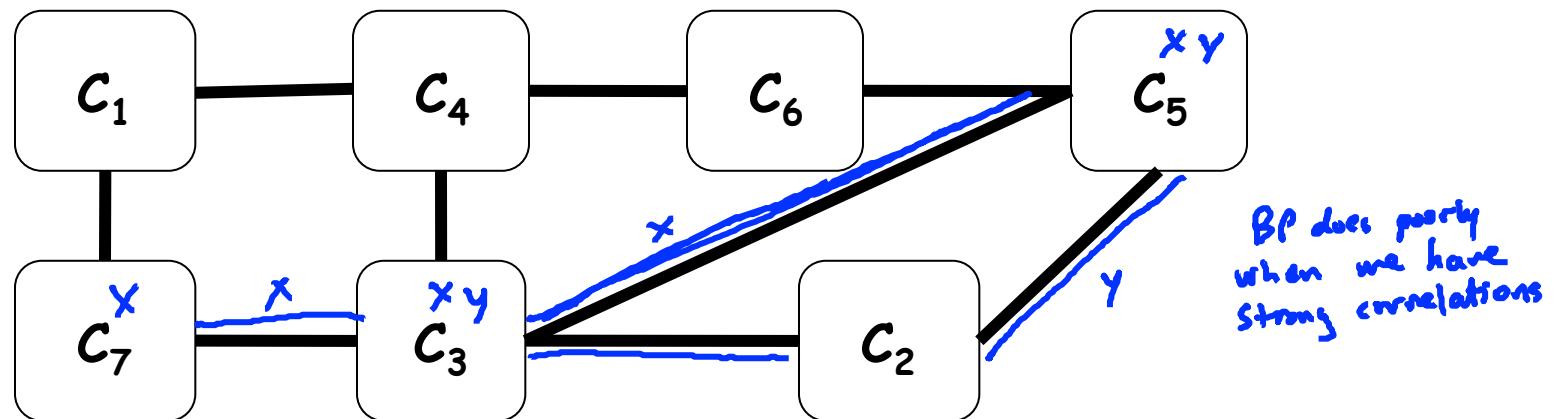
## Family Preservation

- Given set of factors  $\Phi$ , we assign each  $\phi_k$  to a cluster  $C_{\alpha(k)}$  s.t.  $\text{Scope}[\phi_k] \subseteq C_{\alpha(k)}$
- For each factor  $\phi_k \in \Phi$ , there exists a cluster  $C_i$  s.t.  $\text{Scope}[\phi_k] \subseteq C_i$   $\leftarrow$  accommodates  $\phi_k$

# Running Intersection Property

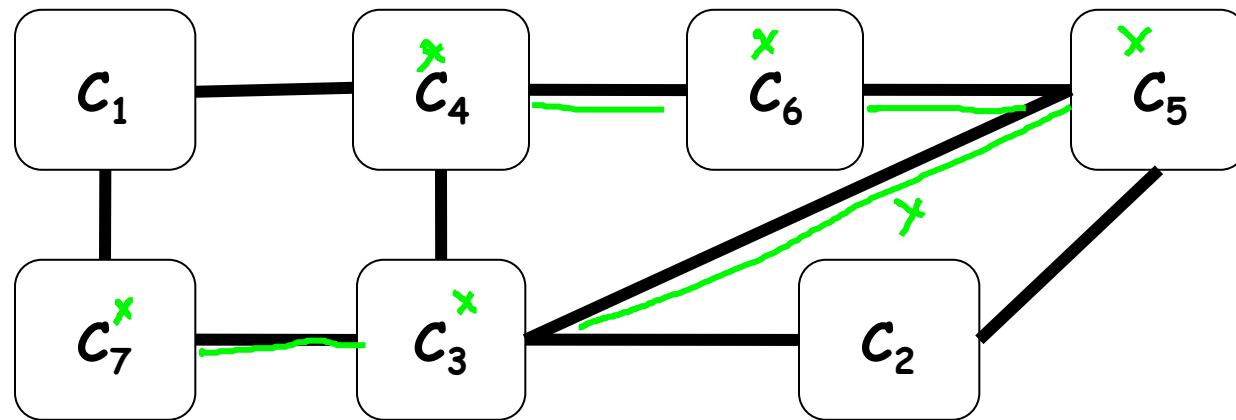
- For each pair of clusters  $C_i, C_j$  and variable  $X \in C_i \cap C_j$  there exists a unique path between  $C_i$  and  $C_j$  for which all clusters and sepsets contain  $X$

*x and Y that are very strongly correlated*



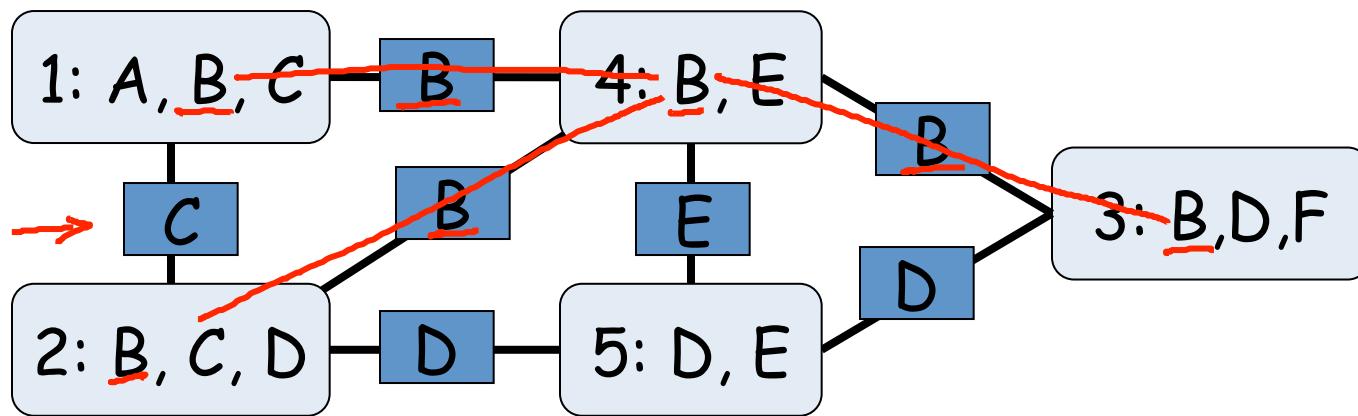
# Running Intersection Property

- Equivalently: For any  $X$ , the set of clusters and sepsets containing  $X$  forms a tree

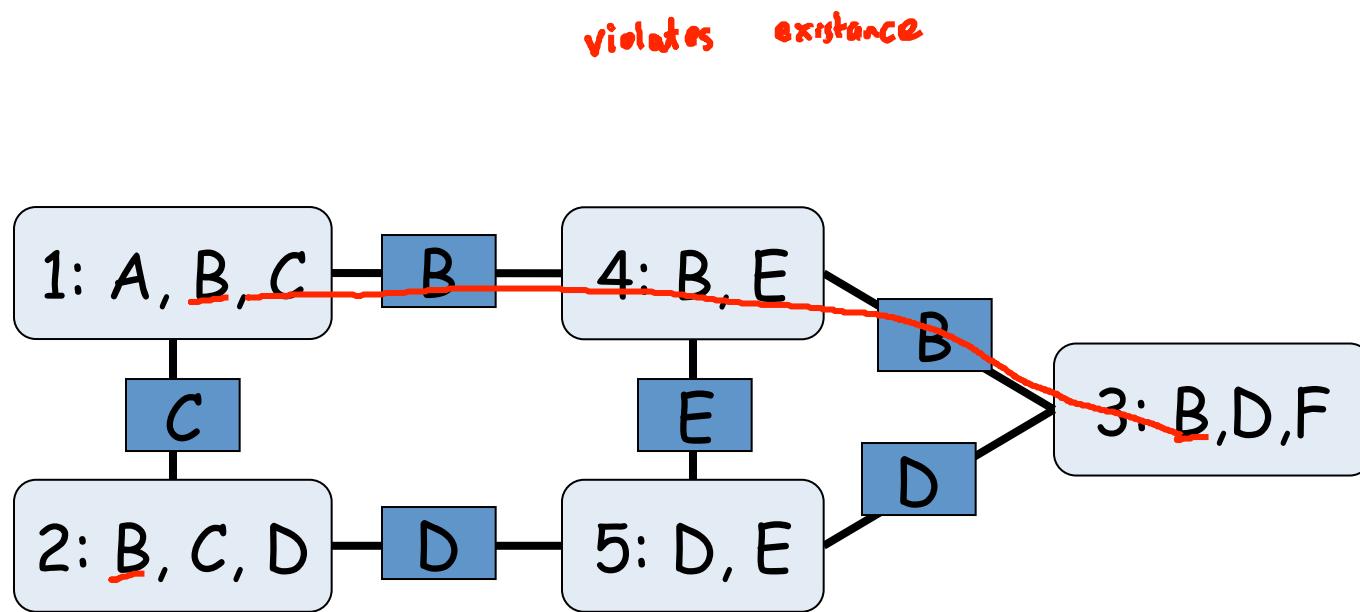


Daphne Koller

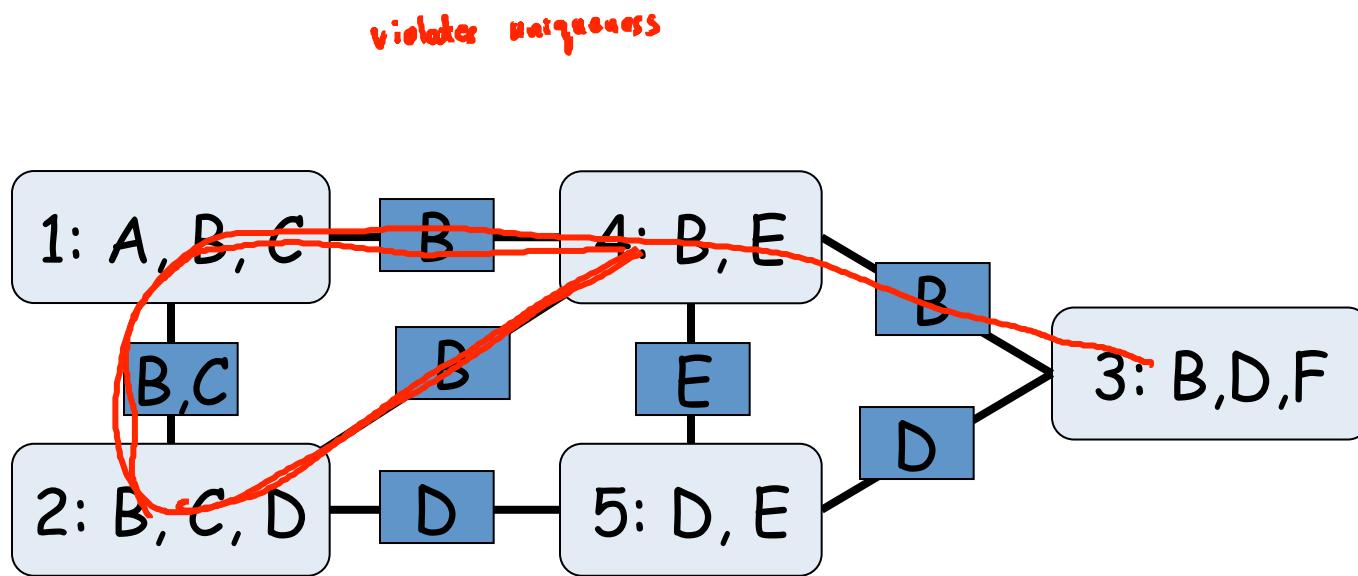
# Example Cluster Graph



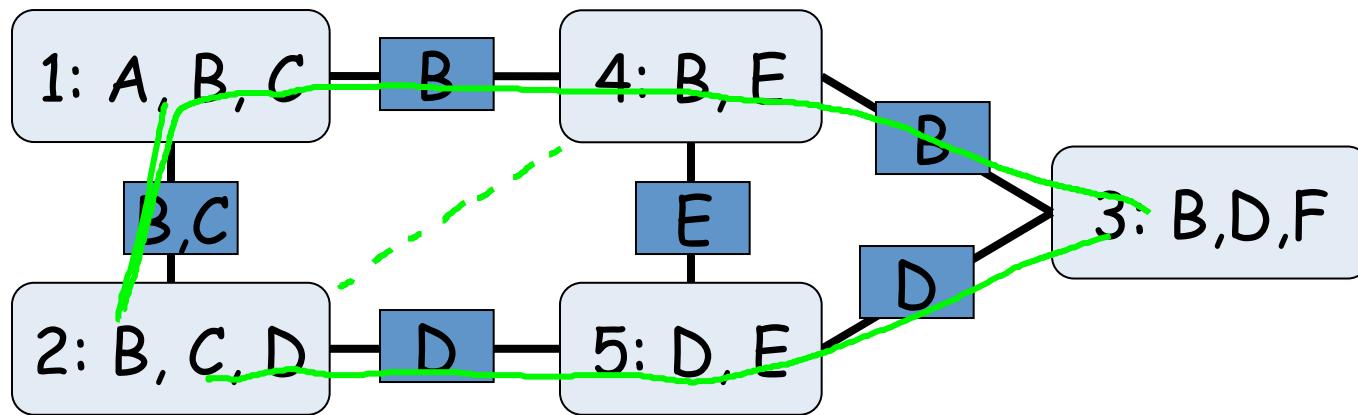
# Illegal Cluster Graph I



# Illegal Cluster Graph II



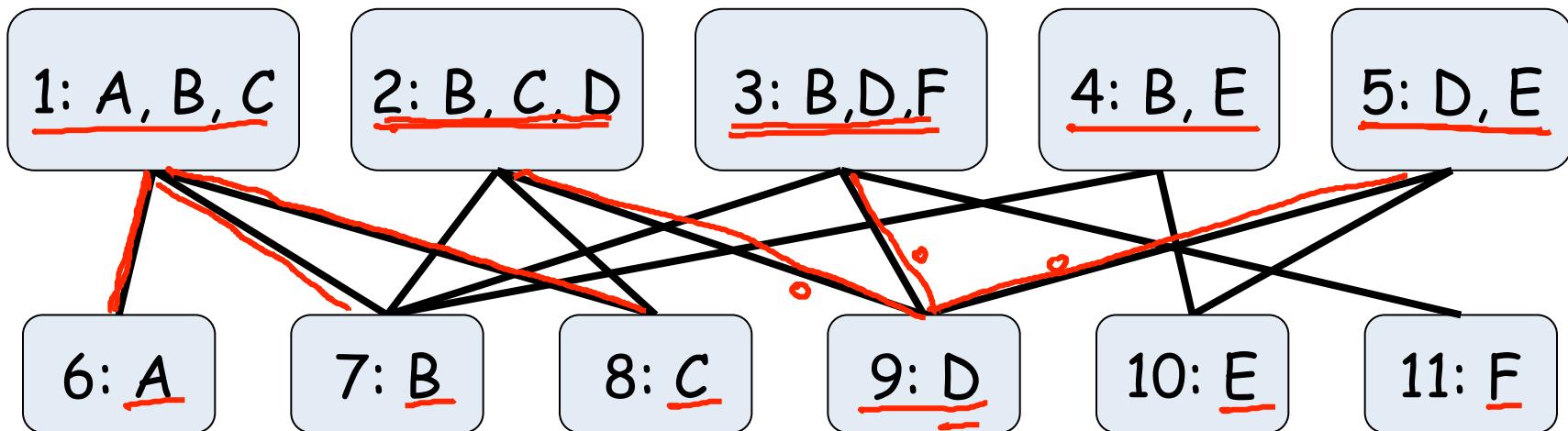
# Alternative Legal Cluster Graph



# Bethe Cluster Graph

*big clusters = factor in  $\Phi$*

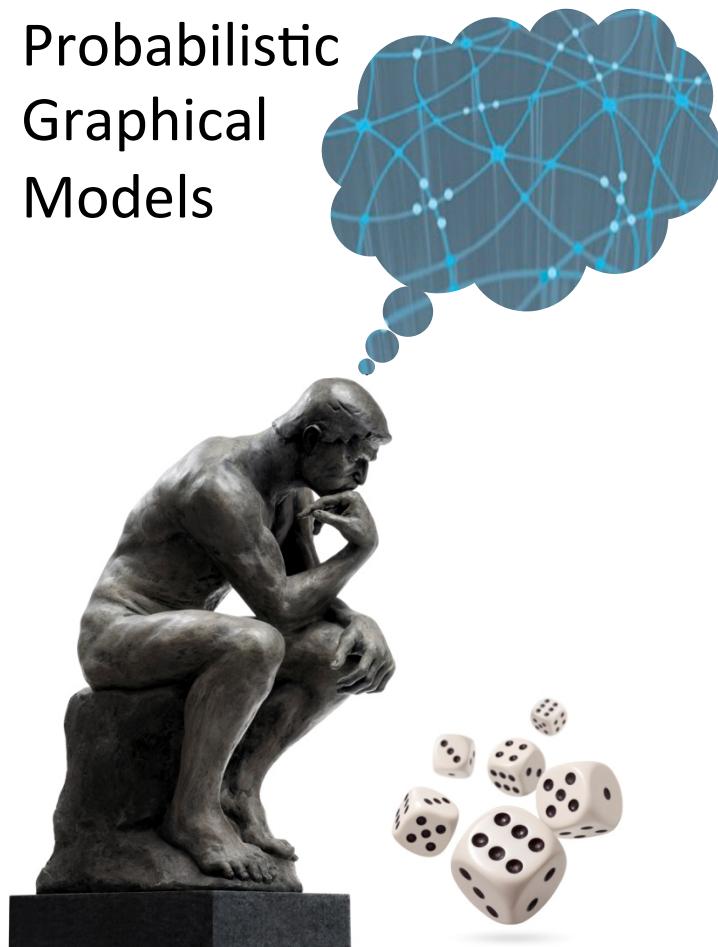
- For each  $\phi_k \in \Phi$ , a factor cluster  $C_k = \text{Scope}[\phi_k]$
- For each  $X_i$  a singleton cluster  $\{X_i\}$
- Edge  $C_k — X_i$  if  $X_i \in C_k$



# Summary

- Cluster graph must satisfy two properties
  - family preservation: allows  $\Phi$  to be encoded
  - running intersection: connects all information about any variable, but without feedback loops
- Bethe cluster graph is often first default
- Richer cluster graph structures can offer different tradeoffs wrt computational cost and preservation of dependencies

Probabilistic  
Graphical  
Models



Inference

---

Message Passing

---

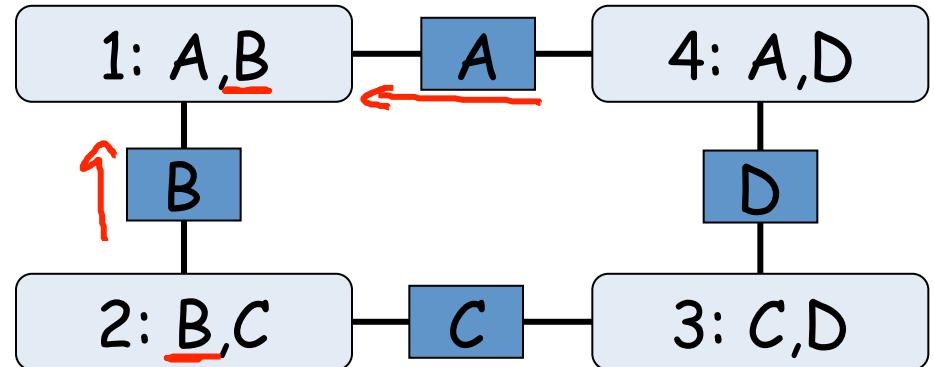
# Properties of BP Algorithm

# Calibration

$$\beta_1(A, B) = \underline{\psi_1(A, B)} \times \delta_{4 \rightarrow 1}(A) \times \delta_{2 \rightarrow 1}(B)$$

- Cluster beliefs:

$$\underline{\beta_i(C_i)} = \underline{\psi_i} \times \prod_{k \in \mathcal{N}_i} \underline{\delta_{k \rightarrow i}}$$



- A cluster graph is calibrated if every pair of adjacent clusters  $C_i, C_j$  agree on their sepset  $S_{i,j}$

$$\sum_{C_i - S_{i,j}} \underline{\beta_i(C_i)} = \sum_{C_j - S_{i,j}} \underline{\beta_j(C_j)}$$

sepset  $S_{i,j}$

# Convergence $\Rightarrow$ Calibration

- Convergence:

$$\delta_{i \rightarrow j}(S_{i,j}) = \delta'_{i \rightarrow j}(S_{i,j})$$

$$\beta_i(C_i) = \psi_i \times \prod_{k \in \mathcal{N}_i} \delta_{k \rightarrow i}$$

$$\delta'_{i \rightarrow j}(S_{i,j}) = \sum_{C_i - S_{i,j}} \left( \psi_i \times \prod_{k \in (\mathcal{N}_i - \{j\})} \delta_{k \rightarrow i} \right) = \sum_{C_i - S_{i,j}} \frac{\beta_i(C_i)}{\delta_{j \rightarrow i}(S_{i,j})} =$$

*all msgs*

$$\delta_{j \rightarrow i}(S_{i,j}) \delta_{i \rightarrow j}(S_{i,j}) = \sum_{C_i - S_{i,j}} \underline{\beta_i(C_i)}$$

*calibration*

$$\delta_{j \rightarrow i}(S_{i,j}) \delta_{i \rightarrow j}(S_{i,j}) = \sum_{C_j - S_{i,j}} \underline{\beta_j(C_j)}$$

*subset beliefs*

$$\mu_{i,j}(S_{i,j}) = \delta_{j \rightarrow i} \delta_{i \rightarrow j} = \sum_{C_j - S_{i,j}} \beta_j(C_j)$$

Daphne Koller

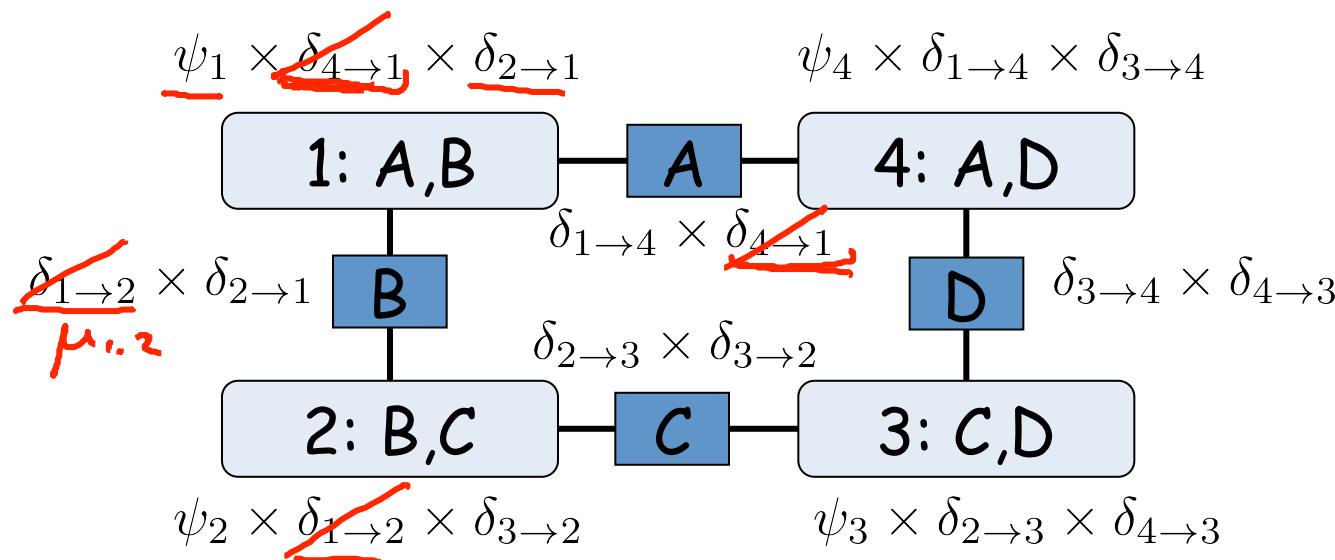
# Reparameterization

$$\beta_i(C_i) = \psi_i \times \prod_{k \in \mathcal{N}_i} \delta_{k \rightarrow i}$$

$$\mu_{i,j}(S_{i,j}) = \delta_{j \rightarrow i} \delta_{i \rightarrow j}$$

*separat beliefs*

$$\frac{\prod_i \beta_i}{\prod_{i,j} \mu_{i,j}}$$



# Reparameterization

no information loss

$$\beta_i(C_i) = \psi_i \times \prod_{k \in \mathcal{N}_i} \delta_{k \rightarrow i} \quad \mu_{i,j}(S_{i,j}) = \delta_{j \rightarrow i} \delta_{i \rightarrow j}$$

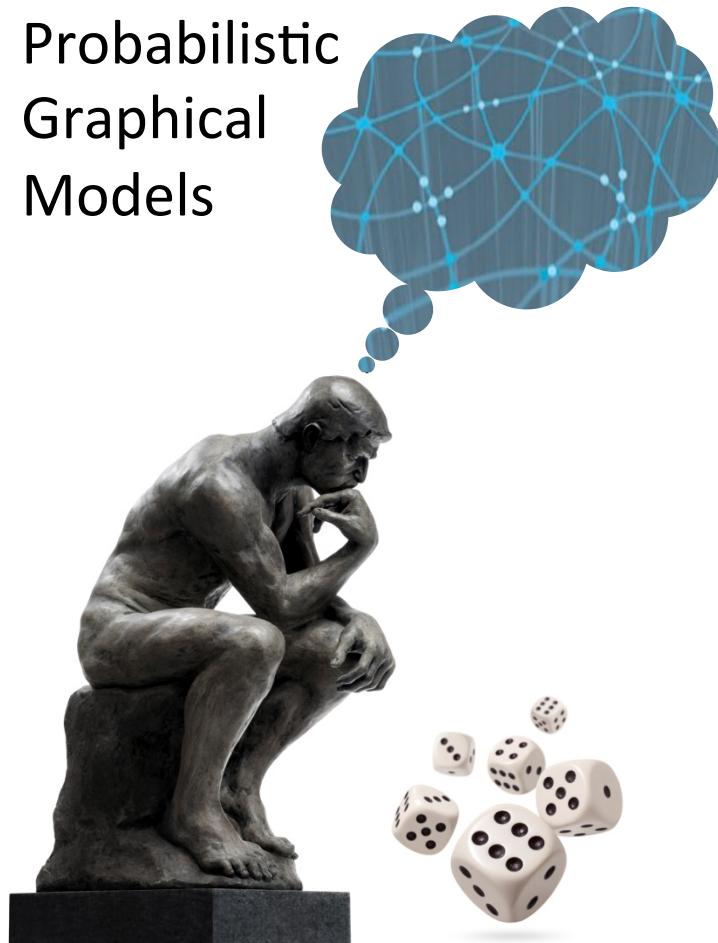
$$\begin{aligned} \frac{\prod_i \beta_i}{\prod_{i,j} \mu_{i,j}} &= \frac{\prod_i (\psi_i \prod_{j \in \mathcal{N}_i} \delta_{j \rightarrow i})}{\prod_{i,j} \delta_{i \rightarrow j}} \\ &= \prod_i \psi_i = \tilde{P}_{\Phi}(X_1, \dots, X_n) \end{aligned}$$

unnormalized measure

# Summary

- At convergence of BP, cluster graph beliefs are calibrated:
  - beliefs at adjacent clusters agree on sepsets
- Cluster graph beliefs are an alternative, calibrated parameterization of the original unnormalized density
  - No information is lost by message passing

Probabilistic  
Graphical  
Models



Inference

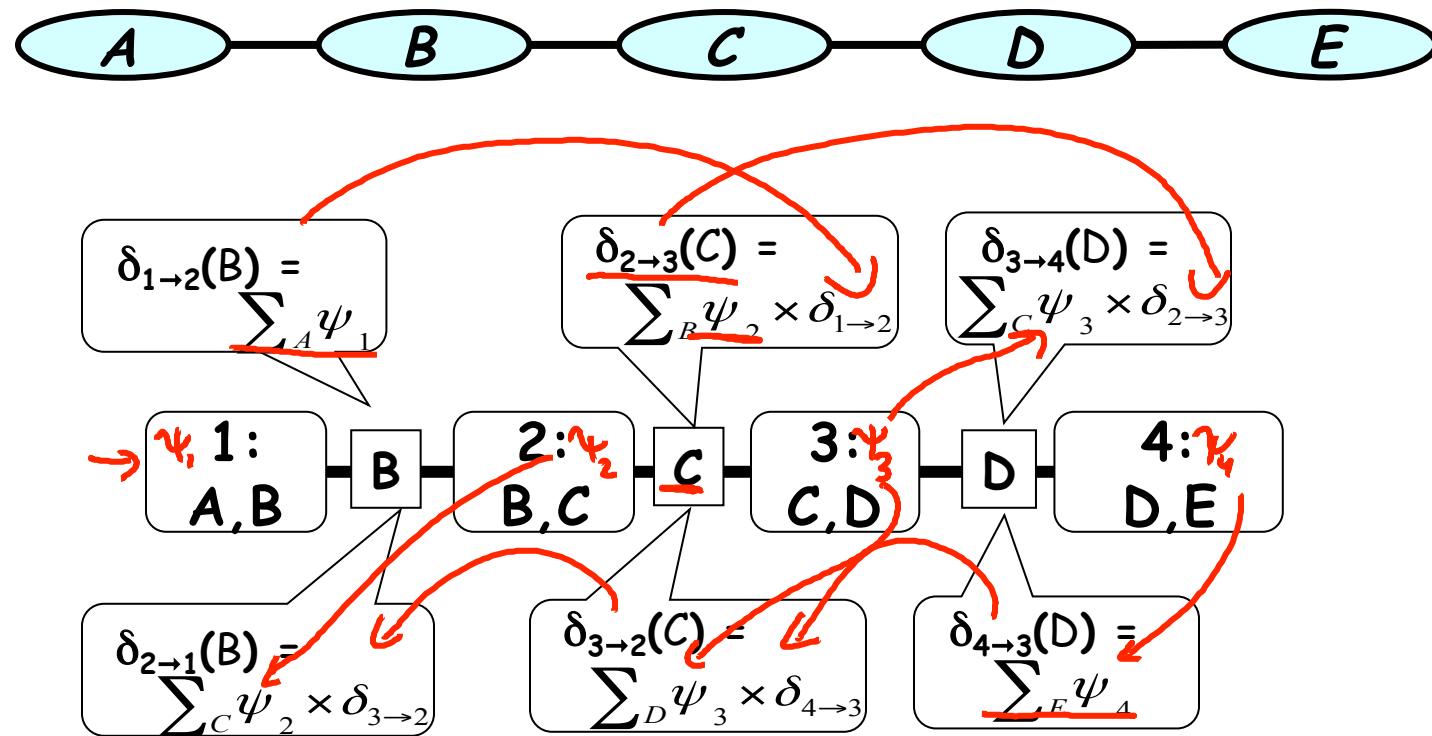
---

Message Passing

---

# Clique Tree Algorithm & Correctness

# Message Passing in Trees



# Correctness

$$\begin{aligned} \delta_{1 \rightarrow 2}(B) &= \sum_A \psi_1 \\ \delta_{2 \rightarrow 3}(C) &= \sum_B \psi_2 \times \delta_{1 \rightarrow 2}(B) \end{aligned}$$

1: A, B      2: B, C      3: C, D      4: D, E

*legal order of operations*

$$\begin{aligned} \beta_3(C, D) &= \psi_3 \times \delta_{2 \rightarrow 3} \times \delta_{4 \rightarrow 3} \\ &= \psi_3 \times \left( \sum_B (\psi_2 \times \delta_{1 \rightarrow 2}) \right) \times \sum_E \psi_4 \\ &= \psi_3 \times \left( \sum_B \psi_2 \times \sum_A \psi_1 \right) \times \sum_E \psi_4 \end{aligned}$$

*product of factors marginalized out unnecessary variables*

Daphne Koller

# Clique Tree

- Undirected tree such that:
  - nodes are clusters  $C_i \subseteq \{X_1, \dots, X_n\}$
  - edge between  $C_i$  and  $C_j$  associated with sepset  $S_{i,j} = C_i \cap C_j$

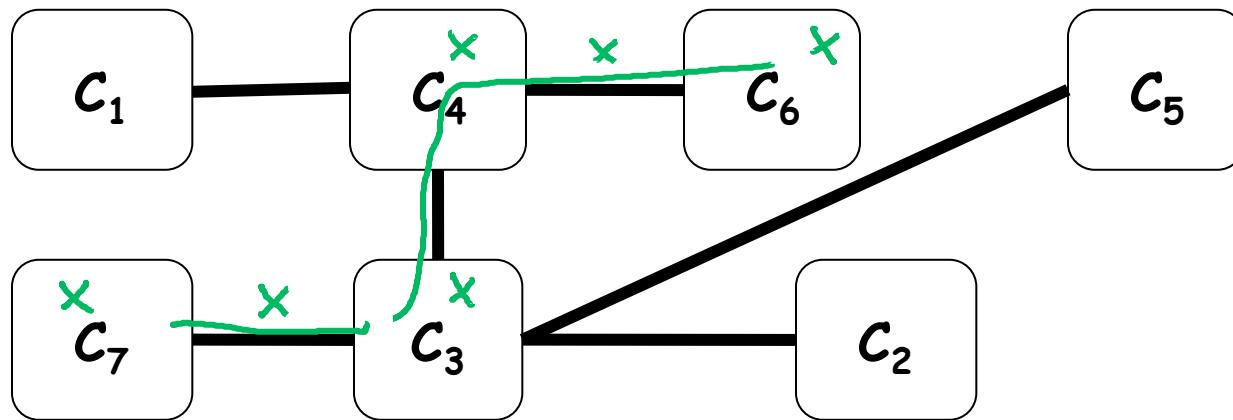
# Family Preservation

- Given set of factors  $\Phi$ , we assign each  $\phi_k \in \Phi$  to a cluster  $C_{\alpha(k)}$  s.t.  $\text{Scope}[\phi_k] \subseteq C_{\alpha(k)}$
- For each factor  $\phi_k \in \Phi$ , there exists a cluster  $C_i$  s.t.  $\text{Scope}[\phi_k] \subseteq C_i$

# Running Intersection Property

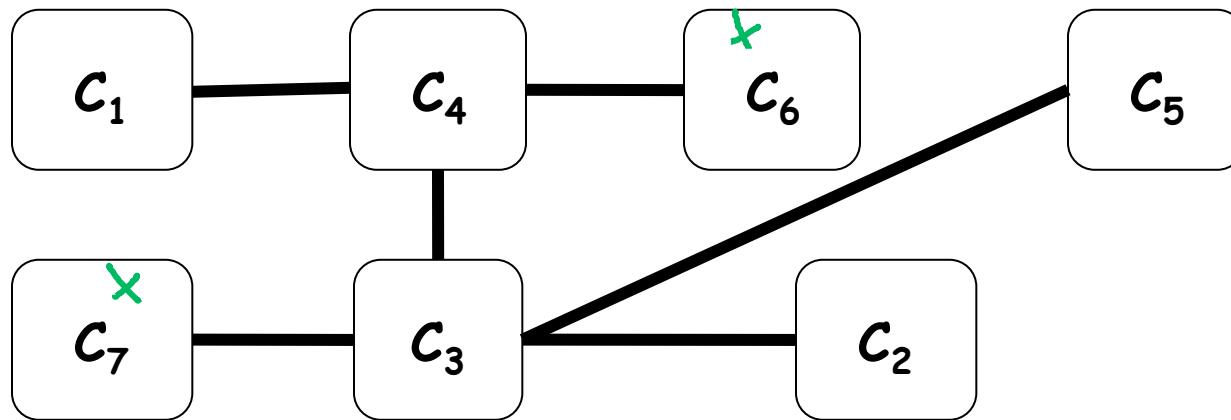
*Cluster graph variant*

- For each pair of clusters  $C_i, C_j$  and variable  $X \in C_i \cap C_j$  there exists a unique path between  $C_i$  and  $C_j$  for which all clusters and sepsets contain  $X$

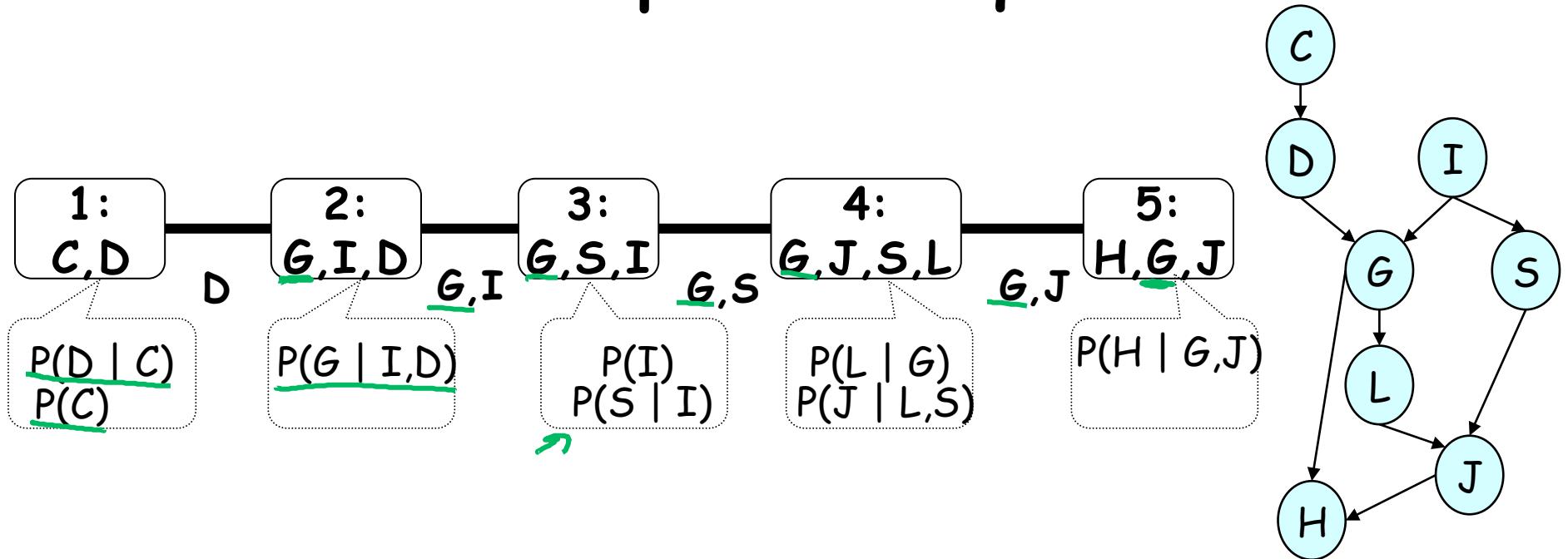


# Running Intersection Property

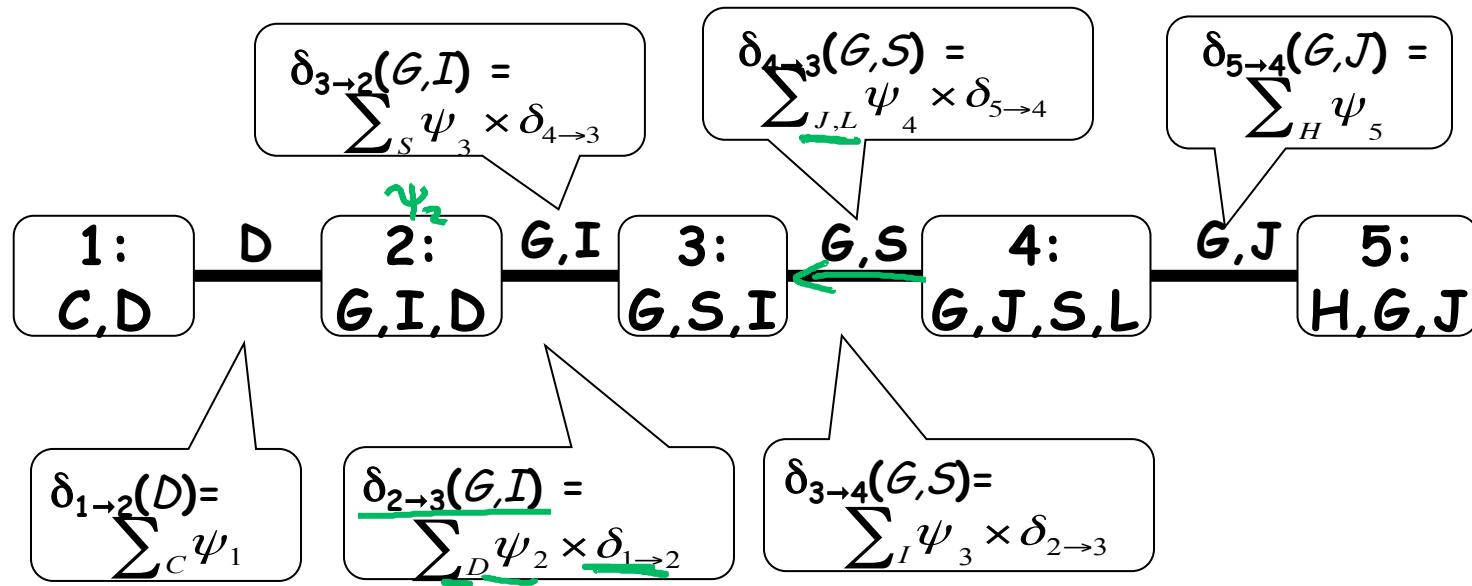
- For each pair of clusters  $C_i, C_j$  and variable  $X \in C_i \cap C_j$ , in the unique path between  $C_i$  and  $C_j$ , all clusters and sepsets contain  $X$



# More Complex Clique Tree

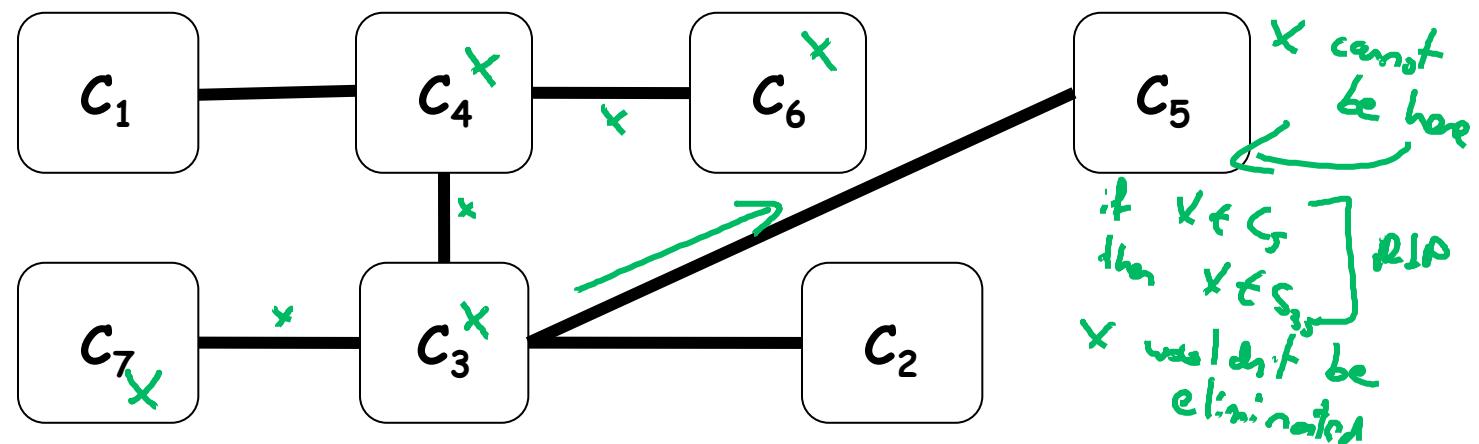


# Clique Tree Message Passing

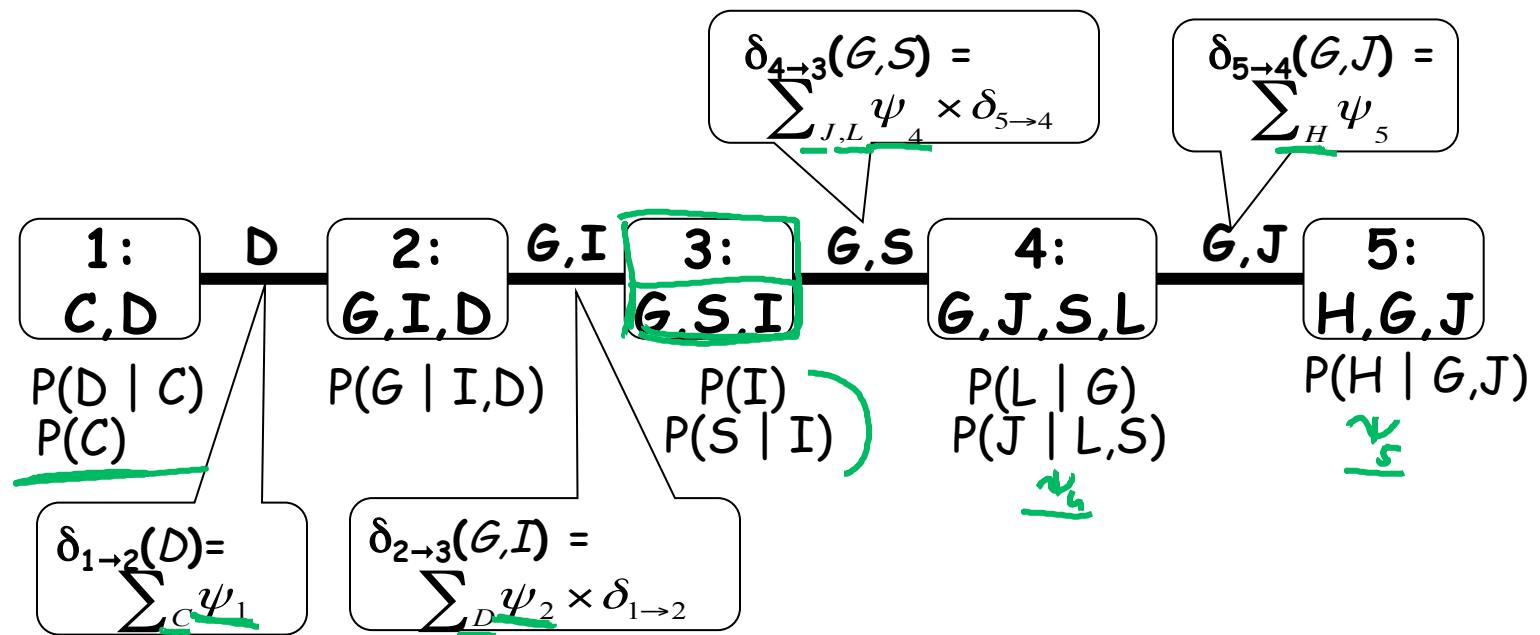


# RIP $\Rightarrow$ Clique Tree Correctness

- If  $X$  is eliminated when we pass the message  $C_i \rightarrow C_j$
- Then  $X$  does not appear in the  $C_j$  side of the tree



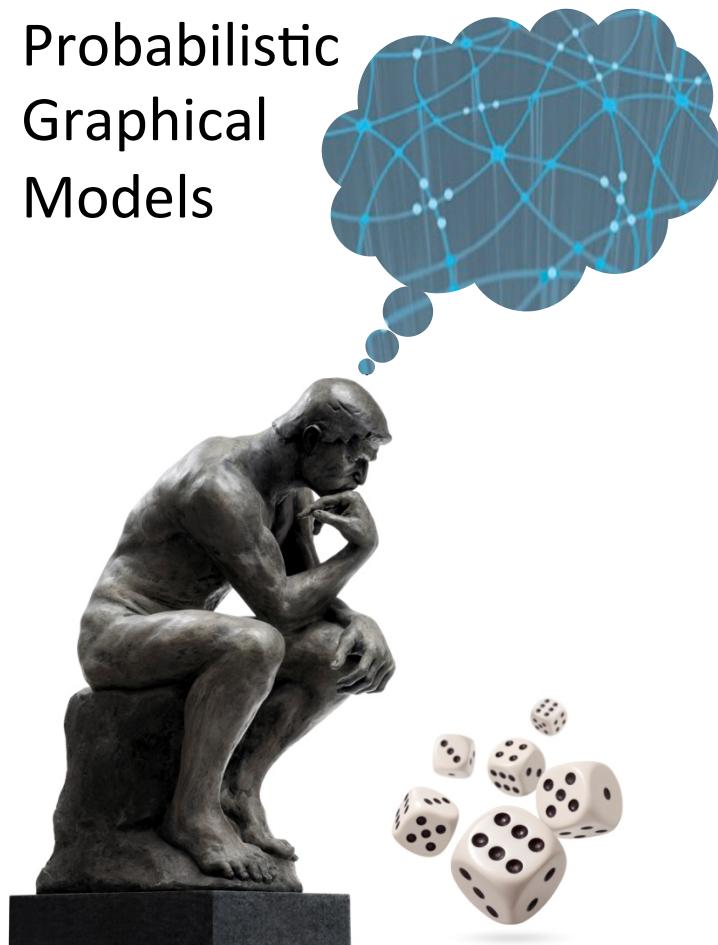
# Clique Tree Correctness



# Summary

- Belief propagation can be run over a tree-structured cluster graph
- In this case, computation is a variant of variable elimination
- Resulting beliefs are guaranteed to be correct marginals  $\pi_1 \dots \pi_5$

Probabilistic  
Graphical  
Models



Inference

---

Message Passing

---

# Clique Tree Algorithm: Computation

# Message Passing in Trees



*once computed  
never changes*

$$\delta_{1 \rightarrow 2}(B) = \sum_A \psi_1$$

1:  
A, B

*wait for  $\delta_{1 \rightarrow 2}$*

$$\delta_{2 \rightarrow 3}(C) = \sum_B \psi_2 \times \delta_{1 \rightarrow 2}$$

2:  
B, C

*wait for  $\delta_{2 \rightarrow 3}$*

$$\delta_{3 \rightarrow 4}(D) = \sum_C \psi_3 \times \delta_{2 \rightarrow 3}$$

3:  
C, D

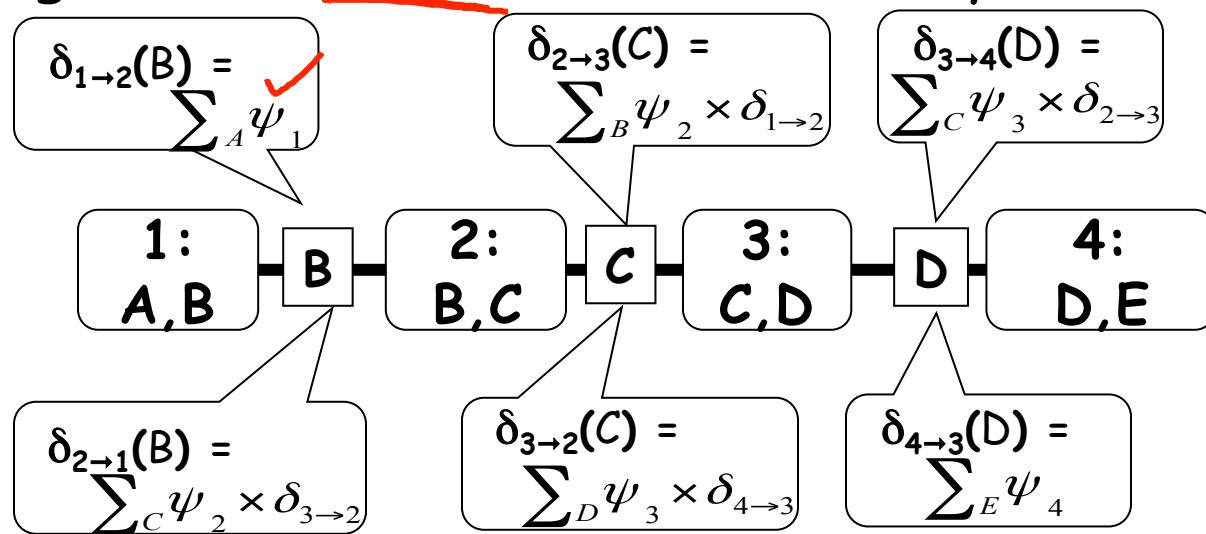
*Beliefs that  
are  $\pi_D(C_i)$*

$$\delta_{4 \rightarrow 3}(D) = \sum_E \psi_4$$

*converges  
instantly*

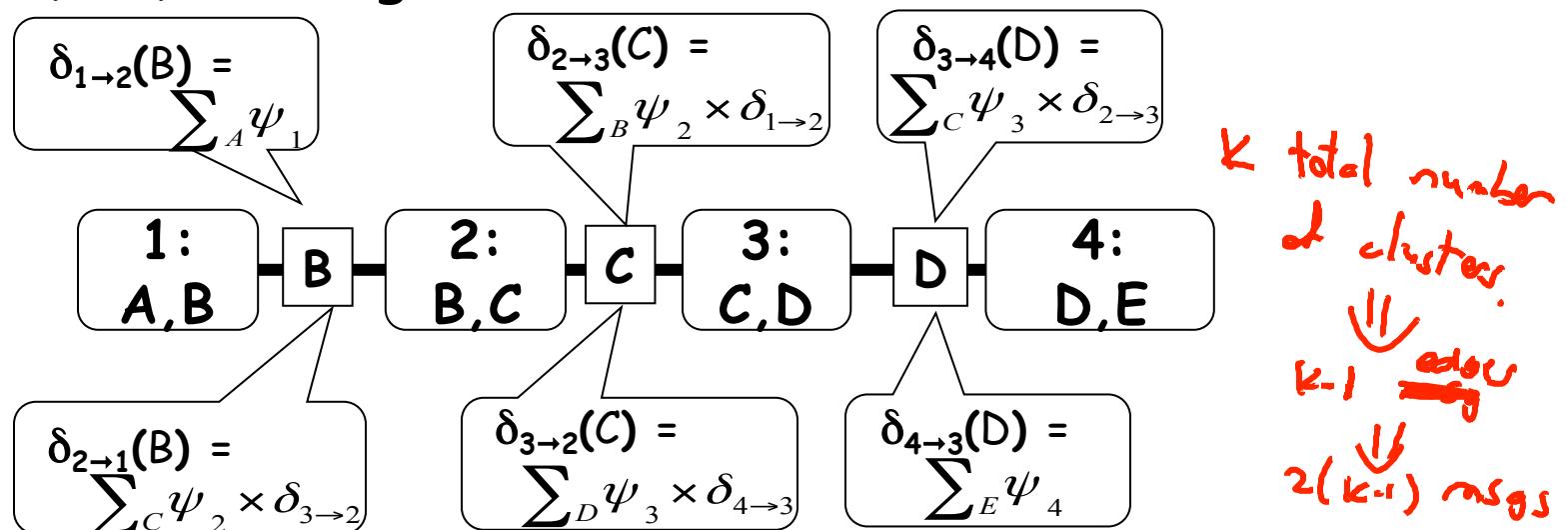
# Convergence of Message Passing

- Once  $C_i$  receives a final message from all neighbors except  $C_j$ , then  $\delta_{i \rightarrow j}$  is also final (will never change)
- Messages from leaves are immediately final

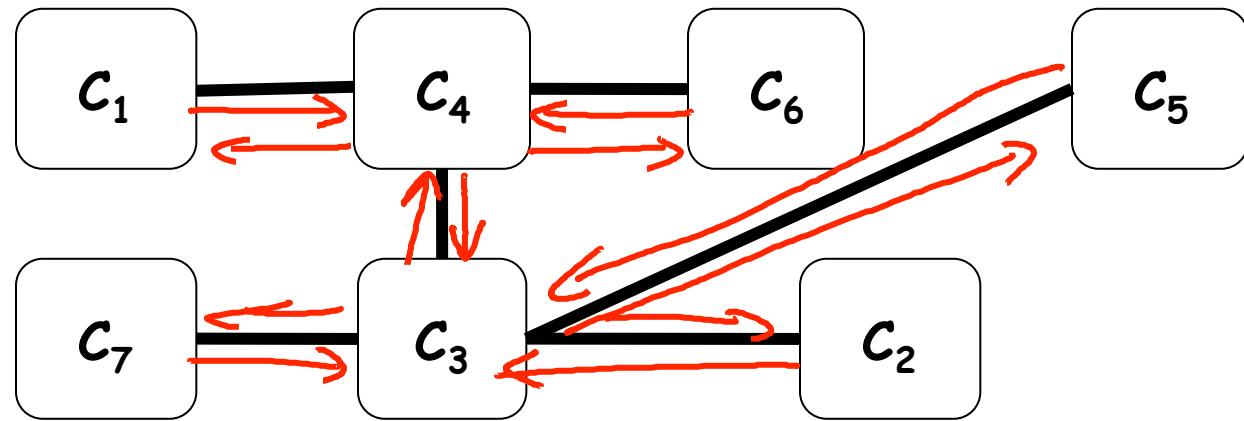


# Convergence of Message Passing

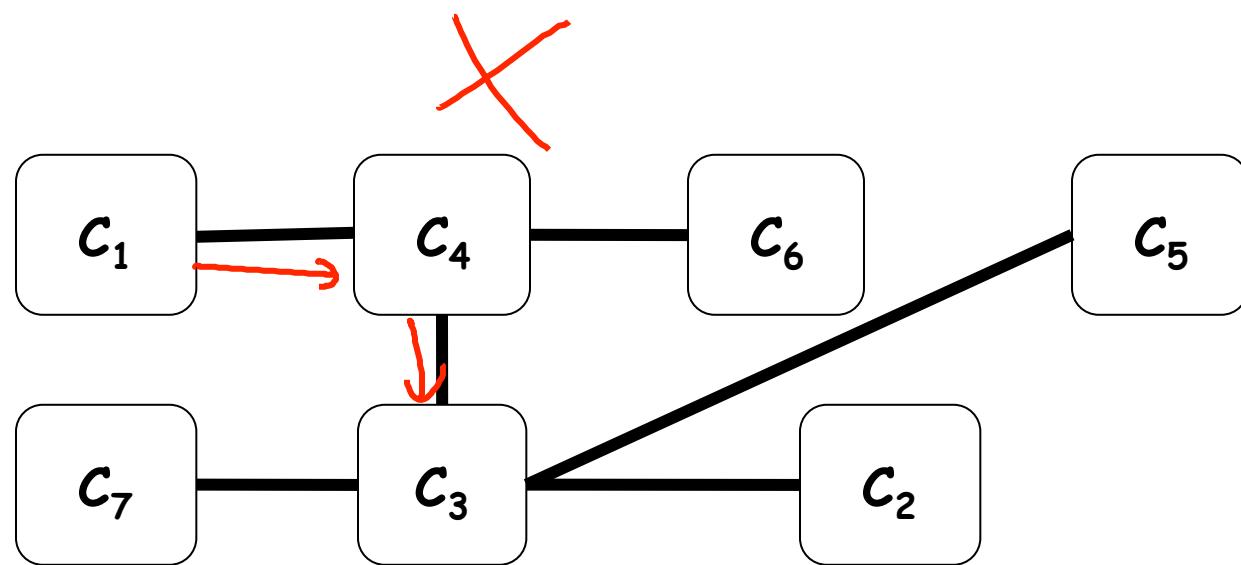
- Can pass messages from leaves inward
- If messages are passed in the right order, only need to pass  $2(K-1)$  messages



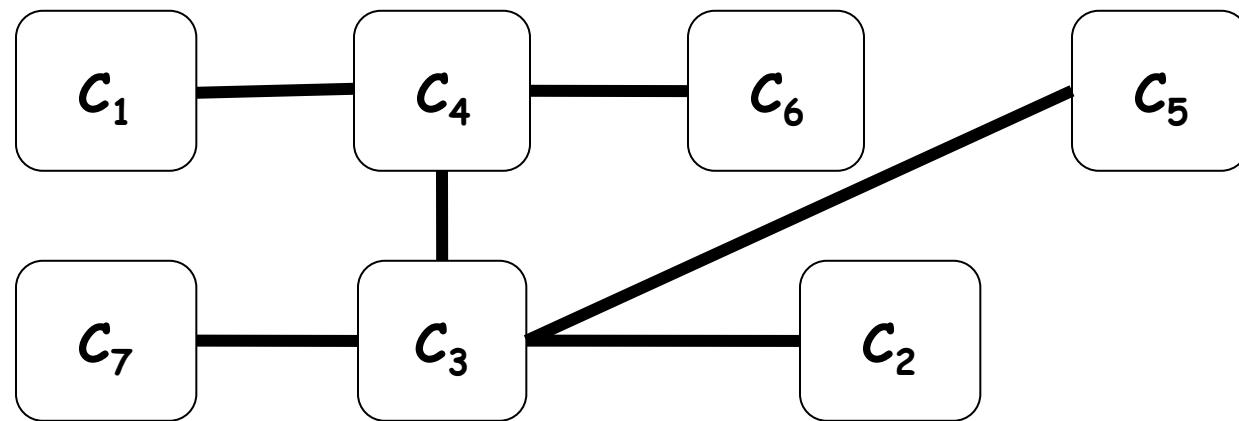
# Message Passing Order I



# Message Passing Order II



# Message Passing Order III

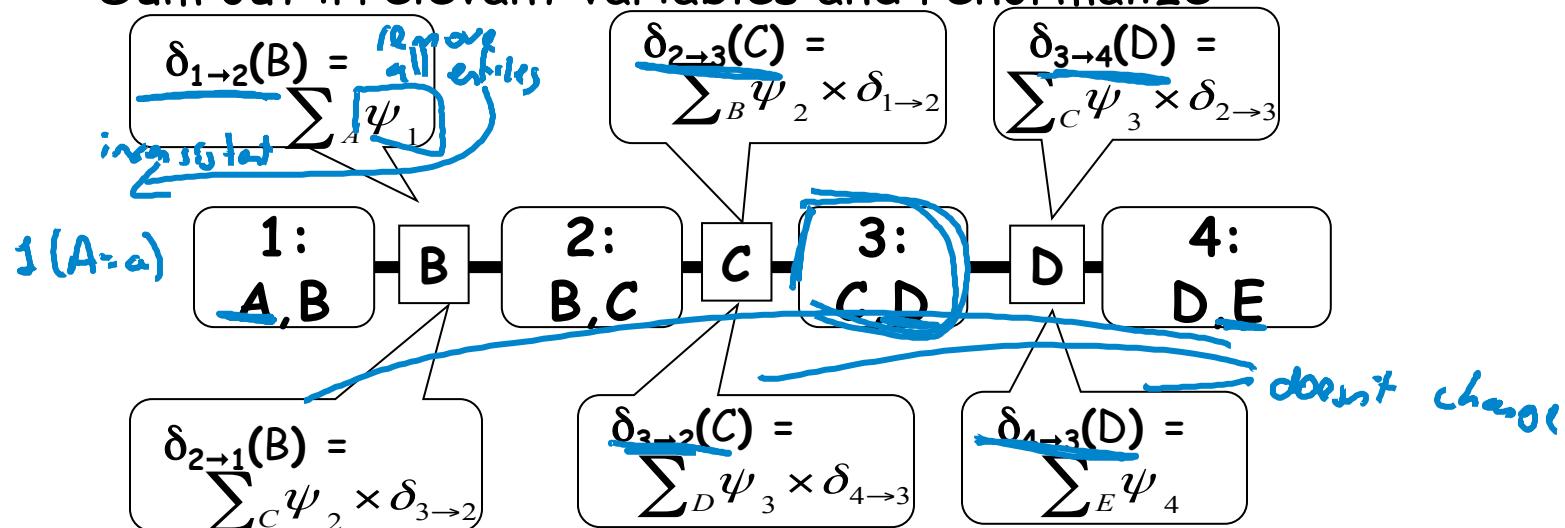


# Answering Queries

- Posterior distribution queries on variables that appear together in clique
  - Sum out irrelevant variables from any clique containing those variables  $p_{\phi}$  renormalize
- Introducing new evidence  $Z=z$  and querying  $X$ 
  - If  $X$  appears in clique with  $Z$  incremental inference
    - Multiply clique that contains  $X$  and  $Z$  with indicator function  $1(Z=z)$  reduce clique  $p_{\phi}(z, X)$
    - Sum out irrelevant variables and renormalize

# And More Queries

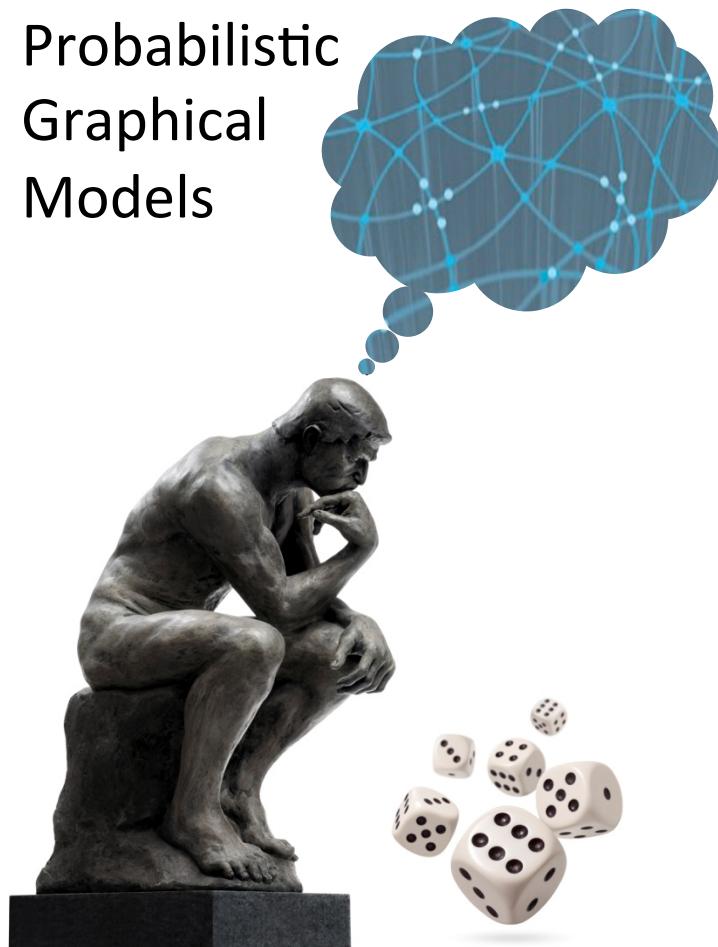
- Introducing new evidence  $Z=z$  and querying  $X$  if  $X$  does not share a clique with  $Z$ 
  - Multiply  $\mathbf{1}(Z=z)$  into some clique containing  $Z$  *reduction of factor*
  - Propagate messages along path to clique containing  $X$
  - Sum out irrelevant variables and renormalize



# Summary

- In clique tree with  $K$  cliques, if messages are passed starting at leaves,  $2(K-1)$  messages suffice to compute all beliefs
- Can compute marginals over all variables at only twice the cost of variable elimination
- By storing messages, inference can be reused in incremental queries

Probabilistic  
Graphical  
Models



Inference

---

Message Passing

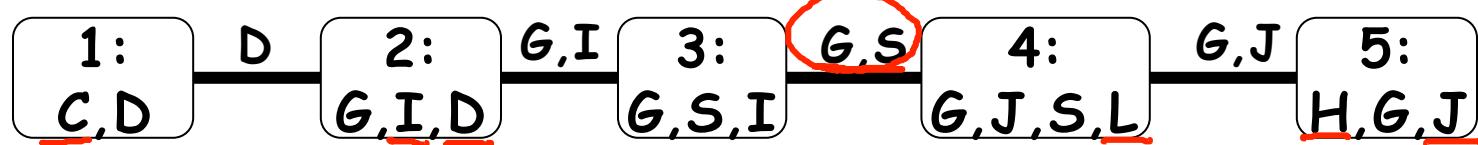
---

Clique Tree &  
Independence

# RIP and Independence

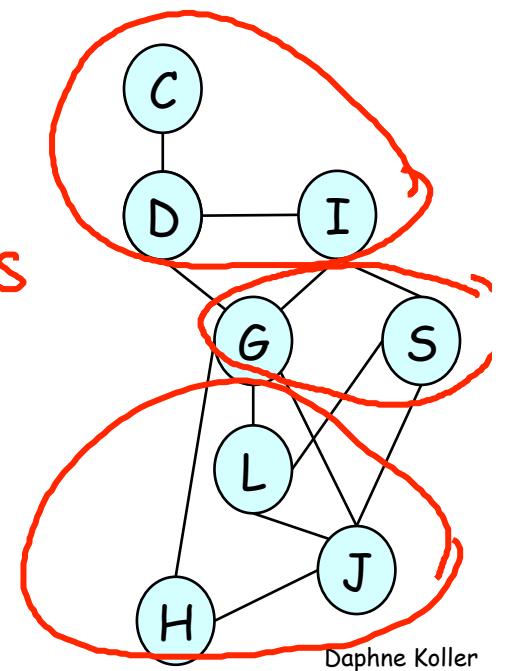
- For an edge  $(i,j)$  in  $T$ , let:
  - $W_{<(i,j)}$  = all variables that appear only on  $C_i$  side of  $T$
  - $W_{<(j,i)}$  = all variables that appear only on  $C_j$  side
  - Variables on both sides are in the sepset  $S_{i,j}$
- Theorem:  $T$  satisfies RIP if and only if, for every  $\underline{(i,j)}$   $P_\Phi \models (W_{<(i,j)} \perp W_{<(j,i)} \mid S_{i,j})$

# RIP and Independence



$$P_\Phi \models (\{\underline{C}, \underline{I}, \underline{D}\} \perp \underline{\{J, L, H\}} \mid \underline{\{G, S\}})$$

*C, I, D separated from H, L, J given G, S*

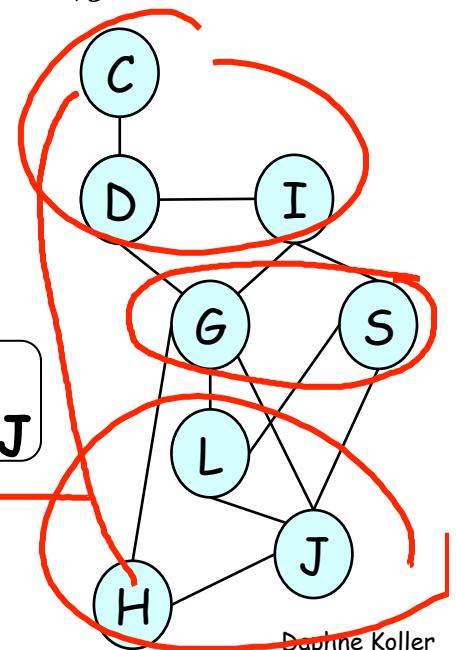
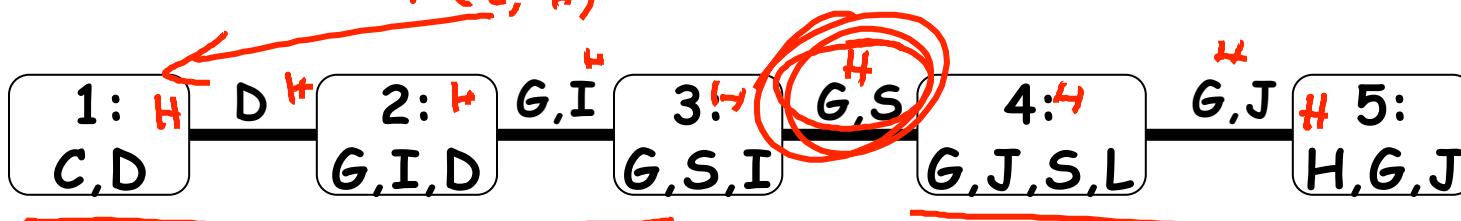


# RIP and Independence

- Theorem: T satisfies RIP if and only if, for every edge  $(i,j)$   $P_\Phi \models (W_{<(i,j)} \perp W_{<(j,i)} \mid S_{i,j})$

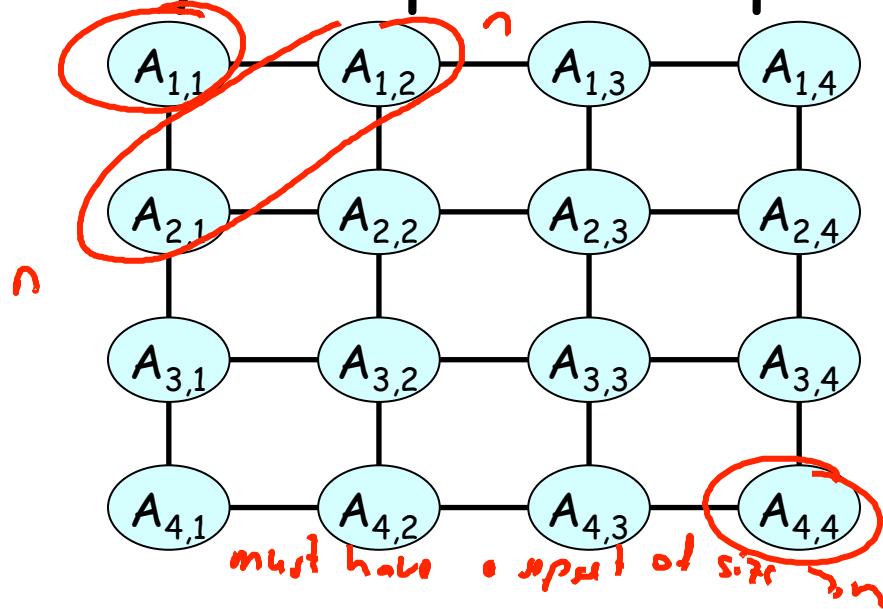
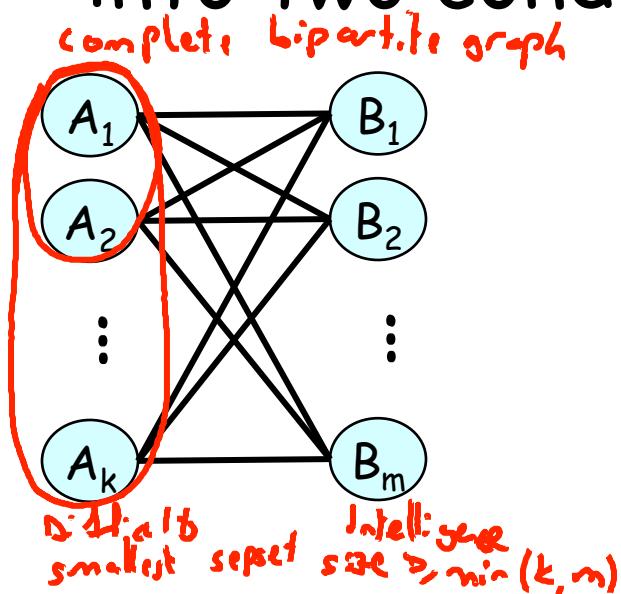
Assume otherwise  $\Rightarrow \exists$  path in induced Markov network between  $W_{<(i,j)}$   $W_{<(j,i)}$  that doesn't go through  $S_{i,j}$

Factor  $\phi(c, h)$



# Implications

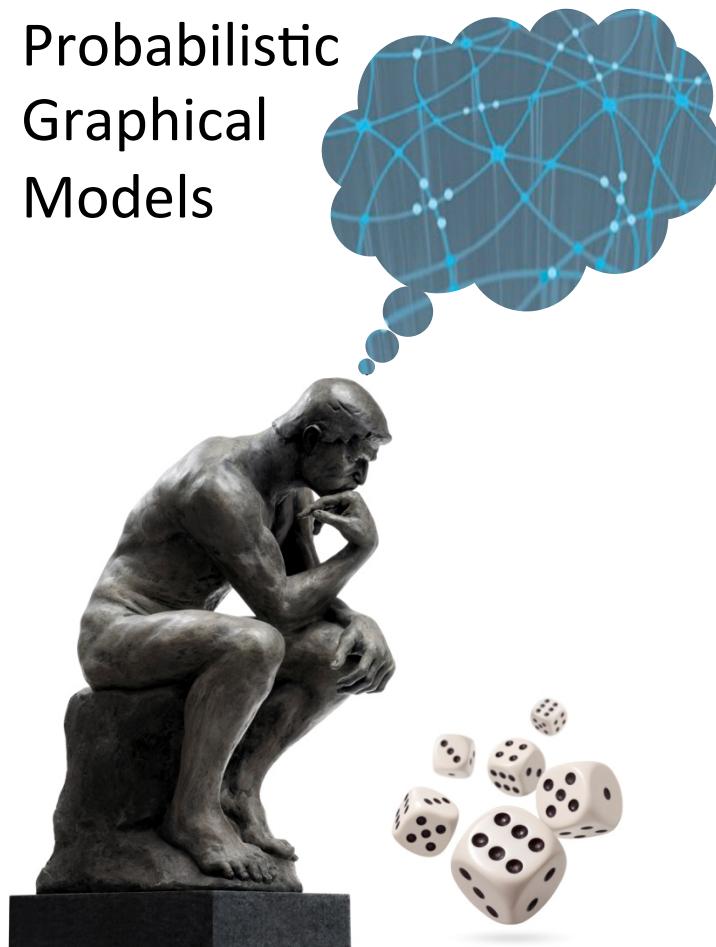
- Each sepset needs to separate graph into two conditionally independent parts



# Summary

- Correctness of clique tree inference relies on running intersection property
- Running intersection property implies separation in original distribution
- Implies minimal complexity incurred by any clique tree:
  - Related to minimal induced width of graph

Probabilistic  
Graphical  
Models



Inference

---

Message Passing

---

Clique Tree  
and VE

# Variable Elimination & Clique Trees

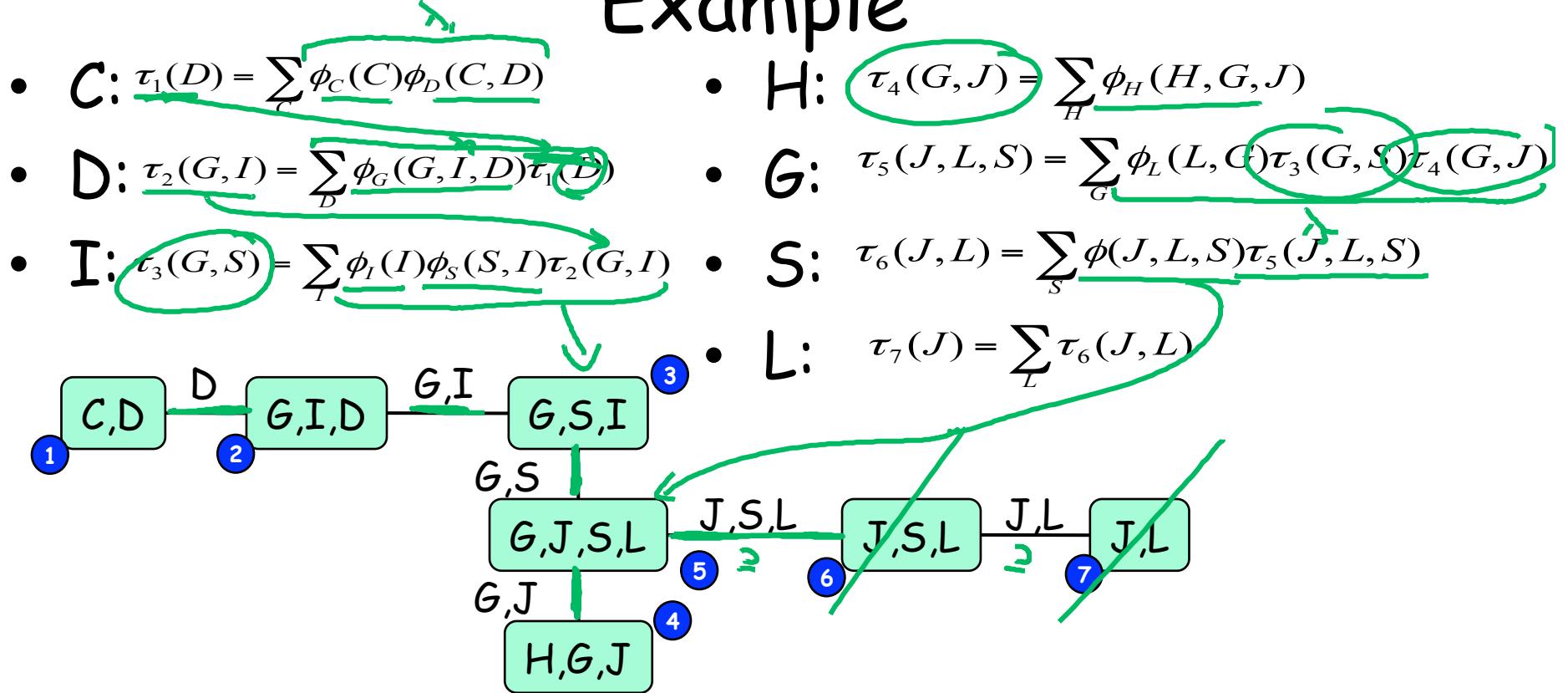
- Variable elimination
  - Each step creates a factor  $\lambda_i$  through factor product
  - A variable is eliminated in  $\lambda_i$  to generate new factor  $\tau_i$
  - $\tau_i$  is used in computing other factors  $\lambda_j$
- Clique tree view
  - Intermediate factors  $\lambda_i$  are cliques
  - $\tau_i$  are "messages" generated by clique  $\lambda_i$  and transmitted to another clique  $\lambda_j$

# Clique Tree from VE

- VE defines a graph
  - Cluster  $C_i$  for each factor  $\lambda_i$  used in the computation
  - Draw edge  $C_i - C_j$  if the factor generated from  $\lambda_i$  is used in the computation of  $\lambda_j$



# Example



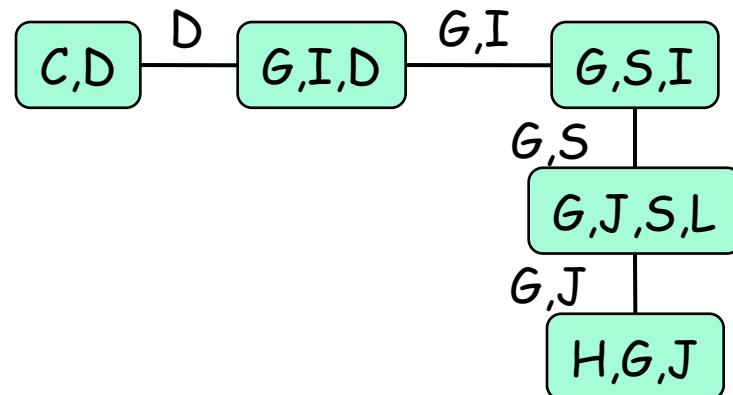
Remove redundant cliques:  
those whose scope is a subset of adjacent clique's scope

# Properties of Tree

- VE process induces a tree  $T_i$ 
  - In VE, each intermediate factor is used only once
  - Hence, each cluster “passes” a factor (message) to exactly one other cluster *(every cluster has at most one parent)*
- Tree is family preserving:  $\phi \in \Phi$ 
  - Each of the original factors must be used in some elimination step
  - And therefore contained in scope of associated  $\psi_i$   
*Scope that contains  $\text{Scop}(\phi)$*

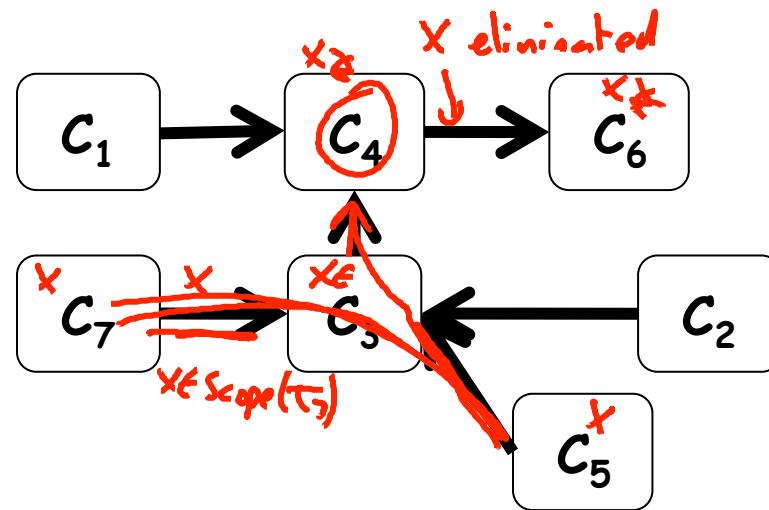
# Properties of Tree

- Tree obeys running intersection property
  - If  $\underline{X} \in C_i$  and  $\underline{X} \in C_j$  then  $X$  is in each cluster in the (unique) path between  $C_i$  and  $C_j$



# Running Intersection Property

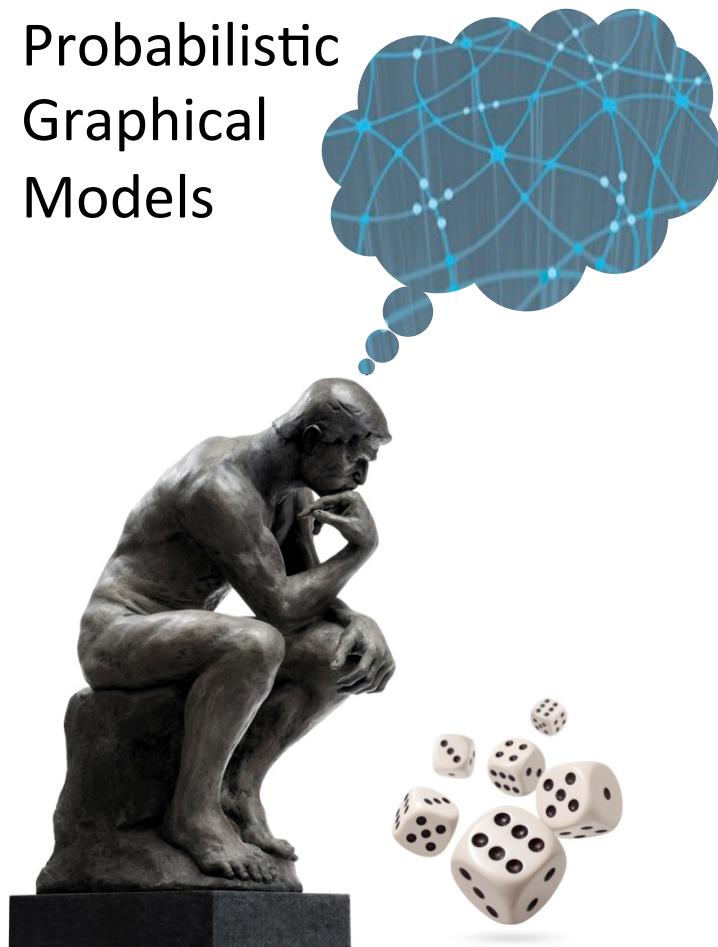
- **Theorem:** If  $T$  is a tree of clusters induced by VE, then  $T$  obeys RIP



# Summary

- A run of variable elimination implicitly defines a correct clique tree
  - We can “simulate” a run of VE to define cliques and connections between them
- Cost of variable elimination is ~ the same as passing messages in one direction in tree
- Clique trees use dynamic programming (storing messages) to compute marginals over all variables at only twice the cost of VE

Probabilistic  
Graphical  
Models



Inference

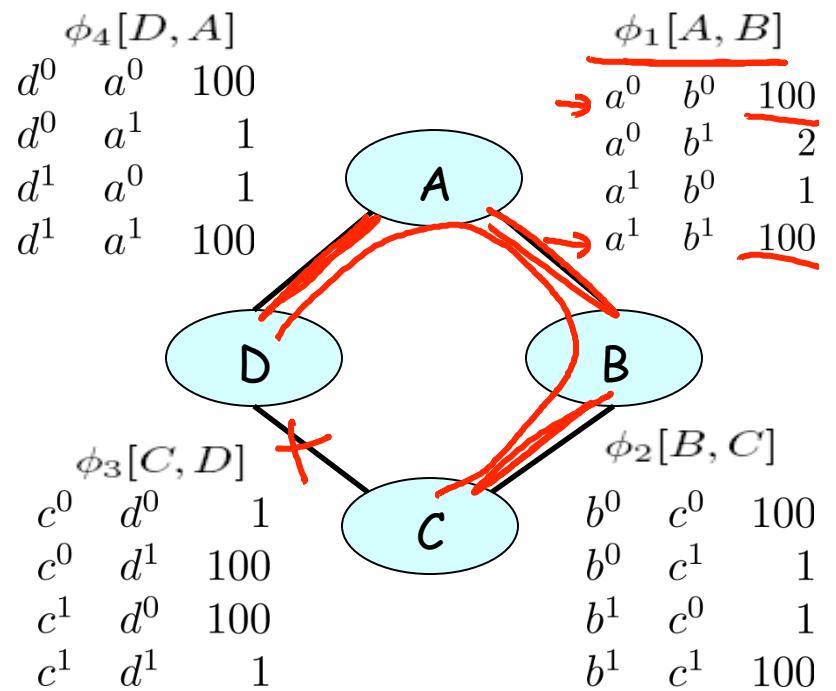
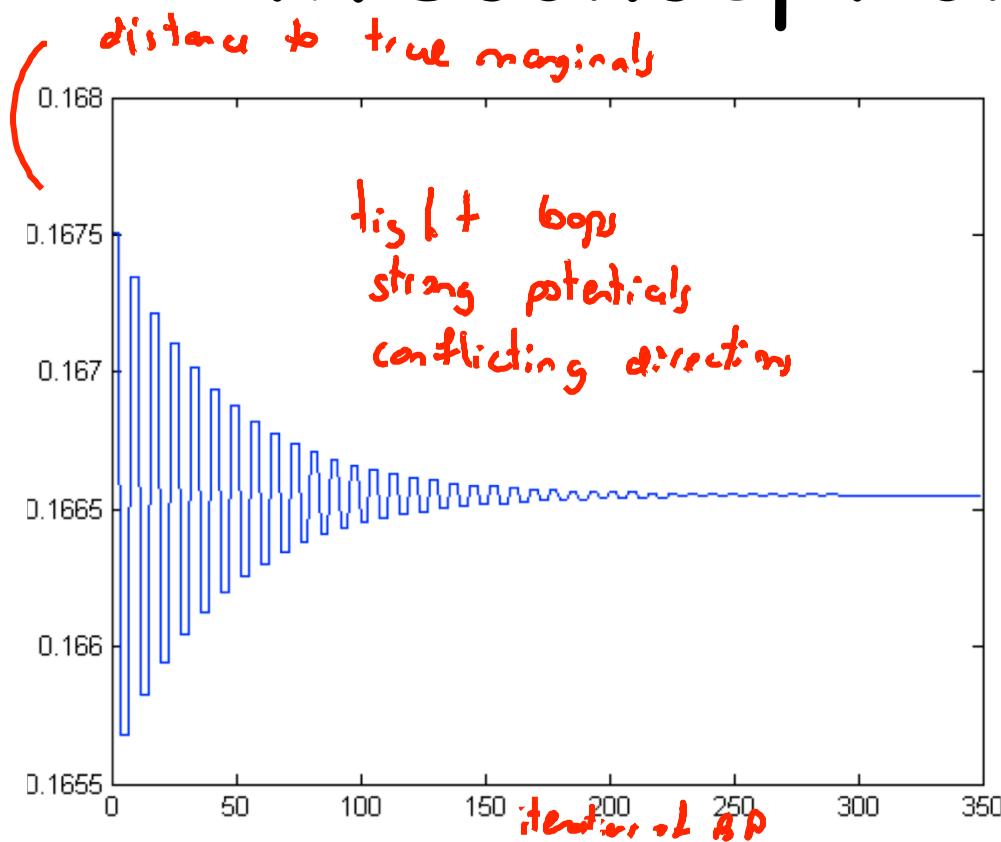
---

Message Passing

---

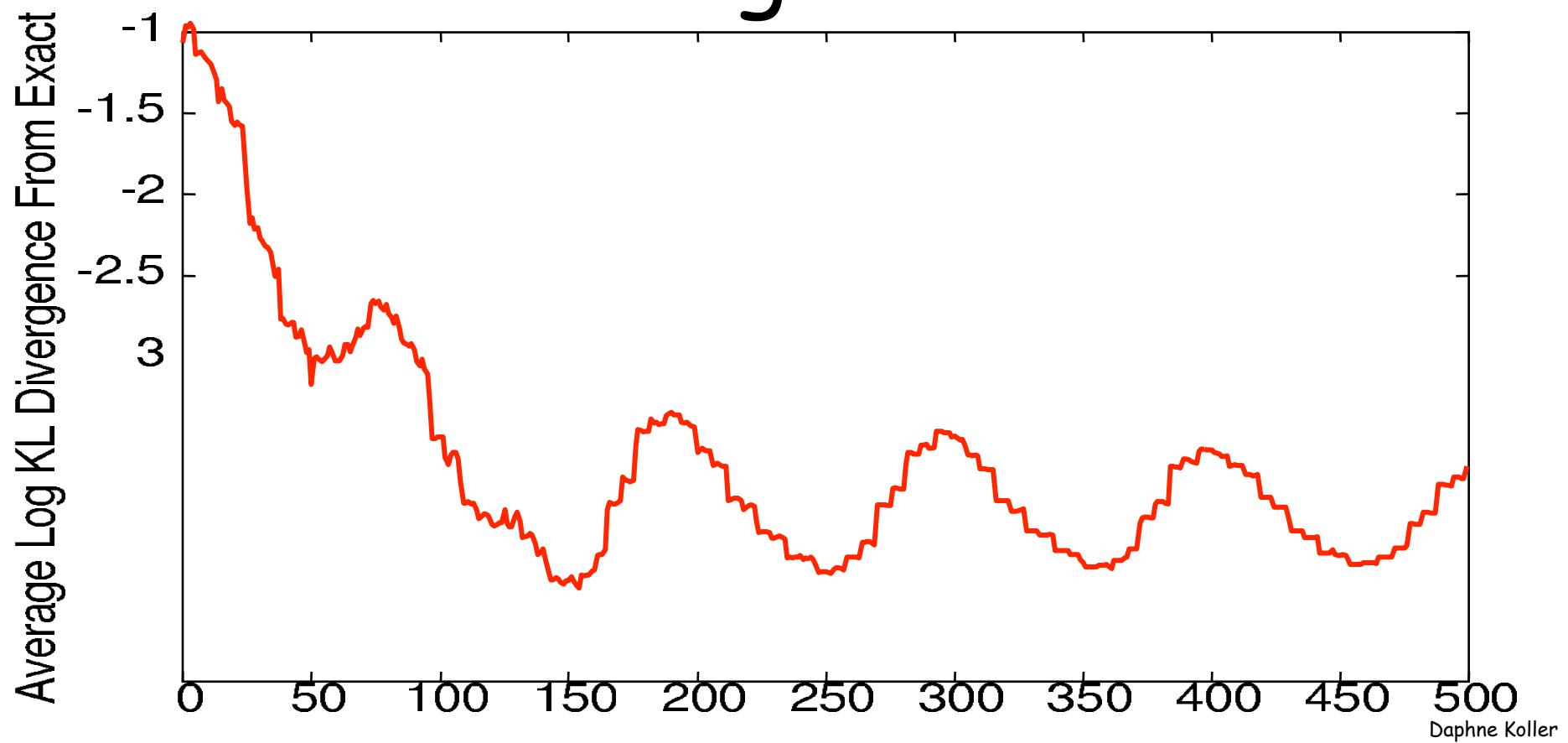
BP In Practice

# Misconception Revisited



Daphne Koller

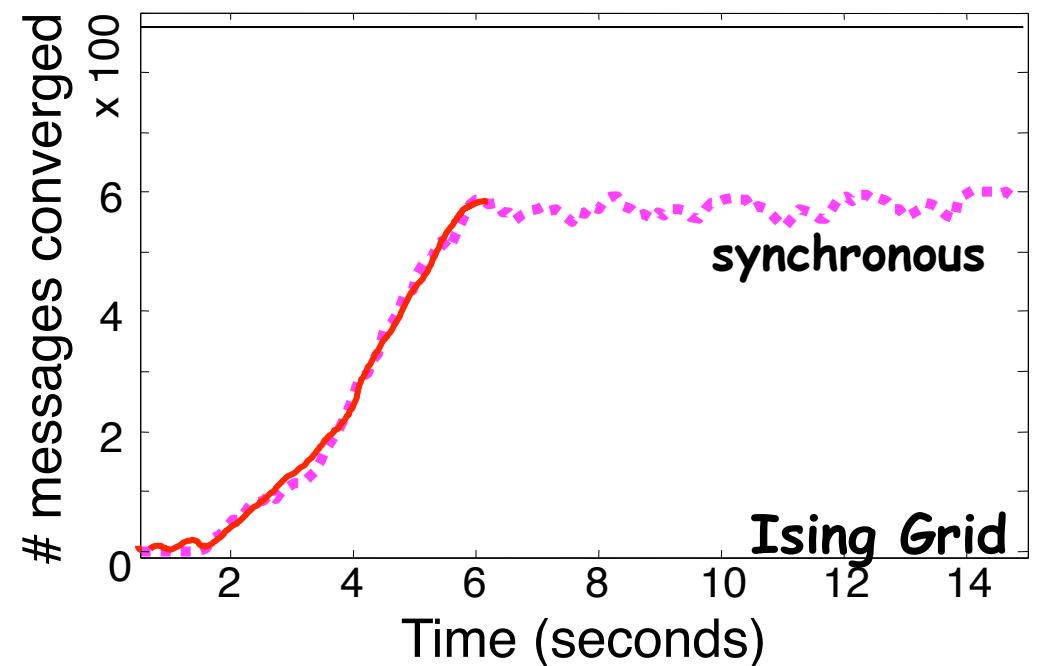
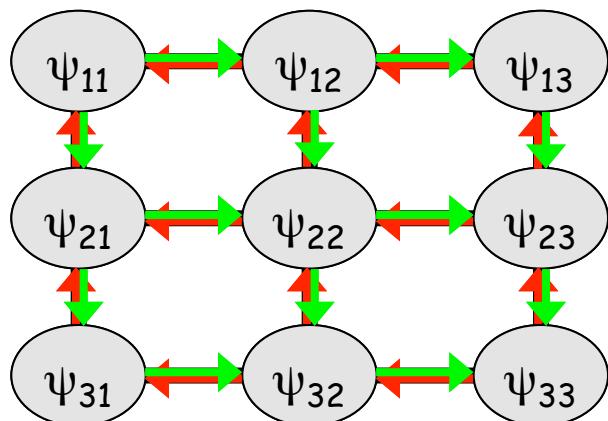
# Nonconvergent BP Run



Daphne Koller

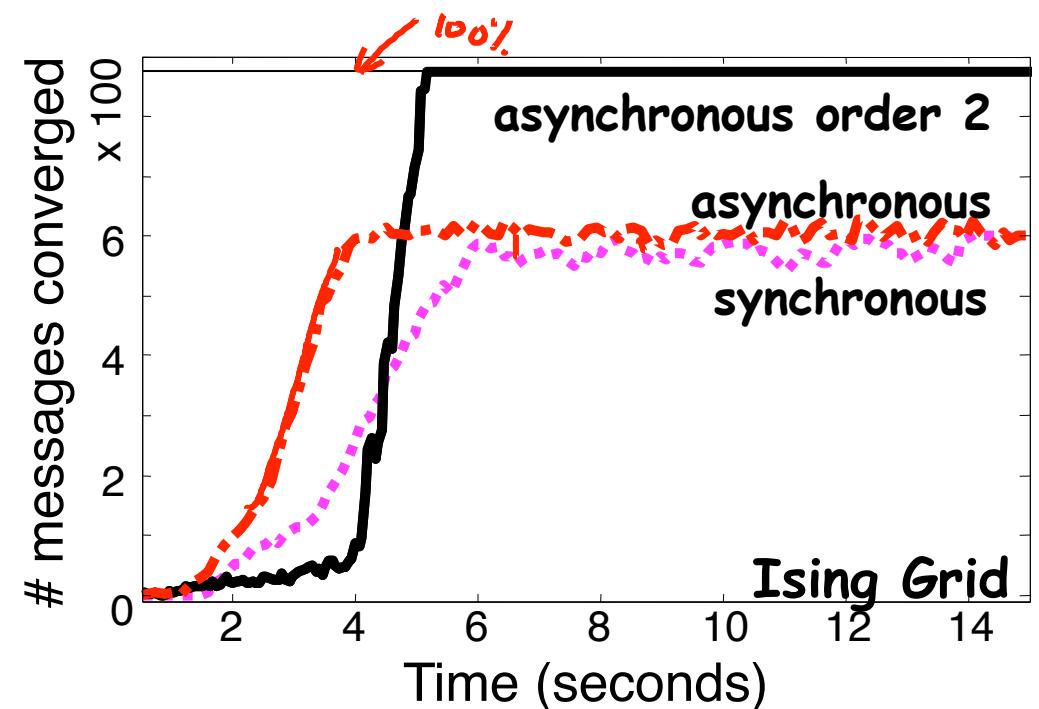
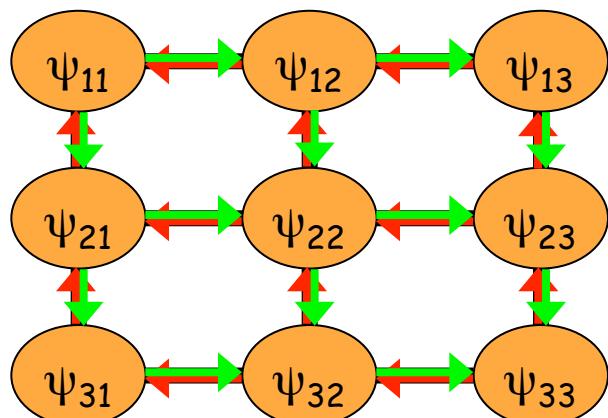
# Different Variants of BP

Synchronous BP:  
all messages are  
updated in parallel



# Different Variants of BP

**Asynchronous BP:**  
Messages are updated  
one at a time

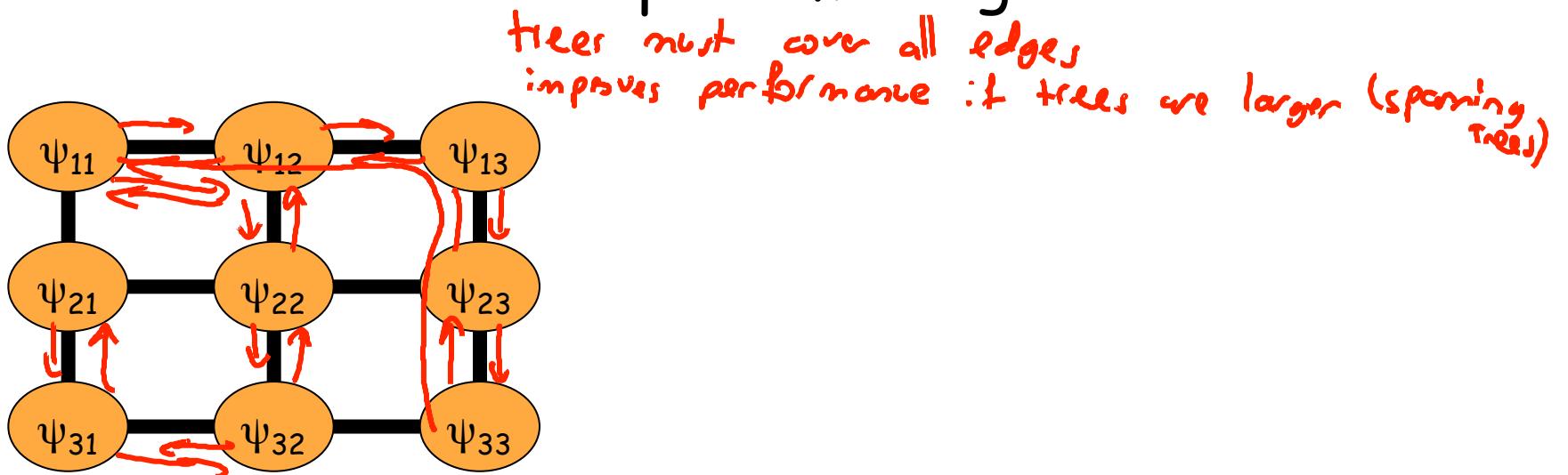


# Observations

- Convergence is a local property:
  - some messages converge soon
  - others may never converge
- Synchronous BP converges considerably worse than asynchronous
- Message passing order makes a difference to extent and rate of convergence

# Informed Message Scheduling

- Tree reparameterization (TRP)
  - Pick a tree and pass messages to calibrate



Daphne Koller

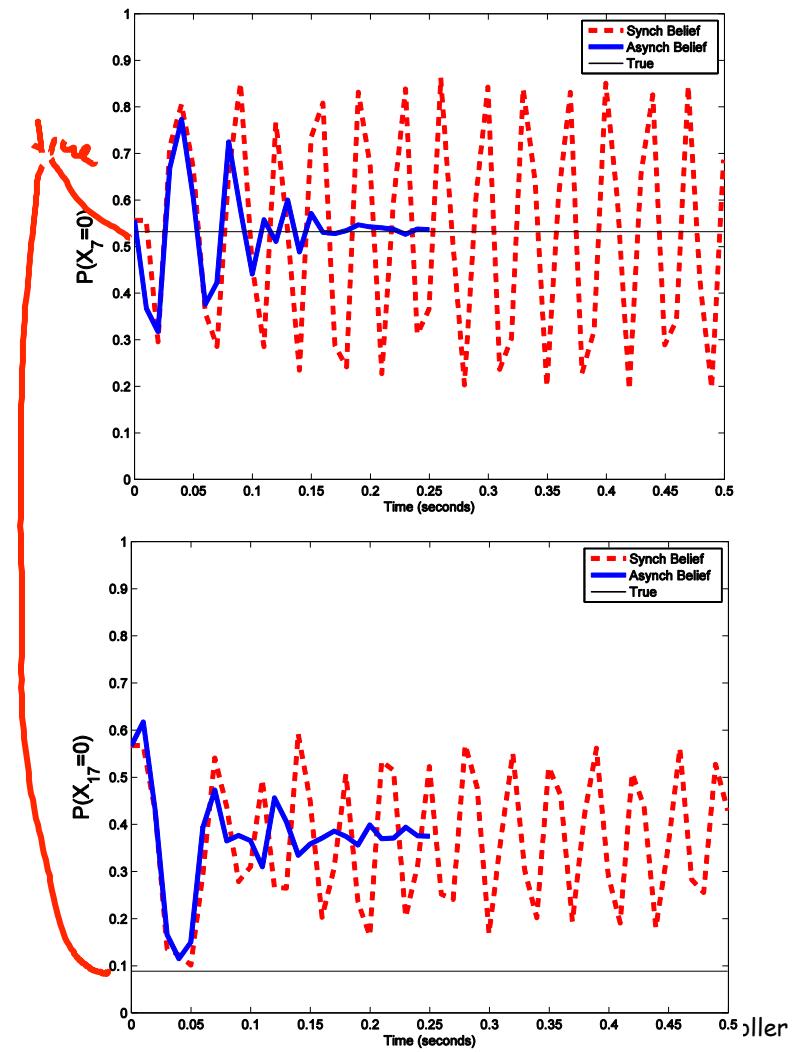
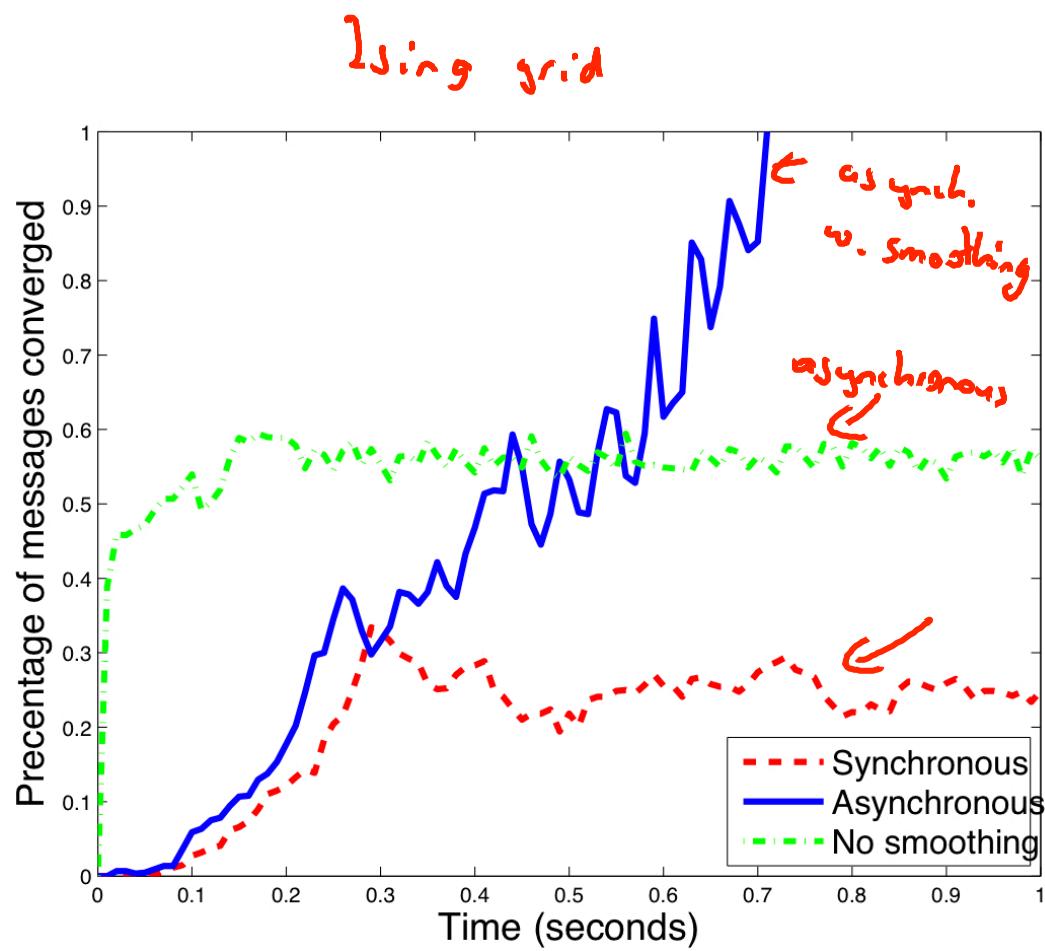
# Informed Message Scheduling

- Tree reparameterization (TRP)
  - Pick a tree and pass messages to calibrate
- Residual belief propagation (RBP)
  - Pass messages between two clusters whose beliefs over the sepset disagree the most  
*priorly in w of edges*

# Smoothing (Damping) Messages

$$\delta_{i \rightarrow j} \leftarrow \underbrace{\sum_{C_i - S_{i,j}} \psi_i \prod_{k \neq j} \delta_{k \rightarrow i}}_{\text{new msg}}$$
$$\delta_{i \rightarrow j} \leftarrow \underbrace{\lambda}_{\text{new msg}} \left( \sum_{C_i - S_{i,j}} \psi_i \prod_{k \neq j} \delta_{k \rightarrow i} \right) + (1 - \lambda) \underbrace{\delta_{i \rightarrow j}^{\text{old}}}_{\text{old msg}}$$

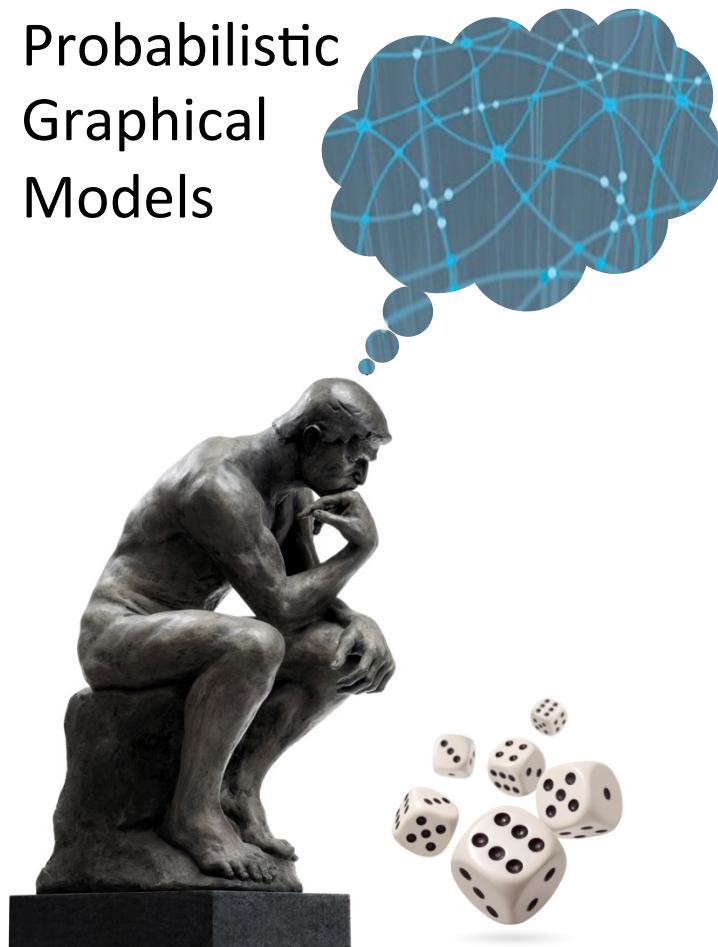
- Dampens oscillations in messages



# Summary

- To achieve BP convergence, two main tricks
  - Damping
  - Intelligent message ordering
- Convergence doesn't guarantee correctness
- Bad cases for BP – both convergence & accuracy:
  - Strong potentials pulling in different directions
  - Tight loops
- Some new algorithms have better convergence:
  - Optimization-based view to inference

Probabilistic  
Graphical  
Models



Inference

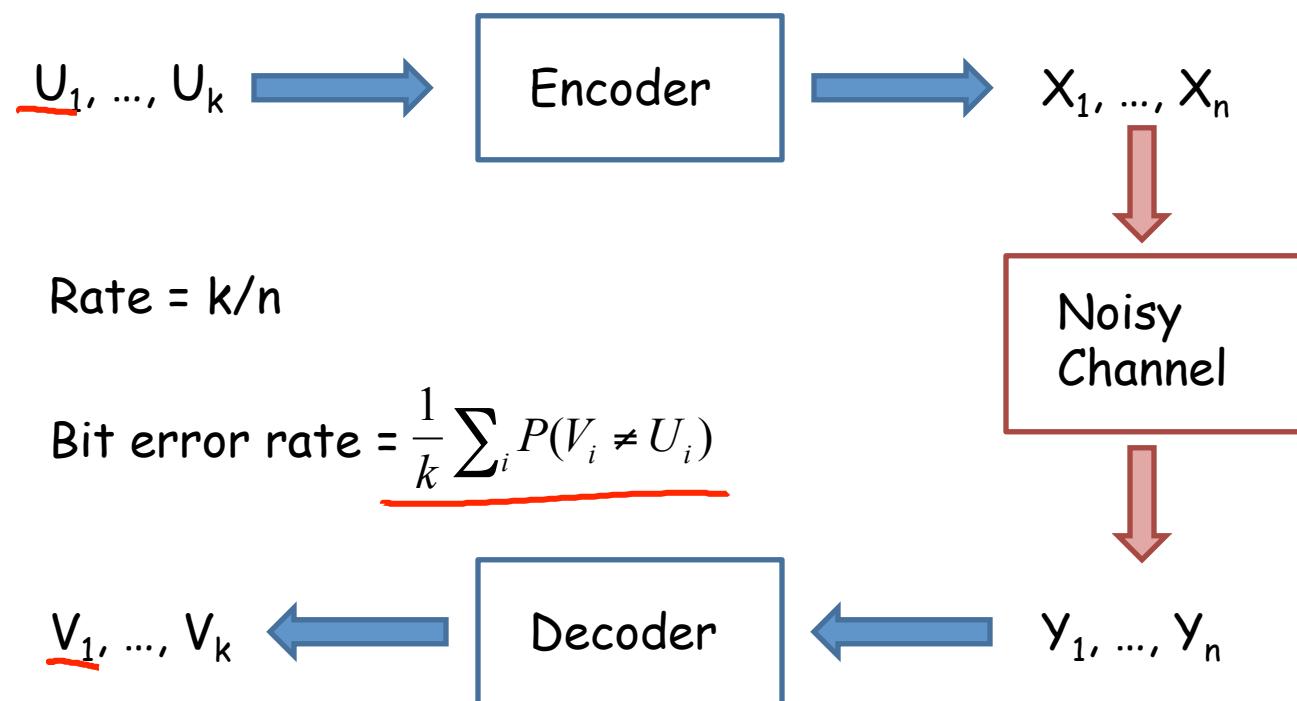
---

Message Passing

---

Loopy BP and  
Message  
Decoding

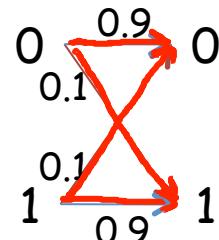
# Message Coding & Decoding



Noisy  
Channel

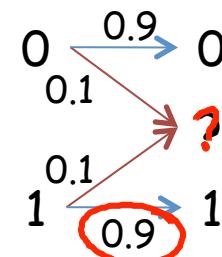
# Channel Capacity

Binary  
symmetric  
channel

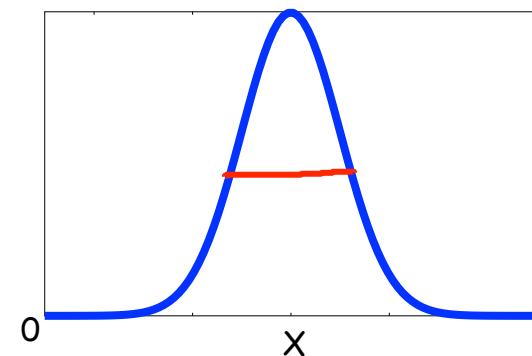


$$\text{capacity} = \underline{0.531}$$

Binary  
erasure  
channel

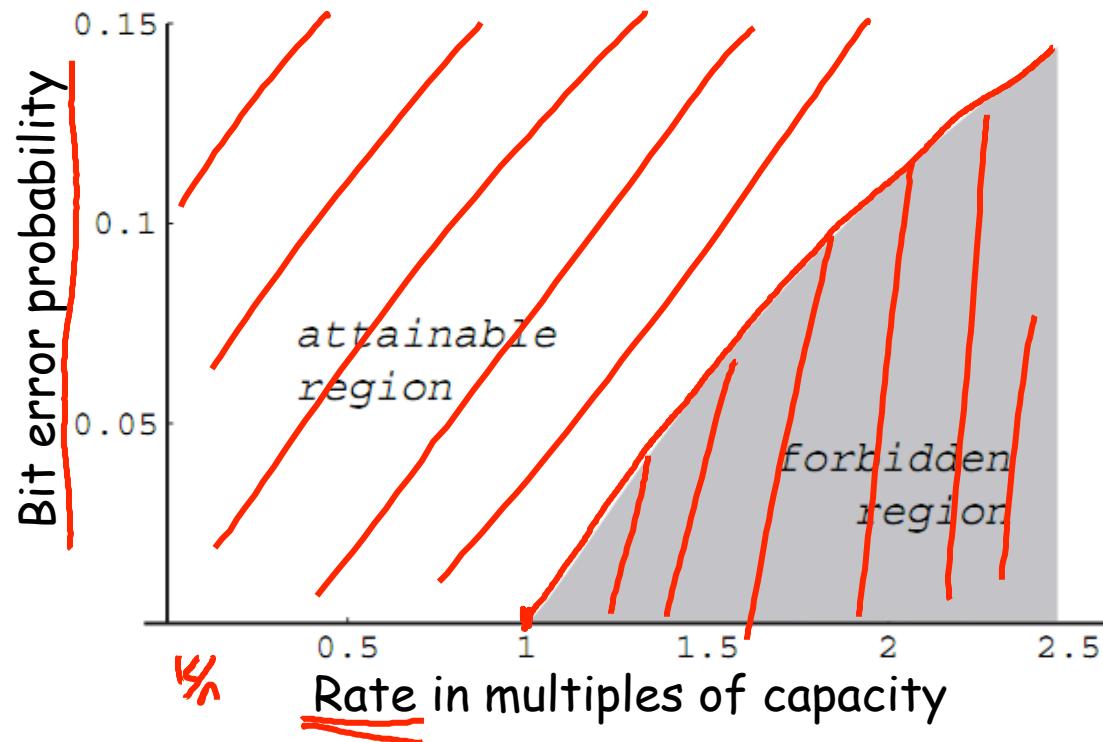


$$\text{capacity} = \underline{0.9}$$



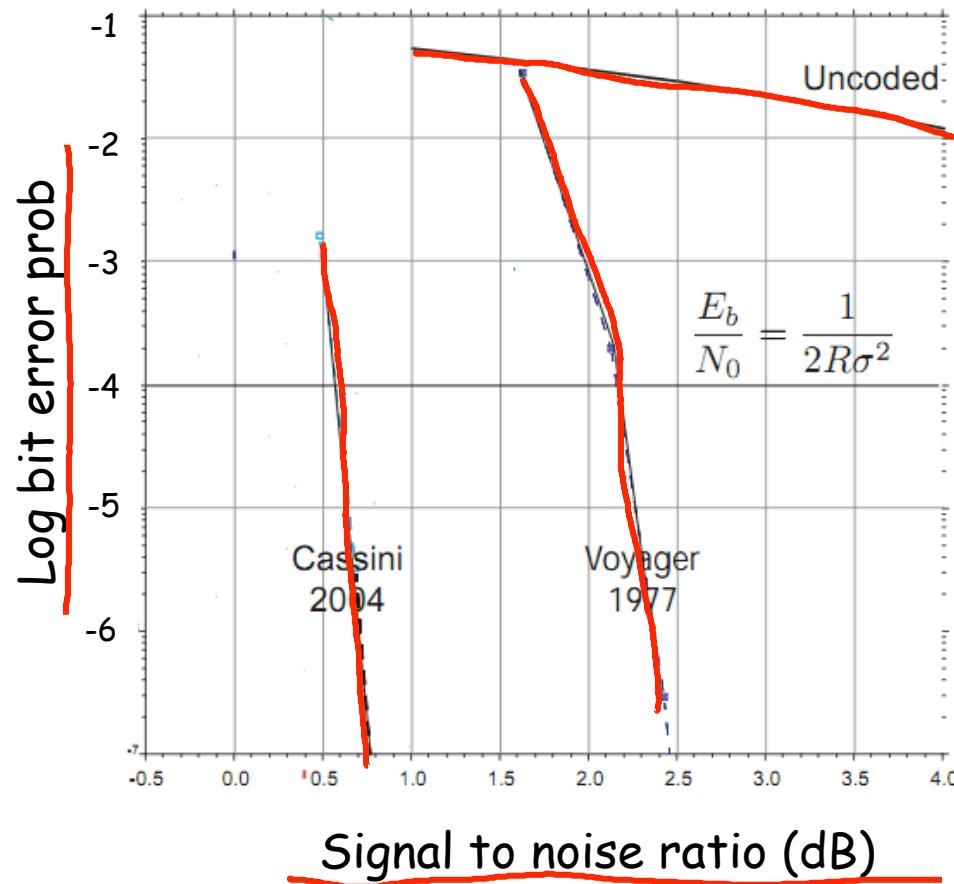
$$\text{capacity} = \underline{\frac{1}{2} \log \left( 1 + \frac{E(X^2)}{\sigma^2} \right)}$$

# Shannon's Theorem



McEliece

# How close to C can we get?



Daphne Koller

# Turbocodes (May 1993)

**NEAR SHANNON LIMIT ERROR - CORRECTING  
CODING AND DECODING : TURBO-CODES (1)**

Claude Berrou, Alain Glavieux and Punya Thitimajshima

Claude Berrou, Integrated Circuits for Telecommunication Laboratory

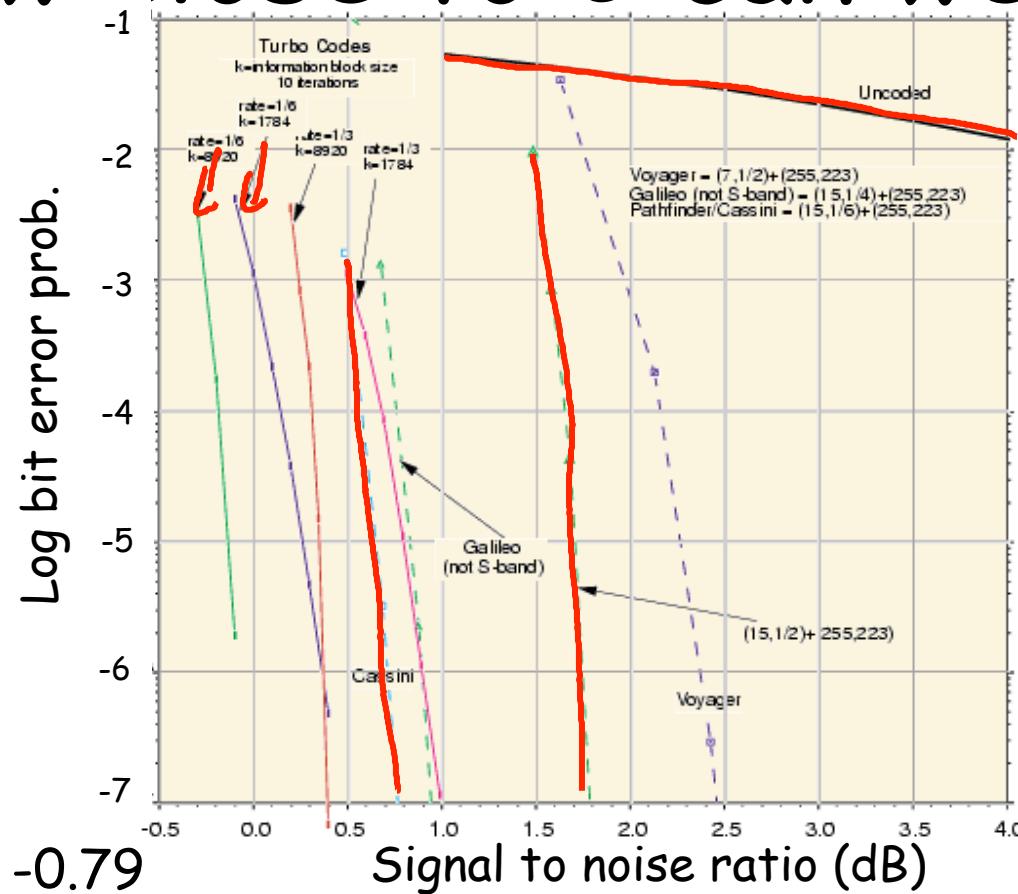
Alain Glavieux and Punya Thitimajshima, Digital Communication Laboratory

Ecole Nationale Supérieure des Télécommunications de Bretagne, France

(1) Patents N° 9105279 (France), N° 92460011.7 (Europe), N° 07/870,483 (USA)

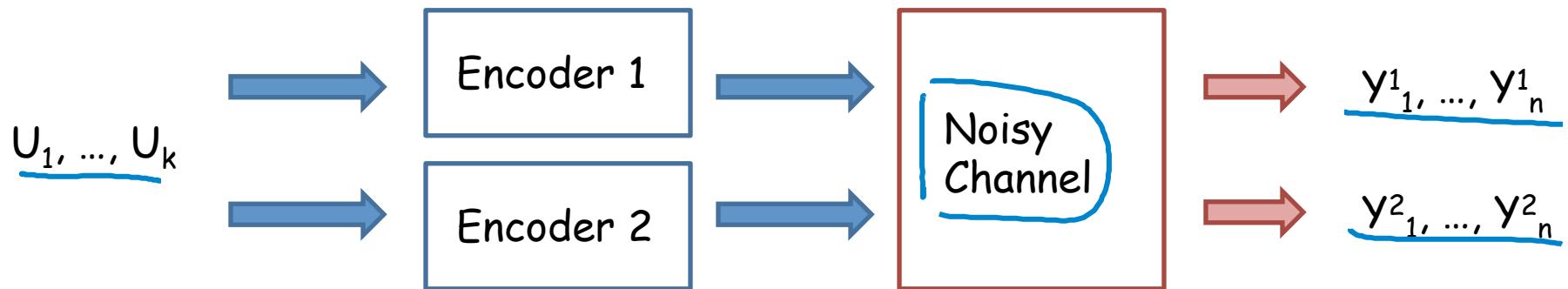
McEliece

# How close to C can we get?

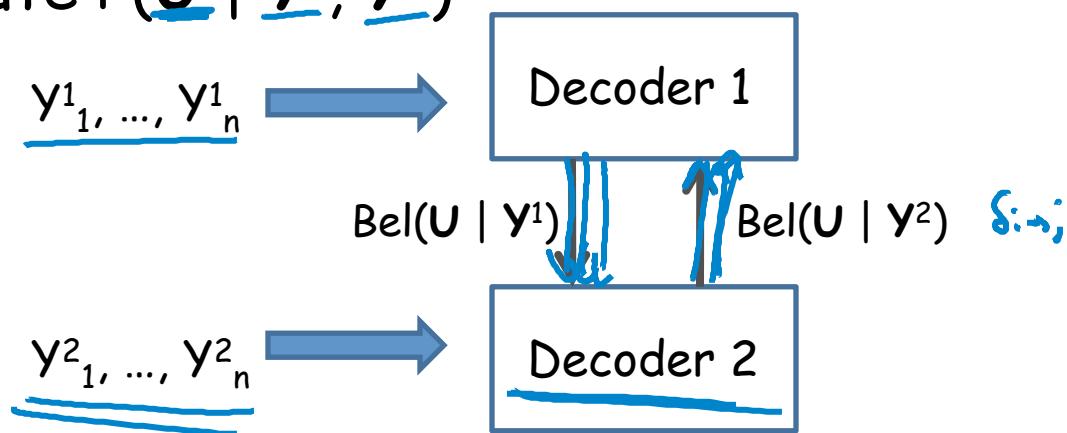


Daphne Koller

# Turbocodes: The Idea

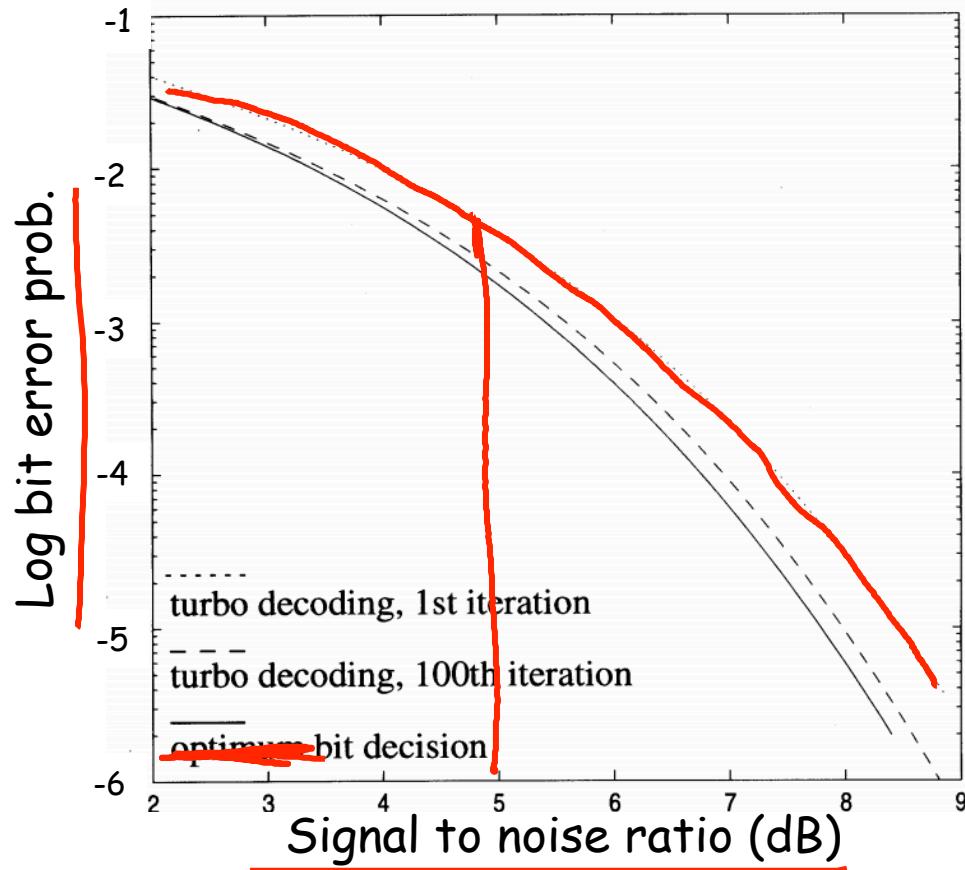


Compute  $P(\underline{U} \mid \underline{y^1}, \underline{y^2})$



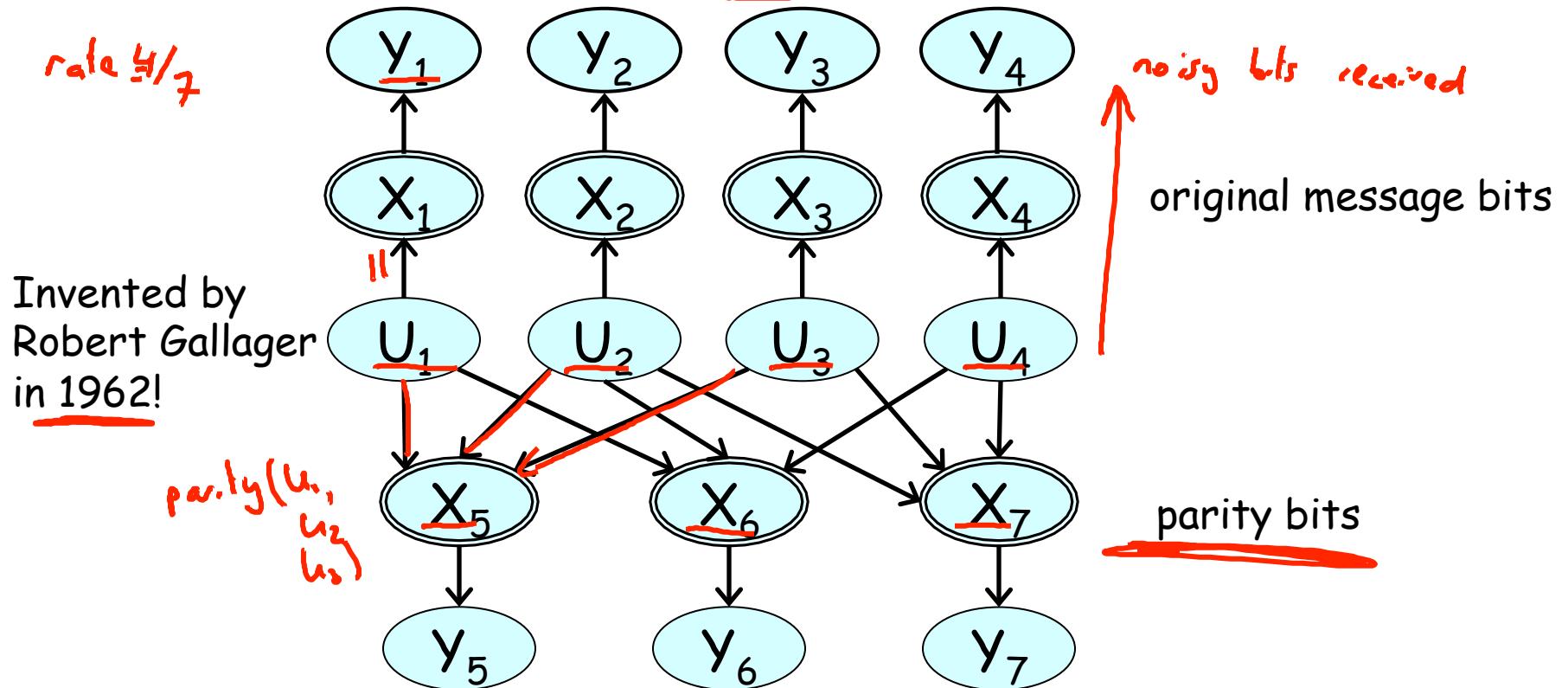
McEliece

# Iterations of Turbo Decoding

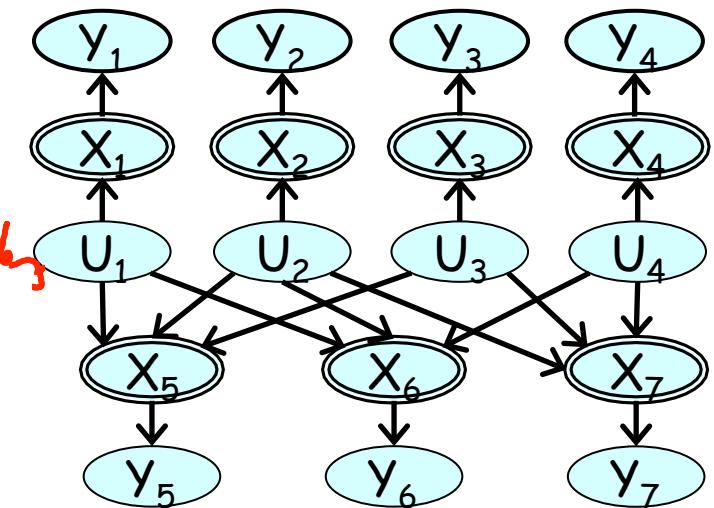
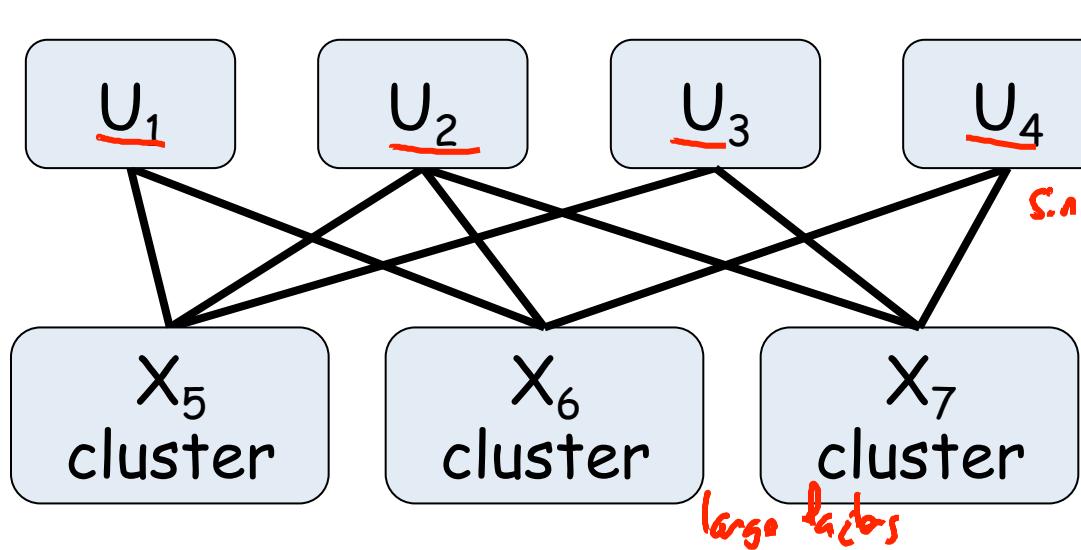


Daphne Koller

# Low-Density Parity Checking Codes



# Decoding as Loopy BP



# Turbo-Codes & LDPCs

- 3G and 4G mobile telephony standards
- Mobile television system from Qualcomm
- Digital video broadcasting
- Satellite communication systems
- New NASA missions (e.g., Mars Orbiter)
- Wireless metropolitan network standard

# Summary

- Loopy BP rediscovered by coding practitioners
- Understanding turbocodes as loopy BP led to development of many new and better codes
  - Current codes coming closer and closer to Shannon limit
- Resurgence of interest in BP led to much deeper understanding of approximate inference in graphical models
  - Many new algorithms