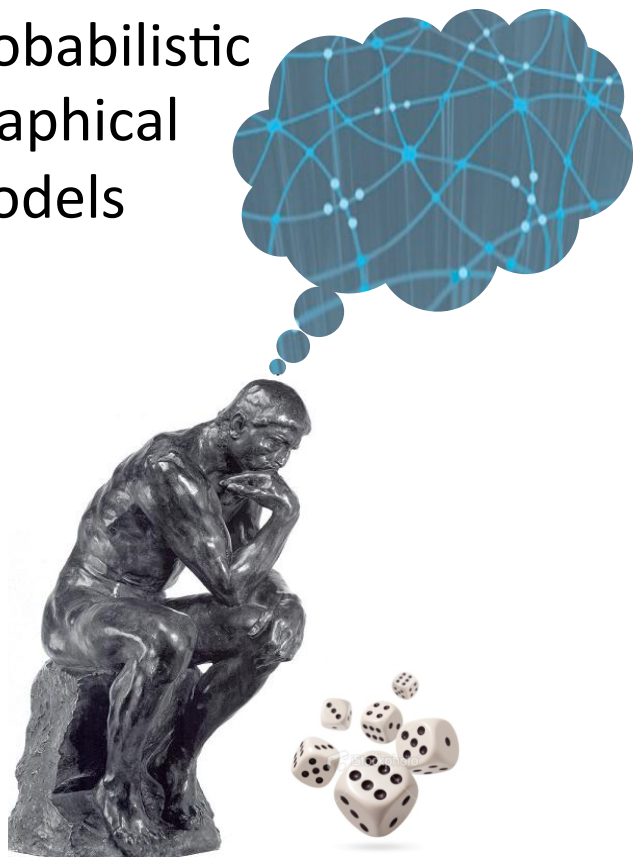


Probabilistic
Graphical
Models



Inference

Message Passing

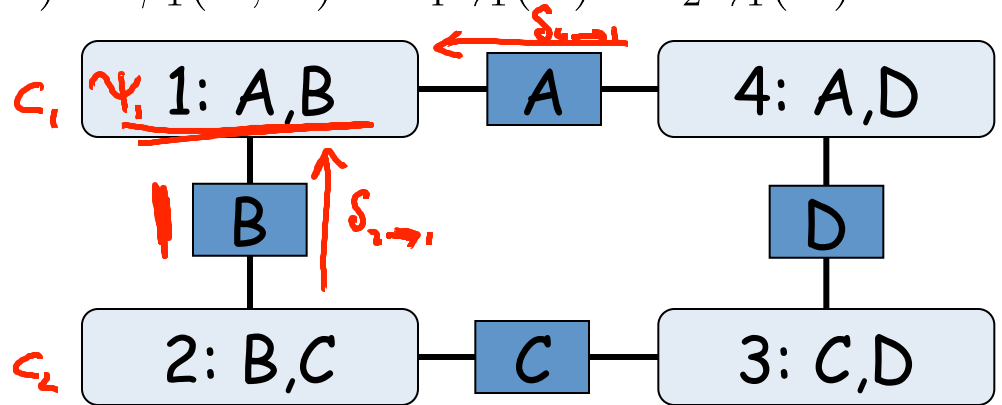
Properties of BP Algorithm

Calibration

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

- Cluster beliefs:

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



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

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

Convergence \Rightarrow Calibration

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

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

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

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

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

Reparameterization

- Sepset marginals:

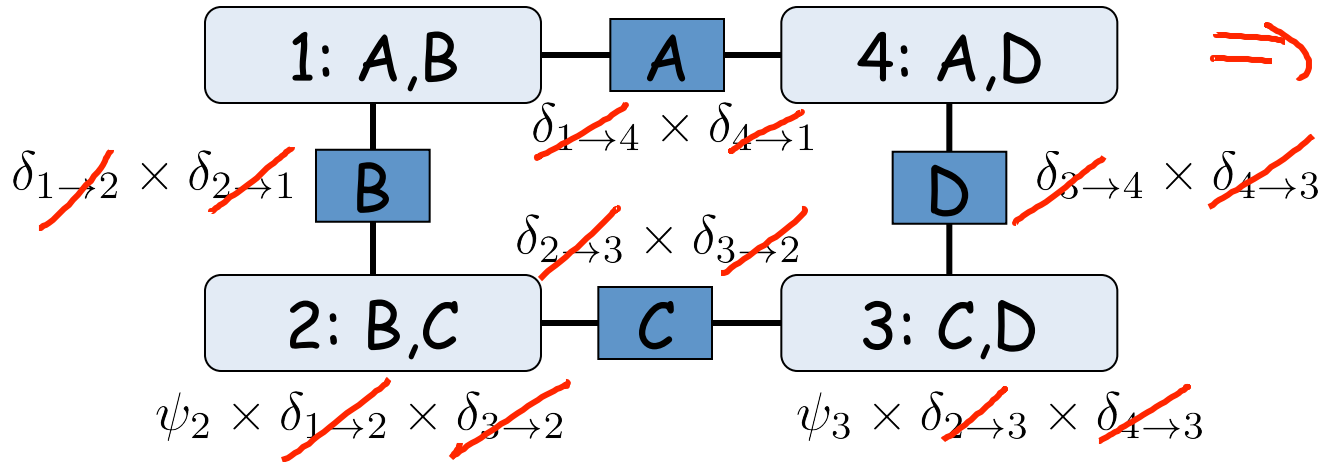
multiply all the beliefs
multiply sepsets

$$\rightarrow \psi_1 \times \cancel{\delta_{4 \rightarrow 1}} \times \cancel{\delta_{2 \rightarrow 1}}$$

$$\psi_4 \times \cancel{\delta_{1 \rightarrow 4}} \times \cancel{\delta_{3 \rightarrow 4}}$$

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

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



Reparameterization

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

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

$$\frac{\prod_i \beta_i}{\prod_{i,j} \mu_{i,j}} = \frac{\prod_i \psi_i \prod_{j \in \mathcal{N}_i} \delta_{j \rightarrow i}}{\prod_{i,j} \delta_{i \rightarrow j}}$$

initial potentials (pointing to ψ_i)
all msgs (pointing to $\delta_{j \rightarrow i}$)
all msgs (pointing to $\delta_{i \rightarrow j}$)

$$= \prod_i \psi_i = \tilde{P}_\Phi(X_1, \dots, X_n)$$

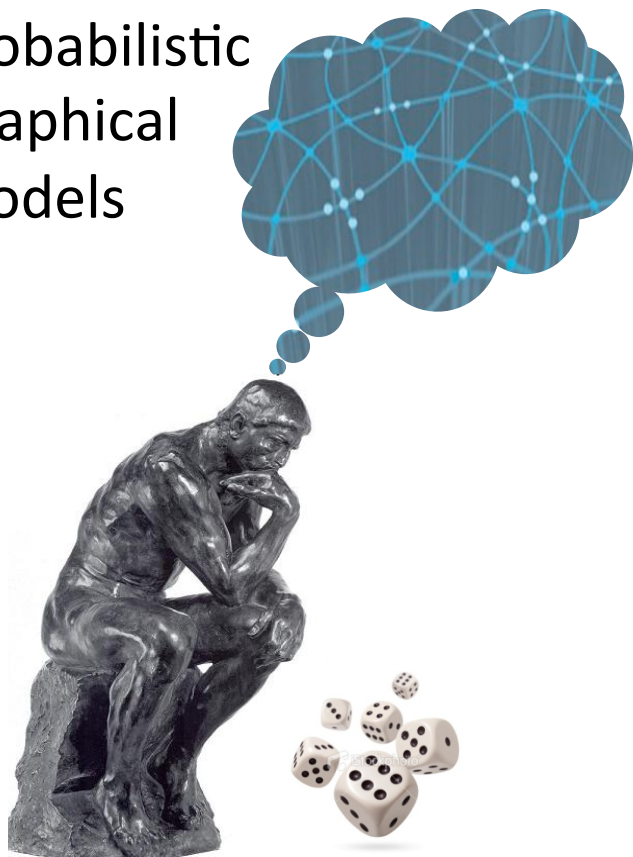
initial factors ψ

*different parameterization
 $\beta_i, \mu_{i,j}$
 (calibrated)*

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



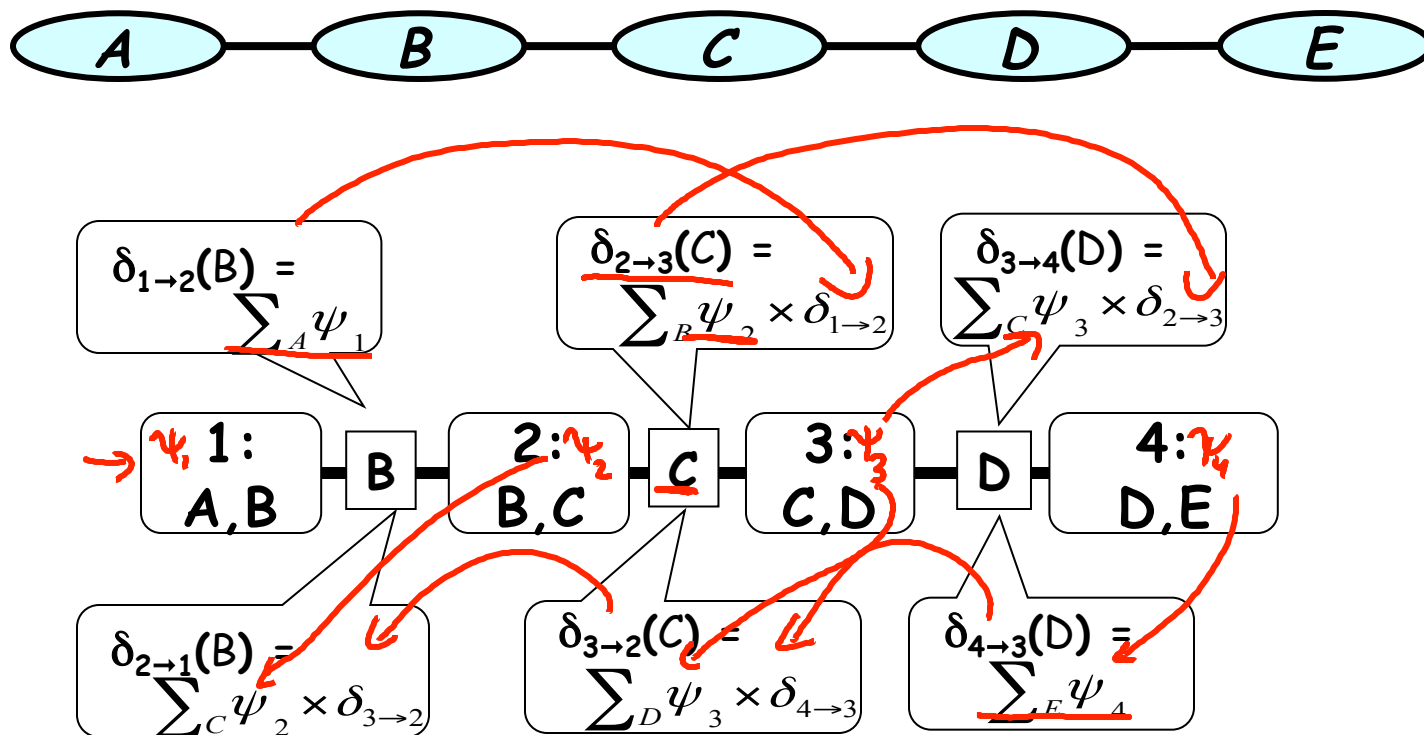
Inference

Message Passing

Clique Tree

Algorithm & Correctness

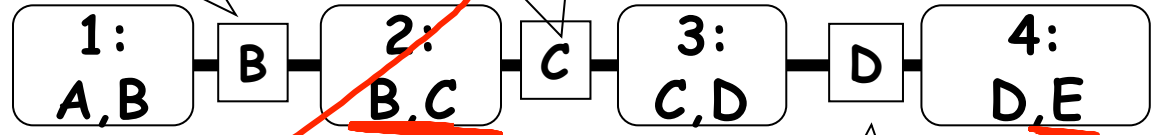
Message Passing in Trees



Correctness

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

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



$$\beta_3(C, D) = \psi_3 \times \delta_{2 \rightarrow 3} \times \delta_{4 \rightarrow 3}$$

$$= \psi_3 \times \left(\sum_B (\psi_2 \times \delta_{1 \rightarrow 2}) \right) \times \sum_E \psi_4$$

$$= \psi_3 \times \left(\sum_B \psi_2 \times \sum_A \psi_1 \right) \times \sum_E \psi_4$$

legal order of operations

→

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

product of factors
marginalized out unnecessary
variables

Clique Tree

- Undirected tree such that:
 - nodes are clusters $\mathcal{C}_i \subseteq \{X_1, \dots, X_n\}$
 - edge between \mathcal{C}_i and \mathcal{C}_j associated with sepset $\mathcal{S}_{i,j} = \mathcal{C}_i \cap \mathcal{C}_j$

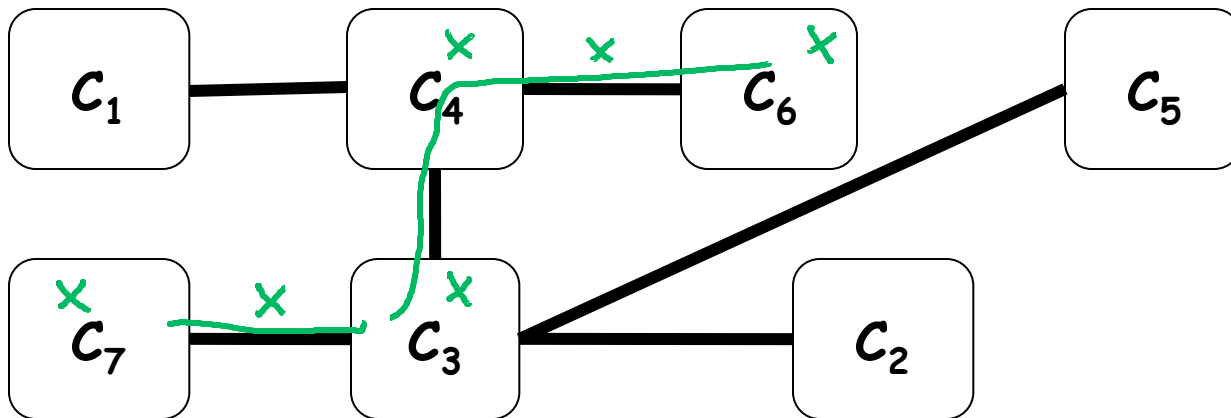
Family Preservation

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

Running Intersection Property

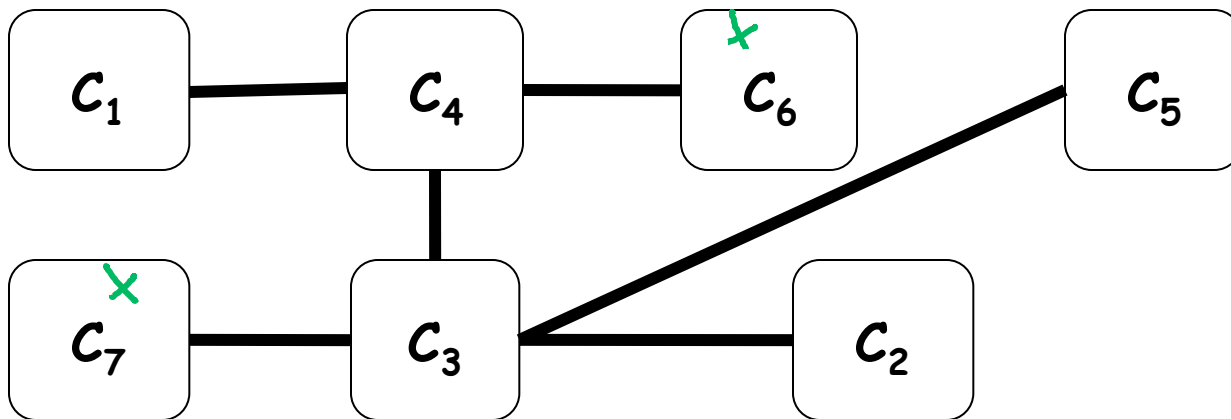
Cluster graph variant

- 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

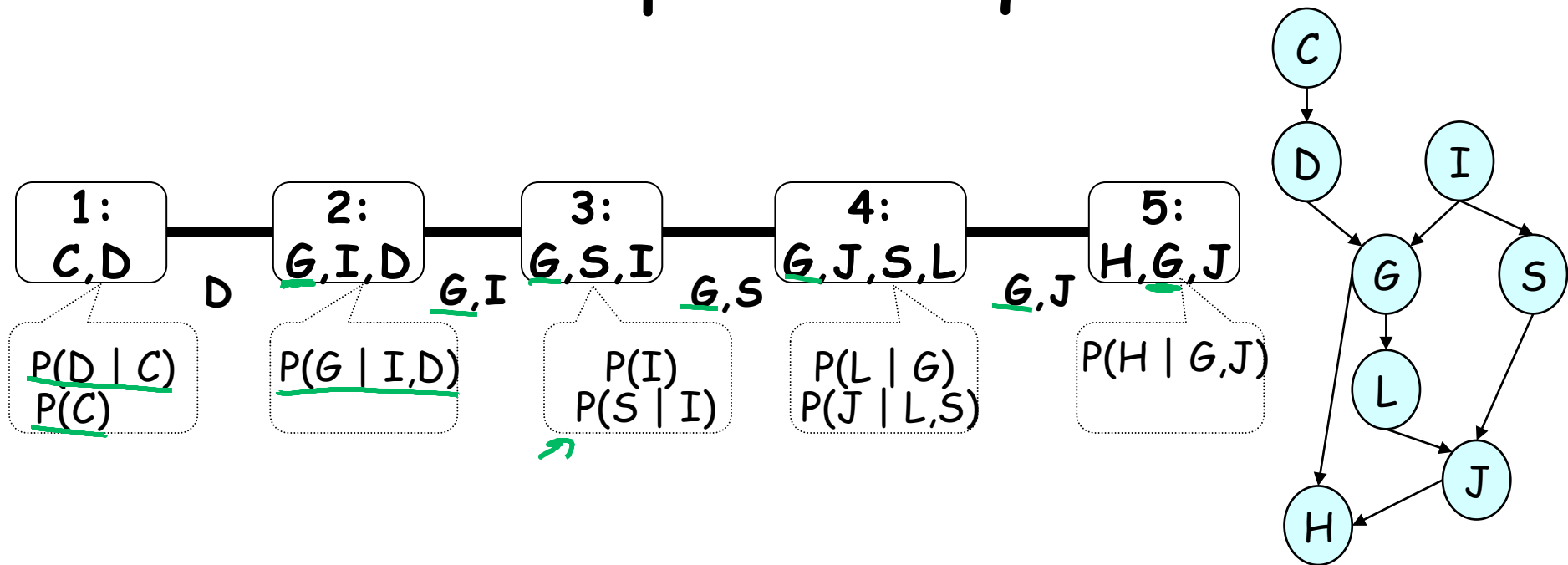


Running Intersection Property

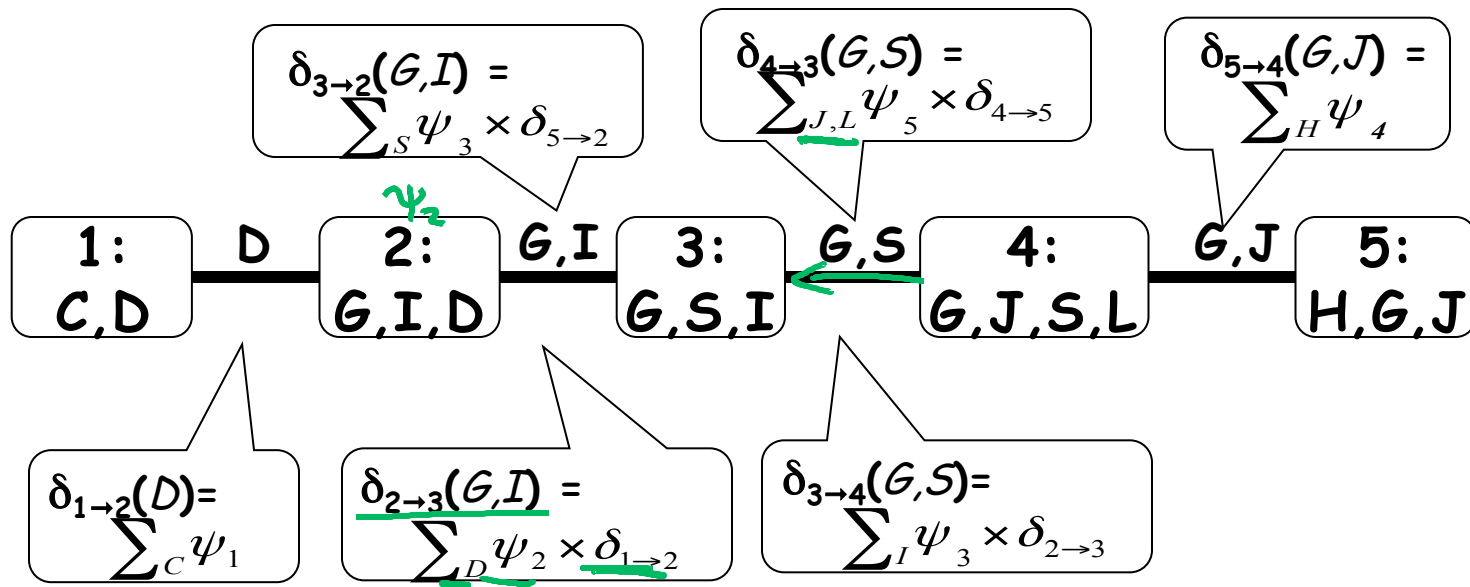
- For each pair of clusters C_i, C_j and variable $X \in C_i \cap C_j$, in the unique path between C_i and C_j , all clusters and sepsets contain X



More Complex Clique Tree

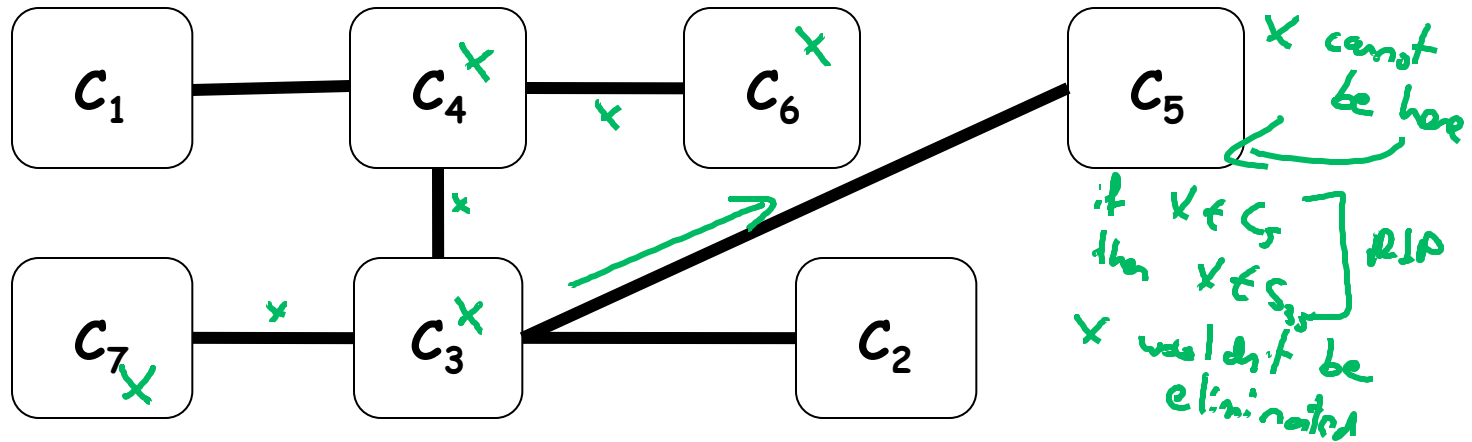


Clique Tree Message Passing

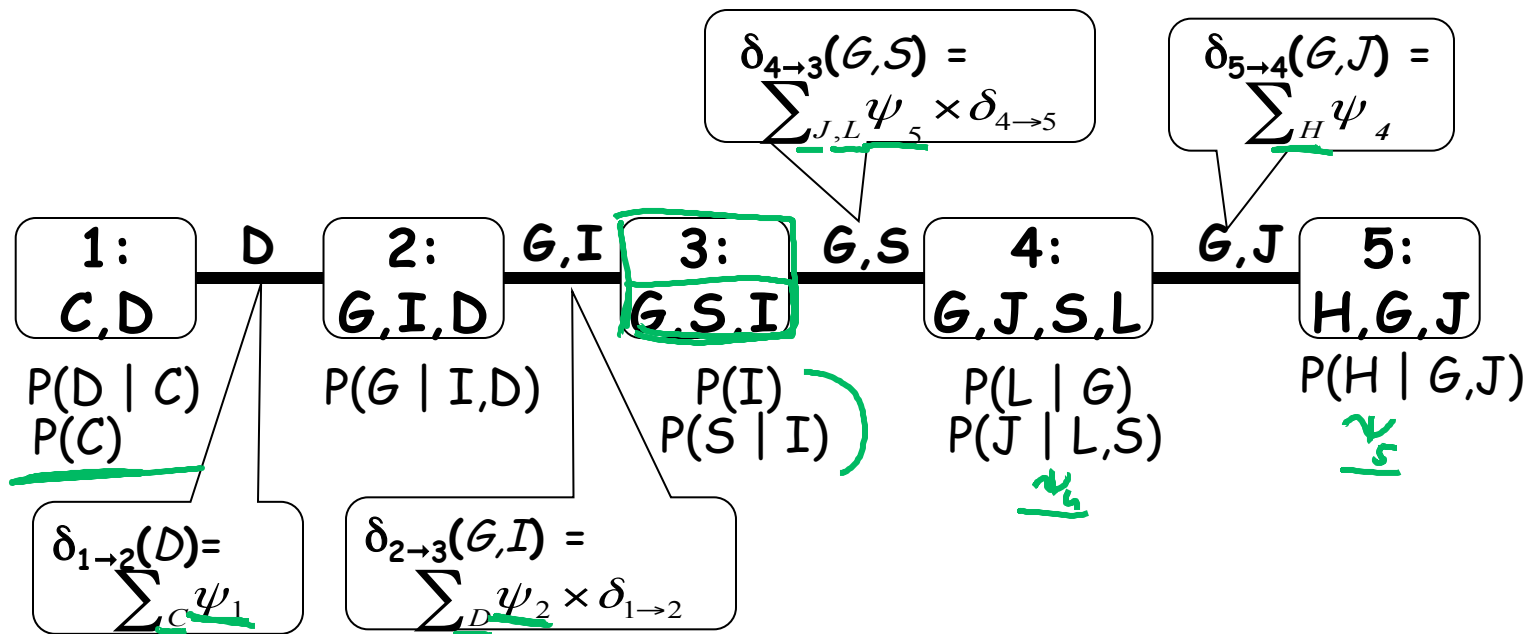


RIP \Rightarrow Clique Tree Correctness


- 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



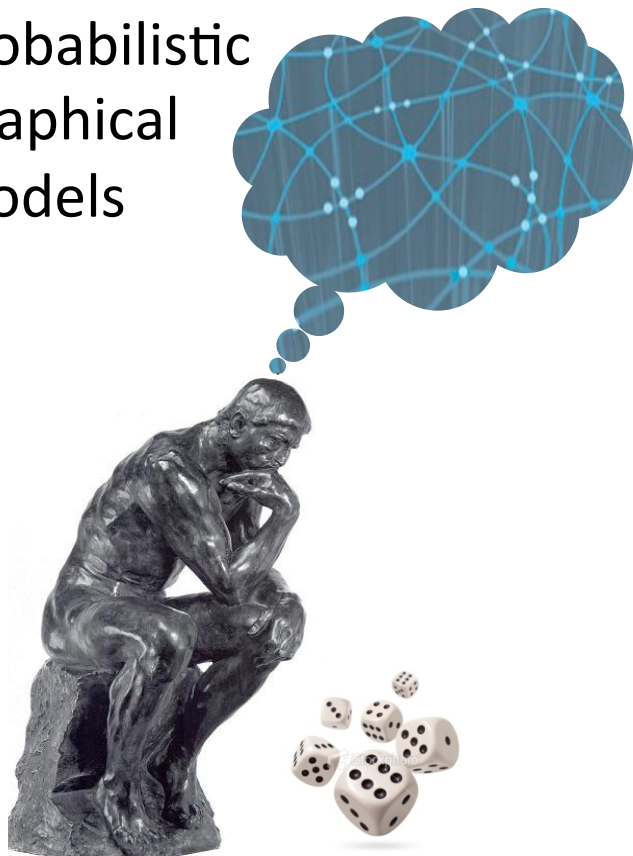
Clique Tree Correctness



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 

Probabilistic
Graphical
Models

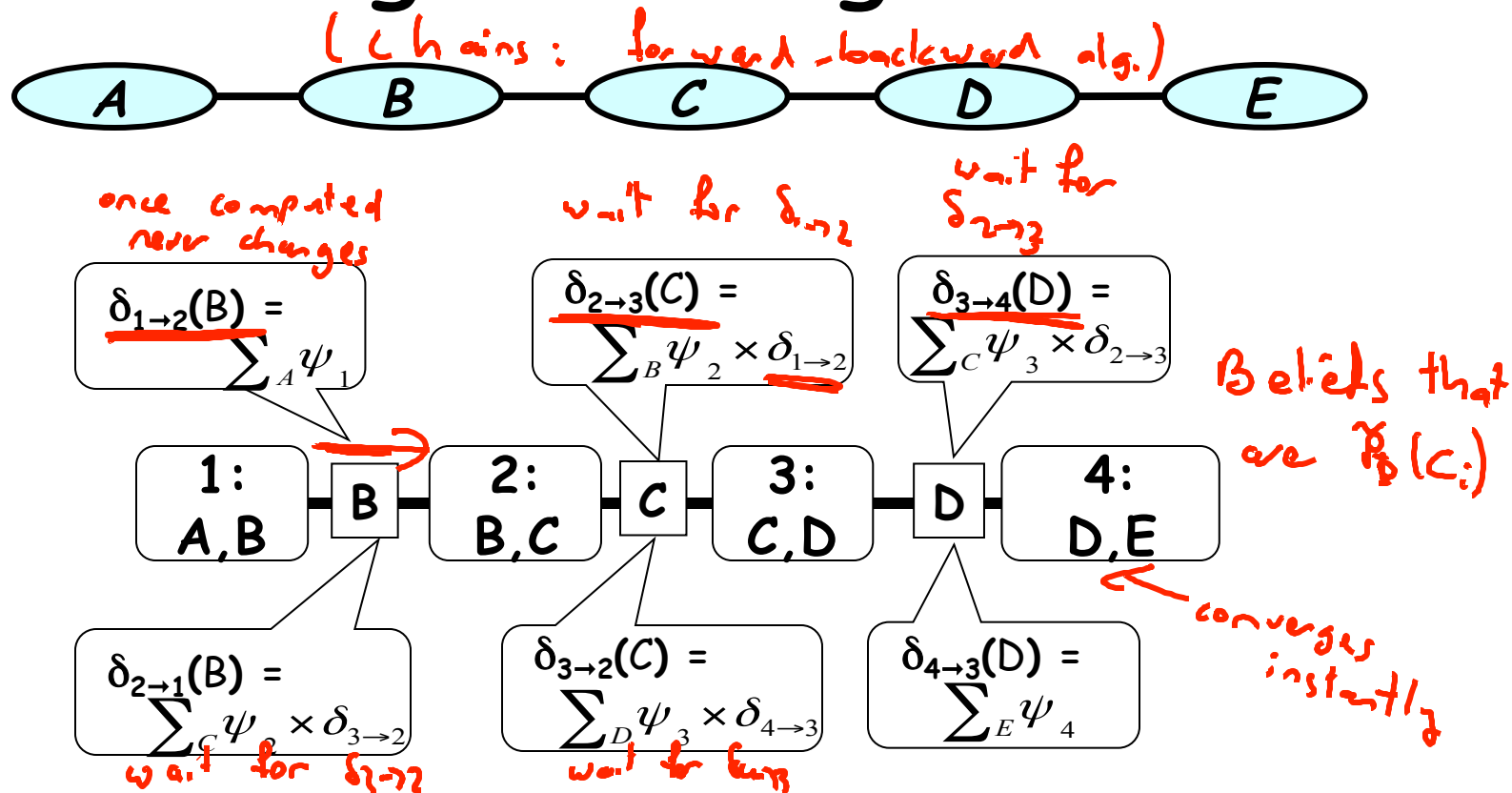


Inference

Message Passing

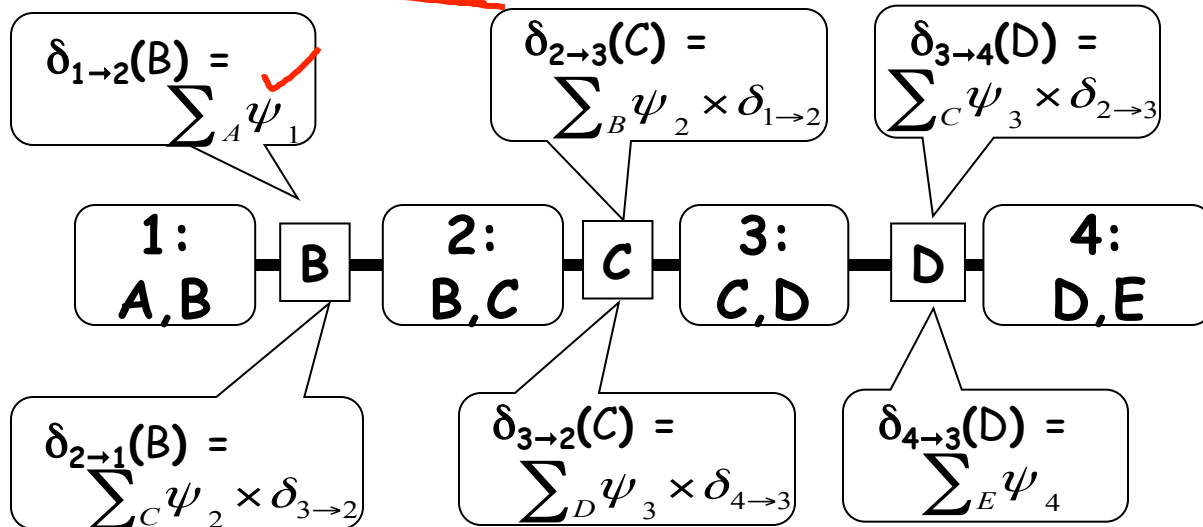
Clique Tree Algorithm: Computation

Message Passing in Trees



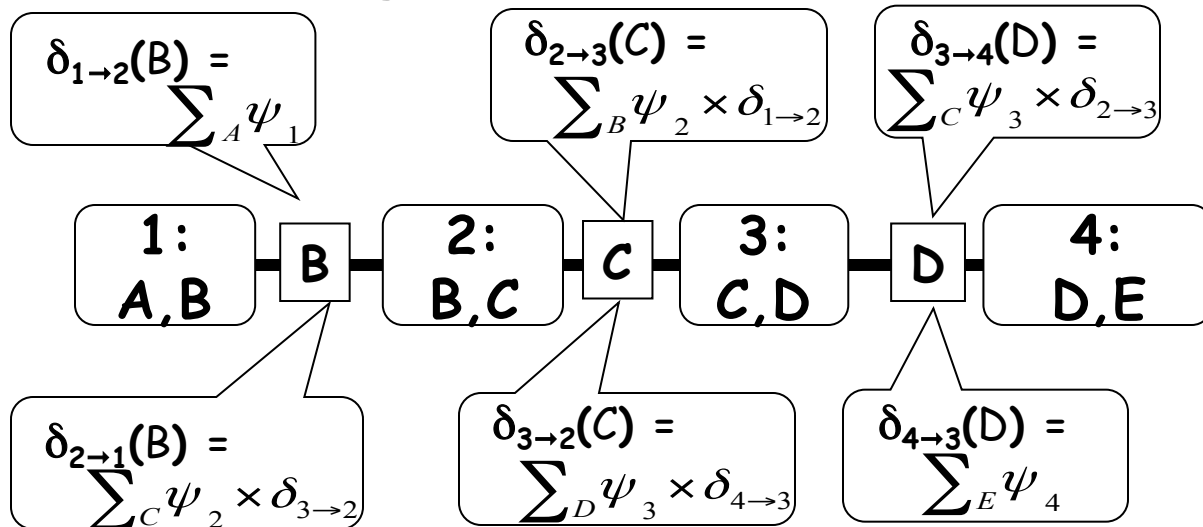
Convergence of Message Passing

- 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



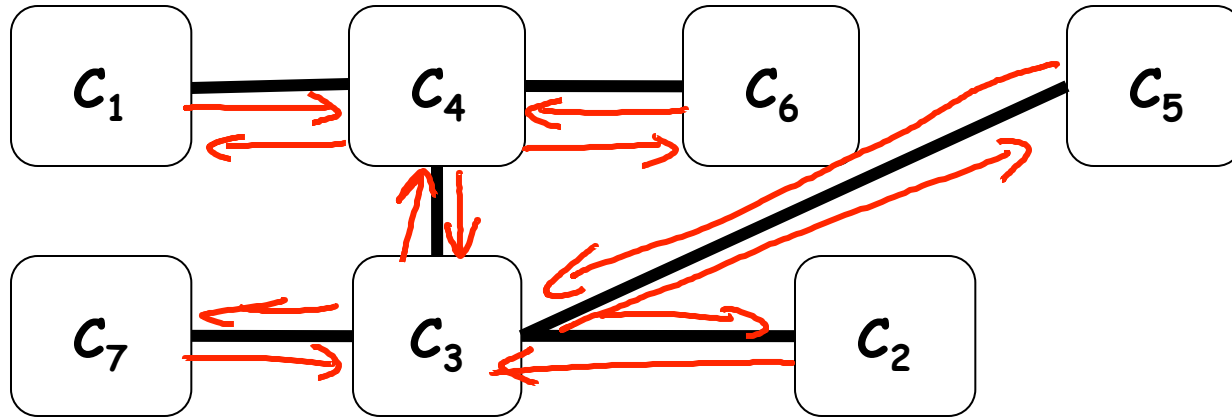
Convergence of Message Passing

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

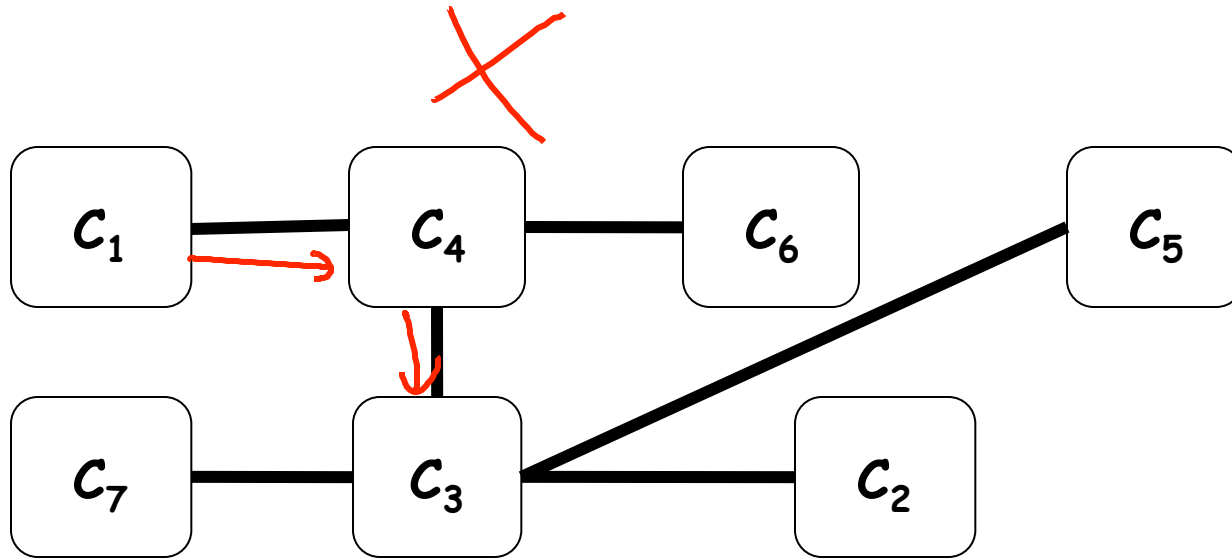


K total number of clusters.
 \Downarrow
 $K-1$ edges
 \Downarrow
 $2(K-1)$ msgs

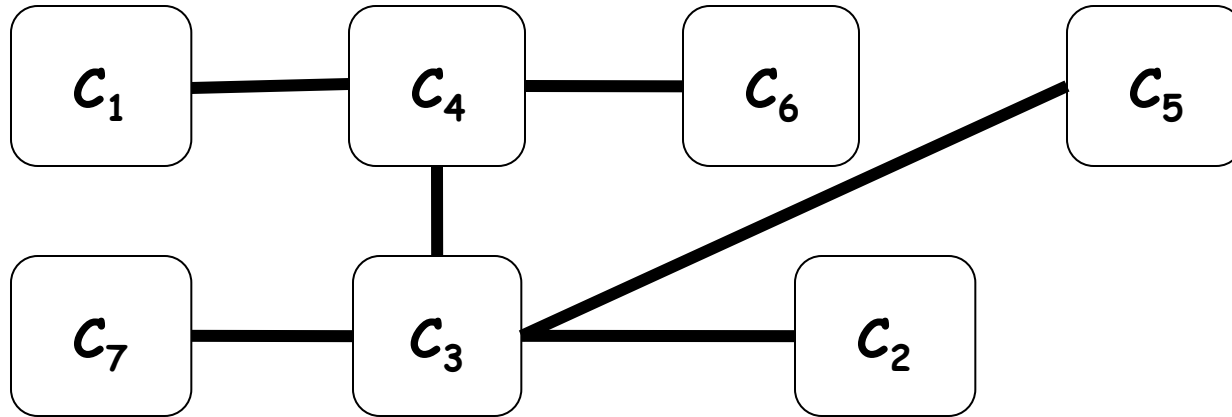
Message Passing Order I




Message Passing Order II



Message Passing Order III

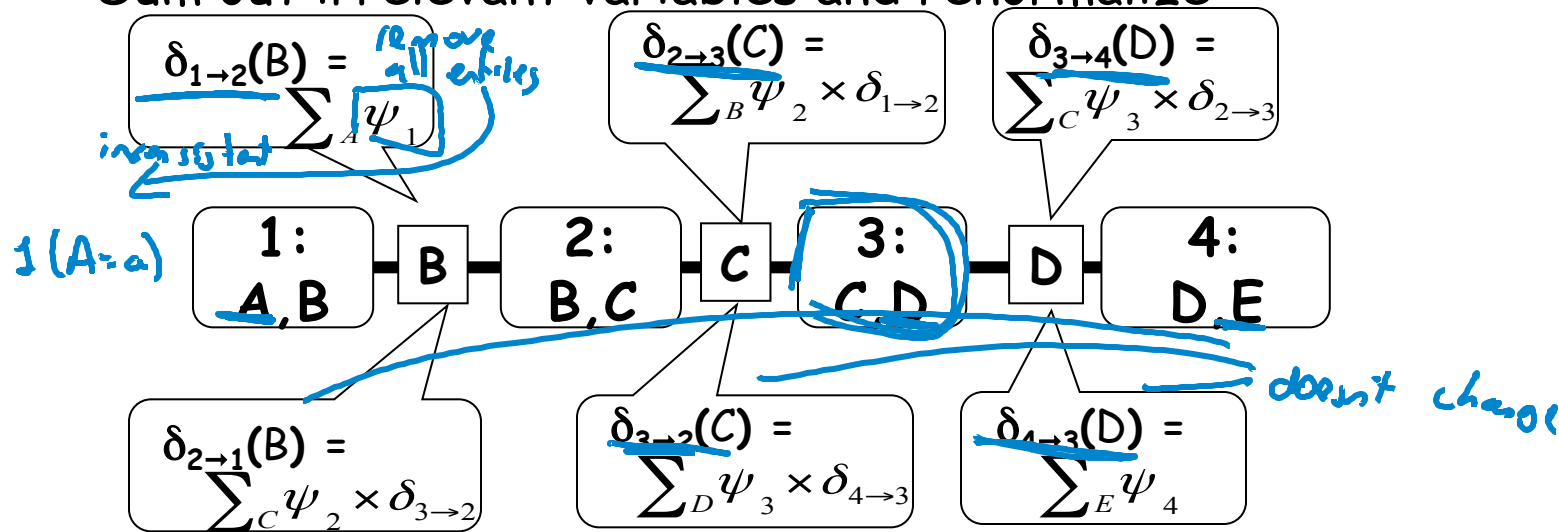


Answering Queries

- Posterior distribution queries on variables that appear together in clique
 - Sum out irrelevant variables from any clique containing those variables \mathcal{P}_ϕ renormalize
- Introducing new evidence $Z=\underline{z}$ 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 $\mathcal{P}_\phi(z, X)$ 
 - Sum out irrelevant variables and renormalize

And More Queries

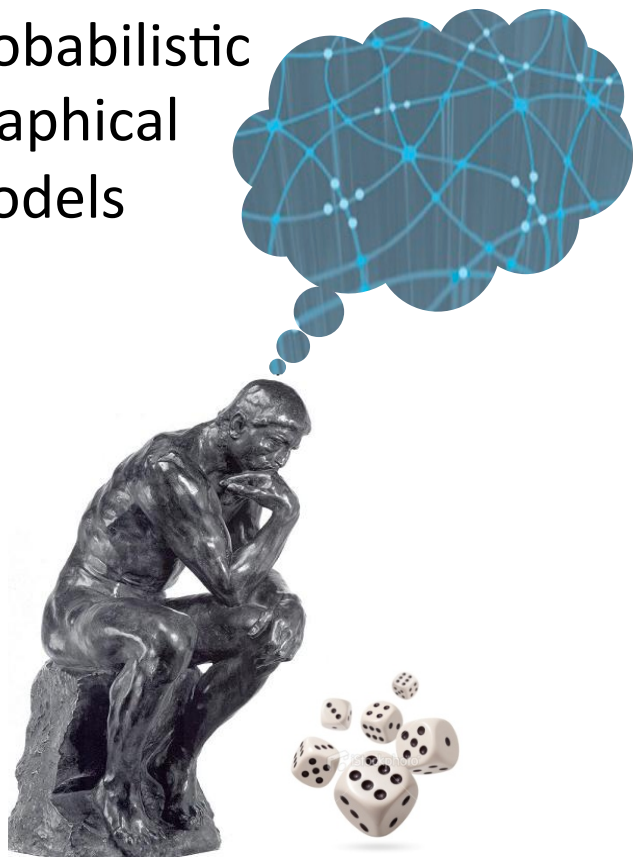
- Introducing new evidence $Z=z$ and querying X if X does not share a clique with Z
 - Multiply $1(Z=z)$ into some clique containing Z *reduction of factor*
 - Propagate messages along path to clique containing X
 - Sum out irrelevant variables and renormalize



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



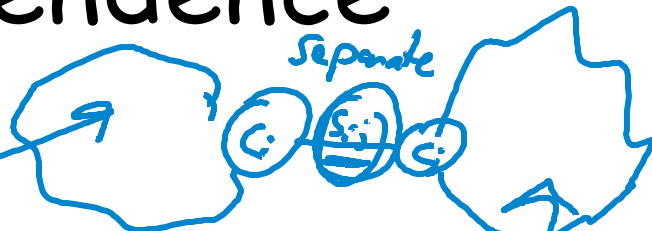
Inference

Message Passing

Clique Tree & Independence

RIP and Independence

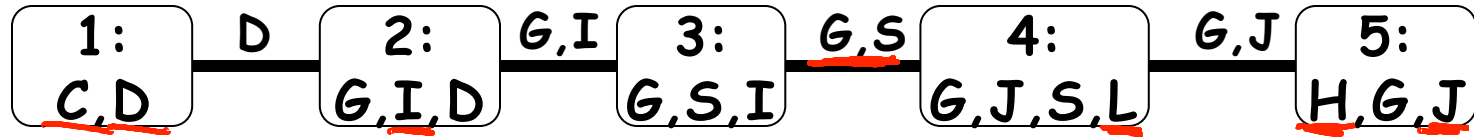
- For an edge (i,j) in T , let:



- $\mathbf{W}_{\langle(i,j)}$ = all variables that appear only on C_i side of T
- $\mathbf{W}_{\langle(j,i)}$ = all variables that appear only on C_j side
- Variables on both sides are in the sepset $S_{i,j}$

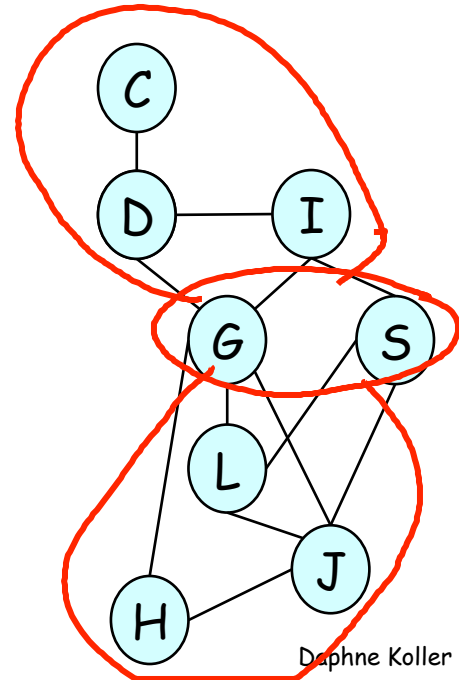
- Theorem:** T satisfies RIP if and only if, for every (i,j) $P_\Phi \models (\mathbf{W}_{\langle(i,j)} \perp \mathbf{W}_{\langle(j,i)} \mid \mathbf{S}_{i,j})$

RIP and Independence



$P_{\mathcal{G}} \models (\{G, I, D\} \perp \{J, L, H\} \mid \{G, S\})$
separate C, I, D from H, G, J

$\Leftrightarrow (C, I, D \perp H, G, J \mid C, D)$

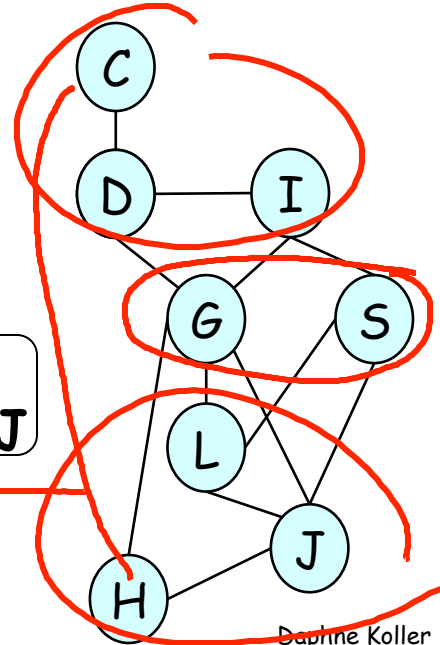


RIP and Independence

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

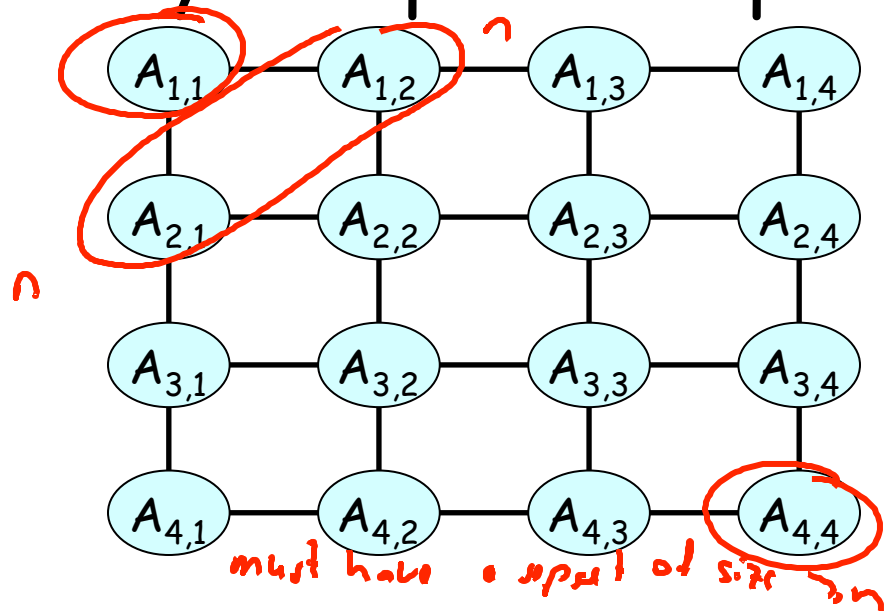
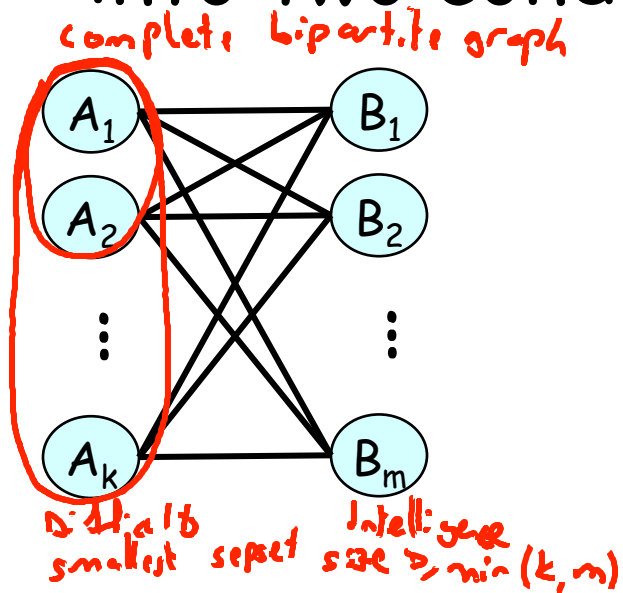
Assume otherwise $\Rightarrow \exists$ path in induced Markov network between $W_{<(i,j)}$ and $W_{<(j,i)}$ that doesn't go through $S_{i,j}$

Factor $\phi(C, H)$



Implications

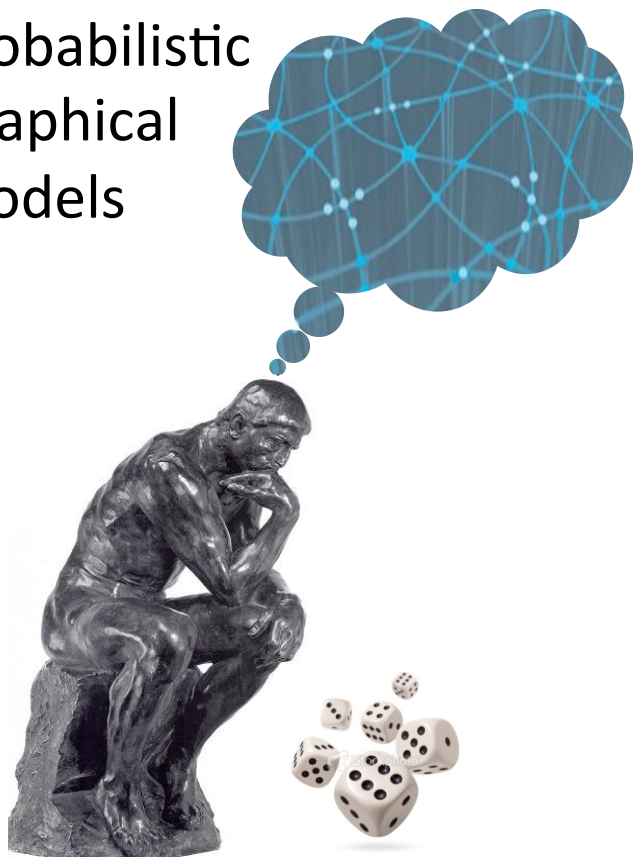
- Each sepset needs to separate graph into two conditionally independent parts



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: *separoids*
 - Related to *cliques* minimal induced width of graph

Probabilistic
Graphical
Models



Inference

Message Passing

Clique Tree and VE

Variable Elimination & Clique Trees

- Variable elimination

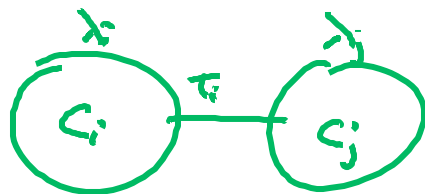
- Each step creates a factor λ_i through factor product
- A variable is eliminated in λ_i to generate new factor τ_i
- τ_i is used in computing other factors λ_j

- Clique tree view

- Intermediate factors λ_i are cliques
- τ_i are "messages" generated by clique λ_i and transmitted to another clique λ_j

Clique Tree from VE

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



Example

- $C: \tau_1(D) = \sum_C \phi_C(C) \phi_D(C, D)$

- $D: \tau_2(G, I) = \sum_D \phi_G(G, I, D) \tau_1(D)$

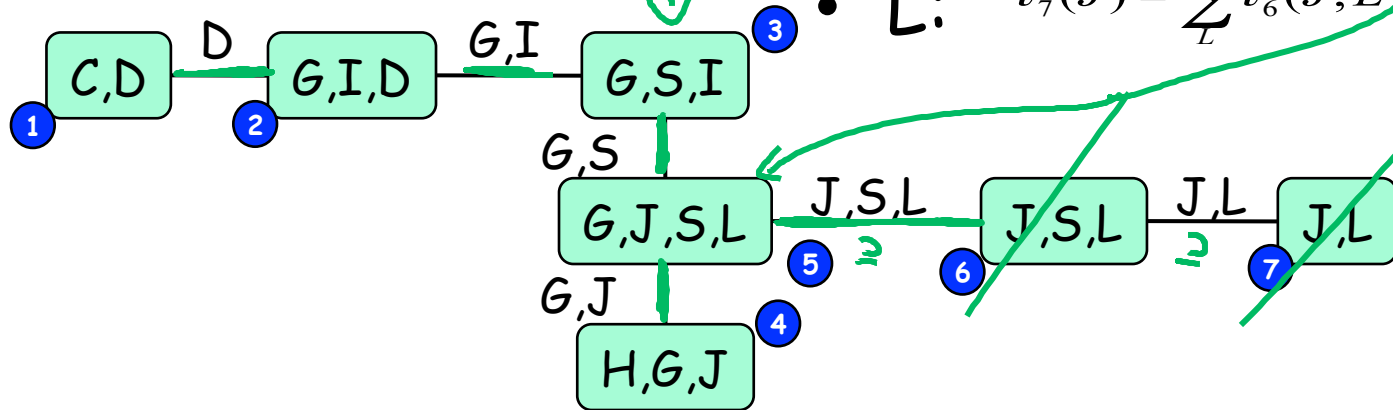
- $I: \tau_3(G, S) = \sum_I \phi_I(I) \phi_S(S, I) \tau_2(G, I)$

- $H: \tau_4(G, J) = \sum_H \phi_H(H, G, J)$

- $G: \tau_5(J, L, S) = \sum_G \phi_L(L, G) \tau_3(G, S) \tau_4(G, J)$

- $S: \tau_6(J, L) = \sum_S \phi(J, L, S) \tau_5(J, L, S)$

- $L: \tau_7(J) = \sum_L \tau_6(J, L)$



Remove redundant cliques:

those whose scope is a subset of adjacent clique's scope

Properties of Tree

- VE process induces a tree

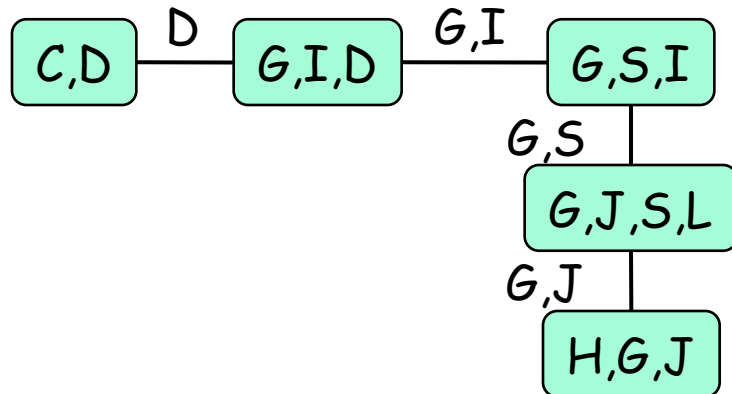
- In VE, each intermediate factor τ_i is used only once
- Hence, each cluster "passes" a factor (message) to exactly one other cluster (every cluster has at most one parent)

- Tree is family preserving: $\phi \in \mathcal{D}$

- Each of the original factors must be used in some elimination step
- And therefore contained in scope of associated ϕ_i (scope that contains $\text{Scope}(\phi)$)

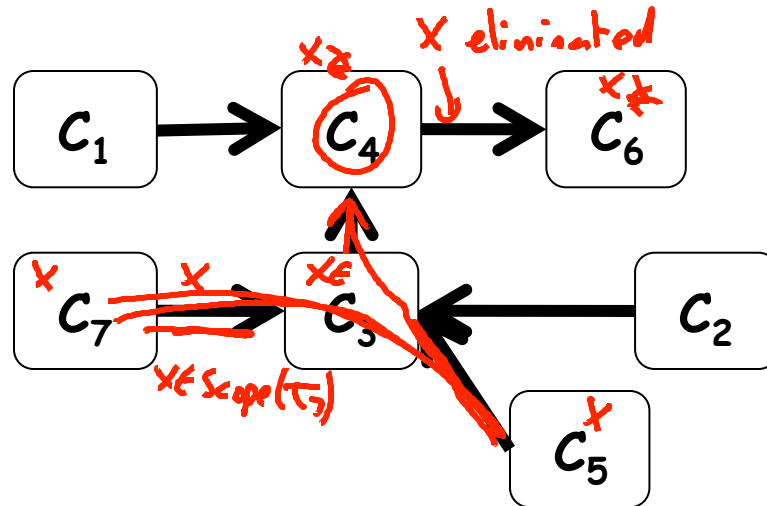
Properties of Tree

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



Running Intersection Property

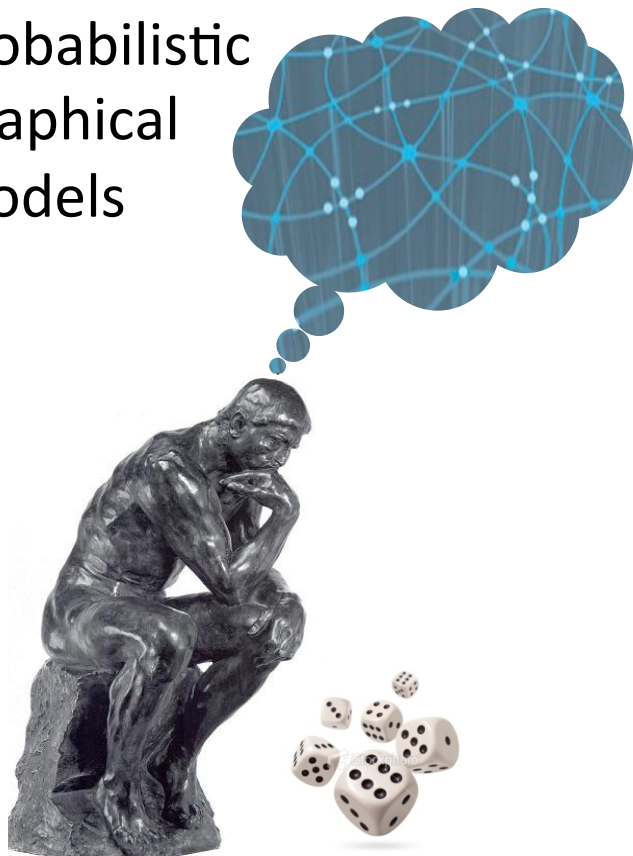
- Theorem:** If T is a tree of clusters induced by VE, then T obeys RIP



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 \sim 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

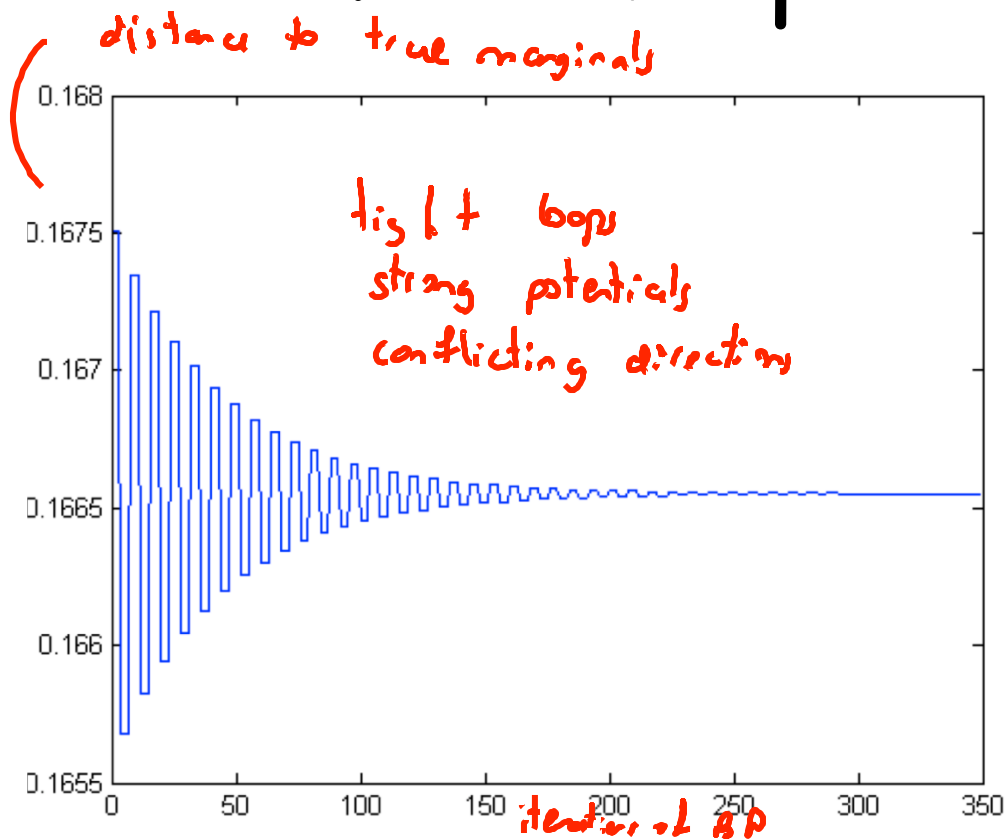


Inference

Message Passing

BP in Practice

Misconception Revisited



$\phi_4[D, A]$

d^0	a^0	100
d^0	a^1	1
d^1	a^0	1
d^1	a^1	100

$\phi_1[A, B]$

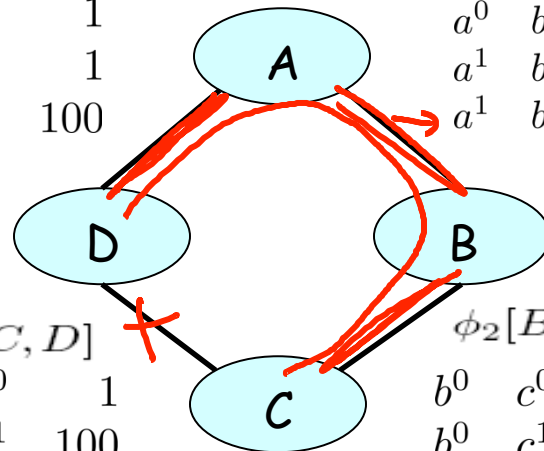
a^0	b^0	100
a^0	b^1	2
a^1	b^0	1
a^1	b^1	100

$\phi_3[C, D]$

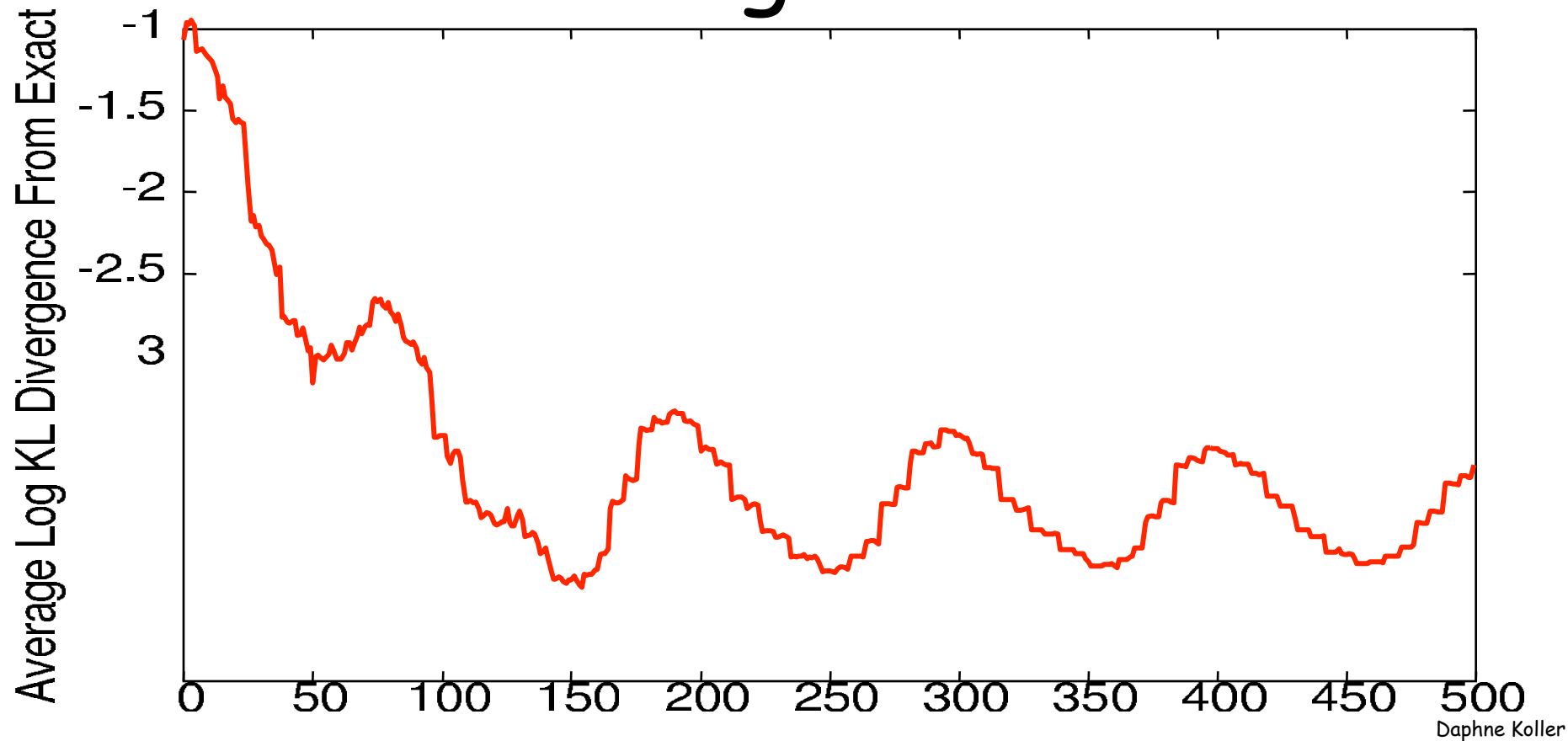
c^0	d^0	1
c^0	d^1	100
c^1	d^0	100
c^1	d^1	1

$\phi_2[B, C]$

b^0	c^0	100
b^0	c^1	1
b^1	c^0	1
b^1	c^1	100



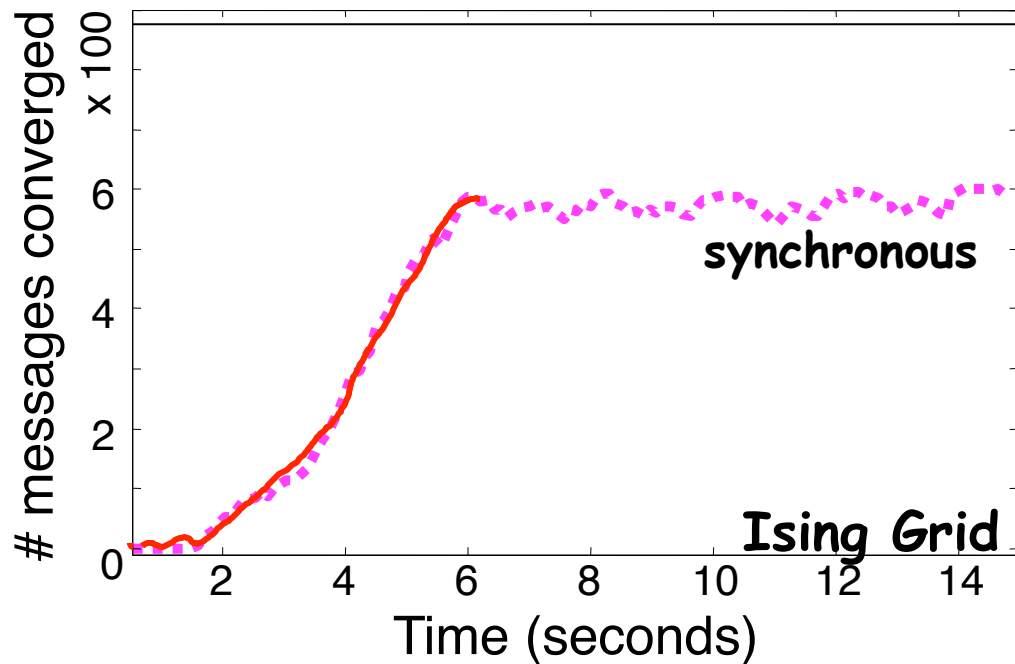
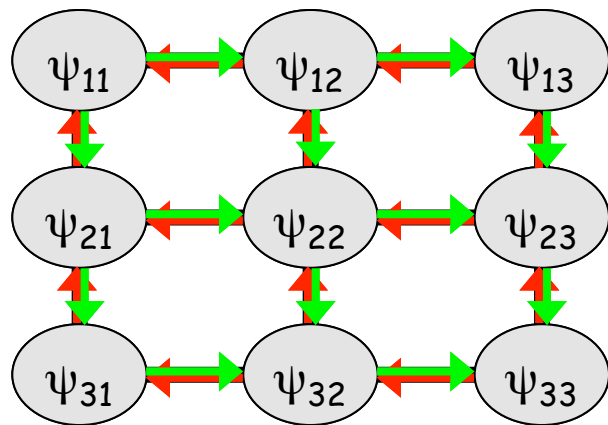
Nonconvergent BP Run



Different Variants of BP

Synchronous BP:

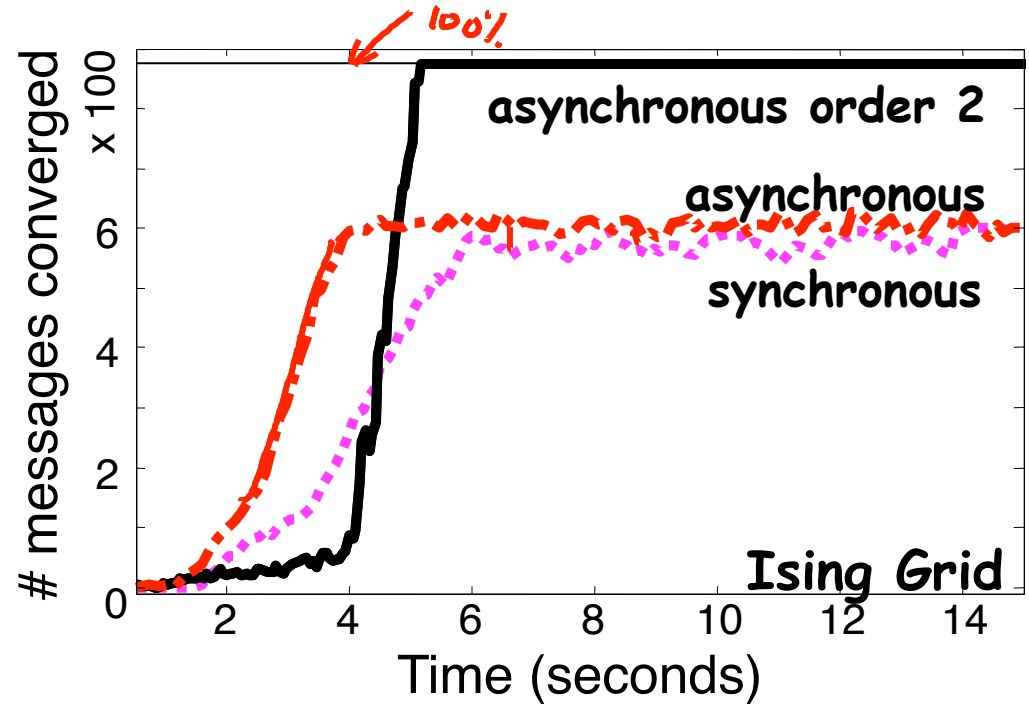
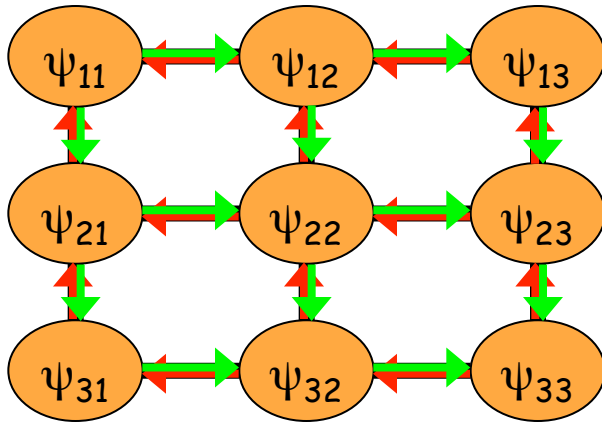
all messages are updated in parallel



Different Variants of BP

Asynchronous BP:

Messages are updated one at a time



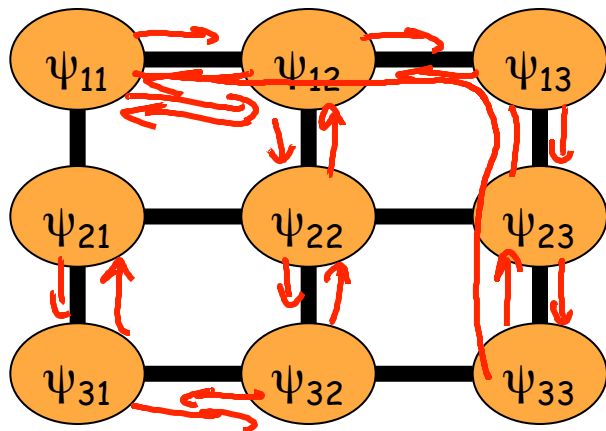
Observations

- 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

*trees must cover all edges
improves performance if trees are larger (spanning trees)*



Informed Message Scheduling

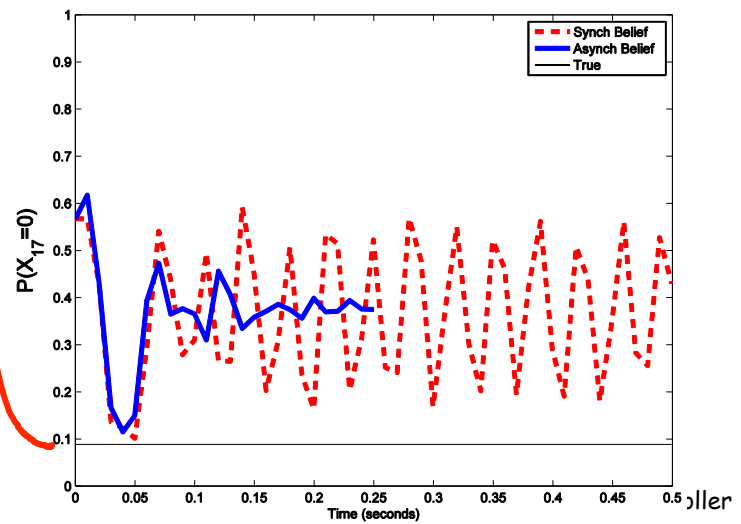
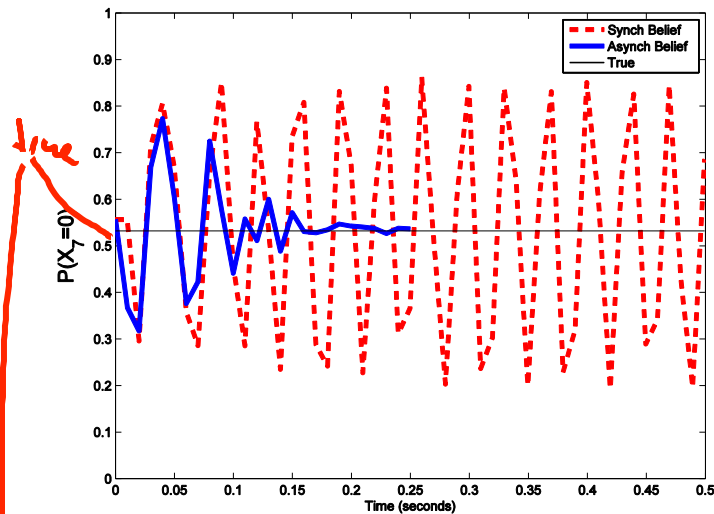
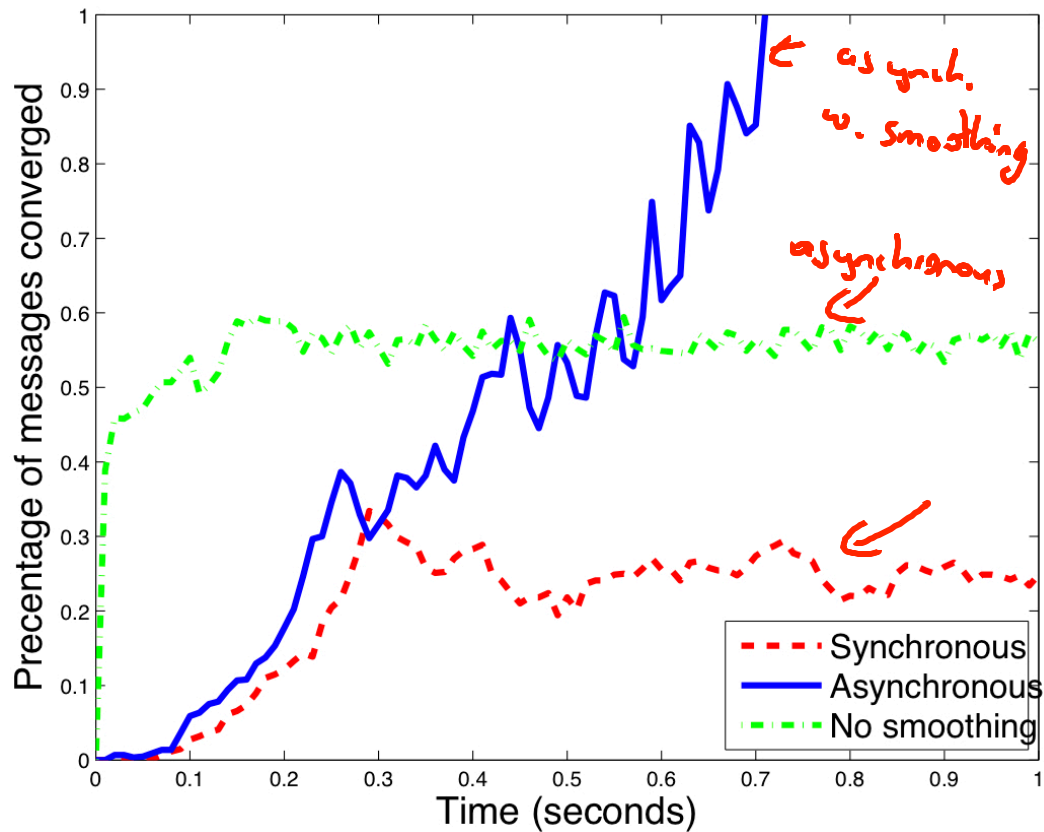
- Tree reparameterization (TRP)
 - Pick a tree and pass messages to calibrate
- Residual belief propagation (RBP)
 - Pass messages between two clusters whose beliefs over the sepset disagree the most
priority now of edges

Smoothing (Damping) Messages

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

- Dampens oscillations in messages

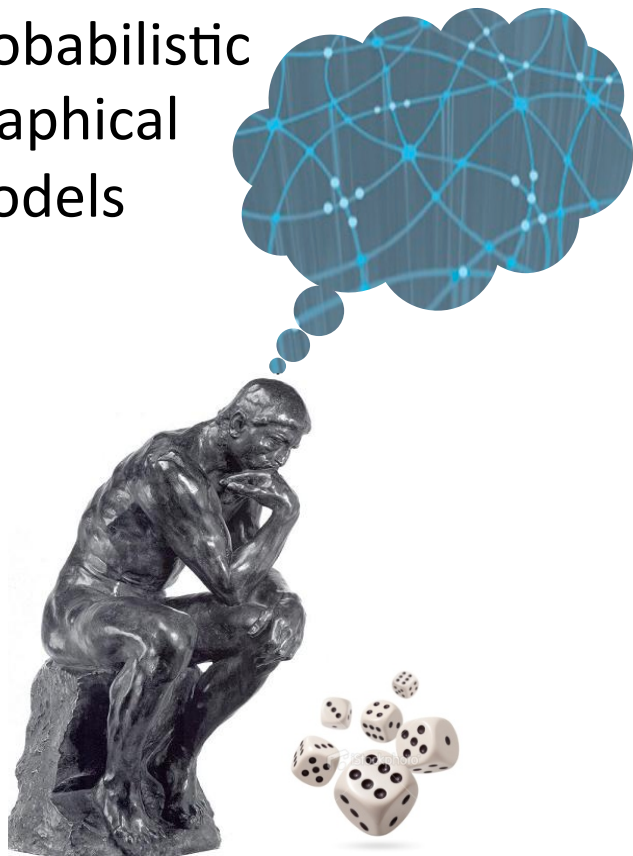
Using grid



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

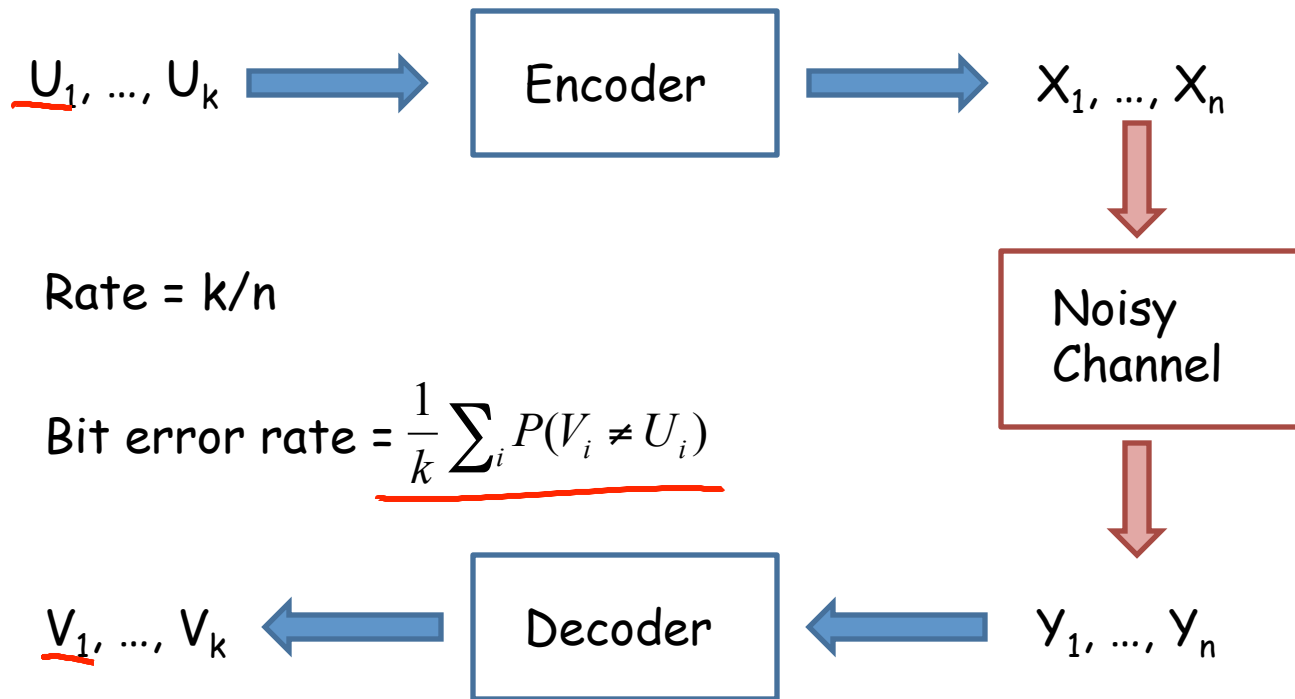


Inference

Message Passing

Loopy BP and Message Decoding

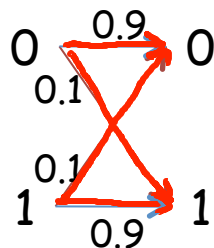
Message Coding & Decoding



Noisy
Channel

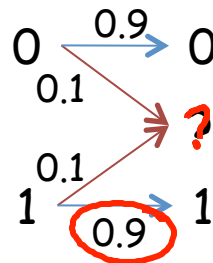
Channel Capacity

Binary
symmetric
channel

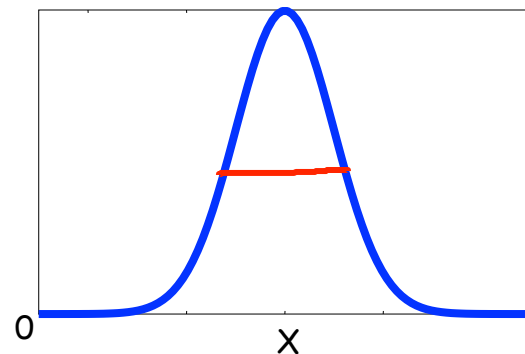


capacity = 0.531

Binary
erasure
channel

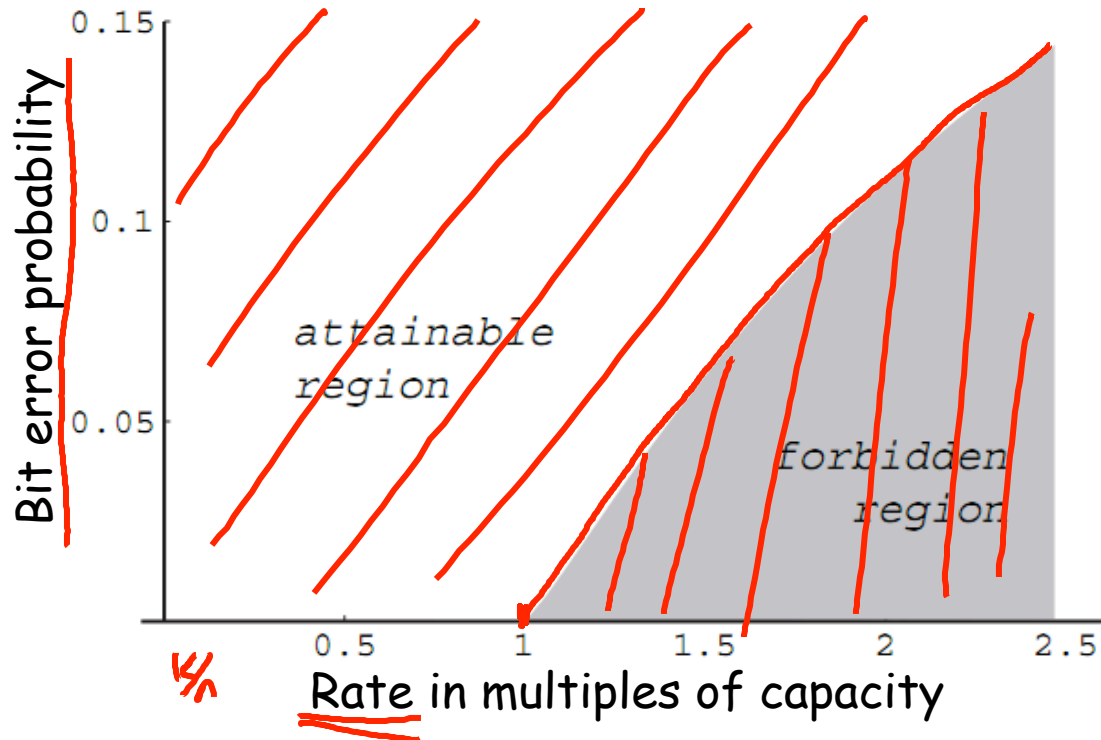


capacity = 0.9

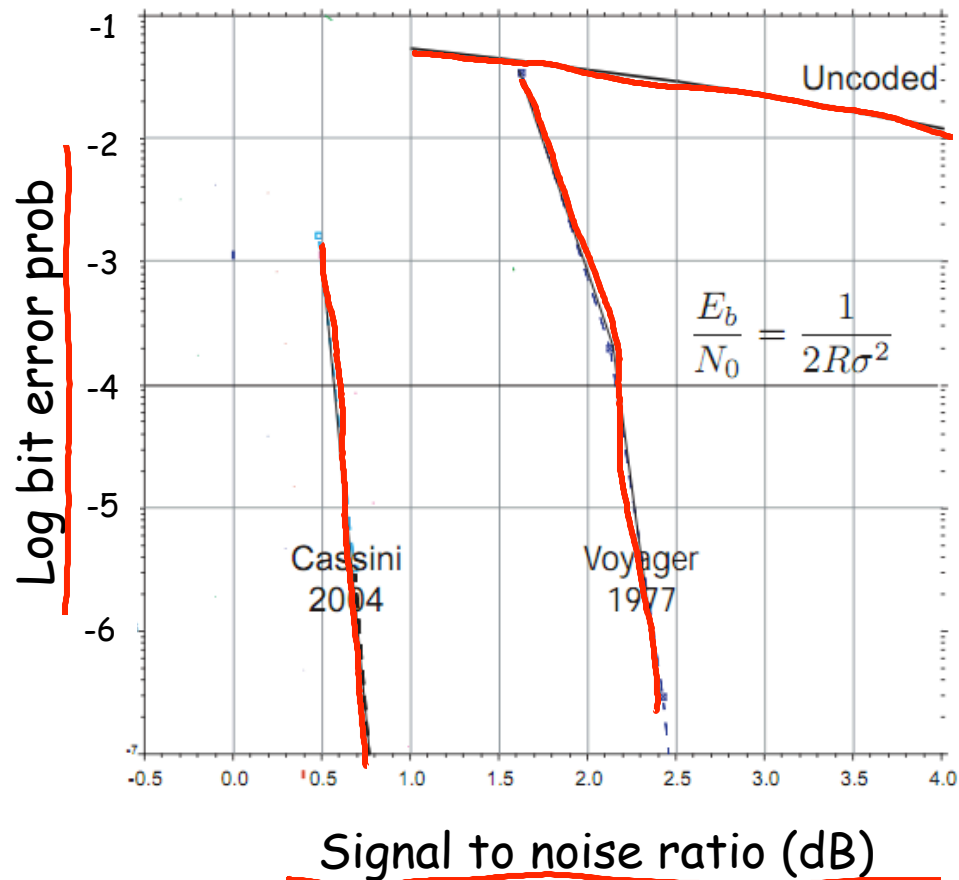


capacity = $\frac{1}{2} \log \left(1 + \frac{E(X^2)}{\sigma^2} \right)$

Shannon's Theorem



How close to C can we get?



Turbocodes (May 1993)

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

Claude Berrou, Alain Glavieux and Punya Thitimajshima

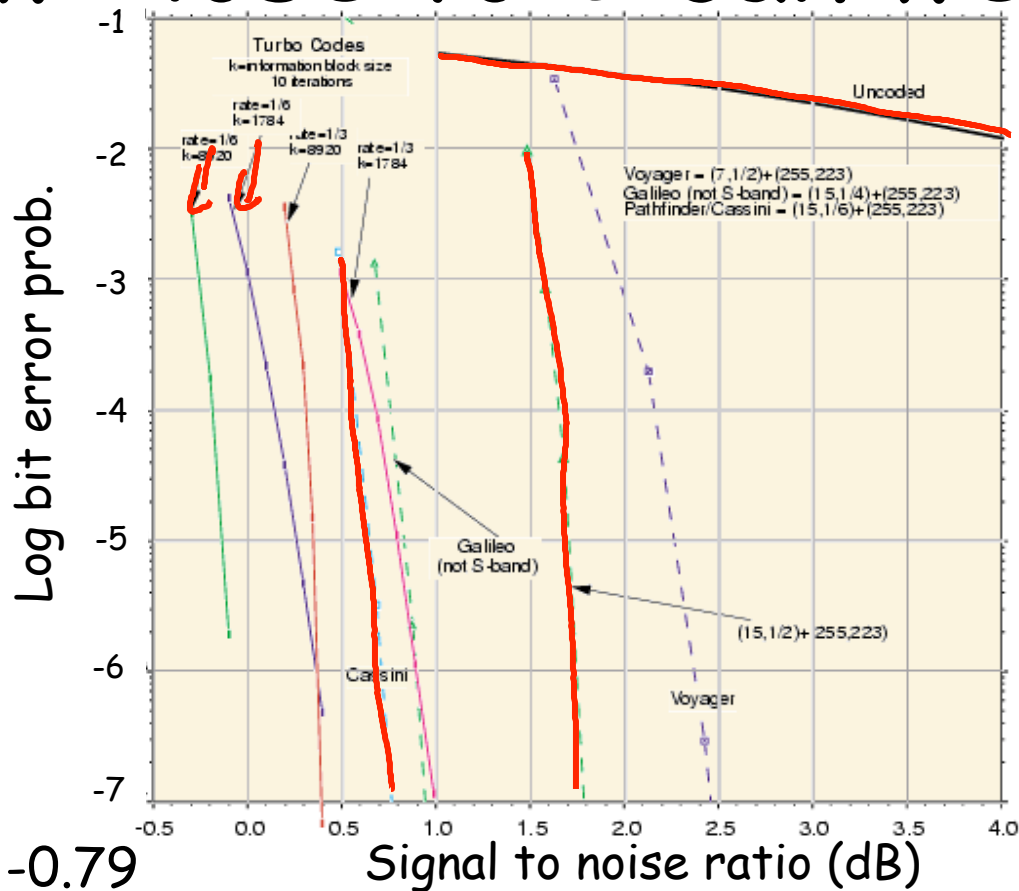
Claude Berrou, Integrated Circuits for Telecommunication Laboratory

Alain Glavieux and Punya Thitimajshima, Digital Communication Laboratory

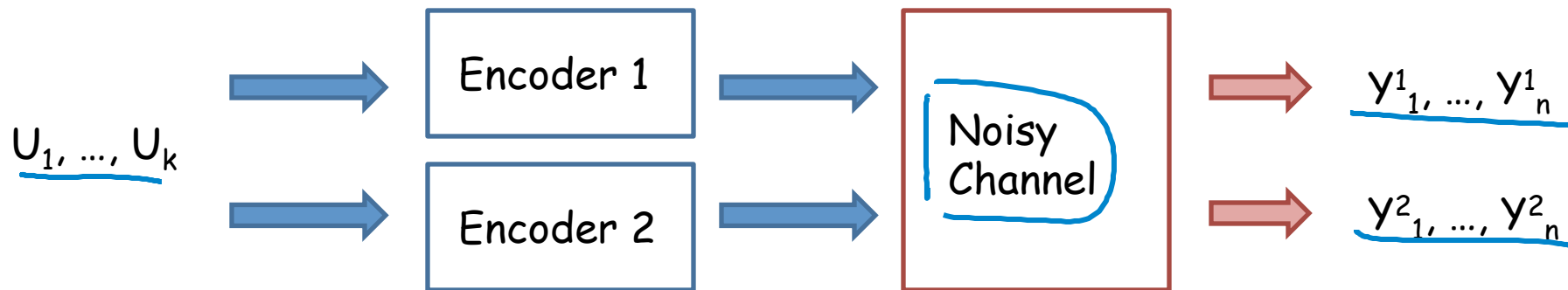
Ecole Nationale Supérieure des Télécommunications de Bretagne, France

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

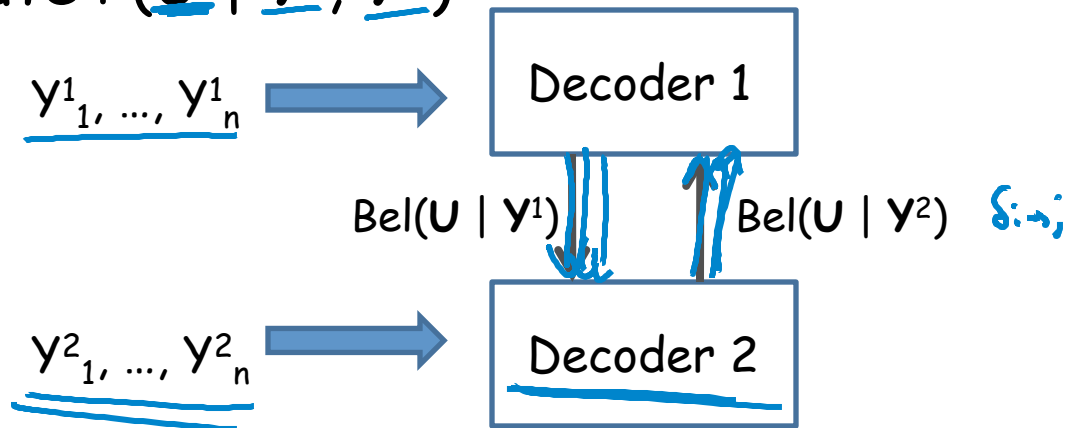
How close to C can we get?



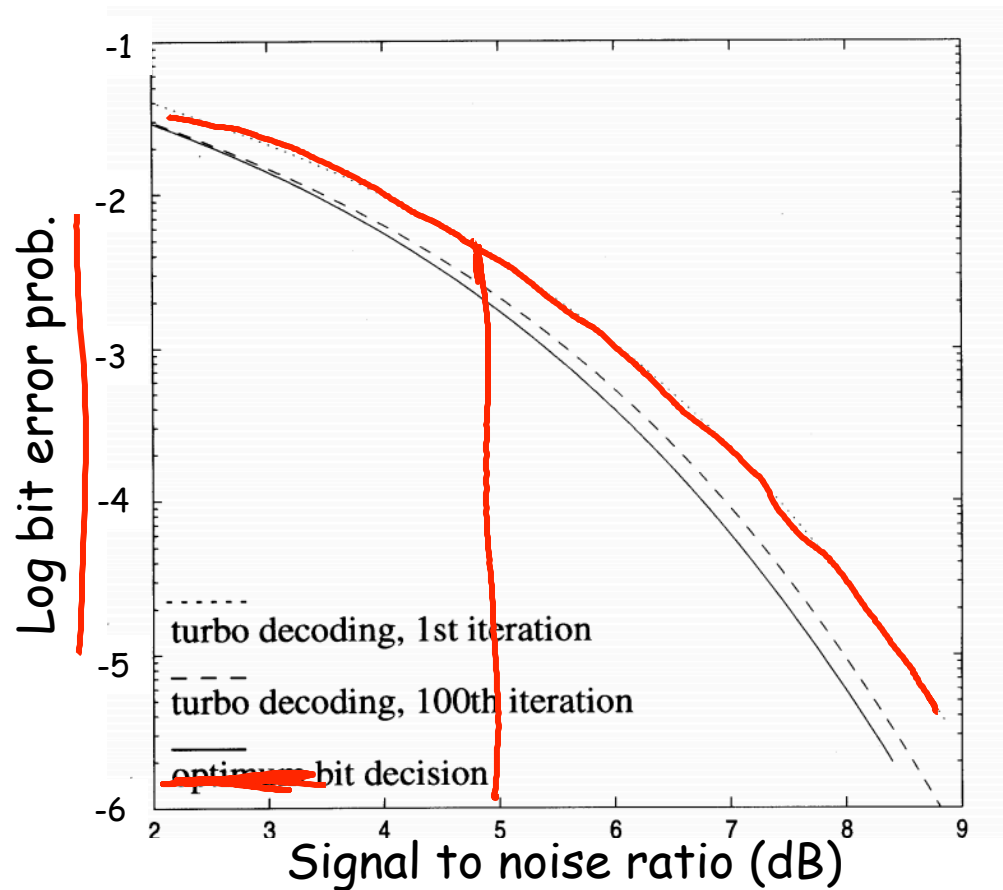
Turbocodes: The Idea



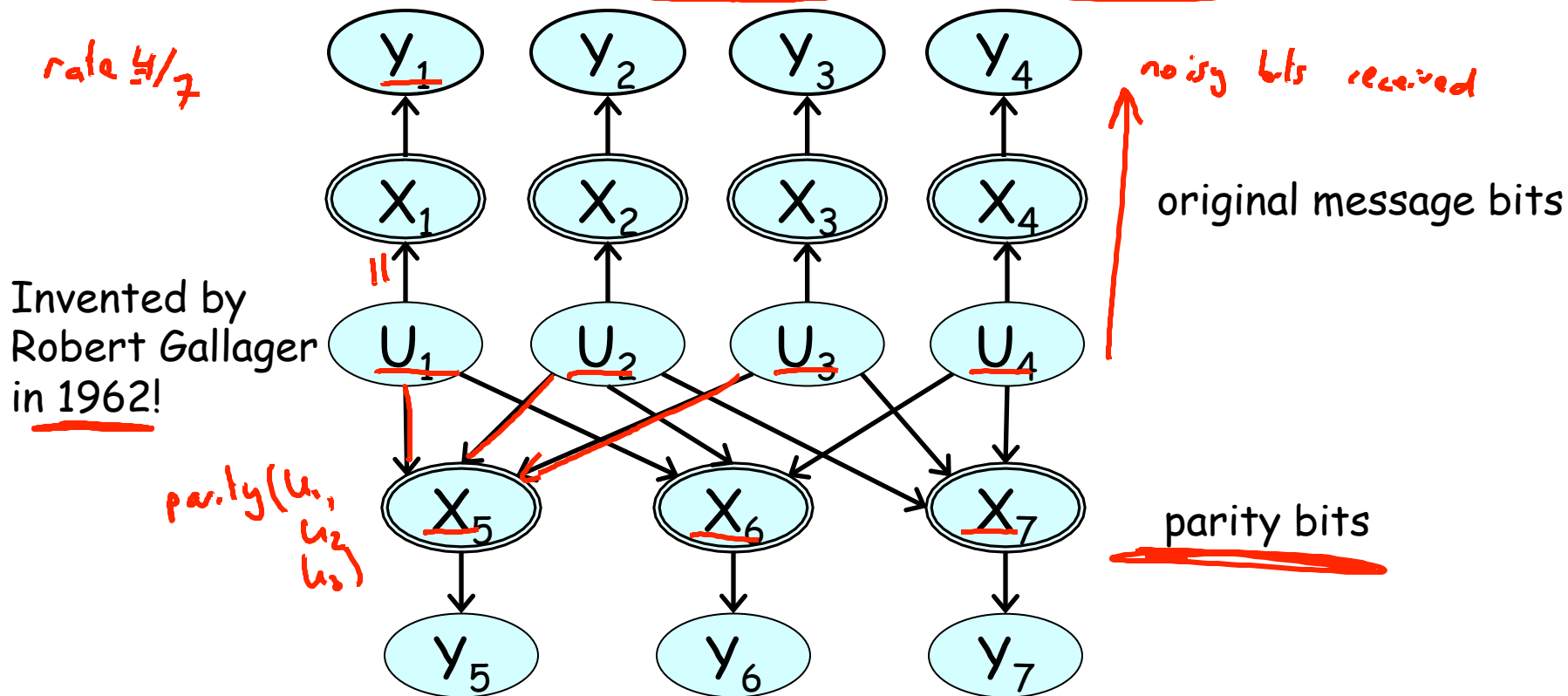
Compute $P(\underline{U} \mid \underline{y}^1, \underline{y}^2)$



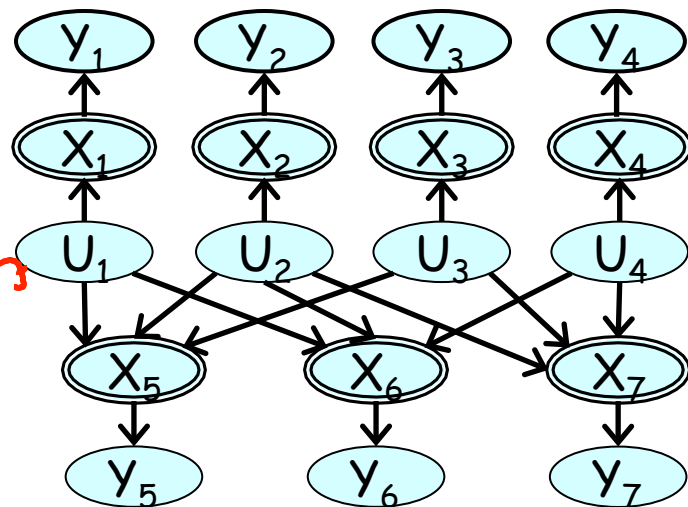
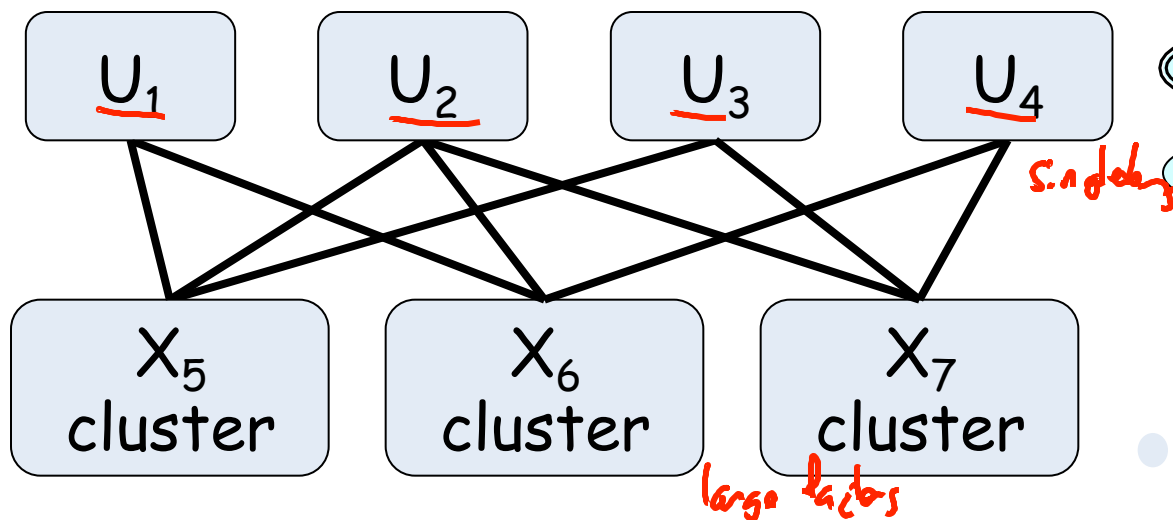
Iterations of Turbo Decoding



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