

Bounded inference non-iteratively; Mini-bucket elimination

COMPSCI 276, Spring 2017

Set 8: Rina Dechter

Reading: Class Notes (8), Darwiche chapters 14

Agenda

- Mini-bucket elimination
- Weighted Mini-bucket
- Mini-clustering
- Iterative Belief propagation
- Iterative-join-graph propagation

Probabilistic Inference Tasks

- Belief updating:

$$\text{BEL}(X_i) = P(X_i = x_i \mid \text{evidence})$$

- Finding most probable explanation (MPE)

$$\bar{x}^* = \underset{\bar{x}}{\operatorname{argmax}} P(\bar{x}, e)$$

- Finding maximum a-posteriori hypothesis

$$(a_1^*, \dots, a_k^*) = \underset{\bar{a}}{\operatorname{argmax}} \sum_{X/A} P(\bar{x}, e) \quad \begin{matrix} A \subseteq X : \\ \text{hypothesis variables} \end{matrix}$$

- Finding maximum-expected-utility (MEU) decision

$$(d_1^*, \dots, d_k^*) = \underset{\bar{d}}{\operatorname{argmax}} \sum_{X/D} P(\bar{x}, e) U(\bar{x}) \quad \begin{matrix} D \subseteq X : \text{decision variables} \\ U(\bar{x}) : \text{utility function} \end{matrix}$$

Queries

- Probability of evidence (or partition function)

$$P(e) = \sum_{X-\text{var}(e)} \prod_{i=1}^n P(x_i | pa_i)|_e \quad Z = \sum_X \prod_i \psi_i(C_i)$$

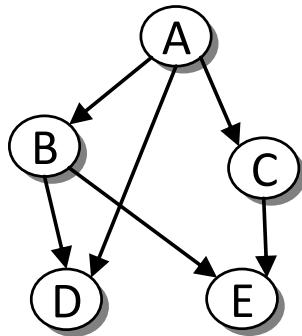
- Posterior marginal (beliefs):

$$P(x_i | e) = \frac{P(x_i, e)}{P(e)} = \frac{\sum_{X-\text{var}(e)-X_i} \prod_{j=1}^n P(x_j | pa_j)|_e}{\sum_{X-\text{var}(e)} \prod_{j=1}^n P(x_j | pa_j)|_e}$$

- Most Probable Explanation

$$\bar{x}^* = \operatorname{argmax}_{\bar{x}} P(\bar{x}, e)$$

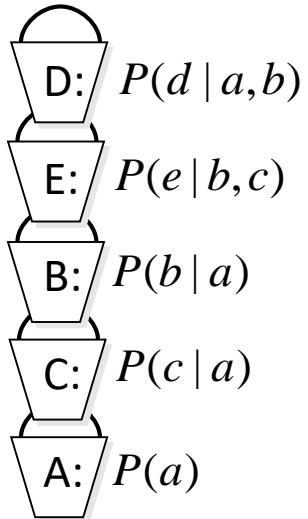
Bucket Elimination



Query: $P(a | e = 0) \propto P(a, e = 0)$ Elimination Order: d,e,b,c

$$\begin{aligned} P(a, e = 0) &= \sum_{c,b,e=0,d} P(a)P(b | a)P(c | a)P(d | a,b)P(e | b,c) \\ &= P(a) \sum_c P(c | a) \sum_b P(b | a) \sum_{e=0} P(e | b,c) \sum_d P(d | a,b) \end{aligned}$$

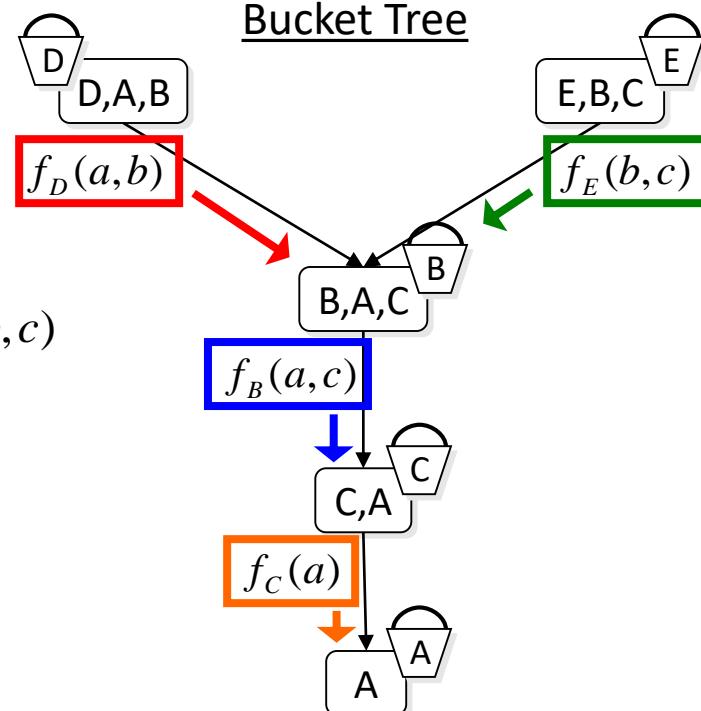
Original Functions



Messages

$$\begin{aligned} f_D(a, b) &= \sum_d P(d | a, b) \\ f_E(b, c) &= P(e = 0 | b, c) \\ f_B(a, c) &= \sum_b P(b | a) f_D(a, b) f_E(b, c) \\ f_C(a) &= \sum_c P(c | a) f_B(a, c) \\ P(a, e = 0) &= p(A) f_C(a) \end{aligned}$$

Bucket Tree



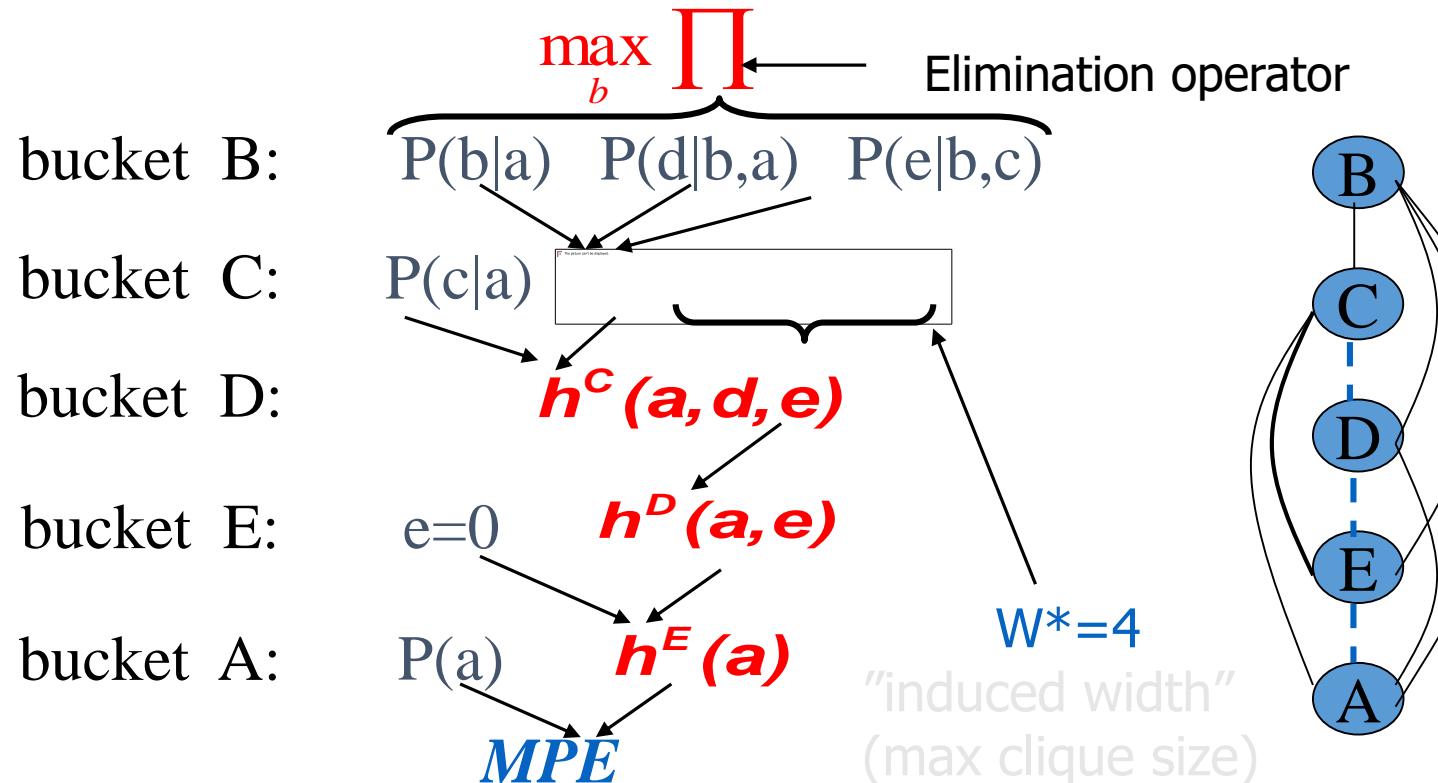
Finding

$$MPE = \max_{\bar{x}} P(\bar{x})$$

Algorithm *elim-mpe* (Dechter 1996)

\sum is replaced by **max** :

$$MPE = \max_{a,e,d,c,b} P(a)P(c|a)P(b|a)P(d|b,a)P(e|b,c)$$



Generating the MPE-tuple

$$5. \ b' = \arg \max_{\substack{b \\ a'}} P(b | a') \times P(d' | b, a') \times P(e' | b, c')$$

$$4. \ c' = \arg \max_{\substack{c \\ a'}} P(c | a') \times h^B(a', d', c, e')$$

$$3. \ d' = \arg \max_d h^C(a', d, e')$$

$$2. \ e' = 0$$

$$1. \ a' = \arg \max_a P(a) \cdot h^E(a)$$

$$B: \quad P(b|a) \quad P(d|b,a) \quad P(e|b,c)$$

$$C: \quad P(c|a) \quad h^B(a, d, c, e)$$

$$D: \quad h^C(a, d, e)$$

$$E: \quad e=0 \quad h^D(a, e)$$

$$A: \quad P(a) \quad h^E(a)$$

Return (a', b', c', d', e')

Approximate Inference

- Metrics of evaluation
- **Absolute error:** given $e > 0$ and a query $p = P(x|e)$, an estimate r has absolute error e iff $|p-r| < e$
- **Relative error:** the ratio r/p in $[1-e, 1+e]$.
- Dagum and Luby 1993: approximation up to a relative error is NP-hard.
- Absolute error is also NP-hard if error is less than .5

Mini-Buckets: “Local Inference”

- Computation in a bucket is time and space exponential in the number of variables involved
- Therefore, partition functions in a bucket into “mini-buckets” on smaller number of variables

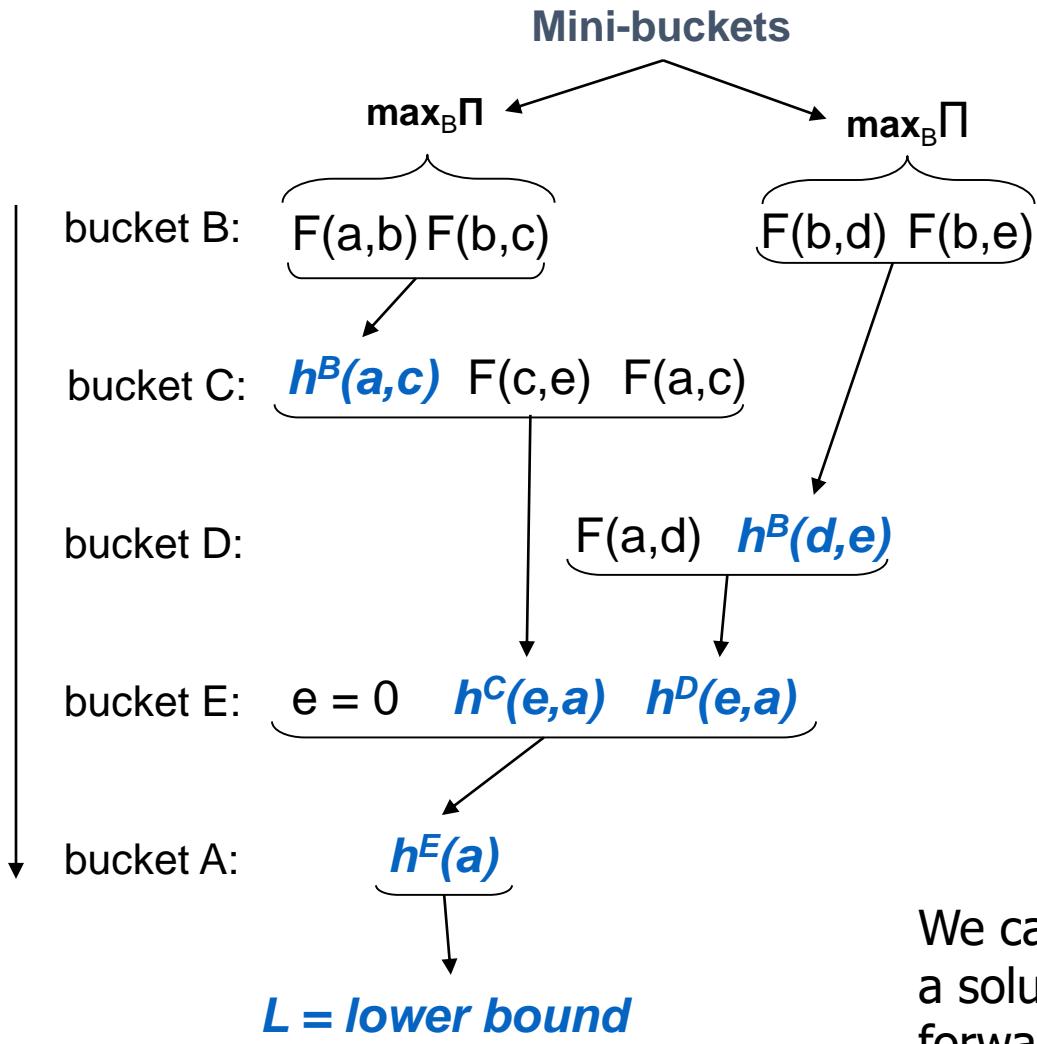
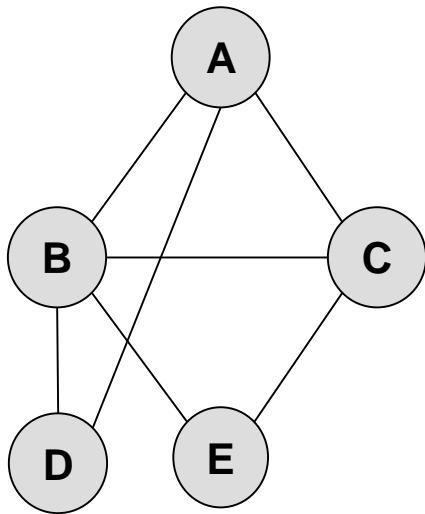
Mini-Bucket Approximation: MPE task

Split a bucket into mini-buckets => bound complexity

$$\begin{aligned} \textbf{bucket}(\mathbf{X}) &= \underbrace{\{ \mathbf{h}_1, \dots, \mathbf{h}_r, \mathbf{h}_{r+1}, \dots, \mathbf{h}_n \}}_{h^X = \max_X \prod_{i=1}^n h_i} \\ &\quad \swarrow \qquad \searrow \\ \underbrace{\{ \mathbf{h}_1, \dots, \mathbf{h}_r \}}_{g^X = (\max_X \prod_{i=1}^r h_i) \cdot (\max_X \prod_{i=r+1}^n h_i)} &\qquad \underbrace{\{ \mathbf{h}_{r+1}, \dots, \mathbf{h}_n \}} \\ &\quad \downarrow \\ \mathbf{h}^X &\leq \mathbf{g}^X \end{aligned}$$

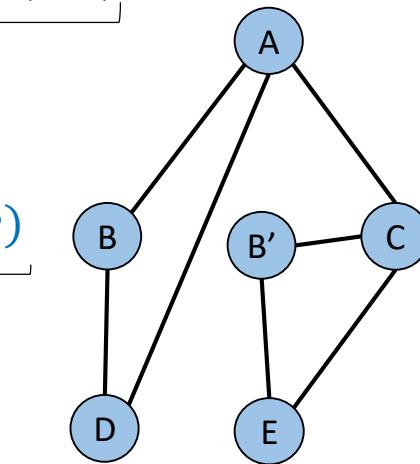
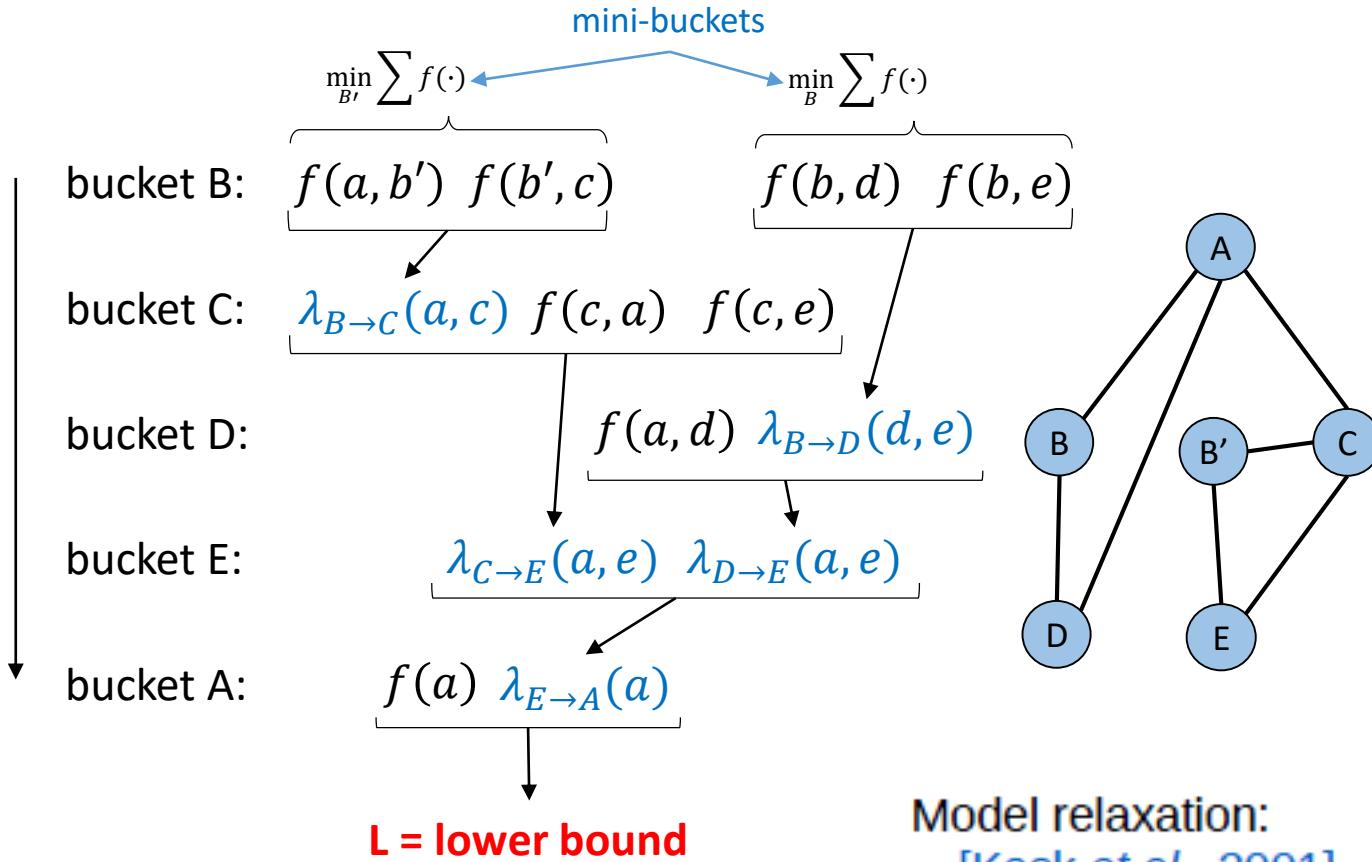
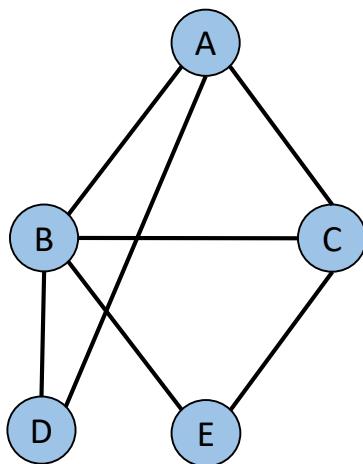
Exponential complexity decrease : $O(e^n) \rightarrow O(e^r) + O(e^{n-r})$

Mini-Bucket Elimination



We can generate
a solution s going
forward as before
 $U = F(s)$

Mini-Bucket Elimination

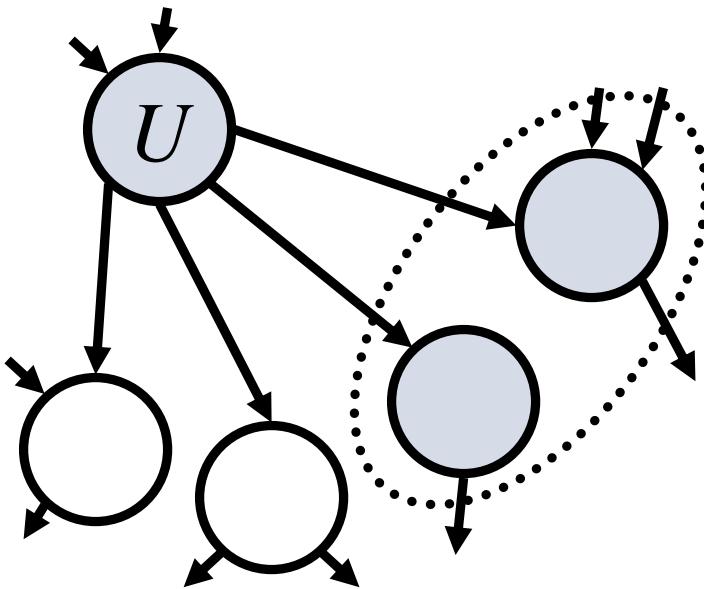


Model relaxation:
[Kask et al., 2001]
[Geffner et al., 2007]
[Choi et al., 2007]
[Johnson et al. 2007]

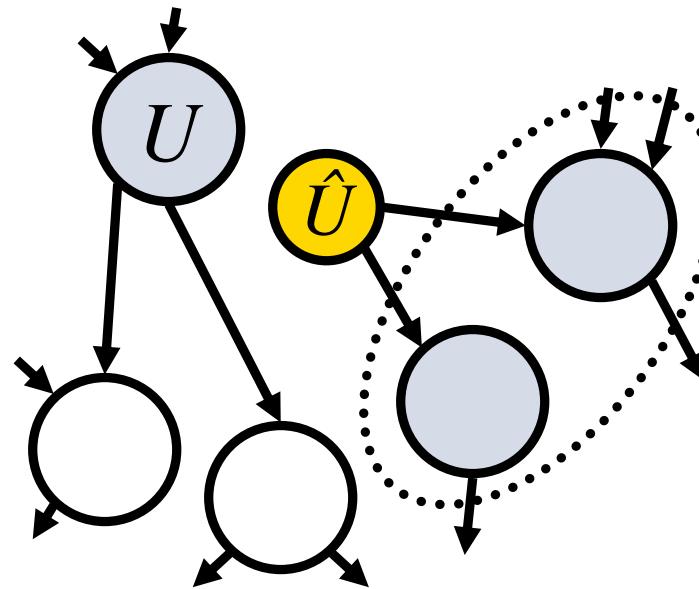
Semantics of Mini-Bucket: Splitting a Node

Variables in different buckets are renamed and duplicated
(Kask *et. al.*, 2001), (Geffner *et. al.*, 2007), (Choi, Chavira, Darwiche , 2007)

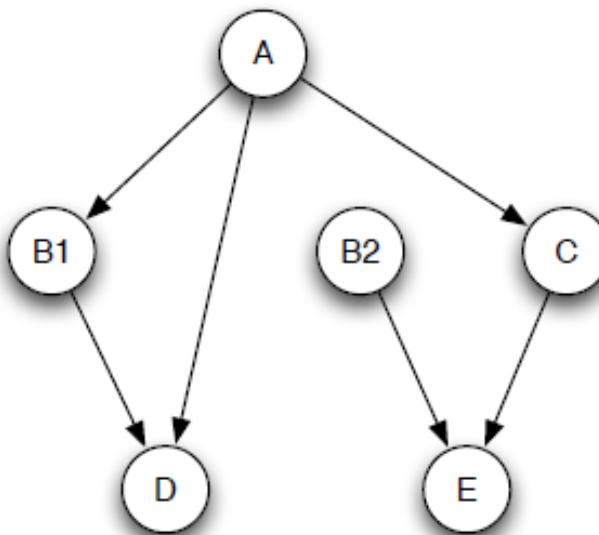
Before Splitting:
Network N



After Splitting:
Network N'



Relaxed Network Example



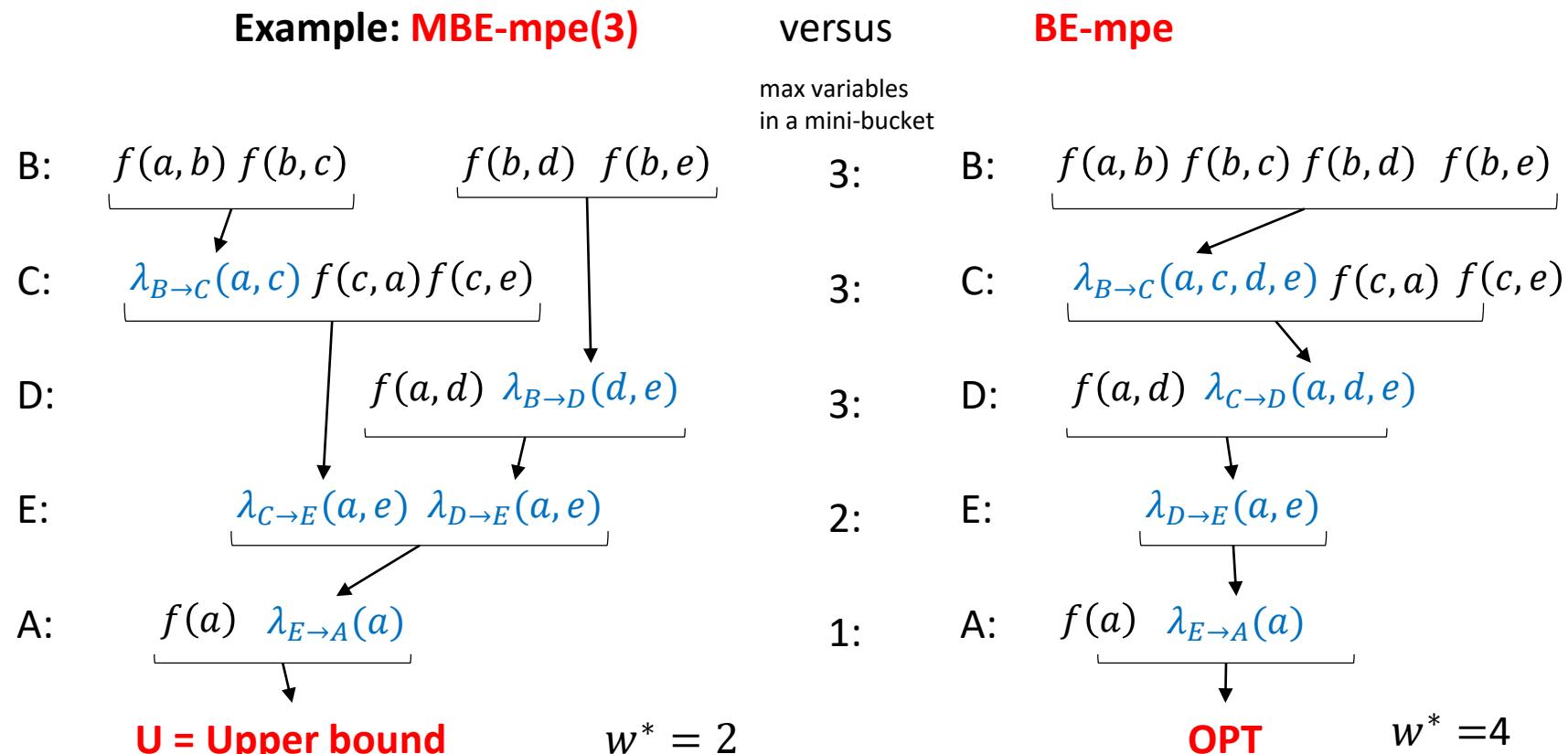
(a)

B1: $P(b1|a), P(d|b1,a)$
B2: $P(e|b2,c)$
C: $P(c|a)$
D:
E: $E=e$
A: $P(a)$

(b)

MBE-MPE(i): Algorithm MBE-mpe

- **Input:** l – max number of variables allowed in a mini-bucket
- **Output:** [lower bound (P of suboptimal solution), upper bound]



[Dechter and Rish, 1997]

Mini-Bucket Decoding

$$\hat{b} = \arg \min_b f(\hat{a}, b) + f(b, \hat{c}) + f(b, \hat{d}) + f(b, \hat{e})$$

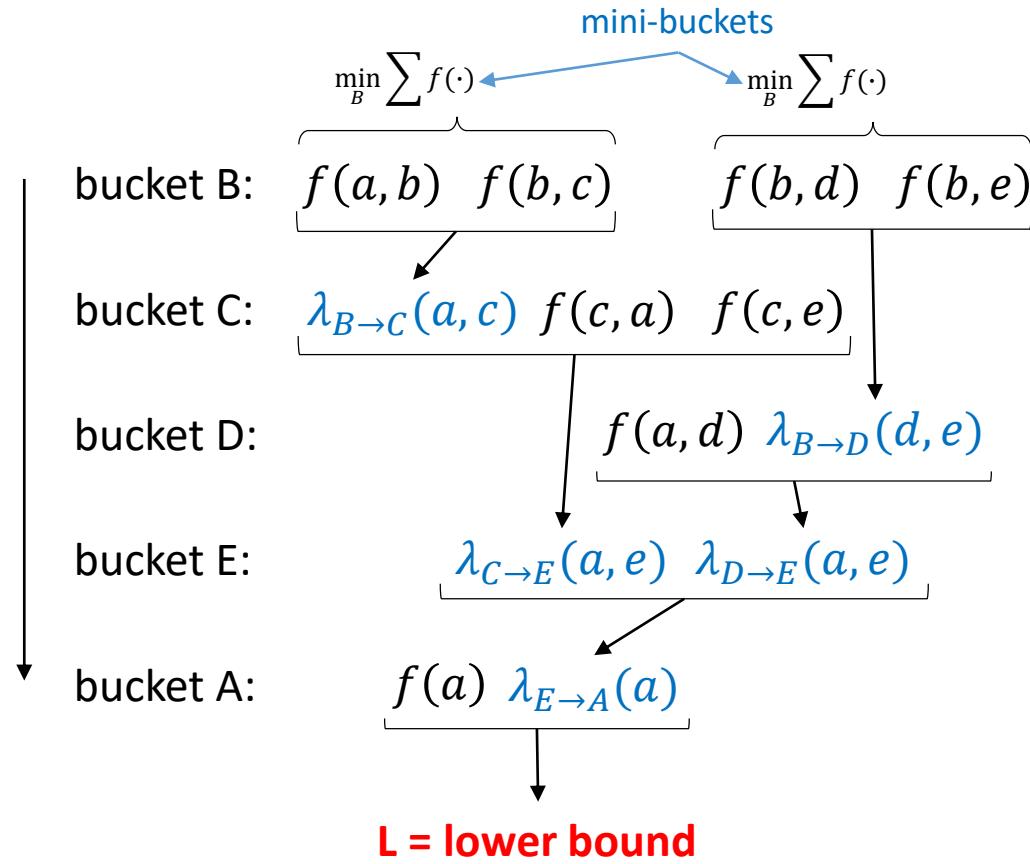
$$\hat{c} = \arg \min_c \lambda_{B \rightarrow C}(\hat{a}, c) + f(c, \hat{a}) + f(c, \hat{e})$$

$$\hat{d} = \arg \min_d f(\hat{a}, d) + \lambda_{B \rightarrow D}(d, \hat{e})$$

$$\hat{e} = \arg \min_e \lambda_{C \rightarrow E}(\hat{a}, e) + \lambda_{D \rightarrow E}(\hat{a}, e)$$

$$\hat{a} = \arg \min_a f(a) + \lambda_{E \rightarrow A}(a)$$

Greedy configuration = upper bound



[Dechter and Rish, 2003]

(i,m)-Patitionings

Definition 7.1.1 ((i,m)-partitioning) Let H be a collection of functions h_1, \dots, h_t defined on scopes S_1, \dots, S_t , respectively. We say that a function f is subsumed by a function h if any argument of f is also an argument of h . A partitioning of h_1, \dots, h_t is canonical if any function f subsumed by another function is placed into the bucket of one of those subsuming functions. A partitioning Q into mini-buckets is an (i, m) -partitioning if and only if (1) it is canonical, (2) at most m non-subsumed functions are included in each mini-bucket, (3) the total number of variables in a mini-bucket does not exceed i , and (4) the partitioning is refinement-maximal, namely, there is no other (i, m) -partitioning that it refines.

$\text{MBE}(i,m)$, $\text{MBE}(i)$

- Input: Belief network (P_1, \dots, P_n)
- Output: upper and lower bounds
- Initialize: put functions in buckets along ordering
- Process each bucket from $p=n$ to 1
 - Create (i,m) -partitions
 - Process each mini-bucket
- (For mpe): assign values in ordering d
- Return: mpe-configuration, upper and lower bounds

Algorithm MBE-mpe(i,m)

Input: A belief network $\mathcal{B} = \langle X, D, G, \mathcal{P} \rangle$, where $\mathcal{P} = \{P_1, \dots, P_n\}$; an ordering of the variables, $d = X_1, \dots, X_n$; observations e .

Output: An upper bound U and a lower bound L on the most probable configuration given the evidence. A suboptimal solution \bar{x}^a that provides the lower bound $L = P(\bar{x}^a)$.

1. **Initialize:** Generate an ordered partition of the conditional probability function, $bucket_1, \dots, bucket_n$, where $bucket_i$ contains all functions whose highest variable is X_i . Put each observed variable in its bucket. Let ψ_i be the product of input function in a bucket and let h_i be the messages in the bucket.
2. **Backward:** For $p \leftarrow n$ downto 1, do
for all the functions h_1, h_2, \dots, h_j in $bucket_p$, do

- If (observed variable) $bucket_p$ contains $X_p = x_p$, assign $X_p = x_p$ to each function and put each in appropriate bucket.
- else, Generate an an (i, m) -partitioning, $Q' = \{Q_1, \dots, Q_r\}$ of h_1, h_2, \dots, h_t in $bucket_p$.
- for each $Q_l \in Q'$ containing h_{l_1}, \dots, h_{l_t} , do

$$h_l \leftarrow \max_{X_p} \prod_{j=1}^t h_{l_j} \quad (1.1)$$

Add h_l to the bucket of the largest-index in $scope(h_l)$. Put constants in $bucket_1$.

3. **Forward:**

- Generate an mpe upper bound cost by maximizing over X_1 , the product in $bucket_1$. Namely $U \leftarrow \max_{X_1} \psi_1 \prod_j h_{1_j}$.
- (Generate an approximate mpe tuple): Given x_1^a, \dots, x_{p-1}^a , assign x_p^a to X_p that maximizes the product of all functions in $bucket_p$. $L \leftarrow P(x_1^a, \dots, x_n^a)$

4. **Return** U and L and configuration: $\bar{x}^a = (x_1^a, \dots, x_n^a)$

Figure 1.2: Algorithm $MBE\text{-}mpe(i,m)$.

Partitioning, Refinements

Clearly, as the mini-buckets get smaller, both complexity and accuracy decrease.

Definition 7.1.4 *Given two partitionings Q' and Q'' over the same set of elements, Q' is a refinement of Q'' if and only if for every set $A \in Q'$ there exists a set $B \in Q''$ such that $A \subseteq B$.*

It is easy to see that:

Proposition 7.1.5 *If Q'' is a refinement of Q' in bucket_p, then $h^p \leq g_{Q'}^p \leq g_{Q''}^p$.*

Remember that *mbe-mpe* computes the bounds on $MPE = \max_{\bar{x}} P(\bar{x}, \bar{e})$, rather than on $M = \max_{\bar{x}} P(\bar{x}|\bar{e}) = MPE/P(\bar{e})$. Thus

$$\frac{L}{P(\bar{e})} \leq M \leq \frac{U}{P(\bar{e})}$$

Properties of MBE-mpe(i)

- **Complexity:** $O(r \exp(i))$ time and $O(r \exp(i))$ space.
- **Accuracy:** determined by upper/lower (U/L) bound.
- As i increases, both accuracy and complexity increase.
- Possible use of mini-bucket approximations:
 - As [anytime algorithms](#)
 - As [heuristics](#) in best-first search

Anytime Approximation

Algorithm anytime-mpe(ϵ)

Input: Initial values of i and m , i_0 and m_0 ; increments i_{step} and m_{step} , and desired approximation error ϵ .

Output: U and L

1. **Initialize:** $i = i_0, m = m_0$.
2. **do**
3. run $mbe\text{-}mpe(i, m)$
4. $U \leftarrow$ upper bound of $mbe\text{-}mpe(i, m)$
5. $L \leftarrow$ lower bound of $mbe\text{-}mpe(i, m)$
6. Retain best bounds U, L , and best solution found so far
7. **if** $1 \leq U/L \leq 1 + \epsilon$, return solution
8. **else** increase i and m : $i \leftarrow i + i_{step}$ and $m \leftarrow m + m_{step}$
9. **while** computational resources are available
10. **Return** the largest L
and the smallest U found so far.

MBE for Belief Updating and for Probability of Evidence or Partition Function

- Idea mini-bucket is the same:

$$\sum_X f(x) \bullet g(x) \leq \sum_X f(x) \bullet \sum_X g(x)$$

$$\sum_X f(x) \bullet g(x) \leq \sum_X f(x) \bullet \max_X g(X)$$

- So we can apply a sum in each mini-bucket, or better, one sum and the rest max, or min (for lower-bound)
- **MBE-bel-max(i,m), MBE-bel-min(i,m)** generating upper and lower-bound on beliefs approximates BE-bel
- MBE-map(i,m): max buckets will be maximized, sum buckets will be sum-max. Approximates BE-map.

Algorithm MBE-bel-max(i,m)

Input: A belief network $\mathcal{B} = \langle \mathbf{X}, \mathbf{D}, \mathbf{P}_G, \Pi \rangle$, an ordering $d = (X_1, \dots, X_n)$; evidence e

Output: an upper bound on $P(X_1, \bar{e})$ and an upper bound on $P(e)$.

1. Initialize: Partition $P = \{P_1, \dots, P_n\}$ into buckets $bucket_1, \dots, bucket_n$, where $bucket_k$ contains all CPTs h_1, h_2, \dots, h_t whose highest-index variable is X_k .

2. Backward: for $k = n$ to 2 do

- If X_p is observed ($X_k = a$), assign $X_k \leftarrow a$ in each h_j and put the result in the highest-variable bucket of its scope (put constants in $bucket_1$).
- Else for h_1, h_2, \dots, h_t in $bucket_k$ Generate an (i, m) -partitioning, $Q' = \{Q_1, \dots, Q_r\}$. For each $Q_l \in Q'$, containing h_{l_1}, \dots, h_{l_t} , do

$$h_l \leftarrow \sum_{X_k} \Pi_{j=1}^t h_{l_j}, \text{ if } l = 1$$

$$h_l \leftarrow \max_{X_k} \Pi_{j=1}^t h_{l_j}, \text{ if } k \neq 1$$

Add h_l to the bucket of the highest-index variable in its scope $\bigcup_{j=1}^t scope(h_{l_j}) - \{X_k\}$. (put constant functions in $bucket_1$).

3. Return $P'(\bar{x}_1, e) \leftarrow$ the product of functions in the bucket of X_1 , which is an upper bound on $P(x_1, \bar{e})$.

$P'(e) \leftarrow \sum_{x_1} P'(\bar{x}_1, e)$, which is an upper bound on probability of evidence.

Figure 8.5: Algorithm MBE-bel-max(i, m).

Normalization

- MBE-bel computes upper/lower bound on the joint marginal distributions.

Alternatively, let U_i and L_i be the upper bound and lower bounding functions on $P(X_1 = x_i, \bar{e})$ obtained by *mbe-bel-max* and *mbe-bel-min*, respectively. Then,

$$\frac{L_i}{P(\bar{e})} \leq P(x_i|\bar{e}) \leq \frac{U_i}{P(\bar{e})}$$

We sometime use normalization of the approximation, but then no guarantee. The problem is that we have to approximate also $P(e)$.

Empirical Evaluation

(Dechter and Rish, 1997; Rish thesis, 1999)

- Randomly generated networks
 - Uniform random probabilities
 - Random noisy-OR
- CPCS networks
- Probabilistic decoding

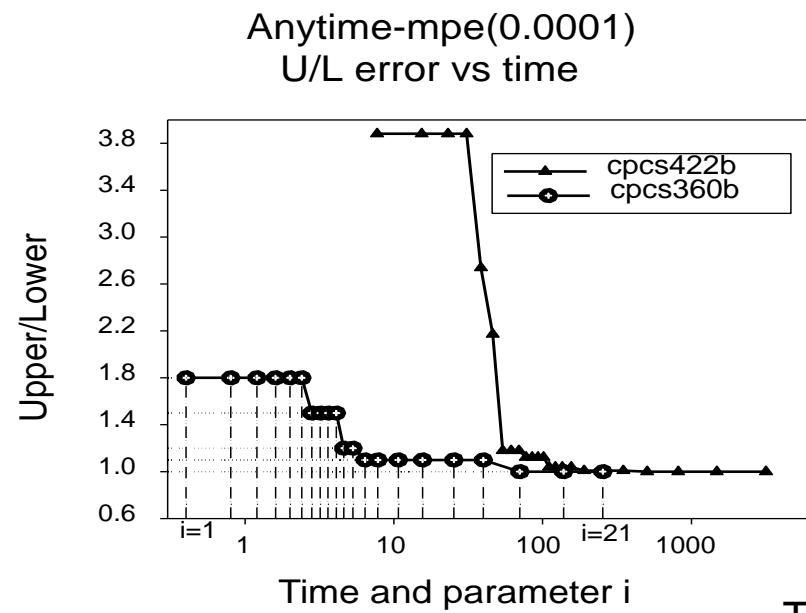
Comparing MBE-mpe and anytime-mpe
versus BE-mpe

Methodology for Empirical Evaluation (for mpe)

- U/L –accuracy
- Better (U/mpe) or mpe/L
- Benchmarks: Random networks
 - Given n, e, v generate a random DAG
 - For x_i and parents generate table from uniform [0,1], or noisy-or
- Create k instances. For each, generate random evidence, likely evidence
- Measure averages

CPCS Networks – Medical Diagnosis (noisy-OR model)

Test case: no evidence



Algorithm	cpcs360	cpcs422
elim-mpe	115.8	1697.6
anytime-mpe(ϵ), $\epsilon = 10^{-4}$	70.3	505.2
anytime-mpe(ϵ), $\epsilon = 10^{-1}$	70.3	110.5

Agenda

- Mini-bucket elimination
- Weighted Mini-bucket
- Mini-clustering
- Iterative Belief propagation
- Iterative-join-graph propagation

The Power Sum and Holder Inequality

The power sum is defined as follows:

$$\sum_x^w f(x) = \left(\sum_x f(x)^{\frac{1}{w}} \right)^w \quad (1.2)$$

where w is a non-negative weight. The power sum reduce to a standard summation when $w = 1$ and approaches max when $w \rightarrow 0^+$.

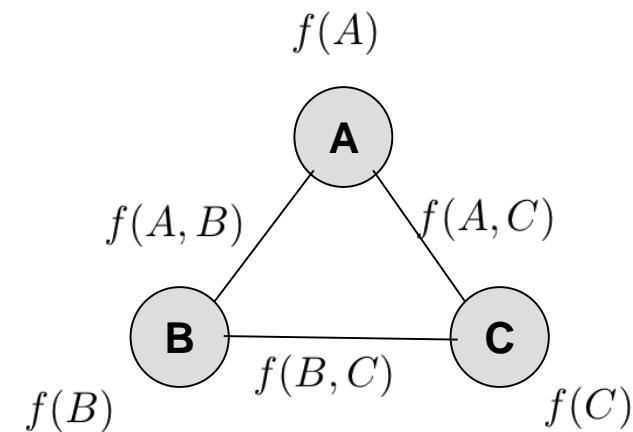
[Holder inequality] Let $f_i(x)$, $i = 1..r$ be a set of functions and w_1, \dots, w_r b a set of of non-zero weights, s.t., $w = \sum_{i=1}^r w_i$ then,

$$\sum_x^w \prod_{i=1}^r f_i(x) \leq \prod_{i=1}^r \sum_x^{w_i} f_i(x)$$

Working Example

- Model:
 - Markov network
- Task:
 - Partition function

$$Z = \sum_{A,B,C} f(A)f(B)f(C)f(A,B)f(A,C)f(B,C)$$



(Qiang Liu slides)

Mini-Bucket (Basic Principles)

- Upper bound

$$\sum_i a_i b_i \leq \left(\sum_i a_i \right) \max_i(b_i)$$

- Lower bound

$$\sum_i a_i b_i \geq \left(\sum_i a_i \right) \min_i(b_i)$$

(Qiang Liu slides)

I am using a_i b_i to represent the general constant.

Hölder Inequality

$$\sum_i a_i b_i \leq \left(\sum_i a_i^{1/w_1} \right)^{w_1} \left(\sum_i b_i^{1/w_2} \right)^{w_2}$$

- Where $a_i > 0, b_i > 0$ and $w_1 + w_2 = 1$ $w_1 > 0, w_2 > 0$
- When $\frac{a_i^{1/w_1}}{\sum a_i^{1/w_1}} = \frac{b_i^{1/w_2}}{\sum b_i^{1/w_2}}$ the equality is achieved.

(Qiang Liu slides)

G. H. Hardy, J. E. Littlewood and G. Pólya, *Inequalities*, Cambridge Univ. Press, London and New York, 1934.

Reverse Holder Inequality

- If $w_1 + w_2 = 1$, but $w_1 < 0, w_2 > 1$
the direction of the inequality reverses.

$$\sum_i a_i b_i \geq \left(\sum_i a_i^{1/w_1} \right)^{w_1} \left(\sum_i b_i^{1/w_2} \right)^{w_2}$$

(Qiang Liu slides)

G. H. Hardy, J. E. Littlewood and G. Pólya, *Inequalities*, Cambridge Univ. Press, London and New York, 1934.

Mini-Bucket as Holder Inequality

$$w_1 = 1; w_2 \rightarrow +0 \quad (\text{mbe-bel-max})$$

$$\sum_i a_i b_i \leq \left(\sum_i a_i \right) \max_i(b_i)$$

► $w_1 = 1, w_2 \rightarrow -0$ (mbe-bel-min)

$$\sum_i a_i b_i \geq \left(\sum_i a_i \right) \min_i(b_i)$$

Algorithm Weighted WMBE(i,m), ($w_1, \dots w_n$)

Input: A belief network $\mathcal{B} = \langle \mathbf{X}, \mathbf{D}, \mathbf{P}_G, \Pi \rangle$, an ordering $d = (X_1, \dots, X_n)$; evidence e

Output: an upper bound on $\sum_{\mathbf{X}}^w \prod_{i=1}^n P_i$

1. **Initialize:** Partition $P = \{P_1, \dots, P_n\}$ into buckets $bucket_1, \dots, bucket_n$, where $bucket_k$ contains all CPTs h_1, h_2, \dots, h_t whose highest-index variable is X_k .

2. **Backward:** for $k = n$ to 1 do

- If X_p is observed ($X_k = a$), assign $X_k \leftarrow a$ in each h_j and put the result in the highest-variable bucket of its scope (put constants in $bucket_1$).
- Else for h_1, h_2, \dots, h_t in $bucket_k$ Generate an (i, m) -partitioning, $Q' = \{Q_1, \dots, Q_r\}$. Select a set of weights $w_1, \dots w_r$ s.t $\sum_l w_l = w$. For each $Q_l \in Q'$, containing $h_{l_1}, \dots h_{l_t}$, do

$$h_l \leftarrow \sum_{X_k}^{w_l} \prod_{j=1}^t h_{l_j} = \left(\sum_{X_k} \prod_{j=1}^t (h_{l_j})^{w_l} \right)^{\frac{1}{w_l}}$$

Add h_l to the bucket of the highest-index variable in its scope (and put constant functions in $bucket_1$).

3. **Return** $U \leftarrow$ the weighted product of functions in the bucket of X_1 , which is an upper bound on $P(x_1, e)$.

Weighted Mini-Bucket

(for summation)

Exact bucket elimination:

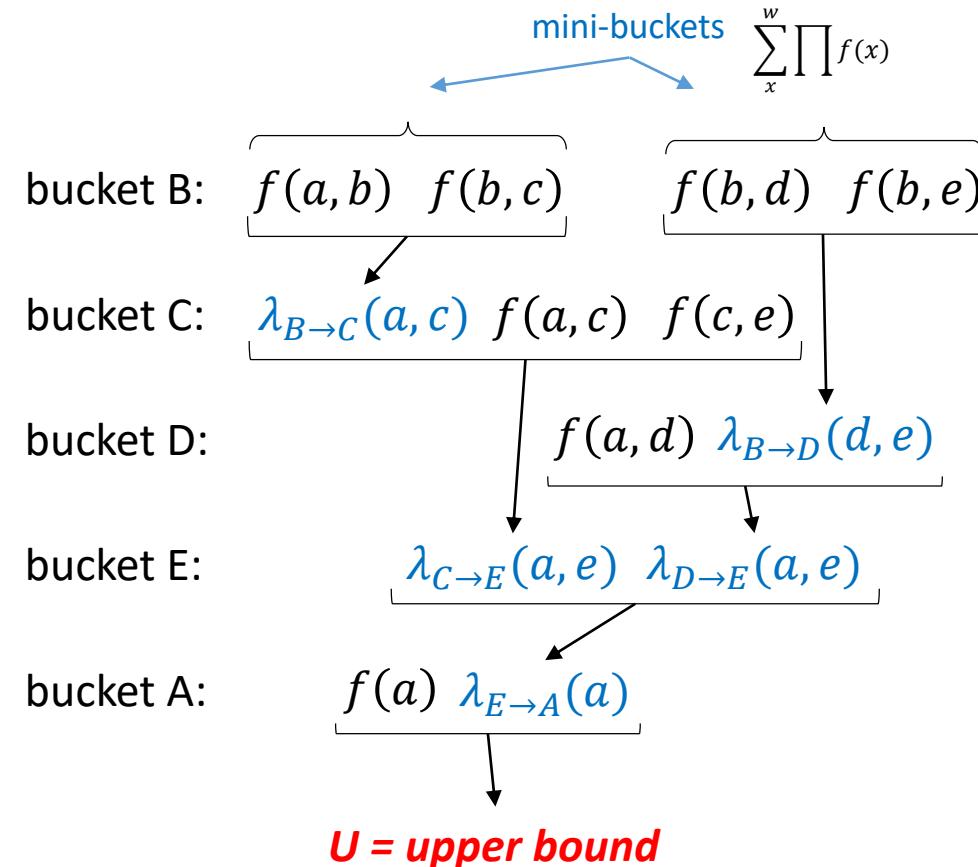
$$\begin{aligned}\lambda_B(a, c, d, e) &= \sum_b [f(a, b) \cdot f(b, c) \cdot f(b, d) \cdot f(b, e)] \\ &\leq \left[\sum_b^{w_1} f(a, b) f(b, c) \right] \cdot \left[\sum_b^{w_2} f(b, d) f(b, e) \right] \\ &= \lambda_{B \rightarrow C}(a, c) \cdot \lambda_{B \rightarrow D}(d, e)\end{aligned}$$

(mini-buckets)

where $\sum_x^w f(x) = \left[\sum_x f(x)^{\frac{1}{w}} \right]^w$
is the weighted or “power” sum operator

$$\sum_x^w f_1(x) f_2(x) \leq \left[\sum_x^{w_1} f_1(x) \right] \left[\sum_x^{w_2} f_2(x) \right]$$

where $w_1 + w_2 = w$ and $w_1 > 0, w_2 > 0$
(lower bound if $w_1 > 0, w_2 < 0$)



[Liu and Ihler, 2011]

Choosing the weights

- $w_1 = 1/2, w_2 = 1/2$ (Cauchy–Schwarz inequality)

$$\sum_i a_i b_i \leq \left(\sum_i a_i^2 \right)^{1/2} \left(\sum_i b_i^2 \right)^{1/2}$$

What is the optimal weights?



(Qiang Liu slides)

Allocating the probabilities

$$\sum_i a_i b_i \leq \left(\sum_i (a_i \xi_i)^{1/w_1} \right)^{w_1} \left(\sum_i (b_i \frac{1}{\xi_i})^{1/w_2} \right)^{w_2}$$

Notice that $a_i b_i = (a_i \xi_i) \cdot (b_i \frac{1}{\xi_i})$

What is the optimal allocation ?



(Qiang Liu slides)

Extention of Mini-Bucket

- Allocation of the probability:

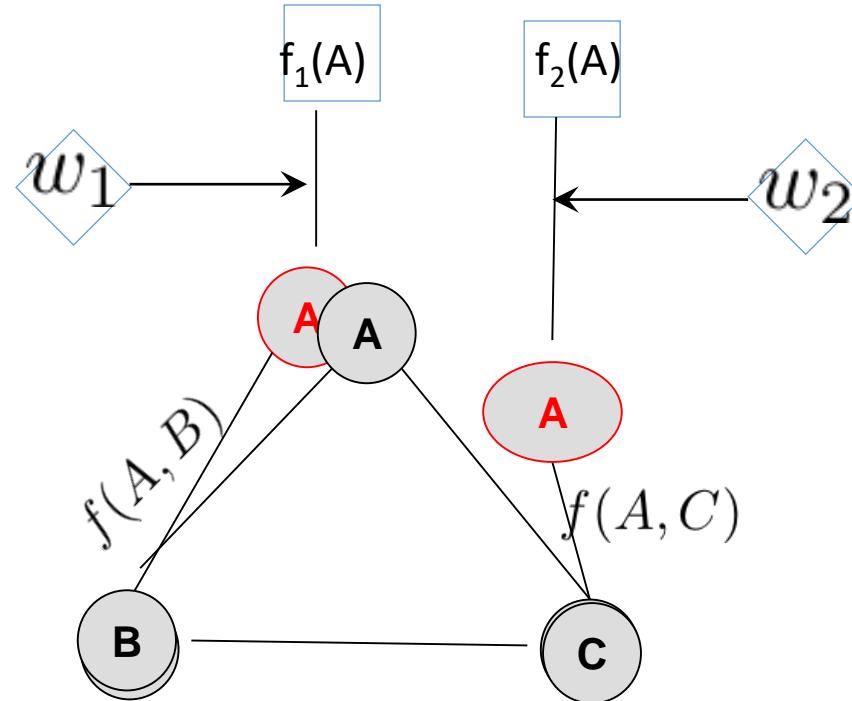
$$f_1(A)f_2(A) = f(A)$$

- weights: w_1 and w_2

$$w_1 + w_2 = 1$$

$$\sum_A f(A)f(A, B)f(A, C)$$

$$\leq \left(\sum_A (f_1(A)f(A, C)^{1/w_1})^{w_1} \right) \left(\sum_A (f_2(A)f(B, C))^{1/w_2} \right)^{w_2}$$



(Qiang Liu slides)

Algorithm MBE-map(i,m)

Input: A Bayesian network $\mathcal{B} = \langle \mathbf{X}, \mathbf{D}, \mathbf{P}_G, \prod \rangle$, $P = \{P_1, \dots, P_n\}$; a subset of hypothesis variables $A = \{A_1, \dots, A_k\}$; an ordering of the variables, d , in which the A 's are first in the ordering; observations e . ψ_i is the product of input function in the bucket of X_i .

Output: An upper bound on the map and a and a suboptimal solution $A = a_k^a$.

1. Initialize: Generate an ordered partition of the conditional probability functions, $bucket_1, \dots, bucket_n$, where $bucket_i$ contains all functions whose highest variable is X_i .

2. Backwards For $p \leftarrow n$ down to 1, do

for all the message functions $\beta_1, \beta_2, \dots, \beta_j$ in $bucket_p$ do

- If (observed variable) $bucket_p$ contains the observation $X_p = x_p$, assign $X_p = x_p$ to each β_i and ψ_p and put each in appropriate bucket.
- Else for h_1, h_2, \dots, h_j in $bucket_p$ generate an (i, m) -partitioning, Q' mini-buckets Q_1, \dots, Q_r .
- If $X_p \notin A$ (not a hypothesis variable)
foreach $Q_l \in Q'$, containing h_{l_1}, \dots, h_{l_t} , do

$$h_l \leftarrow \sum_{X_p} \Pi_{i=1}^t h_{l_i}, \text{ if } l = 1$$

$$h_l \leftarrow \max_{X_p} \Pi_{i=1}^t h_{l_i}, \text{ if } l \neq 1$$

- Else, $(X_p \in A)$ (a hypothesis variable)
for each $Q_l \in Q'$ containing h_{l_1}, \dots, h_{l_t}

$$h_l \leftarrow \max_{X_p} \Pi_{i=1}^t h_{l_i}$$

Add h^l to the bucket of the highest-index variable in $U_l \leftarrow \bigcup_{i=1}^t scope(h_{l_i}) - \{X_p\}$, (put constants in $bucket_1$).

3. Forward: for $p = 1$ to k , given $A_1 = a_1^a, \dots, A_{p-1} = a_{p-1}^a$, assign a value a_p^a to A_p that maximizes the product of all functions in $bucket_p$.

4. Return An upper bound $U = \max_{a_1} \prod_{h_i \in bucket_1} h_i$ on map, computed in the first bucket, and the assignment $a_k^a = (a_1^a, \dots, a_k^a)$.

Figure 8.8: Algorithm $MBE-map(i, m)$.

MBE-map

Process max buckets
 With max mini-buckets
 And sum buckets with sum
 Mini-bucket and max
 mini-buckets

Algorithm MBE-map(i,m)

Input: A Bayesian network $\mathcal{B} = \langle \mathbf{X}, \mathbf{D}, \mathbf{P}_G, \prod \rangle$, $P = \{P_1, \dots, P_n\}$; a subset of hypothesis variables $A = \{A_1, \dots, A_k\}$; an ordering of the variables, d , in which the A 's are first in the ordering; observations e .

Output: An upper bound on the map and a suboptimal solution $A = \bar{a}_k^a$.

1. **Initialize:** Partition $P = \{P_1, \dots, P_n\}$ into $bucket_1, \dots, bucket_n$, where $bucket_i$ contains all functions whose highest variable is X_i .

2. **Backwards** For $p \leftarrow n$ downto 1, do

for all the functions h_1, h_2, \dots, h_j in $bucket_p$ do

- If (observed variable) $bucket_p$ contains the observation $X_p = x_p$, assign $X_p = x_p$ to each h_i and put each in appropriate bucket.
- Else for h_1, h_2, \dots, h_j in $bucket_p$ generate an (i, m) -partitioning, $Q' = \{Q_1, \dots, Q_r\}$.
- If $X_p \notin A$ assign $w_p = 1$, otherwise $w_p = 0$. Select weights for the mini-buckets in X_p bucket: w_{p1}, \dots, w_{pr} . s.t $\sum_i w_{pi} = w_p$.
 foreach $Q_l \in Q'$, containing h_{l1}, \dots, h_{lt} , do

$$h_l \leftarrow \sum_{X_k}^{w_{pl}} \prod_{j=1}^t h_{lj} = \left(\sum_{X_k} \left(\prod_{j=1}^t h_{lj} \right)^{w_{pl}} \right)^{\frac{1}{w_{pl}}}$$

Add h_l to the bucket of the highest-index variable in its scope.

3. **Forward:** for $p = 1$ to k , given $A_1 = a_1^a, \dots, A_{p-1} = a_{p-1}^a$, assign a value a_p^a to A_p that maximizes the product of all functions in $bucket_p$. conditioned on earlier assignments.

4. **Return** An upper bound $U = \max_{a_1} \prod_{h_i \in bucket_1} h_i$ on the map value, computed in the first bucket, and the assignment $\bar{a}_k^a = (a_1^a, \dots, a_k^a)$.

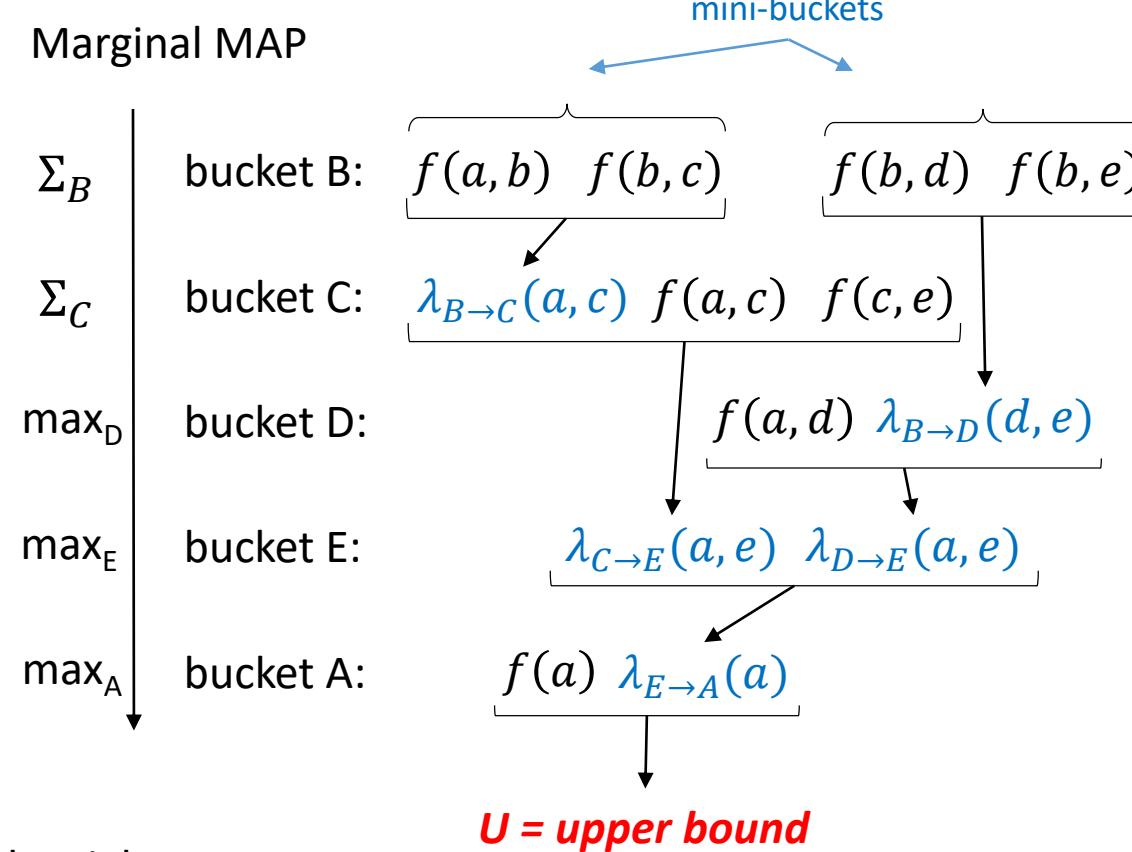
Figure 8.7: Algorithm MBE-map(i, m).

MB and WMB for Marginal MAP

$$\max_Y \sum_{X \setminus Y} \prod_j P_j$$

$$\begin{aligned}\lambda_{B \rightarrow C}(a, c) &= \sum_b^{w_1} f(a, b)f(b, c) \\ \lambda_{B \rightarrow D}(d, e) &= \sum_b^{w_2} f(b, d)f(b, e) \\ &\vdots \\ \lambda_{E \rightarrow A}(a) &= \max_e \lambda_{C \rightarrow E}(a, e)\lambda_{D \rightarrow E}(a, e)\end{aligned}$$

$$U = \max_a f(a)\lambda_{E \rightarrow A}(a)$$

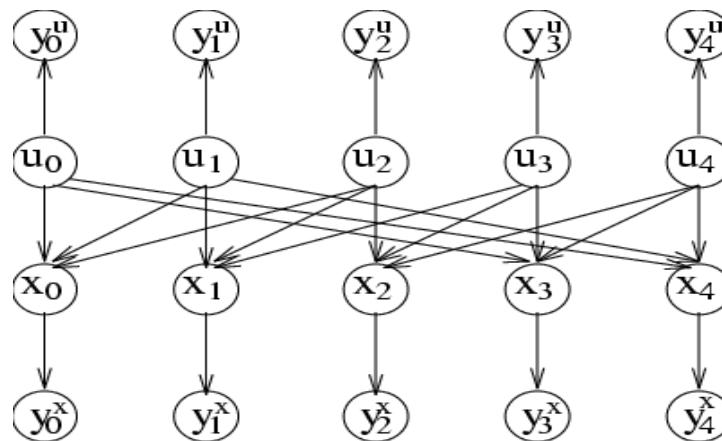


Can optimize over cost-shifting and weights
(single pass “MM” or iterative message passing)

[Liu and Ihler, 2011; 2013]
[Dechter and Rish, 2003]

Probabilistic decoding

Error-correcting linear block code



State-of-the-art:
approximate algorithm – iterative belief propagation (IBP)
(Pearl's poly-tree algorithm applied to loopy networks)

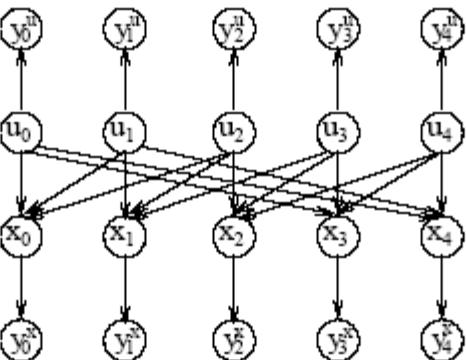


Figure 7.7: Belief network for a linear block code.

Example 7.3.1 We will next demonstrate the mini-bucket approximation for MAP on an example of *probabilistic decoding* (see Chapter 2) Consider a belief network which describes the decoding of a *linear block code*, shown in Figure 7.7. In this network, U_i are *information bits* and X_j are *code bits*, which are functionally dependent on U_i . The vector (U, X) , called the channel input, is transmitted through a noisy channel which adds Gaussian noise and results in the channel output vector $Y = (Y^u, Y^x)$. The decoding task is to assess the most likely values for the U 's given the observed values $Y = (\bar{y}^u, \bar{y}^x)$, which is the MAP task where U is the set of hypothesis variables, and $Y = (\bar{y}^u, \bar{y}^x)$ is the evidence. After processing the observed buckets we get the following bucket configuration (lower case y 's are observed values):

$$\begin{aligned} \text{bucket}(X_0) &= P(y_0^x|X_0), P(X_0|U_0, U_1, U_2), \\ \text{bucket}(X_1) &= P(y_1^x|X_1), P(X_1|U_1, U_2, U_3), \\ \text{bucket}(X_2) &= P(y_2^x|X_2), P(X_2|U_2, U_3, U_4), \\ \text{bucket}(X_3) &= P(y_3^x|X_3), P(X_3|U_3, U_4, U_0), \\ \text{bucket}(X_4) &= P(y_4^x|X_4), P(X_4|U_4, U_0, U_1), \\ \text{bucket}(U_0) &= P(U_0), P(y_0^u|U_0), \\ \text{bucket}(U_1) &= P(U_1), P(y_1^u|U_1), \\ \text{bucket}(U_2) &= P(U_2), P(y_2^u|U_2), \\ \text{bucket}(U_3) &= P(U_3), P(y_3^u|U_3), \\ \text{bucket}(U_4) &= P(U_4), P(y_4^u|U_4). \end{aligned}$$

Initial
partitioning

Processing by *mbe-map(4,1)* of the first top five buckets by summation and the rest by maximization, results in the following mini-bucket partitionings and function generation:

$$\begin{aligned}
& \text{bucket}(X_0) = \{P(y_0^x|X_0), P(X_0|U_0, U_1, U_2)\}, \\
& \text{bucket}(X_1) = \{P(y_1^x|X_1), P(X_1|U_1, U_2, U_3)\}, \\
& \text{bucket}(X_2) = \{P(y_2^x|X_2), P(X_2|U_2, U_3, U_4)\}, \\
& \text{bucket}(X_3) = \{P(y_3^x|X_3), P(X_3|U_3, U_4, U_0)\}, \\
& \text{bucket}(X_4) = \{P(y_4^x|X_4), P(X_4|U_4, U_0, U_1)\}. \\
& \text{bucket}(U_0) = \boxed{\{P(U_0), P(y_0^u|U_0), h^{X_0}(U_0, U_1, U_2)\}}, \boxed{\{h^{X_3}(U_3, U_4, U_0)\}} \boxed{\{h^{X_4}(U_4, U_0, U_1)\}}, \\
& \text{bucket}(U_1) = \{P(U_1), P(y_1^u|U_1), h^{X_1}(U_1, U_2, U_3), h^{U_0}(U_1, U_2)\}, \{h^{U_0}(U_4, U_1)\}, \\
& \text{bucket}(U_2) = \{P(U_2), P(y_2^u|U_2), h^{X_2}(U_2, U_3, U_4), h^{U_1}(U_2, U_3)\}, \\
& \text{bucket}(U_3) = \{P(U_3), P(y_3^u|U_3), h^{U_0}(U_3, U_4), h^{U_1}(U_3, U_4), h^{U_2}(U_3, U_4)\}, \\
& \text{bucket}(U_4) = \{P(U_4), P(y_4^u|U_4), h^{U_1}(U_4), h^{U_3}(U_4)\}.
\end{aligned}$$

The first five buckets are not partitioned at all and are processed as full buckets, since in this case a full bucket is a (4,1)-partitioning. This processing generates five new functions, three are placed in bucket U_0 , one in bucket U_1 and one in bucket U_2 . Then bucket U_0 is partitioned into three mini-buckets processed by maximization, creating two functions placed in bucket U_1 and one function placed in bucket U_3 . Bucket U_1 is partitioned into two mini-buckets, generating functions placed in bucket U_2 and bucket U_3 . Subsequent buckets are processed as full buckets. Note that the scope of recorded functions is bounded by 3.

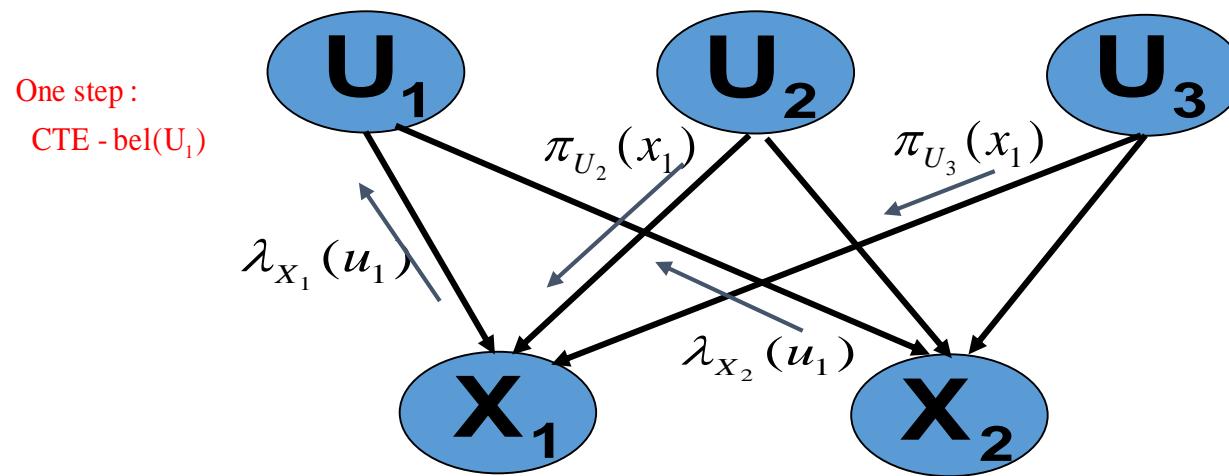
In the bucket of U_4 we get an upper bound U satisfying $U \geq \text{MAP} = P(U, \bar{y}^u, \bar{y}^x)$ where \bar{y}^u and \bar{y}^x are the observed outputs for the U 's and the X 's bits transmitted. In order to bound $P(U|\bar{e})$, where $\bar{e} = (\bar{y}^u, \bar{y}^x)$, we need $P(\bar{e})$ which is not available. Yet, again, in most cases we are interested in the ratio $P(U = \bar{u}_1|\bar{e})/P(U = \bar{u}_2|\bar{e})$ for competing hypotheses $U = \bar{u}_1$ and $U = \bar{u}_2$ rather than in the absolute values. Since $P(U|\bar{e}) = P(U, \bar{e})/P(\bar{e})$ and the probability of the evidence is just a constant factor independent of U , the ratio is equal to $P(U_1, \bar{e})/P(U_2, \bar{e})$. \square

Complexity and Tractability of MBE(i,m)

Theorem 7.6.1 *Algorithm mbe(i,m) takes $O(r \cdot \exp(i))$ time and space, where r is the number of input functions², and where $|F|$ is the maximum scope of any input function, $|F| \leq i \leq n$. For $m = 1$, the algorithm is time and space $O(r \cdot \exp(|F|))$.*

Iterative Belief Propagation

- Belief propagation is exact for poly-trees
- IBP - applying BP iteratively to cyclic networks



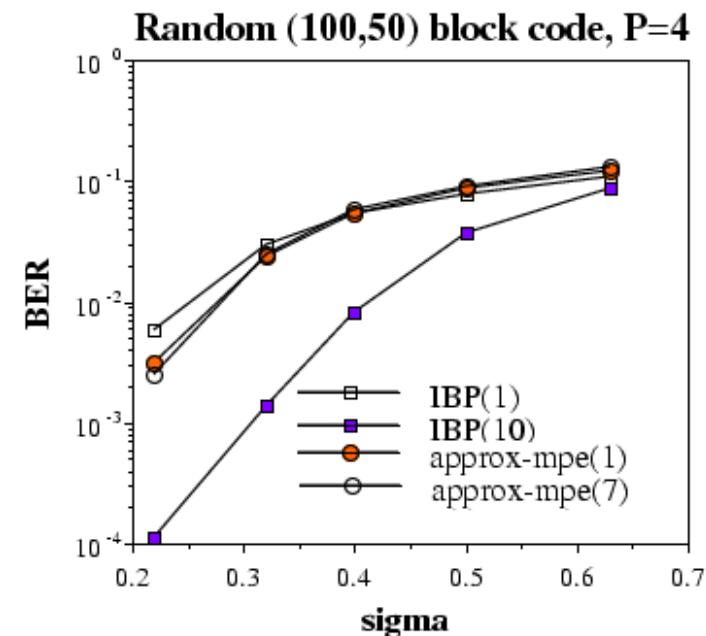
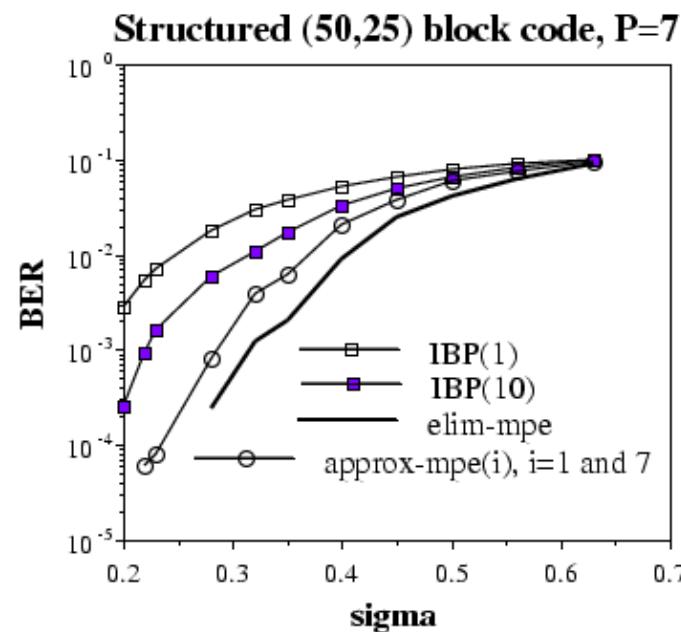
- No guarantees for convergence
- Works well for many coding networks

MBE-mpe vs. IBP

mbe - mpe is better on low - w * codes

IBP is better on randomly generated (high - w *) codes

Bit error rate (BER) as a function of noise (sigma):



Mini-Buckets: Summary

- Mini-buckets – local inference approximation
- Idea: bound size of recorded functions
- MBE-mpe(i) - mini-bucket algorithm for MPE
 - Better results for noisy-OR than for random problems
 - Accuracy increases with decreasing noise in coding
 - Accuracy increases for likely evidence
 - Sparser graphs -> higher accuracy
 - Coding networks: MBE-mpe outperforms IBP on low-induced width codes

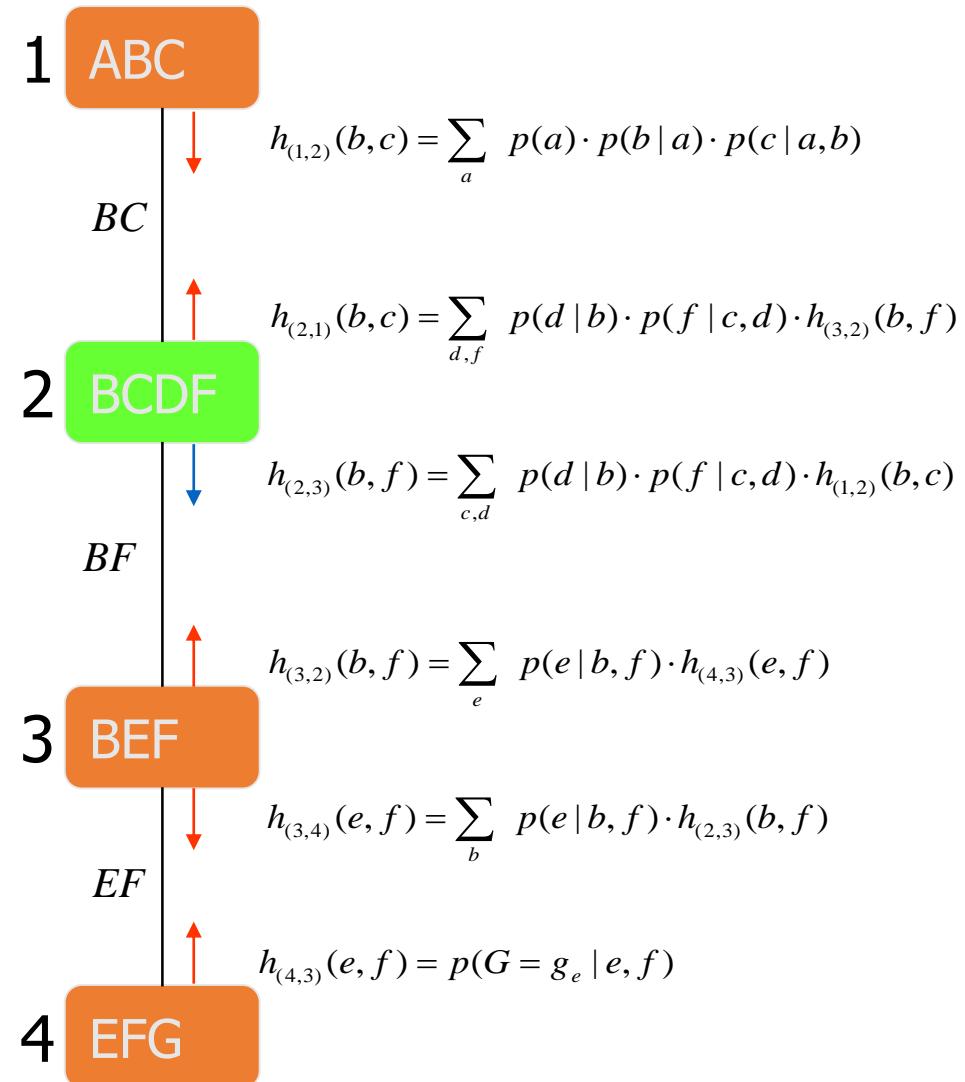
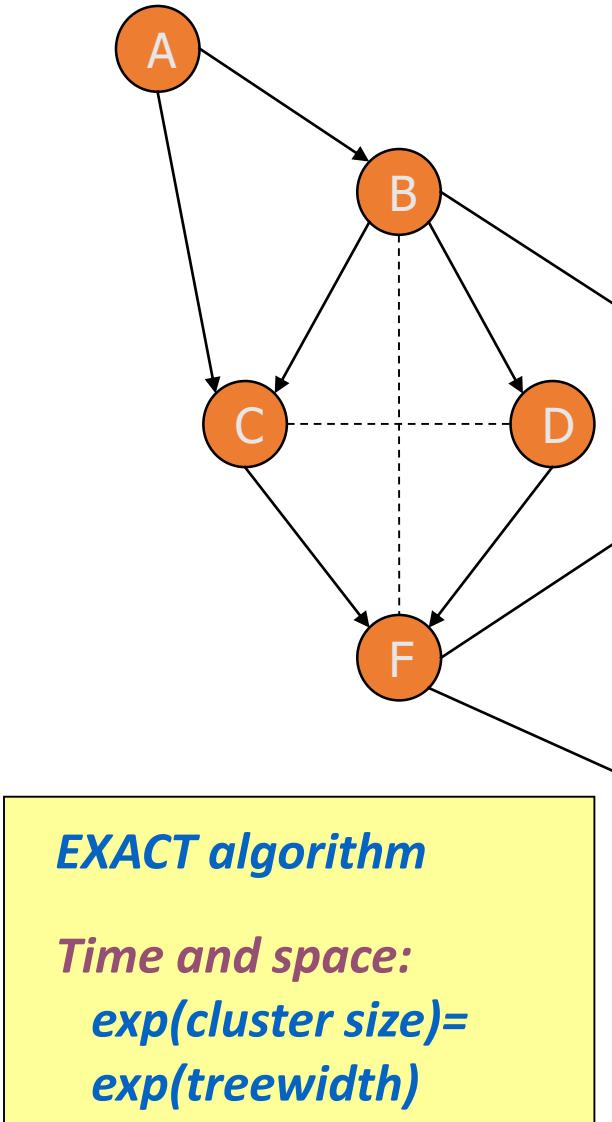
Agenda

- Mini-bucket elimination
- **Mini-clustering**
- Iterative Belief propagation
- Iterative-join-graph propagation

Cluster Tree Elimination - properties

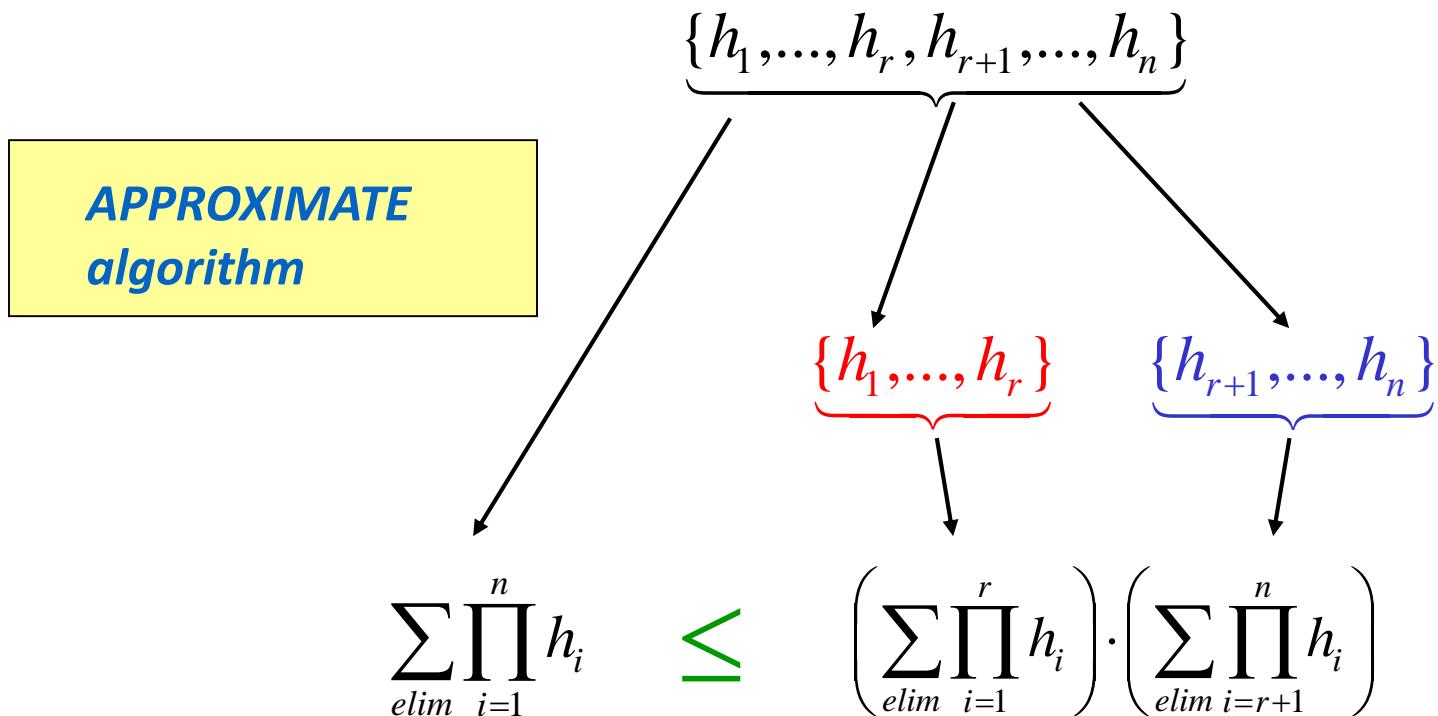
- Correctness and completeness: Algorithm CTE is correct, i.e. it computes the exact joint probability of a single variable and the evidence.
- Time complexity: $O(\deg \times (n+N) \times d^{w^*+1})$
- Space complexity: $O(N \times d^{sep})$
 - where
 - \deg = the maximum degree of a node
 - n = number of variables (= number of CPTs)
 - N = number of nodes in the tree decomposition
 - d = the maximum domain size of a variable
 - w^* = the induced width
 - sep = the separator size

Join-Tree Clustering (Cluster-Tree Elimination)



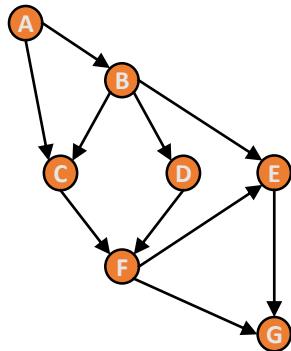
Mini-Clustering

Split a cluster into mini-clusters => bound complexity



Exponential complexity decrease $O(e^n) \rightarrow O(e^{\text{var}(r)}) + O(e^{\text{var}(n-r)})$

Mini-Clustering, i-bound=3



1

A B C
 $p(a), p(b|a), p(c|a,b)$

2

B C D
 $p(d|b), h_{(1,2)}^1(b,c)$

C D F
 $p(f|c,d)$

3

B E F
 $p(e|b,f),$
 $h_{(2,3)}^1(b), h_{(2,3)}^2(f)$

4

E F G
 $p(g|e,f)$

$$h_{(1,2)}^1(b,c) = \sum_a p(a) \cdot p(b|a) \cdot p(c|a,b)$$

$$h_{(2,3)}^1(b) = \sum_{c,d} p(d|b) \cdot h_{(1,2)}^1(b,c)$$

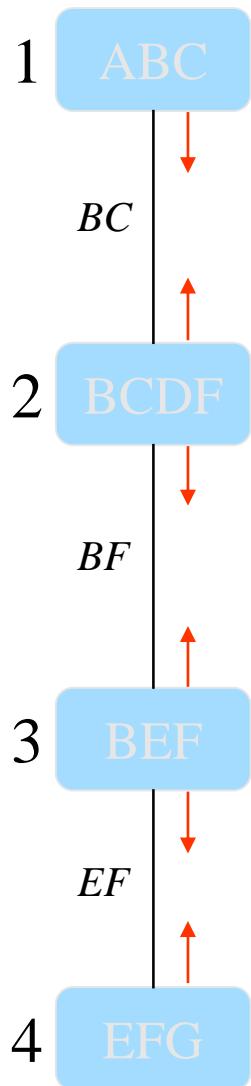
$$h_{(2,3)}^2(f) = \max_{c,d} p(f|c,d)$$

APPROXIMATE algorithm

Time and space:
 $\exp(i\text{-bound})$

Number of variables in a mini-cluster

Mini-Clustering - example



$$H_{(1,2)} \quad h_{(1,2)}^1(b, c) := \sum_a p(a) \cdot p(b | a) \cdot p(c | a, b)$$

$$H_{(2,1)} \quad h_{(2,1)}^1(b) := \sum_{d,f} p(d | b) \cdot h_{(3,2)}^1(b, f)$$

$$h_{(2,1)}^2(c) := \max_{d,f} p(f | c, d)$$

$$H_{(2,3)} \quad h_{(2,3)}^1(b) := \sum_{c,d} p(d | b) \cdot h_{(1,2)}^1(b, c)$$

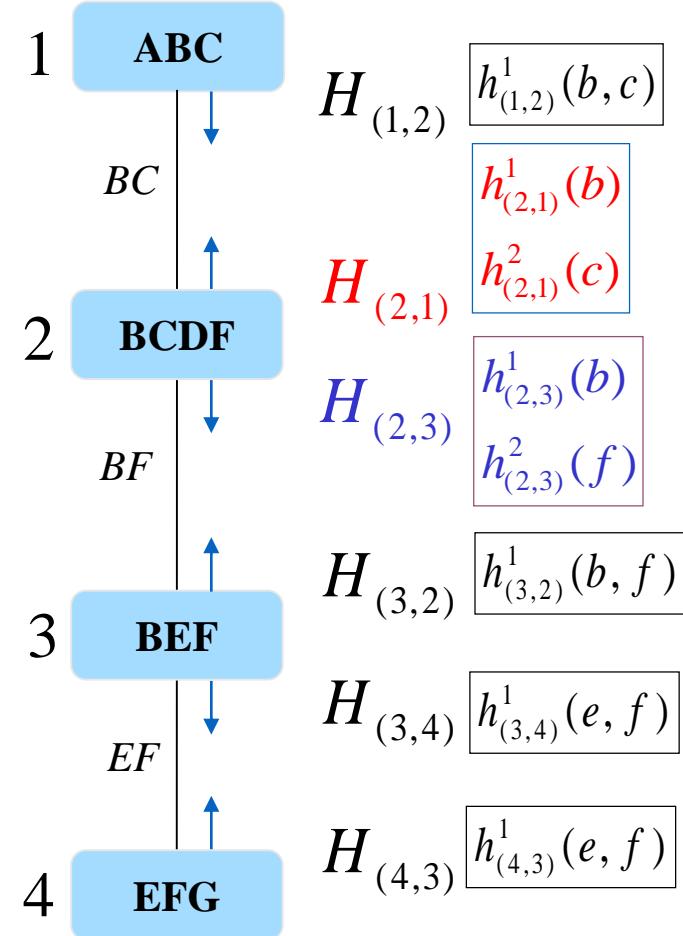
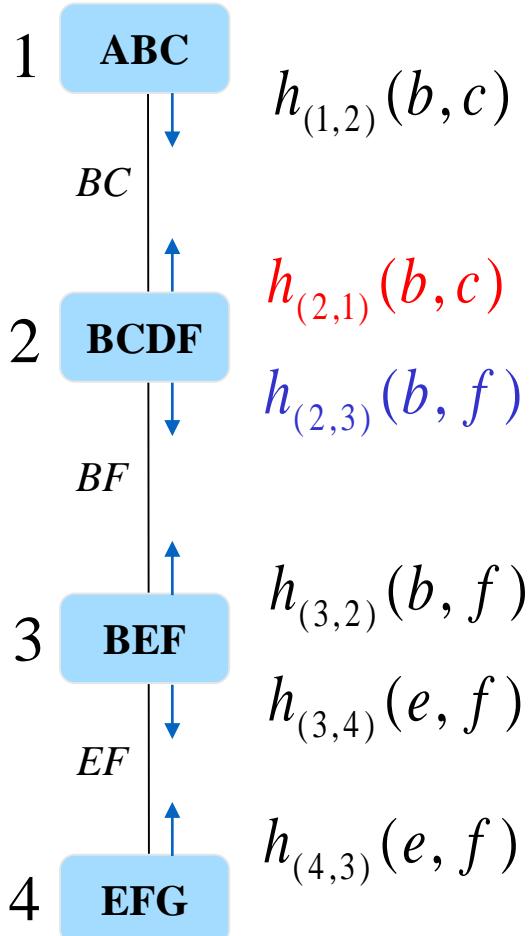
$$h_{(2,3)}^2(f) := \max_{c,d} p(f | c, d)$$

$$H_{(3,2)} \quad h_{(3,2)}^1(b, f) := \sum_e p(e | b, f) \cdot h_{(4,3)}^1(e, f)$$

$$H_{(3,4)} \quad h_{(3,4)}^1(e, f) := \sum_b p(e | b, f) \cdot h_{(2,3)}^1(b) \cdot h_{(2,3)}^2(f)$$

$$H_{(4,3)} \quad h_{(4,3)}^1(e, f) := p(G = g_e | e, f)$$

Cluster Tree Elimination vs. Mini-Clustering



Semantic of node duplication for mini-clustering

- We can have a different duplication of nodes going up and down. Example: going down.

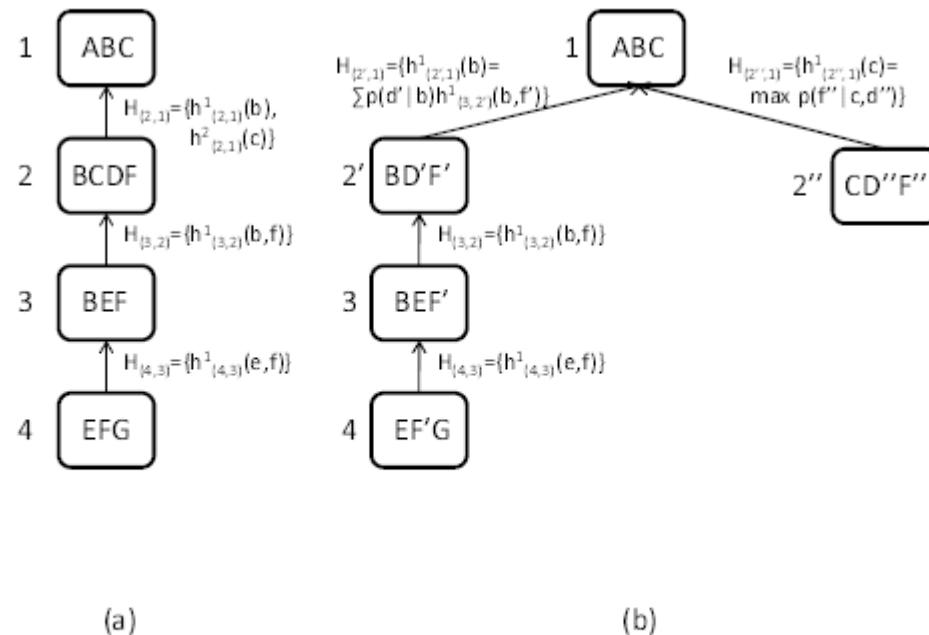


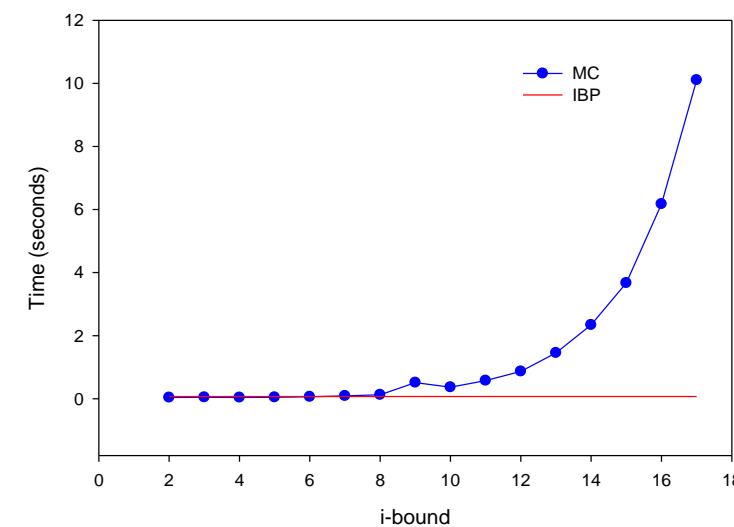
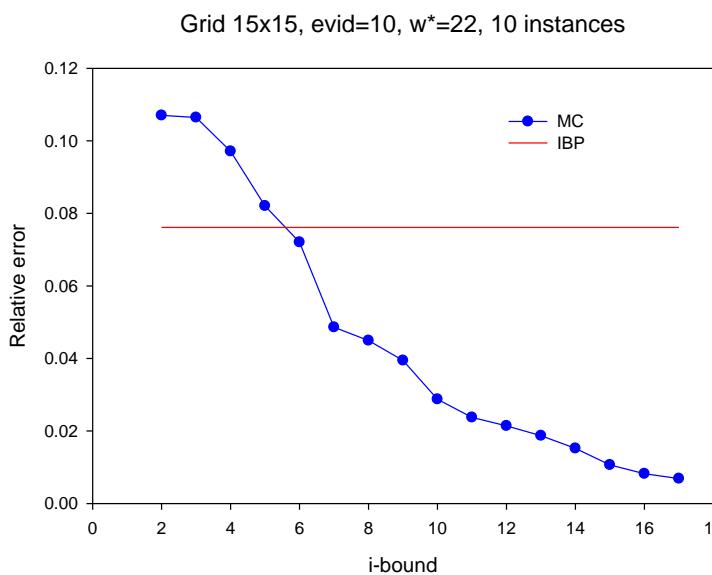
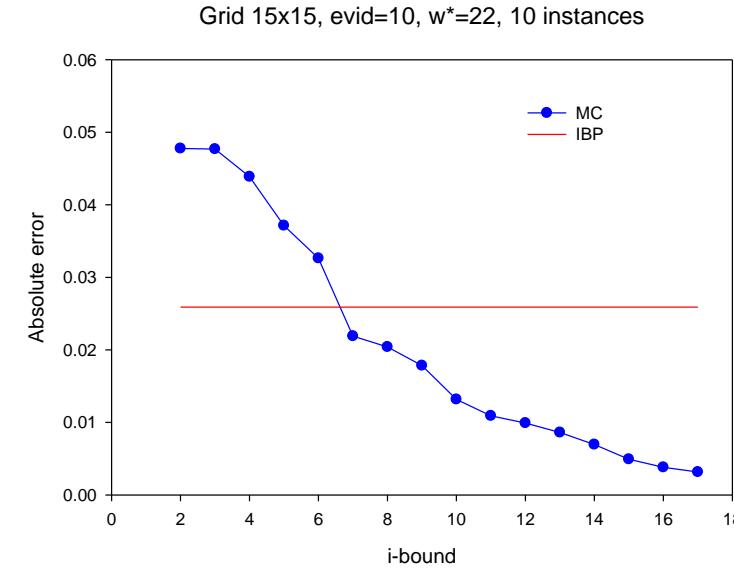
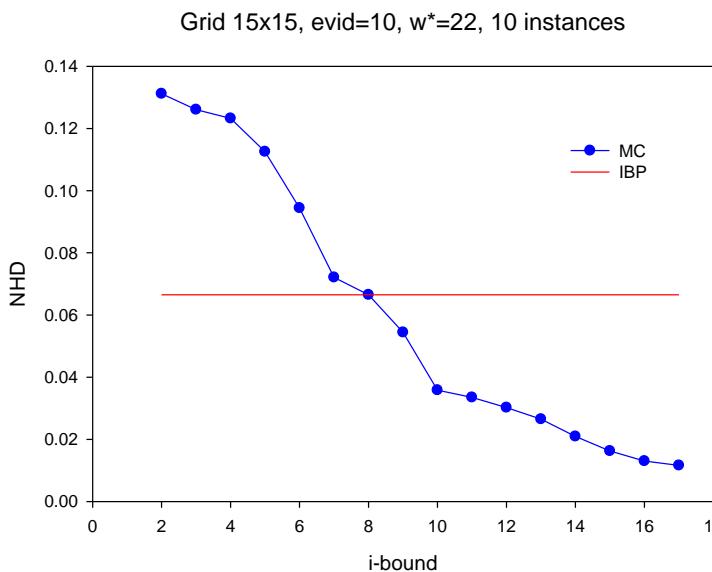
Figure 1.14: Node duplication semantics of MC: (a) trace of MC-BU(3); (b) trace of CTE-BU.

Lower Bounds and Mean Approximations

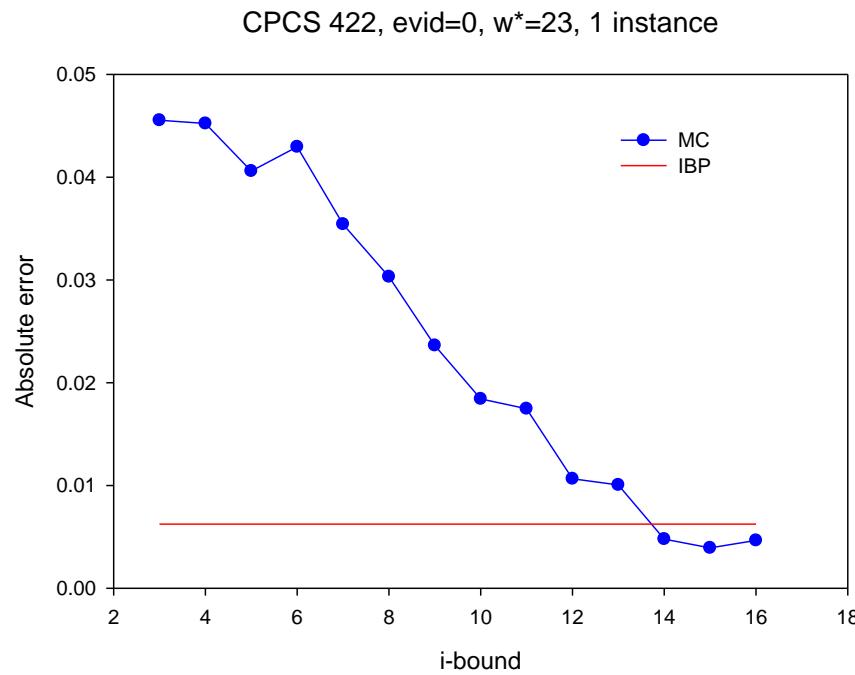
We can replace *max* operator by

- *min* => lower bound on the joint
- *mean* => approximation of the joint

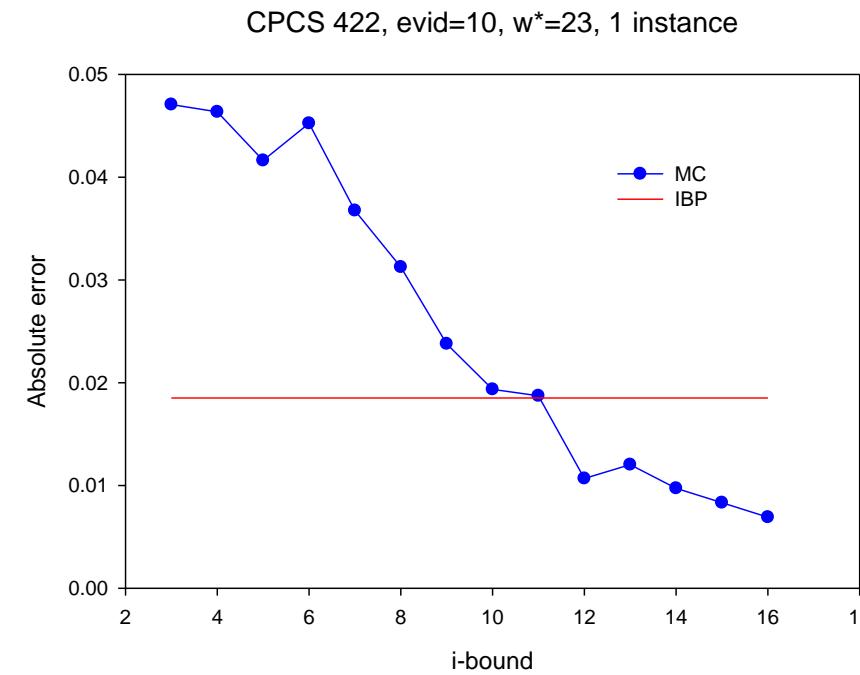
Grid 15x15 - 10 evidence



CPCS422 - Absolute error

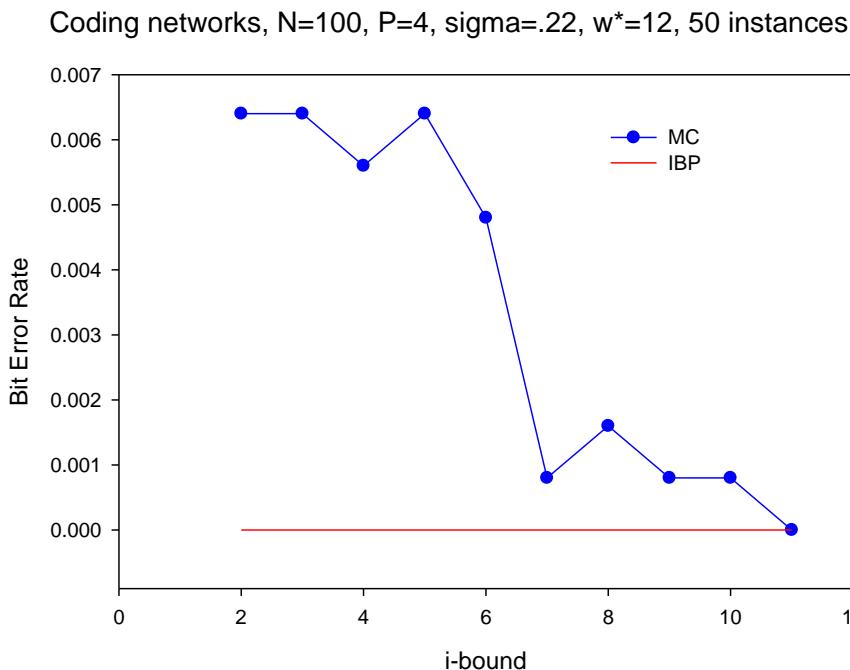


evidence=0

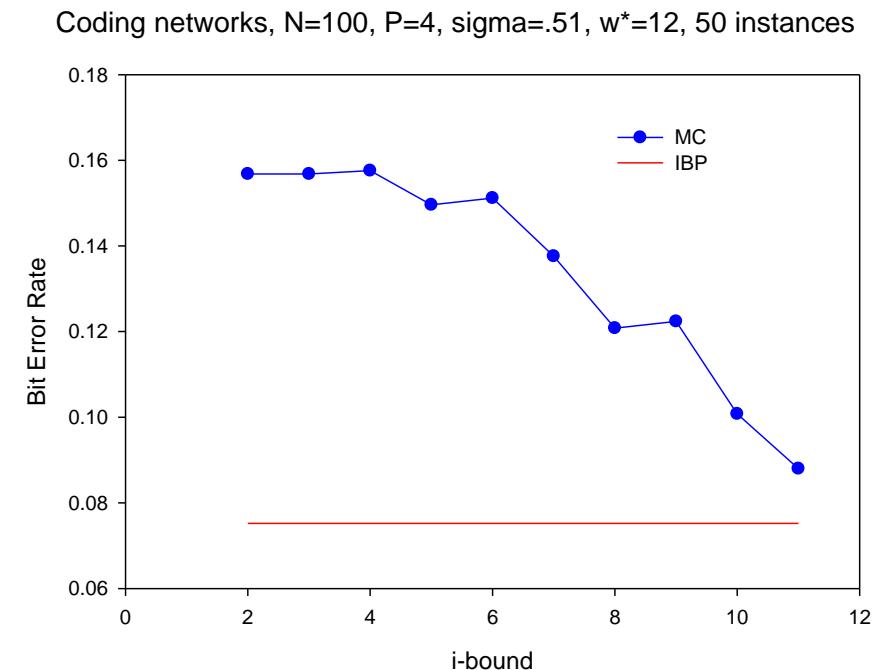


evidence=10

Coding networks - Bit Error Rate

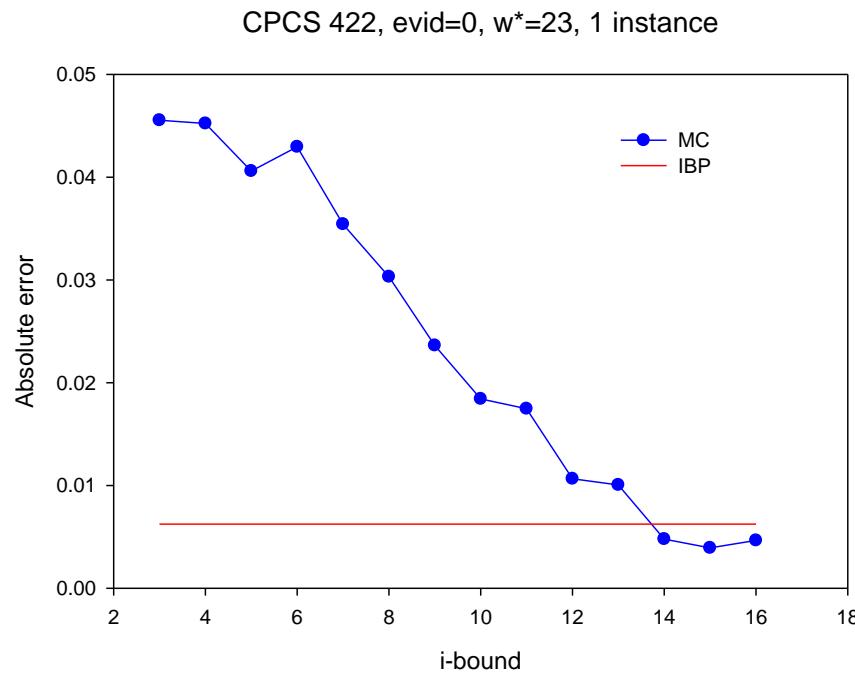


sigma=0.22

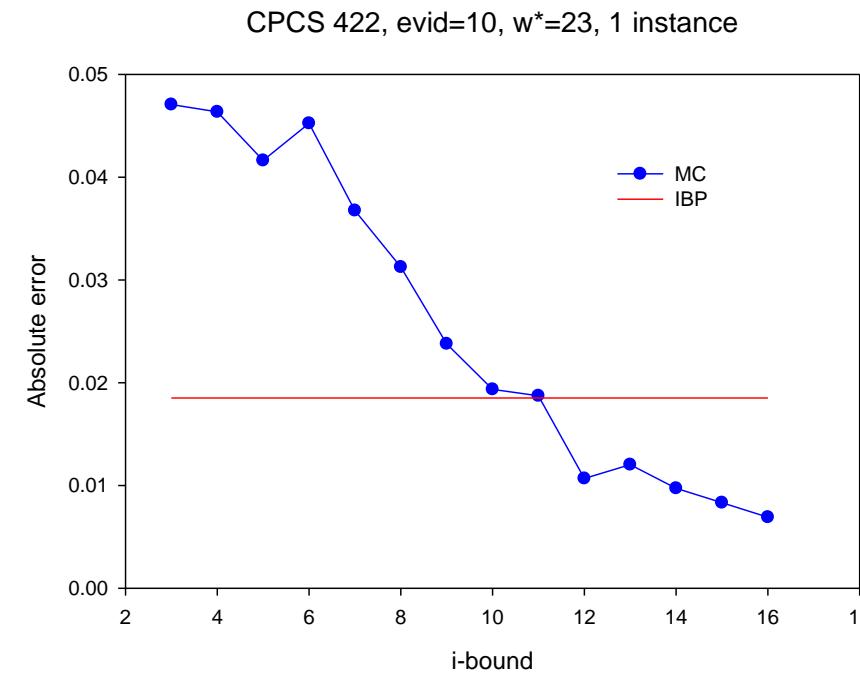


sigma=.51

CPCS422 - Absolute error

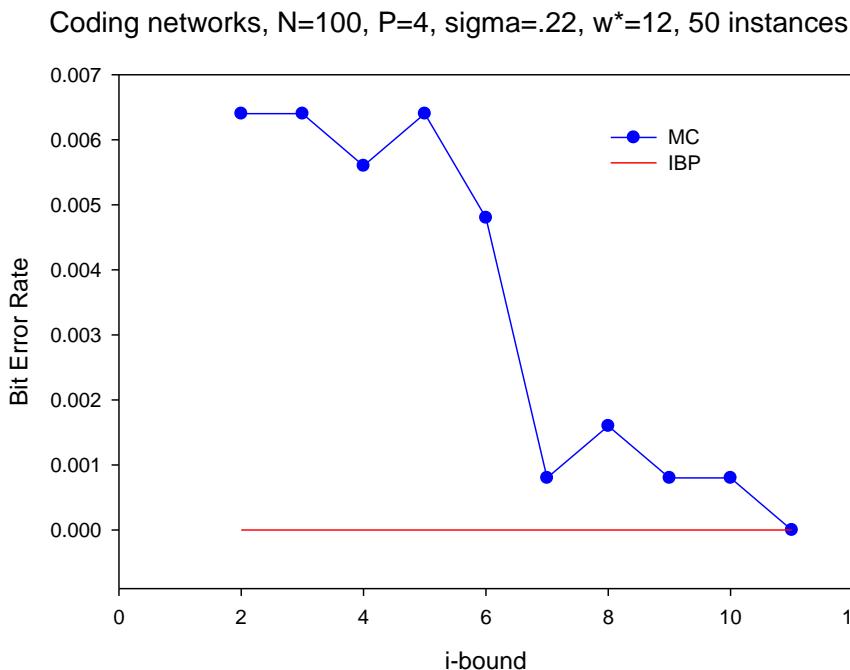


evidence=0

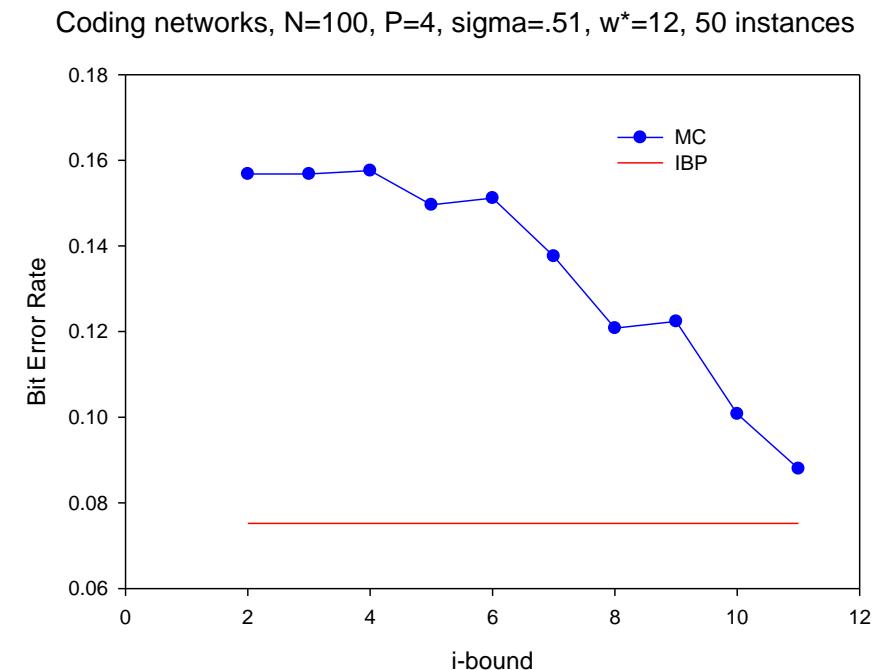


evidence=10

Coding networks - Bit Error Rate



sigma=0.22



sigma=.51

Heuristic for Partitioning

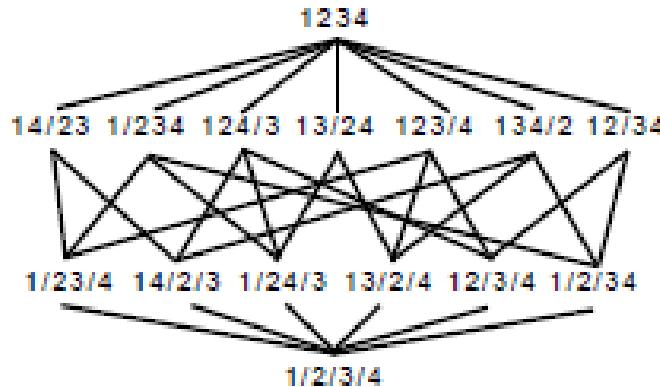
Scope-based Partitioning Heuristic. The *scope-based* partition heuristic (SCP) aims at minimizing the number of mini-buckets in the partition by including in each minibucket as many functions as possible as long as the i bound is satisfied. First, single function mini-buckets are decreasingly ordered according to their arity. Then, each minibucket is absorbed into the left-most minibucket with whom it can be merged.

The time and space complexity of $\text{Partition}(B, i)$, where B is the partitioned bucket, using the SCP heuristic is $O(|B| \log (|B|) + |B|^2)$ and $O(\exp(i))$, respectively.

The scope-based heuristic is quite fast, its shortcoming is that it does not consider the actual information in the functions.

Content-based heuristics

(Rollon and Dechter 2010)



Partitioning lattice of bucket $\{f_1, f_2, f_3, f_4\}$.

- Log relative error:

$$RE(f, h) = \sum_t (\log(f(t)) - \log(h(t)))$$

- Max log relative error:

$$MRE(f, h) = \max_t \{\log(f(t)) - \log(h(t))\}$$

Use greedy heuristic derived from a distance function to decide which functions go into a single mini-bucket

Greedy Partition as a function of a distance function h

```
function GreedyPartition( $B, i, h$ )
1. Initialize  $Q$  as the bottom partition of  $B$ ;
2. While  $\exists Q' \in ch(Q)$  which is a  $i$ -partition
 $Q \leftarrow \arg \min_{Q'} \{h(Q \rightarrow Q')\}$  among child  $i$ -partitions of  $Q$ ;
3.Return  $Q$ ;
```

Figure 8.13: Greedy partitioning

Proposition 8.6.5 *The time complexity of GreedyPartition is $O(|B| \times T)$ where $O(T)$ is the time complexity of selecting the min child partition according to h .*

Agenda

- Mini-bucket elimination
- Mini-clustering
- Reparameterization, cost-shifting
- Iterative Belief propagation
- Iterative-join-graph propagation

Cost-Shifting

(Reparameterization)

$+ \lambda(B)$

A	B	$f(A,B)$
b	b	6 + 3
b	g	0 - 1
g	b	0 + 3
g	g	6 - 1

$- \lambda(B)$

B	C	$f(B,C)$
b	b	6 - 3
b	g	0 - 3
g	b	0 + 1
g	g	6 + 1

B	$\lambda(B)$
b	3
g	-1

A	B	C	$f(A,B,C)$
b	b	b	12
b	b	g	6
b	g	b	0
b	g	g	6
g	b	b	6
g	b	g	0
g	g	b	6
g	g	g	12

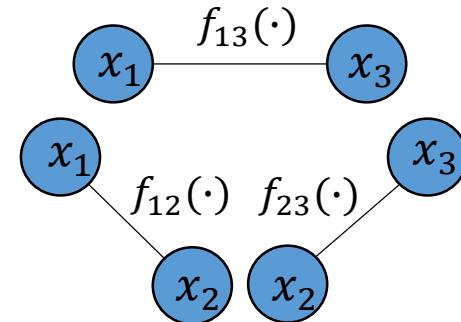
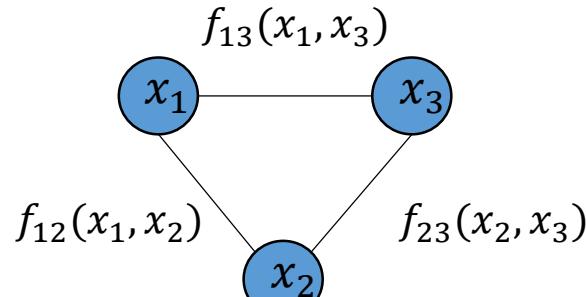
$$= 0 + 6$$

Modify the individual functions

- but -

keep the sum of functions the same

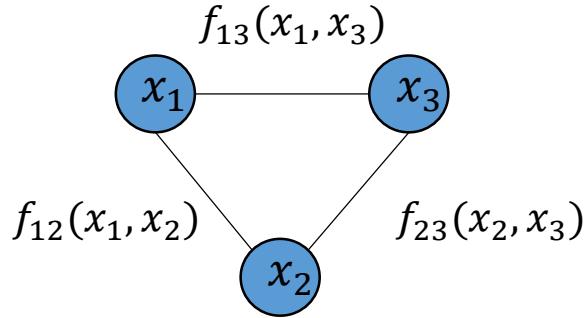
Dual Decomposition



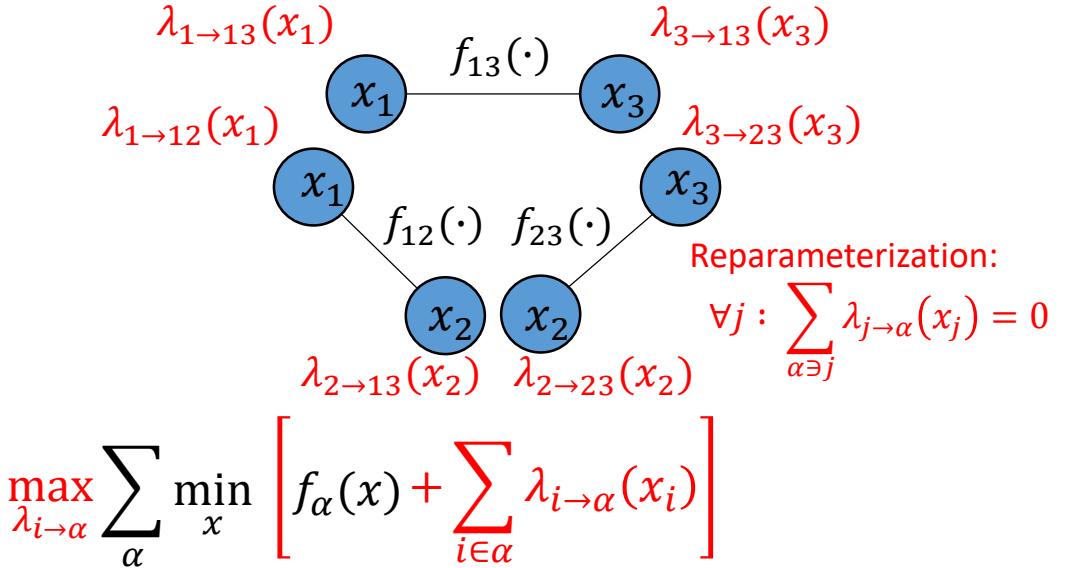
$$F^* = \min_x \sum_{\alpha} f_{\alpha}(x) \geq \sum_{\alpha} \min_x f_{\alpha}(x)$$

- Bound solution using decomposed optimization
- Solve independently: optimistic bound

Dual Decomposition

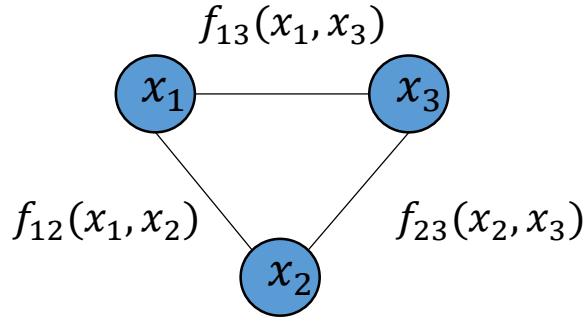


$$F^* = \min_x \sum_{\alpha} f_{\alpha}(x)$$

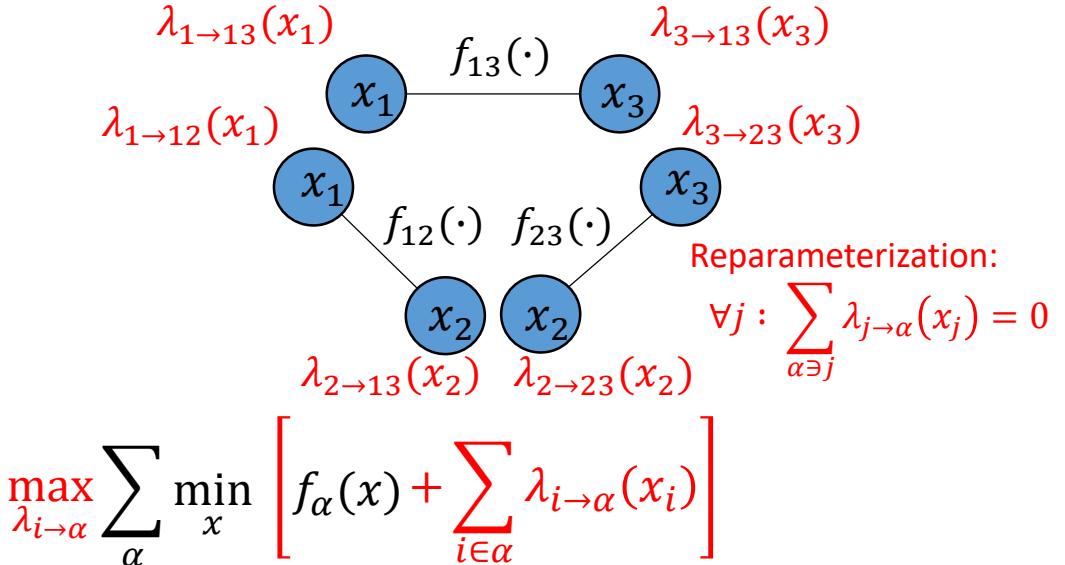


- Bound solution using decomposed optimization
- Solve independently: optimistic bound
- Tighten the bound by reparameterization
 - Enforce lost equality constraints via Lagrange multipliers

Dual Decomposition



$$F^* = \min_x \sum_{\alpha} f_{\alpha}(x)$$

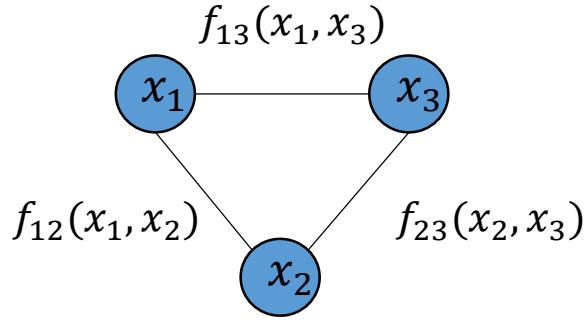


$$\geq \max_{\lambda_{i \rightarrow \alpha}} \sum_{\alpha} \min_x \left[f_{\alpha}(x) + \sum_{i \in \alpha} \lambda_{i \rightarrow \alpha}(x_i) \right]$$

Many names for the same class of bounds:

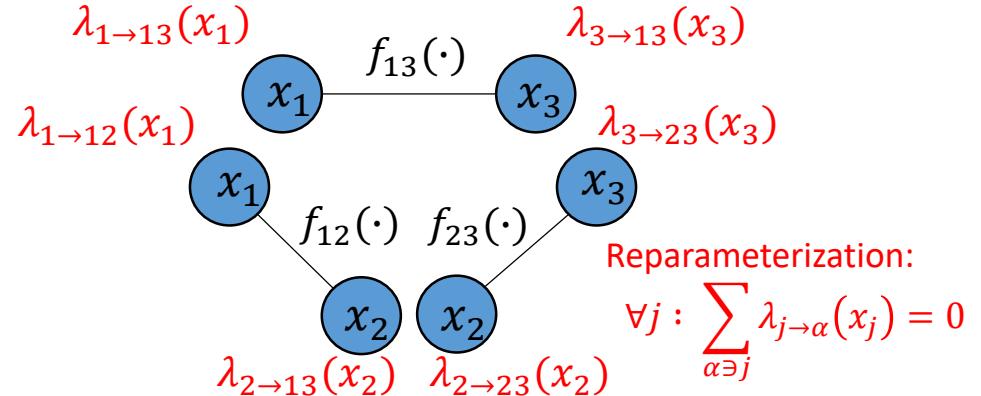
- Dual decomposition [Komodakis et al. 2007]
- TRW, MPLP [Wainwright et al. 2005; Globerson & Jaakkola, 2007]
- Soft arc consistency [Cooper & Schiex, 2004]
- Max-sum diffusion [Warner 2007]

Dual Decomposition



$$F^* = \min_x \sum_{\alpha} f_{\alpha}(x)$$

$$\geq \max_{\lambda_{i \rightarrow \alpha}} \sum_{\alpha} \min_x \left[f_{\alpha}(x) + \sum_{i \in \alpha} \lambda_{i \rightarrow \alpha}(x_i) \right]$$



Many ways to optimize the bound:

- Sub-gradient descent
- Coordinate descent
- Proximal optimization
- ADMM

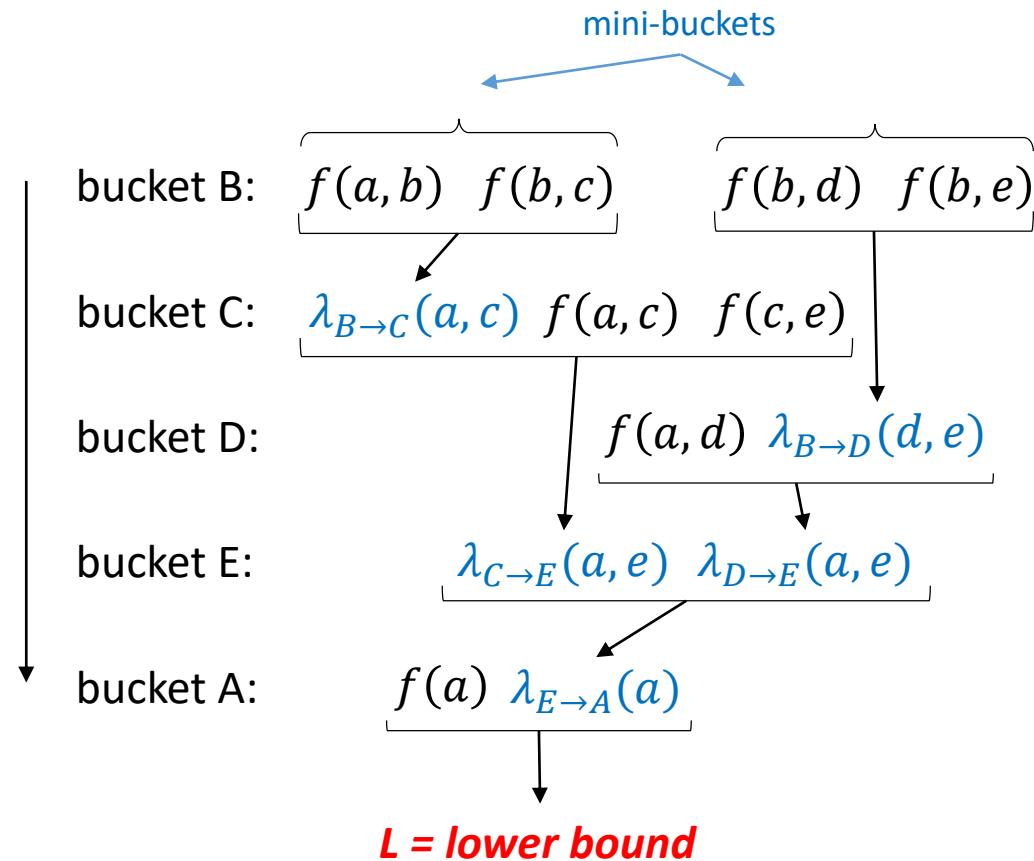
[Komodakis et al. 2007; Jojic et al. 2010]

[Warner 2007; Globerson & Jaakkola 2007; Sontag et al. 2009; Ihler et al. 2012]

[Ravikumar et al., 2010]

[Meshi & Globerson 2011; Martins et al. 2011; Forouzan & Ihler 2013]

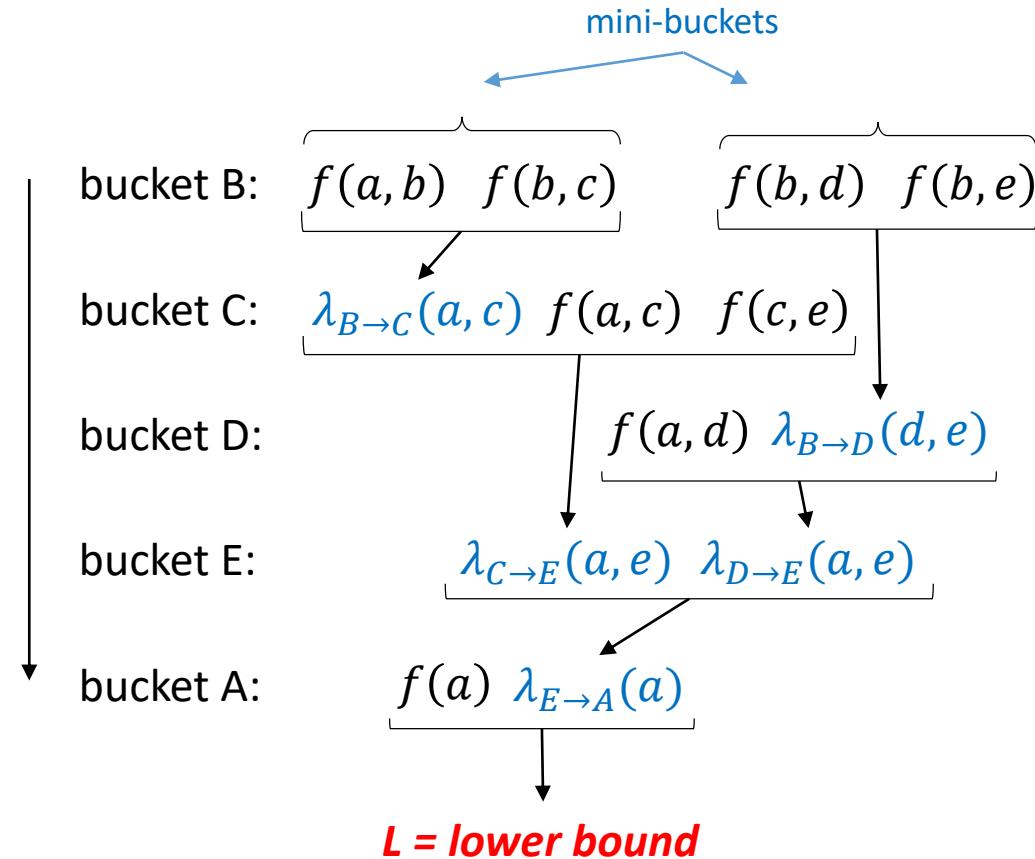
Mini-Bucket as Dual Decomposition



Mini-Bucket as Dual Decomposition

$$\min_{a,c,b} [f(a, b) + f(b, c) - \lambda_{B \rightarrow C}(a, c)] = 0$$

$$\min_{d,e,b} [f(b, d) + f(b, e) - \lambda_{B \rightarrow D}(d, e)] = 0$$

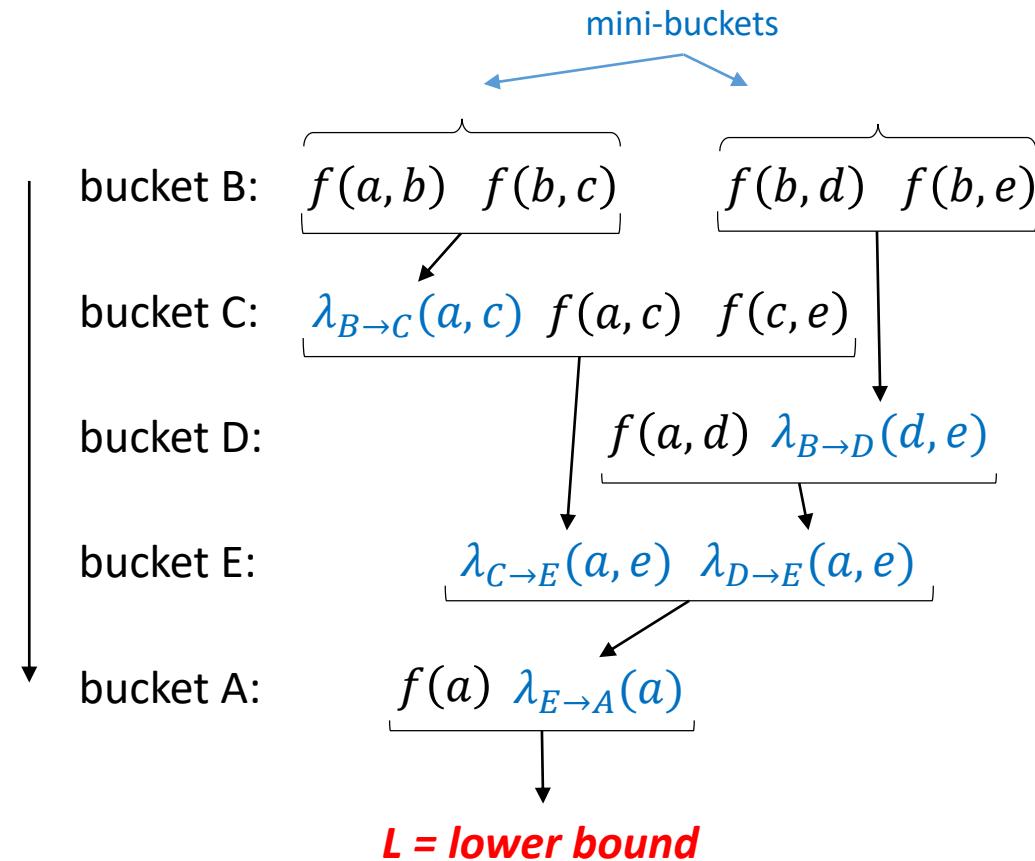


Mini-Bucket as Dual Decomposition

$$\min_{a,c,b} [f(a,b) + f(b,c) - \lambda_{B \rightarrow C}(a,c)] = 0$$

$$\min_{d,e,b} [f(b,d) + f(b,e) - \lambda_{B \rightarrow D}(d,e)] = 0$$

$$\min_{a,e,c} [\lambda_{B \rightarrow C}(a,c) + f(a,c) + f(c,e) - \lambda_{C \rightarrow E}(a,e)] = 0$$



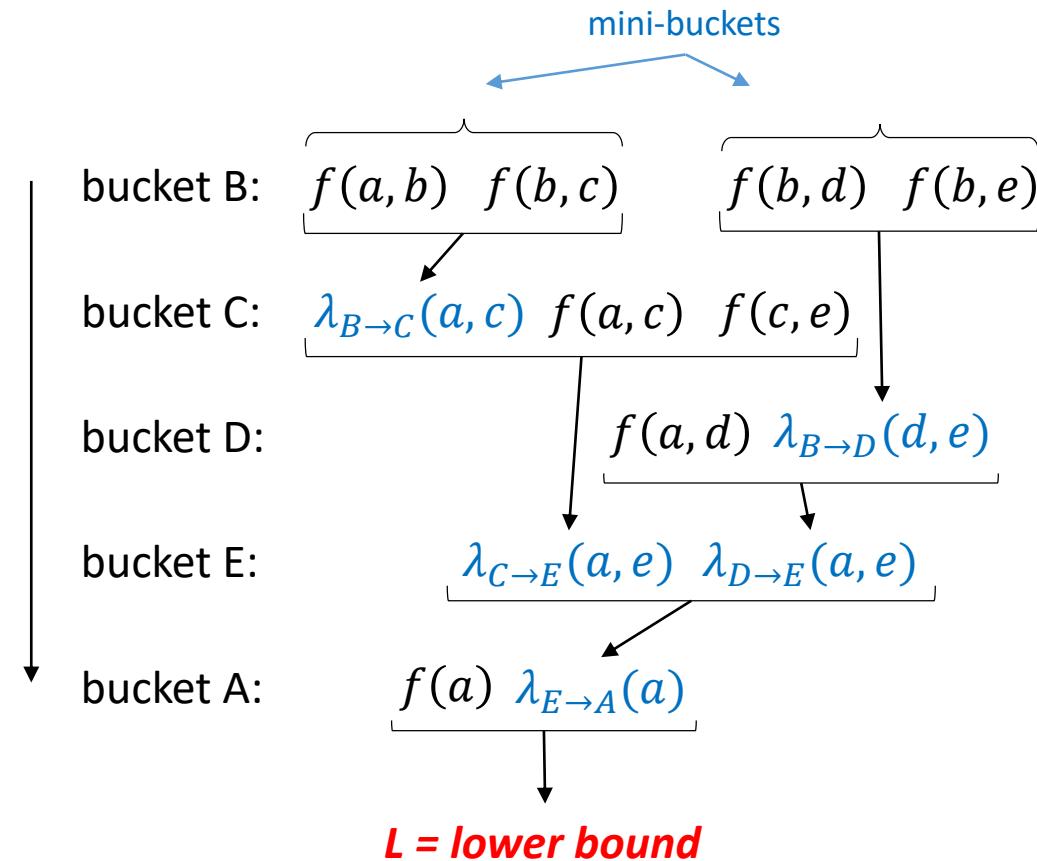
Mini-Bucket as Dual Decomposition

$$\min_{a,c,b} [f(a,b) + f(b,c) - \lambda_{B \rightarrow C}(a,c)] = 0$$

$$\min_{d,e,b} [f(b,d) + f(b,e) - \lambda_{B \rightarrow D}(d,e)] = 0$$

$$\min_{a,e,c} [\lambda_{B \rightarrow C}(a,c) + f(a,c) + f(c,e) - \lambda_{C \rightarrow E}(a,e)] = 0$$

$$\min_{a,d} [f(a,d) + \lambda_{B \rightarrow D}(d,e) - \lambda_{D \rightarrow E}(a,e)] = 0$$



Mini-Bucket as Dual Decomposition

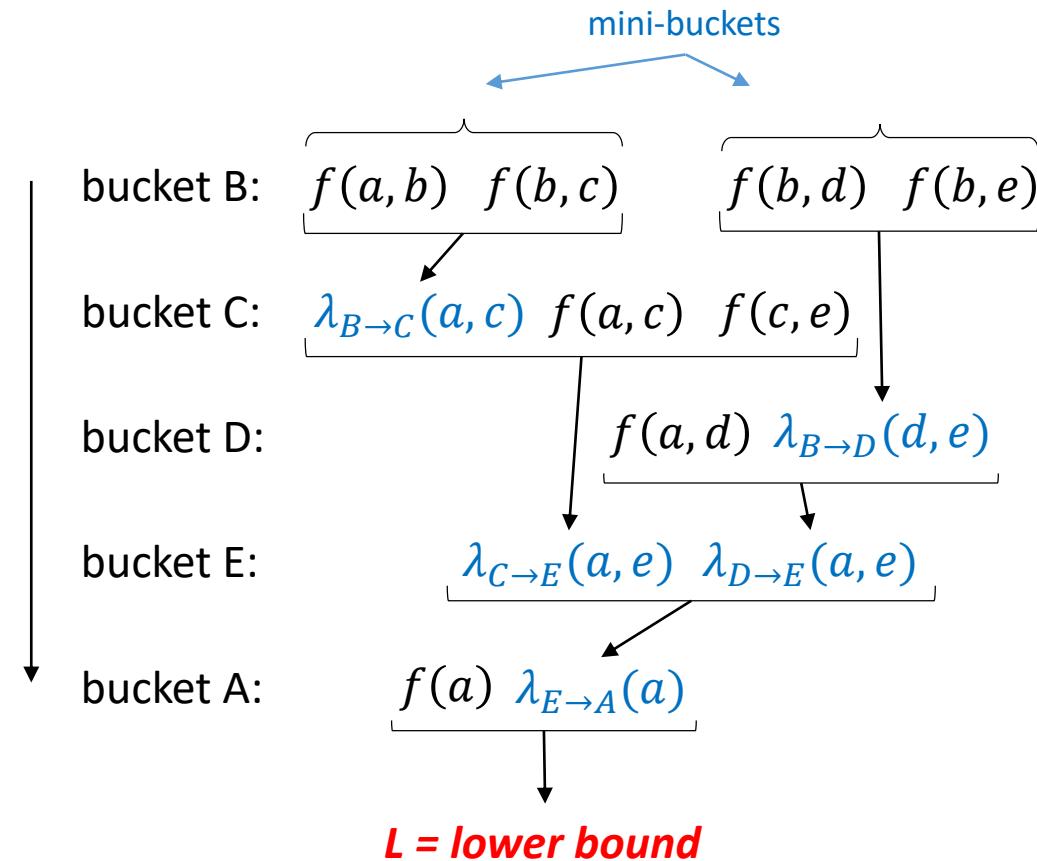
$$\min_{a,c,b} [f(a,b) + f(b,c) - \lambda_{B \rightarrow C}(a,c)] = 0$$

$$\min_{d,e,b} [f(b,d) + f(b,e) - \lambda_{B \rightarrow D}(d,e)] = 0$$

$$\min_{a,e,c} [\lambda_{B \rightarrow C}(a,c) + f(a,c) + f(c,e) - \lambda_{C \rightarrow E}(a,e)] = 0$$

$$\min_{a,d} [f(a,d) + \lambda_{B \rightarrow D}(d,e) - \lambda_{D \rightarrow E}(a,e)] = 0$$

$$\min_{a,e} [\lambda_{C \rightarrow E}(a,e) + \lambda_{D \rightarrow E}(a,e) - \lambda_{E \rightarrow A}(a)] = 0$$



Mini-Bucket as Dual Decomposition

$$\min_{a,c,b} [f(a,b) + f(b,c) - \lambda_{B \rightarrow C}(a,c)] = 0$$

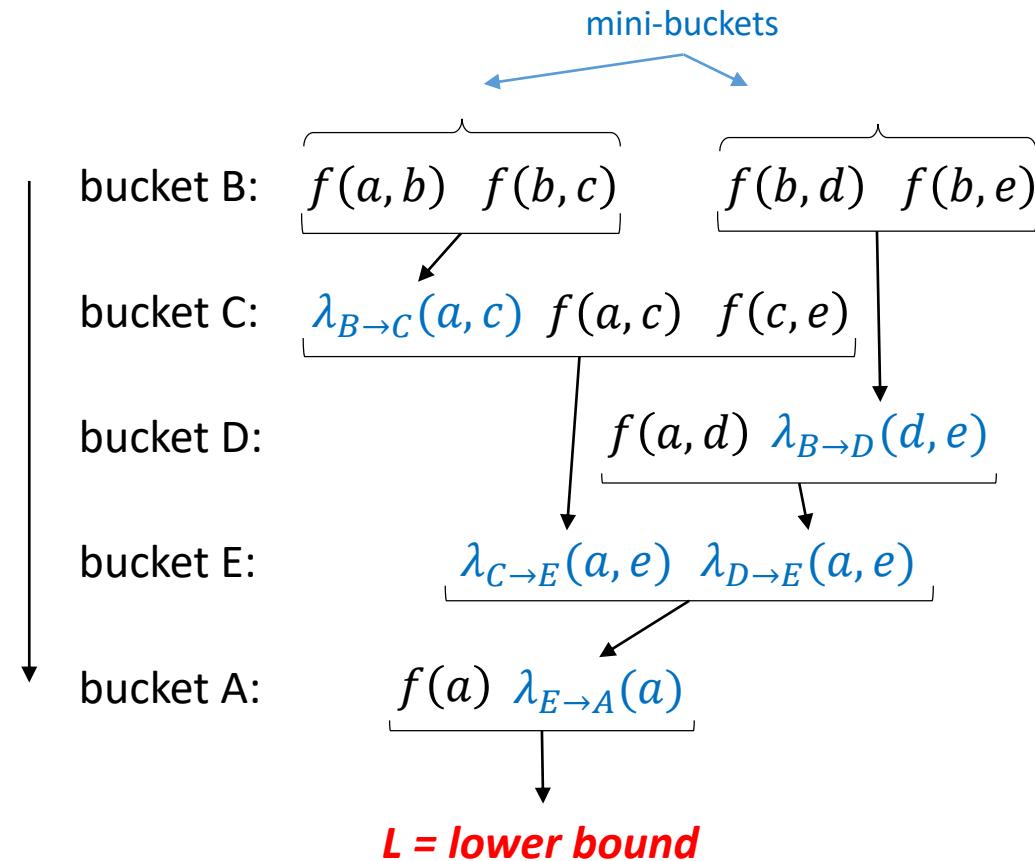
$$\min_{d,e,b} [f(b,d) + f(b,e) - \lambda_{B \rightarrow D}(d,e)] = 0$$

$$\min_{a,e,c} [\lambda_{B \rightarrow C}(a,c) + f(a,c) + f(c,e) - \lambda_{C \rightarrow E}(a,e)] = 0$$

$$\min_{a,d} [f(a,d) + \lambda_{B \rightarrow D}(d,e) - \lambda_{D \rightarrow E}(a,e)] = 0$$

$$\min_{a,e} [\lambda_{C \rightarrow E}(a,e) + \lambda_{D \rightarrow E}(a,e) - \lambda_{E \rightarrow A}(a)] = 0$$

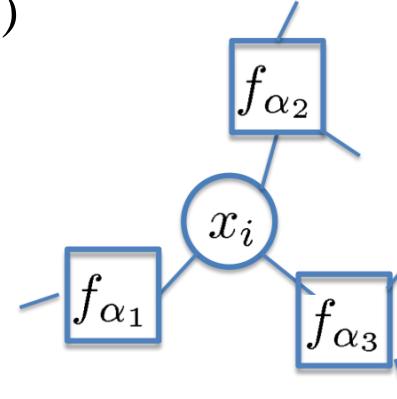
$$\min_a [f(a) + \lambda_{E \rightarrow A}(a)] = L$$



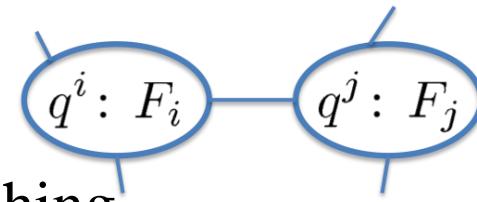
Various Update Schemes

- Can use any decomposition updates
 - (message passing, subgradient, augmented, etc.)

- **FGLP:** Update the original factors



- **JGLP:** Update clique function of the join graph



- **MBE-MM:** Mini-bucket with moment matching
 - Apply cost-shifting within each bucket only

FGLP (Factor Graph Linear Programming)

Algorithm 1: Factor graph LP (FGLP)

[1] **Input:** Graphical Model $\langle X, D, F, \sum \rangle$, where f_α is a function defined on variables X_α

Output: Re-parameterized factors F' , upper bound on optimum value.

1. Iterate until convergence:

2. For each variable X_i do:

3. Let $F_i = \{\alpha : i \in \alpha\}$ with X_i in their scope

4. $\forall \alpha \in F_i$, compute max-marginals: $\lambda_\alpha(X_i) = \max_{X_\alpha} f_\alpha(X_\alpha)$

5. $\forall \alpha \in F_i$, update parameterization: $f_\alpha(X_\alpha) \leftarrow f_\alpha(X_\alpha) - \lambda_\alpha(X_i) + \frac{1}{|F_i|} \sum_{\beta \in F_i} \lambda_\beta(X_i)$

6. **Return:** Reparameterized factors F' and bound $\sum_{\alpha \in F} \max_X f_\alpha(X_\alpha)$

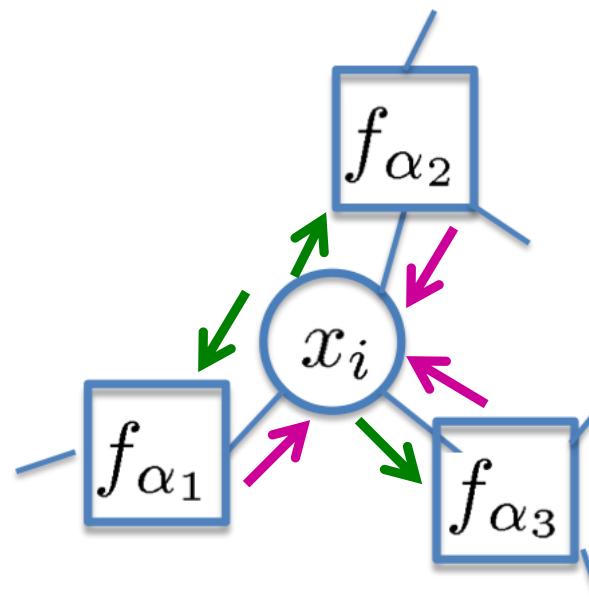
$$C^* \leq \min_{\lambda \in \Lambda} \sum_{(ij) \in F} \max_X (f_{ij}(X_i, X_j) + \lambda_{ij}(X_i) + \lambda_{ji}(X_j)).$$

The new functions $\tilde{f}_{ij} = f_{ij}(X_i, X_j) + \lambda_{ij}(X_i) + \lambda_{ji}(X_j)$ define a *re-parameterization* or cost-shifting of the original distribution. The λ_{ij} can be interpreted in various ways [81, 77] and they can be tightened by different methods [38, 75]. MPLP [25], soft-arc consistency, and many other LP algorithms operate directly

Factor graph Linear Programming

- Update the original factors (FGLP)
 - Tighten all factors over x_i simultaneously
 - Compute **max-marginals** $\forall \alpha, \gamma_\alpha(x_i) = \max_{x_\alpha \setminus x_i} f_\alpha$
 - & **update**:

$$\forall \alpha, f_\alpha(x_\alpha) \leftarrow f_\alpha(x_\alpha) - \gamma_\alpha(x_i) + \frac{1}{|F_i|} \sum_{\beta} \gamma_\beta(x_i)$$



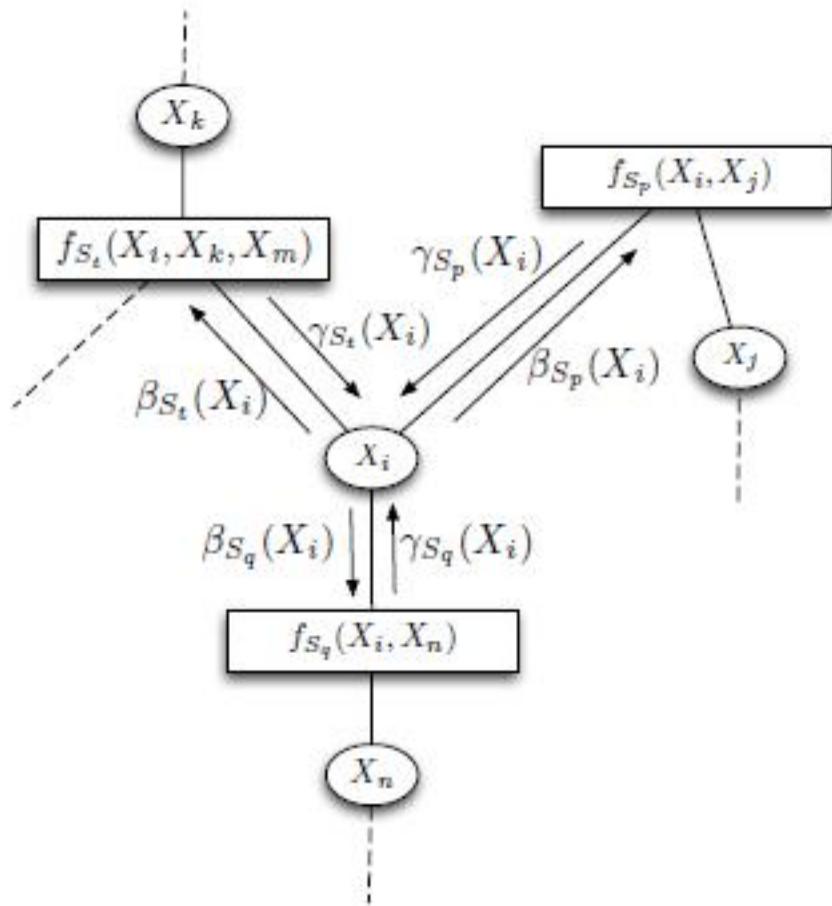


Figure 5.2: FGLP example: local messages involving variable X_i .

Factor Graph Linear Programming

Algorithm 23: Factor Graph Linear Programming (FGLP, based on [40])

Input: A graphical model $M = (\mathbf{X}, \mathbf{D}, \mathbf{F}, \sum)$, variable ordering σ

Output: Upper bound on the optimum value of MPE cost

```
1 while NOT converged do
2   for each variable  $X_i$  do
3     Get factors  $F_i = f_{\mathbf{S}_k} : X_i \in \mathbf{S}_k$  with  $X_i$  in their scope;
4     //for each function compute max-marginals  $\gamma$  marginalizing out all variables except for  $X_i$ :
5      $\forall f_{\mathbf{S}_k} \gamma_{\mathbf{S}_k}(X_i) = \max_{\mathbf{S}_k \setminus X_i} f_{\mathbf{S}_k};$ 
6     // compute messages  $\beta_{\mathbf{S}_k}(X_i)$  from  $X_i$  back to a function  $f_{\mathbf{S}_k}$  correcting for the function's
7     own max-marginal  $\gamma_{\mathbf{S}_k}$ :
8      $\forall f_{\mathbf{S}_k} \beta_{\mathbf{S}_k} = \frac{1}{|F_i|} \sum_{\{\mathbf{S}_j | f_{\mathbf{S}_j} \in F_i\}} \gamma_{\mathbf{S}_j}(X_i) - \gamma_{\mathbf{S}_k}(X_i)$ 
9     // update (re-parametrize) each function:
10     $\forall f_{\mathbf{S}_k}, f_{\mathbf{S}_k} \leftarrow f_{\mathbf{S}_k} + \beta_{\mathbf{S}_k};$ 
```

The update messages from variable X_i back to the functions are:

$$\beta_{\mathbf{S}_q}(X_i) = \frac{1}{3}(\gamma_{\mathbf{S}_q}(X_i) + \gamma_{\mathbf{S}_p}(X_i) + \gamma_{\mathbf{S}_t}(X_i)) - \gamma_{\mathbf{S}_q}(X_i)$$

$$\beta_{\mathbf{S}_p}(X_i) = \frac{1}{3}(\gamma_{\mathbf{S}_q}(X_i) + \gamma_{\mathbf{S}_p}(X_i) + \gamma_{\mathbf{S}_t}(X_i)) - \gamma_{\mathbf{S}_p}(X_i)$$

$$\beta_{\mathbf{S}_t}(X_i) = \frac{1}{3}(\gamma_{\mathbf{S}_q}(X_i) + \gamma_{\mathbf{S}_p}(X_i) + \gamma_{\mathbf{S}_t}(X_i)) - \gamma_{\mathbf{S}_t}(X_i)$$

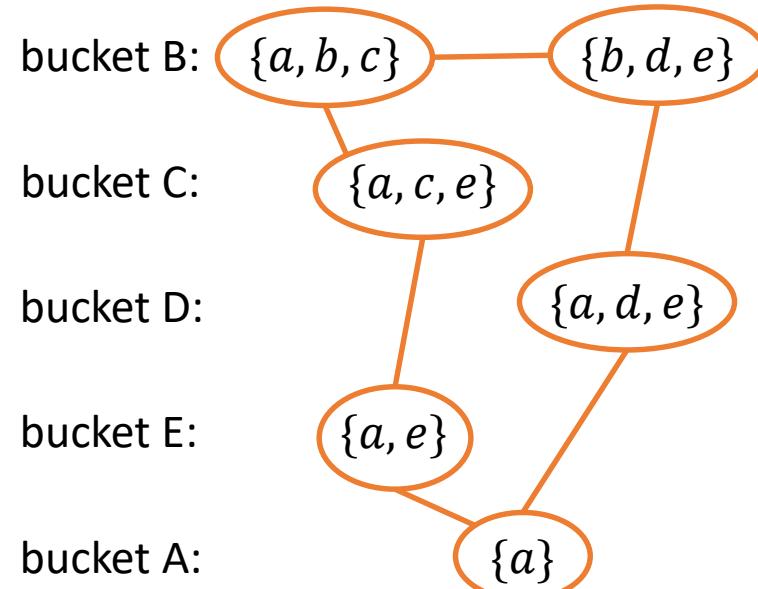
Complexity of FGLP

Theorem 5.1 (Complexity of FGLP). *The total time complexity of a single iteration of FGLP is $O(n \cdot Q \cdot k^{Sc})$, where n is the number of variables in the problem, k is the largest domain size, $|F|$ is the number of functions, Sc bounds the largest scope of the original functions, Q is the largest number of functions having the same variable X_j in their scopes. The space complexity is $O(|F| \cdot k^{Sc})$.*

Mini-Bucket as Dual Decomposition

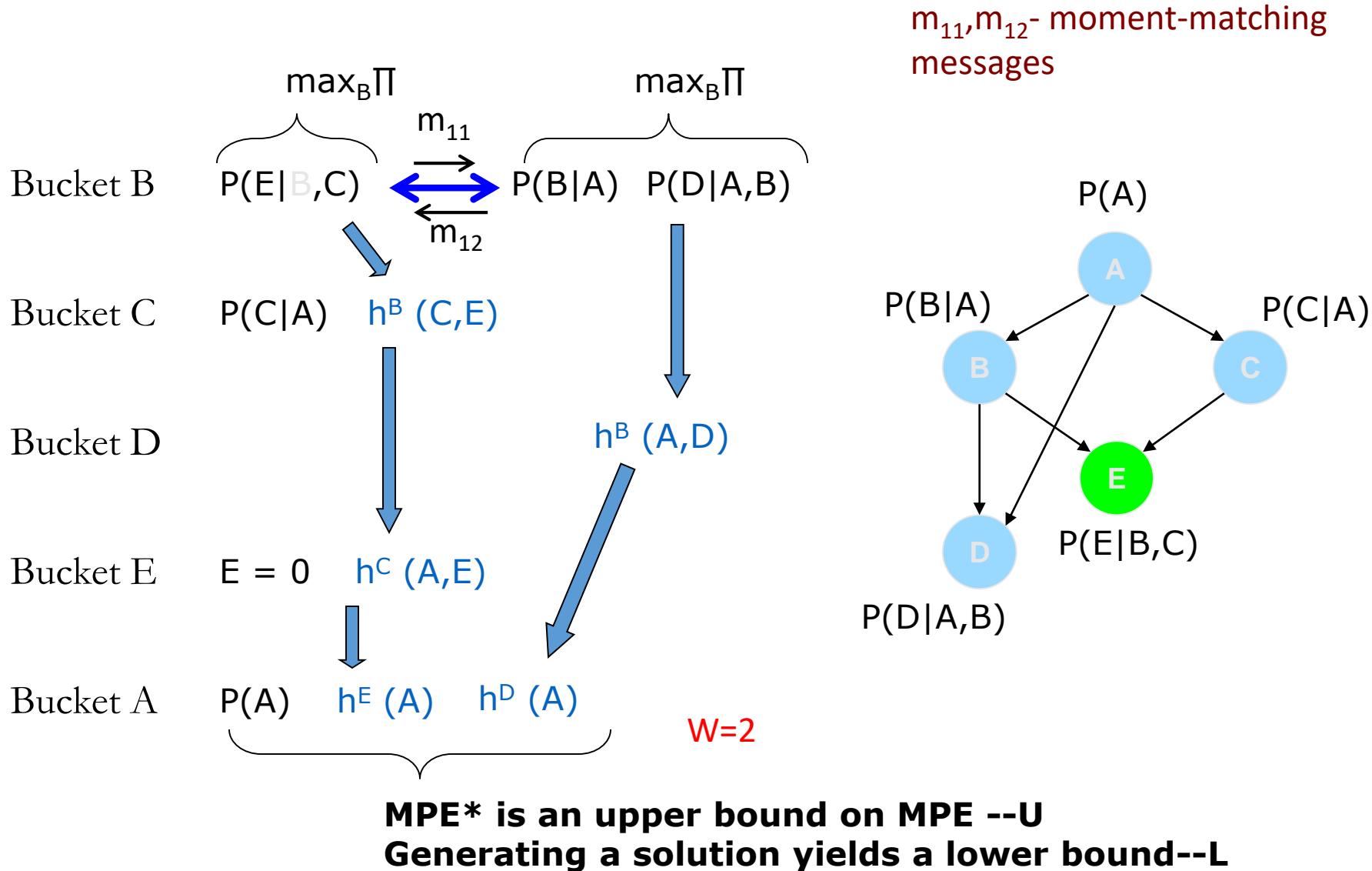
- Downward pass as cost-shifting
- Can also do cost-shifting within mini-buckets
- “Join graph” message passing
- “Moment matching” version: one message update within each bucket during downward sweep

Join graph:



L = lower bound

MBE-MM: MBE with moment matching



MBE-MM (MBE with Moment-Matching)

Algorithm 26: Algorithm MBE-MM

Input: A graphical model $\mathcal{M} = \langle \mathbf{X}, \mathbf{D}, \mathbf{F}, \sum \rangle$, variable order $o = \{X_1, \dots, X_n\}$, i-bound parameter i

Output: Upper bound on the optimum value of MPE cost

//Initialize:

- 1 Partition the functions in \mathbf{F} into B_{X_1}, \dots, B_{X_n} , where B_{X_k} contains all functions f_j whose highest variable is X_k ;
//processing bucket B_{X_k}
 - 2 **for** $k \leftarrow n$ down to 1 **do**
 - 3 Partition functions g (both original and messages generated in previous buckets) in B_{X_k} into the mini-buckets defined $Q_{X_k} = \{q_k^1, \dots, q_k^t\}$, where each q_k^p has no more than $i + 1$ variables;
 - 4 Find the set of variables common to all the mini-buckets of variable X_k :
 $S_k = \text{Scope}(q_k^1) \cap \dots \cap \text{Scope}(q_k^t)$;
Find the function of each mini-bucket
 $q_k^p: F_k^p \leftarrow \prod_{g \in q_k^p} g$;
 - 5 Find the max-marginals of each mini-bucket
 $q_k^p: \gamma_{X_k}^p = \max_{\text{Scope}(q_k^p) \setminus S_k} (F_k^p)$;
 - 6 Update functions of each mini-bucket
 $F_k^p \leftarrow F_k^p - \gamma_{X_k}^p + \frac{1}{t} \sum_{j=1}^t \gamma_{X_k}^j$;
 - 7 Generate messages $h_{X_k \rightarrow X_m}^p = \max_{X_k} F_k^p$ and place each in the bucket of highest in the ordering o variable X_m in $\text{Scope}(q_k^p)$;
 - 8 **return** All the buckets and the cost bound from B_1 ;
-

Theorem 5.3 (Complexity of MBE-MM). *Given a problem with n variables having domain of size k and an i-bound i , the worst-case time complexity of MBE-MM is $O(n \cdot Q \cdot k^{i+1})$ and its space complexity is $O(n \cdot k^i)$, where Q bounds the number of functions having the same variable X_i in their scopes.*

Anytime Approximation

anytime - mpe(ε)

Initialize : $i = i_0$

While time and space resources are available

$$i \leftarrow i + i_{step}$$

$U \leftarrow$ upper bound computed by $approx\text{-}mpe(i)$

$L \leftarrow$ lower bound computed by $approx\text{-}mpe(i)$

keep the best solution found so far

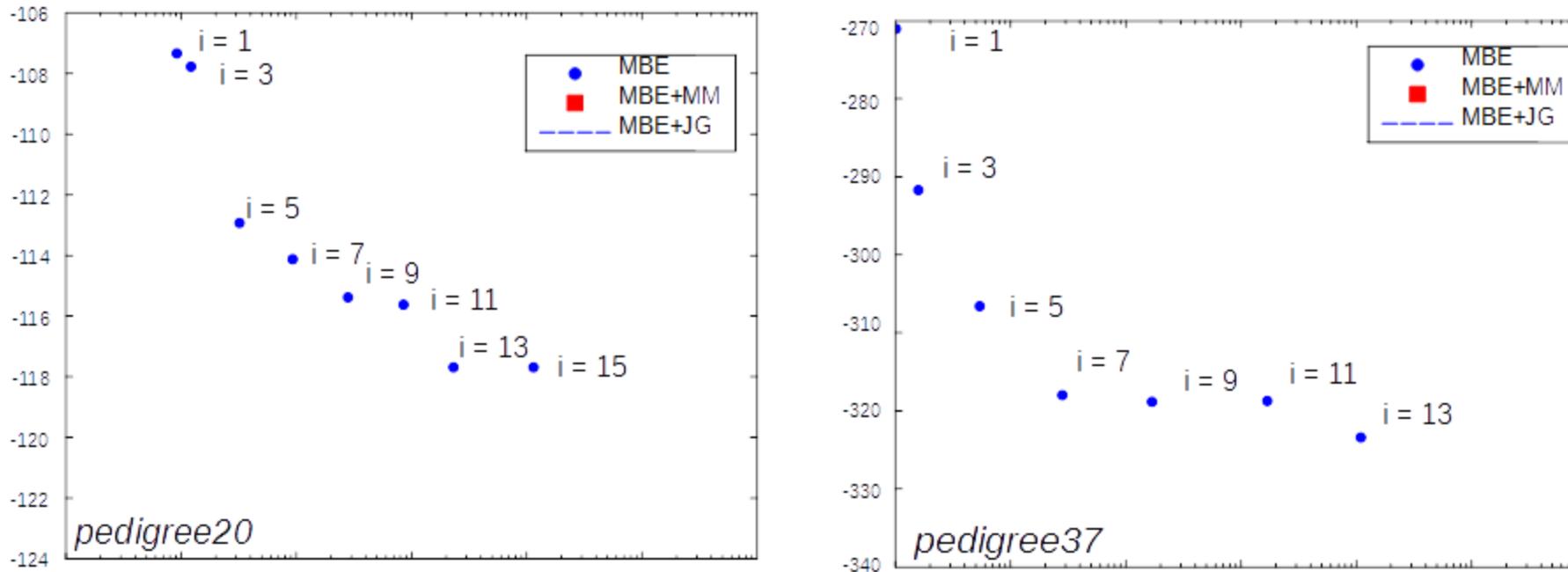
if $1 \leq \frac{U}{L} \leq 1 + \varepsilon$, return solution

end

return the largest L and the smallest U

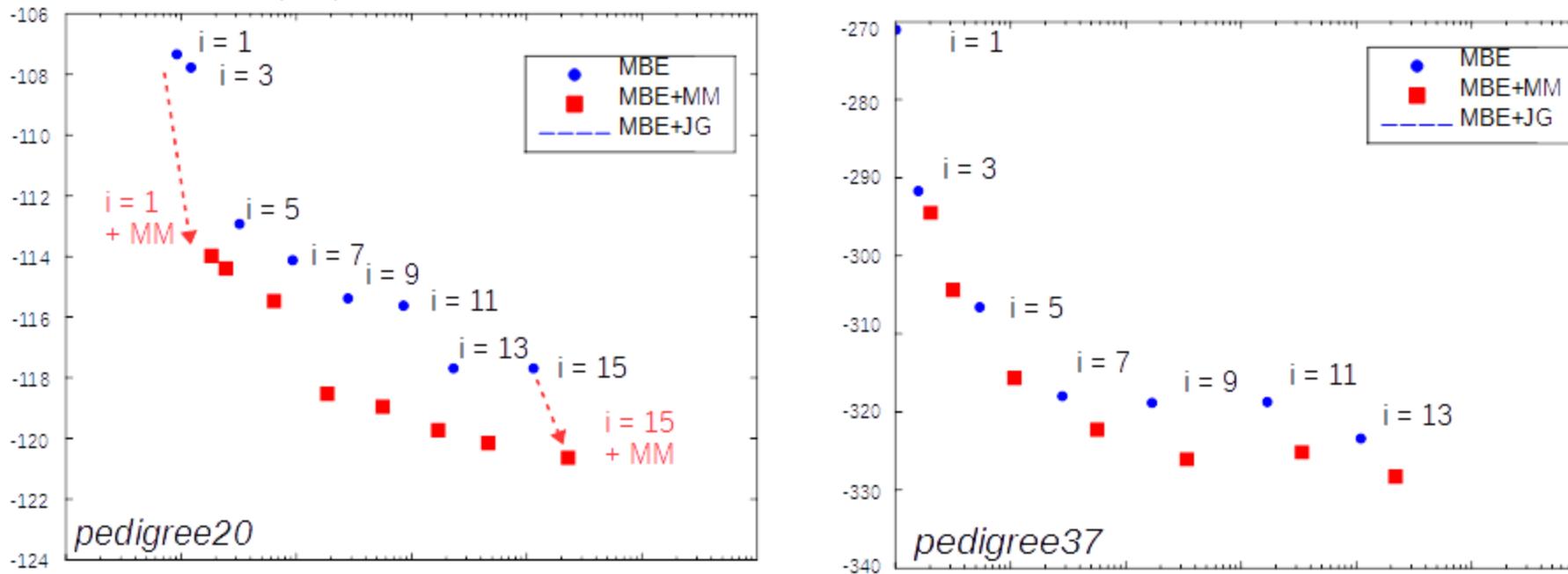
[Dechter and Rish, 2003]

Anytime Approximation



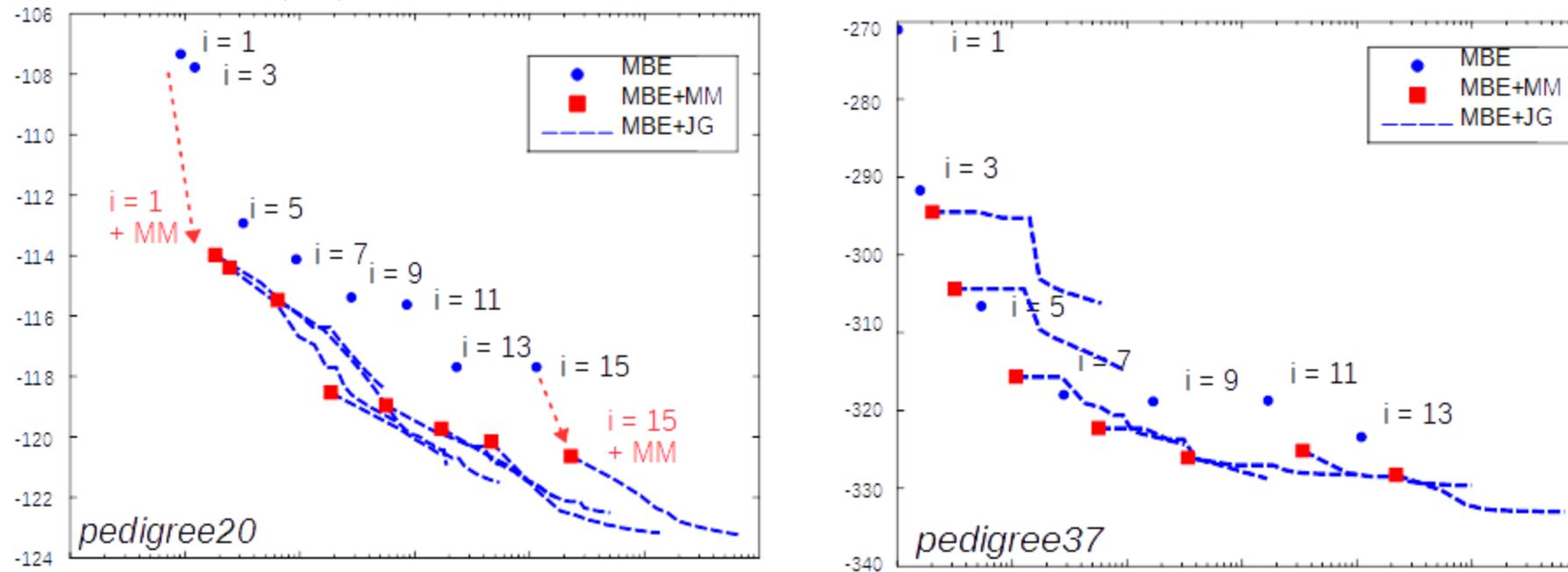
- Can tighten the bound in various ways
 - Cost-shifting (improve consistency between cliques)
 - Increase i-bound (higher order consistency)
- Simple moment-matching step improves bound significantly

Anytime Approximation



- Can tighten the bound in various ways
 - Cost-shifting (improve consistency between cliques)
 - Increase i-bound (higher order consistency)
- Simple moment-matching step improves bound significantly

Anytime Approximation



- Can tighten the bound in various ways
 - Cost-shifting (improve consistency between cliques)
 - Increase i-bound (higher order consistency)
- Simple moment-matching step improves bound significantly

Weighted Mini-Bucket

(for summation)

Exact bucket elimination:

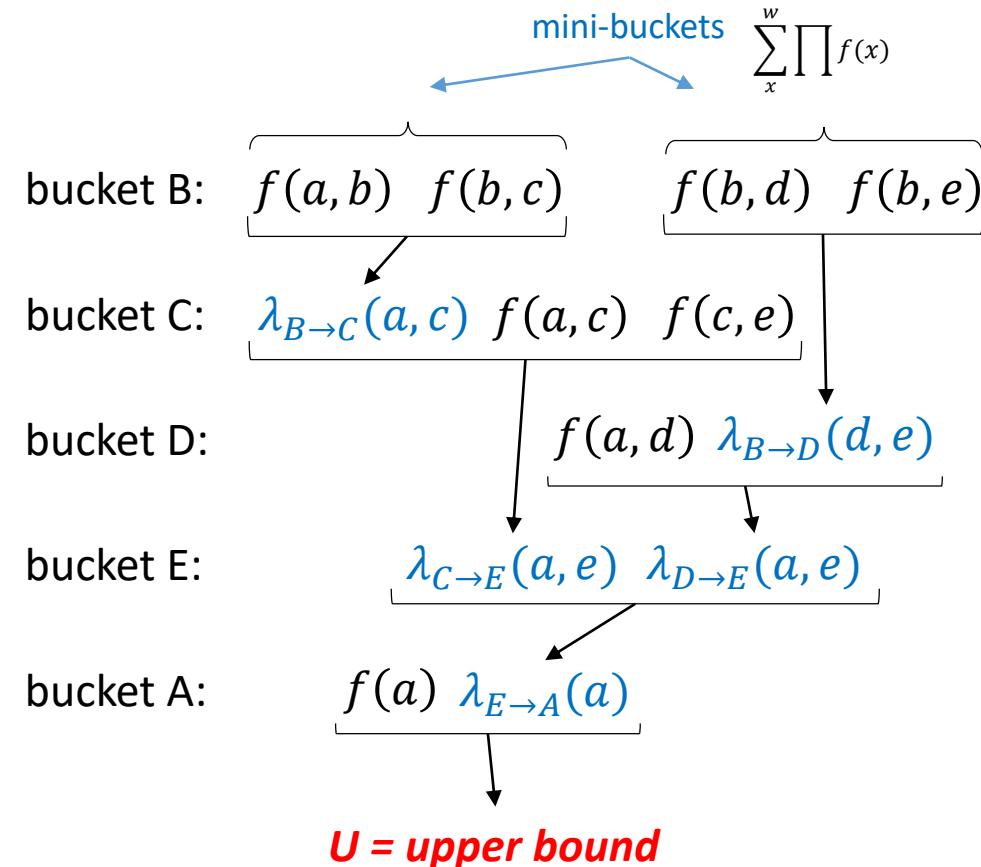
$$\begin{aligned}\lambda_B(a, c, d, e) &= \sum_b [f(a, b) \cdot f(b, c) \cdot f(b, d) \cdot f(b, e)] \\ &\leq \left[\sum_b^{w_1} f(a, b) f(b, c) \right] \cdot \left[\sum_b^{w_2} f(b, d) f(b, e) \right] \\ &= \lambda_{B \rightarrow C}(a, c) \cdot \lambda_{B \rightarrow D}(d, e)\end{aligned}$$

(mini-buckets)

where $\sum_x^w f(x) = \left[\sum_x f(x)^{\frac{1}{w}} \right]^w$
is the weighted or “power” sum operator

$$\sum_x^w f_1(x) f_2(x) \leq \left[\sum_x^{w_1} f_1(x) \right] \left[\sum_x^{w_2} f_2(x) \right]$$

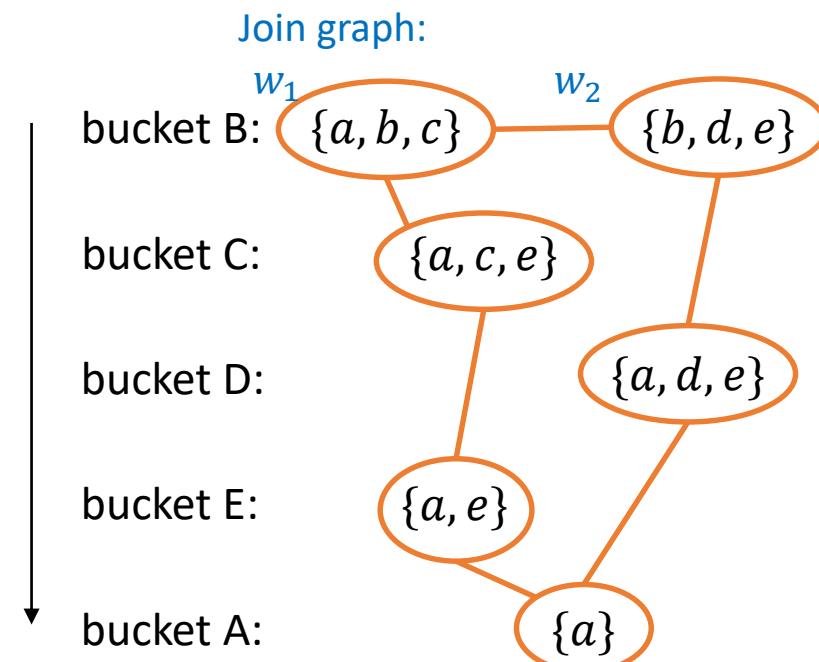
where $w_1 + w_2 = w$ and $w_1 > 0, w_2 > 0$
(lower bound if $w_1 > 0, w_2 < 0$)



[Liu and Ihler, 2011]

Weighted Mini-Bucket

- Related to conditional entropy decomposition but, with an efficient “primal” bound form
- Can optimize the bound over:
 - Cost-shifting
 - Weights
- Again, involves message passing over JG
- Similar, one-pass “moment-matching” variant



$U = \text{upper bound}$

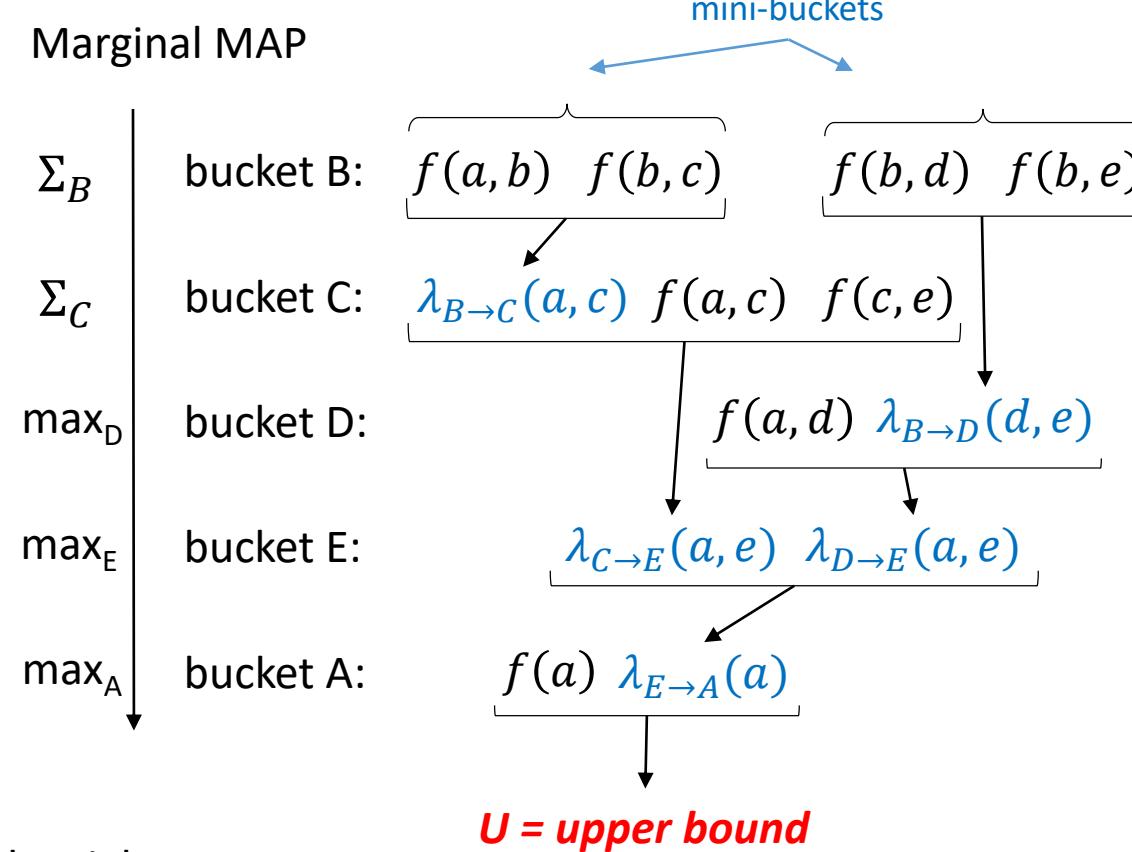
[Liu and Ihler, 2011]

MB and WMB for Marginal MAP

$$\max_Y \sum_{X \setminus Y} \prod_j P_j$$

$$\begin{aligned}\lambda_{B \rightarrow C}(a, c) &= \sum_b^{w_1} f(a, b)f(b, c) \\ \lambda_{B \rightarrow D}(d, e) &= \sum_b^{w_2} f(b, d)f(b, e) \\ &\vdots \\ \lambda_{E \rightarrow A}(a) &= \max_e \lambda_{C \rightarrow E}(a, e)\lambda_{D \rightarrow E}(a, e)\end{aligned}$$

$$U = \max_a f(a)\lambda_{E \rightarrow A}(a)$$



Can optimize over cost-shifting and weights
(single pass “MM” or iterative message passing)

U = upper bound

[Liu and Ihler, 2011; 2013]
[Dechter and Rish, 2003]