
Subproblem-Tree Calibration: A Unified Approach to Max-Product Message Passing

Huayan Wang

HUAYANW@CS.STANFORD.EDU

Computer Science Department, Stanford University, Palo Alto, CA 94305 USA

Daphne Koller

KOLLER@CS.STANFORD.EDU

Computer Science Department, Stanford University, Palo Alto, CA 94305 USA

Abstract

Max-product (max-sum) message passing algorithms are widely used for MAP inference in MRFs. It has many variants sharing a common flavor of passing “messages” over some graph-object. Recent advances revealed that its convergent versions (such as MPLP, MSD, TRW-S) can be viewed as performing block coordinate descent (BCD) in a dual objective. That is, each BCD step achieves dual-optimal w.r.t. a block of dual variables (messages), thereby decreases the dual objective monotonically. However, most existing algorithms are limited to updating blocks selected in rather restricted ways. In this paper, we show a “unified” message passing algorithm that: (a) subsumes MPLP, MSD, and TRW-S as special cases when applied to their respective choices of dual objective and blocks, and (b) is able to perform BCD under much more flexible choices of blocks (including very large blocks) as well as the dual objective itself (that arise from an arbitrary dual decomposition).

1. Introduction

MAP-MRF (finding the most probable assignments for MRFs) is one of the most important components in learning and applying structured probabilistic models. In general this problem is NP-hard (Shimony, 1994). Many different methods have been proposed to approximate or solve it under specific circumstances. A large family of these methods is based on solving

a dual problem of an LP relaxation. Different duals of different LP relaxations (with different tightness) have been used. These methods are usually formulated as max-product (max-sum) message passing over the MRF or its cluster (region) graph.

Recent advances revealed that convergent versions of these algorithms can often be interpreted as block coordinate descent (BCD) in the dual (Meltzer et al., 2009; Sontag et al., 2011). However, most of them operate on local (small) blocks, such as MPLP (Globerson & Jaakkola, 2007) and its generalizations (Sontag et al., 2008), max-sum diffusion (MSD) (Werner, 2007), and TRW-S (Kolmogorov, 2006). Given a block of dual variables (messages), these algorithms work by enforcing some *consistency constraint* over the block, thereby achieve dual-optimal w.r.t. these dual variables. Meltzer et al. (2009) noted that it was difficult to generalize them to larger blocks while enforcing the same consistency constraints. Sontag et al. (2009) proposed a method (tree-BCD) that updates much larger blocks for a specific choice of dual. But it was not clear how to apply it to duals of tighter LP relaxations (e.g., a dual decomposition with cycle subproblems).

We observe that the difficulties in generalizing these methods arise from the fact that they all impose too strong consistency constraints, which are *sufficient but not necessary* for the dual objective to be optimal on the blocks being updated. By losing these constraints, we are able to perform BCD on much larger blocks.

Indeed, we show that dual-optimality on blocks (of messages) can be established on a much broader basis—with quite general choices of the dual objective itself as well as the blocks to be updated. Specifically, we illustrate this by deriving a “unified” message passing algorithm in the most general setup—an

arbitrary dual decomposition. The resulted algorithm (*subproblem-tree calibration*, or STC) has the following properties:

1. It is formulated as message passing on a graph-object (*subproblem multi-graph*, or SMG) that generalizes traditional cluster (region) graphs.
2. It subsumes MPLP, MSD and TRW-S as special cases when applied to their respective choices of dual objectives and blocks. We will precisely characterize the choices made by these existing methods under our framework.
3. It achieves dual-optimal on blocks that can be chosen in an extremely flexible manner, including very large blocks. For example, we can choose a block that spans all subproblems in an arbitrary dual decomposition.

In other words, our framework attempts to characterize the *degrees of freedom* we have in designing a message passing algorithm (for MAP inference). Understanding these flexibilities could help us better understand existing methods, as well as design more powerful ones. In practice, we observe that we often get stuck in sub-optimal dual states when only updating blocks chosen in a restricted manner, whereas being able to choose blocks with more flexibility could lead to better dual (and primal) states.

2. Background

2.1. MAP inference, LP relaxation, and dual decomposition

The MAP inference problem over $\mathbf{X} = \{X_{1:N}\}$ and graph structure $G = \{V, E\}$ can be formulated as:

$$\underset{\mathbf{X}}{\text{maximize}} \quad \Theta(\mathbf{X}) \quad (1)$$

where $\Theta(\mathbf{X}) = \sum_{\alpha \in \mathcal{A}} \theta_{\alpha}(X_{\alpha})$; \mathcal{A} is the set of MRF cliques. Without loss of generality we choose $\mathcal{A} = V \cup E$ for a parametrization with unary and pairwise potentials¹. We use lowercase $x_i \in \text{Val}(X_i)$ and $\mathbf{x} = \{x_{1:N}\}$ to denote assignments to the variables. Problem (1) is NP-hard in general (Shimony, 1994).

A large family of MAP inference methods builds on solving a linear programming (LP) relaxation (Wainwright & Jordan, 2008; Koller & Friedman, 2009):

$$\underset{\mu \in \mathcal{M}}{\text{maximize}} \quad \Theta \cdot \mu \quad (2)$$

¹A high-order potential can be converted into unary and pairwise potentials by introducing auxiliary variables.

where $\mu = \{\mu_i(x_i), \mu_{ij}(x_i, x_j) \mid \forall i, x_i, (i, j), (x_i, x_j)\}$; Θ is all MRF parameters $\{\theta_i, \theta_{ij}\}$ concatenated in same ordering as μ . Choosing different polytopes \mathcal{M} results in different LP relaxations. One choice is the *marginal polytope*:

$$\begin{aligned} \mathcal{M}_G = \{ \mu \mid \exists p(\mathbf{X}), \sum_{\mathbf{x} \in \text{Val}(\mathbf{X})} p(\mathbf{x}) = 1, p(\mathbf{x}) \geq 0, \\ p(X_i) = \mu_i, p(X_i, X_j) = \mu_{ij}, \forall i, (i, j) \} \end{aligned} \quad (3)$$

where $p(X_i)$ and $p(X_i, X_j)$ denote marginal distributions of $p(\mathbf{X})$. When $\mathcal{M} = \mathcal{M}_G$ the LP relaxation (2) is equivalent to the original problem (1). Therefore it is also NP-hard—the marginal polytope is defined by exponentially many faces (constraints).

A widely used choice for \mathcal{M} is the *local polytope*:

$$\begin{aligned} \mathcal{M}_L = \{ \mu \mid \mu_{ij} \geq 0, \sum_{x_j} \mu_{ij}(X_i, x_j) = \mu_i(X_i), \\ \sum_{x_i} \mu_i(x_i) = 1, \forall i \in V, (i, j) \in E \} \end{aligned} \quad (4)$$

which only has a polynomial number of constraints, and the LP relaxation (2) with \mathcal{M}_L is tractable. In general \mathcal{M}_L is a loose outer bound of \mathcal{M}_G . Any solution to (2) with $\mathcal{M} = \mathcal{M}_L$ is a vertex of \mathcal{M}_L . If it happens to be a vertex of \mathcal{M}_G , we have the exact solution to (1). This is usually not the case except for several special families such as the problems with tree structures or sub-modular potentials etc.

Solving (2) using a standard LP solver is very inefficient (Yanover et al., 2006) even with a tractable choice of \mathcal{M} . In practice we usually solve its dual LP. Indeed, it is more convenient to start with a *dual decomposition*, which directly gives us a dual LP (of some underlying primal LP that we do not explicitly deal with).

Specifically, consider a decomposition of $\Theta(\mathbf{X})$ into subproblems $c \in \mathcal{C}$, parametrized by $\{\Theta^c\}$, which altogether give rise to a reparametrization of $\Theta(\mathbf{x})$:

$$\forall \mathbf{x}, \sum_{c \in \mathcal{C}} \Theta^c(\mathbf{x}|_c) = \Theta(\mathbf{x}) \quad (5)$$

where $\mathbf{x}|_c$ denotes restricting the joint assignment to the scope of subproblem c . The subproblems could be single nodes, edges, trees, or cycles that are assumed to be tractable.

Note that (5) has exponentially many constraints. A succinct way to enforce them is to express the reparametrization in terms of messages:

$$\Theta^c = \Theta_0^c + \sum_{c': \mathbf{x}_c \cap \mathbf{x}_{c'} \neq \emptyset} \delta_{c' \rightarrow c}(\mathbf{X}_c \cap \mathbf{X}_{c'}) \quad (6)$$

where the messages satisfy $\delta_{c' \rightarrow c} = -\delta_{c \rightarrow c'}$. The initial subproblem potentials Θ_0^c are assumed to satisfy (5). They can be constructed, for example, by splitting the original potentials. It is easy to check that subproblem potentials defined by (6) satisfy (5) for any choice of the messages.

Note that the scopes of messages are defined by $\mathbf{X}_c \cap \mathbf{X}_{c'}$, which is the intersections among \mathbf{X}_c (with a subscript): the original *MRF variables* covered by subproblem c . This is to distinguish from \mathbf{X}^c (with a superscript), which are *subproblem variables*. For example, if subproblems c and c' both cover X_1 (original MRF variable), we would have X_1^c and $X_1^{c'}$: two different variables with independent assignments.

Since each subproblem has its own copy of variables \mathbf{X}^c , summing over their optimal values gives rise to an upper bound of (1) because each subproblem is maximized independently. This upper bound is a function of the messages:

$$\mathcal{D}(\{\delta_{c' \rightarrow c}\}) = \sum_{c \in \mathcal{C}} \max_{\mathbf{X}^c} \Theta^c(\mathbf{X}^c). \quad (7)$$

It turns out that (7) is a dual problem of (2) under some choice of \mathcal{M} . The tightness of the underlying \mathcal{M} depends on our choice of the subproblems. It has been shown (Komodakis et al., 2011) that choosing only tree structured subproblems leads to a dual of (2) with $\mathcal{M} = \mathcal{M}_L$. Duals of tightened LP relaxations can be constructed by introducing complex subproblems such as cycles.

Given a decomposition, our goal is to minimize (7) by rearranging subproblem potentials (via message passing), thereby solving the underlying LP relaxation.

2.2. Bethe cluster (region) graph

In principle the messages $\delta_{c' \rightarrow c}$ in (6) can be defined between all pairs of overlapping subproblems. However, a rather restricted form has been dominantly used in defining such messages: the Bethe cluster graph. It has a bipartite structure with one layer of “factor” nodes and one layer of small (usually unary) nodes encoding intersections between the factors. Specifically, consider using the set of subproblem $\mathcal{C} = \{i\} \cup \mathcal{F}$ consisting of unary subproblems $\{i\}$ and larger subproblems (factors) \mathcal{F} . So the dual (7) takes a restricted form:

$$\mathcal{D}(\{\delta_{f \rightarrow i}\}) = \sum_i \max_{X^i} \theta^i(X^i) + \sum_f \max_{\mathbf{X}^f} \Theta^f(\mathbf{X}^f) \quad (8)$$

where the messages are only defined *between* the two layers because of the bipartite structure. This

construction is shown in Fig. 1 (a) and (c). Note that we used superscripts (instead of subscripts) in θ^i and X^i to indicate that these are subproblem potentials and variables to distinguish from the original MRF potentials θ_i and variables X_i .

It has been shown (Meltzer et al., 2009; Sontag & Jaakkola, 2009; Sontag et al., 2011) that many existing algorithms, including (Globerson & Jaakkola, 2007; Sontag et al., 2008; 2011; Werner, 2007; Sontag & Jaakkola, 2009; Tarlow et al., 2011; Komodakis & Paragios, 2008; Meltzer et al., 2009), can be interpreted as passing messages on the Bethe cluster graph and thereby optimizing the dual objective (8). This “restricted” design had arisen from the historical concern of satisfying the *running intersection property*, thus alleviating double-counting in loopy BP (Weiss, 2000). However, this is no longer relevant under the modern view of message passing as a BCD algorithm. To this end, we will introduce a more general graph-object: *subproblem multi-graph* (Sec. 3.1) to serve the role of traditional cluster graphs. This more general setup allows us to pass messages in more flexible ways and achieve better dual (and primal) states in many situations (as shown in experiments).

3. The STC Framework

3.1. Subproblem multi-graph and subproblem trees

For notation simplicity, in this paper we treat graphs \mathcal{G} and trees \mathcal{T} as *sets* consisting of nodes and edges. So we will use $\mathcal{T} \subset \mathcal{G}$ to denote that \mathcal{T} is a subgraph of \mathcal{G} . And we will use $e \in \mathcal{T}$ or $c \in \mathcal{T}$ to denote that edge e or node c is in the tree.

Given a dual decomposition with subproblems \mathcal{C} , we build a graph-object as follows:

Definition 1 (Subproblem Multi-Graph/Tree). *Given \mathcal{C} , the **subproblem multi-graph (SMG)** $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ has one node for each $c \in \mathcal{C}$, and one edge between c and c' for each tuple (c, c', φ) , where $\varphi \in V \cup E$ is shared by c and c' . A **subproblem tree** is a tree $\mathcal{T} \subset \mathcal{G}$.*

Note that we use (V, E) for the MRF graph and $(\mathcal{V}, \mathcal{E})$ for SMG.

This construction is illustrated in Fig. 1 (b) where we decompose the MRF (Fig. 1 (a)) into four subproblems.

Note that if we include all unary subproblems into the decomposition, we would get a SMG similar to Fig. 1 (c) but with extra edges among the non-unary

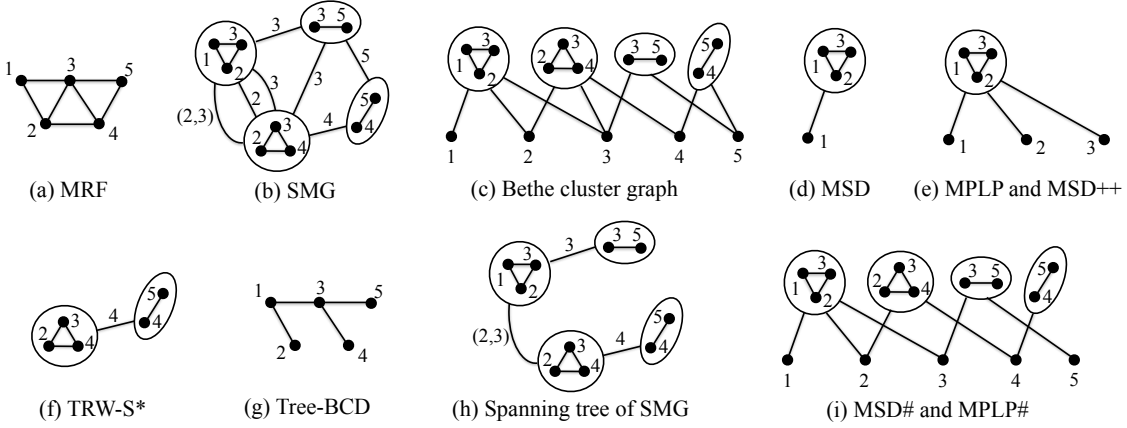


Figure 1. Examples of subproblem multi-graph and subproblem trees. See text for explanation.

subproblems. So a tree in the Bethe cluster graph (which we call a *Bethe tree*) is also a *subproblem tree* by definition. Therefore all conclusions in this paper on subproblem trees apply to Bethe trees.

Fig. 1 (d) to (i) are examples of subproblem trees. They correspond to blocks updated by different algorithms as explained later. Note that (d) (e) and (i) are Bethe trees corresponding to the Bethe cluster graph (c); (f) and (h) are trees in the SMG (b); (g) is a Bethe tree of a different dual decomposition (with all edges and nodes as subproblems).

For each SMG edge $(c, c', \varphi) \in \mathcal{E}$, we have messages² $\delta_{c' \rightarrow c}(X_\varphi) = -\delta_{c \rightarrow c'}(X_\varphi)$. Therefore the block (of dual variables) associated with subproblem tree \mathcal{T} is given by

$$\mathcal{B}^{\mathcal{T}} = \{\delta_{c' \rightarrow c}(X_\varphi) : (c, c', \varphi) \in \mathcal{T}\}. \quad (9)$$

3.2. Max-consistency and dual-optimal on trees

Given a block $\mathcal{B}^{\mathcal{T}}$ associated with some subproblem tree \mathcal{T} , we want to achieve dual-optimal w.r.t. that block:

Definition 2 (Dual-optimal on \mathcal{T}). *The subproblem potentials $\{\Theta^c\}$ are dual-optimal on \mathcal{T} if we can not further decrease the dual objective by changing messages in $\mathcal{B}^{\mathcal{T}}$.*

A message passing algorithm achieves this by enforcing some *consistency constraint*. We first identify a constraint that is equivalent to dual-optimal on \mathcal{T} .

²The messages defined in (6) appear to have larger scopes. However, if the MRF is parametrized over $\mathcal{A} = V \cup E$, we can always reparametrize the messages in (6) using messages over unary or pairwise scopes. For high-order MRFs we could allow φ to be larger scopes as well.

Definition 3 (Assignments agree on \mathcal{T}). *Assignments to all subproblems $\{\mathbf{x}^c\}_{c \in \mathcal{T}}$ agree on \mathcal{T} , denoted as $\{\mathbf{x}^c\} \sim \mathcal{T}$, if for $\forall (c, c', \varphi) \in \mathcal{T}$ we have $x_\varphi^c = x_\varphi^{c'}$.*

Definition 4 (Weak max-consistency on \mathcal{T}). *$\{\Theta^c\}_{c \in \mathcal{T}}$ satisfies weak max-consistency if*

$$\sum_{c \in \mathcal{T}} \max_{\mathbf{X}^c} \Theta^c(\mathbf{X}^c) = \max_{\{\mathbf{x}^c\} \sim \mathcal{T}} \sum_{c \in \mathcal{T}} \Theta^c(\mathbf{X}^c) \quad (10)$$

That is, maximizing each subproblem independently gets to the same optimal value as maximizing them while requiring the assignments to agree on the tree. This condition turns out to be equivalent to dual optimal on \mathcal{T} (Proposition 1). We will show that our message passing algorithm (Sec. 3.3) enforces this constraint for arbitrary \mathcal{T} .

We now identify two other stronger constraints (enforced by existing algorithms) that are *sufficient but not necessary* for dual-optimal on \mathcal{T} .

Let M_φ^c be the (log)-max-marginals of c on φ :

$$M_\varphi^c(x_\varphi) = \max_{\mathbf{X}^c|_\varphi = x_\varphi} \Theta^c(\mathbf{X}^c), \quad (11)$$

Note that if $\varphi = (i, j) \in E$, $\mathbf{X}^c|_\varphi = x_\varphi$ means $X_i^c = x_i$ and $X_j^c = x_j$.

The following consistency constraint requires the subproblem to agree on their max-marginals over the tree:

Definition 5 (Strong max-consistency³ on \mathcal{T}). *$\{\Theta^c\}_{c \in \mathcal{T}}$ satisfies strong max-consistency if $M_\varphi^c = M_\varphi^{c'}, \forall (c, c', \varphi) \in \mathcal{T}$.*

Another consistency constraint requires that the

³This corresponds to “max-consistency” in (Meltzer et al., 2009).

subproblem potentials match the max-marginals of the tree distribution:

Definition 6 (MPLP max-consistency on \mathcal{T}). For Bethe tree \mathcal{T} with N unary clusters, $\{\theta^i, \Theta^f\}_{i,f \in \mathcal{T}}$ satisfies MPLP max-consistency if $\theta^i = \frac{1}{N} \mu_i^{\mathcal{T}}$ and $\Theta^f = \mu_f^{\mathcal{T}} - \sum_{i \in f, i \in \mathcal{T}} \frac{N_i^f}{N} \mu_i^{\mathcal{T}}$ where $\mu_i^{\mathcal{T}}$ and $\mu_f^{\mathcal{T}}$ are max-marginals⁴ of the tree \mathcal{T} , and N_i^f is the number of unary clusters in the subtree rooted by i on the opposite side of f .

Note that MPLP max-consistency is only defined for Bethe trees. With a slight abuse of notation, we used $i \in f$ in the above definition to denote that the scope of unary subproblem i is contained in the scope of factor subproblem f in a Bethe cluster graph.

The relations among these consistency constraints are given below.

Proposition 1. For any Bethe tree \mathcal{T} ,

MPLP max-consistency \implies weak max-consistency

For any subproblem tree \mathcal{T} (including Bethe trees),

strong max-consistency \implies weak max-consistency.

weak max-consistency \iff dual-optimal on \mathcal{T} .

Proof for this (as well as all other propositions in the rest of this paper) are given in Appendix (supplement material), where we also give an example that satisfies weak max-consistency but not the other two.

The two stronger constraints are enforced by existing methods:

Proposition 2. Max-sum diffusion (MSD) (Werner, 2007) performs BCD on the Bethe dual (8). Each BCD step enforces strong max-consistency for a Bethe tree consisting of one $f \in \mathcal{F}$ and one $i \in f$ (as in Fig. 1 (d)).

Proposition 3. TRW-S (Kolmogorov, 2006) performs BCD on the dual (7) where each subproblem $c \in \mathcal{C}$ is a tree (of the MRF). Each BCD step enforces strong max-consistency for a subproblem tree \mathcal{T} with nodes $\mathcal{V}_{\mathcal{T}} = \{c : \text{scope}(c) \supset \varphi\}$ for a given $\varphi \in V \cup E$, and edges (c, c', φ) that constitute a tree over $\mathcal{V}_{\mathcal{T}}$ ⁵.

Proposition 4. MPLP (Globerson & Jaakkola, 2007; Sontag et al., 2011) performs BCD on the Bethe dual

⁴ Max-marginals of subproblem tree \mathcal{T} over \tilde{c} is defined as $\mu_{\tilde{c}}^{\mathcal{T}}(\mathbf{x}_{\tilde{c}}) = \max_{\{\mathbf{x}^c\} \sim \mathcal{T}, \mathbf{x}^{\tilde{c}} = \mathbf{x}_{\tilde{c}}} \sum_{c \in \mathcal{T}} \Theta^c(\mathbf{X}^c)$.

⁵In TRW-S the subproblems (MRF trees) can be connected in arbitrary ways into \mathcal{T} . Because all edges in \mathcal{T} have the same scope φ , strong max-consistency essentially requires the max-marginals on φ to be the same for all $c \in \mathcal{T}$.

(8). Each BCD step enforces MPLP max-consistency for a Bethe tree consisting of one $f \in \mathcal{F}$ and all $i \in f$ (as in Fig. 1 (e)).

Proposition 5. Tree-BCD (Sontag & Jaakkola, 2009) performs BCD on the Bethe dual (8) where all f are pairwise factors (i, j) . Each BCD step enforces MPLP max-consistency for a Bethe tree corresponding to a spanning tree of the MRF (as in Fig. 1 (g)).

Note that a spanning tree of the MRF (with $|V| - 1$ edges) is much smaller than a spanning tree of the SMG (of a dual decomposition into all edges and nodes). The latter contains all (possibly $O(|V|^2)$) edges of the MRF.

3.3. The STC algorithm

Now we show how to attain weak max-consistency for any subproblem tree (such as Fig. 1 (h) and (i)). To express the algorithm concisely, we use $x \xrightarrow{m} y$ to denote two updates: $x \leftarrow x - m$ and $y \leftarrow y + m$. We assume subproblem solvers for all c that output \mathbf{M}_{φ}^c . For tree-structured subproblems this is straightforward. For cycle subproblems we use (Felzenszwalb & McAuley, 2011) which provides a fast way of computing the junction-tree messages.

Our algorithm (Alg. 1) calibrates a subproblem-tree by an upstream pass and a downstream pass; both update subproblem potentials “in place” without storing any message (although conceptually it can be viewed as updating the messages in (7) as well). The

Algorithm 1 Subproblem tree calibration (STC)

- 1: INPUT: \mathcal{T} ; allocation weight a_c for each $c \in \mathcal{T}$ satisfying $a_c \geq 0, \sum_{c \in \mathcal{T}} a_c = 1$.
 - 2: /* upstream pass */
 - 3: **for** edges (c, c', φ) in upstream order **do**
 - 4: $\Theta^c \xrightarrow{\mathbf{M}_{\varphi}^c} \Theta^{c'}$, where c is a child of c'
 - 5: /* downstream pass */
 - 6: $w_c \leftarrow \sum_{\tilde{c} \in \mathcal{T}_{\tilde{c}}} a_{\tilde{c}}$ for all $c \in \mathcal{T}$, where $\mathcal{T}_{\tilde{c}}$ is the subtree rooted by \tilde{c}
 - 7: **for** nodes c' in downstream order **do**
 - 8: (P-STC only) precompute all $\mathbf{M}_{\varphi}^{c'}$ used below
 - 9: **for** all c (children of c') **do**
 - 10: $\Theta^{c'} \xrightarrow{(w_c/w_{c'}) \cdot \mathbf{M}_{\varphi}^{c'}} \Theta^c$
 - 11: (S-STC only) $w_{c'} \leftarrow w_{c'} - w_c$
-

downstream pass of STC can be performed for each children of a node either *sequentially* or in *parallel*. They differ in line 8 and 11 of Alg. 1. We call these two alternatives S-STC and P-STC, respectively. Any statement about STC in the following, if not specified, applies to both P-STC and S-STC.

Overall we repeatedly choose different trees and

perform the calibration.

Proposition 6. *STC enforces weak max-consistency for any \mathcal{T} and any allocation weights a_c satisfying $a_c \geq 0$, $\sum_{c \in \mathcal{T}} a_c = 1$.*

Note that this (together with Proposition 1) implies monotonicity and convergence of the overall algorithm.

We now clarify the relations between STC and existing methods.

Proposition 7. *Applying P-STC to the Bethe tree specified in Proposition 4 with allocation weights $a_i = \frac{1}{|\mathcal{F}|}$, $a_f = 0$ is equivalent to MPLP.*

Proposition 8. *Applying STC to the Bethe tree specified in Proposition 2 with allocation weights $a_f = 0.5$, $a_i = 0.5$ is equivalent to MSD.*

Proposition 9. *Applying STC to the subproblem tree specified in Proposition 3 with uniform allocation weights $a_c = \frac{1}{|\mathcal{T}|}$ is equivalent to TRW-S.*

This implies that STC actually enforces strong max-consistency when applied to settings of MSD and TRW-S. Indeed, we have:

Proposition 10. *If all edges of \mathcal{T} have same scope φ , P-STC is equivalent to S-STC. Both achieve strong max-consistency when applied with uniform allocation weights $a_c = \frac{1}{|\mathcal{T}|}$*

This condition (all edges of \mathcal{T} have the same scope) is satisfied in MSD and TRW-S.

Given Proposition 9, we can easily generalize TRW-S beyond tree subproblems⁶. We call this generalized algorithm TRW-S*: applying STC to a subproblem tree consists of all subproblems sharing a given MRF node or edge, as in Fig. 1 (f). It will appear in experiments.

3.4. Fixed-point characterization

In this part we give a fixed-point characterization of the STC algorithm, and we show that it is in fact equivalent to the weak tree agreement (WTA) condition (Kolmogorov, 2006). Indeed, the fact that TRW-S (Kolmogorov, 2006) is a special case of STC suggests that the latter should be at least as powerful as the former.

Definition 7 (WTA2). *Subproblem potentials $\{\Theta^c\}$ satisfy weak tree-agreement-two (WTA2) if for any subproblem tree \mathcal{T} , $\exists\{\mathbf{x}^c\} \sim \mathcal{T}$ that are optimal for each $c \in \mathcal{T}$ individually.*

⁶However, the acceleration trick by monotonic chains is not applicable in general.

Proposition 11 (STC Fixed-Point). *If $\{\Theta^c\}$ do not satisfy WTA2, we can always find \mathcal{T} such that applying STC to \mathcal{T} decreases the dual objective. If $\{\Theta^c\}$ satisfy WTA2, applying STC to any \mathcal{T} does not change the dual objective and preserves the WTA2 condition.*

We re-state the WTA condition from (Kolmogorov, 2006) in a slightly generalized form (allowing c to be arbitrary subproblems, not restricted to trees).

Definition 8 (WTA). *Subproblem potentials $\{\Theta^c\}$ satisfy weak tree-agreement (WTA) if we can find a subset of optimal assignments for each c , denoted as $OPT(c)$, such that: for any small subproblem tree \mathcal{T} consists of two nodes c and c' connected by (c, c', φ) , we have: for $\forall \mathbf{x}^c \in OPT(c)$, $\exists \mathbf{x}^{c'} \in OPT(c')$, $\{\mathbf{x}^c, \mathbf{x}^{c'}\} \sim \mathcal{T}$.*

Comparing to WTA, WTA2 appears to assert a weaker constraint on more general \mathcal{T} . This originates from the fact that STC calibrates more general \mathcal{T} and enforces a weaker consistency constraint than that of MSD and TRW-S. However asymptotically they are equally powerful:

Proposition 12. *WTA2 is equivalent to WTA.*

3.5. Choosing allocation weights

In this part we clarify the role of allocation weights a_c , which are responsible for the message coefficients in the downstream pass of STC. As we have seen in Proposition 8 and 7, there are subtle (but important) differences between MPLP and MSD in this aspect. (This has also been noted in (Sontag et al., 2011).)

In order to understand the role of allocation weights a_c , we first show some detailed characterization of STC.

Proposition 13 (STC Allocation). *After STC, for each subproblem c we have:*

$$\max_{\mathbf{X}^c} \Theta^c(\mathbf{X}^c) = a_c \cdot \max_{\{\mathbf{X}^{\tilde{c}}\} \sim \mathcal{T}} \sum_{\tilde{c} \in \mathcal{T}} \Theta^{\tilde{c}}(\mathbf{X}^{\tilde{c}}) \quad (12)$$

Intuitively, the downstream pass allocate “energy” to all subproblems according to their allocation weights. We can further show that:

Proposition 14 (Detailed Monotonicity). *Let \mathcal{D}_0 , \mathcal{D}_1 , \mathcal{D}_2 be the dual objective value before STC, after upstream pass, and after downstream pass, respectively. We have $\mathcal{D}_0 \geq \mathcal{D}_1 = \mathcal{D}_2$. Moreover, the dual objective remains constant in each single step of the downstream pass.*

That is, the downstream pass essentially moves around in a plateau of the dual, preparing for the next

downhill move. Therefore the allocation weights determine where to settle on that plateau.

When \mathcal{T} is a Bethe tree, two styles of allocation have been used. One is the MPLP-style allocation, which assigns zero weights to all non-unary subproblems, and uniformly allocate among unary subproblems. This is the case for MPLP and Tree-BCD⁷. The other is the MSD-style allocation, which assigns uniform allocation weights to all subproblems. This is the case for MSD and its generalized version MSD++ (Sontag et al., 2011).

Given our framework, it is straightforward to further generalize these two styles of allocation to spanning trees of the Bethe cluster graph (as in Fig. 1 (i)). We call these two algorithms MPLP# and MSD#, respectively⁸. They will appear in experiments.

3.6. Generate primal solutions

Given subproblem potentials, solutions to the original MAP inference problem can be constructed in different ways (Komodakis et al., 2011). For example, when unary subproblems are present, we could simply take the assignments that minimize each unary subproblem. Or we could visit all subproblems in turn, and for each one commit the variables in its scope to the subproblem solution.

In order to better leverage “beliefs” of different subproblems as well as “smoothness” across subproblems, we propose the following method, which works better in practice comparing to other heuristics. (The comparison is not shown due to space limit.) In practice, one could construct multiple assignments from different heuristics and choose the one with the best objective value.

We visit the variables (in the original MRF) in some ordering, for example, X_1, X_2, \dots, X_N . And for X_i we choose the assignment:

$$x_i = \operatorname{argmax}_{c:i \in \text{scope}(c)} \sum_{\mathbf{X}^c \setminus X_i} \max_{\mathbf{X}^c \setminus X_i} \Theta^c(\mathbf{X}^c | X_j = x_j, \forall j < i) \quad (13)$$

That is, when visiting each X_i , we choose its assignment to maximize the sum of all max-marginals from all subproblems covering X_i . And then we fix $X_i = x_i$ in all subproblems (this will affect subsequent visits to these subproblems).

⁷Although Tree-BCD is not a special case of STC. We can define its choice of allocation weights according to the distribution of “energy” in its calibrated tree.

⁸Specifically, MPLP# means applying STC with MPLP-style allocation weights to randomly selected spanning trees as in Fig. 1 (i), and similarly for MSD#.

4. Experimental Results

[Experiment Setup] We used three MAP inference tasks in experiments: (1) The protein design benchmark (Yanover et al., 2006). We use the 20 largest problems from that dataset, with number of variables from 101 to 180, number of edges from 1,973 to 3,005, variable cardinality up to 154. (2) Synthetic 20-by-20 grid with variable cardinality of 100. Potentials are drawn from $\mathcal{N}(0, 1)$. (3) The “object detection” task from PIC-2011⁹. There are 37 problem instances, each has 60 variables and 1,770 edges, variable cardinality from 11 to 21. Each task corresponds to one row in Fig. 2.

For each tasks, we use two settings: with or without cycle subproblems (loose or tightened LP relaxation). The former corresponds to the left two columns in Fig. 2, and the latter corresponds to the right two columns. Each setting for each task is shown by a primal(right)-dual(left) pair of figures (averaged over all problem instances). Each primal-dual pair share one legend.

The cycles in dual decomposition are selected in a static manner by applying the criterion of (Sontag et al., 2008) to the original edge potentials. We did this (instead of dynamically adding cycles) because our goal is to compare different dual methods and we want them to always operate on the same dual problem. For protein design and PIC, we selected 500 triangles for each problem instance. For grid, we used all 381 squares of size four.

[Methods Compared] Among the methods compared, MPLP (Globerson & Jaakkola, 2007), MSD++ (Sontag et al., 2011), Tree-BCD (Sontag & Jaakkola, 2009), TRW-S MonoChain (Kolmogorov, 2006) are existing methods. TRW-S* is a straightforward generalization of a existing method. MPLP#, MSD#, and STC are new methods derived from our framework. Here STC means calibrating randomly chosen spanning trees of the SMG without unary subproblems (as in Fig. 1 (h)). For TRW-S MonoChain we used the code of (Szeliski et al., 2008). Note that it is only applicable to the grid problem with long monotonic chains. For TRW-S* we use MRF edges (and cycles for the tightened setting) as subproblems. All other methods have been explained in earlier sections.

[Result Analysis] When cycle subproblems are present, we observe that STC and TRW-S* performs significantly better in all tasks. The crucial difference

⁹<http://www.cs.huji.ac.il/project/PASCAL/index.php>

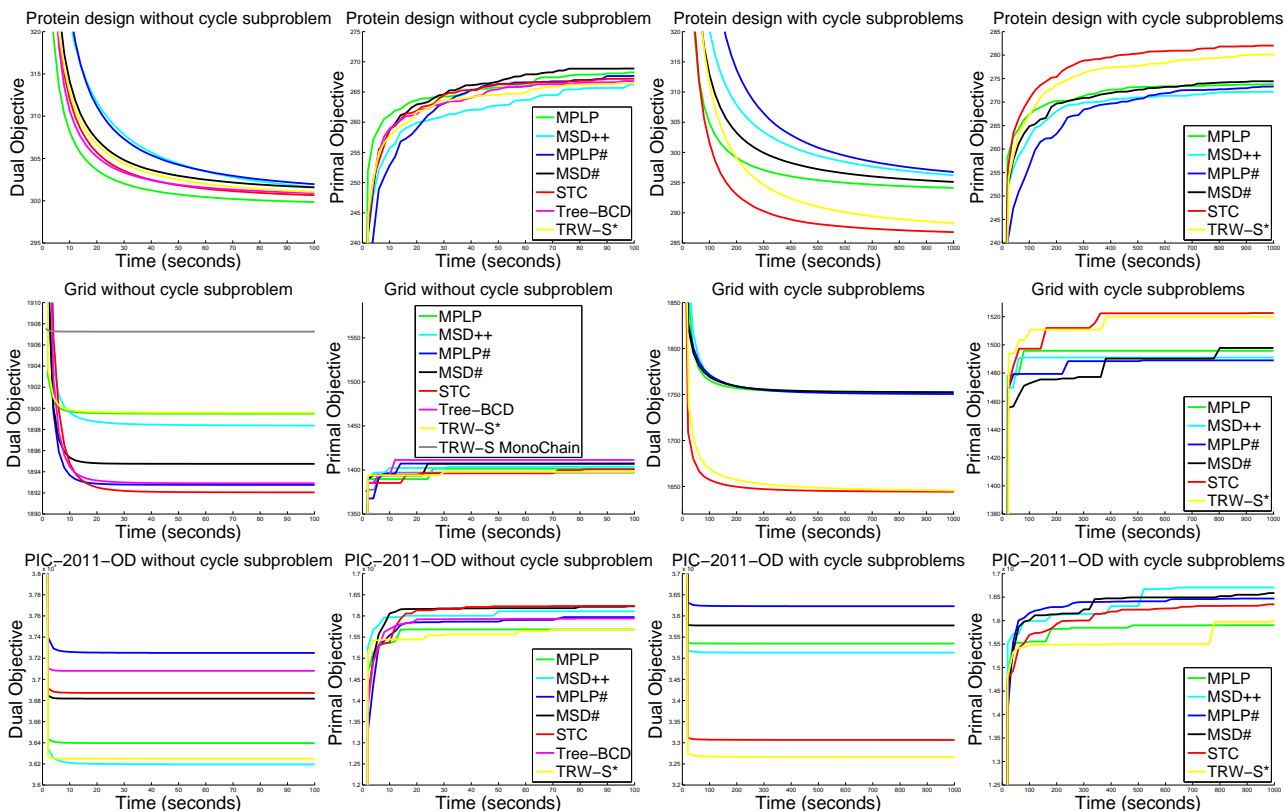


Figure 2. Experimental results. This figure is better viewed in color.

is that: these two methods allow the subproblems communicating through MRF edges, whereas all other methods are restricted to passing messages over MRF nodes (they use the Bethe cluster graph as in Fig. 1 (c)). This reveals the limitation of this bipartite construction.

For sparse MRFs (the grid case), getting more subproblems involved in each BCD steps turns out to be important. As we can see from the second plot in the first column of Fig. 2. The methods that updates “global” blocks (STC, Tree-BCD, MPLP#, MSD#) performs significantly better than the ones with “local” blocks (MPLP, MSD++, and TRW-S*). Note that TRW-S MonoChain get stuck at an even worse dual state than the “local” methods. This is because its block selection (and ordering) is even more restrictive due to using monotonic chains.

Overall, we observe that different methods tend to “converge” to different dual objectives, Even though the dual objectives in each plot should have exactly the same optimal value. Note that all these methods perform BCD—they achieve dual optimal on a block in each step. Therefore choosing blocks (as well as choosing plateau dual states) is very important to the performance (and final result) of a message passing

algorithm. Generally speaking, being able to choose blocks and plateau states in a more flexible manner (as our framework reveals) could help us get to better dual states, whereas making these choices in a restricted manner (as most existing methods do) could easily get stuck.

5. Conclusion

Our framework revealed two dimensions of flexibility in designing a message passing algorithm for MAP inference: choosing blocks to update, and choosing a dual state on a plateau in each BCD step. The STC algorithm can be applied with extreme flexibility in these choices. Although these choices appear to be important to performance, any known fixed strategy in making them does not seem to be optimal across different scenarios. If we could find principled and adaptive strategies in making these choices, we will be able to design much more powerful message passing algorithms.

Acknowledgement

This work has been supported by the Max Planck Center for Visual Computing and Communication.

References

- Felzenszwalb, P. F. and McAuley, J. J. Fast inference with min-sum matrix product. *PAMI*, 33(12):2549–2554, 2011.
- Globerson, Amir and Jaakkola, Tommi. Fixing max-product: Convergent message passing algorithms for MAP LP-relaxations. In *NIPS*, 2007.
- Koller, D. and Friedman, N. *Probabilistic Graphical Models: Principles and Techniques*. MIT Press, 2009.
- Kolmogorov, V. Convergent tree-reweighted message passing for energy minimization. *PAMI*, 28(10):1568–1583, October 2006.
- Komodakis, N. and Paragios, N. Beyond loose LP-relaxations: Optimizing MRFs by repairing cycles. In *ECCV*, pp. III: 806–820, 2008.
- Komodakis, Nikos, Paragios, Nikos, and Tziritas, Georgios. MRF energy minimization and beyond via dual decomposition. *PAMI*, 33(3):531–552, 2011.
- Meltzer, Talya, Globerson, Amir, and Weiss, Yair. Convergent message passing algorithms - a unifying view. In *UAI*, pp. 393–401, 2009.
- Shimony. Finding MAPs for belief networks is NP-hard. *AIJ: Artificial Intelligence*, 68, 1994.
- Sontag, David and Jaakkola, Tommi. Tree block coordinate descent for MAP in graphical models. *AISTATS*, 5:544–551, 2009.
- Sontag, David, Meltzer, Talya, Globerson, Amir, Jaakkola, Tommi, and Weiss, Yair. Tightening LP relaxations for MAP using message passing. In *UAI*, pp. 503–510, 2008.
- Sontag, David, Globerson, Amir, and Jaakkola, Tommi. Introduction to dual decomposition for inference. In Sra, Suvrit, Nowozin, Sebastian, and Wright, Stephen J. (eds.), *Optimization for Machine Learning*. MIT Press, 2011.
- Szeliski, R. S., Zabih, R., Scharstein, D., Veksler, O., Kolmogorov, V., Agarwala, A., Tappen, M., and Rother, C. A comparative study of energy minimization methods for markov random fields with smoothness-based priors. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 30(6):1068–1080, June 2008.
- Tarlow, Daniel, Batra, Dhruv, Kohli, Pushmeet, and Kolmogorov, Vladimir. Dynamic tree block coordinate ascent. In *ICML*, 2011.
- Wainwright, Martin J. and Jordan, Michael I. Graphical models, exponential families, and variational inference. *Foundations and Trends in Machine Learning*, 1(1-2):1–305, 2008.
- Weiss, Y. Correctness of local probability propagation in graphical models with loops. *Neural Computation*, 12:1–41, 2000.
- Werner, T. A linear programming approach to max-sum problem: A review. *PAMI*, 29(7):1165–1179, July 2007.
- Yanover, Chen, Meltzer, Talya, and Weiss, Yair. Linear programming relaxations and belief propagation - an empirical study. *JMLR*, 7:1887–1907, 2006.

Appendix

In this section we prove all propositions in the paper. Instead of showing them in their original order, we follow a logical ordering to allow building later derivations on earlier conclusions.

Lemma 1. *Weak max-consistency is equivalent to: $\exists\{\mathbf{x}^c\} \sim \mathcal{T}$ that maximizes each subproblem individually.*

Proof. If $\exists\{\mathbf{x}^c\} \sim \mathcal{T}$ that maximizes each subproblem, we have:

$$\sum_{c \in \mathcal{T}} \max_{\mathbf{x}^c} \Theta^c(\mathbf{X}^c) = \sum_{c \in \mathcal{T}} \Theta^c(\mathbf{x}^c) \leq \max_{\{\mathbf{x}^c\} \sim \mathcal{T}} \sum_{c \in \mathcal{T}} \Theta^c(\mathbf{X}^c) \quad (14)$$

However, we also have $\sum_{c \in \mathcal{T}} \max_{\mathbf{x}^c} \Theta^c(\mathbf{X}^c) \geq \max_{\{\mathbf{x}^c\} \sim \mathcal{T}} \sum_{c \in \mathcal{T}} \Theta^c(\mathbf{X}^c)$ because the LHS has less constraint in maximization. Therefore weak max-consistency holds.

In the other direction, if we have weak max-consistency, i.e. $\sum_{c \in \mathcal{T}} \max_{\mathbf{x}^c} \Theta^c(\mathbf{X}^c) = \max_{\{\mathbf{x}^c\} \sim \mathcal{T}} \sum_{c \in \mathcal{T}} \Theta^c(\mathbf{X}^c)$. We claim that the maximizer of RHS, denoted as $\{\mathbf{x}^c\}$, must be optimal for each c individually, otherwise the LHS can be increased by changing \mathbf{x}^c that is not optimal for the corresponding subproblem c . \square

Lemma 2. *For any $\{\mathbf{x}^c\} \sim \mathcal{T}$, we have $\sum_{c \in \mathcal{T}} \Theta^c(\mathbf{x}^c) = \sum_{c \in \mathcal{T}} \Theta_0^c(\mathbf{x}^c)$ if Θ^c is obtained from Θ_0^c by only varying block $\mathcal{B}^{\mathcal{T}}$.*

Proof. For any $\{\mathbf{x}^c\} \sim \mathcal{T}$

$$\sum_{c \in \mathcal{T}} \Theta^c(\mathbf{x}^c) = \sum_{c \in \mathcal{T}} [\Theta_0^c(\mathbf{x}^c) + \sum_{c': (c, c', \varphi) \in \mathcal{T}} \delta_{c' \rightarrow c}(\mathbf{x}^c|_{\varphi})] \quad (15)$$

So it suffice to show that

$$\sum_{c \in \mathcal{T}} \sum_{c': (c, c', \varphi) \in \mathcal{T}} \delta_{c' \rightarrow c}(\mathbf{x}^c|_{\varphi}) = 0 \quad (16)$$

The follows from the fact that $\delta_{c' \rightarrow c} = -\delta_{c \rightarrow c'}$ (by definition), and $\mathbf{x}^c|_{\varphi} = \mathbf{x}^{c'}|_{\varphi}$ (because $\{\mathbf{x}^c\} \sim \mathcal{T}$). \square

Proof of Proposition 1

Proposition. *For any Bethe tree \mathcal{T} ,*

MPLP max-consistency \implies weak max-consistency

For any subproblem tree \mathcal{T} (including Bethe trees),

strong max-consistency \implies weak max-consistency.

weak max-consistency \iff dual-optimal on \mathcal{T} .

Proof. We first show that MPLP max-consistency implies weak max-consistency. Note that according to the definition of $\mu_i^{\mathcal{T}}$ and $\mu_f^{\mathcal{T}}$

$$\begin{aligned} \max_{\{X_i, \mathbf{X}_f\} \sim \mathcal{T}} [\sum_{i \in \mathcal{T}} \theta^i(X_i) + \sum_{f \in \mathcal{T}} \Theta^f(\mathbf{X}_f)] = \\ \max_{x_i} \mu_i^{\mathcal{T}}(x_i) = \max_{\mathbf{x}_f} \mu_f^{\mathcal{T}}(\mathbf{x}_f) \end{aligned} \quad (17)$$

We denote this quantity (maximum of tree) as $M^{\mathcal{T}}$. To validate the weak max-consistency condition (10), we first note that the RHS of (10) equals $M^{\mathcal{T}}$. For the LHS (10) of we have:

$$\begin{aligned} \sum_{i \in \mathcal{T}} \max_{X_i} \theta^i(X_i) + \sum_{f \in \mathcal{T}} \max_{\mathbf{X}_f} \Theta^f(\mathbf{X}_f) \\ = \sum_{i \in \mathcal{T}} \max_{X_i} \frac{1}{N} \mu_i^{\mathcal{T}}(X_i) \\ + \sum_{f \in \mathcal{T}} \max_{\mathbf{X}_f} [\mu_f^{\mathcal{T}} - \sum_{i \in f, i \in \mathcal{T}} \frac{N_i^f}{N} \mu_i^{\mathcal{T}}] \end{aligned} \quad (18)$$

Note that N is the number of unary subproblems in \mathcal{T} , so we have $\sum_{i \in \mathcal{T}} \max_{X_i} \frac{1}{N} \mu_i^{\mathcal{T}}(X_i) = M^{\mathcal{T}}$. It remains to show that all the remaining terms are zero. We note that

$$\mu_i^{\mathcal{T}}(x_i) \geq \mu_f^{\mathcal{T}}(\mathbf{x}_f) \quad (19)$$

for any $i \in f$, $i \in \mathcal{T}$, and assignments x_i consistent with \mathbf{x}_f . So we have

$$\mu_f^{\mathcal{T}} - \sum_{i \in f, i \in \mathcal{T}} \frac{N_i^f}{N} \mu_i^{\mathcal{T}} \leq 0 \quad (20)$$

by noting that $\sum_{i \in f, i \in \mathcal{T}} \frac{N_i^f}{N} = 1$. Also note that (20) equals zero with the assignments corresponds to the maximizer of the LHS of (17). Therefore all terms in the second sum of (18) are zero, which proves that MPLP max-consistency implies weak max-consistency.

Now we show that strong max-consistency implies weak max-consistency. Using Lemma 1, we can always find such $\{\mathbf{x}^c\} \sim \mathcal{T}$ by following these step: choose arbitrary maximizer \mathbf{x}^c for subproblem c ; choose maximizers for neighbors of c that agrees with \mathbf{x}^c ; and so on. We can find $\{\mathbf{x}^c\} \sim \mathcal{T}$ (that maximizes each subproblem) in this manner because of the tree-structure and the fact that all adjacent subproblems in the tree \mathcal{T} have identical max-marginals on their sepsets (strong max-consistency).

Now we show that weak max-consistency is equivalent to dual-optimal on \mathcal{T} . We have

$$\begin{aligned} \sum_{c \in \mathcal{T}} \max_{\mathbf{X}^c} \Theta^c(\mathbf{X}^c) &\geq \max_{\{\mathbf{x}^c\} \sim \mathcal{T}} \sum_{c \in \mathcal{T}} \Theta^c(\mathbf{X}^c) \\ &= \max_{\{\mathbf{x}^c\} \sim \mathcal{T}} \sum_{c \in \mathcal{T}} \Theta_0^c(\mathbf{X}^c) \end{aligned} \quad (21)$$

where the second equality follows from Lemma 2. The first expression in (21) is the dual objective (excluding terms outside \mathcal{T} , that do not change). The last expression is a constant; therefore weak max-consistency (equal sign holds for (21)) implies dual-optimal on \mathcal{T} .

If weak max-consistency is not satisfied, from (21) we must have:

$$\sum_{c \in \mathcal{T}} \max_{\mathbf{X}^c} \Theta^c(\mathbf{X}^c) > \max_{\{\mathbf{X}^c\} \sim \mathcal{T}} \sum_{c \in \mathcal{T}} \Theta_0^c(\mathbf{X}^c) \quad (22)$$

However, from Proposition 6¹⁰ we know that the value on the RHS can be attained by enforcing weak consistency with STC on \mathcal{T} (that only changes dual variables in $\mathcal{B}^{\mathcal{T}}$). So we have not attained dual optimal on \mathcal{T} yet. Therefore dual optimal on \mathcal{T} implies weak max-consistency.

Now we give an example to show that weak max-consistency is *strictly* weaker than the other two. Consider a subproblem tree consists of three nodes (X_1, X_2) (as root), (X_1) , and (X_2) , with the subproblem potentials given by:

$$\begin{aligned} \Theta^{12} &= \begin{pmatrix} -0.05 & 0.05 \\ 0 & 0.2 \end{pmatrix} \\ \Theta^1 &= (0.1, 0.2) \quad \Theta^2 = \begin{pmatrix} 0.05 \\ 0.2 \end{pmatrix} \end{aligned} \quad (23)$$

$X_1, X_2 \in \{0, 1\}$ are binary variables. It is easy to verify that they satisfy weak max-consistency (and dual optimal), because the assignment $X_1 = 1, X_2 = 1$ maximizes each subproblem individually. However, the max-marginal of Θ^{12} on X_1 is $(0, 0.2)$, which does not agree with Θ^1 . Therefore strong max-consistency does not hold. It is also easy to verify that MPLP max-consistency does not hold. \square

We use the following lemmas to characterize the behavior STC. Let \mathcal{T}^c be the subtree of \mathcal{T} rooted by c , and $c \downarrow$ be the set of all children of c .

Lemma 3. *If for all $c \in c' \downarrow$ we have that:*

$$\Theta^c(\mathbf{x}^c) = \max_{\substack{\{\mathbf{X}^{\tilde{c}}\} \sim \mathcal{T}^c \\ \mathbf{X}^c = \mathbf{x}^c}} \sum_{\tilde{c} \in \mathcal{T}^c} \Theta^{\tilde{c}}(\mathbf{X}^{\tilde{c}}), \quad \forall \mathbf{x}^c \quad (24)$$

$$\max_{\mathbf{X}^{\tilde{c}} |_{\varphi_{\tilde{c}}} = \mathbf{x}^c |_{\varphi_{\tilde{c}}}} \Theta^{\tilde{c}}(\mathbf{X}^{\tilde{c}}) = 0, \quad \forall \mathbf{x}^c, \tilde{c} \in c \downarrow \quad (25)$$

then (24, 25) also holds for c' after performing upstream updates for all $c \in c' \downarrow$.

¹⁰Here we used Proposition 6, which is proved later. This is justified because proving Proposition 6 will not use Proposition 1.

Proof. Note that $\varphi_{\tilde{c}}$ denote the sepset of edge $(\tilde{c}, c, \varphi_{\tilde{c}})$; it is indexed by the children node \tilde{c} because each node in the tree (except for root), as children node, is associated with exactly one edge. This lemma claims that the properties described by (24, 25) propagates in the upstream direction of the tree. It will be used as an induction step to characterize all subproblems when the upstream pass finishes.

We use $\Theta_{new}^{c'}$ to denote the new potentials of c' after performing upstream pass along edge (c, c', φ_c) for all $c \in c' \downarrow$.

$$\begin{aligned} & \Theta_{new}^{c'}(\mathbf{x}^{c'}) \\ &= \Theta^{c'}(\mathbf{x}^{c'}) + \sum_{c \in c' \downarrow} M_{\varphi_c}^c(\varphi_c \sim \mathbf{x}^{c'}) \\ &= \Theta^{c'}(\mathbf{x}^{c'}) + \sum_{c \in c' \downarrow} \max_{\mathbf{x}^c |_{\varphi_c} = \mathbf{x}^{c'} |_{\varphi_c}} \Theta^c(\mathbf{x}^c) \\ &= \Theta^{c'}(\mathbf{x}^{c'}) + \sum_{c \in c' \downarrow} \max_{\substack{\{\mathbf{X}^{\tilde{c}}\} \sim \mathcal{T}^c \\ \mathbf{X}^c |_{\varphi_c} = \mathbf{x}^{c'} |_{\varphi_c}}} \sum_{\tilde{c} \in \mathcal{T}^c} \Theta^{\tilde{c}}(\mathbf{X}^{\tilde{c}}) \\ &= \Theta^{c'}(\mathbf{x}^{c'}) + \max_{\substack{\forall c \in c' \downarrow, \{\mathbf{X}^{\tilde{c}}\} \sim \mathcal{T}^c \\ \mathbf{X}^c |_{\varphi_c} = \mathbf{x}^{c'} |_{\varphi_c}}} \sum_{c \in c' \downarrow} \sum_{\tilde{c} \in \mathcal{T}^c} \Theta^{\tilde{c}}(\mathbf{X}^{\tilde{c}}) \\ &= \max_{\substack{\{\mathbf{X}^{\tilde{c}}\} \sim \mathcal{T}^{c'} \\ \mathbf{X}^{c'} = \mathbf{x}^{c'}}} \sum_{\tilde{c} \in \mathcal{T}^{c'}} \Theta^{\tilde{c}}(\mathbf{X}^{\tilde{c}}), \text{ which proves (24) for } c' \end{aligned} \quad (26)$$

where the first equality directly follows from upstream update rule; the second equality follows from the definition of max-marginals; the third equality uses the induction hypothesis (24); the fourth equality is because the trees rooted by different c do not overlap, therefore performing the maximization jointly or separately makes no difference; the last equality is moving the constant into the minimization and re-interpreting the sum and maximization scope.

For each children $c \in c' \downarrow, \forall \mathbf{x}^{c'}$ we have

$$\begin{aligned} & \max_{\mathbf{x}^c |_{\varphi_c} = \mathbf{x}^{c'} |_{\varphi_c}} \Theta_{new}^c(\mathbf{x}^c) \\ &= \max_{\mathbf{x}^c |_{\varphi_c} = \mathbf{x}^{c'} |_{\varphi_c}} [(\Theta^c(\mathbf{x}^c) - M_{\varphi_c}^c)] \\ &= \max_{\mathbf{x}^c |_{\varphi_c} = \mathbf{x}^{c'} |_{\varphi_c}} [\Theta^c(\mathbf{x}^c) - \max_{\mathbf{X}^c |_{\varphi_c} = \mathbf{x}^c |_{\varphi_c}} \Theta^c(\mathbf{X}^c)] \\ &= 0, \text{ which proves (25) for } c' \end{aligned} \quad (27)$$

where the first equality directly follows from the upstream update rule; the second equality follows from the definition of max-marginals; the last equality is because the expression inside the outer max is by definition non-positive and it attains zero when $\mathbf{x}^c = \operatorname{argmax}_{\mathbf{X}^c |_{\varphi_c} = \mathbf{x}^c |_{\varphi_c}} \Theta^c(\mathbf{X}^c)$. \square

Lemma 4. *After upstream pass of STC, we have:*

$$\max_{\mathbf{X}^c |_{\varphi_c = \mathbf{x}^{c'} |_{\varphi_c}}} \Theta^c(\mathbf{X}^c) = 0, \forall c \neq r, \forall \mathbf{x}^{c'} \quad (28)$$

where c' is the parent of c .

$$\max_{\mathbf{X}^r} \Theta^r(\mathbf{X}^r) = \max_{\{\mathbf{X}^c\} \sim \mathcal{T}} \sum_{c \in \mathcal{T}} \Theta^c(\mathbf{X}^c) \quad (29)$$

where r is the root of \mathcal{T} .

Proof. We first note that (24,25) are trivial facts if c is a leaf node of \mathcal{T} . Using Lemma 3, by induction we know that (24) hold for r when we have finished the upstream pass. That is:

$$\Theta^r(\mathbf{x}^r) = \max_{\substack{\{\mathbf{X}^c\} \sim \mathcal{T} \\ \mathbf{X}^r = \mathbf{x}^r}} \sum_{c \in \mathcal{T}} \Theta^{\tilde{c}}(\mathbf{X}^{\tilde{c}}), \quad \forall \mathbf{x}^r \quad (30)$$

Maximizing it over \mathbf{x}^r on both sides we get (29).

We also have that (25) holds for c after the upstream pass from c to its parent c' . All subsequent upstream updates will not touch c again. So (25) holds for all $c \neq r$ after the upstream pass has finished. This directly gives us (28). \square

The lemmas above have fully characterized the STC upstream pass: *the root subproblem collects all “energy” over the tree, leaving a zero max-marginal for each non-root subproblem at their upstream sepsets.* Now we consider the downstream pass through edge (c, c', φ_c) , where c' is the parent of c . There are two different schedules of downstream updates: sequential and parallel, leading to S-STC, and P-STC, respectively. In S-STC we process different children of one parent c' one-by-one, for each children c we split a portion of the max-marginal of c' and add to c . In P-STC we process all children of c' altogether: splitting a portion of max-marginal of c' (possibly at different sepsets) at the same time and adding to its children. They are different because in S-STC the parent and its max-marginals have changed in processing each children before processing subsequent ones.

We use the following lemma to characterize the basic operations of P-STC downstream pass. The basic operations of S-STC is a special case of it.

Lemma 5. *Let $\max_{\mathbf{X}^{c'}} \Theta^{c'}(\mathbf{X}^{c'}) = \eta$. If for a given subset set of children of c' , denoted as $\{c' \downarrow\}$, we have:*

$$\max_{\mathbf{X}^c |_{\varphi_c = \mathbf{x}^{c'} |_{\varphi_c}}} \Theta^c(\mathbf{X}^c) = 0, \forall \mathbf{x}^{c'}. \quad (31)$$

Then after performing the update

$$\Theta^{c'} \xrightarrow{\gamma_c \cdot M_{\varphi_c}^{c'}} \Theta^c, \quad (32)$$

in parallel for all $c \in \{c' \downarrow\}$, with $\forall \gamma_c \geq 0$, $\sum_{c \in \{c' \downarrow\}} \gamma_c \leq 1$, we have:

$$\max_{\mathbf{X}^{c'}} \Theta_{new}^{c'}(\mathbf{X}^{c'}) = (1 - \sum_{c \in \{c' \downarrow\}} \gamma_c) \eta, \quad (33)$$

$$\max_{\mathbf{X}^c} \Theta_{new}^c(\mathbf{X}^c) = \gamma_c \eta, \quad \forall c \in \{c' \downarrow\} \quad (34)$$

Note that this holds for any subset of children $\{c' \downarrow\}$. We can pick $\{c' \downarrow\}$ to be a single children to represent the basic downstream operation of S-STC.

Proof. For the new potentials of parent subproblem c' we have

$$\begin{aligned} & \max_{\mathbf{X}^{c'}} \Theta_{new}^{c'}(\mathbf{X}^{c'}) \\ &= \max_{\mathbf{X}^{c'}} [\Theta^{c'}(\mathbf{X}^{c'}) - \sum_{c \in \{c' \downarrow\}} M_{\varphi_c}^{c'}] \\ &= \max_{\mathbf{X}^{c'}} [(1 - \sum_{c \in \{c' \downarrow\}} \gamma_c) \cdot \Theta^{c'}(\mathbf{X}^{c'}) + \\ & \quad \sum_{c \in \{c' \downarrow\}} \gamma_c (\Theta^{c'}(\mathbf{X}^{c'}) - M_{\varphi_c}^{c'})] \\ &\leq (1 - \sum_{c \in \{c' \downarrow\}} \gamma_c) \eta + \sum_{c \in \{c' \downarrow\}} \gamma_c \max_{\mathbf{X}^{c'}} (\Theta^{c'}(\mathbf{X}^{c'}) - M_{\varphi_c}^{c'}) \\ &= (1 - \sum_{c \in \{c' \downarrow\}} \gamma_c) \eta, \end{aligned} \quad (35)$$

The first equality follows directly from the update rule; the second equality simply splits the first term; the third inequality uses the fact that $\gamma_c \geq 0$, $\sum_{c \in \{c' \downarrow\}} \gamma_c \leq 1$; the last equality follows from the definition of max-marginals. Also note that

$$\Theta_{new}^{c'}(\mathbf{x}^{c'}) = (1 - \sum_{c \in \{c' \downarrow\}} \gamma_c) \eta \quad (36)$$

for $\mathbf{x}^{c'} = \operatorname{argmax}_{\mathbf{X}^{c'}} \Theta^{c'}(\mathbf{X}^{c'})$. Therefore we have proved (33).

For the new potentials of children subproblem c we have

$$\begin{aligned} & \max_{\mathbf{X}^c} \Theta_{new}^c(\mathbf{X}^c) \\ &= \max_{\mathbf{X}^c} [\Theta^c(\mathbf{X}^c) + \gamma_c \cdot M_{\varphi_c}^{c'}] \\ &= \max_{\mathbf{x}_{\varphi_c}} [\max_{\mathbf{X}^c |_{\varphi_c = \mathbf{x}_{\varphi_c}}} \Theta^c(\mathbf{X}^c) + \gamma_c \cdot M_{\varphi_c}^{c'}(\mathbf{x}_{\varphi_c})] \quad (37) \\ &= \max_{\mathbf{x}_{\varphi_c}} \gamma_c \cdot M_{\varphi_c}^{c'}(\mathbf{x}_{\varphi_c}) \\ &= \gamma_c \eta, \quad \text{which proves (34) for } \forall c \in \{c' \downarrow\} \end{aligned}$$

where the first equality follows directly from the update rule; the second equality re-interprets the maximization; the third equality uses the assumption (31); the last equality follows from the definition of max-marginals and the fact that $\gamma_c \geq 0$. \square

Proof of Proposition 13

Proposition (STC Allocation). *After STC, for each subproblem c we have:*

$$\max_{\mathbf{X}^c} \Theta^c(\mathbf{X}^c) = a_c \cdot \max_{\{\mathbf{X}^{\tilde{c}}\} \sim \mathcal{T}} \sum_{\tilde{c} \in \mathcal{T}} \Theta^{\tilde{c}}(\mathbf{X}^{\tilde{c}}) \quad (38)$$

Proof. Lemma 5 has established that the downstream pass of STC allocates “energy” from parent subproblems to children subproblems in the same way as splitting real numbers. Specifically, when the downstream pass has reached c' but not its children. We have (for both S-STC and P-STC):

$$\max_{\mathbf{X}^{c'}} \Theta_{new}^{c'}(\mathbf{X}^{c'}) = \left(\sum_{c \in \mathcal{T}^{c'}} a_c \right) \cdot \max_{\{\mathbf{X}^c\} \sim \mathcal{T}} \sum_{c \in \mathcal{T}} \Theta^c(\mathbf{X}^c) \quad (39)$$

After performing downstream update for all children of c' we have

$$\max_{\mathbf{X}^{c'}} \Theta_{new}^{c'}(\mathbf{X}^{c'}) = a_{c'} \cdot \max_{\{\mathbf{X}^c\} \sim \mathcal{T}} \sum_{c \in \mathcal{T}} \Theta^c(\mathbf{X}^c) \quad (40)$$

We can use induction over the tree (with initial condition given by Lemma 4, induction step given by Lemma 5) to show that (40) holds for all subproblems in the tree after the downstream pass finishes, which proves the proposition. \square

Proof of Proposition 6

Proposition. *STC enforces weak max-consistency for any \mathcal{T} and any allocation weights a_c satisfying $a_c \geq 0$, $\sum_{c \in \mathcal{T}} a_c = 1$.*

Proof. This directly follows from Proposition 13 by summing over $c \in \mathcal{T}$ on both sides of (38). \square

Proof of Proposition 14

Proposition (Detailed Monotonicity). *Let \mathcal{D}_0 , \mathcal{D}_1 , \mathcal{D}_2 be the dual objective value before STC, after upstream pass, and after downstream pass, respectively. We have $\mathcal{D}_0 \geq \mathcal{D}_1 = \mathcal{D}_2$. Moreover, the dual objective remains constant in each single step of the downstream pass.*

Proof. Let Θ_0^c , Θ_1^c and Θ_2^c be the subproblem potentials before STC, after upstream pass, and after

downstream pass, respectively. Note that

$$\begin{aligned} & \sum_{c \in \mathcal{T}} \max_{\mathbf{X}^c} \Theta_1^c(\mathbf{X}^c) \\ &= \max_{\{\mathbf{X}^c\} \sim \mathcal{T}} \sum_{c \in \mathcal{T}} \Theta_1^c(\mathbf{X}^c) \\ &= \max_{\{\mathbf{X}^c\} \sim \mathcal{T}} \sum_{c \in \mathcal{T}} \Theta_0^c(\mathbf{X}^c) \\ &\leq \sum_{c \in \mathcal{T}} \max_{\mathbf{X}^c} \Theta_0^c(\mathbf{X}^c) \end{aligned} \quad (41)$$

where the first equation follows from Lemma 4; the second equation follows from Lemma 2; the last inequality is obtained by eliminating the constraint and pushing the max into the sum.

Therefore we have $\mathcal{D}_1 \leq \mathcal{D}_0$ by noting that the two ends of (41) represent the dual objectives excluding the terms outside \mathcal{T} (that are not changed).

For \mathcal{D}_2 we also have

$$\sum_{c \in \mathcal{T}} \max_{\mathbf{X}^c} \Theta_2^c(\mathbf{X}^c) = \max_{\{\mathbf{X}^c\} \sim \mathcal{T}} \sum_{c \in \mathcal{T}} \Theta_0^c(\mathbf{X}^c) \quad (42)$$

by using Proposition 13 and following the same reasoning as for \mathcal{D}_1 . Therefore we have $\mathcal{D}_2 = \mathcal{D}_1$. From Lemma 5 we can see that the dual objective remains constant in each step of the downstream pass. \square

Proof of Proposition 12

Proposition. *WTA2 is equivalent to WTA.*

Proof. If $\{\Theta^c\}$ satisfies WTA, we show how to find $\{\mathbf{x}^c\}$ for $\forall \mathcal{T}$ to satisfy WTA2. For any \mathcal{T} , we start from any $\mathbf{x}^c \in OPT(c)$. Because of WTA, we can always expand it to neighboring subproblems in \mathcal{T} , and so on, until we have $\{\mathbf{x}^c\} \sim \mathcal{T}$ that minimize each subproblem individually.

If $\{\Theta^c\}$ satisfies WTA2, we show how to find $\{OPT(c)\}$ to satisfy WTA. We initialize $\{OPT(c)\}$ to contain all optimal assignments for each subproblem c , and repeatedly perform the following steps:

1. Find $\mathbf{x}^c \in OPT(c)$ that cannot be expanded to one of its neighbor c' maintaining consistency over all edges between c and c' . That is, $\exists(c, c', \varphi) \in \mathcal{G}$, such that $\forall \mathbf{x}^{c'} \in OPT(c')$, $\mathbf{x}^c|_{\varphi} \neq \mathbf{x}^{c'}|_{\varphi}$.
2. Assert that \mathbf{x}^c is not the only assignment in $OPT(c)$.
3. Remove \mathbf{x}^c from $OPT(c)$.

We can always make the assertion in Step 2 because of the WTA2 condition: if \mathbf{x}^c is the only assignment left in $OPT(c)$, and it satisfies the condition in Step 1, consider the subproblem tree \mathcal{T} with only one edge (c, c', φ) ; we have the WTA2 condition on \mathcal{T} violated. Also note that this process must end, because the total number of assignments in $\{OPT(c)\}$ is finite. Therefore it must end in a state where we cannot find a \mathbf{x}^c in Step 1, and each $OPT(c)$ is not empty. This $\{OPT(c)\}$ satisfies WTA. \square

Proof of Proposition 5

Proposition (STC Fixed-Point). *If $\{\Theta^c\}$ do not satisfy WTA2, we can always find \mathcal{T} such that applying STC to \mathcal{T} decreases the dual objective. If $\{\Theta^c\}$ satisfy WTA2, applying STC to any \mathcal{T} does not change the dual objective and preserves the WTA2 condition.*

Proof. If $\{\Theta^c\}$ do not satisfy WTA2, by definition we can find \mathcal{T} such that there is no $\{\mathbf{x}^c\} \sim \mathcal{T}$ that are optimal for each $c \in \mathcal{T}$. We therefore claim that:

$$\sum_{c \in \mathcal{T}} \max_{\mathbf{X}^c} \Theta^c(\mathbf{X}^c) > \max_{\{\mathbf{x}^c\} \sim \mathcal{T}} \sum_{c \in \mathcal{T}} \Theta^c(\mathbf{X}^c) \quad (43)$$

this is because if the equal sign holds, the maximizer of RHS must be optimal for each $c \in \mathcal{T}$ (Lemma 1). So applying STC to \mathcal{T} must increase the dual object because it does not change the sum over subproblems outside \mathcal{T} , and it decreases the sum over subproblem within \mathcal{T} . To see this we denote the new subproblem potentials as $\{\Theta_{new}^c\}$, and we have

$$\begin{aligned} \sum_{c \in \mathcal{T}} \max_{\mathbf{X}^c} \Theta_{new}^c(\mathbf{X}^c) &= \max_{\{\mathbf{x}^c\} \sim \mathcal{T}} \sum_{c \in \mathcal{T}} \Theta_{new}^c(\mathbf{X}^c) \\ &= \max_{\{\mathbf{x}^c\} \sim \mathcal{T}} \sum_{c \in \mathcal{T}} \Theta^c(\mathbf{X}) < \sum_{c \in \mathcal{T}} \max_{\mathbf{X}^c} \Theta^c(\mathbf{X}^c) \end{aligned} \quad (44)$$

where the first equality follows from Proposition 6; the second equality follows from Lemma 2; the last inequality uses (43).

To show that STC preserves the dual objective when WTA2 holds, let us assume that $\{\Theta^c\}$ satisfy WTA2. Consider applying STC to any \mathcal{T} . From the assumption we know that $\exists \{\mathbf{x}^c\} \sim \mathcal{T}$ that are optimal for each $c \in \mathcal{T}$ individually. Apply Lemma 1 we know that

$$\sum_{c \in \mathcal{T}} \max_{\mathbf{X}^c} \Theta^c(\mathbf{X}^c) = \max_{\{\mathbf{x}^c\} \sim \mathcal{T}} \sum_{c \in \mathcal{T}} \Theta^c(\mathbf{X}^c) \quad (45)$$

Again, let $\{\Theta_{new}^c\}$ be the subproblem potentials after applying STC to \mathcal{T} . It has the same dual objective

because:

$$\begin{aligned} \sum_{c \in \mathcal{T}} \max_{\mathbf{X}^c} \Theta_{new}^c(\mathbf{X}^c) &= \max_{\{\mathbf{x}^c\} \sim \mathcal{T}} \sum_{c \in \mathcal{T}} \Theta_{new}^c(\mathbf{X}^c) \\ &= \max_{\{\mathbf{x}^c\} \sim \mathcal{T}} \sum_{c \in \mathcal{T}} \Theta^c(\mathbf{X}) = \sum_{c \in \mathcal{T}} \max_{\mathbf{X}^c} \Theta^c(\mathbf{X}^c) \end{aligned} \quad (46)$$

following the same reasoning as above.

To show that STC preserves the WTA2 condition, we instead show that STC preserves the WTA condition (which is equivalent to WTA2 according to Proposition 12). Assuming $\{\Theta^c\}$ satisfy WTA with respect to the optimal assignment sets $\{OPT(c)\}$, it is easy to check that any single message passing step in the STC algorithm will not make any assignment in $\{OPT(c)\}$ non-optimal. Therefore WTA still holds with respect to the same $\{OPT(c)\}$ after STC. \square

Proof of Proposition 10

Proposition. *If all edges of \mathcal{T} have same scope φ , P-STC is equivalent to S-STC. Both achieve strong max-consistency when applied with uniform allocation weights $a_c = \frac{1}{|\mathcal{T}|}$*

Proof. It suffice to show that if all edges of \mathcal{T} have the same scope φ , and the allocation weights are uniform $a_c = \frac{1}{|\mathcal{T}|}$, after running STC we have

$$M_\varphi^c = \frac{1}{|\mathcal{T}|} \mu_\varphi^\mathcal{T} \quad (47)$$

where the max-marginal of subproblem tree \mathcal{T} over φ is defined as

$$\mu_\varphi^\mathcal{T}(\mathbf{x}_\varphi) = \max_{\forall c \in \mathcal{T}, \mathbf{X}_c|_\varphi = \mathbf{x}_\varphi} \sum_{c \in \mathcal{T}} \Theta^c(\mathbf{X}_c). \quad (48)$$

Note that the condition $\forall c \in \mathcal{T}, \mathbf{X}_c|_\varphi = \mathbf{x}_\varphi$ implies $\{\mathbf{X}_c\} \sim \mathcal{T}$ given that all edges of \mathcal{T} have the same scope φ , so this definition of tree max-marginal is consistent with the definition used in Def. 6.

We note that in this special case the STC algorithm is equivalent to the following:

1. Compute M_φ^c for all $c \in \mathcal{T}$. Note that all these max-marginals have identical scope.
2. Compute the average of these max-marginals \overline{M}_φ .
3. Add $\overline{M}_\varphi - M_\varphi^c$ to each subproblem.

This is because if we were to run MAP inference for each subproblem individually using variable elimination with φ being the variable(s) last eliminated, all

subproblem inference steps before reaching φ will not change with any change in \mathcal{B}^T . Therefore both S-STC and P-STC in this case essentially averages the max-marginals on φ . So the max-marginals on φ for all subproblems become \overline{M}_φ .

It remains to show that $\overline{M}_\varphi = \frac{1}{|\mathcal{T}|} \mu_\varphi^T$. This follows directly from Proposition 6 (the weak max-consistency property enforced by STC). \square

Proof of Proposition 8

Proposition. *Applying STC to the Bethe tree specified in Proposition 2 with allocation weights $a_f = 0.5$, $a_i = 0.5$ is equivalent to MSD.*

Proof. To prove the equivalence we use the update rules of MSD given in (Sontag et al., 2011). MSD computes messages $\delta_{fi}(X_i)$ and defines the reparametrization as follows:

$$\Theta^f = \Theta_0^f - \sum_{i \in f} \delta_{fi}(X_i) \quad (49)$$

$$\theta^i = \theta_0^i + \sum_{f \ni i} \delta_{fi}(X_i) \quad (50)$$

Note that these correspond to (6) in our notation. The MDS updates are given by

$$\delta_{fi} = -\frac{1}{2} \delta_i^{-f} + \frac{1}{2} \max_{\mathbf{X}^{f \setminus i}} [\Theta_0^f - \sum_{j \in f \setminus i} \delta_{fj}] \quad (51)$$

where $\delta_i^{-f} = \theta_0^i + \sum_{f' \ni i, f' \neq f} \delta_{f'i}$.

In STC we consider calibrating the Bethe tree specified in the proposition, and in MSD we consider updating δ_{fi} for a given f and a given $i \in f$. We show that they result in the same updates to subproblem potentials Θ^f and θ^i .

We denote MSD messages and subproblem potentials before updates as $\widehat{\delta}_{fi}$, $\widehat{\delta}_i^{-f}$, $\widehat{\theta}^i$ and $\widehat{\Theta}^f$. The STC upstream pass subtracts min-marginals from unary subproblems $i \in f$ and add to f . That is:

$$\theta^i \leftarrow 0, \quad (52)$$

$$\begin{aligned} \Theta^f &\leftarrow \Theta_0^f - \sum_{i \in f} \widehat{\delta}_{fi} + \theta_0^i + \sum_{f' \ni i} \widehat{\delta}_{f'i} \\ &= \Theta_0^f - \sum_{j \in f \setminus i} \widehat{\delta}_{fj} + \widehat{\delta}_i^{-f} \end{aligned} \quad (53)$$

Then we perform the downstream pass with allocation weights $a_i = 0.5$ and $a_f = 0.5$:

$$\begin{aligned} \theta^i &\leftarrow 0 + \frac{1}{2} \max_{\mathbf{X}^{f \setminus i}} [\Theta_0^f - \sum_{j \in f \setminus i} \widehat{\delta}_{fj} + \widehat{\delta}_i^{-f}] \\ &= \frac{1}{2} \widehat{\delta}_i^{-f} + \frac{1}{2} \max_{\mathbf{X}^{f \setminus i}} [\Theta_0^f - \sum_{j \in f \setminus i} \widehat{\delta}_{fj}] \end{aligned} \quad (54)$$

$$\begin{aligned} \Theta^f &\leftarrow \Theta_0^f - \sum_{j \in f \setminus i} \widehat{\delta}_{fj} + \widehat{\delta}_i^{-f} - \frac{1}{2} \max_{\mathbf{X}^{f \setminus i}} [\Theta_0^f - \sum_{j \in f \setminus i} \widehat{\delta}_{fj} + \widehat{\delta}_i^{-f}] \\ &= \Theta_0^f - \sum_{j \in f \setminus i} \widehat{\delta}_{fj} + \frac{1}{2} \widehat{\delta}_i^{-f} - \frac{1}{2} \max_{\mathbf{X}^{f \setminus i}} [\Theta_0^f - \sum_{j \in f \setminus i} \widehat{\delta}_{fj}] \end{aligned} \quad (55)$$

Now consider MSD. For both Θ^f and θ^i , only one message in their definitions is changed. Therefore

$$\begin{aligned} \theta^i &= \widehat{\theta}^i - \widehat{\delta}_{fi} + \delta_{fi} \\ &= \theta_0^i + \sum_{f' \ni i, f' \neq f} \widehat{\delta}_{f'i} - \frac{1}{2} \widehat{\delta}_i^{-f} + \frac{1}{2} \max_{\mathbf{X}^{f \setminus i}} [\Theta_0^f - \sum_{j \in f \setminus i} \widehat{\delta}_{fj}] \\ &= \frac{1}{2} \widehat{\delta}_i^{-f} + \frac{1}{2} \max_{\mathbf{X}^{f \setminus i}} [\Theta_0^f - \sum_{j \in f \setminus i} \widehat{\delta}_{fj}] \end{aligned} \quad (56)$$

which is the same as (54).

$$\begin{aligned} \Theta^f &= \widehat{\Theta}^f + \widehat{\delta}_{fi} - \delta_{fi} \\ &= \Theta_0^f - \sum_{j \in f} \widehat{\delta}_{fj} + \widehat{\delta}_{fi} + \frac{1}{2} \widehat{\delta}_i^{-f} - \frac{1}{2} \max_{\mathbf{X}^{f \setminus i}} [\Theta_0^f - \sum_{j \in f \setminus i} \widehat{\delta}_{fj}] \\ &= \Theta_0^f - \sum_{j \in f \setminus i} \widehat{\delta}_{fj} + \frac{1}{2} \widehat{\delta}_i^{-f} - \frac{1}{2} \max_{\mathbf{X}^{f \setminus i}} [\Theta_0^f - \sum_{j \in f \setminus i} \widehat{\delta}_{fj}] \end{aligned} \quad (57)$$

which is the same as (55). Hence we have proved the equivalence. \square

Proof of Proposition 2

Proposition. *Max-sum diffusion (MSD) (Werner, 2007) performs BCD on the Bethe dual (8). Each BCD step enforces strong max-consistency for a Bethe tree consisting of one $f \in \mathcal{F}$ and one $i \in f$.*

Proof. This directly follows from Proposition 10 and Proposition 8. \square

Proof of Proposition 9

Proposition. *Applying STC to the subproblem tree specified in Proposition 3 with uniform allocation weights $a_c = \frac{1}{|\mathcal{T}|}$ is equivalent to TRW-S.*

Proof. In proving Proposition 10, we noted that STC when applied to \mathcal{T} with all edges sharing the same scope φ is equivalent to averaging the max-marginals of all subproblems (trees in this case) on φ . This is exactly what TRW-S does as described in (Kolmogorov, 2006). \square

Proof of Proposition 3

Proposition. *TRW-S (Kolmogorov, 2006) performs BCD on the dual (7) where each subproblem $c \in \mathcal{C}$ is a tree (of the MRF). Each BCD step enforces strong max-consistency for a subproblem tree \mathcal{T} with nodes $\mathcal{V}_{\mathcal{T}} = \{c : \text{scope}(c) \supset \varphi\}$ for a given $\varphi \in V \cup E$, and edges (c, c', φ) that constitute a tree over $\mathcal{V}_{\mathcal{T}}$.*

Proof. This directly follows from Proposition 10 and Proposition 9. \square

Proof of Proposition 7

Proposition. *Applying P-STC to the Bethe tree specified in Proposition 4 with allocation weights $a_i = \frac{1}{|f|}$, $a_f = 0$ is equivalent to MPLP.*

Proof. Note that the equivalence only holds for P-STC but not S-STC in this case.

Max-product linear programming (MPLP) was originally proposed in (Globerson & Jaakkola, 2007). We use its generalized form as described in (Sontag et al., 2011). MPLP maintains messages $\delta_{fi}(X_i)$ to define the same reparametrization (49,50) as in MSD.

The MPLP updates are given by

$$\delta_{fi} = -\delta_i^{-f} + \frac{1}{|f|} \max_{\mathbf{X}^{f \setminus i}} [\Theta_0^f + \sum_{j \in f} \delta_j^{-f}] \quad (58)$$

where $\delta_i^{-f} = \theta_0^i + \sum_{f' \ni i, f' \neq f} \delta_{f'i}$.

In STC we consider calibrating the Bethe tree specified in the proposition, and in MPLP we consider updating δ_{fi} for a given f and all $i \in f$. We show that they result in the same updates to subproblem potentials Θ^f and θ^i .

We denote MPLP messages and dual variables before updates as $\widehat{\delta}_{fi}$, $\widehat{\delta}_i^{-f}$, $\widehat{\theta}^i$ and $\widehat{\Theta}^f$. The STC upstream pass subtracts min-marginals from unary subproblems $i \in f$ and add to f . That is:

$$\theta^i \leftarrow 0, \quad \forall i \in f \quad (59)$$

$$\begin{aligned} \Theta^f &\leftarrow \Theta_0^f - \sum_{i \in f} \widehat{\delta}_{fi} + \sum_{i \in f} [\theta_0^i + \sum_{f' \ni i} \widehat{\delta}_{f'i}] \\ &= \Theta_0^f + \sum_{i \in f} [\theta_0^i + \sum_{f' \ni i, f' \neq f} \widehat{\delta}_{f'i}] \\ &= \Theta_0^f + \sum_{i \in f} \widehat{\delta}_i^{-f} \end{aligned} \quad (60)$$

Then we perform the downstream pass (of P-STC) with allocation weights $a_i = 1/|f|$ and $a_f = 0$:

$$\theta^i \leftarrow 0 + \frac{1}{|f|} \max_{\mathbf{X}^{f \setminus i}} [\Theta_0^f + \sum_{j \in f} \widehat{\delta}_j^{-f}], \quad \forall i \in f \quad (61)$$

$$\Theta^f \leftarrow \Theta_0^f + \sum_{i \in f} \widehat{\delta}_i^{-f} - \sum_{i \in f} \frac{1}{|f|} \max_{\mathbf{X}^{f \setminus i}} [\Theta_0^f + \sum_{j \in f} \widehat{\delta}_j^{-f}] \quad (62)$$

Now consider MPLP. Substituting (58) into (50) we can see that the update to Θ^f is exactly (62). For θ^i we note that only one term in RHS of (49) has changed for all $i \in f$. Therefore

$$\begin{aligned} \theta^i &= \widehat{\theta}^i - \widehat{\delta}_{fi} + \delta_{fi} \\ &= \theta_0^i + \sum_{f' \ni i} \widehat{\delta}_{f'i} - \widehat{\delta}_{fi} - \delta_i^{-f} + \frac{1}{|f|} \max_{\mathbf{X}^{f \setminus i}} [\Theta_0^f + \sum_{j \in f} \widehat{\delta}_j^{-f}] \\ &= (\theta_0^i + \sum_{f' \ni i, f' \neq f} \widehat{\delta}_{f'i}) - \delta_i^{-f} + \frac{1}{|f|} \max_{\mathbf{X}^{f \setminus i}} [\Theta_0^f + \sum_{j \in f} \widehat{\delta}_j^{-f}] \\ &= \frac{1}{|f|} \max_{\mathbf{X}^{f \setminus i}} [\Theta_0^f + \sum_{j \in f} \widehat{\delta}_j^{-f}] \end{aligned} \quad (63)$$

where the last equation uses the fact that $\delta_i^{-f} = \widehat{\delta}_i^{-f}$ since the messages for $f' \neq f$ are not changed. Note that the result is the same as (61) so we have proved the equivalence. \square

Proof of Proposition 4

Proposition. *MPLP (Globerson & Jaakkola, 2007; Sontag et al., 2011) performs BCD on the Bethe dual (8). Each BCD step enforces MPLP max-consistency for a Bethe tree consisting of one $f \in \mathcal{F}$ and all $i \in f$.*

Proof. We have shown that MPLP is a special case of P-STC. So it suffice to show that P-STC, when applied to the setting specified in Proposition 7, enforces MPLP max-consistency.

Applying Def. 6 to the Bethe tree of MPLP, it suffice to show that

$$\theta^i = \frac{1}{|f|} \mu_i^T, \quad \Theta^f = \mu_f^T - \sum_{i \in f} \frac{1}{|f|} \mu_i^T \quad (64)$$

The result of applying P-STC to this setting is given by (61)(62). Comparing them to (64) and using the definition of tree max-marginals μ , it suffice to show that $\Theta_0^f + \sum_{i \in f} \widehat{\delta}_i^{-f}$ equals $\widehat{\Theta}^f + \sum_{i \in f} \widehat{\theta}^i$ (the ‘‘tree distribution’’ before updates). Indeed,

$$\begin{aligned} & \widehat{\Theta}^f + \sum_{i \in f} \widehat{\theta}^i \\ &= \Theta_0^f - \sum_{i' \in f} \widehat{\delta}_{f i'} + \sum_{i \in f} [\theta_0^i + \sum_{f' \ni i} \widehat{\delta}_{f' i}] \\ &= \Theta_0^f + \sum_{i \in f} [\theta_0^i + \sum_{f' \ni i, f' \neq f} \widehat{\delta}_{f' i}] \\ &= \Theta_0^f + \sum_{i \in f} \widehat{\delta}_i^{-f} \end{aligned} \quad (65)$$

□

Proof of Proposition 5

Proposition. *Tree-BCD (Sontag & Jaakkola, 2009) performs BCD on the Bethe dual (8) where all f are pairwise factors (i, j) . Each BCD step enforces MPLP max-consistency for a Bethe tree corresponding to a spanning tree of the MRF.*

Proof. In (Sontag & Jaakkola, 2009), the tree-BCD algorithm was described as updating the beliefs $f_T = \{f_i(x_i), f_{ij}(x_i, x_j) : ij \in T\}$ on a spanning tree T of the MRF as follows:

$$f_i(x_i) = \frac{1}{n} \log \mu_i(x_i) \quad (66)$$

$$\begin{aligned} f_{ij}(x_i, x_j) &= \log \mu_{ij}(x_i, x_j) - \frac{n_{j \rightarrow i}}{n} \log \mu_i(x_i) \\ &\quad - \frac{n_{i \rightarrow j}}{n} \log \mu_j(x_j) \end{aligned} \quad (67)$$

where $\log \mu_i$, $\log \mu_{ij}$ are (log)-max-marginals of the tree distribution over T , and $n_{j \rightarrow i}$ is the number of nodes in the subtree of node i with parent j . Note that T is a special case of subproblem tree when the dual decomposition has MRF edges and nodes as subproblems; $\log \mu_i$, $\log \mu_{ij}$, and $n_{j \rightarrow i}$ correspond to μ_i^T , μ_f^T , and N_i^f in Def. 6, respectively, Therefore the proposition follows directly from Def. 6. □