# Fast, Approximately Optimal Solutions for Single and Dynamic MRFs<sup>\*</sup>

Nikos Komodakis, Georgios Tziritas University of Crete, Computer Science Department

{komod,tziritas}@csd.uoc.gr

Abstract

A new efficient MRF optimization algorithm, called Fast-PD, is proposed, which generalizes  $\alpha$ -expansion. One of its main advantages is that it offers a substantial speedup over that method, e.g. it can be at least 3-9 times faster than  $\alpha$ -expansion. Its efficiency is a result of the fact that Fast-PD exploits information coming not only from the original MRF problem, but also from a dual problem. Furthermore, besides static MRFs, it can also be used for boosting the performance of dynamic MRFs, i.e. MRFs varying over time. On top of that, Fast-PD makes no compromise about the optimality of its solutions: it can compute exactly the same answer as  $\alpha$ -expansion, but, unlike that method, it can also guarantee an almost optimal solution for a much wider class of NP-hard MRF problems. Results on static and dynamic MRFs demonstrate the algorithm's efficiency and power. E.g., Fast-PD has been able to compute disparity for stereoscopic sequences in real time, with the resulting disparity coinciding with that of  $\alpha$ -expansion.

## 1. Introduction

Discrete MRFs are ubiquitous in computer vision, and thus optimizing them is a problem of fundamental importance. According to it, given a weighted graph  $\mathcal{G}$  (with nodes  $\mathcal{V}$ , edges  $\mathcal{E}$  and weights  $w_{pq}$ ), one seeks to assign a label  $x_p$  (from a discrete set of labels  $\mathcal{L}$ ) to each  $p \in \mathcal{V}$ , so that the following cost is minimized:

$$\sum_{p \in \mathcal{V}} \mathfrak{c}_p(x_p) + \sum_{(p,q) \in \mathcal{E}} w_{pq} d(x_p, x_q).$$
(1)

Here,  $c_p(\cdot)$ ,  $d(\cdot, \cdot)$  determine the singleton and pairwise MRF potential functions respectively.

Up to now, graph-cut based methods, like  $\alpha$ -expansion [3], have been very effective in MRF optimization, generating solutions with good optimality properties [8]. However, besides solutions' optimality, another important issue is that of computational efficiency. In fact, this issue has recently been looked at for the special case of dynamic MRFs [5, 4], *i.e.* MRFs varying over time. Thus, trying to concentrate on both of these issues here, we raise the following questions: can there be a graph-cut based method, which will be more efficient, but equally (or even more) powerful, than  $\alpha$ -expansion, for the case of single MRFs? Furthermore,

Nikos Paragios MAS, Ecole Centrale de Paris

can that method also offer a computational advantage for the case of dynamic MRFs? With respect to the questions raised above, this work makes the following contributions.

**Efficiency for single MRFs:**  $\alpha$ -expansion works by solving a series of max-flow problems. Its efficiency is thus largely determined from the efficiency of these max-flow problems, which, in turn, depends on the number of augmenting paths per max-flow. Here, we build upon recent work of [6], and propose a new primal-dual MRF optimization method, called Fast-PD. This method, like [6] or  $\alpha$ -expansion, also ends up solving a max-flow problem for a series of graphs. However, unlike these techniques, the graphs constructed by Fast-PD ensure that the number of augmentations per max-flow decreases dramatically over time, thus boosting the efficiency of MRF inference. To show this, we prove a generalized relationship between the number of augmentations and the so-called *primal-dual* gap associated with the original MRF problem and its dual. Furthermore, to fully exploit the above property, 2 new extensions are also proposed: an *adapted max-flow algorithm*, as well as an *incremental graph construction* method.

**Optimality properties:** Despite its efficiency, our method also makes no compromise regarding the optimality of its solutions. So, if  $d(\cdot, \cdot)$  is a metric, Fast-PD is as powerful as  $\alpha$ -expansion, *i.e.* it computes exactly the same solution, but with a substantial speedup. Moreover, it applies to a much wider class of MRFs<sup>1</sup>, *e.g.* even with a non-metric  $d(\cdot, \cdot)$ , while still guaranteeing an almost optimal solution.

Efficiency for dynamic MRFs: Furthermore, our method can also be used for boosting the efficiency of dynamic MRFs (introduced to computer vision in [5]). Two works have been proposed in this regard recently [5, 4]. These methods can be applied to dynamic MRFs that are binary or have convex priors. On the contrary, Fast-PD naturally handles a much wider class of dynamic MRFs, and can do so by also exploiting information from a problem, which is dual to the original MRF problem. Fast-PD can thus be thought of as a generalization of previous techniques.

The rest of the paper is organized as follows. In sec. 2, we briefly review the work of [6] about using the primaldual schema for MRF optimization. The Fast-PD algorithm is then described in sec. 3. Its efficiency for optimizing

<sup>\*</sup>This work was partially supported from the French ANR-Blanc grant SURF (2005-2008) and Platon (2006-2007).

<sup>&</sup>lt;sup>1</sup>Fast-PD requires only  $d(a, b) \ge 0, \ d(a, b) = 0 \Leftrightarrow a = b$ 

$1: [\mathbf{x}, \mathbf{y}] \leftarrow \texttt{INIT\_DUALS\_PRIMALS(}); \mathbf{x}_{old} \leftarrow \mathbf{x}$	6: $\mathbf{x} \leftarrow \mathbf{x}'; \ \mathbf{y} \leftarrow \mathbf{y}';$
2: for each label $c$ in $\mathcal{L}$ do	7: end for
3: $\mathbf{y} \leftarrow \text{PREEDIT}_\text{DUALS}(c, \mathbf{x}, \mathbf{y});$	8: if $\mathbf{x} \neq \mathbf{x}_{old}$ then
4: $[\mathbf{x}', \mathbf{y}'] \leftarrow \text{UPDATE}_\text{DUALS}_\text{PRIMALS}(c, \mathbf{x}, \mathbf{y});$	9: $\mathbf{x}_{old} \leftarrow \mathbf{x}; \text{ goto } 2;$
5: $\mathbf{y}' \leftarrow \text{POSTEDIT}_\text{DUALS}(c, \mathbf{x}', \mathbf{y}');$	10: end if

Fig. 1: The primal dual schema for MRF optimization.

single MRFs is further analyzed in sec. 4, where related results and some important extensions of Fast-PD are presented as well. Sec. 5 explains how Fast-PD can boost the performance of dynamic MRFs, and also contains more experimental results. Finally, we conclude in section 6.

#### 2. Primal-dual MRF optimization algorithms

In this section, we review very briefly the work of [6]. Consider the primal-dual pair of linear programs, given by:

PRIMAL: min 
$$\mathbf{c}^T \mathbf{x}$$
 DUAL: max  $\mathbf{b}^T \mathbf{y}$   
s.t.  $\mathbf{A}\mathbf{x} = \mathbf{b}, \mathbf{x} \ge \mathbf{0}$  s.t.  $\mathbf{A}^T \mathbf{y} \le \mathbf{c}$ 

One seeks an optimal primal solution, with the extra constraint of  $\mathbf{x}$  being integral. This makes for an NP-hard problem, and so one can only hope for finding an approximate solution. To this end, the following schema can be used:

**Theorem 1 (Primal-Dual schema).** *Keep generating pairs* of integral-primal, dual solutions  $(\mathbf{x}^{\mathbf{k}}, \mathbf{y}^{\mathbf{k}})$ , until the elements of the last pair, say  $\mathbf{x}, \mathbf{y}$ , are both feasible and have costs that are close enough, e.g. their ratio is  $\leq f_{app}$ :

$$\mathbf{c}^T \mathbf{x} \le f_{\mathrm{app}} \cdot \mathbf{b}^T \mathbf{y} \tag{2}$$

Then  $\mathbf{x}$  is guaranteed to be an  $f_{app}$ -approximate solution to the optimal integral solution  $x^*$ , i.e.  $\mathbf{c}^T \mathbf{x} \leq f_{app} \cdot \mathbf{c}^T \mathbf{x}^*$ .

The above schema has been used in [6], for deriving approximation algorithms for a very wide class of MRFs. To this end, MRF optimization was first cast as an equivalent integer program and then, as required by the primal-dual schema, its linear programming relaxation and its dual were derived. Based on these LPs, the authors then show that, for Theorem 1 to be true with  $f_{app} = 2 \frac{d_{max}^2}{d_{min}}^2$ , it suffices that the next (so-called *relaxed complementary slackness*) conditions hold true for the resulting primal and dual variables:

$$h_p(x_p) = \min_{a \in \mathcal{L}} h_p(a), \quad \forall p \in \mathcal{V}$$
(3)

$$y_{pq}(x_p) + y_{qp}(x_q) = w_{pq}d(x_p, x_q), \quad \forall pq \in \mathcal{E}$$
(4)

$$y_{pq}(a) + y_{qp}(b) \le 2w_{pq}d_{\max}, \quad \forall pq \in \mathcal{E}, a \in \mathcal{L}, b \in \mathcal{L}$$
 (5)

In these formulas, the primal variables, denoted by  $\mathbf{x} = \{x_p\}_{p \in \mathcal{V}}$ , determine the labels assigned to nodes (called *active labels* hereafter), *e.g.*  $x_p$  is the active label of node *p*. Whereas, the dual variables are divided into *balance* and *height* variables. There exist 2 balance variables  $y_{pq}(a), y_{qp}(a)$  per edge (p,q) and label *a*, as well as 1 height variable  $h_p(a)$  per node *p* and label *a*. Variables  $y_{pq}(a), y_{qp}(a)$  are also called *conjugate* and, for the dual solution to be feasible, these must be set opposite to each other, *i.e.*:  $y_{qp}(\cdot) \equiv -y_{pq}(\cdot)$ . Furthermore, the height variables are always defined in terms of the balance variables as follows:

$$h_p(\cdot) \equiv \mathfrak{c}_p(\cdot) + \sum_{q:qp \in \mathcal{E}} y_{pq}(\cdot).$$
 (6)

Note that, due to (6), only the vector y (of all balance variables) is needed for specifying a dual solution. In addition, for simplifying conditions (4),(5), one can also define:

10

$$\operatorname{pad}_{pq}(a,b) \equiv y_{pq}(a) + y_{qp}(b).$$
(7)

The primal-dual variables are iteratively updated until all conditions (3)-(5) hold true. The basic structure of a primal-dual algorithm can be seen in Fig. 1. During an inner *c*-iteration (lines 3-6 in Fig. 1), a label *c* is selected and a new primal-dual pair of solutions  $(\mathbf{x}', \mathbf{y}')$  is generated based on the current pair  $(\mathbf{x}, \mathbf{y})$ . To this end, among all balance variables  $y_{pq}(.)$ , only the balance variables of *c*-labels (*i.e.*  $y_{pq}(c)$ ) are updated during a *c*-iteration.  $|\mathcal{L}|$  such iterations (*i.e.* one *c*-iteration per label *c* in  $\mathcal{L}$ ) make up an outer iteration (lines 2-7 in Fig. 1), and the algorithm terminates if no change of label takes place at the current outer iteration.

During an inner iteration, the main update of the primal and dual variables takes place inside UP-DATE\_DUALS\_PRIMALS, and (as shown in [6]) this update reduces to solving a max-flow problem in an appropriate graph  $\mathcal{G}^c$ . Furthermore, the routines PREEDIT\_DUALS and POSTEDIT\_DUALS simply apply corrections to the dual variables before and after this main update, *i.e.* to variables **y** and **y'** respectively. Also, for simplicity's sake, note that we will hereafter refer to only one of the methods derived in [6], and this will be the so-called PD3<sub>a</sub> method.

#### 3. Fast primal-dual MRF optimization

The complexity of the  $PD3_a$  primal-dual method largely depends on the complexity of all max-flow instances (one instance per inner-iteration), which, in turn, depends on the number of augmentations per max-flow. So, for designing faster primal-dual algorithms, we first need to understand how the graph  $\mathcal{G}^c$ , associated with the max-flow problem at a *c*-iteration of  $PD3_a$ , is constructed. To this end, we also have to recall the following intuitive interpretation of the dual variables [6]: for each node p, a separate copy of all labels in  $\mathcal{L}$  is considered, and all these labels are represented as balls, which float at certain heights relative to a reference plane. The role of the height variables  $h_p(\cdot)$  is then to determine the balls' height (see Figure 2(a)). E.g. the height of label a at node p is given by  $h_p(a)$ . Also, expressions like "label a at p is below/above label b" imply  $h_p(a) \leq h_p(b)$ . Furthermore, balls are not static, but may move in pairs through updating pairs of conjugate balance variables. E.g., in Figure 2(a), label c at p is raised by  $+\delta$  (due to adding  $+\delta$ to  $y_{pq}(c)$ ), and so label c at q has to move down by  $-\delta$  (due to adding  $-\delta$  to  $y_{qp}(c)$  so that condition  $y_{pq}(c) = -y_{qp}(c)$ still holds). Therefore, the role of balance variables is to raise or lower labels. In particular, the value of balance variable  $y_{pa}(a)$  represents the partial raise of label a at p due to edge pq, while (by (6)) the total raise of a at p equals the sum of partial raises from all edges of  $\mathcal{G}$  incident to p.

 $<sup>^{2}</sup>d_{\max} \equiv \max_{a \neq b} d(a, b), \ d_{\min} \equiv \min_{a \neq b} d(a, b)$ 



Fig. 2: (a) Dual variables' visualization for a simple MRF with 2 nodes  $\{p, q\}$  and 2 labels  $\{a, c\}$ . A copy of labels  $\{a, c\}$  exists for every node, and all these labels are represented by balls floating at certain heights. The role of the *height variables*  $h_{p}(\cdot)$  is to specify exactly these heights. Furthermore, balls are not static, but may move (i.e. change their heights) in pairs by updating conjugate *balance variables.* E.g., here, ball c at p is pulled up by  $+\delta$  (due to increasing  $y_{pq}(c)$  by  $+\delta$ ) and so ball c at q moves down by  $-\delta$  (due to decreasing  $y_{qp}(c)$  by  $-\delta$ ). Active labels are drawn with a thicker circle. (b) If label c at p is below  $x_p$ , then (due to (3)) we want label c to raise and reach  $x_p$ . We thus connect node p to the source  $\mathfrak{s}$  with an edge  $\mathfrak{s}p$  (*i.e.* p is an  $\mathfrak{s}$ -linked node), and flow  $f_{sp}$  represents the total raise of c (we also set  $\operatorname{cap}_{\mathfrak{s}p} = h_p(x_p) - h_p(c)$ . (c) If label c at p is above  $x_p$ , then (due to (3)) we want label c not to go below  $x_p$ . We thus connect node p to the sink t with edge pt (i.e. p is a t-linked node), and flow  $f_{pt}$  represents the total decrease in the height of c (we also set  $\operatorname{cap}_{pt} = h_p(c) - h_p(x_p)$  so that c will still remain above  $x_p$ ).

Hence,  $PD3_a$  tries to iteratively move labels up or down, until all conditions (3)-(5) hold true. To this end, it uses the following strategy: it ensures that conditions (4)-(5) hold at each iteration (which is always easy to do) and is just left with the main task of making the labels' heights satisfy condition (3) as well in the end (which is the most difficult part, requiring each active label  $x_p$  to be the lowest label for p). For this purpose, labels are moved in groups. In particular, during a c-iteration, only the c-labels are allowed to move. Furthermore, it was shown in [6] that the movement of all *c*-labels (*i.e.* the update of dual variables  $y_{pq}(c)$  and  $h_p(c)$ for all p, q can be simulated by pushing the maximum flow through a directed graph  $\mathcal{G}^c$  (which is constructed based on the current primal-dual pair  $(\mathbf{x}, \mathbf{y})$  at a *c*-iteration). The nodes of  $\mathcal{G}^c$  consist of all nodes of graph  $\mathcal{G}$  (the *internal* nodes), plus 2 external nodes, the source s and the sink t. In addition, all nodes of  $\mathcal{G}^c$  are connected by two types of edges: interior and exterior edges. Interior edges come in pairs pq, qp (with one such pair for every 2 neighbors p, q in  $\mathcal{G}$ ), and are responsible for updating the balance variables. In particular, the flows  $f_{pq}/f_{qp}$  of these edges represent the increase/decrease of balance variable  $y_{pq}(c)$ , *i.e.*  $y'_{pq}(c) =$  $y_{pq}(c) + f_{pq} - f_{qp}$ . Also, as we shall see, the capacities of interior edges are used together with PREEDIT\_DUALS, POSTEDIT\_DUALS to impose conditions (4), (5).

But for now, in order to understand how to make a faster primal-dual method, it is the exterior edges (which are in charge of the update of height variables), as well as their capacities (which are used for imposing the remaining condition (3)), that are of interest to us. The reason is that these edges determine the number of s-linked nodes, which, in turn, affects the number of augmenting paths per max-flow. In particular, each internal node connects to either the source  $\mathfrak{s}$  (*i.e.* it is an  $\mathfrak{s}$ -linked node) or to the sink t (i.e. it is a t-linked node) through one of these exterior edges, and this is done (with the goal of ensuring (3)) as follows: if label c at p is above  $x_p$  during a c-iteration (*i.e.*  $h_p(c) > h_p(x_p)$ ), then label c should not go below  $x_p$ , or else (3) will be violated for p. Node p thus connects to t through directed edge pt (i.e. p becomes t-linked), and flow  $f_{pt}$  represents the total decrease in the height of c after UPDATE\_DUALS\_PRIMALS, *i.e.*  $h'_p(c) = h_p(c) - f_{pt}$  (see Fig. 2(c)). Furthermore, the capacity of pt is set so that label cwill still remain above  $x_p$ , *i.e.*  $\operatorname{cap}_{pt} = h_p(c) - h_p(x_p)$ . On the other hand, if label c at p is below active label  $x_p$  (*i.e.*  $h_p(c) < h_p(x_p)$ ), then (due to (3)) label c should raise so as to reach  $x_p$ , and so p connects to s through edge  $\mathfrak{s}p$  (*i.e.* p becomes s-linked), while flow  $f_{sp}$  represents the total raise of ball c, i.e.  $h'_p(c) = h_p(c) + f_{\mathfrak{s}p}$  (see Fig. 2(b)). In this case, we also set  $\operatorname{cap}_{\mathfrak{s}p} = h_p(x_p) - h_p(c)$ .

This way, by pushing flow through the exterior edges of  $\mathcal{G}^c$ , all *c*-labels that are strictly below an active label try to raise and reach that label during UPDATE\_DU-ALS\_PRIMALS<sup>3</sup>. Not only that, but the fewer are the *c*-labels below an active label (*i.e.* the fewer are the *s*-linked nodes), the fewer will be the edges connected to the source, and thus the less will be the number of possible augmenting paths. In fact, the algorithm terminates when, for any label *c*, there are no more *c*-labels strictly below an active label (*i.e.* no *s*-linked nodes exist and thus no augmenting paths may be found), in which case condition (3) will finally hold true, as desired. Put another way, UPDATE\_DUALS\_PRIMALS tries to push *c*-labels (which are at a low height) up, so that the number of *s*-linked nodes is reduced and thus fewer augmenting paths may be possible for the next iteration.

However, although UPDATE\_DUALS\_PRIMALS tries to reduce the number of  $\mathfrak{s}$ -linked nodes (by pushing the maximum amount of flow), PREEDIT\_DUALS or POSTEDIT\_DU-ALS very often spoil that progress. As we shall see later, this occurs because, in order to restore condition (4) (which is their main goal), these routines are forced to apply corrections to the dual variables (*i.e.* to the labels' height). This is abstractly illustrated in Figure 3, where, as a result of pushing flow, a *c*-label initially managed to reach an active label  $x_p$ , but it again dropped below  $x_p$ , due to some correction applied by these routines. In fact, as one can show, the only point where a new  $\mathfrak{s}$ -linked node can be created is during either PREEDIT\_DUALS or POSTEDIT\_DUALS.

<sup>&</sup>lt;sup>3</sup>Equivalently, if *c*-label at *p* cannot raise high enough to reach  $x_p$ , UPDATE\_DUALS\_PRIMALS then assigns that *c*-label as the new active label of *p* (*i.e.*  $x'_p = c$ ), thus effectively making the active label go down. This helps condition (3) to become true, and forms the main rationale behind the update of the primal variables **x** in UPDATE\_DUALS\_PRIMALS.



**Fig. 3:** (a) Label c at p is below  $x_p$ , and thus label c is allowed to raise itself in order to reach  $x_p$ . This means that p will be an s-linked node of graph  $\mathcal{G}^c$ , *i.e.*  $\operatorname{cap}_{sp} > 0$ , and thus a non-zero flow  $f_{sp}$  (representing the total raise of label c) may pass through edge sp. Therefore, in this case, edge sp may become part of an augmenting path during max-flow. (b) After UPDATE\_DUALS\_PRIMALS, label c has managed to raise by  $f_{sp}$ and reach  $x_p$ . Since it cannot go higher than that, no flow can pass through edge sp, *i.e.*  $\operatorname{cap}_{sp} = 0$ , and so no augmenting path may traverse that edge thereafter. (c) However, due to some correction applied to c-label's height, label c has dropped below  $x_p$  once more and p has become an s-linked node again (*i.e.*  $\operatorname{cap}_{sp} > 0$ ). Edge sp can thus be part of an augmenting path again (as in (a)).

To fix this problem, we will redefine PREEDIT\_DUALS, POSTEDIT\_DUALS so that they can now ensure condition (4) by using just a minimum amount of corrections for the dual variables, (*e.g.* by touching these variables only rarely). To this end, however, UPDATE\_DUALS\_PRIMALS needs to be modified as well. The resulting algorithm, called Fast-PD, carries the following main differences over PD3<sub>a</sub> during a *c*-iteration (its pseudocode appears in Fig. 4):

- the new PREEDIT\_DUALS modifies a pair  $y_{pq}(c), y_{qp}(c)$ of dual variables only when absolutely necessary. So. whereas the previous version modified these variables (thereby changing the height of a c-label) whenever  $c \neq x_p$ ,  $c \neq x_q$  (which could happen extremely often), a modification is now applied only if  $load_{pq}(c, x_q) > w_{pq}d(c, x_q)$  or  $load_{pq}(x_p, c) > w_{pq}d(x_p, c)$  (which, in practice, happens much more rarely). In this case, a modification is needed (see code in Fig. 4), because the above inequalities indicate that condition (4) will be violated if either  $(c, x_q)$  or  $(x_p, c)$ become the new active labels for p, q. On the contrary, no modification is needed if the following inequalities are true:  $load_{pq}(c, x_q) < w_{pq}d(c, x_q), \ load_{pq}(x_p, c) < w_{pq}d(x_p, c),$ because then, as we shall see below, the new UP-DATE\_DUALS\_PRIMALS can always restore (4) (i.e. even if  $(c, x_q)$  or  $(x_p, c)$  are the next active labels - e.g., see (12)). In fact, the modification to  $y_{pq}(c)$  that is occasionally applied by the new PREEDIT\_DUALS can be shown to be the minimal correction that restores exactly the above inequalities (assuming, of course, this restoration is possible).

- Similarly, the new POSTEDIT\_DUALS modifies<sup>4</sup> balance variables  $y'_{pq}(x'_p)$  (with  $x'_p = c$ ) and  $y'_{qp}(x'_q)$  (with  $x'_q = c$ ) only if the inequality  $\log d'_{pq}(x'_p, x'_q) > w_{pq}d(x'_p, x'_q)$  holds, in which case POSTEDIT\_DUALS simply has to



## Fig. 4: Fast-PD's pseudocode.

reduce  $\operatorname{load}'_{pq}(x'_p, x'_q)$  for restoring (4). However, this inequality will hold true very rarely (*e.g.* for a metric  $d(\cdot, \cdot)$ , one may show that it can never hold), and so POSTEDIT\_DU-ALS will modify a *c*-balance variable (thereby changing the height of a *c*-label) only in very seldom occasions.

- But, to allow for the above changes, we also need to modify the construction of graph  $\mathcal{G}^c$  in UPDATE\_DU-ALS\_PRIMALS. In particular, for  $c \neq x_p$  and  $c \neq x_q$ , the capacities of interior edges pq, qp must now be set as follows:<sup>5</sup>

$$\operatorname{cap}_{pq} = \left[ w_{pq} d(c, x_q) - \operatorname{load}_{pq}(c, x_q) \right]^{+}, \qquad (8)$$

$$\operatorname{cap}_{qp} = \left\lfloor w_{pq}d(x_p,c) - \operatorname{load}_{pq}(x_p,c) \right\rfloor^{+}, \qquad (9)$$

where  $[x]^+ \equiv \max(x, 0)$ . Besides ensuring (5) (by not letting the balance variables increase too much), the main rationale behind the above definition of interior capacities is to also ensure that (after max-flow) condition (4) will be met by most pairs (p,q), even if  $(c, x_q)$  or  $(x_p, c)$  are the next labels assigned to them (which is a good thing, since we will thus manage to avoid the need for a correction by POSTEDIT\_DUALS for all but a few p,q). For seeing this, the crucial thing to observe is that if, say,  $(c, x_q)$  are the next labels for p and q, then capacity  $\operatorname{cap}_{pq}$  can be shown to represent the increase of  $\operatorname{load}_{pq}(c, x_q)$  after max-flow, *i.e.*:

$$\operatorname{load}_{pq}'(c, x_q) = \operatorname{load}_{pq}(c, x_q) + \operatorname{cap}_{pq}.$$
 (10)

Hence, if the following inequality is true as well:

1

$$\operatorname{oad}_{pq}(c, x_q) \le w_{pq} d(c, x_q) , \qquad (11)$$

then condition (4) will do remain valid after max-flow, as the following trivial derivation shows:

$$\log d'_{pq}(c, x_q) \stackrel{(10),(8)}{=} \log d_{pq}(c, x_q) + [w_{pq}d(c, x_q) - \log d_{pq}(c, x_q)]^{-1}$$

$$\stackrel{(11)}{=} w_{pq}d(c, x_q)$$
(12)

But this means that a correction may need to be applied by POSTEDIT\_DUALS only for pairs p, q violating (11) (before max-flow). However, such pairs tend to be very rare in practice (*e.g.*, as one can prove, no such pairs exist when  $d(\cdot, \cdot)$  is a metric), and thus very few corrections need to take place.

Fig. 5 summarizes how Fast-PD sets the capacities for all edges of  $\mathcal{G}^c$ . As already explained, the interior capacities (with the help of PREEDIT\_DUALS, POSTEDIT\_DUALS

<sup>&</sup>lt;sup>4</sup>We recall that POSTEDIT\_DUALS may modify only dual solution  $\mathbf{y}'$ . For that solution, we define  $\log d'_{pq}(a,b) \equiv y'_{pq}(a) + y'_{qp}(b)$ , as in (7).

<sup>&</sup>lt;sup>5</sup>If  $c = x_p$  or  $c = x_q$ , then  $cap_{pq} = cap_{qp} = 0$  as before, *i.e.* as in PD3<sub>a</sub>.

in a few cases) allow UPDATE\_DUALS\_PRIMALS to impose conditions (4),(5), while the exterior capacities allow UP-DATE\_DUALS\_PRIMALS to impose condition (3). As a result, the next theorem holds (see [1] for a complete proof):

**Theorem 2.** The last primal-dual pair  $(\mathbf{x}, \mathbf{y})$  of Fast-PD satisfies (3)-(5), and so  $\mathbf{x}$  is an  $f_{app}$ -approximate solution.

In fact, Fast-PD maintains all good optimality properties of the PD3<sub>a</sub> method. *E.g.*, for a metric  $d(\cdot, \cdot)$ , Fast-PD proves to be as powerful as  $\alpha$ -expansion (see [1]):

**Theorem 3.** If  $d(\cdot, \cdot)$  is a metric, then the Fast-PD algorithm computes the best *c*-expansion after any *c*-iteration.

#### 4. Efficiency of Fast-PD for single MRFs

But, besides having all these good optimality properties, a very important advantage of Fast-PD over all previous primal-dual methods, as well as  $\alpha$ -expansion, is that it proves to be much more efficient in practice.

In fact, the computational efficiency for all methods of this kind is largely determined from the time taken by each max-flow problem, which, in turn, depends on the number of augmenting paths that need to be computed. For the case of Fast-PD, the number of augmentations per inner-iteration decreases dramatically, as the algorithm progresses. E.g. Fast-PD has been applied to the problem of image restoration, and fig. 7 contains a related result about the denoising of a corrupted (with gaussian noise) "penguin" image (256 labels and a truncated quadratic distance  $d(a,b) = \min(|a-b|^2, D)$  - where D = 200 - has been used in this case). Also, fig. 8(a) shows the corresponding number of augmenting paths per outer-iteration (*i.e.* per group of  $|\mathcal{L}|$  inner-iterations). Notice that, for both  $\alpha$ -expansion, as well as  $PD3_a$ , this number remains very high (*i.e.* almost over  $2 \cdot 10^6$  paths) throughout all iterations. On the contrary, for the case of Fast-PD, it drops towards zero very quickly, e.g. only 4905 and 7 paths had to be found during the 8<sup>th</sup> and last outer-iteration respectively (obviously, as also shown in Fig. 9(a), this directly affects the total time needed per outer-iteration). In fact, for the case of Fast-PD, it is very typical that, after very few inner-iterations, no more than 10 or 20 augmenting paths need to be computed per max-flow, which really boosts the performance in this case.

This property can be explained by the fact that Fast-PD maintains both a primal, as well as a dual solution throughout its execution. Fast-PD then manages to effectively use the dual solutions of previous inner iterations, so as to reduce the number of augmenting paths for the next inneriterations. Intuitively, what happens is that Fast-PD ultimately wants to close the gap between the primal and the

exterior capacities	interior capacities		
$\operatorname{cap}_{\mathfrak{sp}} = [h_p(x_p) - h_p(c)]^+$	$x_p \neq c cap_{pq} = [w_{pq}d(c,x_q)-load_{pq}(c,x_q)]^+$	$x_{p} = c cap_{pq}$	= (
$cap_{pt} = [h_p(c) - h_p(x_p)]^+$	$x_q \neq c$ cap <sub>qp</sub> =[ $w_{pq}d(x_p,c)$ -load <sub>pq</sub> ( $x_p,c$ )] <sup>+</sup>	$x_q = c cap_{qp}$	= (

**Fig. 5:** Capacities of graph  $\mathcal{G}^c$ , as set by Fast-PD.



**Fig. 6:** (a) Fast-PD generates pairs of primal-dual solutions iteratively, with the goal of always reducing the primal-dual gap (*i.e.* the gap between the resulting primal and dual costs). But, for the case of Fast-PD, this gap can be viewed as a rough estimate for the number of augmentations, and so this number is forced to reduce over time as well. (b) On the contrary,  $\alpha$ -expansion works only in the primal domain (*i.e.* it is as if a fixed dual cost is used at the start of each new iteration) and thus the primal-dual gap can never become small enough. Therefore, no significant reduction in the number of augmentations takes place as the algorithm progresses.

dual cost (see Theorem 1), and, for this, it iteratively generates primal-dual pairs, with the goal of decreasing the size of this gap (see Fig. 6(a)). But, for Fast-PD, the gap's size can be thought of as, roughly speaking, an upper-bound for the number of augmenting paths per inner-iteration. Since, furthermore, Fast-PD manages to reduce this gap at any time throughout its execution, the number of augmenting paths is forced to decrease over time as well.

On the contrary, a method like  $\alpha$ -expansion, that works only in the primal domain, ignores dual solutions completely. It is, roughly speaking, as if  $\alpha$ -expansion is resetting the dual solution to zero at the start of each inner-iteration, thus effectively forgetting that solution thereafter (see Fig. 6(b)). For this reason, it fails to reduce the primal-dual gap and thus also fails to achieve a reduction in path augmentations over time, *i.e.* across inneriterations. But the  $PD3_a$  algorithm as well fails to mimic Fast-PD's behavior (despite being a primal-dual method). As explained in sec. 3, this happens because, in this case, PREEDIT\_DUAL and POSTEDIT\_DUAL temporarily destroy the gap just before the start of UPDATE\_DUALS\_PRIMALS, *i.e.* just before max-flow is about to begin computing the augmenting paths. (Note, of course, that this destruction is only temporary, and the gap is restored again after the execution of UPDATE\_DUALS\_PRIMALS).

The above mentioned relationship between primal-dual gap and number of augmenting paths is formally described in the next theorem (see [1] for a complete proof):

**Theorem 4.** For Fast-PD, the primal-dual gap at the current inner-iteration forms an approximate upper bound for the number of augmenting paths at each iteration thereafter. **Sketch of proof.** During a *c*-iteration, it can be shown that dual-cost  $\leq \sum_{p} \min(h_p(c), h_p(x_p))$ , whereas primal-cost=  $\sum_{p} h_p(x_p)$ , and so the primal-dual gap upper-bounds the following quantity:  $\sum_{p} [h_p(x_p) - h_p(c)]^+ = \sum_{p} \operatorname{cap}_{sp}$ .



**Fig. 7: Left:** "Tsukuba" image and its disparity by Fast-PD. **Middle:** a "SRI tree" image and corresponding disparity by Fast-PD. **Right:** noisy "penguin" image and its restoration by Fast-PD.

But this quantity obviously forms an upper-bound on the maximum flow, which, in turn, upper-bounds the number of augmentations (assuming integral flows).

Due to the above mentioned property, the time per outer-iteration decreases dramatically over time. This has been verified experimentally with virtually all problems that Fast-PD has been tested on. E.g. Fast-PD has been also applied to the problem of stereo matching, and fig. 7 contains the resulting disparity (of size  $384 \times 288$  with 16 labels) for the well-known "Tsukuba" stereo pair, as well as the resulting disparity (of size 256×233 with 10 labels) for an image pair from the well-known "SRI tree" sequence (in both cases, a truncated linear distance  $d(a,b) = \min(|a-b|, D)$  - with D=2 and D=5 - has been used, while the weights  $w_{pq}$  were allowed to vary based on the image gradient at p). Figures 9(b), 9(c) contain the corresponding running times per outer iteration. Notice how much faster the outer-iterations of Fast-PD become as the algorithm progresses, e.g. the last outer-iteration of Fast-PD (for the "SRI-tree" example) lasted less than 1 msec (since, as it turns out, only 4 augmenting paths had to be found during that iteration). Contrast this with the behavior of either the  $\alpha$ -expansion or the PD3<sub>a</sub> algorithm, which both require an almost constant amount of time per outer-iteration, e.g. the last outer-iteration of  $\alpha$ -expansion needed more than 0.4 secs to finish (i.e. it was more than 400 times slower than Fast-PD's iteration!). Similarly, for the "Tsukuba" example,  $\alpha$ -expansion's last outer-iteration was more than 2000 times slower than Fast-PD's iteration.

**Max-flow algorithm adaptation:** However, for fully exploiting the decreasing number of path augmentations and reduce the running time, we had to properly adapt the max-flow algorithm. To this end, the crucial thing to observe was that the decreasing number of augmentations was directly related to the decreasing number of  $\mathfrak{s}$ -linked nodes, as already explained in sec. 3. *E.g.* fig. 8(b) shows how the number of  $\mathfrak{s}$ -linked nodes varies per outer-iteration for the "penguin" example (with a similar behavior being observed for the other examples as well). As can be seen, this number decreases drastically over time. In fact, as



**Fig. 8:** (a) Number of augmenting paths per outer iteration for the "penguin" example (similar results hold for the other examples as well). Only in the case of Fast-PD, this number decreases dramatically over time. (b) This property of Fast-PD is directly related to the decreasing number of  $\mathfrak{s}$ -linked nodes per outer-iteration (this number is shown here for the same example as in (a)).



**Fig. 9:** Total time per outer iteration for the (**a**) "penguin", (**b**) "Tsukuba" and (**c**) "SRI tree" examples. (**d**) Total running times. For all experiments of this paper, a 1.6GHz laptop has been used.

implied by condition (3), no  $\mathfrak{s}$ -linked nodes will finally exist upon the algorithm's termination. Any augmentation-based max-flow algorithm striving for computational efficiency, should certainly exploit this property when trying to extract its augmenting paths. The most efficient of these algorithms [2] maintains 2 search trees for the fast extraction of these paths, a source and a sink tree. Here, the source tree will start growing by exploring non-saturated edges that are adjacent to s-linked nodes, whereas the sink tree will grow starting from all t-linked nodes. Of course, the algorithm terminates when no adjacent unsaturated edges can be found any more. However, in our case, maintaining the sink tree is completely inefficient and does not exploit the much smaller number of s-linked nodes. We thus propose maintaining only the source tree during max-flow, which will be a much cheaper thing to do here (e.g., in many inner iterations, there can be fewer than 10 s-linked nodes, but many thousands of t-linked nodes). Moreover, due to the small size of the source tree, detecting the termination of the max-flow procedure can now be done a lot faster, *i.e.* with-



**Fig. 10:** Suboptimality bounds per inner iteration (for "Tsukuba" and "penguin"). These bounds drop to 1 very fast, meaning that the corresponding solutions have become almost optimal very early.

out having to fully expand the large sink tree (which is a very costly operation), thus giving a substantial speedup. In addition to that, for efficiently building the source tree, we keep track of all  $\mathfrak{s}$ -linked nodes and don't recompute them from scratch each time. In our case, this tracking can be done without cost, since, as explained in sec. 3, an  $\mathfrak{s}$ -linked node can be created only inside the PREEDIT\_DUALS or the POSTEDIT\_DUALS routine, and thus can be easily detected. The above simple strategy has been extremely effective for boosting the performance of max-flow, especially when a small number of augmentations were needed.

**Incremental graph construction:** But besides the maxflow algorithm adaptation, we may also modify the way graph  $\mathcal{G}^c$  is constructed. *I.e.* instead of constructing the capacitated graph  $\mathcal{G}^c$  from scratch each time, we also propose an incremental way of setting its capacities. The following lemma turns out to be crucial in this regard:

**Lemma 1.** Let  $\mathcal{G}^c$ ,  $\mathcal{G}^c$  be the graphs for the current and previous c-iteration. Let also p, q be 2 neighboring MRF nodes. If, during the interval from the previous to the current c-iteration, no change of label took place for p and q, then the capacities of the interior edges pq, qp in  $\mathcal{G}^c$  and of the exterior edges  $\mathfrak{s}p, \mathfrak{p}t, \mathfrak{s}q, qt$  in  $\mathcal{G}^c$  equal the residual capacities of the corresponding edges in  $\overline{\mathcal{G}}^c$ .

The proof follows directly from the fact that if no change of label took place for p, q, then none of the height variables  $h_p(x_p), h_q(x_q)$  or the balance variables  $y_{pq}(x_p), y_{qp}(x_q)$ could have changed. Due to lemma 1, for building graph  $\mathcal{G}^c$ , we can simply reuse the residual graph of  $\overline{\mathcal{G}}^c$  and only recompute those capacities of  $\mathcal{G}^c$  for which the above lemma does not hold, thus speeding-up the algorithm even further.

**Combining speed with optimality:** Fig. 9(d) contains the running times of Fast-PD for various MRF problems. As can be seen from that figure, Fast-PD proves to be much faster than either the  $\alpha$ -expansion<sup>6</sup> or the PD3<sub>a</sub> method, *e.g.* Fast-PD has been more than 9 times faster than  $\alpha$ -expansion for the case of the "penguin" image (17.44 secs vs 173.1 secs). In fact, this behavior is a typical one, since Fast-PD has consistently provided at least a 3-9 times speedup for all the problems it has been tested on. However, besides its efficiency, Fast-PD does not make any compromise regarding the optimality of its solutions. On one hand, this is ensured by theorems 2, 3. On the other hand, Fast-PD, like

$[\mathbf{x}, \mathbf{y}] \leftarrow \text{INIT_DUALS_PRIMALS}(\bar{\mathbf{x}}, \bar{\mathbf{y}}): \mathbf{x} \leftarrow \bar{\mathbf{x}}; \mathbf{y} \leftarrow \bar{\mathbf{y}};$	
$\forall pq, y_{pq}(x_p) \mathrel{+}= w_{pq}d(x_p, x_q) - \bar{w}_{pq}\bar{d}(x_p, x_q);$	
$\forall p, h_p(\cdot) \mathrel{+}= \mathfrak{c}_p(\cdot) - \overline{\mathfrak{c}}_p(\cdot);$	

Fig. 11: Fast-PD's new pseudocode for dynamic MRFs.

any other primal-dual method, can also tell for free how well it performed by always providing a per-instance suboptimality bound for its solution. This comes at no extra cost, since any ratio between the cost of a primal solution and the cost of a dual solution can form such a bound. *E.g.* fig. 10 shows how these ratios vary per inner-iteration for the "tsukuba" and "penguin" problems (with similar results holding for the other problems as well). As one can notice, these ratios drop to 1 very quickly, meaning that an almost optimal solution has already been estimated even after just a few iterations (and despite the problem being NP-hard).

## 5. Dynamic MRFs

But, besides single MRFs, Fast-PD can be easily adapted to also boost the efficiency for dynamic MRFs [5], *i.e.* MRFs varying over time, thus showing the generality and power of the proposed method. In fact, Fast-PD fits perfectly to this task. The implicit assumption here is that the change between successive MRFs is small, and so, by initializing the current MRF with the final (primal) solution of the previous MRF, one expects to speed up inference. A significant advantage of Fast-PD in this regard, however, is that it can exploit not only previous MRF's primal solution (say  $\bar{x}$ ), but also its dual solution (say  $\bar{y}$ ). And this, for initializing current MRF's both primal and dual solutions (say x, y).

Obviously, for initializing x, one can simply set  $x = \bar{x}$ . Regarding the initialization of y, however, things are slightly more complicated. For maintaining Fast-PD's optimality properties, it turns out that, after setting  $y = \bar{y}$ , a slight correction still needs to be applied to y. In particular, Fast-PD requires its initial solution y to satisfy condition (4), *i.e.*  $y_{pq}(x_p) + y_{qp}(x_q) = w_{pq}d(x_p, x_q)$ , whereas  $\bar{\mathbf{y}}$  satisfies  $\bar{y}_{pq}(x_p) + \bar{y}_{qp}(x_q) = \bar{w}_{pq}\bar{d}(x_p, x_q)$ , *i.e.* condition (4) with  $w_{pq}d(\cdot,\cdot)$  replaced by the pairwise potential  $\bar{w}_{pq}\bar{d}(\cdot,\cdot)$  of the previous MRF. The solution for fixing that is very simple: e.g. we can simply set  $y_{pq}(x_p) + = w_{pq}d(x_p, x_q) - \bar{w}_{pq}\bar{d}(x_p, x_q)$ . Finally, for taking into account the possibly different singleton potentials between successive MRFs, the new heights will obviously need to be updated as  $h_p(\cdot) + = \mathfrak{c}_p(\cdot) - \overline{\mathfrak{c}}_p(\cdot)$ , where  $\overline{\mathfrak{c}}_p(\cdot)$  are the singleton potentials of the previous MRF. These are the only changes needed for the case of dynamic MRFs, and thus the new pseudocode appears in Fig. 11.

As expected, for dynamic MRFs, the speedup provided by Fast-PD is even greater than single MRFs. *E.g.* Fig. 12(a) shows the running times per frame for the "SRI tree" image sequence. Fast-PD proves to be be more than 10 times faster than  $\alpha$ -expansion in this case (requiring on average 0.22 secs per frame, whereas  $\alpha$ -expansion required 2.28 secs on average). Fast-PD can thus run on about 5

<sup>&</sup>lt;sup>6</sup>Since  $\alpha$ -expansion cannot be used if  $d(\cdot, \cdot)$  is not a metric, the method proposed in [7] had to be used for the cases of a non-metric  $d(\cdot, \cdot)$ .



**Fig. 12:** Statistics for the "SRI tree" sequence.

frames/sec, *i.e.* it can do stereo matching almost in real time for this example (in fact, if successive MRFs bear greater similarity, even much bigger speedups can be achieved). Furthermore, fig. 12(b) shows the corresponding number of augmenting paths per frame for the "SRI tree" image sequence (for both  $\alpha$ -expansion and Fast-PD). As can be seen from that figure, a substantial reduction in the number of augmenting paths is achieved by Fast-PD, which helps that algorithm to reduce its running time.

This same behavior has been observed in all other dynamic problems that Fast-PD has been tested on as well. Intuitively, what happens is illustrated in Fig. 13(a). Fast-PD has already managed to close the gap between the final primal-dual costs primal<sub>x</sub>, dual<sub>y</sub> of the previous MRF. However, due to the possibly different singleton (*i.e.*  $c_p(\cdot)$ ) or pairwise (*i.e.*  $w_{pq}d(\cdot, \cdot)$ ) potentials of the current MRF, these costs need to be perturbed to generate the new initial costs primal<sub>x</sub>, dual<sub>y</sub>. Nevertheless, as only slight perturbations take place, the new primal-dual gap (*i.e.* between primal<sub>x</sub>, dual<sub>y</sub>) will still be close to the previous gap (*i.e.* between primal<sub>x</sub>, dual<sub>y</sub>). As a result, the new gap will remain small. Few augmenting paths will therefore have to be found for the current MRF, and thus the algorithm's performance is boosted.

Put otherwise, for the case of dynamic MRFs, Fast-PD manages to boost performance, *i.e.* reduce number of augmenting paths, across two different "axes". The first axis lies along the different inner-iterations of the same MRF (*e.g.* see red arrows in Fig. 13(b)), whereas the second axis extends across time, *i.e.* across different MRFs (*e.g.* see blue arrow in Fig. 13(b), connecting the last iteration of MRF<sup>t-1</sup> to the first iteration of MRF<sup>t</sup>).



**Fig. 13:** (a) The final costs  $\operatorname{primal}_{\bar{\mathbf{x}}}$ ,  $\operatorname{dual}_{\bar{\mathbf{y}}}$  of the previous MRF are slightly perturbed to give the initial costs  $\operatorname{primal}_{\mathbf{x}}$ ,  $\operatorname{dual}_{\mathbf{y}}$  of the current MRF. Therefore, the initial primal-dual gap of the current MRF will be close to the final primal-dual gap of the previous MRF. Since the latter is small, so will be the former, and thus few augmenting paths will need to be computed for the current MRF. (b) Fast-PD reduces the number of augmenting paths in 2 ways: internally, *i.e.* across iterations of the same MRF (see red arrows), as well as externally, *i.e.* across different MRFs (see blue arrow).

#### **6.** Conclusions

In conclusion, a new graph-cut based method for MRF optimization has been proposed. It generalizes  $\alpha$ -expansion, while it also manages to be substantially faster than this state-of-the-art technique. Hence, regarding optimization of static MRFs, this method provides a significant speedup. In addition to that, however, it can also be used for boosting the performance of dynamic MRFs. In both cases, its efficiency comes from the fact that it exploits information not only from the "primal" problem (*i.e.* the MRF optimization problem), but also from a "dual" problem. Moreover, despite its speed, the proposed method can nevertheless guarantee almost optimal solutions for a very wide class of NP-hard MRFs. Due to all of the above, and given the ubiquity of MRFs, we strongly believe that Fast-PD can prove to be an extremely useful tool for many problems in computer vision in the years to come.

#### References

- N. Komodakis, G. Tziritas and N. Paragios. Fast Primal-Dual Strategies for MRF Optimization. Technical report, 2006. 5
- [2] Y. Boykov and V. Kolmogorov. An experimental comparison of min-cut/max-flow algorithms for energy minimization in vision. *PAMI*, 26(9), 2004. 6
- [3] Y. Boykov, O. Veksler, and R. Zabih. Fast approximate energy minimization via graph cuts. *PAMI*, 23(11), 2001. 1
- [4] O. Juan and Y. Boykov. Active graph cuts. In CVPR, 2006. 1
- [5] P. Kohli and P. H. Torr. Efficiently solving dynamic markov random fields using graph cuts. In *ICCV*, 2005. 1, 7
- [6] N. Komodakis and G. Tziritas. A new framework for approximate labeling via graph-cuts. In *ICCV*, 2005. 1, 2, 3
- [7] C. Rother, S. Kumar, V. Kolmogorov, and A. Blake. Digital tapestry. In CVPR, 2005. 7
- [8] R. Szeliski, *et al.* A comparative study of energy minimization methods for markov random fields. In *ECCV*, 2006. 1