

# Continuous 3D Label Stereo Matching using Local Expansion Moves

Tatsunori Taniai, Yasuyuki Matsushita, Yoichi Sato, and Takeshi Naemura,

**Abstract**—We present an accurate stereo matching method using *local expansion moves* based on graph cuts. This new move-making scheme is used to efficiently infer per-pixel 3D plane labels on a pairwise Markov random field (MRF) that effectively combines recently proposed slanted patch matching and curvature regularization terms. The local expansion moves are presented as many  $\alpha$ -expansions defined for small grid regions. The local expansion moves extend traditional expansion moves by two ways: localization and spatial propagation. By localization, we use different candidate  $\alpha$ -labels according to the locations of local  $\alpha$ -expansions. By spatial propagation, we design our local  $\alpha$ -expansions to propagate currently assigned labels for nearby regions. With this localization and spatial propagation, our method can efficiently infer MRF models with a continuous label space using randomized search. Our method has several advantages over previous approaches that are based on fusion moves or belief propagation; it produces *submodular moves* deriving a *subproblem optimality*; it helps find good, smooth, piecewise linear disparity maps; it is suitable for parallelization; it can use cost-volume filtering techniques for accelerating the matching cost computations. Even using a simple pairwise MRF, our method is shown to have best performance in the Middlebury stereo benchmark V2 and V3.

**Index Terms**—Stereo Vision, 3D Reconstruction, Graph Cuts, Markov Random Fields, Discrete-Continuous Optimization.



## 1 INTRODUCTION

STEREO vision often struggles with a bias toward reconstructing fronto-parallel surfaces, which can stem from matching cost, smoothness regularization, or even inference [5], [45].

Segment-based stereo [4] that represents disparity maps by disparity plane segments, could ease the bias issue but recovered surfaces are constrained to be piecewise planar. Recently, two breakthroughs have been independently made that overcome the fronto-parallel bias while surpassing the limitation of segment-based stereo.

One breakthrough is matching cost using slanted patch matching [5]. In this approach, the disparity  $d_p$  of each pixel  $p$  is over-parameterized by a local disparity plane

$$d_p = a_p u + b_p v + c_p \quad (1)$$

defined on the image domain  $(u, v)$ , and the triplet  $(a_p, b_p, c_p)$  is estimated for each pixel  $p$  instead of directly estimating  $d_p$ . The matching window is then slanted according to Eq. (1), which produces linearly-varying disparities within the window and thereby measures patch dissimilarity accurately even with large matching windows.

Another key invention is curvature regularization [27], [33] by tangent plane labels. Unlike second-order smoothness [45] that forms a higher-order term, this curvature regularization is nicely represented by pairwise terms like conventional (fronto-parallel) linear [19] and truncated linear models, and can handle smooth surfaces beyond planes.

Given slanted patch matching [5] and tangent-based curvature regularization [33], the use of 3D plane labels allows us to establish a new stereo model using a pairwise Markov random field (MRF) [12] that is free from the

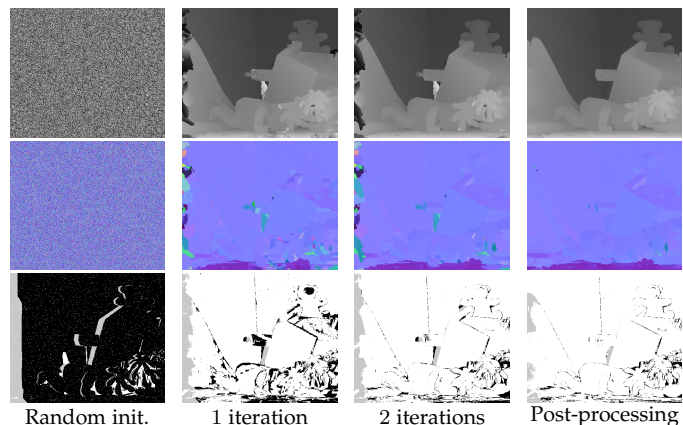


Fig. 1. Evolution of our stereo matching estimates. From top to bottom, we show disparity maps, normal maps of disparity planes, and error maps with 0.5 pixel threshold where ground truth is given. In our framework, we start with random disparities that are represented by per-pixel 3D planes, *i.e.*, disparities (top) and normals (middle). We then iteratively apply our local expansion moves using GC (middle) to update and propagate local disparity planes. Finally, the resulting disparity map is further refined at a post-processing stage using left-right consistency check and weighted median filtering (rightmost).

fronto-parallel bias. However, while stereo with standard 1D discrete disparity labels [8], [23], [42] can be directly solved by discrete optimizers such as graph cuts (GC) [7], [24] and belief propagation (BP) [11], [47], such approaches cannot be directly used for continuous 3D labels due to the huge or infinite label space  $(a, b, c) \in \mathbb{R}^3$ .

To efficiently infer 3D labels, recent successful methods [3], [5], [14], [32] use *PatchMatch* [1], [2], an inference algorithm using spatial label propagation. In *PatchMatch*, each pixel is updated in raster-scan order and its refined label is propagated to next pixels as their candidate labels.

- T. Taniai, T. Naemura, and Y. Sato are with the University of Tokyo, Japan.
- Y. Matsushita is with Osaka University, Japan.

Report submitted Mar. 28, 2016; revised Jul. 5 and Oct. 17, 2017.

Further in [3], this sequential algorithm is combined with BP yielding an efficient optimizer PMBP for pairwise MRFs. In terms of MRF optimization, however, BP is considered a *sequential optimizer*, which improves each variable individually keeping others conditioned at the current state. In contrast, GC improves all variables simultaneously by accounting for interactions across variables, and this global property helps optimization avoid bad local minimums [38], [45]. In order to take advantage of this and efficiently infer 3D planes by GC, it is important to use spatial propagation. Nevertheless, incorporating such spatial propagation into GC-based optimization is not straightforward, because inference using GC rather processes *all nodes simultaneously*, not *one-by-one sequentially* like PatchMatch and BP.

In this paper, we introduce a new move making scheme, *local expansion moves*, that enables spatial propagation in GC optimization. The local expansion moves are presented as many  $\alpha$ -expansions [8] defined for a small extent of regions at different locations. Each of this small or local  $\alpha$ -expansion tries to improve the current labels in its local region in an energy minimization manner using GC. Here, those current labels are allowed to move to a candidate label  $\alpha$ , which is given uniquely to each local  $\alpha$ -expansion in order to achieve efficient label searching. At the same time, this procedure is designed to propagate a current label in a local region for nearby pixels. For natural scenes that often exhibit locally planar structures, the joint use of local expansion moves and GC has a useful property. It allows multiple pixels in a local region to be assigned the same disparity plane *by a single min-cut* in order to find a smooth solution. Being able to simultaneously update multiple variables also helps to avoid being trapped at a bad local minimum.

For continuous MRF inference, fusion-based approaches using GC [25] are often employed with some heuristics to generate disparity map proposals [33], [34], [45]. Our work bridges between apparently different inference approaches of PatchMatch and GC, and can take advantage of those heuristics used in both PatchMatch and fusion based approaches leading to higher efficiency and accuracy.

The advantages of our method are as follows. 1) Our local expansion move method produces *submodular moves* that guarantee the optimal labeling at each min-cut (subproblem optimal), which in contrast is not guaranteed in general fusion moves [25]. 2) This optimality and spatial propagation allow randomized search, rather than employ external methods to generate plausible initial proposals as done in fusion approaches [25], [45], which may limit the possible solutions. 3) Our method achieves greater accuracy than BP [3] thanks to the good properties of GC and local expansion moves. 4) Unlike other PatchMatch based methods [3], [5], [14], our method can incorporate the fast cost filtering technique of [32]. In this manner, we can efficiently reduce the computation complexity of unary terms from  $O(|W|)$  to approximately  $O(1)$ , removing dependency from support window size  $|W|$ . 5) Unlike PMBP [3], our method is well suited for parallelization using both CPU and GPU.<sup>1</sup> With multiple CPU cores, each of our local  $\alpha$ -expansions

(*i.e.*, min-cut computations) can be individually performed in parallel. With a GPU implementation, computations of unary terms can be more efficiently performed in a parallel manner for further acceleration.

This paper is an extended version of our conference paper [39]. The extensions are summarized as follows. We add theoretical verifications on the preferability of our method for piecewise planar scenes in Sec. 3.1, and also on the subproblem optimality of our algorithm in Sec. 3.3. Furthermore, the efficiency of our algorithm is improved by two ways; In Sec. 3.3 we show the parallelizability of our local expansion move algorithm; In Sec. 3.5 we show that the fast cost filtering technique of [32] can be used in our method. The effectiveness of both extensions is thoroughly evaluated in the experiments, and we show that even a CPU implementation of the proposed method achieves about 2.1x faster running times than our previous GPU implementation [39], with comparable or even greater accuracy. Finally, we add experimental evaluations in Sec. 4 showing that our method outperforms a fusion-based approach [33] and all the state-of-the-art methods registered in the latest Middlebury benchmark V3. Our code is publicly available online (<https://github.com/t-taniai/LocalExpStereo>).

## 2 RELATED WORK

### 2.1 MRF stereo methods

MRF stereo methods can be categorized into three approaches: discrete stereo, segment-based stereo, and continuous stereo.

Discrete stereo [8], [23], [42] formulates stereo matching as a discrete multi-labeling problem, where each pixel is individually assigned one of pre-defined discrete disparity values. For this problem, many powerful discrete optimizers, such as BP [11], [47], TRW [21], and GC [7], [24], can be directly used. Successful results are shown using GC with expansion moves [8], [38]. In expansion moves, the multi-labeling problem is reduced to a sequence of binary-labeling problems, each of which can be exactly solved by GC, if only pairwise potentials  $\psi$  meet the following submodularity of expansion moves [8], [22]:

$$\psi(\alpha, \alpha) + \psi(\beta, \gamma) \leq \psi(\beta, \alpha) + \psi(\alpha, \gamma). \quad (2)$$

Segment-based stereo [4], [15], [41], [43] assigns a 3D disparity plane for each of over-segmented image regions. Although this approach yields continuous-valued disparities, it strictly limits the reconstruction to a piecewise planar representation and is subject to the quality of initial segmentation. More recent methods alleviate these limitations by using a complex layered MRF [6], multi-scale segmentation [26], or jointly estimating segmentation [46].

The last group, to which our method belongs, is continuous stereo [3], [5], [14], [32], [33], [45], where each pixel is assigned a distinct continuous disparity value. Some methods [33], [45] use fusion moves [25], an operation that combines two solutions to make better one (binary fusion) by solving a non-submodular binary-labeling problem using QPBO-GC [22]. In this approach, a number of continuous-valued disparity maps (or so-called *proposals* in the literature [25]) are first generated by other external methods (*e.g.*, segment-based stereo [45]), which are then combined as a

1. Although BP is usually parallelizable, PMBP differs from BP's standard settings in that it defines label space *dynamically* for each pixel and *propagates* it; both make parallelization indeed non-trivial.

sequence of binary fusions. Our method differs from this fusion-based approach in that we use spatial propagation and randomization search during inference, by which we only require a randomized initial solution instead of those generated by external methods. Also, binary energies produced in our method are always submodular allowing exact inference via GC (subproblem optimal). More importantly, we rather provide an efficient inference mechanism that is conceptually orthogonal to proposal generation. Thus, conventional proposal generation schemes used in fusion approaches [33], [34], [45] can be incorporated in our method for further improvements. A stereo method by Bleyer *et al.* [5] proposes accurate photo-consistency measures using 3D disparity planes that are inferred by PatchMatch [1], [2]. Heise *et al.* [14] incorporate Huber regularization into [5] using convex optimization. Besse *et al.* [3] point out a close relationship between PatchMatch and BP and present a unified method called PatchMatch BP (PMBP) for pairwise continuous MRFs. PMBP is probably the closest approach to ours in spirit, but we use GC instead of BP for the inference. Therefore, our method is able to take advantage of better convergence of GC [38] for achieving greater accuracy. In addition, our method allows efficient parallel computation of unary matching costs and even min-cut operations.

## 2.2 Cost-volume filtering

Patch-based stereo methods often use cost-volume filtering for fast implementations. Generally, computing a matching cost  $C$  for a patch requires  $O(|W|)$  of computation, where  $|W|$  is the size of the patch. However, given a cost-volume slice  $\rho_d(p)$  that represents pixelwise raw matching costs  $\|I(p) - I'(p - d)\|$  for a certain disparity label  $d$ , the patch-based matching costs can be efficiently computed by applying a filtering to the cost map as  $C_d(p) = \sum_q \omega_{pq} \rho_d(q)$ . Here, the filter kernel  $\omega_{pq}$  represents the matching window at  $p$ . If we use a constant-time filtering, each matching cost  $C_d(p)$  is efficiently computed in  $O(1)$ .

The box filtering can achieve  $O(1)$  by using integral image but such simple filtering flattens object boundaries. For this boundary issue, Yoon and Kweon [48] propose an adaptive support-window technique that uses the joint bilateral filtering [35] for cost filtering. Although this adaptive window technique successfully deals with the boundary issue [16], it involves  $O(|W|)$  of computation because of the complexity of the bilateral filtering. Recently, He *et al.* [13] propose a constant-time edge-aware filtering named the *guided image filtering*. This filtering is employed in a cost-volume filtering method of [17], [36], achieving both edge-awareness and  $O(1)$  of matching-cost computation.

In principle, stereo methods using PatchMatch inference [3], [5], [14] cannot take advantage of the cost filtering acceleration, since in those methods the candidate disparity labels are given dynamically to each pixel and we cannot make a consistent-label cost-volume slice  $\rho_d(p)$ . To this issue, Lu *et al.* [32] extend PatchMatch [5] to use super-pixels as a unit of cost calculations. In their method, called PatchMatch filter (PMF), fast cost-volume filtering of [17], [36] is applied in small subregions and that approximately achieves  $O(1)$  of complexity. We will show in Sec. 3.5 that their subregion filtering technique can be effectively incor-

porated into our method, and we achieve greater accuracy than their local stereo method [32].

## 3 PROPOSED METHOD

This section describes our proposed method. Given two input images  $I_L$  and  $I_R$ , our purpose is to estimate the disparities of both images.

In Sec. 3.1, we first discuss a geometric property of our model. We then define our energy function in Sec. 3.2, and describe the fundamental idea of our optimizing strategy and its properties in Sec. 3.3. The whole optimization procedure is presented in Sec. 3.4, and we further discuss a fast implementation in Sec. 3.5.

### 3.1 Geometric interpretation to slanted patch matching

Slanted patch matching [5] by Eq. (1) implicitly assumes that the true disparity maps are approximately piecewise linear. While it is not discussed in [5], linear disparity is exactly related to motion by *planar surfaces* [4], [33]. Formally, if there exists a plane in the 3D world coordinates  $(x, y, z) \in \mathbb{R}^3$

$$a'_p x + b'_p y + c'_p z = h'_p \quad (3)$$

parameterized by  $(a'_p, b'_p, c'_p, h'_p)$ , then motion due to this geometric plane is represented by a disparity plane

$$d(u, v) = \frac{B}{h'_p} (a'_p u + b'_p v + f c'_p), \quad (4)$$

where  $B$  and  $f$  are the baseline and focal length of the stereo cameras.<sup>2</sup> This piecewise planar assumption is reasonable as many natural scenes approximately exhibit locally planar surfaces. As we will discuss later, this fact also supports our method design as our model and inference both softly prefer piecewise linear disparity maps while having the ability to express arbitrary curved surfaces.

### 3.2 Formulation

We use the 3D plane label formulation to take advantage of the powerful slanted patch matching [5] and curvature regularization [33]. Here, each pixel  $p$ 's disparity  $d_p$  is over-parameterized by a 3D plane  $d_p = a_p u + b_p v + c_p$ . Therefore, we seek a mapping  $f_p = f(p) : \Omega \rightarrow \mathcal{L}$  that assigns a disparity plane  $f_p = (a_p, b_p, c_p) \in \mathcal{L}$  for every pixel  $p$  in the left and right images. Following conventional MRF stereo methods [8], [23], [33], [42], we estimate  $f$  by minimizing the following energy function based on a pairwise MRF.

$$E(f) = \sum_{p \in \Omega} \phi_p(f_p) + \lambda \sum_{(p, q) \in \mathcal{N}} \psi_{pq}(f_p, f_q). \quad (5)$$

The first term, called the *data term* or *unary term*, measures the photo-consistency between matching pixels. The disparity plane  $f_p$  defines a warp from a pixel  $p$  in one image to its correspondence in the other image. The second term is called the *smoothness term* or *pairwise term*, which penalizes discontinuity of disparities between neighboring pixel pairs  $(p, q) \in \mathcal{N}$ . We define these terms below.

2. Derived easily by using the perspective projection  $x = uz/f$  and  $y = vz/f$ , and the depth-disparity relation  $z = Bf/d$  in rectified stereo.

### Data term

To measure photo-consistencies, we use the slanted patch matching term that has been recently proposed by [5]. The data term of  $p$  in the left image is defined as

$$\phi_p(f_p) = \sum_{s \in W_p} \omega_{ps} \rho(s|f_p). \quad (6)$$

Here,  $W_p$  is a square window centered at  $p$ . The weight  $\omega_{ps}$  implements the adaptive support window proposed in [48]. For this we replace the bilateral filter weight used in [5], [48] with more advanced guided image filtering [13] as below.

$$\omega_{ps} = \frac{1}{|W'|^2} \sum_{k:(p,s) \in W'_k} \left(1 + (I_p - \mu_k)^T (\Sigma_k + e)^{-1} (I_s - \mu_k)\right) \quad (7)$$

Here,  $I_p = I_L(p)/255$  is a normalized color vector,  $\mu_k$  and  $\Sigma_k$  are the mean and co-variance matrix of  $I_p$  in a local regression window  $W'_k$ , and  $e$  is an identity matrix with a small positive coefficient for avoiding over-fitting. We will discuss how to efficiently compute this filtering in Sec. 3.5.

Given a disparity plane  $f_p = (a_p, b_p, c_p)$ , the function  $\rho(s|f_p)$  in Eq. (6) measures the pixel dissimilarity between a support pixel  $s = (s_u, s_v)$  in the window  $W_p$  and its matching point in the right image

$$s' = s - (a_p s_u + b_p s_v + c_p, 0) \quad (8)$$

as

$$\rho(s|f_p) = (1 - \alpha) \min(\|I_L(s) - I_R(s')\|_1, \tau_{\text{col}}) + \alpha \min(|\nabla_x I_L(s) - \nabla_x I_R(s')|, \tau_{\text{grad}}). \quad (9)$$

Here,  $\nabla_x I$  represents the  $x$ -component of the gray-value gradient of image  $I$ , and  $\alpha$  is a factor that balances the weights of color and gradient terms. The two terms are truncated by  $\tau_{\text{col}}$  and  $\tau_{\text{grad}}$  to increase the robustness for occluded regions. We use linear interpolation for  $I_R(s')$ , and a sobel filter kernel of  $[-0.5 \ 0 \ 0.5]$  for  $\nabla_x$ . When the data term is defined on the right image, we swap  $I_L$  and  $I_R$  in Eqs. (7) and (9), and add the disparity value in Eq. (8).

### Smoothness term

For the smoothness term, we use a curvature-based, second-order smoothness regularization term [33] defined as

$$\psi_{pq}(f_p, f_q) = \max(w_{pq}, \epsilon) \min(\bar{\psi}_{pq}(f_p, f_q), \tau_{\text{dis}}). \quad (10)$$

Here,  $w_{pq}$  is a contrast-sensitive weight defined as

$$w_{pq} = e^{-\|I_L(p) - I_L(q)\|_1 / \gamma}, \quad (11)$$

where  $\gamma$  is a user-defined parameter. The  $\epsilon$  is a small constant value that gives a lower bound to the weight  $w_{pq}$  to increase the robustness to image noise. The function  $\bar{\psi}_{pq}(f_p, f_q)$  penalizes the discontinuity between  $f_p$  and  $f_q$  in terms of disparity as

$$\bar{\psi}_{pq}(f_p, f_q) = |d_p(f_p) - d_p(f_q)| + |d_q(f_q) - d_q(f_p)|, \quad (12)$$

where  $d_p(f_q) = a_q p_u + b_q p_v + c_q$ . The first term in Eq. (12) measures the difference between  $f_p$  and  $f_q$  by their disparity values at  $p$ , and the second term is defined similarly at  $q$ . We visualize  $\bar{\psi}_{pq}(f_p, f_q)$  as red arrows in Fig. 2 (a). The  $\bar{\psi}_{pq}(f_p, f_q)$  is truncated at  $\tau_{\text{dis}}$  to allow sharp jumps in disparity at depth edges.

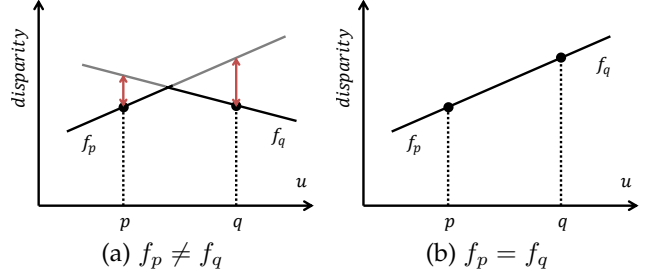


Fig. 2. Illustration of the smoothness term proposed in [33]. (a) The smoothness term penalizes the deviations of neighboring disparity planes shown as red arrows. (b) When neighboring pixels are assigned the same disparity plane, it gives no penalty; thus, it enforces second order smoothness for the disparity maps.

Notice that  $\bar{\psi}_{pq}(f_p, f_q) = 2|c_p - c_q|$  when  $a = b = 0$  is forced; therefore, the smoothness term  $\psi_{pq}(f_p, f_q)$  naturally extends the traditional truncated linear model [8], although the latter has a fronto-parallel bias and should be avoided [33], [45]. Also, as shown in Fig. 2 (b), this term becomes zero when  $f_p = f_q$ . This enforces piecewise linear disparity maps and so piecewise planar object surfaces. Furthermore, this term satisfies the following property for taking advantage of GC.

**Lemma 1.** *The term  $\psi_{pq}(f_p, f_q)$  in Eq. (10) satisfies the submodularity of expansion moves in Eq. (2).*

*Proof.* See [33] and also Appendix A.  $\square$

### 3.3 Local expansion moves

In this section, we describe the fundamental idea of our method, local expansion moves, as the main contribution of this paper. We first briefly review the original expansion move algorithm [8], and then describe how we extend it for efficiently optimizing continuous MRFs.

The expansion move algorithm [8] is a discrete optimization method for pairwise MRFs of Eq. (5), which iteratively solves a sequence of the following binary labeling problems

$$f^{(t+1)} = \underset{f'}{\operatorname{argmin}} E(f' | f_p' \in \{f_p^{(t)}, \alpha\}) \quad (13)$$

for all possible candidate labels  $\forall \alpha \in \mathcal{L}$ . Here, the binary variable  $f_p'$  for each pixel  $p$  is assigned either its current label  $f_p^{(t)}$  or a candidate label  $\alpha$ . If all the pairwise terms in  $E(f)$  meet the condition of Eq. (2), then the binary energies  $E(f')$  in Eq. (13) are submodular and this minimization can thus be exactly solved via GC [8] (subproblem optimal). Here, it is guaranteed that the energy does not increase:  $E(f^{(t+1)}) \leq E(f^{(t)})$ . However, the label space  $\mathcal{L}$  in our setting is a three dimensional continuous space  $(a, b, c)$ ; therefore, such an exhaustive approach cannot be employed.

Our local expansion moves extend traditional expansion moves by two ways; *localization* and *spatial propagation*. By localization, we use different candidate labels  $\alpha$  depending on the locations of pixels  $p$  in Eq. (13), rather than using the same  $\alpha$  label for all pixels. This is reasonable because the distributions of disparities should be different from location to location, therefore the selection of candidate labels  $\alpha$  should be accordingly different. By spatial propagation, we incorporate label propagation similar to the PatchMatch



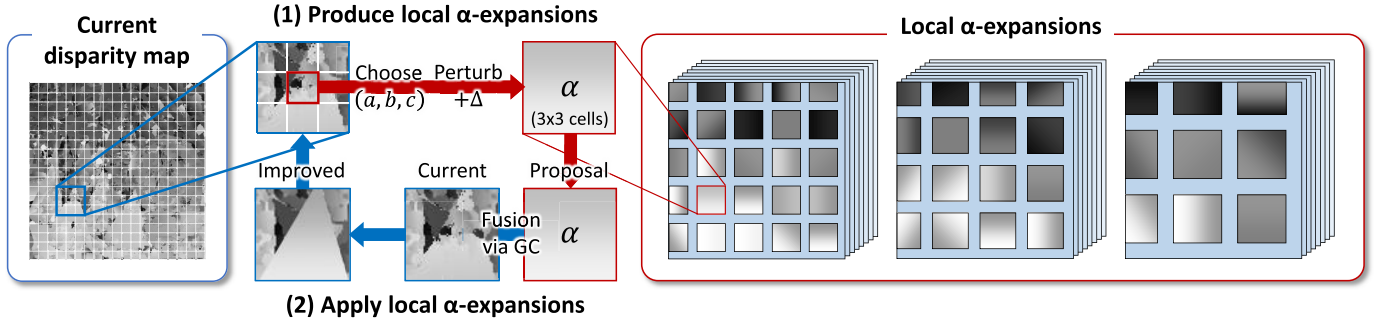


Fig. 3. Illustration of the proposed local expansion moves. The local expansion moves consist of many small  $\alpha$ -expansions (or *local  $\alpha$ -expansions*), which are defined using grid structures such shown in the left figure. These local  $\alpha$ -expansions are defined at each grid-cell and applied for  $3 \times 3$  neighborhood cells (or *expansion regions*). In the middle part, we illustrate how each of local  $\alpha$ -expansions works. (1) The candidate label  $\alpha$  (i.e.,  $\alpha = (a, b, c)$  representing a disparity plane  $d = au + bw + c$ ) is produced by randomly choosing and perturbing one of the currently assigned labels in its center cell. (2) The current labels in the expansion region are updated by  $\alpha$  in an energy minimization manner using GC. Consequently, a current label in the center cell can be propagated for its surrounding cells. In the right part, local  $\alpha$ -expansions are visualized as small patches on stacked layers with three different sizes of grid structures. As shown here, using local  $\alpha$ -expansions we can localize the scopes of label searching by their locations. Each layer represents a group of mutually-disjoint local  $\alpha$ -expansions, which are performed individually in a parallel manner.

inference [1], [2], [5] into GC optimization, and propagate currently assigned labels to nearby pixels via GC. The assumption behind this propagation is that, if a good label is assigned to a pixel, this label is likely a good estimate for nearby pixels as well. The localization and spatial propagation together make it possible for us to use a powerful randomized search scheme, where we no longer need to produce initial solution proposals as usually done in the fusion based approach [25]. Below we provide the detailed descriptions of our algorithm.

#### Local $\alpha$ -expansions for spatial propagation

We first define a grid structure that divides the image domain  $\Omega$  into grid regions  $C_{ij} \subset \Omega$ , which are indexed by 2D integer coordinates  $(i, j) \in \mathbb{Z}^2$ . We refer to each of these grid regions as a *cell*. We in this paper simply assume regular square cells, but this could be extended to use superpixels [18], [29], [40]. At a high level, the size of cells balances between the level of localization and the range of spatial propagation. Smaller sizes of cells can achieve finer localization but result in shorter ranges of spatial propagation. Later in Sec 3.4 we introduce different sizes of multiple grid structures for balancing these two factors well, but for now let us focus on using one grid structure.

Given a grid structure, we define a *local  $\alpha$ -expansion* at each cell  $(i, j)$ , which we specifically denote as  $\alpha_{ij}$ -*expansion*. We further define two types of regions for each  $\alpha_{ij}$ -expansion: its *center region*  $C_{ij}$  and *expansion region*

$$R_{ij} = C_{ij} \cup \left\{ \bigcup_{(m,n) \in \mathcal{N}(i,j)} C_{mn} \right\}, \quad (14)$$

i.e.,  $3 \times 3$  cells consisting of the center region  $C_{ij}$  and its eight neighbor cells.

In the middle part of Fig. 3, we focus on an expansion region and illustrate how an  $\alpha_{ij}$ -expansion works. We first randomly select a pixel  $r$  from the center region  $C_{ij}$ , and take its currently assigned label as  $(a, b, c) = f_r$ . We then make a candidate label  $\alpha_{ij}$  by perturbing this current label as  $\alpha_{ij} = (a, b, c) + \Delta$ . Finally, we update the current labels of pixels  $p$  in the expansion region  $R_{ij}$ , by choosing either their current labels  $f_p$  or the candidate label  $\alpha_{ij}$ . Here, similarly

---

#### Algorithm 1: ITERATIVE $\alpha_{ij}$ -EXPANSION

---

**input** : current  $f$ , target cell  $(i, j)$ , perturbation size  $|\Delta'|$   
**output**: updated  $f$  (only labels  $f_p$  at  $p \in R_{ij}$  are updated)

- 1 **repeat expansion proposer (propagation)**:
- 2    $\alpha_{ij} \leftarrow f_r$  with randomly chosen  $r \in C_{ij}$ ;
- 3    $f \leftarrow \operatorname{argmin} E(f' | f'_p \in \{f_p, \alpha_{ij}\}, p \in R_{ij})$ ;
- 4 **until**  $K_{prop}$  times;
- 5 **repeat RANSAC proposer (optional)**:
- 6    $\alpha_{ij} \leftarrow \operatorname{RANSAC}(f, C_{ij})$ ;
- 7    $f \leftarrow \operatorname{argmin} E(f' | f'_p \in \{f_p, \alpha_{ij}\}, p \in R_{ij})$ ;
- 8 **until**  $K_{RANS}$  times;
- 9 **repeat randomization proposer (refinement)**:
- 10    $\alpha_{ij} \leftarrow f_r$  with randomly chosen  $r \in C_{ij}$ ;
- 11    $\alpha_{ij} \leftarrow \alpha_{ij} + \Delta'$ ;
- 12    $f \leftarrow \operatorname{argmin} E(f' | f'_p \in \{f_p, \alpha_{ij}\}, p \in R_{ij})$ ;
- 13    $|\Delta'| \leftarrow |\Delta'|/2$ ;
- 14 **until**  $K_{rand}$  times (or  $|\Delta'|$  is sufficiently small);

---

to Eq. (13), we update the partial labeling by minimizing  $E(f')$  with binary variables:  $f'_p \in \{f_p, \alpha_{ij}\}$  for  $p \in R_{ij}$ , and  $f'_p = f_p$  for  $p \notin R_{ij}$ . Consequently, we obtain an improved solution as its minimizer with a lower or equal energy.

Notice that making the expansion region  $R_{ij}$  larger than the label selection region  $C_{ij}$  is the key idea for achieving spatial propagation. We can see this since, in an  $\alpha_{ij}$ -expansion without perturbation ( $\Delta = 0$ ), a current label  $f_p$  in the center region  $C_{ij}$  can be propagated for its nearby pixels in  $R_{ij}$  as the candidate label  $\alpha_{ij}$ .

We use this  $\alpha_{ij}$ -expansion iteratively as shown in Algorithm 1. Similarly to the PatchMatch algorithm [1], this iterative algorithm has propagation (lines 1–4) and refinement (lines 9–14) steps. In the propagation step, we apply  $\alpha_{ij}$ -expansions without perturbation to propagate labels from  $C_{ij}$  for  $R_{ij}$ . In the refinement step, we apply  $\alpha_{ij}$ -expansions with an exponentially decreasing perturbation-size to refine the labels. Here, propagation and refinement can be seen as heuristic proposers of a candidate label  $\alpha_{ij}$  that is expectedly a good estimate for the region  $C_{ij}$  or even its neighborhood  $R_{ij}$ . Our local expansion move method is conceptually orthogonal to any such heuristics for generating  $\alpha_{ij}$ . As such an example, we in lines 5–8 use a RANSAC proposer by following conventional proposal

generation schemes used in fusion-based methods [33], [45]. Here, we generate  $\alpha_{ij}$  by fitting a plane to the current disparity map  $f$  in the region  $C_{ij}$  using LO-RANSAC [9]. We perform this iterative  $\alpha_{ij}$ -expansion at every cell  $(i, j)$ .

This local  $\alpha$ -expansion method has the following useful properties. **Piecewise linearity:** It helps to obtain smooth solutions. In each  $\alpha_{ij}$ -expansion, multiple pixels in the expansion region  $R_{ij}$  are allowed to move-at-once to the same candidate label  $\alpha_{ij}$  at one binary-energy minimization, which contrasts to BP that updates only one pixel at once. Since a label represents a disparity plane here, our method helps to obtain piecewise linear disparity maps and thus piecewise planar surfaces as we have discussed in Sec 3.1. **Cost filtering acceleration:** We can accelerate the computation of matching costs  $\phi_p(f_p)$  in Eq. (6) by using cost-volume filtering techniques. We discuss this more in Sec 3.5. **Optimality and parallelizability:** With our energy formulation, it is guaranteed that each binary-energy minimization in Algorithm 1 can be optimally solved via GC. In addition, we can efficiently perform many  $\alpha_{ij}$ -expansions in a parallel manner. We discuss these matters in the following sections.

#### Mutually-disjoint local $\alpha$ -expansions

While the previous section shows how each local  $\alpha$ -expansion behaves, here we discuss the scheduling of local  $\alpha$ -expansions. We need a proper scheduling, because local  $\alpha$ -expansions cannot be simultaneously performed due to the overlapping expansion regions.

To efficiently perform local  $\alpha$ -expansions, we divide them into groups such that the local  $\alpha$ -expansions in each group are mutually disjoint. Specifically, we assign each  $\alpha_{ij}$ -expansion a group index  $k$  given by

$$k = 4(j \bmod 4) + (i \bmod 4), \quad (15)$$

and perform the iterative  $\alpha_{ij}$ -expansions in one group and another. As illustrated in Fig. 4, this grouping rule picks  $\alpha_{ij}$ -expansions at every four vertical and horizontal cells into the same group, and it amounts to 16 groups of mutually-disjoint local  $\alpha$ -expansions. We visualize a group of disjoint local  $\alpha$ -expansions as orange regions in Fig. 5, and also as a single layer of stacks in the right part of Fig. 3 with three different grid structures.

Notice that in each group, we leave *gaps* of one cell width between neighboring local  $\alpha$ -expansions. These gaps are to guarantee submodularity and independence for local  $\alpha$ -expansions. By the submodularity, we can show that our local  $\alpha$ -expansions always produce submodular energies and can thus be optimally solved via GC. Because of this submodularity, we can use a standard GC algorithm [7] instead of an expensive QPBO-GC algorithm [22], which is usually required in the fusion based approach. By the independence, we can show that the local  $\alpha$ -expansions in the same group do not interfere with each other. Hence, we can perform them simultaneously in a parallel manner. The parallelization of GC algorithms are of interest in computer vision [30], [37]. Our scheme is simple and can directly use existing GC implementations. The proof of this submodularity and independence is presented in the next section.

---

#### Algorithm 2: OPTIMIZATION PROCEDURE

---

```

1 Define three levels of grid structures:  $H = \{h_1, h_2, h_3\}$ ;
2 Initialize the current solution  $f$  randomly;
3 Initialize the perturbation size  $|\Delta|$ ;
4 repeat
5   foreach grid structure  $h \in H$  do
6     foreach disjoint group  $k = 0, 1, \dots, 15$  do
7       foreach cell  $(i, j)$  in the group  $k$  do [in parallel]
8         Do iterative  $\alpha_{ij}$  expansion  $(f, (i, j), |\Delta|)$ ;
9       end
10    end
11  end
12   $|\Delta| \leftarrow |\Delta|/2$ ;
13 until convergence;
14 Do post processing;
```

---

#### Submodularity and independence

To formally address the submodularity and independence of local  $\alpha$ -expansions, we discuss it using the form of fusion moves [25]. Let us assume a current solution  $f$  and a group of mutually-disjoint  $\alpha_{ij}$ -expansions to be applied. Simultaneously applying these  $\alpha_{ij}$ -expansions is equivalent to the following fusion-energy minimization:

$$f^* = \operatorname{argmin} E(f' | f'_p \in \{f_p, g_p\}), \quad (16)$$

where the proposal solution  $g$  is set to  $g_p = \alpha_{ij}$  if  $p$  belongs to any of expansion regions  $R_{ij}$  in this group; otherwise,  $p$  is in gaps so we assign  $\phi_p(g_p)$  an infinite unary cost for forcing  $f'_p = f_p$ . This proposal disparity map  $g$  is visualized as a single layer of stacks in the right part of Fig. 3. We prove the following lemmas:

**Lemma 2.** *Submodularity: the binary energies in Eq. (16) are submodular, i.e., all the pairwise interactions in Eq. (16) meet the following submodularity of fusion moves [22], [25]:*

$$\psi_{pq}(g_p, g_q) + \psi_{pq}(f_p, f_q) \leq \psi_{pq}(f_p, g_q) + \psi_{pq}(g_p, f_q). \quad (17)$$

*Proof.* Omitted.<sup>3</sup>  $\square$

**Lemma 3.** *Independence: the assignments to  $f'_p$  and  $f'_q$  do not influence each other, if  $p$  and  $q$  are in different expansion regions.*

*Proof.* The  $f'_p$  and  $f'_q$  have interactions if and only if there exists a chain of valid pairwise interactions  $C = \{\psi(f'_{s_0}, f'_{s_1}), \dots, \psi(f'_{s_{n-1}}, f'_{s_n})\}$  connecting  $p = s_0$  and  $q = s_n$ . But there is no such chain because  $C$  inevitably contains constant or unary terms at a gap ( $\psi_{cd}$  and  $\psi_{ac}$  in Fig. 5).  $\square$

### 3.4 Optimization procedure

Using the local expansion moves shown in the previous section, we present the overall optimization procedure in this section and summarize it in Algorithm 2.

This algorithm begins with defining grid structures. Here, we use three different sizes of grid structures for better balancing localization and spatial propagation.

<sup>3</sup> We outline proof. For  $(p, q)$  inside an expansion region such as  $(a, b)$  in Fig. 5, Eq. (17) is relaxed to Eq. (2). For the other cases, e.g.,  $(a, c)$  or  $(c, d)$ , pairwise terms  $\psi_{ac}(f'_a, f'_c)$  and  $\psi_{cd}(f'_c, f'_d)$  become unary or constant terms. Therefore, when implementing an  $\alpha_{ij}$ -expansion using min-cut in  $R_{ij}$ , we create nodes for  $\forall p \in R_{ij}$  and add  $\psi_{ac}(f'_a, f'_c)$  as unary potentials of nodes  $a$  at the inner edge of  $R_{ij}$ . See also [8] for the conversion of  $\psi_{ab}(f'_a, f'_b)$  into edge capacities under expansion moves.

		Cell index $i$							
		0	1	2	3	4	5	6	7
Cell index $j$	0	0	1	2	3	0	1	2	3
	1	4	5	6	7	4	5	6	7
	2	8	9	10	11	8	9	10	11
	3	12	13	14	15	12	13	14	15
	4	0	1	2	3	0	1	2	3
	5	4	5	6	7	4	5	6	7
	6	8	9	10	11	8	9	10	11
	7	12	13	14	15	12	13	14	15

Fig. 4. Group index for  $\alpha_{ij}$ -expansions. We perform  $\alpha_{ij}$ -expansions in the same group in parallel.

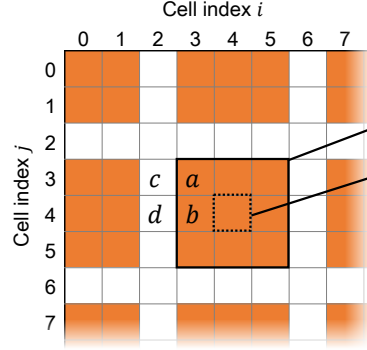


Fig. 5. Expansion regions of mutually disjoint  $\alpha_{ij}$ -expansions (group index  $k = 0$ ). We leave white gaps between neighbors.

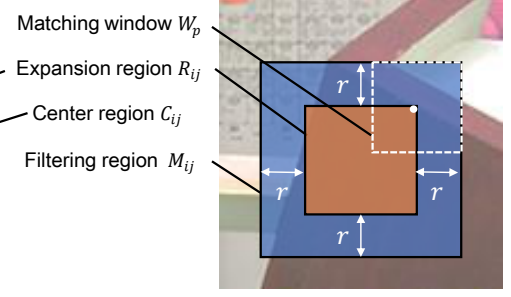


Fig. 6. Filtering region  $M_{ij}$ . The margin width  $r$  corresponds with the radius of the matching window  $W_p$ .

At line 2 of Algorithm 2, the solution  $f$  is randomly initialized. To evenly sample the allowed solution space, we take the initialization strategy described in [5]. Specifically, for each  $f_p = (a_p, b_p, c_p)$  we select a random disparity  $z_0$  in the allowed disparity range  $[0, \text{dispmax}]$ . Then, a random unit vector  $n = (n_x, n_y, n_z)$  and  $z_0$  are converted to the plane representation by  $a_p = -n_x/n_z$ ,  $b_p = -n_y/n_z$ , and  $c_p = -(n_x p_u + n_y p_v + n_z z_0)/n_z$ .

In the main loop through lines 4–13, we select one grid level  $h$  from the pre-defined grid structures (line 5), and apply the iterative  $\alpha_{ij}$  expansions of Algorithm 1 for each cell of the selected structure  $h$  (lines 6–10). As discussed in Sec. 3.3, we perform the iterative  $\alpha_{ij}$  expansions by dividing into disjoint groups defined by Eq. (15). Because of this grouping, the  $\alpha_{ij}$  expansions in the loop at lines 7–9 are mutually independent and can be performed in parallel.

The perturbation at line 9 of Algorithm 1 is implemented as described in [5]. Namely, each candidate label  $\alpha_{ij} = (a, b, c)$  is converted to the form of a disparity  $d$  and normal vector  $n$ . We then add a random disparity  $\Delta'_d \in [-r_d, r_d]$  and a random vector  $\Delta'_n$  of size  $\|\Delta'_n\|_2 = r_n$  to them, respectively, to obtain  $d'$  and  $n'$ . Finally,  $d'$  and  $n'/|n'|$  are converted to the plane representation  $\alpha_{ij} \leftarrow (a', b', c')$  to obtain a perturbed candidate label. The values  $r_d$  and  $r_n$  define an allowed change of disparity planes. We initialize them by setting  $r_d \leftarrow \text{dispmax}/2$  and  $r_n \leftarrow 1$  at line 3 of Algorithm 2, and update them by  $r_d \leftarrow r_d/2$  and  $r_n \leftarrow r_n/2$  at line 12 of Algorithm 2 and line 13 of Algorithm 1.

Finally, we perform post-processing using left-right consistency check and weighted median filtering as described in [5] for further improving the results. This step is widely employed in recent methods [3], [5], [14], [32].

Note that there are mainly two differences between this algorithm and our previous version [39]. For one thing, we have removed a per-pixel label refinement step of [39]. This step is required for updating a special label space structure named *locally shared labels* (LSL) used in [39], but can be removed in our new algorithm by using local  $\alpha$ -expansions instead. For the other, the previous algorithm proceeds rather in a batch cycle; *i.e.*, it produces all local  $\alpha$ -expansions with all grid structures at once and then applies them to the current solution  $f$  in one iteration. This leads to slower convergence than our new algorithm that produces each local  $\alpha$ -expansion always from the latest current solution  $f$ .

### 3.5 Fast implementation

Our method has two major computation parts: the calculations of matching costs  $\phi_p(f_p)$  of Eq. (6), and application of GC in Algorithm 1. In Sec. 3.3 and Sec. 3.4, we have shown that the latter part can be accelerated by performing disjoint local  $\alpha$  expansions in parallel. In this section, we discuss the former part.

The calculations of the matching costs  $\phi_p(f_p)$  are very expensive, since they require  $O(|W|)$  of computation for each term, where  $|W|$  is the size of the matching window. In our previous algorithm [39], we accelerate this part by using GPU. The use of GPU is reasonable for our method because  $\phi_p(f_p)$  of all pixels can be individually computed in parallel during the inference, which contrasts to PMBP [3] that can only sequentially process each pixel. Still, the computation complexity is  $O(|W|)$  and it becomes very inefficient if we only use CPU [39].

In the following part, we show that we can approximately achieve  $O(1)$  of complexity for computing each  $\phi_p(f_p)$ . The key observation here is that, we only need to compute this matching cost  $\phi_p(f_p)$  during the  $\alpha_{ij}$ -expansions in Algorithm 1, where  $\phi_p(f_p)$  is computed as  $\phi_p(\alpha_{ij})$  for all  $p \in R_{ij}$ . With this consistent-label property, we can effectively incorporate the fast subregion cost-filtering technique used in [32], by which  $\phi_p(\alpha_{ij})$  for all  $p \in R_{ij}$  are efficiently computed at once.

To more specifically discuss it, we first separate the computation of  $\phi_p(f_p)$  into two steps: the calculations of raw matching costs

$$\rho_{f_p}(s) = \rho(s|f_p) \quad (18)$$

for the support pixels  $s \in W_p$ , and the aggregation of the raw matching costs using an edge-aware filter kernel  $\omega_{ps}$

$$\phi_{f_p}(p) = \sum_{s \in W_p} \omega_{ps} \rho_{f_p}(s). \quad (19)$$

We also define a filtering region  $M_{ij}$  as the joint matching windows in  $R_{ij}$  as

$$M_{ij} = \bigcup_{p \in R_{ij}} W_p. \quad (20)$$

As shown in Fig. 6, this  $M_{ij}$  is typically a square region margining  $R_{ij}$  with  $r$  pixels of width around  $R_{ij}$ , where  $r$  is the radius of the matching windows.

In the raw matching part of Eq. (18), the calculations of the raw costs  $\rho_{f_p}(s)$ ,  $\forall s \in W_p$  for all  $p \in R_{ij}$  generally require  $O(|W||R_{ij}|)$  of total computation. However, with the consistent-label property (*i.e.*,  $f_p = \alpha_{ij}$  for all  $p \in R_{ij}$ ), they can be computed at once in  $O(|M_{ij}|)$  by computing  $\rho_{\alpha_{ij}}(s)$  for  $\forall s \in M_{ij}$ . Here, the computation complexity for each unary term is  $O(|M_{ij}|/|R_{ij}|)$ . Therefore, if  $|M_{ij}| \simeq |R_{ij}|$ , we can approximately achieve  $O(|M_{ij}|/|R_{ij}|) \simeq O(1)$  [32].

Similarly, the cost aggregation part of Eq. (19) can be done in approximately  $O(1)$ , if we apply a constant-time edge-aware filtering  $\omega_{ps}$  to  $\rho_{\alpha_{ij}}(s)$ ,  $\forall s \in M_{ij}$  [32]. As such an example we have chosen guided image filtering [13] for Eq. (7) but could use more sophisticated filtering [31].

Note that as the consistent-label property is the key to being able to use the filtering techniques, this scheme cannot be used in other PatchMatch based methods [3], [5], [14] except for [32] that uses superpixels as computation units.

## 4 EXPERIMENTS

In the experiments, we first evaluate our method on the Middlebury benchmark V2 and V3. While the latest V3 benchmark allows us to compare with state-of-the-art methods, we also use the older V2 as there are more related methods [3], [5], [14], [32], [39] registered only in this version.

We further assess the effect of sizes of grid-cells and also the effectiveness of the proposed acceleration schemes in comparison to our previous method [39]. Our method is further compared with the PMBP [3], PMF [32] and Olsson *et al.* [33] methods that are closely related to our approach. Additional results for these analyses and comparisons on more image pairs are provided in the supplementary material.

### Settings

We use the following settings throughout the experiments. We use a desktop computer with a Core i7 CPU (3.5 GHz  $\times$  4 physical cores) and NVIDIA GeForce GTX Titan Black GPU. All methods are implemented using C++ and OpenCV.

The parameters of our data term are set as  $\{e, \tau_{\text{col}}, \tau_{\text{grad}}, \alpha\} = \{0.01^2, 10, 2, 0.9\}$  following [32], [36] and [5]. The size of matching windows  $W_p$  is set to  $41 \times 41$  (*i.e.*, the size of local regression windows  $W'_k$  in Eq. (7) is set to  $21 \times 21$ ), which is the same setting with PMBP [3] and [5], [14]. For the smoothness term, we use  $\{\lambda, \tau_{\text{dis}}, \epsilon, \gamma\} = \{1, 1, 0.01, 10\}$  and eight neighbors for  $\mathcal{N}$ .

For our algorithm, we use three grid structures with cell sizes of  $5 \times 5$ ,  $15 \times 15$ , and  $25 \times 25$  pixels. The iteration numbers  $\{K_{\text{prop}}, K_{\text{rand}}\}$  in Algorithm 1 are set to  $\{1, 7\}$  for the first grid structure, and  $\{2, 0\}$  (only propagation step) for the other two. Except in Sec. 4.2 and 4.7 we disable the RANSAC proposer for fair comparisons with related methods or for avoiding cluttered analysis. When enabled, we set  $K_{\text{RANS}} = 1$  for all the grid structures. We iterate the main loop ten times. We use a GC implementation of [7].

We use two variants of our method. **LE-GF** uses guided image filtering [13] for  $w_{ps}$  as proposed in Sec. 3.2. We only use a CPU implementation for this method. **LE-BF** uses bilateral filtering for  $\omega_{ps}$  defined similarly to Eq. (11) and  $\lambda = 20$  so that the energy function corresponds to one we used in [39]. The computation of matching costs is thus as slow as  $O(|W|)$ . In a GPU implementation, we accelerate

this calculation by computing each unary term individually in parallel on GPU. For both LE-BF and LE-GF, we perform disjoint local  $\alpha$ -expansions in parallel using four CPU cores.

### 4.1 Evaluation on the Middlebury benchmark V2

We show in Table 1 selected rankings on the Middlebury stereo benchmark V2 for 0.5-pixel accuracy, comparing with other PatchMatch-based methods [3], [5], [14], [32], [39]. The proposed LE-GF method achieves the current best average rank (3.9) and bad-pixel-rate (5.97%) amongst more than 150 stereo methods including our previous method (GC+LSL) [39]. Even without post-processing, our LE-GF method still outperforms the other methods in average rank, despite that all the methods in Table 1 use the post-processing. On the other hand, the proposed LE-BF method achieves comparable accuracy with our previous algorithm [39], since both methods optimize the same energy function in very similar ways.

### 4.2 Evaluation on the Middlebury benchmark V3

As major differences from the older version, the latest Middlebury benchmark introduces more challenging difficulties such as high-resolution images, different exposure and lighting between an image pair, imperfect rectification, etc. Since our model, especially our data term borrowed from [5], does not consider such new difficulties, we here slightly modify our model for adapting it to the latest benchmark. To this end, we incorporate the state-of-the-art CNN-based matching cost function by Zbontar and LeCun [49], by following manners of current top methods [10], [20], [28], [29] on the latest benchmark. Here, we only replace our pixelwise raw matching cost function  $\rho(\cdot)$  of Eq. (9) with their matching cost function as follows.

$$\rho'(s|f_p) = \min(C_{\text{CNN}}(s, s'), \tau_{\text{CNN}}). \quad (21)$$

The function  $C_{\text{CNN}}(s, s')$  computes a matching cost between left and right image patches of  $11 \times 11$ -size centered at  $s$  and  $s'$ , respectively, using a neural network (called MC-CNN-acrt in [49]).<sup>4</sup> In our method we truncate the matching cost values at  $\tau_{\text{CNN}} = 0.5$  and aggregate them for the support window pixels  $s \in W_p$  using slanted patch matching of Eq. (6). In addition to the change in Eq. (9), the widths of square cells of three grid structures are changed proportionally to 1%, 3%, and 9% of the image width. Since only one submission is allowed for the test data by the new benchmark regulation, we here only use our LE-GF method with  $\lambda = 0.5$  and  $K_{\text{RANS}} = 1$ . We keep all the other parameters (*i.e.*,  $\{e, W_p, \tau_{\text{dis}}, \epsilon, \gamma, K_{\text{prop}}, K_{\text{rand}}\}$ ) as default.

Table 2 shows selected rankings for the *bad 2.0* metric (*i.e.*, percentage of bad pixels by the error threshold of 2.0 pixels at full resolution) for non-occluded regions as the default metric. Our method outperforms all existing 63 methods on the table not only for the default metric but for all combinations of  $\{\textit{bad } 0.5, 1.0, 2.0, 4.0\} \times \{\textit{nonocc}, \textit{all}\}$

4. Using the authors' code and pre-trained model, we pre-compute the matching costs  $C_{\text{CNN}}$  at all allowed integer disparities using fronto-parallel patches. During the inference, we linearly interpolate the function  $C_{\text{CNN}}(s, s')$  at non-integer coordinates  $s'$ . MC-CNN learns a patch similarity metric invariantly to the left-right patch distortion. Thus, the bias due to using fronto-parallel patches in  $C_{\text{CNN}}$  is small.



TABLE 1

Middlebury benchmark V2 for 0.5-pixel accuracy. Our method using the guided image filtering (GF) achieves the current best average rank 3.9. Error percents for *all* pixels, non-occluded pixels (*nonocc*), and pixels around depth discontinuities (*dis*) are shown. (Snapshot on April 22, 2015)

Algorithm	Avg. Rank	Tsukuba			Venus			Teddy			Cones			Average Percent Bad Pixels
		nonocc	all	disc	nonocc	all	disc	nonocc	all	disc	nonocc	all	disc	
<b>1. PROPOSED (GF)</b>	<b>3.9</b>	4.08 2	4.71 2	9.71 4	0.35 2	0.56 2	3.30 2	5.16 5	7.73 3	14.2 5	3.46 5	8.65 6	9.72 9	5.97
2. GC+LSL [39]	6.2	5.04 3	5.56 3	14.0 13	0.66 6	0.88 6	5.82 8	<b>4.20 1</b>	7.12 2	12.9 3	3.77 8	9.16 9	10.4 13	6.63
3. PM-Huber [14]	8.8	7.12 11	7.80 13	13.7 11	1.00 12	1.40 13	7.80 19	5.53 8	9.36 5	15.9 9	<b>2.70 1</b>	7.90 2	<b>7.77 1</b>	7.33
7. PMF [32]	12.5	11.0 39	11.4 36	16.0 32	0.72 8	0.92 7	5.27 7	4.45 3	9.44 7	13.7 4	2.89 2	8.31 3	8.22 2	7.69
11. PMBP [3]	19.7	11.9 52	12.3 48	17.8 60	0.85 10	1.10 8	6.45 11	5.60 9	12.0 12	15.5 6	3.48 6	8.88 8	9.41 6	8.77
14. PatchMatch [5]	28.2	15.0 74	15.4 73	20.3 89	1.00 13	1.34 12	7.75 17	5.66 10	11.8 10	16.5 10	3.80 9	10.2 11	10.2 11	9.91
Reference Evaluations														
<b>BF w/ post-proc</b>	<b>6.6</b>	5.48 3	6.07 3	14.5 15	0.82 9	1.08 7	6.32 10	<b>4.05 1</b>	7.24 3	12.4 3	3.69 7	9.12 8	9.96 10	6.73
<b>BF w/o post-proc</b>	<b>8.4</b>	5.36 3	5.99 3	14.1 14	0.82 9	1.11 9	6.33 10	4.59 4	10.8 8	14.0 5	3.90 9	10.3 13	10.6 14	7.32
<b>GF w/ post-proc</b>	<b>3.9</b>	4.08 2	4.71 2	9.71 4	0.35 2	0.56 2	3.30 2	5.16 5	7.73 3	14.2 5	3.46 5	8.65 6	9.72 9	5.97
<b>GF w/o post-proc</b>	<b>3.9</b>	4.07 2	4.79 2	9.78 4	0.41 2	0.78 2	4.01 2	4.26 2	9.19 4	13.6 4	3.37 5	9.63 9	9.59 9	6.13

TABLE 2

Middlebury benchmark V3 for the *bad 2.0 nonocc* metric. Our method achieves the current best average error rate among 64 existing algorithms. We here list current top methods [10], [20], [28], [29], [49] that use the matching costs of MC-CNN-acrt [49] like ours. (Snapshot on July 4, 2017)

bad 2.0 (%)		Weight	Austr	AustrP	Bicyc2	Class	ClassE	Compu	Crusa	CrusaP	Djemb	DjembL	Hoops	Livgrm	Nkuba	Plants	Stairs
<b>LocalExp (ours)</b>	H	<b>5.43 1</b>	3.65 2	2.87 3	<b>2.98 1</b>	<b>1.99 1</b>	<b>5.59 1</b>	<b>3.37 1</b>	3.48 2	<b>3.35 1</b>	<b>2.05 1</b>	10.3 2	9.75 2	8.57 4	14.4 8	5.40 3	9.55 5
3DMST [28]	H	5.92 2	3.71 3	2.78 2	4.75 2	2.72 4	7.36 4	4.28 2	<b>3.44 1</b>	3.76 2	2.35 2	12.6 5	11.5 4	8.56 3	14.0 7	5.35 2	8.87 4
MC-CNN+TDSR [10]	F	6.35 3	5.45 8	4.45 12	6.80 13	3.46 10	10.7 10	6.05 7	5.01 7	5.19 8	2.62 6	10.8 3	<b>9.62 1</b>	<b>6.59 1</b>	<b>11.4 1</b>	6.01 6	<b>7.04 1</b>
PMSC [29]	H	6.71 4	<b>3.46 1</b>	<b>2.68 1</b>	6.19 9	2.54 2	6.92 2	4.54 3	3.96 3	4.04 4	2.37 3	13.1 7	12.3 5	12.2 6	16.2 13	5.88 5	10.8 8
NTDE [20]	H	7.44 8	5.72 12	4.36 11	5.92 7	2.83 5	10.4 7	5.71 5	5.30 8	5.54 9	2.40 4	13.5 8	14.1 9	12.6 8	13.9 6	6.39 8	12.2 13
MC-CNN-acrt [49]	H	8.08 9	5.59 11	4.55 15	5.96 8	2.83 5	11.4 14	5.81 6	8.32 12	8.89 16	2.71 7	16.3 12	14.1 10	13.2 10	13.0 3	6.40 9	11.1 10

except for *bad 4.0 - all*. In Table 3 we analyze effects of our post-processing and RANSAC proposer using training data. Again, our method is ranked first even without post-processing for the default and all other *bad nonocc* metrics. Also, the RANSAC proposer reduces errors by one point (see Sec. 4.7 for more analysis). The avr and stdev show that our inference is stable by different random initializations.

### 4.3 Effect of grid-cell sizes

To observe the effect of grid-cell sizes, we use the three sizes of grid-cells,  $5 \times 5$  pixels (denoted as “S”mall),  $15 \times 15$  pixels (denoted as “M”edium),  $25 \times 25$  pixels (denoted as “L”arge), in different combinations and assess the performance using the following five different settings; (S, S, S): the small-size cells for all grid-structures; (M, M, M): the medium-size cells for all grid-structures; (L, L, L): the large-size cells for all grid-structures; (S, M, M): the small and medium-size cells for the first and the other two grid-structures, respectively; (S, M, L): the small, medium, and large-size cells for the first, second, and third grid-structures, respectively (the default setting for our method described above). Here, the iteration numbers for the first to third grid-structures are kept as default. We use the LE-GF method so as to access the effect on cost filtering acceleration as well. We use  $\lambda = 0.5$  but keep the other parameters as default. Using these settings, we observe the performance variations by estimating the disparities of only the left image of the *Reindeer* dataset without post-processing.

The plots in Fig. 7 (a) show the transitions of the energy function values over iterations. The plots in Fig. 7 (b)

show the temporal transitions of error rates with subpixel accuracy. As shown, the joint use of multiple cell sizes improves the performance in both energy reduction and accuracy. Although (S, M, M) and (S, M, L) show almost the same energy transitions, the proposed combination (S, M, L) shows faster and better convergence in accuracy.

Figure 8 compares the results of the five settings with error maps for (S, M, M) and (S, M, L). The use of larger grid-cells helps to obtain smoother disparities, and it is especially effective for occluded regions.

Comparing the running times in Fig. 7 (b), the use of the small-size cells is inefficient due to the increased overhead in cost filtering, whereas the use of the medium and large-size cells achieves almost the same efficiency.

### 4.4 Efficiency evaluation in comparison to our previous algorithm [39]

We evaluate the effectiveness of the three acceleration techniques: 1) parallelization of disjoint  $\alpha$ -expansions, and 2) acceleration of unary cost computation by GPU and 3) by cost filtering. For this, we compare following six variants of our method: LE-BF and LE-GF using one or four CPU cores (denoted as CPUx1 and CPUx4), and LE-BF using GPU and one or four CPU cores (denoted as GPU+CPUx1 and GPU+CPUx4). Additionally, we compare with our previous algorithm of [39] with CPU and GPU implementations (denoted as LSL with CPUx1 and GPU+CPUx4). We use the *Rocks1* dataset by estimating the disparities of only the left image without post-processing. Figures 9 (a)–(c) show



TABLE 3

Effect of our post-processing (PP) and RANSAC proposer (RP). We show weighted average bad 2.0 scores of our LE-GF method with 3DMST [28] that ranks second following ours. We use 15 training image pairs from Middlebury V3. We also show average (avr) and standard deviation (stdev) by ten sets of different random initializations. See also Sec. 4.7 for more analysis on the RANSAC proposer.

	PP	RP	nonocc	all
Ours	✓	✓	6.52	12.1
		✓	6.65	13.6
			7.72	14.6
avr ± stdev	✓	✓	6.63	12.3
			±0.12	±0.2
ref. 3DMST [28]	✓		7.08	12.9

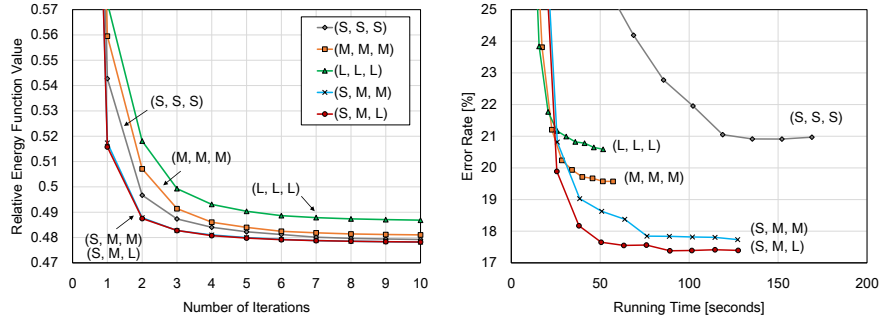


Fig. 7. Effect of grid-cell sizes. We use LE-GF with different combinations of grid structures. The S, M, and L denote small, medium, and large grid-cells, respectively. The joint use of different sizes of grid-cells improves the performance. See also Fig. 8 for visual comparison.

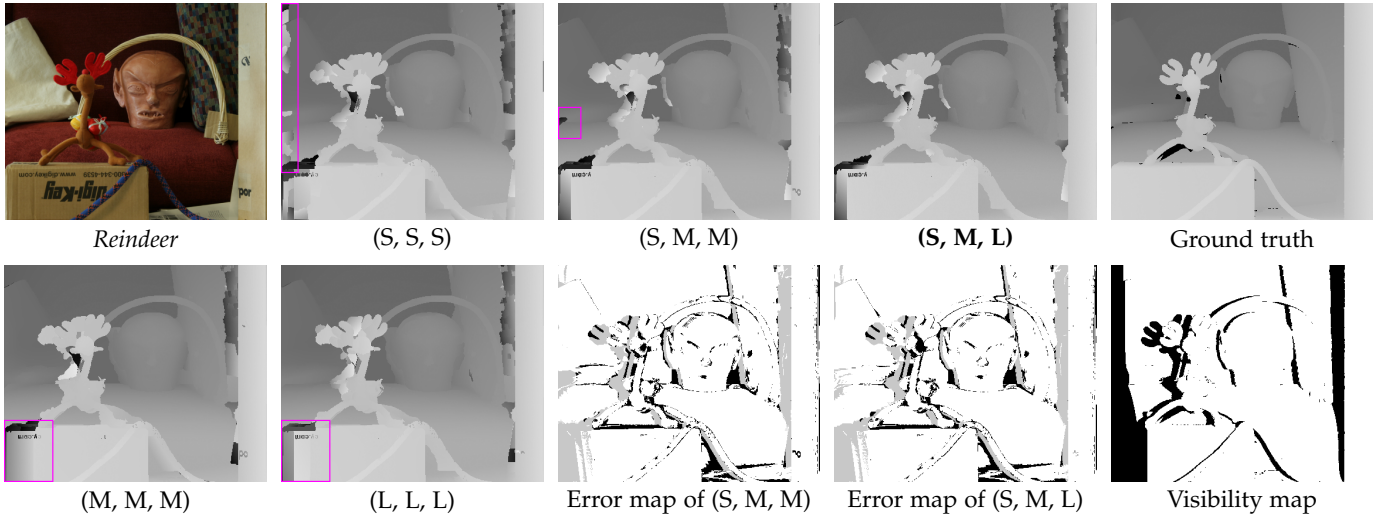


Fig. 8. Visual effect of grid-cell sizes. The use of larger grid-cells leads to smoother solutions and effective for occluded regions. The proposed combination (S, M, L) well balances localization and spatial propagation and performs best. These are all raw results without post-processing.

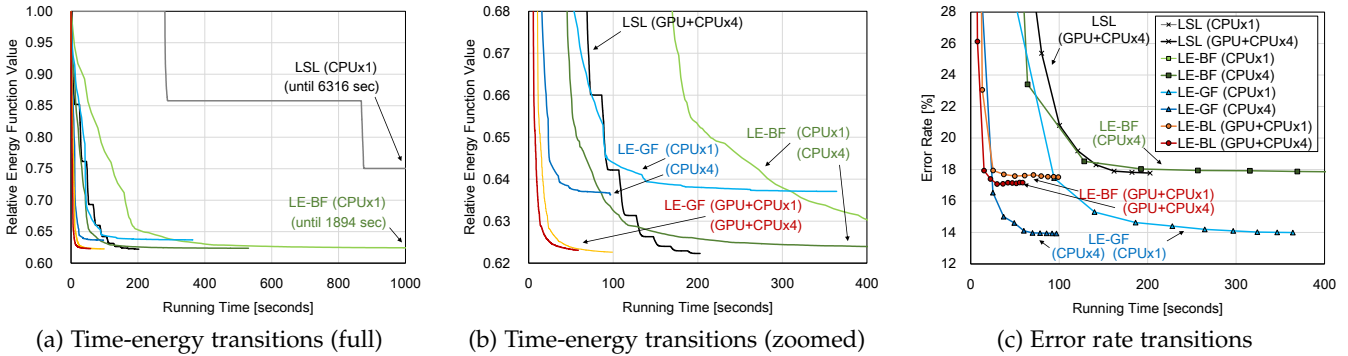


Fig. 9. Efficiency evaluation in comparison to our previous algorithm (LSL) [39]. Accuracies are evaluated for all-regions at each iteration.

the temporal transitions of the energy values in the full and zoomed scales, and the subpixel error rates, respectively. Note that energy values in Fig. 9 are evaluated by the energy function of LE-BF and LSL for better visualization.

**Parallelization of local  $\alpha$ -expansions on CPU:** Comparing CPUx1 and CPUx4, we observe about 3.5x of speed-up for both LE-BF and LE-GF. Note that the converged energy values of LE-GF are relatively higher than the others because it optimizes a different energy function. However, if we see Fig. 9 (c), it actually finds better solutions in this example.

On the other hand, speed-up for LE-BF from GPU+CPUx1 to GPU+CPUx4 is limited to about 1.7x. This is because the unary cost computation is already parallelized by GPU and this speed-up is accounted for by the parallelization of other parts, *e.g.*, the computation of pairwise terms and min-cut.

**Parallelization of unary cost computation on GPU:** By comparing GPU+CPUx1 and CPUx1 of LE-BF, we obtain about 19x of speed-up. This is relatively smaller than 32x of speed-up in our previous LSL method, mainly due to the larger overhead of GPU data transferring.

**Fast cost filtering:** By comparing LE-GF and LE-BF methods, we observe about 5.3x speed-up for both CPUx1 and CPUx4 versions. Notably, even a CPU implementation of LE-GF (CPUx4) achieves almost comparable efficiency with a GPU implementation of LE-BF (GPU+CPUx1).

**Comparison to our previous method [39]:** In comparison to our previous LSL method [39], our LE-BF shows much faster convergence than LSL in the same single-core results (CPUx1). We consider there are mainly three factors contributing to this speed-up; First, the batch-cycle algorithm adopted in [39] makes its convergence slower; Second, we use cost filtering and compute its first step of Eq. (18) in  $O(1)$ , while [39] computes each unary term individually in  $O(|W|)$ ; Finally, we have removed a per-pixel label refinement step of [39] which requires additional unary-cost computations. Similar speed-up can be observed from LSL (GPU+CPUx4) to LE-BF (GPU+CPUx4).

#### 4.5 Comparison with PMBP [3]

We compare our method with PMBP [3] that is the closest method to ours. For a fair comparison, we use four neighbors for  $\mathcal{N}$  in Eq. (5) as the same setting with PMBP. For a comparable smoothness weight with the default setting (eight-neighbor  $\mathcal{N}$ ), we use  $\lambda = 40$  for LE-BF and  $\lambda = 2$  for LE-GF, and keep the other parameters as default. For PMBP, we use the same model as ours; the only difference from the original PMBP is the smoothness term, which does not satisfy the submodularity of Eq. (2). In PMBP, it defines  $K$  candidate labels for each pixel, for which we set  $K = 1$  and  $K = 5$  (original paper uses  $K = 5$ ). We show the comparison using the *Cones* dataset by estimating the disparities of only the left image without post-processing.

Figures 10 (a)–(c) show the temporal transitions of the energy values in the full and zoomed scales, and the sub-pixel error rates, respectively. We show the performance of our method using its GPU and CPU (one or four CPU cores) implementations. For PMBP, we also implemented the unary cost computation on GPU, but it became rather slow, due to the overhead of data transfer. Efficient GPU implementations for PMBP are not available in literature.<sup>5</sup> Therefore, the plots show PMBP results that use a single CPU core. Figures 10 (a) and (b) show that, even with a single CPU-core implementation, our LE-BF and LE-GF show comparable or even faster convergence than PMBP. With CPU and GPU parallelization, our methods achieve much faster convergence than PMBP. Furthermore, our methods reach lower energies with greater accuracies than PMBP at the convergence. In Fig. 11 we compare the results of our LE-GF method and PMBP with  $K = 5$ . While PMBP yields noisy disparities, our method finds smoother and better disparities around edges and at occluded regions.

#### 4.6 Comparison with PMF [32]

We also compare our LE-GF method with PMF [32] using the same data term as ours. For PMF, the number of super-

pixels  $K$  is set to 300, 500, and 700 as used in [32] ( $K = 500$  is the default in [32]), and we sufficiently iterate 30 times. Both LE-GF and PMF are run using a single CPU core.

Figures 12 (a)–(c) show the temporal transitions of the energy values in the full and zoomed scales, and the sub-pixel error rates, respectively. As shown in Figs. 12 (a) and (b), PMF converges at higher energies than ours, since it cannot explicitly optimize pairwise smoothness terms because it is a local method. Furthermore, although energies are reduced almost monotonically in PMF, the transitions of accuracies are not stable and even degrade after many iterations. This is also because of the lack of explicit smoothness regularizer, and PMF converges at a bad local minimum. Figure 13 compares the results of our method and PMF with  $K = 500$ . Again, our methods find smoother and better disparities around edges and at occluded regions.

#### 4.7 Comparison with Olsson *et al.* [33]

We compare with a stereo method by Olsson *et al.* [33], from which we borrow our smoothness term. As differences from our method, their data term does not use accurate slanted patch matching and they use conventional fusion-based optimization using planar proposals generated by RANSAC-based local plane fitting (see [33] for details). To compare optimization approaches of ours and [33], we replace the data term of [33] with ours so that the two methods use the same energy function as ours. We here use two variants of our LE-GF method that remove the RANSAC proposer and the regularization term, which are components closely related to [33]. We disable post-processing in all methods.

Table 4 shows error rates for *Teddy* from Middlebury V2 and three more datasets from V3, where the proposed method outperforms [33] and others. As shown in Figs. 14 and 15, our method has faster convergence and produces qualitatively better smooth surfaces than [33]. Spatial label propagation not only accelerates inference but also introduces an implicit smoothness effect [5], which helps our method to find smoother surfaces. Our method mostly performs comparably well regardless of the RANSAC proposer, outperforming [33]. On the other hand, the RANSAC proposer promotes convergence (see Fig. 14) and also effectively recovers large texture-less planes in *Vintage* (Fig. 16). Using a more advanced proposer [34] may lead to further improvements. Also note that, except for *Vintage* in Table 4, our method performs well for non-occluded regions even without regularization (as comparable to [33]). The slanted patch matching term alone has an implicit second-order smoothness effect [5], and the regularizer of [33] further enforces it especially for occluded or texture-less regions.

## 5 CONCLUSIONS

In this paper, we have presented an accurate and efficient stereo matching method for continuous 3D plane label estimation. Unlike previous approaches that use fusion moves [25], [33], [45], our method is subproblem optimal and only requires a randomized initial solution. By comparing with a recent continuous MRF stereo method, PMBP [3], our method has shown an advantage in efficiency and comparable or greater accuracy. The use of a GC-based optimizer makes our method advantageous.

5. GPU-parallelization schemes of BP are not directly applicable due to PMBP’s unique settings. The “jump flooding” used in the original PatchMatch [1] reports 7x speed-ups by GPU. However, because it propagates candidate labels to distant pixels, it is not applicable to PMBP that must propagate messages to *neighbors*, and is not as efficient as our 19x, either.

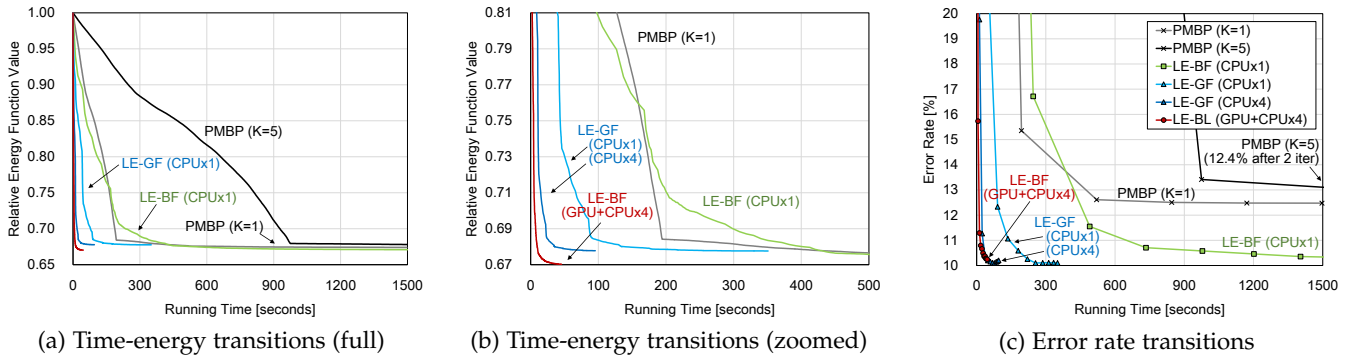


Fig. 10. Efficiency and accuracy comparison with PMBP [3]. Our methods achieve much faster convergence, reaching lower energies and better accuracies at the convergence. Accuracies are evaluated for all-regions at each iteration. See also Fig. 11 for visual comparison.

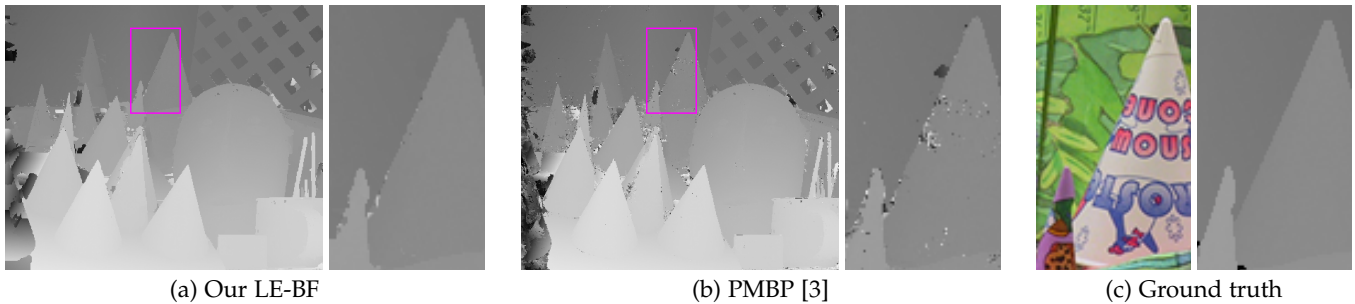


Fig. 11. Visual comparison with PMBP [3]. Our method finds smoother and better disparities around edges and at occluded regions.

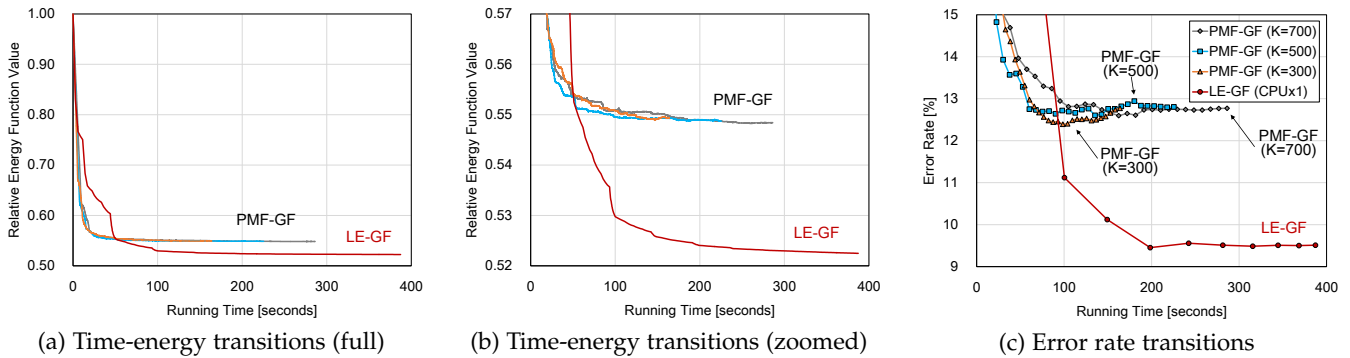


Fig. 12. Efficiency and accuracy comparison with PMF [32]. Our method stably improves the solution and reaches a lower energy with greater accuracy at the convergence. Accuracies are evaluated for all-regions at each iteration. See also Fig. 13 for visual comparison.

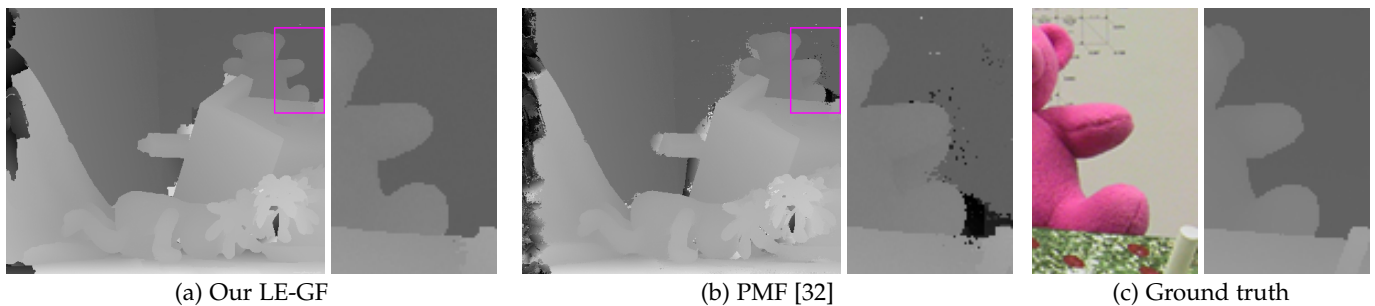


Fig. 13. Visual comparison with PMF [32]. With explicit smoothness regularization, our method finds smoother and better disparities around edges and at occluded regions.

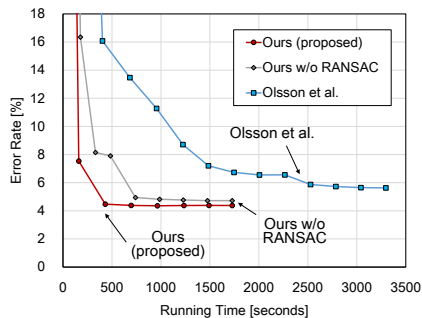
Furthermore, by using the subregion cost-filtering scheme developed in a local stereo method PMF [32], we achieve a fast algorithm and greater accuracy than PMF. As shown in [32], our method will be even faster by using a more light-weight constant-time filtering such as [31].

Follow-up works [18], [40] suggest that our optimization strategy can be used for more general correspondence field estimation than stereo. We also believe that occlusion handling schemes using GC [23], [44] can be incorporated into our framework, which may yield even greater accuracy.

TABLE 4

Accuracy comparison with Olsson *et al.* [33]. We use our energy function using 3D patch matching for our method and [33]. *Teddy* dataset is from Middlebury V2 (Sec. 4.1) and the other three from V3 (Sec. 4.2). Our local expansion move method (LE) outperforms fusion-based optimization of [33]. While our method has comparably good accuracies regardless of the RANSAC proposer, it is effective for large texture-less planes in *Vintage*.

	Reg.	Options Optimization	Teddy (from V2)			Adirondack		ArtL		Vintage	
			nonocc	all	disc	nonocc	all	nonocc	all	nonocc	all
Ours (proposed LE-GF)	[33]	LE+RANSAC	<b>3.98</b>	9.46	<b>12.8</b>	<b>1.19</b>	<b>4.36</b>	<b>3.97</b>	<b>15.5</b>	<b>5.65</b>	<b>11.7</b>
Ours w/o RANSAC	[33]	LE	4.26	<b>9.19</b>	13.6	<b>1.20</b>	4.63	<b>3.98</b>	<b>15.3</b>	13.7	18.7
Ours w/o regularization	None	LE+RANSAC	5.47	13.0	16.0	2.56	7.74	4.12	17.6	22.8	27.9
Olsson <i>et al.</i> [33] + 3DPM	[33]	Fusion+RANSAC	5.21	9.98	15.8	2.74	5.42	4.41	15.8	5.94	<b>11.8</b>



(For *Adirondack*, by a single CPU core)

Fig. 14. Efficiency and accuracy comparison with Olsson *et al.* [33]. Our method is faster than [33] even without parallelization. The RANSAC proposer promotes the inference.

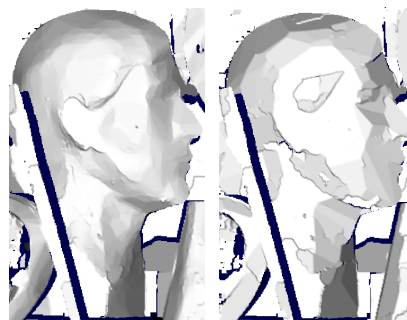


Fig. 15. 3D visual comparison with Olsson *et al.* [33] for *ArtL* dataset. Although we use the same energy function, the result by [33] is strongly biased to piecewise planar surfaces.

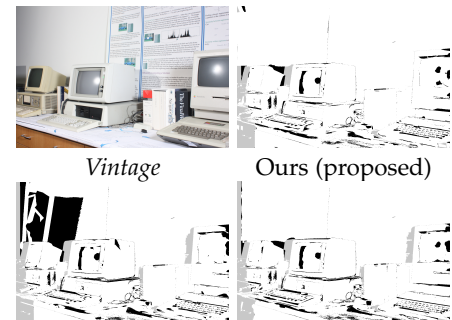


Fig. 16. Effectiveness of the RANSAC-based plane proposer for large texture-less regions. We show error maps of our method and Olsson *et al.* [33] for *Vintage* dataset.

## ACKNOWLEDGMENTS

The authors would like to thank the anonymous reviewers including those of the conference paper. This work was supported by JSPS KAKENHI Grant Number 14J09001 and Microsoft Research Asia Ph.D. Fellowship.

## REFERENCES

- [1] C. Barnes, E. Shechtman, A. Finkelstein, and D. B. Goldman. PatchMatch: a randomized correspondence algorithm for structural image editing. *Proc. of SIGGRAPH (ACM Trans. on Graph.)*, 28(3):24:1–24:11, 2009.
- [2] C. Barnes, E. Shechtman, D. B. Goldman, and A. Finkelstein. The generalized patchmatch correspondence algorithm. In *Proc. of European Conf. on Computer Vision (ECCV)*, pages 29–43, 2010.
- [3] F. Besse, C. Rother, A. Fitzgibbon, and J. Kautz. PMBP: PatchMatch Belief Propagation for Correspondence Field Estimation. In *Proc. of British Machine Vision Conf. (BMVC)*, pages 132.1–132.11, 2012.
- [4] S. Birchfield and C. Tomasi. Multiway cut for stereo and motion with slanted surfaces. In *Proc. of Int'l Conf. on Computer Vision (ICCV)*, volume 1, pages 489–495, 1999.
- [5] M. Bleyer, C. Rhemann, and C. Rother. PatchMatch Stereo - Stereo Matching with Slanted Support Windows. In *Proc. of British Machine Vision Conf. (BMVC)*, pages 14.1–14.11, 2011.
- [6] M. Bleyer, C. Rother, and P. Kohli. Surface stereo with soft segmentation. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 1570–1577, 2010.
- [7] Y. Boykov and V. Kolmogorov. An experimental comparison of min-cut/max-flow algorithms for energy minimization in vision. *IEEE Trans. Pattern Anal. Mach. Intell.*, 26(9):1124–1137, 2004.
- [8] Y. Boykov, O. Veksler, and R. Zabih. Fast approximate energy minimization via graph cuts. *IEEE Trans. Pattern Anal. Mach. Intell.*, 23(11):1222–1239, 2001.
- [9] O. Chum, J. Matas, and J. Kittler. Locally Optimized RANSAC. In *DAGM-Symposium*, volume 2781, pages 236–243, 2003.
- [10] S. Drouyer, S. Beucher, M. Bilodeau, M. Moreaud, and L. Sorbier. Sparse Stereo Disparity Map Densification using Hierarchical Image Segmentation. In *Proc. of Int'l Symp. on Mathematical Morphology*, pages 172–184, 2017.
- [11] P. Felzenszwalb and D. Huttenlocher. Efficient belief propagation for early vision. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, volume 1, pages 261–268, 2004.
- [12] S. Geman and D. Geman. Stochastic Relaxation, Gibbs Distributions, and the Bayesian Restoration of Images. *IEEE Trans. Pattern Anal. Mach. Intell.*, 6(6):721–741, 1984.
- [13] K. He, J. Sun, and X. Tang. Guided image filtering. *IEEE Trans. Pattern Anal. Mach. Intell.*, 35(6):1397–1409, 2013.
- [14] P. Heise, S. Klose, B. Jensen, and A. Knoll. PM-Huber: PatchMatch with Huber Regularization for Stereo Matching. In *Proc. of Int'l Conf. on Computer Vision (ICCV)*, pages 2360–2367, 2013.
- [15] L. Hong and G. Chen. Segment-based stereo matching using graph cuts. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, volume 1, pages 74–81, 2004.
- [16] A. Hosni, M. Gelautz, and M. Bleyer. Accuracy-Efficiency Evaluation of Adaptive Support Weight Techniques for Local Stereo Matching. In *Proc. DAGM/OAGM Symposium*, pages 337–346, 2012.
- [17] A. Hosni, C. Rhemann, M. Bleyer, C. Rother, and M. Gelautz. Fast cost-volume filtering for visual correspondence and beyond. *IEEE Trans. Pattern Anal. Mach. Intell.*, 35(2):504–511, 2013.
- [18] J. Hur and S. Roth. MirrorFlow: Exploiting Symmetries in Joint Optical Flow and Occlusion Estimation. In *Proc. of Int'l Conf. on Computer Vision (ICCV)*, 2017. (to appear).
- [19] H. Ishikawa. Exact optimization for Markov random fields with convex priors. *IEEE Trans. Pattern Anal. Mach. Intell.*, 25(10):1333–1336, 2003.
- [20] K. R. Kim and C. S. Kim. Adaptive smoothness constraints for efficient stereo matching using texture and edge information. In *Proc. of Int'l Conf. on Image Processing (ICIP)*, pages 3429–3433, 2016.
- [21] V. Kolmogorov. Convergent Tree-Reweighted Message Passing for Energy Minimization. *IEEE Trans. Pattern Anal. Mach. Intell.*, 28(10):1568–1583, 2006.
- [22] V. Kolmogorov and C. Rother. Minimizing Nonsubmodular Functions with Graph Cuts – A Review. *IEEE Trans. Pattern Anal. Mach. Intell.*, 29(7):1274–1279, 2007.
- [23] V. Kolmogorov and R. Zabih. Computing visual correspondence with occlusions using graph cuts. In *Proc. of Int'l Conf. on Computer Vision (ICCV)*, volume 2, pages 508–515, 2001.
- [24] V. Kolmogorov and R. Zabih. What energy functions can be minimized via graph cuts? *IEEE Trans. Pattern Anal. Mach. Intell.*, 26(2):147–159, 2004.
- [25] V. Lempitsky, C. Rother, S. Roth, and A. Blake. Fusion Moves

- for Markov Random Field Optimization. *IEEE Trans. Pattern Anal. Mach. Intell.*, 32(8):1392–1405, 2010.
- [26] A. Li, D. Chen, Y. Liu, and Z. Yuan. Coordinating multiple disparity proposals for stereo computation. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 4022–4030, 2016.
- [27] G. Li and S. W. Zucker. Differential geometric inference in surface stereo. *IEEE Trans. Pattern Anal. Mach. Intell.*, 32(1):72–86, 2010.
- [28] L. Li, X. Yu, S. Zhang, X. Zhao, and L. Zhang. 3d cost aggregation with multiple minimum spanning trees for stereo matching. *Appl. Opt.*, 56(12):3411–3420, 2017.
- [29] L. Li, S. Zhang, X. Yu, and L. Zhang. PMSC: PatchMatch-Based Superpixel Cut for Accurate Stereo Matching. *IEEE Transactions on Circuits and Systems for Video Technology*, PP(99):1–1, 2017.
- [30] J. Liu and J. Sun. Parallel graph-cuts by adaptive bottom-up merging. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 2181–2188, 2010.
- [31] J. Lu, K. Shi, D. Min, L. Lin, and M. Do. Cross-based local multipoint filtering. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 430–437, 2012.
- [32] J. Lu, H. Yang, D. Min, and M. N. Do. Patch Match Filter: Efficient Edge-Aware Filtering Meets Randomized Search for Fast Correspondence Field Estimation. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 1854–1861, 2013.
- [33] C. Olsson, J. Ulen, and Y. Boykov. In Defense of 3D-Label Stereo. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 1730–1737, 2013.
- [34] C. Olsson, J. Uln, and A. Eriksson. Local Refinement for Stereo Regularization. In *Proc. of Int’l Conf. on Pattern Recognition (ICPR)*, pages 4056–4061, 2014.
- [35] G. Petschnigg, R. Szeliski, M. Agrawala, M. Cohen, H. Hoppe, and K. Toyama. Digital photography with flash and no-flash image pairs. *ACM Trans. on Graph.*, 23(3):664–672, 2004.
- [36] C. Rhemann, A. Hosni, M. Bleyer, C. Rother, and M. Gelautz. Fast cost-volume filtering for visual correspondence and beyond. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 3017–3024, 2011.
- [37] P. Strandmark and F. Kahl. Parallel and distributed graph cuts by dual decomposition. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 2085–2092, 2010.
- [38] R. Szeliski, R. Zabih, D. Scharstein, O. Veksler, V. Kolmogorov, A. Agarwala, M. Tappen, and C. Rother. A Comparative Study of Energy Minimization Methods for Markov Random Fields with Smoothness-Based Priors. *IEEE Trans. Pattern Anal. Mach. Intell.*, 30(6):1068–1080, 2008.
- [39] T. Taniai, Y. Matsushita, and T. Naemura. Graph Cut based Continuous Stereo Matching using Locally Shared Labels. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 1613–1620, 2014.
- [40] T. Taniai, S. N. Sinha, and Y. Sato. Joint Recovery of Dense Correspondence and Cosegmentation in Two Images. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 4246–4255, 2016.
- [41] H. Tao, H. Sawhney, and R. Kumar. A global matching framework for stereo computation. In *Proc. of Int’l Conf. on Computer Vision (ICCV)*, volume 1, pages 532–539, 2001.
- [42] L. Wang and R. Yang. Global stereo matching leveraged by sparse ground control points. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 3033–3040, 2011.
- [43] Z.-F. Wang and Z.-G. Zheng. A region based stereo matching algorithm using cooperative optimization. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2008.
- [44] Y. Wei and L. Quan. Asymmetrical Occlusion Handling Using Graph Cut for Multi-View Stereo. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 902–909, 2005.
- [45] O. Woodford, P. Torr, I. Reid, and A. Fitzgibbon. Global Stereo Reconstruction under Second-Order Smoothness Priors. *IEEE Trans. Pattern Anal. Mach. Intell.*, 31(12):2115–2128, 2009.
- [46] K. Yamaguchi, T. Hazan, D. McAllester, and R. Urtasun. Continuous markov random fields for robust stereo estimation. In *Proc. of European Conf. on Computer Vision (ECCV)*, pages 45–58, 2012.
- [47] J. S. Yedidia, W. T. Freeman, and Y. Weiss. Generalized Belief Propagation. In *Advances in Neural Information Processing Systems (NIPS)*, volume 13, pages 689–695, 2000.
- [48] K.-J. Yoon and I.-S. Kweon. Locally adaptive support-weight approach for visual correspondence search. In *Proc. of IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, volume 2, pages 924–931, 2005.
- [49] J. Zbontar and Y. LeCun. Stereo matching by training a convolutional neural network to compare image patches. *Journal of Machine Learning Research*, 17:1–32, 2016.



## APPENDIX A

### PROOF OF LEMMA 1

*Proof.* For the completeness, we repeat the original proof given in [33] with our notations. Regarding  $\bar{\psi}_{pq}(f_p, f_q)$ ,

$$\begin{aligned}
 \bar{\psi}_{pq}(\alpha, \alpha) + \bar{\psi}_{pq}(\beta, \gamma) &= \bar{\psi}_{pq}(\beta, \gamma) \\
 &= |d_p(\beta) - d_p(\gamma)| + |d_q(\beta) - d_q(\gamma)| \\
 &= |(d_p(\beta) - d_p(\alpha)) - (d_p(\gamma) - d_p(\alpha))| + |(d_q(\beta) - d_q(\alpha)) - (d_q(\gamma) - d_q(\alpha))| \\
 &\leq |d_p(\beta) - d_p(\alpha)| + |d_p(\gamma) - d_p(\alpha)| + |d_q(\beta) - d_q(\alpha)| + |d_q(\gamma) - d_q(\alpha)| \\
 &= \bar{\psi}_{pq}(\beta, \alpha) + \bar{\psi}_{pq}(\alpha, \gamma).
 \end{aligned}$$

Thus,  $\bar{\psi}_{pq}(f_p, f_q)$  satisfies the submodularity of expansion moves. For its truncated function,

$$\begin{aligned}
 \min(\bar{\psi}_{pq}(\alpha, \alpha), \tau) + \min(\bar{\psi}_{pq}(\beta, \gamma), \tau) &= \min(\bar{\psi}_{pq}(\beta, \gamma), \tau) \\
 &\leq \min(\bar{\psi}_{pq}(\beta, \alpha) + \bar{\psi}_{pq}(\alpha, \gamma), \tau) \\
 &\leq \min(\bar{\psi}_{pq}(\beta, \alpha), \tau) + \min(\bar{\psi}_{pq}(\alpha, \gamma), \tau).
 \end{aligned}$$

Therefore,  $\psi_{pq}(f_p, f_q)$  also satisfies the submodularity.  $\square$

## APPENDIX B

### ADDITIONAL RESULTS

We here provide additional results for the analyses on effects of grid-cell sizes (Sec. 4.3) and parallelization (Sec. 4.4) as well as the comparisons with the methods of PMBP [3] (Sec. 4.5), PMF [32] (Sec. 4.6) and Olsson *et al.* [33] (Sec. 4.7). We evaluate performances by showing plots of error rate transitions over running time, where error rates are evaluated by the *bad 2.0* metric for all regions. To provide diverse examples, we use 15 image pairs from the Middlebury V3 training dataset.

Unless noted otherwise, all methods are run on a single CPU core to optimize the same energy function that we use in Sec. 4.2 for the benchmark V3 evaluation. Settings of our method are also the same as we use in Sec. 4.2. Post-processing is disabled in all methods.

#### *Effect of grid-cell sizes*

Figure A1 shows analyses on effects of grid-cell sizes. Similarly to Sec. 4.3, we compare the proposed grid-cell combination (S, M, L) with other four combinations (L, L, L), (M, M, M), (S, S, S), and (S, M, M). Here, the sizes of “S”mall, “M”edium, and “L”arge grid-cells are proportionally set to 1%, 3% and 9% of the image width, respectively. For most of the image pairs, the proposed combination (S, M, L) performs best among the five combinations.

#### *Efficiency comparison by parallelization*

Figure A2 shows analyses on parallelization of our local expansion move algorithm. When our method is performed in parallel on four CPU cores, we observe about 3.8x of average speed-up among the 15 image pairs.

#### *Comparison with PMBP [3]*

Figure A3 compares our method with PMBP [3]. Similarly to Sec. 4.4, we compare PMBP (using  $K = 1$  and 5) with two variants of our method using guided image filtering (GF) and bilateral filtering (BF). Because PMBP cannot use the cost-map filtering acceleration of GF, we use only BF for PMBP. Therefore, our method using BF optimizes the same energy function with PMBP, while ours using GF optimizes a different function. All methods use four neighbors of pairwise terms as the same setting with PMBP.

As shown in Fig. A3, our methods both using GF and BF consistently outperform PMBP in accuracy at convergence, while ours using GF performs best and shows much faster convergence than the others. Our method can be further accelerated by parallelization as demonstrated in Fig. A2.

#### *Comparison with PMF [32]*

Figure A4 compares our method with PMF [32]. Both methods use the same energy function but PMF can only optimize its data term. With explicit regularization, our method consistently outperforms PMF for all the image pairs.

#### *Comparison with Olsson *et al.* [33]*

Figure A5 compares our method with the method by Olsson *et al.* [33]. Both methods use the same energy function but optimize it differently. For most of the image pairs, our method is about 6x faster to reach comparable or better accuracies than [33]. Our method can be further accelerated by parallelization as demonstrated in Fig. A2.

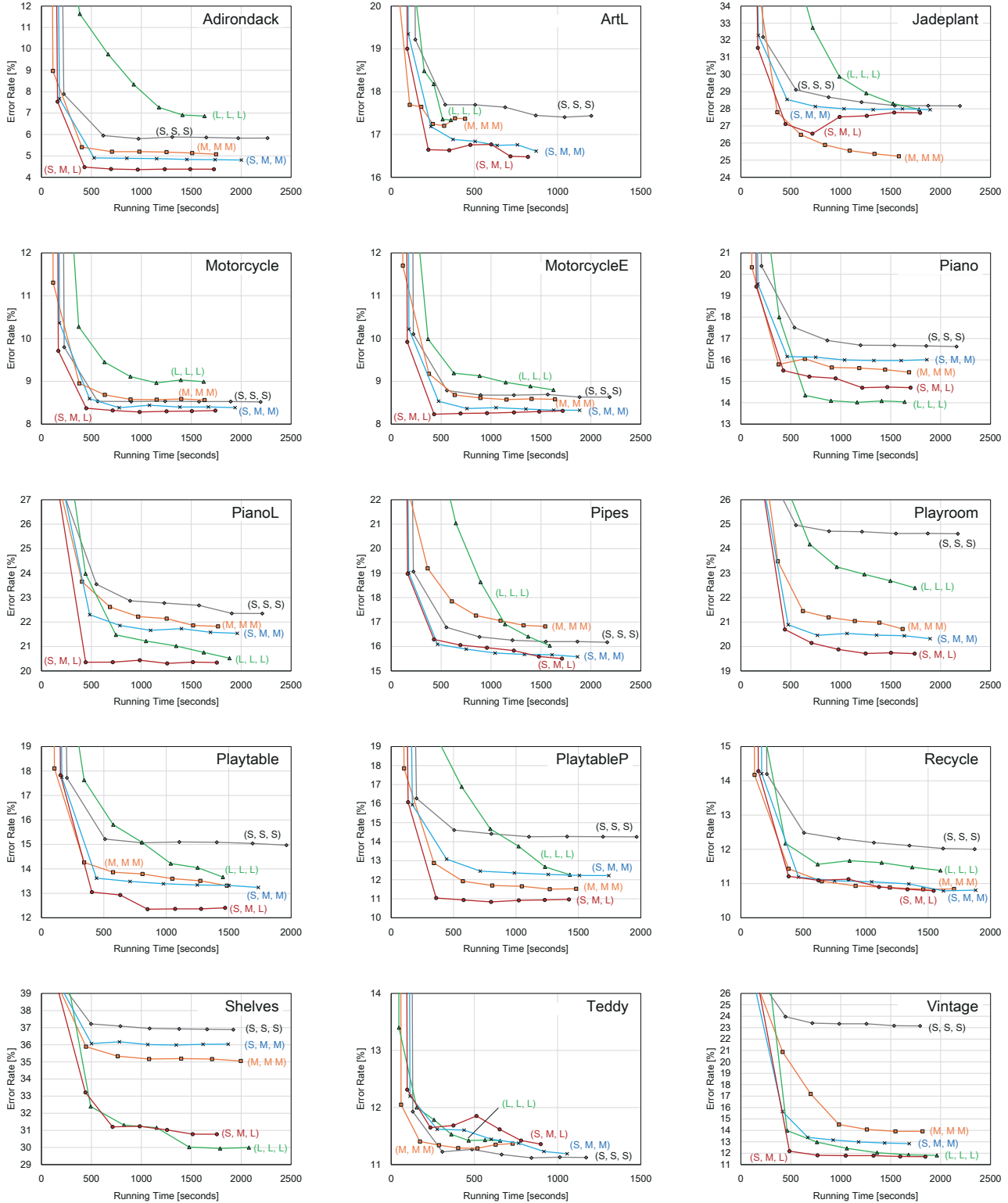


Fig. A1. Convergence comparison for different grid-cell sizes on 15 image pairs from the Middlebury V3 training dataset. Error rates are evaluated by the *bad 2.0* metric for all regions. For most of the image pairs, the proposed grid-cell combination (S, M, L) outperforms the other four combinations (L, L, L), (M, M, M), (S, S, S), and (S, M, M). Here, the sizes of “S” small, “M” medium, and “L” large grid-cells are proportionally set to 1%, 3% and 9% of the image width, respectively. All methods are run on a single CPU core to optimize the same energy function used in Sec 4.2 without post-processing.

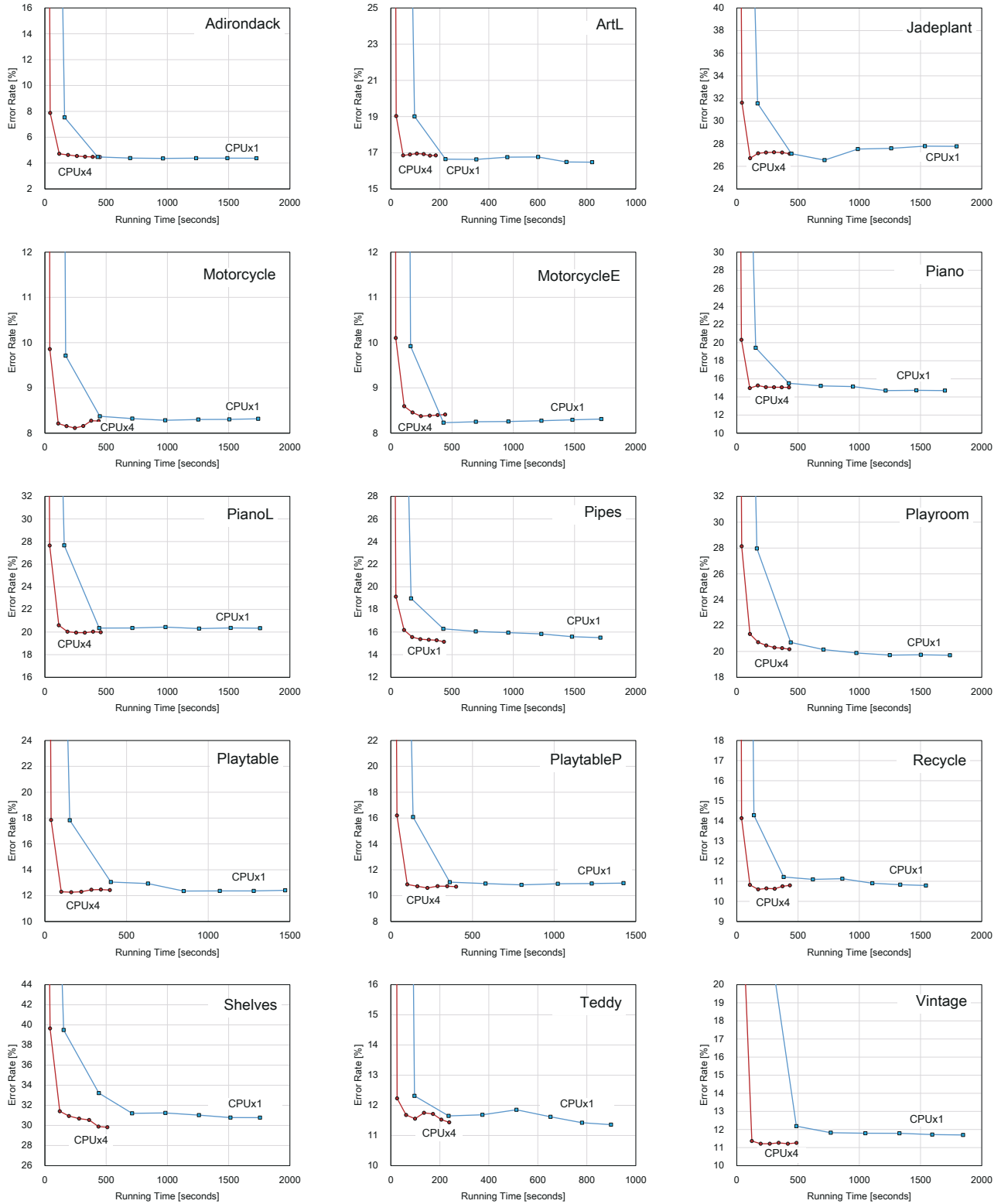


Fig. A2. Convergence comparison for parallelization on 15 image pairs from the Middlebury V3 training dataset. Error rates are evaluated by the *bad 2.0* metric for all regions. By performing local expansion moves in parallel on four CPU cores, we observe about 3.8x of average speed-up. Both methods optimize the same energy function used in Sec 4.2 without post-processing.

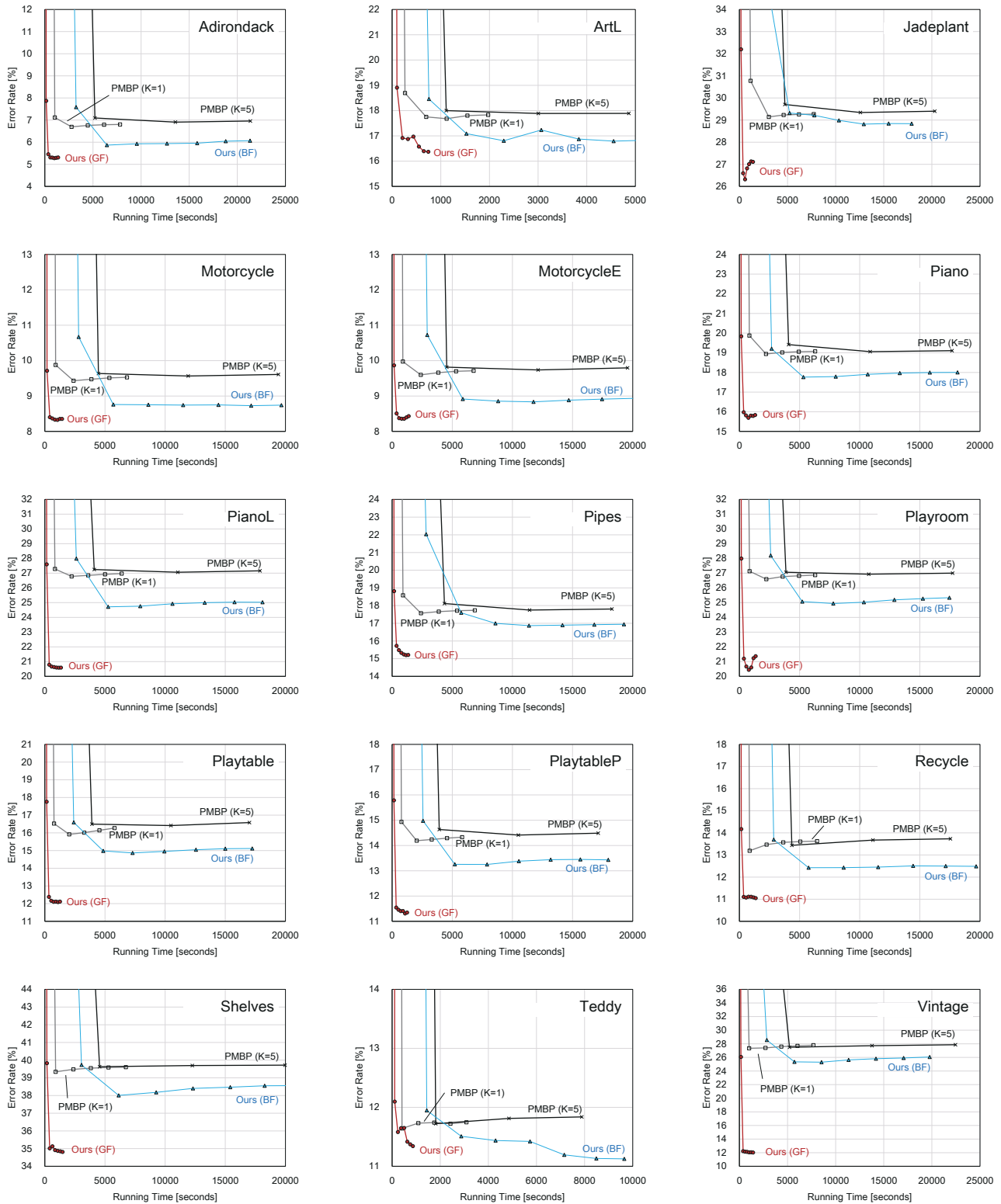


Fig. A3. Convergence comparison with PMBP [3] on 15 image pairs from the Middlebury V3 training dataset. Error rates are evaluated by the *bad* 2.0 metric for all regions. Our methods both using guided image filter (GF) and bilateral filter (BF) consistently outperform PMBP in accuracy at convergence. For most of the image pairs, our method using GF performs best, showing much faster convergence than the others. Note that both PMBP and ours using BF optimize the same energy function, but ours using GF optimizes a different function. All methods are run on a single CPU core without post-processing. Our method can be further accelerated by parallelization as demonstrated in Fig. A2.

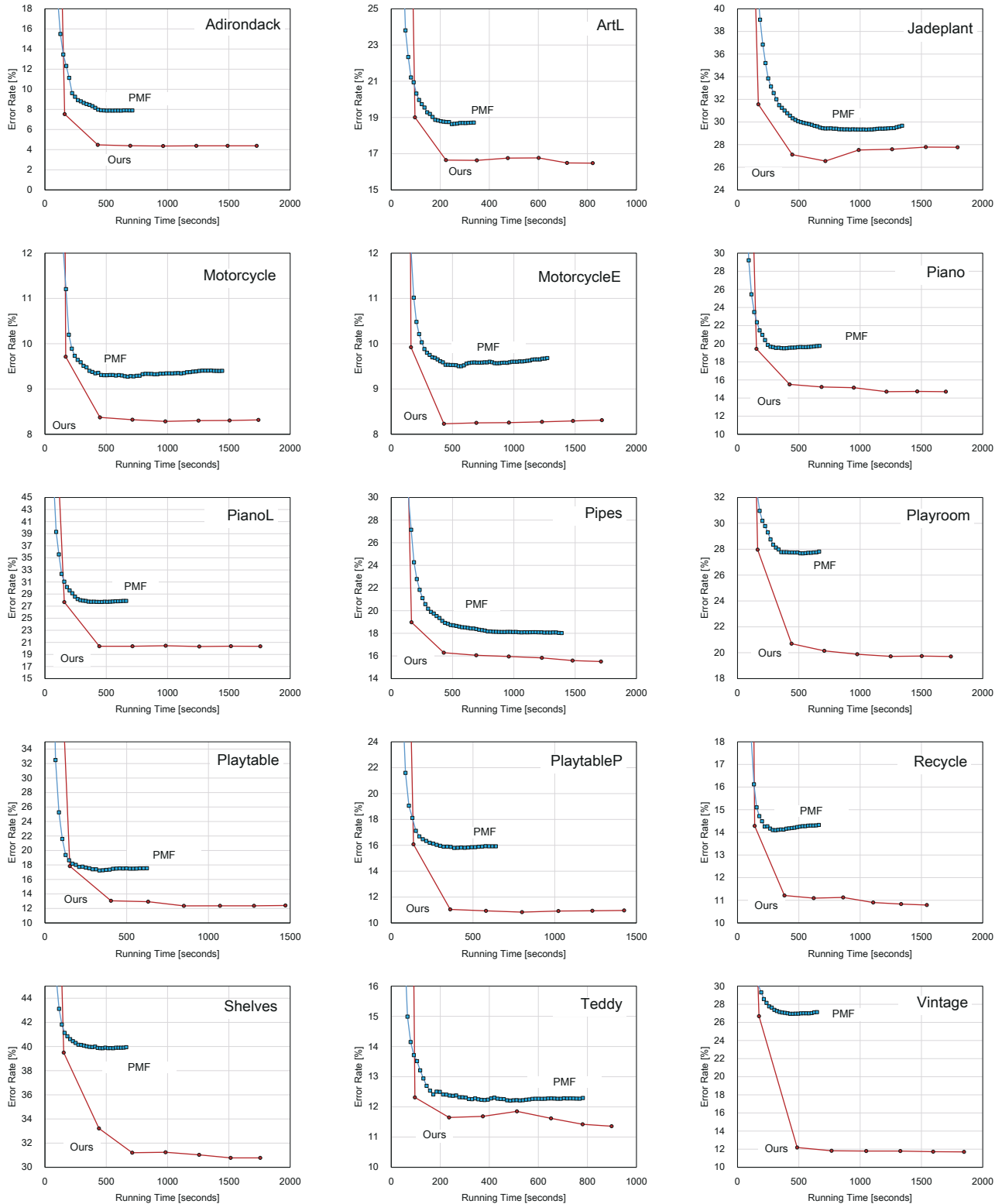


Fig. A4. Convergence comparison with PMF [32] on 15 image pairs from the Middlebury V3 training dataset. Error rates are evaluated by the *bad 2.0* metric for all regions. Our method consistently outperforms PMF. Because performances of PMF by  $K = 300, 500, 700$  had very similar trends, we here only show results by  $K = 500$  (default setting) to avoid cluttered profiles. Both methods are run on a single CPU core using the same energy function in Sec 4.2 without post-processing, but PMF can optimize only its data term. Note that while PMF mostly converged in 30 iterations, we did 60 iterations for some cases (*Jadeplant, Motorcycle, MotorcycleE, Pipes, and Teddy*).



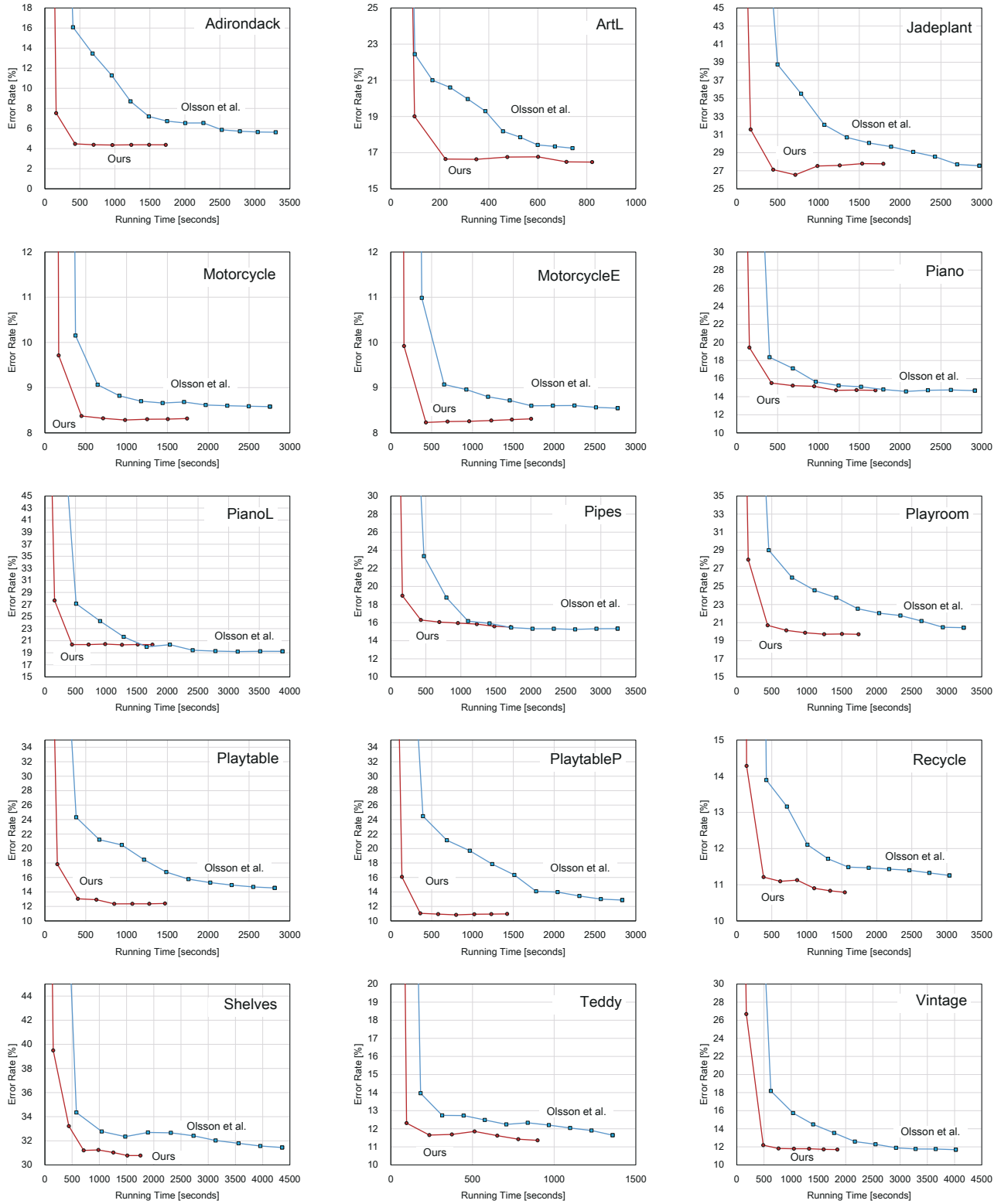


Fig. A5. Convergence comparison with the optimization method used in Olsson *et al.* [33] on 15 image pairs from the Middlebury V3 training dataset. Error rates are evaluated by the *bad 2.0* metric for all regions. For most of the image pairs, our method is about 6x faster to reach comparable or better accuracies than Olsson *et al.* [33]. Both methods are run on a single CPU core to optimize the same energy function used in Sec 4.2 without post-processing. Our method can be further accelerated by parallelization as demonstrated in Fig. A2.