

Continuous Markov Random Fields for Robust Stereo Estimation

Koichiro Yamaguchi^{1,2}, Tamir Hazan¹, David McAllester¹,
and Raquel Urtasun¹

¹ TTI Chicago

² Toyota Central R&D Labs., Inc.

Abstract. In this paper we present a novel slanted-plane model which reasons jointly about occlusion boundaries as well as depth. We formulate the problem as one of inference in a hybrid MRF composed of both continuous (i.e., slanted 3D planes) and discrete (i.e., occlusion boundaries) random variables. This allows us to define potentials encoding the ownership of the pixels that compose the boundary between segments, as well as potentials encoding which junctions are physically possible. Our approach outperforms the state-of-the-art on Middlebury high resolution imagery [1] as well as in the more challenging KITTI dataset [2], while being more efficient than existing slanted plane MRF methods, taking on average 2 minutes to perform inference on high resolution imagery.

1 Introduction

Over the past few decades we have witnessed a great improvement in performance of stereo algorithms. Most modern approaches frame the problem as inference on a Markov random field (MRF) and utilize global optimization techniques such as graph cuts or message passing [3] to reason jointly about the depth of each pixel in the image.

A leading approach to stereo vision uses slanted-plane MRF models which were introduced a decade ago [4]. Most methods [5–10] assume a fixed set of superpixels on a reference image, say the left image, and model the surface under each superpixel as a slanted plane. The MRF typically has a robust data term scoring the assigned plane in terms of a matching score induced by the plane on the pixels contained in the superpixel. This data term often incorporates an explicit treatment of occlusion — pixels in one image that have no corresponding pixel in the other image [11, 12, 8, 13]. Slanted-plane models also typically include a robust smoothness term expressing the belief that the planes assigned to adjacent superpixels should be similar.

A major issue with slanted-plane stereo models is their computational complexity. For example, [13] reports an average of approximately one hour of computation for each low-resolution Middlebury stereo pair. This makes these approaches impractical for applications such as robotics or autonomous driving. A main source of difficulty is the fact that each plane is defined by three continuous parameters and inference for continuous MRFs with non-convex energies is computationally challenging.

This paper contains two contributions. First, we introduce the use of junction potentials, described below, into this class of models. Second, we show that particle methods can achieve state-of-the-art performance with reasonable inference times on very challenging high-resolution imagery, i.e., KITTI dataset [2].

Junction potentials originate in early line labeling algorithms [14, 15]. These algorithms assign labels to the lines of a line drawing where the label indicate whether the line represents a discontinuity due to changes in depth (an occlusion), surface orientation (a corner), lighting (a shadow) or albedo (paint). A *junction* is a place where three lines meet. Only certain combinations of labels are physically realizable at junctions. The constraints on label combinations at junctions often force the labeling of the entire line drawing [14]. Here, following work on monocular image interpretation [16–18], we label the boundaries between image segments with labels –“left occlusion”, “right occlusion”, “hinge” or “coplanar”. In our model the occlusion labels play a role in the data term, where they are interpreted as expressing ownership of the pixels that compose the boundary between segments — an occlusion boundary is “owned” by the foreground object.

Our second contribution is to show that particle methods can be used to perform inference in high resolution imagery with reasonable running times. Particle methods avoid premature commitment to any fixed quantization of continuous variables and hence allow a precise exploration of the continuous space. Our particle inference method is based on the recently developed particle convex belief propagation (PCBP) [19]. Furthermore, we learn the contribution of each potential via the primal-dual framework of [20].

In the remainder of the paper we first review related work. We then introduce our continuous MRF model for stereo and show how to do learning and inference in this model. Finally, we demonstrate the effectiveness of our approach and show that it outperforms the state-of-the-art for high resolution Middlebury imagery [1] as well as in the more challenging KITTI dataset [2].

2 Related Work

In the past few years much progress has been made towards solving the stereo problem, as evidenced by Scharstein et al. overview [21]. Local methods typically aggregate image statistics in a small window, thus imposing smoothness implicitly. Optimization is usually performed using a winner-takes-all strategy, which selects for each pixel the disparity with the smallest value under some distance metric [21]. Traditional local methods [22] often suffer from border bleeding effects or struggle with correspondence ambiguities. Approaches based on adaptive support windows [23, 24] adjust their computations locally to improve performance, especially close to border discontinuities. This results in better performance at the price of more computation.

Hirschmüller proposed semi-global matching [25], an approach which extends polynomial time 1D scan-line methods to propagate information along 16 orientations. This reduces streaking artifacts and improves accuracy compared to

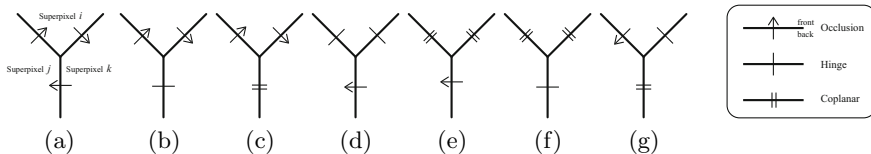


Fig. 1. Impossible cases of 3-way junctions. (a) 3 cyclic occlusions, (b) hinge and 2 occlusion with opposite directions, (c) coplanar and 2 occlusion with opposite directions, (d) 2 hinge and occlusion, (e) 2 coplanar and occlusion, (f) 2 coplanar and hinge, (g) hinge, coplanar, and occlusion (supapixel with coplanar boundary is in front).

traditional methods. In this paper we employ this technique to compute a disparity map from which we build our potentials. In [26, 27] disparities are ‘grown’ from a small set of initial correspondence seeds. Though these methods produce accurate results and can be faster than global approaches, they do not provide dense matching and struggle with textureless and distorted image areas. Approaches to reduce the search space have been investigated for global stereo methods [28, 29] as well as local methods [30].

Dense and accurate matching can be obtained by global methods, which enforce smoothness explicitly by minimizing an MRF-based energy function. These MRFs can be formulated at the pixel level [31], however, the smoothness is then defined very locally. Slatend-plane MRF models for stereo vision were introduced in [4] and have been since very widely used [5–7, 9, 10, 13]. In the context of this literature, our work has several distinctive features. First, we use a novel model involving “boundary labels”, “junction potentials”, and “edge ownership”. Second, for inference we employ the convex form of the particle norm-product belief propagation [32], which we refer to as particle convex belief propagation (PCBP) [19]. In contrast, some previous works used particle belief propagation (PBP) [33, 34, 10] which correspond to non-convex norm-product with the Bethe entropy approximation. The efficiency and convexity of PCBP makes it possible to evaluate our approach on hundreds of high-resolution images [2], whereas previous empirical evaluations of slanted-plane models have largely been restricted to the low-resolution versions of the small number of highly controlled Middlebury images. Third, we use a training algorithm based on primal-dual approximate inference [35] which allow us to effectively learn the importance of each potential.

3 Continuous MRFs for Stereo

In this section we describe our approach to joint reasoning of boundary labels and depth. We reason at the segment level, employing a richer representation than a discrete disparity label. In particular, we formulate the problem as inference in a hybrid conditional random field, which contains continuous and discrete random variables. The continuous random variables represent, for each segment, the disparities of all pixels contained in that segment in the form of a 3D slanted plane. The discrete random variables indicate for each pair of neighboring segments,

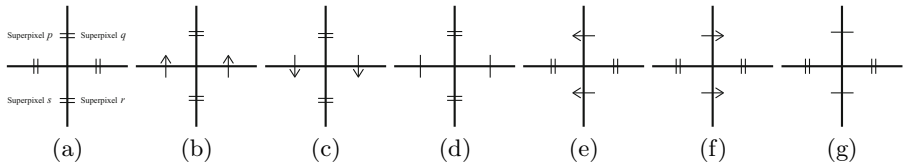


Fig. 2. Valid 4-way junctions. (a) 4 coplanar boundaries, (b)-(d) 2 coplanar vertical boundaries and 2 occlusion/hinge horizontal boundaries, (e)-(g) 2 coplanar horizontal boundaries and 2 vertical occlusion/hinge boundaries. A 4-way junction only appears in a region of uniform color.

whether they are co-planar, they form a hinge or there is a depth discontinuity (indicating which plane is in front of which).

More formally, let $y_i = (\alpha_i, \beta_i, \gamma_i) \in \mathfrak{R}^3$ be a random variable representing the i -th slanted 3D plane. We can compute the disparities of pixel $\mathbf{p} \in S_i$, where S_i is the set of pixels belonging to the i -th segment, as follows

$$\hat{d}_i(\mathbf{p}, \mathbf{y}_i) = \alpha_i(u - c_{ix}) + \beta_i(v - c_{iy}) + \gamma_i \quad (1)$$

with $\mathbf{p} = (u, v)$, and $\mathbf{c}_i = (c_{ix}, c_{iy})$ the center of the i -th segment. We have defined γ_i to be the disparity in the segment center as it improves the efficiency of PCBP inference. Let $o_{i,j} \in \{co, hi, lo, ro\}$ be a discrete random variable representing whether two neighboring planes are coplanar, form a hinge or an occlusion boundary. Here, *lo* implies that plane i occludes plane j , and *ro* represents that plane j occludes plane i .

We define our hybrid conditional random field as follows

$$p(\mathbf{y}, \mathbf{o}) = \frac{1}{Z} \prod_{\vartheta} \psi_{\vartheta}(\mathbf{y}_{\vartheta}) \prod_{\zeta} \psi_{\zeta}(\mathbf{o}_{\zeta}) \prod_{\tau} \psi_{\tau}(\mathbf{y}_{\tau}, \mathbf{o}_{\tau}) \quad (2)$$

where \mathbf{y} represents the set of all 3D slanted planes, \mathbf{o} the set of all discrete random variables, and $\psi_{\vartheta}, \psi_{\zeta}, \psi_{\tau}$ encode potential functions over sets of continuous, discrete or mixture of both types of variables. Note that \mathbf{y} contains three random variables for every segment in the image, and there is a random variable $o_{i,j}$ for each pair of neighboring segments.

In the following, we describe the different potentials we employed for our joint occlusion boundary and depth reasoning. For clarity, we describe the potentials in the log domain. Each type of potential will have a weight associated. All the weights \mathbf{w} will be learned using structure prediction methods [20].

3.1 Occlusion Boundary and Segmentation Potentials

Our approach takes as input a disparity image computed by any matching algorithm. In this paper we employ semi-global block matching [25]. Most matching methods return estimated disparity values on a subset of pixels. Let \mathcal{F} be the set of all pixels whose initial disparity has been estimated, and let $\mathcal{D}(\mathbf{p})$ be the

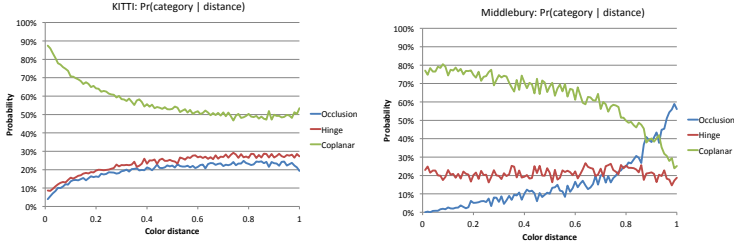


Fig. 3. Color statistics on (left) KITTI dataset, (right) Middlebury high-resolution

Table 1. SGBM matching functions in the validation set of KITTI

	> 2 pixels		> 3 pixels		> 4 pixels		> 5 pixels	
	Non-Occ	All	Non-Occ	All	Non-Occ	All	Non-Occ	All
Census	13.03 %	14.31 %	9.84 %	10.92 %	8.11 %	9.02 %	6.94 %	7.71 %
Gradient	8.68 %	9.82 %	6.04 %	6.96 %	4.73 %	5.47 %	3.90 %	4.51 %
Gradient + Census	8.19 %	9.39 %	5.55 %	6.52 %	4.26 %	5.05 %	3.47 %	4.13 %

disparity of pixel $\mathbf{p} \in \mathcal{F}$. Our model jointly reasons about segmentation in the form of occlusion boundaries as well as depth. We define potentials for each of these tasks individually as well as potentials which link both tasks. We start by defining truncated quadratic potentials, which we will employ in the definition of some of our potentials, i.e.,

$$\phi_i^{TP}(\mathbf{p}, \mathbf{y}_i, K) = \min \left(\left| \mathcal{D}(\mathbf{p}) - \hat{d}_i(\mathbf{p}, \mathbf{y}_i) \right|, K \right)^2, \quad \mathbf{p} \in S_i \cap \mathcal{F} \quad (3)$$

with K a constant threshold, and $\hat{d}_i(\mathbf{p}, \mathbf{y}_i)$ the disparity of pixel \mathbf{p} estimated as in Eq. 1. Note that we have made the quadratic potential robust via the min function. We now describe each of the potentials employed in more details.

Disparity Potential: We define truncated quadratic potentials for each segment encoding that the plane should agree with the results of the matching,

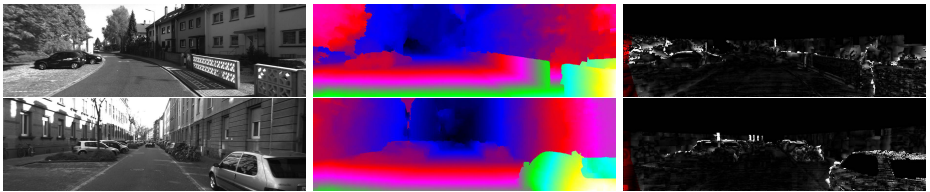
$$\phi_i^{\text{seg}}(\mathbf{y}_i) = \sum_{\mathbf{p} \in S_i \cap \mathcal{F}} \phi_i^{TP}(\mathbf{p}, \mathbf{y}_i, K)$$

Boundary Potential: We employ 3-way potentials linking our discrete and continuous variables. In particular, these potentials express the fact that when two neighboring planes are hinge or coplanar they should agree on the boundary, and when a segment occludes another segment, the boundary should be explained by the occluder. We thus define

$$\phi_{ij}^{\text{bdy1}}(o_{ij}, \mathbf{y}_i, \mathbf{y}_j) = \begin{cases} \sum_{\mathbf{p} \in B_{ij} \cap \mathcal{F}} \phi_i^{TP}(\mathbf{p}, \mathbf{y}_i, K) & \text{if } o_{ij} = lo \\ \sum_{\mathbf{p} \in B_{ij} \cap \mathcal{F}} \phi_j^{TP}(\mathbf{p}, \mathbf{y}_j, K) & \text{if } o_{ij} = ro \\ \frac{1}{2} \sum_{\mathbf{p} \in B_{ij} \cap \mathcal{F}} \phi_i^{TP}(\mathbf{p}, \mathbf{y}_i, K) + \phi_j^{TP}(\mathbf{p}, \mathbf{y}_j, K) & \text{if } o_{ij} = hi \vee co \end{cases}$$

Table 2. Comparison with the state of the art on the test set of KITTI [2]

	> 2 pixels		> 3 pixels		> 4 pixels		> 5 pixels	
	Non-Occ	All	Non-Occ	All	Non-Occ	All	Non-Occ	All
GC+occ [36]	39.76 %	40.97 %	33.50 %	34.74 %	29.86 %	31.10 %	27.39 %	28.61 %
OCV-BM [37]	27.59 %	28.97 %	25.39 %	26.72 %	24.06 %	25.32 %	22.94 %	24.14 %
CostFilter [38]	25.85 %	27.05 %	19.96 %	21.05 %	17.12 %	18.10 %	15.51 %	16.40 %
GCS [26]	18.99 %	20.30 %	13.37 %	14.54 %	10.40 %	11.44 %	8.63 %	9.55 %
GCSF [39]	20.75 %	22.69 %	13.02 %	14.77 %	9.48 %	11.02 %	7.48 %	8.84 %
SDM [27]	15.29 %	16.65 %	10.98 %	12.19 %	8.81 %	9.87 %	7.44 %	8.39 %
ELAS [30]	10.95 %	12.82 %	8.24 %	9.95 %	6.72 %	8.22 %	5.64 %	6.94 %
OCV-SGBM [25]	10.58 %	12.20 %	7.64 %	9.13 %	6.04 %	7.40 %	5.04 %	6.25 %
ITGV [40]	8.86 %	10.20 %	6.31 %	7.40 %	5.06 %	5.97 %	4.26 %	5.01 %
Ours	6.25 %	7.78 %	4.13 %	5.45 %	3.18 %	4.32 %	2.66 %	3.66 %

**Fig. 4.** KITTI examples. (Left) Original. (Middle) Disparity. (Right) Disparity errors.

where B_{ij} is the set of pixels around the boundary (within 2 pixels of the boundary) between segments i and j .

Compatibility Potential: We introduce an additional potential which ensures that the discrete occlusion labels match well the disparity observations. We do so by penalizing occlusion boundaries that are not supported by the data

$$\phi_{ij}^{\text{occ}}(\mathbf{y}_{\text{front}}, \mathbf{y}_{\text{back}}) = \begin{cases} \lambda_{\text{imp}} & \text{if } \exists \mathbf{p} \in B_{ij} : \hat{d}_i(\mathbf{p}, \mathbf{y}_{\text{front}}) < \hat{d}_j(\mathbf{p}, \mathbf{y}_{\text{back}}) \\ 0 & \text{otherwise} \end{cases}$$

We also define $\phi_{ij}^{\text{neg}}(\mathbf{y}_i)$ to be a function which penalizes negative disparities

$$\phi_{ij}^{\text{neg}}(\mathbf{y}_i) = \begin{cases} \lambda_{\text{imp}} & \text{if } \min_{\mathbf{p} \in B_{ij}} \hat{d}_i(\mathbf{p}, \mathbf{y}_i) < 0 \\ 0 & \text{otherwise} \end{cases}$$

We impose a regularization on the type of occlusion boundary, where we prefer simpler explanations (i.e., coplanar is preferable than hinge which is more desirable than occlusion). We encode this preference by defining $\lambda_{\text{occ}} > \lambda_{\text{hinge}} > 0$. We thus define our computability potential

$$\phi_{ij}^{\text{bdy2}}(o_{ij}, \mathbf{y}_i, \mathbf{y}_j) = \begin{cases} \lambda_{\text{occ}} + \phi_{ij}^{\text{neg}}(\mathbf{y}_i) + \phi_{ij}^{\text{neg}}(\mathbf{y}_j) + \phi_{ij}^{\text{occ}}(\mathbf{y}_i, \mathbf{y}_j) & \text{if } o_{ij} = lo \\ \lambda_{\text{occ}} + \phi_{ij}^{\text{neg}}(\mathbf{y}_i) + \phi_{ij}^{\text{neg}}(\mathbf{y}_j) + \phi_{ij}^{\text{occ}}(\mathbf{y}_j, \mathbf{y}_i) & \text{if } o_{ij} = ro \\ \lambda_{\text{hinge}} + \phi_{ij}^{\text{neg}}(\mathbf{y}_i) + \phi_{ij}^{\text{neg}}(\mathbf{y}_j) + \frac{1}{|B_{ij}|} \sum_{\mathbf{p} \in B_{ij}} \Delta d_{ij} & \text{if } o_{ij} = hi \\ \phi_{ij}^{\text{neg}}(\mathbf{y}_i) + \phi_{ij}^{\text{neg}}(\mathbf{y}_j) + \frac{1}{|S_i \cup S_j|} \sum_{\mathbf{p} \in S_i \cup S_j} \Delta d_{i,j} & \text{if } o_{ij} = co \end{cases}$$

with $\Delta d_{i,j} = (\hat{d}_i(\mathbf{p}, \mathbf{y}_i) - \hat{d}_j(\mathbf{p}, \mathbf{y}_j))^2$.

Table 3. Comparison with the state-of-the-art on Middlebury high-resolution imagery

	> 1 pixel		> 2 pixels		> 3 pixels		> 4 pixels		> 5 pixels	
	N.-Occ	All	N.-Occ	All	N.-Occ	All	N.-Occ	All	N.-Occ	All
GC+occ [36]	23.8 %	-	16.6 %	-	13.9 %	-	12.5 %	-	11.5 %	-
EBP [41]	14.3 %	-	10.3 %	-	9.4 %	-	9.0 %	-	8.7 %	-
GCS [26]	13.2 %	-	9.0 %	-	7.4 %	-	6.5 %	-	5.9 %	-
SDM [27]	12.8 %	-	9.3 %	-	8.2 %	-	7.7 %	-	7.3 %	-
ELAS [30]	7.1 %	17.0 %	4.7 %	11.7 %	3.9 %	9.2 %	3.5 %	7.9 %	3.2 %	7.3 %
OCV-SGBM [25]	7.0 %	14.6 %	5.9 %	12.5 %	5.5 %	11.5 %	5.3 %	10.9 %	5.2 %	10.4 %
Ours	4.4 %	11.2 %	2.8 %	8.1 %	2.4 %	6.9 %	2.1 %	6.3 %	2.0 %	5.8 %

Table 4. Difference in performance when employing different segmentation methods to compute superpixels on the validation set of KITTI. When employing an intersection of SLIC+UCM superpixels works best for the same amount of superpixels.

	Super-pixels	> 2 pixels		> 3 pixels		> 4 pixels		> 5 pixels	
		N.-Occ	All	N.-Occ	All	N.-Occ	All	N.-Occ	All
UCM	95.3	53.38%	53.88%	48.67%	49.14%	45.59%	46.03%	43.32%	43.73%
SLIC	981.0	7.22%	8.56%	4.92%	6.05%	3.87%	4.84%	3.27%	4.11%
SLIC	1218.2	7.15%	8.63%	4.86%	6.15%	3.81%	4.93%	3.21%	4.20%
UCM+SLIC	1203.2	7.04%	8.42%	4.77%	5.94%	3.74%	4.73%	3.16%	4.01%

Junction Feasibility: Following work on occlusion boundary reasoning [15, 16], we utilize higher order potentials to encode whether a junction of three planes is possible. We refer the reader to Fig. 1 for an illustration of these cases. We thus define the compatibility of a junction $\{i, j, k\}$ to be

$$\phi_{ijk}^{jct}(o_{ij}, o_{jk}, o_{ik}) = \begin{cases} \lambda_{imp} & \text{if impossible case} \\ 0 & \text{otherwise} \end{cases}$$

We also defined a potential encoding the feasibility of a junction of four planes (see Fig. 2) as follows

$$\phi_{pqrs}^{crs}(o_{pq}, o_{qr}, o_{rs}, o_{ps}) = \begin{cases} \lambda_{imp} & \text{if impossible case} \\ 0 & \text{otherwise} \end{cases}$$

Note that, although these potentials are high order, they only involve variables with small number of states, i.e., 4 states.

Potential for Color Similarity: Finally, we employ a simple color potential to reason about segmentation, which is defined in terms of the χ -squared distance between color histograms of neighboring segments. This potential encodes the fact that we expect segments which are coplanar to have similar color statistics (i.e., histograms), while the entropy of this distribution is higher when the planes form an occlusion boundary or a hinge. This trend is shown in Fig. 3 (left) for KITTI [2]¹. We reflect these statistics in the following potential

$$\phi_{ij}^{col}(o_{ij}) = \begin{cases} \min(\kappa \cdot \chi^2(h_i, h_j), \lambda_{col}) & \text{if } o_{ij} = co \\ \lambda_{col} & \text{otherwise} \end{cases}$$

¹ The statistics are less meaningful in the case of the Middlebury high resolution imagery [1], as this dataset is captured in a control environment.

with κ a scalar and $\chi^2(h_i, h_j)$ the χ -squared distance between the color histograms of segments i and j .

3.2 Inference in Continuous MRFs

Now that we have defined the model, we can turn our attention to inference, which is defined as computing the MAP estimate, i.e., $\arg \max_{\mathbf{y}, \mathbf{o}} p(\mathbf{y}, \mathbf{o})$, with $p(\mathbf{y}, \mathbf{o})$ defined in Eq. 2. Inference in this model is in general NP hard. Our inference is also particularly challenging since, unlike traditional MRF stereo formulations, we have defined a hybrid MRF, which reasons about continuous as well as discrete variables. While there is a vast literature on discrete MRF inference, only a few attempts have focussed on solving the continuous case. The exact MAP solution can only be recovered in very restrictive cases, e.g., when the potentials are quadratic and diagonally dominated, Gaussian Belief propagation [42] returns the optimal solution. For general potentials, one can approximate the messages using mixture models, or via particles.

Here we make use of particle convex belief propagation (PCBP) [19], a technique that is guaranteed to converge and gradually approach the optimum. This works very well in practice, yielding state-of-the-art results. PCBP is an iterative algorithm that works as follows: For each random variable, particles are sampled around the current solution. These samples act as labels in a discretized MRF which is solved to convergence using convex belief propagation [32]. The current solution is then updated with the MAP estimate obtained on the discretized MRF. This process is repeated for a fixed number of iterations. In our implementation, we use the distributed message passing algorithm of [43] to solve the discretized MRF at each iteration. Algorithm 1 depicts PCBP for our formulation. At each iteration, to balance the trade off between exploration and exploitation, we decrease the values of the standard deviations $\sigma_\alpha, \sigma_\beta$ and σ_γ of the normal distributions from which the plane random variables are drawn.

Algorithm 1. PCBP for stereo estimation and occlusion boundary reasoning

```

Set  $N$ 
Initialize slanted planes  $\mathbf{y}_i^0 = (\alpha_i^0, \beta_i^0, \gamma_i^0)$  via local fitting  $\forall i$ 
Initialize  $\sigma_\alpha, \sigma_\beta$  and  $\sigma_\gamma$ 
for  $t = 1$  to #iters do
  Sample  $N$  times  $\forall i$  from  $\alpha_i \sim \mathcal{N}(\alpha_i^{t-1}, \sigma_\alpha)$ ,  $\beta_i \sim \mathcal{N}(\beta_i^{t-1}, \sigma_\beta)$ ,  $\gamma_i \sim \mathcal{N}(\gamma_i^{t-1}, \sigma_\gamma)$ 
   $(\mathbf{o}^t, \mathbf{y}^t) \leftarrow$  Solve the discretized MRF using convex BP
  Update  $\sigma_\alpha^c = \sigma_\beta^c = 0.5 \times \exp(-c/10)$  and  $\sigma_\gamma^c = 5.0 \times \exp(-c/10)$ 
end for
Return  $\mathbf{o}^t, \mathbf{y}^t$ 

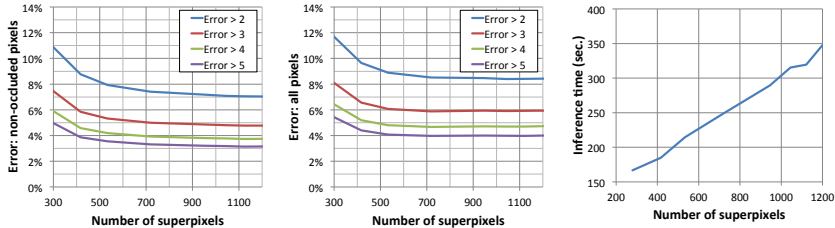
```

3.3 Learning in Continuous MRFs

Given a set of training images and corresponding depth labels, the goal of learning is to estimate the weights which minimize the surrogate loss (e.g., hinge loss

Table 5. Performance changes when employing different segmentation methods to compute superpixels on the Middlebury high resolution imagery

	Super-pixels	> 1 pixel		> 2 pixels		> 3 pixels		> 4 pixels		> 5 pixels	
		N.-Occ	All	N.-Occ	All	N.-Occ	All	N.-Occ	All	N.-Occ	All
UCM	259.0	38.5%	41.7%	30.2%	33.5%	25.7%	28.8%	22.8%	25.6%	20.6%	23.1%
SLIC	1787.6	4.8%	11.8%	3.1%	8.7%	2.7%	7.4%	2.4%	6.8%	2.3%	6.4%
SLIC	2066.1	4.6%	12.0%	3.0%	8.9%	2.6%	7.7%	2.4%	7.0%	2.2%	6.4%
UCM+SLIC	2042.6	4.4%	11.2%	2.8%	8.1%	2.4%	6.9%	2.1%	6.3%	2.0%	5.8%


Fig. 5. Number of superpixels: KITTI validation results as a function of the number of superpixels. Even with a small number our approach still outperforms the baselines. (Right) Inference time scales linearly with the number of superpixels.

for structured SVMs or log-loss for CRFs). In our model, we have a total of 5 weights, associated with ϕ^{seg} , ϕ^{bdy1} , ϕ^{bdy2} , ϕ^{col} , as well as a shared weight for ϕ^{ict} and ϕ^{crs} . We would like to employ the algorithm of [20] for learning. However, our learning problem, as opposed to the one in [20], contains a mixture of continuous and discrete variables. Therefore the surrogate loss in our setting requires to integrate over the continuous variables. We note that our continuous variables have robust quadratic potentials, thus integrating over them can be efficiently approximated by discretizing the continuous variables. In practice, summing over 30 particles gives a good approximation for the integral.

4 Experimental Evaluation

We perform exhaustive experiments on two publicly available datasets: Middlebury high resolution images [1] as well as the more challenging KITTI dataset [2]. The high resolution Middlebury images [1] have an average resolution of 1239.2×1038.0 pixels. We employ 5 images for training (i.e., *Books*, *Laundry*, *Moebius*, *Reindeer*, *Bowling2*) and 9 images for testing (i.e., *Cones*, *Teddy*, *Art*, *Aloe*, *Dolls*, *Baby3*, *Cloth3*, *Lampshade2*, *Rocks2*). We also evaluate our approach on the KITTI dataset [2], which is the only real-world stereo dataset with accurate ground truth. It is composed of 194 training and 195 test high-resolution images (1237.1×374.1 pixels) captured from an autonomous driving platform driving around in a urban environment. The ground truth is generated by means of a Velodyne sensor which is calibrated with the stereo pair. This results in semi-dense ground truth covering

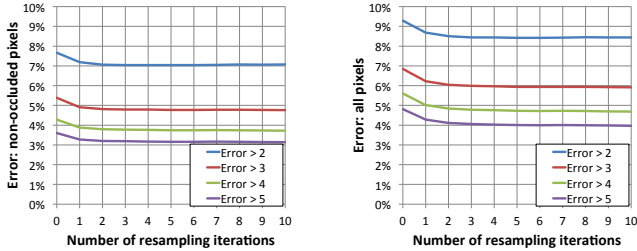


Fig. 6. Importance of re-sampling iterations: on KITTI validation set

Table 6. Training set size: Estimation errors as a function of the training set size on the validation set of KITTI. Very few images are needed to learn good parameters.

Number of training images	> 2 pixels		> 3 pixels		> 4 pixels		> 5 pixels	
	Non-Occ	All	Non-Occ	All	Non-Occ	All	Non-Occ	All
1	7.06 %	8.45 %	4.80 %	5.98 %	3.77 %	4.76 %	3.19 %	4.04 %
5	7.05 %	8.44 %	4.79 %	5.96 %	3.75 %	4.74 %	3.16 %	4.01 %
10	7.04 %	8.42 %	4.77 %	5.94 %	3.74 %	4.73 %	3.16 %	4.01 %
20	7.04 %	8.42 %	4.78 %	5.94 %	3.75 %	4.72 %	3.17 %	4.00 %

approximately 30 % of the pixels. We employ 20 images for training, and utilize the remaining 174 images for validation purposes.

For all experiments, we employ the same parameters which have been validated on the training set. We use a disparity difference threshold $K = 5.0$ pixels, and set $\lambda_{\text{occ}} = 15$, $\lambda_{\text{hinge}} = 3$, $\lambda_{\text{imp}} = 30$ and $\lambda_{\text{col}} = 30$. For the color potential, we use a color histogram with 64 bins and set $\kappa = 60$. Unless otherwise stated, we use 10 training images learning, 10 particles and 5 iterations of re-sampling for PCBP [19], and run each iteration of convex BP to convergence. For learning, we use a value of C equal to the number of examples and unless otherwise stated use a CRF, i.e., $\epsilon = 1$. We learned the importance of each potential, thus 6 parameters. We employ two different metrics. The first one measures the average number of non-occluded pixels which error is bigger than a fixed threshold. To test the extrapolation capabilities of the different approaches, the second metric computes the same metric, but including the occluding pixels as well.

Robust SGBM: We begin our experimentation by developing a new criteria for semi-global block matching which is more robust and accurate. It uses gradients as well as Census transform [44]. Left-right consistency check is performed by computing both left and right disparity maps. Table 1 shows the performance improvement. Matching is performed on average in only 3.6s for each KITTI image. We utilize this more robust matching criteria to create our potentials.

Comparison with the State-of-the-Art: Table 2 and 3 depict results of our approach and the baselines in terms of the two metrics for the KITTI and high resolutions Middlebury datasets respectively. Note that our approach significantly outperforms all the baselines in all settings (i.e., thresholds bigger than 2, 3, 4 and 5 pixels). Fig. 4 depicts an illustrative set of KITTI examples. Despite the challenges, our approach does a good job at estimating disparities.

Table 7. Oracle performance: Oracle, our approach and initial fit on KITTI val

	> 2 pixels		> 3 pixels		> 4 pixels		> 5 pixels	
	N.-Occ	All	N.-Occ	All	N.-Occ	All	N.-Occ	All
Oracle	1.38%	1.70%	1.03%	1.27%	0.90%	1.10%	0.82%	0.99%
Initial fit	7.66%	9.28%	5.38%	6.85%	4.28%	5.61%	3.60%	4.81%
Ours	7.04%	8.42%	4.77%	5.94%	3.74%	4.73%	3.16%	4.01%

Table 8. Oracle performance: Oracle, our approach and initial fit on Middlebury

	> 1 pixel		> 2 pixels		> 3 pixels		> 4 pixels		> 5 pixels	
	N.-Occ	All	N.-Occ	All	N.-Occ	All	N.-Occ	All	N.-Occ	All
Oracle	2.0%	6.2%	1.4%	5.5%	1.3%	5.3%	1.2%	5.2%	1.1%	5.1%
Initial fit	4.9%	13.4%	3.2%	10.5%	2.7%	9.4%	2.5%	8.7%	2.3%	8.2%
Ours	4.4%	11.2%	2.8%	8.1%	2.4%	6.9%	2.1%	6.3%	2.0%	5.8%

Segmentation Strategy: We next investigate how the segmentation strategy affects the stereo estimation. Towards this goal we evaluate the results of our approach when employing UCM segments [45], SLIC superpixels [46] or the intersection of both as input. Table 4 depicts results on the KITTI dataset. UCM performs very poorly as the number of superpixels on average is very small, i.e., a single 3D plane is a poor representation for the disparities in the large segments. SLIC performs quite well, but the intersection of SLIC and UCM superpixels outperforms the other strategies. This is also expected, as UCM respects the boundaries much better than SLIC. Note that as shown in Table 5 similar results are observed for the Middlebury dataset.

Number of Superpixels: We next investigate how well our approach scales with the number of superpixels in terms of computational complexity as well as accuracy. Fig. 5 shows results for the KITTI dataset when varying the number of superpixels. Our approach reduces performance gracefully when reducing the amount of superpixels. Note that inference scales linearly with the number of superpixels, taking on average 5.5 minutes per high resolution image when employing 1200 superpixels and 2.5 minutes when using 300.

Number of Re-sampling Iterations: We evaluate the effects of varying the number of resampling iterations on the performance of our approach. As shown in Fig. 6, our approach converges to a good local optima after only 2 resampling iterations. This reduces the inference cost from 5.5 minutes per high-resolution image for 5 iterations to 2.2 minutes for 2 iterations.

Training Set Size: We evaluate the effect of increasing the training set size in Table 6. Even when training with a single image we outperform all baselines.

Oracle Performance: We evaluate the best performance that our model can achieve, by fitting the model to the ground truth disparities. This is an upper-bound on the performance that our method could ever achieve if we were able to learn an energy that has its MAP at the ground truth, and if we were able to solve the NP-hard inference problem. Tables 7 and 8 depict the oracle performance in terms of both the occluded and non-occluded pixels for both datasets. Note that as KITTI does not release the test ground truth, we compute this values using

Table 9. Importance of potential functions: on the validation set of KITTI

	> 2 pixels		> 3 pixels		> 4 pixels		> 5 pixels	
	Non-Occ	All	Non-Occ	All	Non-Occ	All	Non-Occ	All
SGBM	8.19 %	9.39 %	5.55 %	6.52 %	4.26 %	5.05 %	3.47 %	4.13 %
Initial fit	7.66 %	9.28 %	5.38 %	6.85 %	4.28 %	5.61 %	3.60 %	4.81 %
MRF plain	7.23 %	8.62 %	5.03 %	6.22 %	3.99 %	5.01 %	3.37 %	4.27 %
MRF +color	7.08 %	8.48 %	4.88 %	6.06 %	3.86 %	4.87 %	3.27 %	4.14 %
MRF +color+junction	7.04 %	8.42 %	4.77 %	5.94 %	3.74 %	4.73 %	3.16 %	4.01 %

Table 10. Robustness to noise: RMS and boundary error as a function of noise

Noise	0	1	2	3	5
RMS (pixels)	0.44	0.80	1.37	2.24	4.40
Boundary error	0.3%	0.6%	1.9%	5.3%	8.9%

10 images for training and the rest of the training set for testing. We also report performance of our initialization which is computed by fitting a local plane to the results of our robust semi-global block matching. Note that the oracle can achieve great performance, showing that the errors due to the 3D slanted plane discretization are negligible.

Importance of Potentials: We evaluate the importance of each potential that our model employs in Table 9. Note that for error > 3 pixels, the contribution of junction potential is 18% of the gain from the initial fit.

Robustness to Noise: We investigate the robustness of our approach to noise by building a synthetic dataset, which is composed of 10 images for training and 90 images for test of resolution 320×240 . The average number of superpixels is 108.0. We create $\mathcal{D}(\mathbf{p})$ by sampling 3 to 5 points at random on the boundaries and generating disparities by corrupting the ground truth with Gaussian noise of varying standard deviation. Table 10 shows RMS errors for disparity as well as percentage of boundary variables wrongly estimated.

5 Conclusion

We have presented a novel stereo slanted-plane MRF model that reasons jointly about occlusion boundaries as well as depth. We have formulated the problem as inference in a hybrid MRF composed of both continuous (i.e., slanted 3D planes) and discrete (i.e., occlusion boundaries) random variables, which we have tackled using particle convex belief propagation. We have demonstrated the effectiveness of our approach on high resolution imagery from Middlebury as well as the more challenging KITTI dataset. In the future we plan to investigate alternative inference algorithms as well as other segmentation potentials.

References

1. Scharstein, D., Szeliski, R.: A taxonomy and evaluation of dense two-frame stereo correspondence algorithms. *IJCV* 47(1,2,3) (2002); Microsoft Research Technical Report MSR-TR-2001-81 (November 2001)

2. Geiger, A., Lenz, P., Urtasun, R.: Are we ready for autonomous driving? In: CVPR (2012)
3. Tappen, M., Freeman, W.: Comparison of graph cuts with belief propagation for stereo, using identical mrf parameters. In: ICCV (2003)
4. Birchfield, S., Tomasi, C.: Multiway cut for stereo and motion with slanted surfaces. In: ICCV (1999)
5. Hong, L., Chen, G.: Segment-based stereo matching using graph cuts. In: CVPR (2004)
6. Bleyer, M., Gelautz, M.: A layered stereo matching algorithm using image segmentation and global visibility constraints. *ISPRS Journal of Photogrammetry and Remote Sensing* 59(3), 128–150 (2005)
7. Klaus, A., Sormann, M., Karner, K.: Segment-based stereo matching using belief propagation and a self-adapting dissimilarity measure. In: ICPR 2006 (2006)
8. Deng, Y., Yang, Q., Lin, X., Tang, X.: A symmetric patch-based correspondence model for occlusion handling. In: ICCV (2005)
9. Yang, Q., Engels, C., Akbarzadeh, A.: Near real-time stereo for weakly-textured scenes. In: British Machine Vision Conference (2008)
10. Trinh, H., McAllester, D.: Unsupervised learning of stereo vision with monocular cues. In: Proc. BMVC, Citeseer (2009)
11. Zitnick, C., Kanade, T.: A cooperative algorithm for stereo matching and occlusion detection. *PAMI* 22(7) (July 2000)
12. Kolmogorov, V., Zabih, R.: Multi-camera Scene Reconstruction via Graph Cuts. In: Heyden, A., Sparr, G., Nielsen, M., Johansen, P. (eds.) *ECCV 2002, Part III*. LNCS, vol. 2352, pp. 82–96. Springer, Heidelberg (2002)
13. Bleyer, M., Rother, C., Kohli, P.: Surface stereo with soft segmentation. In: CVPR (2010)
14. Waltz, D.L.: Generating semantic description from drawings of scenes with shadows. MIT Artificial Intelligence Laboratory (1972) Workin Paper 47
15. Malik, J.: Interpreting line drawings of curved objects. *International Journal of Computer Vision* 1, 73–103 (1987)
16. Hoiem, D., Stein, A.N., Efros, A.A., Hebert, M.: Recovering occlusion boundaries from a single image. In: ICCV (2007)
17. Ashutosh Saxena, J.S., Ng, A.Y.: Depth estimation using monocular and stereo cues. In: IJCAI (2007)
18. Maire, M., Arbelaez, P., Fowlkes, C., Malik, J.: Using contours to detect and localize junctions in natural images. In: CVPR (2008)
19. Peng, J., Hazan, T., McAllester, D., Urtasun, R.: Convex max-product algorithms for continuous mrfs with applications to protein folding. In: International Conference on Machine Learning, ICML (2011)
20. Hazan, T., Urtasun, R.: A primal-dual message-passing algorithm for approximated large scale structured prediction. *Advances in Neural Information Processing Systems* 23, 838–846 (2010)
21. Scharstein, D., Szeliski, R.: A taxonomy and evaluation of dense two-frame stereo correspondence algorithms. In: *IJCV* (2002)
22. Konolige, K.: Small vision system. hardware and implementation. In: International Symposium on Robotics Research, pp. 111–116 (1997)
23. Kanade, T., Okutomi, M.: A stereo matching algorithm with an adaptive window: Theory and experiment. In: *ICRA* (1994)
24. Yoon, K.J., Member, S., Kweon, I.S.: Adaptive support-weight approach for correspondence search. *PAMI* 28, 650–656 (2006)

25. Hirschmueller, H.: Stereo processing by semiglobal matching and mutual information. *PAMI* 30, 328–341 (2008)
26. Cech, J., Sara, R.: Efficient sampling of disparity space for fast and accurate matching. In: *BenCOS* (2007)
27. Kostkova, J.: Stratified dense matching for stereopsis in complex scenes. In: *BMVC* (2003)
28. Wang, L., Jin, H., Yang, R.: Search Space Reduction for MRF Stereo. In: Forsyth, D., Torr, P., Zisserman, A. (eds.) *ECCV 2008, Part I. LNCS*, vol. 5302, pp. 576–588. Springer, Heidelberg (2008)
29. Veksler, O.: Reducing search space for stereo correspondence with graph cuts. In: *BMVC* (2006)
30. Geiger, A., Roser, M., Urtasun, R.: Efficient Large-Scale Stereo Matching. In: Kimmel, R., Klette, R., Sugimoto, A. (eds.) *ACCV 2010, Part I. LNCS*, vol. 6492, pp. 25–38. Springer, Heidelberg (2011)
31. Sun, J., Zheng, N.N., Shum, H.Y.: Stereo matching using belief propagation. *PAMI* 25(7) (2003)
32. Hazan, T., Shashua, A.: Norm-product belief propagation: Primal-dual message-passing for approximate inference. *IEEE Transactions on Information Theory* 56(12), 6294–6316 (2010)
33. Koller, D., Lerner, U., Angelov, D.: A general algorithm for approximate inference and its application to hybrid bayes nets. In: *UAI* (1999)
34. Ihler, A., McAllester, D.: Particle belief propagation. In: *AISTATS 2009* (2009)
35. Hazan, T., Urtasun, R.: A primal-dual message-passing algorithm for approximated large scale structured prediction. In: *NIPS* (2010)
36. Kolmogorov, V., Zabih, R.: Computing visual correspondence with occlusions using graph cuts. In: *ICCV*, pp. 508–515 (2001)
37. Bradski, G.: The opencv library. *Dr. Dobb's Journal of Software Tools* (2000)
38. Rhemann, C., Hosni, A., Bleyer, M., Rother, C., Gelautz, M.: Fast cost-volume filtering for visual correspondence and beyond. In: *CVPR* (2011)
39. Cech, J., Sanchez-Riera, J., Horaud, R.P.: Scene flow estimation by growing correspondence seeds. In: *CVPR* (2011)
40. Ranftl, R., Gehrig, S., Pock, T., Bischof, H.: Pushing the Limits of Stereo Using Variational Stereo Estimation. In: *IEEE Intelligent Vehicles Symposium* (2012)
41. Felzenszwalb, P.F., Huttenlocher, D.P.: Efficient belief propagation for early vision. *IJCV* 70(1) (October 2006)
42. Weiss, Y., Freeman, W.: Correctness of belief propagation in gaussian graphical models of arbitrary topology. *Neural Computation* 13(10) (2001)
43. Schwing, A., Hazan, T., Pollefeys, M., Urtasun, R.: Distributed message passing for large scale graphical models. In: *IEEE Conference on Computer Vision and Pattern Recognition, CVPR* (2011)
44. Zabih, R., Woodfill, J.: Non-parametric local transforms for computing visual correspondence. In: *ECCV* (1994)
45. Arbelaez, P., Maire, M., Fowlkes, C., Malik, J.: Contour detection and hierarchical image segmentation. In: *PAMI* (2011)
46. Achanta, R., Shaji, A., Smith, K., Lucchi, A., Fua, P., Süsstrunk, S.: Slic superpixels. Technical Report 149300 EPFL (June 2010)