

# Local Refinement for Stereo Regularization

Carl Olsson, Johannes Ulén  
Centre for Mathematical Sciences  
Lund University  
`{calle,ulen}@maths.lth.se`

Anders Eriksson  
School of Computer Science  
University of Adelaide  
`anders.eriksson@adelaide.edu.au`

**Abstract**—Stereo matching is an inherently difficult problem due to ambiguous and noisy texture. The non-convexity and non-differentiability makes local linear (or quadratic) approximations poor, thereby preventing the use of standard local descent methods. Therefore recent methods are predominantly based on discretization and/or random sampling of some class of approximating surfaces (e.g. planes). While these methods are very efficient in generating a rough surface estimate, via either fusion of proposals or label propagation, the end result is usually not as smooth as desired. In this paper we show that, if the objective function is decomposed correctly, local refinement of candidate solutions can be performed using an ADMM approach. This allows searching over more general function classes, thereby resulting in visually more appealing smooth surface estimations.

## I. INTRODUCTION

In the last decade considerable progress has been made in dense stereo matching, largely due to the availability of powerful regularizers for handling ambiguous and noisy data. The most common are the first order regularization priors [1], [2], [3]. One reason for their popularity is that robust move-making algorithms such as  $\alpha$ -expansion [1] or fusion moves [4] capable of modifying large numbers of pixels simultaneously can be applied. Such moves are essential for avoiding poor local solutions.

First order methods often implicitly assume fronto-parallel planes. For example, standard piecewise smooth (e.g. truncated linear or quadratic) pairwise regularization potentials assign higher cost to surfaces with larger tilt with respect to the camera [1]. To model surfaces more accurately Birchfeld and Tomasi [5] introduced 3D-labels corresponding to arbitrary 3D planes. However, this approach is limited to piecewise planar scenes. To address more general scenes recent papers use 2nd derivative regularization [6], [7], [8]. There are two ways of modeling such higher order smoothness potentials. Woodford et al. [8] retain the scalar disparity labels while using triple-cliques to penalize 2nd derivatives of the reconstructed surface. This encourages near planar smooth disparity maps. The optimization problem is however made substantially more difficult due to the introduction of non-submodular triple interactions. In contrast, [6], [7] use 3D-labels corresponding to tangent planes to encode 2nd order smoothness as pairwise interactions. It is shown in [6] that in contrast to the triple-cliques used by Woodford et al. [8] the 3D-label formulation is often submodular (or near submodular) making fusion moves easier to solve optimally using standard methods like Roof duality [9].

While these methods are very efficient in generating a rough surface estimate the end result is usually not as smooth

as desired. Since proposals are often planar or piecewise planar [8] approximations of the surface, a very large number of tangents have to be generated in order to achieve a smooth result. In this paper we focus on local methods for minimizing the 2nd order stereo energy. Such approaches can be used to generate proposals (in the form of locally optimal surfaces) which can then be fused with the current solution.

The energies considered have three features that make local optimization difficult; First, the parametrization of the pairwise interaction is non-linear resulting in non-convex least squares terms. Second, the interaction is truncated to preserve discontinuities in the scene. This makes local optimization difficult since the Taylor approximation can be an arbitrarily bad approximation of the function (even locally). And third, the data term is based on photo consistency of local patches, making it highly non-convex with lots of local minima, see Figure 1. In this paper we show that with the right parametrization all of these difficulties can be addressed within an ADMM framework [10]. We show that the use of local optimization in the proposal generation process results in smoother and more visually appealing surface estimates.

## II. STEREO AND FUSION MOVES

### A. Energy

The goal of stereo reconstruction is to compute a depth estimate for every pixel in an image. Doing so requires that every pixel in the image is matched to a corresponding pixel in another image. Due to ambiguous texture this matching is rarely unique and as a result the stereo problem is most often ill posed. Resolving these ambiguities requires adding knowledge of the types of surfaces that we can expect to see in natural scenes, in the form of regularization.

The problem is therefore typically formulated as an energy minimization problem of the form

$$\min_z E(z) \quad (1)$$

where

$$E(z) = \sum_{\mathbf{p}} \sum_{\mathbf{q} \in \mathcal{N}(\mathbf{p})} E_{\mathbf{pq}}(z) + \sum_{\mathbf{p}} E_{\mathbf{p}}(z). \quad (2)$$

The function  $z$  represents the sought depth map and  $\mathcal{N}(\mathbf{p})$  is a predefined neighborhood to  $\mathbf{p}$ . Throughout the paper we will think of our assignments as samples of the underlying depth function. The value  $z(\mathbf{p})$  represents the assignment at pixel  $\mathbf{p}$ . The term  $E_{\mathbf{p}}$  is a data term that only depends on the assignment at the particular pixel  $\mathbf{p}$ . Typically this term

is based on some measure of photo consistency. For a given assignment  $z(\mathbf{p})$  its value can be computed by backprojecting into another image and comparing pixel appearance. In our work we will use normalized cross correlation between a  $3 \times 3$  patch with center at  $\mathbf{p}$  in the original image and a patch with center at the backprojection point.

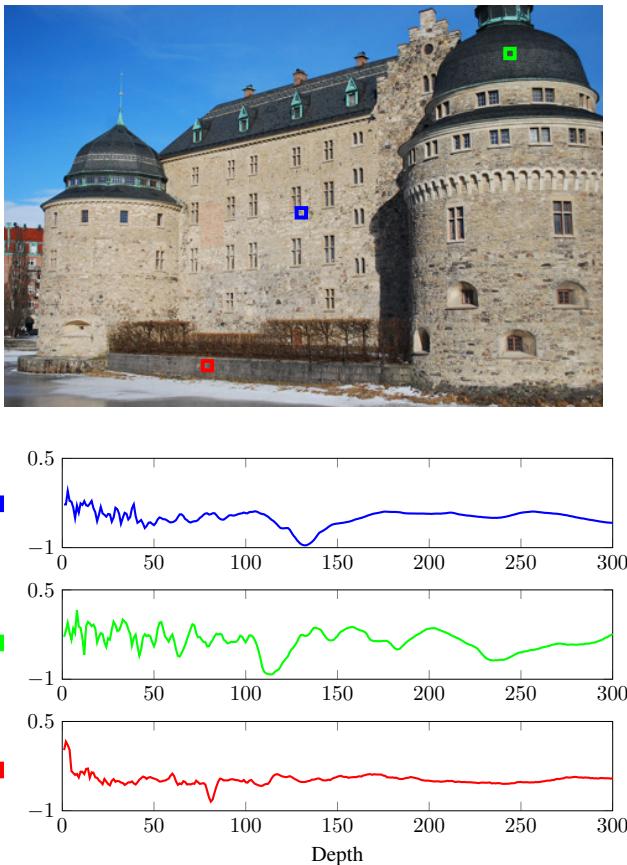


Fig. 1: The dataterm  $E_p$  for three pixels in the Örebro Castle dataset.

Photo consistency based measurements such as this one generally results in noisy and ambiguous functions. Figure 1 shows  $E_p$  for three typical pixels in the Örebro Castle dataset. To handle this the smoothness term  $E_{pq}$  is added to the energy. (By the subscript  $pq$  we mean that it only depends on the assignments at  $\mathbf{p}$  and  $\mathbf{q}$ .) This term is meant to penalize unrealistic assignments such as surfaces with high depth variation. The simplest possible interaction is the so called Potts model, which adds a penalty if neighboring depth assignments are different [1]. In this work we follow [6] which uses a slightly more complicated interaction that penalizes second order smoothness, see Figure 2. In addition to assigning a depth  $z(\mathbf{p})$  we also assign a depth gradient  $\nabla z(\mathbf{p})$  to pixel  $\mathbf{p}$ , which allows us to determine the tangent of  $z$  at  $\mathbf{p}$ . To enforce smoothness we extend the tangent at  $\mathbf{p}$  to  $\mathbf{q}$  and measure the difference  $V_{pq}$  to the depth assignment at  $\mathbf{q}$ . Intuitively, if the function is smooth then the tangent should be a good approximation, and therefore  $V_{pq}$  should be small. Note in particular that planar depth functions do not incur any penalty.

## B. Fusion Moves

Optimizing energies, such as (1), is difficult since  $E_p$  is often non-convex and even non-differentiable. An approach that has been effective are the so called fusion moves [4]. Given two assignments  $z^0$  and  $z^1$  we fuse them into a new one with lower energy by solving

$$\min_{x \in \{0,1\}^n} E(x \cdot z^0 + (1-x) \cdot z^1), \quad (3)$$

where  $\cdot$  is element-wise multiplication. These moves are very effective since they allow changes of a large number of pixels at the same time. In addition they are not local and can therefore escape bad local minimas. In the case of pairwise energies, such as (1), the fusion move can be solved using standard methods [9], [11].

The candidate solutions are typically selected by randomly sampling functions from some low dimensional class, for example [8], [6] use planes and piecewise planes. The motivation for using such a subclass of functions is that sampling can be made very efficient and the resulting solution typically approximates the optimal solution well. However to achieve truly smooth estimations we need to go beyond a fixed class of functions and generate more general surfaces. In this work we show that for the stereo energies that we consider this can be done efficiently using local optimization. Our strategy will be to mix the usage of random sampling of planes with local refinement to generate proposals which are then fused with the current solution.

## III. PARAMETRIZATION OF THE INTERACTION

In order to do local refinement we need to find a parametrization of the distance  $V_{pq}$ . In [6] the tangent plane at  $\mathbf{p}$  is parametrized using a unit normal  $\mathbf{n}_p$  and a scalar  $d_p$  via the affine plane equation  $\mathbf{n}_p^T x - d_p = 0$ . Such a parametrization does however introduce a non-linear constraint ( $\|\mathbf{n}_p\|^2 = 1$ ). In addition this parametrization introduces a singularity. If the tangent plane is selected such that it contains the entire viewing ray, then no unique depth can be determined for that pixel.

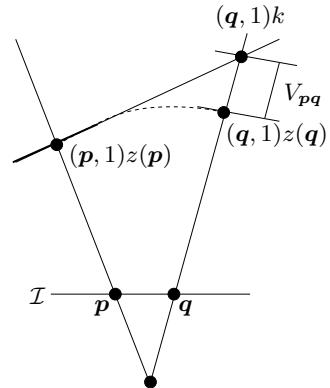


Fig. 2: The regularization interaction  $V_{pq}$  measures the deviation from the neighboring tangent along the viewing ray. Note that interactions are modeled in depth space, working in disparity space would require a modified model.

In the following, we instead derive an expression for the distance  $V_{pq}$  in terms of the assigned depths  $z$  and gradients

$\nabla z$  at the points  $p$  and  $q$ . The point  $(q, 1)k$  is the intersection between the tangent plane and the viewing ray at  $q$ . The line between  $(p, 1)z(p)$  and  $(q, 1)k$  is contained in the tangent plane and can therefore be found by linearizing the curve  $(p + tv, 1)z(p + tv)$ , where  $v$  is a unit vector such that  $q = p + sv$  for some  $s$ . We get

$$l(t) = \begin{pmatrix} p \\ 1 \end{pmatrix} z(p) + t \left( \begin{pmatrix} p \\ 1 \end{pmatrix} z'_v(p) + \begin{pmatrix} v \\ 0 \end{pmatrix} z(p) \right). \quad (4)$$

At the intersection point  $(q, 1)k$ , we have

$$l(t) = k \begin{pmatrix} q \\ 1 \end{pmatrix} = k \left( \begin{pmatrix} p \\ 1 \end{pmatrix} + s \begin{pmatrix} v \\ 0 \end{pmatrix} \right). \quad (5)$$

Identification of the coefficients yields

$$k = z(p) + z'_v(p)t, \quad (6)$$

$$s = \frac{z(p)t}{z(p) + z'_v(p)t} \Leftrightarrow t = \frac{sz(p)}{z(p) - sz'_v(p)}. \quad (7)$$

Therefore we get the residuals

$$k - z(q) = \frac{z(p)^2}{z(p) - sz'_v(p)} - z(q). \quad (8)$$

#### IV. ENERGY AND OPTIMIZATION

Next we formulate the optimization problem. We denote the regularization terms

$$V_{pq}(z) = \frac{z(p)^2}{z(p) - sz'_v(p)} - z(q). \quad (9)$$

The energy consists of a regularization term and a data term

$$\sum_p \sum_{q \in \mathcal{N}(p)} h(V_{pq}(z)) + \lambda \sum_p D_p(z). \quad (10)$$

Here  $D_p$  is a data term that evaluates the cost of the assigned depth for pixel  $p$ . We will assume that this function is densely sampled such that its global minima can be found by simply searching all the sample points. The one dimensional function  $h$  is typically of the form  $h(x) = \min(|x|^p, \tau)$  where  $\tau$  is some threshold level and  $p \in \{1, 2\}$ . The only assumption we make is that it is piecewise differentiable so that the minimum can be found by computing stationary points.

Optimizing energy (10) is typically very challenging since the data term is often non-differentiable with lots of local minima. In addition, the smoothness term is a sum of non-convex functions. To handle these problems we decouple the terms by introducing two new sets of variables;  $x_{pq}$  and  $y_p$ .

We constrain these variables to be  $x_{pq} = V_{pq}(z)$  and  $y_p = z(p)$ . The Augmented Lagrangian [10] is now

$$\begin{aligned} L(x, y, z, \lambda) &= \sum_p \sum_{q \in \mathcal{N}(p)} h(x_{pq}) \\ &+ \sum_p \sum_{q \in \mathcal{N}(p)} \lambda_{pq}(x_{pq} - V_{pq}(z)) \\ &+ \sigma \sum_p \sum_{q \in \mathcal{N}(p)} (x_{pq} - V_{pq}(z))^2 \\ &+ \sum_p (\lambda_p(y_p - z(p)) + \sigma(y_p - z(p))^2) \\ &+ \sum_p D_p(y_p). \end{aligned} \quad (11)$$

When applying ADMM we get the subproblems

$$\min_{x_{pq}} h(x_{pq}) + \lambda_{pq}(x_{pq} - V_{pq}(z)) + \sigma(x_{pq} - V_{pq}(z))^2 \quad (12)$$

$$\begin{aligned} \min_z \sum_p \left( \sum_{q \in \mathcal{N}(p)} \lambda_{pq}(x_{pq} - V_{pq}(z)) + \sigma(x_{pq} - V_{pq}(z))^2 \right. \\ \left. + (\lambda_p(y_p - z(p)) + \sigma(y_p - z(p))^2) \right), \end{aligned} \quad (13)$$

$$\min_{y_p} \lambda_p(y_p - z(p)) + \sigma(y_p - z(p))^2 + D_p(y_p) \quad (14)$$

In addition we obtain the dual update rules (see [10])

$$\lambda_{pq}^{k+1} = \lambda_{pq}^k + \sigma(x_{pq} - V_{pq}(z)), \quad (15)$$

$$\lambda_p^{k+1} = \lambda_p^k + \sigma(y_p - z(p)). \quad (16)$$

This decoupling of terms has the following positive effects: The terms  $h(x_{pq})$  and  $D_p(y_p)$ , that are non-smooth and difficult to approximate locally, end up in two different subproblems both of which are separable and where optimization can be carried out over individual pixels separately, greatly reducing the search space. The coupling between pixels appears in problem (13) where the involved functions are smooth and can be optimized locally using standard descent methods. In the following subsections we outline our solution strategies for the individual problems.

#### A. Problem (12)

To solve (12) we first note that the optimum must be in either a stationary point or in a transition between differentiable segments of the function  $h$ . Since we will be using  $h(x_{pq}) = \min(|x_{pq}|, \tau)$  in the experiments we illustrate the process using this choice. We get four cases:

- 1)  $|x_{pq}| > \tau$ . Taking derivatives of (12) gives

$$\lambda_{pq} + 2\sigma(x_{pq} - V_{pq}(z)) = 0. \quad (17)$$

Solving for  $x_{pq}$  gives the stationary point. Note that the solution may violate  $|x_{pq}| > \tau$ . In this case the solution is false and there is no stationary point in the interval. However since we compare the energies of all "candidate" minimizers we do not have to test for this. We are guaranteed that one of the candidates is the global minimizer of (12).

- 2)  $-\tau < x_{pq} < 0$ . In this case the stationary point is given by

$$-1 + \lambda_{pq} + 2\sigma(x_{pq} - V_{pq}(z)) = 0. \quad (18)$$

- 3)  $0 < x_{pq} < \tau$ . Here the stationary point is given by

$$1 + \lambda_{pq} + 2\sigma(x_{pq} - V_{pq}(z)) = 0. \quad (19)$$

- 4) In addition we need to test the two transition points  $x_{pq} = \pm\tau$  and  $x_{pq} = 0$ .

#### B. Problem (13)

The objective function in (13) is similar to non-linear least squares problem. We will apply a Levenberg-Marquart approach to solve it. We linearize the residual

$$x_{pq} - \left( \frac{z(p)^2}{z(p) - sz'_v(p)} - z(q) \right). \quad (20)$$

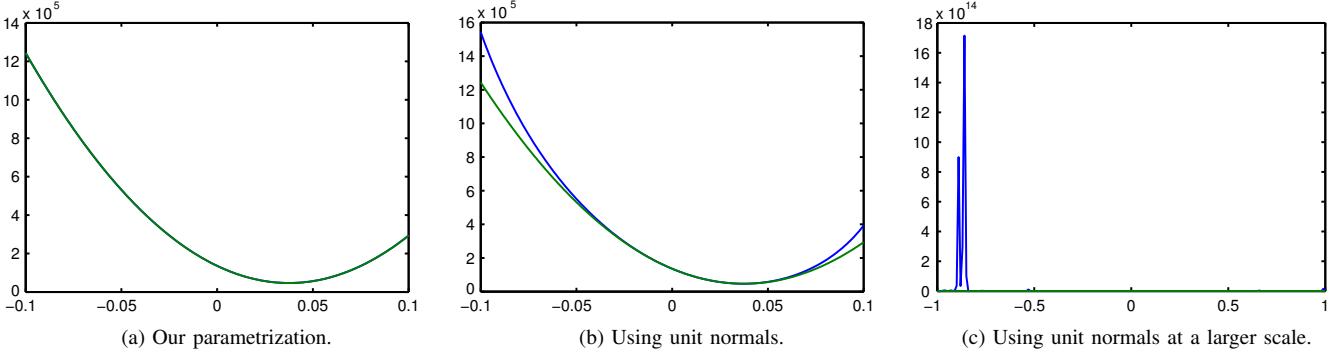


Fig. 3: The effects of parameterizing the problem using the parametrization with unit normals from [6] instead of our parametrization. Green curve shows our approximation and the blue curve the exact error function. For our parametrization the approximation cannot be distinguished from the exact function. Using unit normals at a larger scale gives visible singularities.

Note that  $y_p - z(\mathbf{q})$  is already linear in terms of the unknowns  $(z(\mathbf{p}), z(\mathbf{q}))$  and  $z'_v(\mathbf{p})$ , and does therefore not require modification. The approximation that we make using linearization is in fact in most cases very accurate. This can be heuristically explained by looking at the equivalent expression

$$x_{pq} - \left( z(\mathbf{p}) + sz'_v(\mathbf{p}) + \frac{(sz'_v(\mathbf{p}))^2}{z(\mathbf{p}) - sz'_v(\mathbf{p})} - z(\mathbf{q}) \right). \quad (21)$$

Assuming that  $s$  is small (recall that this is the distance between pixels) the nonlinear term is likely to be neglectable for reasonable values of the derivative  $z'_v(\mathbf{p})$ .

Figure 3 shows an example plot of the exact error function and its least squares approximation, that is, residuals are linearized into  $A\delta - b$  and the non-linear least squares objective is approximated with  $\|A\delta - b\|^2$ . Here  $\delta$  represents the increment in the assignments  $z, \nabla z$ . We plot the functions in the direction  $-2A^T b$  which is the gradient of the approximation. (The data for this figure was taken from the first iteration of local refinement of Örebro Castle, see Figure 5.) For comparison we also plot the approximation obtained when parameterizing with unit normals as in [6]. Note that the two parameterizations are the same in some (possibly very small) local neighborhood. Therefore their linear approximations will be the same. However, it can be seen from Figure 3 that the unit normal parametrization deviates faster from its approximation. In addition it has singularities.

### C. Problem (14)

Since the function  $D_p$  is one dimensional and sampled densely it is easy to optimize it by simply searching the sample values. To solve (14) we simply recompute the samples of  $D_p$  by adding

$$\lambda_p(y_p - z(\mathbf{p})) + \sigma(y_p - z(\mathbf{p}))^2 \quad (22)$$

to  $D_p(y_p)$  and chose the optimum as the best new sample.

## V. IMPLEMENTATION AND RESULTS

In this section we test our proposed approach. We compare two approaches; Fusion moves with sampled planar proposals without local refinement vs. with local refinement. In all of the experiments we use normalized cross correlation (with a minus sign to get a minimization) with patch size  $3 \times 3$  as data term. As regularization term we use the truncated  $L_1$  term  $h(V_{pq}) = \min(|V_{pq}|, \tau)$ . For the specific choices of  $\lambda$  and  $\tau$  see Figure 5. For local refinement we use the ADMM approach. Since the problem is not convex convergence is not guaranteed for fixed  $\sigma$ . Therefore we start with  $\sigma$  at a low value (0.1 in our implementation) and slowly increase it each iteration to a high value (10 seems to be enough for convergence). The specific update rule is

$$\sigma^{k+1} = \eta\sigma^k, \quad (23)$$

where  $\eta$  is determined by ensuring  $\sigma^k$  is 10 in the last iteration.

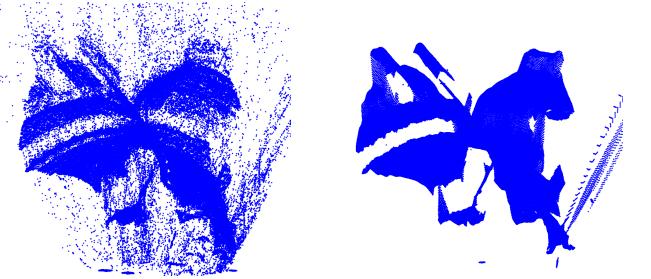


Fig. 4: Comparison of 3D-point positions in unregularized and regularized solutions for the Eglise dataset. Left: unregularized, Right: Regularized. The surface reconstruction (viewed from a different viewpoint) of the regularized solution can also be seen in Figure 5.

In Figure 4 we show the difference between an unregularized solution and one obtained with regularization. We only show the 3D-points since the unregularized solution does not have any normal estimation. (A surface reconstruction of the regularized solution is displayed in the results in Figure 5).

In Figure 5 we show the results of applying the two competing approaches to 4 datasets of varying difficulty; From smooth highly textured surfaces (Église du Dôme) to non-smooth untextured surfaces (Nijo Castle). Note that, in addition to ambiguous texture, the Nijo Castle data set contains people that have walked around between images making the data term incorrect at these positions. We do not handle this in any special way other than applying more regularization. In all of these cases, the fusion moves provide solutions that approximate the underlying surfaces well. However, the planar nature of the proposals gives the appearance of piecewise planarity. In contrast, with local refinement the resulting surfaces have a much smoother appearance and at the same time capture fine details better. In addition the local refinement also repairs some defects, most likely caused by insufficient sampling, such as the hole visible on the roof of the Nijo castle gate.

## VI. ACKNOWLEDGMENTS

This work has been funded by the Swedish Research Council (grants no. 2012-4213 and 2012-4215), the Crafoord Foundation and the Australian Research Councils Discovery Early Career Researcher Award project (DE130101775).

## REFERENCES

- [1] Y. Boykov, O. Veksler, and R. Zabih, “Fast approximate energy minimization via graph cuts,” *IEEE Trans. Pattern Anal. Mach. Intell.*, 2001.
- [2] V. Kolmogorov and R. Zabih, “Multi-camera scene reconstruction via graph cuts,” in *European conf. on Computer Vision*, 2002.
- [3] P. F. Felzenszwalb and D. P. Huttenlocher, “Efficient belief propagation for early vision,” *Int. J. Comput. Vision*, vol. 70, no. 1, pp. 41–54, Oct. 2006.
- [4] V. S. Lempitsky, C. Rother, S. Roth, and A. Blake, “Fusion moves for markov random field optimization,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 32, no. 8, pp. 1392–1405, 2010.
- [5] S. Birchfield and C. Tomasi, “Multiway cut for stereo and motion with slanted surfaces,” in *International Conference on Computer Vision*, 1999.
- [6] C. Olsson, J. Ulén, and Y. Boykov, “In defense of 3d-label stereo,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2013.
- [7] G. Li and S. Zucker, “Differential geometric inference in surface stereo,” *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 32, no. 1, pp. 72–86, 2010.
- [8] O. Woodford, P. Torr, I. Reid, and A. Fitzgibbon, “Global stereo reconstruction under second order smoothness priors,” in *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2009.
- [9] C. Rother, V. Kolmogorov, V. S. Lempitsky, and M. Szummer, “Optimizing binary mrf’s via extended roof duality,” in *IEEE conf. on Computer Vision and Pattern Recognition*, 2007.
- [10] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, “Distributed optimization and statistical learning via the alternating direction method of multipliers,” *Found. Trends Mach. Learn.*, vol. 3, no. 1, pp. 1–122, 2011.
- [11] V. Kolmogorov, “Convergent tree-reweighted message passing for energy minimization,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 28, no. 10, pp. 1568–1583, 2006.

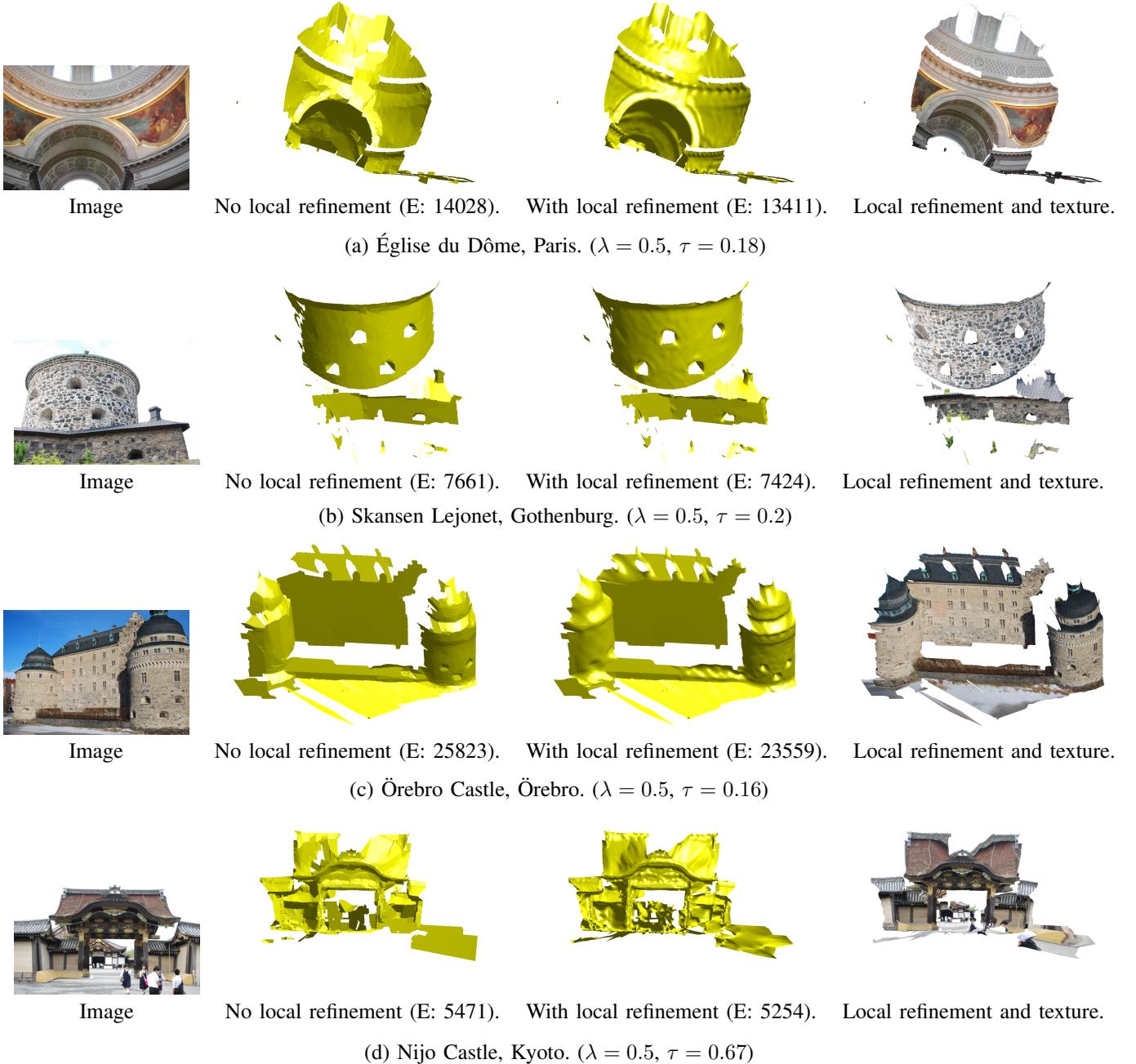


Fig. 5: Reconstructed stereo surfaces, the energy of each solution is given inside parenthesis. After local refinement the energies are lower and the surfaces are more detailed.