

A Method of Estimating Oriented Surface Elements from Stereo Images

Tobias Pietzsch Axel Großmann

Department of Computer Science

Technische Universität Dresden

01062 Dresden, Germany

{Tobias.Pietzsch,axg}@inf.tu-dresden.de

Abstract

In this paper, we describe a technique for the computation of augmented surface patches from stereo images. The stereo vision data is represented as patchlets, which are planar surface elements that have a position, surface normal, size, and confidence measures on the position and normal direction. To estimate the patchlet parameters, we use probabilistic image alignment. Whereas the patchlet generation method by Murray assumes the error on the disparity values to be equal for all pixels, we take the individual intensity patterns into account as well. This way, we obtain better estimates of the confidence measures in the patchlets. We describe the patchlet formation from the disparity and intensity values. To evaluate the quality of the results, we use real and synthetic images as they might be encountered in a typical mobile robot application. For a direct comparison, the images are processed both using Murray's and our approach.

1 Introduction

Stereo vision as a range sensing technology is not new but is gaining popularity in areas such as mobile robotics and augmented reality as the hardware performance has caught up with the processing requirements of this type of sensor. On the one hand, the advantages of stereo vision systems for this kind of applications are obvious. They are relatively cheap, fast, and produce sample points in the full field of view simultaneously. On the other hand, they have a number of drawbacks, too. The smoothing effect of mask-based stereo matching algorithms, occlusions, and poor image contents can make stereo depth extraction incorrect and incomplete. In order to deal with the uncertainty, the representation of stereo data of our choice should contain meaningful confidence measures and facilitate the use of probabilistic methods in the subsequent data interpretation steps.

Recently, Murray proposed a representation of stereo vision data, referred to as patchlets [7]. Patchlets are planar surface elements which have a position, surface normal, size, and confidence on the position and normal direction associated with it. He also proposed a method to form patchlets from stereo disparity images. The algorithm takes into account the uncertainty of the 3D points and propagates them to the patchlet. In the underlying stereo model, the error on the disparity values is assumed to be equal for all pixels, though.

Despite the promising results we obtained using Murray’s patchlet generation technique, we believe that even better estimates of the confidence measures in the patchlets can be found if the algorithm would use the individual intensity patterns as well. Generally speaking, the more salient features we have in the image region used to form the patchlet, the more reliable the estimates of its parameters will be.

If we want to use the intensity patterns to estimate the position and orientation of surface patches, we have to solve an image alignment problem, i.e., bringing into alignment two image regions by finding the parameters of a known function relating position in one region to position in the other, where the parameters are initially only approximately known. Recently, Molton et al. [6] presented the probabilistic inverse compositional image alignment method. In our stereo vision application, we can build upon their technique that combines probabilistic representations with efficient, gradient-based search.

Our contribution is a method for generating augmented surface patches from stereo images. The patchlet parameters are computed from both disparity values and intensity patterns. The method we propose does not necessarily yield better values for the position, surface normal, and size of the surface elements in general, but provides better estimates about the quality of these parameters in terms of confidence values on the position and normal direction. Therefore, the output of our patchlet formation method could be used, for example, to select features which can be tracked reliably, because they have sufficient texture information and are likely to correspond to planar regions in the scene. Reliable confidence estimates could also improve the performance of probabilistic methods for surface segmentation, object detection, and visual mapping. Moreover, we present a detailed empirical comparison with Murray’s patchlet generation method.

In our work, we assume that a mask-based stereo matching algorithm is available to compute the disparity values. As both Murray’s and our implementation is not optimised for speed yet, we do not provide any results about the run-time performance. Both algorithms perform parameter optimisation for each pixel independently and, thus, their theoretical complexity is linear in the number of pixels. Looking for possible applications, we have performed preliminary experiments on surface segmentation. Unfortunately, the presentation of these results would go beyond the scope of this paper.

In the following section, we briefly overview related research on stereo vision and image alignment. In Section 3, we recapitulate the patchlet representation which we use. Section 4 contains a detailed description of the creation of patchlets by image alignment and the underlying sensor model. In Section 5, we evaluate the quality of the results using real and synthetic images as they might be encountered in a typical mobile robot application. Finally, we draw some conclusions in Section 6.

2 Related Work

In this paper, we can draw upon previous work both from stereo vision and image alignment. Murray has obtained promising results for surface segmentation using stereo vision [7]. He introduced the patchlet representation, as described in Section 3, and a plane-fitting technique for patchlet formation. Hattori and Maki [4] proposed a stereo matching algorithm that allows a matching window to locally deform according to the surface orientation. From the intensity gradients within the matching window, they compute surface orientation as well as disparity. However, the rectified stereo setup, which we

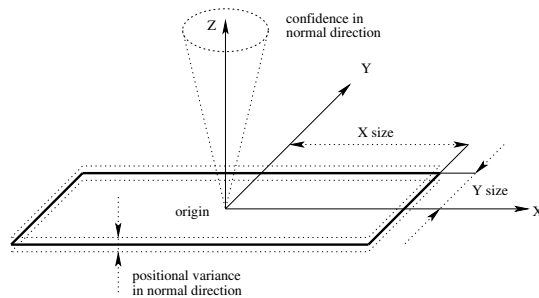


Figure 1: The patchlet model

consider here, is a degenerated configuration for their affine approximation to the image deformation. Devernay and Faugeras [3] used the derivatives of disparity to compute local surface orientation and curvature. The maps for disparity and its derivatives are obtained using *fine correlation*, that is, finding the deformation parameters of local matching windows that maximize correlation. Molton et al. [5] presented planar patch features for a real-time SLAM system, and refine estimates of their surface normals over time. We build upon this work, in that we use their probabilistic formulation [6] of the inverse compositional image alignment algorithm [1] to estimate the orientation of planar surfaces. Our approach differs, in that we simultaneously estimate the depth and orientation of the surface. We also propose a measure of uncertainty that incorporates the intensity differences between the images after converged alignment.

3 Patchlet Surface Elements

Correlation-based stereo matching algorithms typically use a square image mask (in our case, 11x11 pixels) to compute the disparity values. That is, correlation stereo vision is a region matching technique and consequently senses surface patches, not points [7]. Thereby, the surface patch corresponds to the portion of the scene that falls within the neighbourhood of a given stereo pixel as defined by the image mask of the algorithm.

In this paper, we adopt the patchlet model by Murray [7]. One planar surface element, referred to as patchlet, is generated for each valid pixel in the stereo image. The parameters of a patchlet are: a 3D position indicating the centre of the patchlet (origin), a size in both the local X and Y direction, a positional variance in normal direction, λ , and a confidence in normal direction, κ . They are shown in Figure 1.

The size of the patchlet is determined by projecting the pixel on the sensed surface. That is, patchlets that are farther from the camera will be larger in size and are likely to have higher uncertainty values than close patchlets.

4 Estimating Patchlet Parameters Through Image Alignment

Murray's approach to patchlet generation [7] is based on fitting planes to local neighbourhoods of 3D points generated from a disparity map. In contrast, we estimate the patchlet

parameters directly from the deformation of intensity patterns between the images. The parameters of a plane in the scene determine a warp function \mathbf{W} that maps projections of points on the plane between the right and left camera images. We assume that the pixels in a small region around the current point of interest are indeed the projection of a planar surface. The task of computing patchlet position and orientation then becomes to find the warp parameters \mathbf{p} which minimise the dissimilarity between the images

$$\mathbf{p} = \arg \min_{\mathbf{x}} \sum_{\mathbf{x}} [I_L(\mathbf{W}(\mathbf{x}; \mathbf{p})) - I_R(\mathbf{x})]^2. \quad (1)$$

In the above equation, I_L and I_R refer to the intensity in the left and right image respectively. The summation is done over all pixel locations \mathbf{x} in a small window.

Using inverse compositional image alignment [2], the above minimisation problem can be solved by iteratively refining an initial warp estimate. By linearising the intensity function I_R , the parameters of an incremental update \mathbf{W}_i to the current warp \mathbf{W} can be computed as to approximately minimise

$$\sum_{\mathbf{x}} [I_L(\mathbf{W}(\mathbf{x}; \mathbf{p})) - I_R(\mathbf{W}_i(\mathbf{x}; \Delta \mathbf{p}))]^2. \quad (2)$$

The current warp is then updated as $\mathbf{W}_{new} = \mathbf{W} \circ \mathbf{W}_i^{-1}$, and the process is iterated until the parameter estimate converges.

In particular, we employ a probabilistic variation of the inverse compositional algorithm, introduced by Molton et al. [6]. Here, the intensity difference between the images is treated as uncertain measurement of the state $\Delta \mathbf{p}$. The measurement model incorporates zero mean measurement noise with covariance σ_z^2 . The prior estimate of the incremental warp parameters $\Delta \mathbf{p}$ is given by its expected value μ_r and covariance Λ_r , by which we can express the confidence in the initial estimate. By applying Bayes rule, a posterior estimate of $\Delta \mathbf{p}$ can be calculated as

$$\Delta \mathbf{p} = \Lambda_p \left[\Lambda_r^{-1} \mu_r + \sum_{\mathbf{x}} \left(\sigma_z^{-2} \left[\nabla I_R \frac{\partial \mathbf{W}_i}{\partial \Delta \mathbf{p}} \right]^\top (I_L(\mathbf{W}(\mathbf{x}; \mathbf{p})) - I_R(\mathbf{x})) \right) \right] \quad (3)$$

$$\Lambda_p = \left[\Lambda_r^{-1} + \sum_{\mathbf{x}} \left(\sigma_z^{-2} \left[\nabla I_R \frac{\partial \mathbf{W}_i}{\partial \Delta \mathbf{p}} \right]^\top \left[\nabla I_R \frac{\partial \mathbf{W}_i}{\partial \Delta \mathbf{p}} \right] \right) \right]^{-1}. \quad (4)$$

In the present work, we chose the initial warp estimate as follows. First, a standard correlation stereo algorithm [8] is used to compute a dense disparity map. The approximate distance to the plane along a back-projected ray through the origin of the right camera system is parameterised in terms of the disparity d at the patch centre. The prior variance of d is set to a small value according to the expected matching accuracy of the stereo system. Except from the surface being visible in both images, the initial orientation is unknown. By projecting randomly generated planes to a camera, Molton et al. [5] found that the distribution of orientation parameters roughly resembles a Gaussian. We use the same values for orientation uncertainty as they did, and set the initial normal estimate such that the surface is directly facing the right camera. A constant value for all pixels is a sufficient choice for σ_z [6]. In our experiments, the value was 10.0 (intensity levels), which was empirically determined.

In the remainder of this section, we will first introduce the overall and incremental warp functions and their parameterisation. Then, we discuss evaluation of confidence in the warp estimate and finally describe how we compute the patchlet representation from the warp parameters.

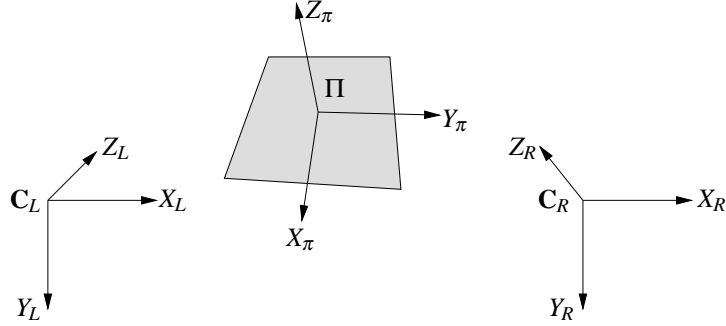


Figure 2: A planar surface is projected in the left and right camera of a stereo system

4.1 The warp functions

It is well known that two projections of a plane are related by a homography that is defined by internal and external camera parameters as well as the parameters of the plane. Given the homography H and the homogeneous coordinates of a pixel $\mathbf{x} = (x, y, 1)^\top$ in the right image the corresponding pixel \mathbf{x}' in the left image can be computed as $\mathbf{x}' = H\mathbf{x}$.

Consider the situation depicted in Figure 2. A planar patch Π is projected in two cameras centred at C_L and C_R . Their optical axis is indicated by the Z_L and Z_R axis of the camera-centred coordinate systems, respectively. We consider a standard rectified stereo setup here, so the relation between the two camera systems is a translation $\mathbf{t}_c = (b, 0, 0)^\top$, where b is the baseline. The coordinate frame attached to the patch is defined such that the plane normal coincides with the Z_π axis. Thus, we have $Z_\pi = 0$ for every point on the plane. Let R_π and \mathbf{t}_π denote the relative rotation and translation between the coordinate systems of the first camera and the plane.

We can decompose the transformation between the images as $H = H_L H_R^{-1}$, where H_L and H_R transform points on the plane to points in the left and right images, respectively. The pinhole model defines the projection of a point in plane coordinates \mathbf{X}_π to the right camera as

$$\mathbf{x} = K [R_\pi | \mathbf{t}_\pi] \mathbf{X}_\pi = KR_\pi [I | R_\pi^\top \mathbf{t}_\pi] \mathbf{X}_\pi. \quad (5)$$

Here, K is the camera calibration matrix which in a rectified stereo setup is identical for both cameras. Because we know that $Z_\pi = 0$ for any point on the plane, Equation (5) simplifies to

$$\mathbf{x} = KR_\pi [I_{12} | R_\pi^\top \mathbf{t}_\pi] \mathbf{x}_\pi = H_R \mathbf{x}_\pi \quad (6)$$

where $\mathbf{x}_\pi = (X_\pi, Y_\pi, 1)^\top$ are the 2D projective coordinates of a point on the plane and I_{12} is the 3×2 matrix consisting of the first two columns of the identity matrix.

Similarly, for the left camera we obtain

$$\mathbf{x}' = KR_\pi [I_{12} | R_\pi^\top (\mathbf{t}_\pi + \mathbf{t}_c)] \mathbf{x}_\pi = H_L \mathbf{x}_\pi. \quad (7)$$

The composition of (6) and (7) yields the required homography $H = H_L H_R^{-1}$ which relates points from the right and left image. Simplification of this homography expression yields (up to scale)

$$H = K [\mathbf{n}^\top \mathbf{t}_\pi I + \mathbf{t}_c \mathbf{n}^\top] K^{-1} \quad (8)$$

where \mathbf{n} is the third column of R_π , the normal vector of the plane. Note that it is not required that \mathbf{n} is a unit vector because H is only defined up to scale.

We choose the origin of the plane coordinate system to lie on the back-projected ray through the centre \mathbf{x}_c of the region of interest, and express the translation $\mathbf{t}_\pi(d)$ as a function of the disparity. Thus, we have the warp function $\mathbf{W}(\mathbf{x}; \mathbf{p}) = H\mathbf{x}$, parameterised in terms of surface normal and disparity.

We parameterise the incremental warp as $\Delta\mathbf{p} = (\alpha, \beta, d_i)^\top$, where d_i is an additive update to the disparity and $(\alpha\mathbf{q}_1 + \beta\mathbf{q}_2)$ is an additive update to the normal. We choose vectors \mathbf{q}_1 and \mathbf{q}_2 perpendicular to each other and \mathbf{t}_π , the normal update is then sufficient to produce any orientation for which the plane is visible in the right camera [5].

The update of the warp estimate is $\mathbf{W}_{new} = \mathbf{W} \circ \mathbf{W}_i^{-1}$, and thus the incremental warp is

$$\begin{aligned} \mathbf{W}_i &= \mathbf{W}_{new}^{-1} \circ \mathbf{W} \\ &= (\mathbf{K} [\mathbf{n}(\alpha, \beta)^\top \mathbf{t}_\pi(d + d_i) + \mathbf{t}_c \mathbf{n}(\alpha, \beta)^\top] \mathbf{K}^{-1})^{-1} (\mathbf{K} [\mathbf{n}^\top \mathbf{t}_\pi(d) + \mathbf{t}_c \mathbf{n}^\top] \mathbf{K}^{-1}) \mathbf{x} \quad (9) \\ &= \mathbf{K} [\mathbf{n}^\top \mathbf{t}_\pi(d + d_i) + \mathbf{t}_c \mathbf{n}(\alpha, \beta)^\top]^{-1} [\mathbf{n}^\top \mathbf{t}_\pi(d) + \mathbf{t}_c \mathbf{n}^\top] \mathbf{K}^{-1} \mathbf{x}, \end{aligned}$$

where $\mathbf{n}(\alpha, \beta) = \mathbf{n} + \alpha\mathbf{q}_1 + \beta\mathbf{q}_2$.

An interesting observation is, that while being homographies for general camera configuration the warps \mathbf{W} and \mathbf{W}_i degenerate to affinities for the rectified scenario. Because the epipolar constraint must hold for all corresponding points, we have $y' = y$. This implies that the partial derivatives of the warps with respect to y are zero, and consequently the intensity gradient in y direction can be excluded from the alignment procedure.

4.2 Evaluation of uncertainty in the parameter estimate

An important feature of the patchlet representation is that it provides measures of confidence in the position and orientation parameters. In Murray's patchlet generation method [7] disparity values are used to triangulate 3D points, each associated with uncertainty according to the expected matching accuracy of the correlation stereo algorithm used to compute the disparity values. A best-fit plane to a local neighbourhood of these points is then computed and the confidence measure defined by propagating the uncertainty from the points to the plane parameters.

From our point of view, a disadvantage of this method is to neglect the fact that the accuracy of correlation matching is influenced by the amount of texture in the pixel windows used. Furthermore, if intensity gradients are used to estimate local surface orientation as well, as in our method, image texture will provide information on the confidence of individual parameters. For instance, an otherwise textureless pixel window containing a single centred black dot contains sufficient information to compute disparity while leaving surface orientation uncertain.

The posterior covariance (4) computed by probabilistic image alignment provides a measure of uncertainty in the converged parameter estimate. Note, however, that only the intensity gradients of the right image play a role in its computation. This is sufficient, if our model of image formation is correct, that is if the considered set of pixels is indeed the image of a planar surface, visible in both views. Frequently, however, these assumptions are violated. For instance, a pixel window in the right image spanning an occlusion

boundary may contain rich texture, but despite the low posterior covariance, we should not put much trust in the estimated plane parameters.

A way to check for violation of our *a priori* assumptions is to analyse the intensity differences at the converged parameter estimate. We expect the intensity errors $I_L(\mathbf{W}(\mathbf{x}; \mathbf{p})) - I_R(\mathbf{x})$ evaluated at the posterior estimate \mathbf{p} to resemble a standard distribution with zero mean and variance σ_z^2 . Pixel windows, for which this is not the case, might be simply discarded as incompatible with our image formation model. Instead, we might assume that the deviation from $N(0, \sigma_z^2)$ increases with the severity of the violation of assumption. A surface with a slight curvature can be expected to be more “well-behaved” than an actual occlusion boundary. In our experiments, we have found that simply using the actual variance of intensity errors $\bar{\sigma}_z^2$ to recompute the posterior covariance as

$$\Lambda_p = \left[\sum_{\mathbf{x}} \left(\max(\bar{\sigma}_z, \sigma_z)^{-2} \left[\nabla I_R \frac{\partial \mathbf{W}_i}{\partial \Delta \mathbf{p}} \right]^\top \left[\nabla I_R \frac{\partial \mathbf{W}_i}{\partial \Delta \mathbf{p}} \right] \right) \right]^{-1}, \quad (10)$$

already gives a reasonable estimate of the uncertainty in the parameters.

4.3 Generation of patchlet parameters

Constructing a patchlet from the parameters of the posterior warp estimate is straightforward and mostly along the same lines as in [7]. The origin of the patchlet is given by the translation $\mathbf{t}_\pi(d)$, and its Z axis by the plane normal \mathbf{n} . The X and Y axes are determined as the major and minor axis of the ellipse that is the intersection of the plane and a cone through the centre pixel of the generating region. Patchlet size along those axes is determined by projecting the pixels boundaries to the patchlet plane. We propagate the uncertainty from the warp parameter space to patchlet coordinate space, to obtain positional and directional confidence.

5 Experimental Evaluation

The performance of stereo vision algorithms typically depends on the structure of the scene, the image information in the textures, and occlusions. We decided to focus on applications in mobile robotics such as constructing geometrical models of indoor environments. That is, planar surfaces are common and the object distances are below 5 metres or so. As the stereo camera is to be mounted on a mobile robot, we used a Point Grey Research Bumblebee camera. The dense disparity image is generated with sum-of-absolute-differences correlation stereo [8]. The experiments were performed both using real and synthetic images. The real images were used to evaluate the patchlet clouds by visual inspections.

Figure 3 shows exemplary results of our algorithm for an outdoor scene. In the top row, the reference (right camera) and disparity images are shown. In the bottom-left, the generated patchlet cloud is rendered from a different point of view. In the close-up view in the bottom-right the individual patchlets are clearly visible.

Visualisations of patchlets generated from various real images indicates that our method can reliably estimate patchlets for regions of the image containing approximately planar scene features. By incorporating surface orientation, our method often improves the disparity results obtained by the correlation stereo method, especially for surfaces tilted away

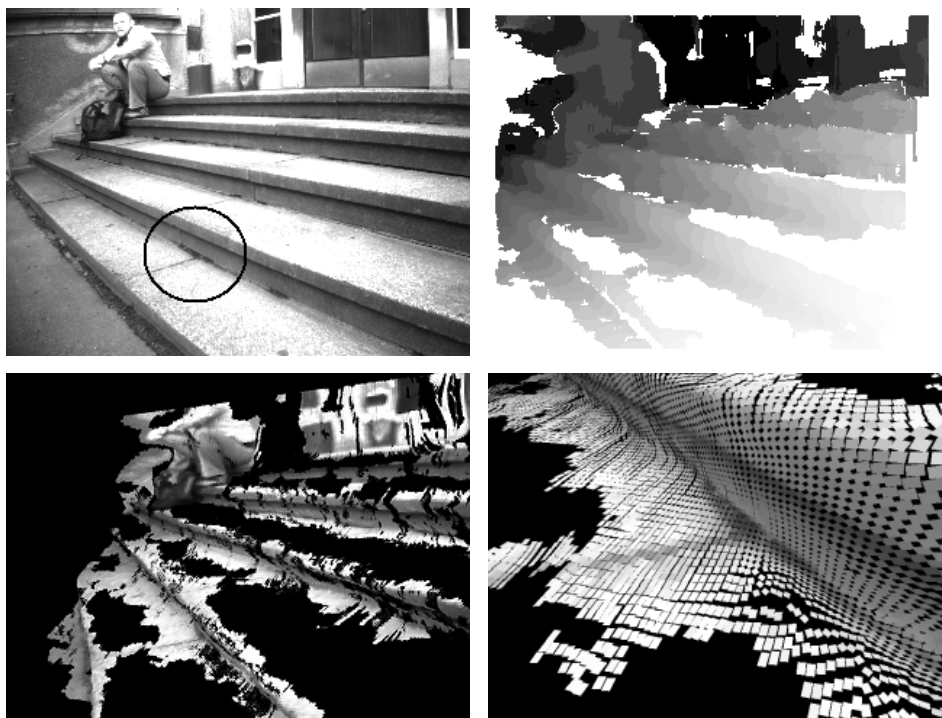


Figure 3: Reference and disparity image and visualisation of the resulting patchlets.

from the camera. We checked that this is not an artifact of the correlation algorithm by excluding the orientation parameters from the alignment procedure.

In a direct comparison, the patchlet clouds generated by Murray’s method often appear smoother, i.e. our method produces gaps between neighbouring patchlets more often. On the other hand, for planar surfaces that are not directly facing the camera our results often look much more even. In general, our results are visually very similar to Murray’s.

The synthetic images were obtained using a 3D robot simulator. The implementation of the Bumblebee camera in the simulator makes use of the OpenGL depth buffer to compute the true 3D position and orientation of the patchlets, in addition to the left and right camera images.

In the simulation, the camera was moving toward a pile of cardboard boxes and then circled around these objects, maintaining a distance of 1.5m. New images were taken at a step size of 10cm. Along the complete trajectory, we recorded 114 images pairs. An example reference image and the corresponding disparity image is given in Figure 4.

For two patchlets a and b , we can specify a position error, a pointing error, and a coordinate origin error, as illustrated in Figure 5 (top-left). The accuracy results of the patchlet estimation methods are strongly affected by outliers. Ideally, we want to consider only those patchlets with high confidence values. For this purpose, the patchlet representation contains the positional variance in normal direction, λ , and the confidence in normal direction, κ . That is, we can sort the computed patchlets for increasing values of λ or decreasing value of κ , respectively, and consider only the first n patchlets of the resulting

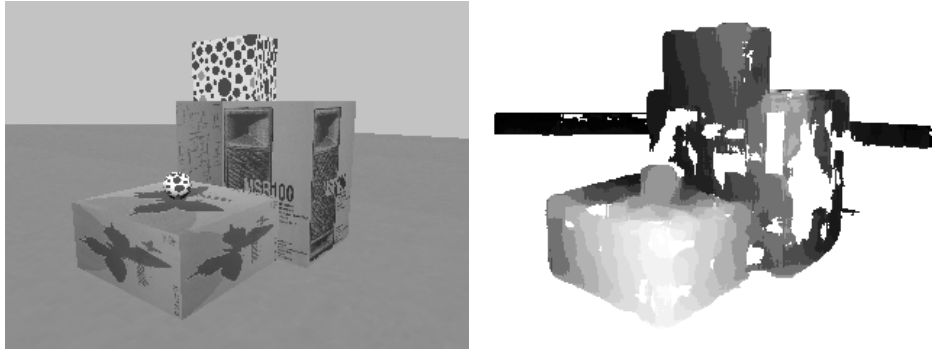


Figure 4: Example reference/disparity image pair of the cardboard boxes.

ordered sequences.

In the given scenario, the stereo matching algorithm computed about 14,000 valid disparity values for a camera image of size 320x240. Each disparity image was processed both by the presented algorithm (denoted by PIA) and Murray's patchlet generation method (denoted by PGR). For 96.2% (PIA) and 99.4% (PGR) of the disparity values, we obtained a corresponding patchlet, respectively. A small number of disparity values was disregarded by the algorithms.

The accuracy comparison of the two algorithms is given in Figure 5. The X-axis denotes the selection ratio based on the parameters λ and κ , respectively. A ratio of 0.1 means that the error was computed for the 10% of patchlets with the smallest value of λ or the largest values of κ , respectively. The larger the selection ratio, the more patchlets were evaluated for each image. If the selection ratio is close to 1.0, we see that PIA and PGR have similar error values. Please note, for better visibility, the error bars in all diagrams are 1/10 of the standard deviation. The error values are the averages obtained for all 114 images.

The diagrams illustrate that PIA yields better estimates of the λ and κ parameters than PGR. For small selection ratios, PIA has got considerably smaller errors. As one would expect, the larger the selection ratio (i.e., an increasing number of 'bad' patchlets are being incorporated), the larger the resulting error. This indicates that the confidence values obtained by PIA reflect the real situation more accurately than those obtained by the original PGR algorithm. In fact, the λ and κ values computed by PIA can be used to differentiate between good and bad patchlet approximations.

6 Conclusions

We expect that stereo vision will become increasingly popular in mobile robotics as a 3D range sensing method. Application for tasks such as navigation and mapping require explicit representation of uncertainty in the sensor data. We believe that the patchlet representation is a useful low-level primitive for this purpose. We have presented a new method to compute patchlets. Our experiments have shown that it is important to incorporate image texture in the computation of confidence measures.

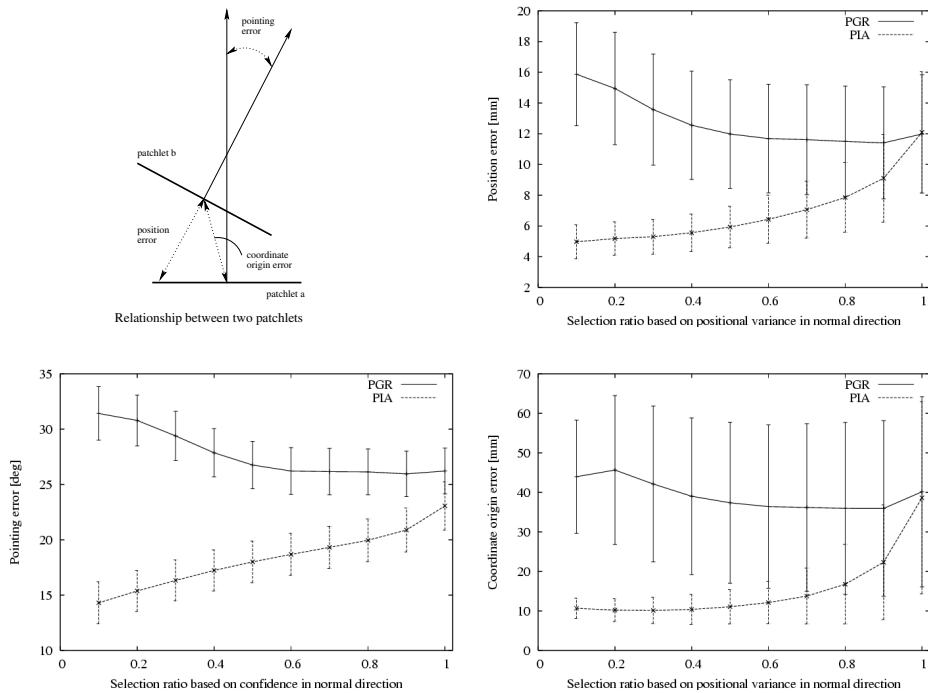


Figure 5: Quantitative comparison of the two patchlet estimation techniques.

References

- [1] Simon Baker, Frank Dellaert, and Iain Matthews. Aligning images incrementally backwards. Technical Report CMU-RI-TR-01-03, Robotics Institute, Carnegie Mellon University, Pittsburgh, PA, USA, 2001.
- [2] Simon Baker and Iain Matthews. Lucas-Kanade 20 years on: A unifying framework. *International Journal of Computer Vision*, 56(3):221–255, March 2004.
- [3] F. Devernay and O. D. Faugeras. Computing differential properties of 3-D shapes from stereoscopic images without 3-D models. In *Proc. of the ICVPR*, 1994.
- [4] H. Hattori and A. Maki. Stereo matching with direct surface orientation recovery. In *Proc. of the BMVC*, 1998.
- [5] N. D. Molton, A. J. Davison, and I. D. Reid. Locally planar patch features for real-time structure from motion. In *Proc. of the BMVC*, 2004.
- [6] N. D. Molton, A. J. Davison, and I. D. Reid. Parameterisation and probability in image alignment. In *Proc. of the Asian Conference on Computer Vision*, 2004.
- [7] Donald R. Murray. *Patchlets: a method of interpreting correlation stereo 3D data*. PhD thesis, The University of British Columbia, 2003.
- [8] Point Grey Research. Triclops Software Development Kit, 2003.