

Segmenting color images into surface patches by exploiting sparse depth data

Babette Dellen Guillem Alenyà Sergi Foix Carme Torras
Institut de Robòtica i Informàtica Industrial (CSIC-UPC), Barcelona, Spain
{bdellen, galenya, sfoix, torras}@iri.upc.edu

Abstract

We present a new method for segmenting color images into their composite surfaces by combining color segmentation with model-based fitting utilizing sparse depth data, acquired using time-of-flight (Swissranger, PMD CamCube) and stereo techniques. The main target of our work is the segmentation of plant structures, i.e., leaves, from color-depth images, and the extraction of color and 3D shape information for automating manipulation tasks. Since segmentation is performed in the dense color space, even sparse, incomplete, or noisy depth information can be used. This kind of data often represents a major challenge for methods operating in the 3D data space directly. To achieve our goal, we construct a three-stage segmentation hierarchy by segmenting the color image with different resolutions - assuming that “true” surface boundaries must appear at some point along the segmentation hierarchy. 3D surfaces are then fitted to the color-segment areas using depth data. Those segments which minimize the fitting error are selected and used to construct a new segmentation. Then, an additional region merging and a growing stage are applied to avoid over-segmentation and label previously unclustered points. Experimental results demonstrate that the method is successful in segmenting a variety of domestic objects and plants into quadratic surfaces. At the end of the procedure, the sparse depth data is completed using the extracted surface models, resulting in dense depth maps. For stereo, the resulting disparity maps are compared with ground truth and the average error is computed.

1. Introduction

The identification and segmentation of 3D surfaces from an image is an important step towards solving robotic object-manipulation tasks as it facilitates object recognition, grasp-point selection, and, in consequence, the execution of appropriate grasping movements. We are interested in developing efficient tools for representing and manipulating domestic objects, including deformable ones. In particular, we aim to segment plants into their composite struc-

tures, i.e., leaves or part of leaves, and to extract color and 3D-shape descriptors to find points of interest with which the robot can interact. However, the characterization and classification of surfaces is only possible if appropriate 3D information is available. Various techniques for 3D data acquisition exist, but there is also always a critical trade-off between the accuracy of the method and its efficiency in terms of computation time, applicability (active versus passive methods), and cost. Laser range scanning for example delivers accurate and dense depth information, but has the drawback that it is very time consuming and thus not practical in the context of our task. Photonic mixer devices (PMDs) or Swissranger cameras deliver depth images in real time using the time-of-flight principle, but are of low resolution and afflicted with uncertainties. Stereo vision has the advantage that it is applicable in most environments but tends to fail in untextured areas (due to correspondence problems) and hence often only delivers sparse and noisy results. This problem motivated us to develop a method for surface segmentation that is also applicable to sparse depth information. It further should allow us to complete missing information.

The segmentation of 3D information is the task of dividing the image into regions so that all the points of the same surface belong to the same region. The regions shall not overlap and, taken together, generate the entire image. Many algorithms for solving this task have been proposed in the past [2, 11, 1, 6, 7, 8]. We roughly distinguish between two main groups: region-based and edge-based segmentation algorithms. Algorithms of the first group segment the image into initial regions which are then merged or extended [1, 7, 8]. Since segmentation is applied at a region level, these methods often produce distorted boundaries. Algorithms of the second group find jump boundaries in the depth image, providing an initial segmentation, which is then refined by fitting quadratic functions to the initial segments [11, 6]. A drawback of these methods is that discontinuities in the depth data are hard to detect and may result in an initial under-segmentation of the image, which cannot be corrected at later steps.

Depth acquisition techniques having the advantage of be-

ing fast and economic (stereo, Swissranger, PMD), or non-invasive (stereo), only provide sparse and noisy depth data, constituting a major problem for the segmentation algorithms described above, which have been developed for 3D information acquired with laser range scanners or structured light. In the case of sparse data, information from other sensors, e.g., color, is required to tackle the given task. Walhoff et al. (2007) combined depth information of a PMD camera with color images to segment objects from their background [12]. In a related work by Bleiweiss and Werman (2009), depth from time-of-flight was fused with color data to improve object segmentation and tracking [3]. Real-time foreground segmentation was also achieved by Crabb et al. (2008), combining range and color imaging [4]. Our goal, however, is the segmentation of images into surface patches, which was not approached by these works.

In this paper, we propose to segment dense color images at different resolutions and to select those segments for which a best fit to the sparse depth data can be obtained using quadratic surface models. In some sense, our method is related to edge-based segmenters, with the important difference that we find “edges” in the color space, not in the depth image. The selected segments are then merged and optionally grown using both their respective surface model and color information. The method is based on the assumption that true surface boundaries must emerge at some point along the segmentation hierarchy. It has the advantage that segmentation can be performed in a dense space even though the data of interest are sparse. Under- and over-segmentation are avoided by considering segmentations at various resolutions. Since segmentation in color space usually provides sharp edges, accurate surface boundaries can be obtained in most cases. Furthermore, the obtained (explicit) surface models can be used to interpolate depth data into regions that are initially undefined.

2. Algorithm

We describe an algorithm for the segmentation of color images into surface patches utilizing sparse depth data. The core idea of the method is based on the notion that surface boundaries are in most cases represented by an edge in the color image (the reverse is often not true). Since we are dealing with sparse depth data, it is further desirable to have as large segments as possible - otherwise model fitting becomes impracticable due to lack of data inside segments. We thus segment the color image with different resolutions (see Fig. 1A). Quadratic surface models are fitted to each segment, and we select those segments from the hierarchy which minimize the total fitting error, while taking into account the hierarchy level, i.e., segments obtained at lower resolutions are given preference to segments at higher ones.

The algorithm thus consists of the following steps (see also Fig. 1B). In step 1, the color image and the correspond-

ing sparse depth data is acquired (see Section 4). Then, in step 2, color segmentation is applied to the color image and segmentations at three different resolutions are obtained, creating a segmentation hierarchy as illustrated in Fig. 1A (see also Section 2.1). In step 3 of the algorithm, surface models (quadratics) are fitted to each segment utilizing the sparse depth data, and, for each segment, the best fitting model is selected (see Section 2.2). Then, in step 4, those segments are selected from the hierarchy that produce the total best fit to the data, as illustrated in detail in Fig. 1C (see also Section 2.3). Segments at lower resolutions are given preference to segments at higher resolutions to avoid over-segmentation. The resulting new segmentation is then further improved by applying an additional region merging step (step 5) through which segments having highly different color values can be merged if they describe the same surface (see Section 2.4). Then, in a final step (6), unclustered points are assigned to the closest surface, using both depth and color information (see Section 2.5).

2.1. Hierarchical color segmentation

The color image is segmented using the method of superparamagnetic clustering of data employed in [5] which allows a segmentation hierarchy to be generated by segmenting with different resolutions. For this purpose, we varied the interaction strength by multiplying the mean distance $\bar{\Delta}$ by a factor between 1.4 and 0.6. Segments smaller than a threshold (here 10 pixels) are considered being unlabeled and are excluded from the steps 3-5 of the algorithm. The segmentation algorithm can be replaced by any other method as long as different resolutions result in a segmentation hierarchy similar to the one illustrated in Fig. 1A.

2.2. Model fitting and selection

For each color segment s_i and model type (see Section 3) we perform a minimization of the mean square distance

$$E_{i,\text{model}} = 1/N \sum_j (z_j - z_{j,m})^2 \quad (1)$$

of measured depth points $z_{j,m}$ from the estimated model depth $z_j = f_{i,\text{model}}(x_j, y_j)$, where $f_{i,\text{model}}$ is the data-model function and N is the number of measured depth points in the area of segment s_i . The optimization is performed with a Nelder-Mead simplex search algorithm provided in MATLAB.

The mean square errors for two different model types, i.e. $E_{i,\text{planar}}$ and $E_{i,\text{curved}}$, are computed. We select the planar model if $E_{i,\text{planar}} < E_{i,\text{curved}} + \tau_1$, and the curved model otherwise.

2.3. Color-segment selection procedure

We define the following *selection procedure* considering first only two levels u and $u + 1$, where a higher level de-

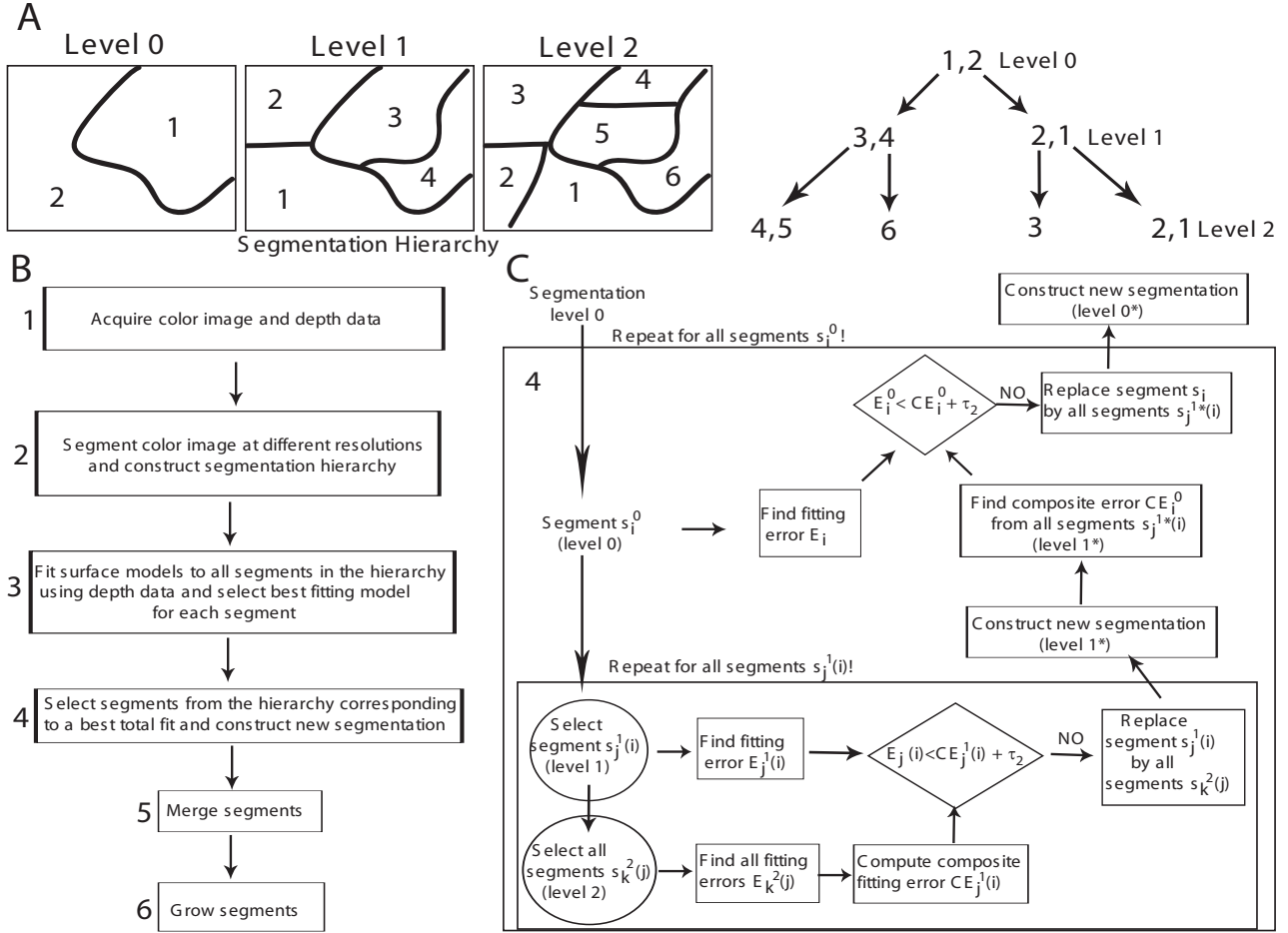


Figure 1. Surface segmentation algorithm. **A** Hierarchical segmentation of an image and resulting segmentation hierarchy. **B** Flow chart of the algorithm. The last step is optional. **C** Detailed schematic of step 4 of the algorithm.

notes a higher resolution. Let s_i^u be a segment at level u of the hierarchy having a fitting error E_i^u (see Section 2.2). At the level $u + 1$ of the hierarchy, segment s_i^u is composed of k segments $s_j^{u+1}(i)$ with respective fitting errors E_j^{u+1} . The composite error of the k segments $s_j^{u+1}(i)$ at level u is then defined as

$$CE_i^u = \sum_j a_j^{u+1} E_j^{u+1} / \sum_j a_j^{u+1} \quad , \quad (2)$$

where a_j^{u+1} is the number of valid depth points in the area of segment s_j^{u+1} . We select s_i^u at level u if

$$E_i^u < CE_i^u + \tau_2 \quad , \quad (3)$$

and the k segments $s_j^{u+1}(i)$ from level $u + 1$ otherwise. Parameter τ_2 is introduced in order to avoid an over-segmentation of the image by preferring segments obtained

at lower resolutions. The procedure is applied to each segment at level u . In this way a new segmentation is constructed which replaces the initial segmentation at u by a segmentation u^* .

Let us now consider a segmentation hierarchy consisting of p levels. We apply the selection procedure to the initial segmentations $u = p - 1$ and $u = p$. The selection procedure is applied to the initial segmentation $u = p - 2$ and $u^* = p - 1$, and so on, until the end of the hierarchy is reached. In this paper, we choose a three-level hierarchy. More levels can be included if desired. The procedure for a three-level hierarchy is illustrated in Fig. 1C.

2.4. Region merging

Let us consider two segments s_i and s_j having fitting errors E_i and E_j . Both segments are merged if $E_{i \cap j} <$

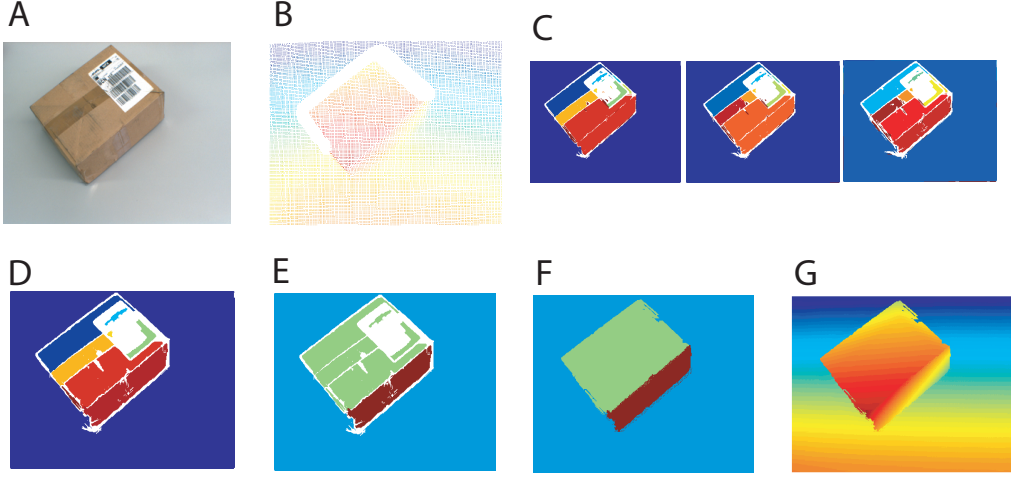


Figure 2. Illustrative example: Segmentation results for a carton box using Swissranger time-of-flight depth. **A** Color image. **B** Sparse depth data, transformed to the color space. **C** Color segmentations at different resolutions (levels 0-2), respectively. **D** Selected segments before merging (after step 4 of the method). **E**. Selected segments after merging (after step 5). **F** Segments after applying region growing (step 6). **G** Fitted depth using surface models. Unclustered points are shown in white.

$(a_i E_i + a_j E_j) / (a_i + a_j) + \tau_2$, where $E_{i \cap j}$ is the fitting error of the merged segments and a_i and a_j is number of valid depth points in the area of the segments s_i and s_j , respectively. This procedure allows segments that have different colors to be merged. The procedure is applied to all segments that are neighbors of each other (i.e. their closest pixels have to be less than 8 pixels apart). When accepting a merge, segments are updated and the new segmentation is used when evaluating the remaining segment pairs.

2.5. Region growing

2.5.1 Time-of-flight

Let p_i be a previously unclustered point with coordinates (x_i, y_i, z_i) and color c_i . We find all segment neighbors of this pixel within a radius of 5 pixels. We compute the distance of p_i to the surface of segment s_j as $\text{dist}_i^j = |z_i - f_j(x_i, y_i)|$, where f_j is the explicit surface-model function of segment s_j , and assign p_i to the closest segment in the neighborhood. For points for which no depth value was originally measured we use the local mean depth value computed over a small area around the point.

2.5.2 Stereo

Let p_i be a previously unclustered point with image coordinates (x_i, y_i) . We find all neighboring segments of this pixel within a radius of 10 pixels. We compute the distance of p_i to the surface of segment s_j as $\text{dist}_i^j = |I_L(x_i) -$

$I_R(x_i - f_j(x_i, y_i))|$, where f_j is the explicit surface-model function of segment s_j , and I_L and I_R are the left and right color images, respectively. According to this distance measure, we assign p_i to the closest segment in the neighborhood.

3. Surface types

We choose two types of surfaces as surface models: planes and a quadratic function, which allows (among others) the modeling of spherical and cylindrical shapes. If desired, higher-order terms can be included. Here we use quadratic functions that allow computing depth z explicitly for the x - y coordinates in the form of $z = f(x, y)$.

3.1. Planes

Planar surfaces are described by three parameters a , b , and c , where the depth z can be expressed as a function of x and y through

$$z = ax + by + c \quad . \quad (4)$$

3.2. Curved surfaces

Curved surfaces are described by five parameters a , b , c , d , and e , where the depth z can be expressed as a function of x and y through

$$z = ax^2 + by^2 + cx + dy + e \quad . \quad (5)$$

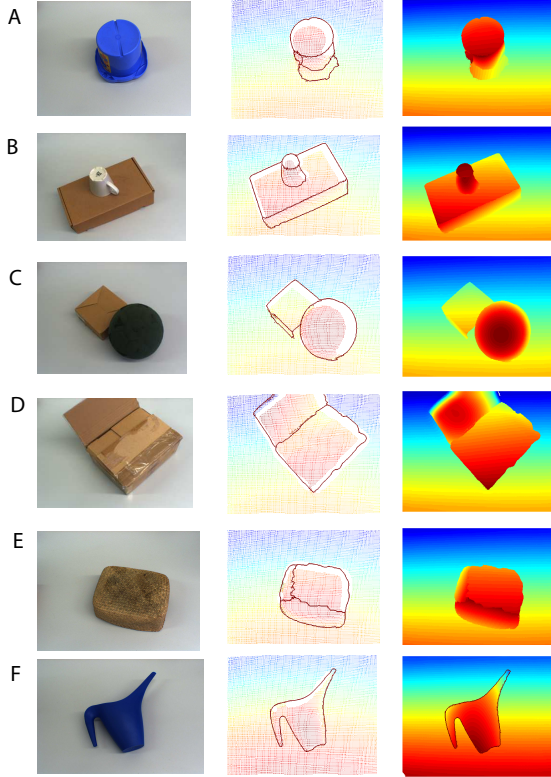


Figure 3. Segmentation results for several domestic objects. Left panels: Original color image. Middle panels: Initial Swissranger sparse depth plotted together with final segment boundaries. Right panels: Fitted depth using segment surface models.

This function has sufficient complexity to allow modelling of basic curved shapes, including spherical, cylindrical, and conical shapes.

4. Data Acquisition

4.1. Time-of-flight depth

We acquire color images together with depth images and combine them by bringing both sensor-pixel matrices in correspondence. For our experiments a Swissranger SR3100 camera and a PMD CamCube camera have been used. The camera is provided with its own illumination system, composed of a set of modulated infra-red LEDs. The SR3100 camera has a low pixel resolution of 176×144 , while the PMD camera has a resolution of 204×204 , which is however still far below the resolutions of standard RGB cameras. Both cameras have a high frame rate average of 25 fps. This high frame rate makes these cameras suitable for real time applications. Depth and RGB color, provided by a PointGrey Flea camera, can be easily combined for a specific depth range thanks to the precise intrinsic and ex-

trinsic camera calibration. To perform a good data registration reducing partial occlusions, both cameras have to be close to each other (few centimeters) and share the same field of view. By means of using multiple images from both cameras with the same calibration pattern, intrinsic and extrinsic parameters for each camera can be computed. Once this information is achieved, both sensor pixel matrices can be put in correspondence using the calibration pattern as a reference frame and hence registering depth and color information.

Due to the different viewpoints of the two cameras, occlusions can occur in particular for close objects. These occlusions are detected and removed using a buffer approach: The 3D point cloud is transferred to the RGB camera reference space using extrinsic parameters and, if several points are detected along the same line of view, the closest point to the camera is selected. Since the ToF camera has a smaller resolution than the RGB camera, the 3D points are collected from a region around the line of view.

4.2. Stereo disparity

We acquire sparse stereo disparity using a recent algorithm proposed in [5], which is applicable to weakly textured images and thus combines well with color-based segmentation. In this method, stereo segments are found for stereo images and segment silhouettes are computed in both frames. Unique correspondences of silhouette-edge points are searched, and the respective disparities are calculated. Occluded edges are identified and removed from the edge disparity map. Additionally, texture inside segments is exploited by applying a window-based stereo algorithm which operates strictly inside stereo segments. Confidence values are computed and only those disparities that have a high confidence are used, resulting in sparse disparity maps. In our method, we use these sparse disparity maps directly as input without employing the interpolation process described in [5]. We further use a phase-based method to establish initial correspondences [9].

5. Results

The algorithm is tested using both time-of-flight depth data (see Section 4) and stereo disparity. In the case of time-of-flight depth, color images and depth images are acquired together and combined by bringing both sensor-pixel matrices in correspondence using the respective camera-calibration information (see also Section 4). The simulations were performed using MATLAB with an Intel Duo Core Processor T2250 of 1.73GHz. Total run time of the algorithm using non-optimized code for segmenting the time-of-flight data in the carton box example (Fig. 2) is 88 s (including computation times for getting color segmentations). Parameters for segmenting time-of-flight data are

$\tau_1 = 2 \text{ cm}^2$ and $\tau_2 = 0.25 \text{ cm}^2$. For stereo disparity, we have $\tau_1 = 0.3 \text{ pixels}^2$ and $\tau_2 = 0.1 \text{ pixels}^2$.

5.1. Time-of-Flight

We first illustrate the algorithm on the example of images of a carton box taken with the RGB and the Swissranger depth camera. In Fig. 2A-B, the color image and the corresponding Swissranger depth are presented, respectively. Points for which no depth is computed (due to occlusions and the smaller size of the depth images compared to the color images) are plotted in white. The results of color segmentation at different resolutions (levels 0-2) are displayed in Fig. 2C. Unclustered points are plotted in white color. The results after step 4 of the algorithm are displayed in Fig. 2D. We observe that those color segments corresponding best to the surface structure of the objects got selected. Remaining over-segmentations are removed by applying step 5 of the algorithm. The result is shown in Fig. 2E. Region growing assigns previously unclustered points to the closest surface (see Fig. 2F). The surface models that have been determined can now be used to fit depth to the segment regions, as shown in Fig. 2G. We obtained a total fitting error of 1.1 cm. Fitting errors of the following examples are in a similar range.

Next, we apply the algorithm to color images with Swissranger depth of several domestic objects, i.e., boxes, a cup, a basket, a ball, and a water can, containing planar, spherical, cylindrical shapes, and other curved shapes (Fig. 3A-F). In the left panels, the original color images are shown. In the middle panel, final segment boundaries (after region growing) are presented together with the initial depth data. In the right panels, the final fitted depth is shown. We observe that the algorithm successfully segments most of the prominent surfaces, including curved ones. Only in areas of highly complex structure or insufficient depth information, failures are observed, e.g., the handle of the cup (Fig. 3B) (which was lost during the procedure) and the edge of the cylinder in Fig. 3A. In the example shown in Fig. 3D, we have some problems with false points in the depth map which suggest the upper segment of the box to be curved even though we would expect it to be planar.

In the next step, we apply the method to color-depth images of plants. Here depth is recorded with a PMD camera. Plant images are challenging because they contain many depth layers and occlusions, caused by overlapping leaves, and weak contrast boundaries separating adjacent leaves. The results are presented in Fig. 4A-E. In the left panels, the original color images are shown. In the middle panels, final segment boundaries (after region growing) are presented together with the initial depth data. In the right panels, the final fitted depths are shown together with the final segment boundaries. Even though plants exhibit complicated shapes and have many occlusions, most of the main surfaces have

Scene	E_{fin}	E_{merge}	ρ_{merge}	E_{pb}
Tsukuba	0.67	0.56	81	0.72
Venus	0.34	0.31	83	0.73
Teddy	1.05	0.75	79	3.61
Cones	1.51	1.1	73	2.3
Babyl	1.61	1.18	65	3.78
Lampshade2	2.95	2.67	87	6.79
Plastic	3.31	3.14	87	6

Table 1. Average error in pixels for Middlebury images of final disparity maps E_{fin} , disparities after merging E_{merge} (with densities ρ_{merge}), and of the phase-based approach E_{pb} [9].

been found, often corresponding to leaves or at least part of leaves, and curved shapes could be modelled correctly in most cases (for example the large leaf at the bottom in Fig. 4A). Basic segment properties such as mean color, size, and mean fitting error are computed, and, based on these criteria, candidate segments (e.g. for robot manipulation) are selected representing leaf structures. An exemplary candidate segment has been marked red for each plant (Fig. 4, left panels). Also the center point of the segments has been marked red.

5.2. Stereo disparity

Next we apply the algorithm to stereo images from the Middlebury stereo database (URL: vision.middlebury.edu/stereo/) [10]. Three examples are shown in Fig. 5A-C, left panels (Lampshade2, Venus, and Teddy). A sparse disparity map of the scene is computed using a recent algorithm proposed in [5]. We use disparity as input for z . The initial sparse disparity maps are shown together with computed final segment boundaries in Fig. 5A-C, middle panels. Points for which no disparity could be found are shown in white color. The extracted surface models of the segments can then be used to create dense disparity maps, as presented in Fig. 5A-C, right panels. The basic structure of the scene could be captured and surfaces identified. Results of our method for Tsukuba, Venus, Teddy, and Cones have been evaluated using the Middlebury evaluation, summarized in Fig. 6 [10]. We further provide the average disparity error for several Middlebury images in Table 1. These results demonstrate that the method can be used to find dense disparity maps even in cases where little information is given at the beginning, e.g., in case of weakly textured images.

6. Discussion

We proposed a novel algorithm for the segmentation of color images into surface patches using sparse depth data, acquired using either time-of-flight or stereo techniques. The color image is segmented with different resolutions, 3D

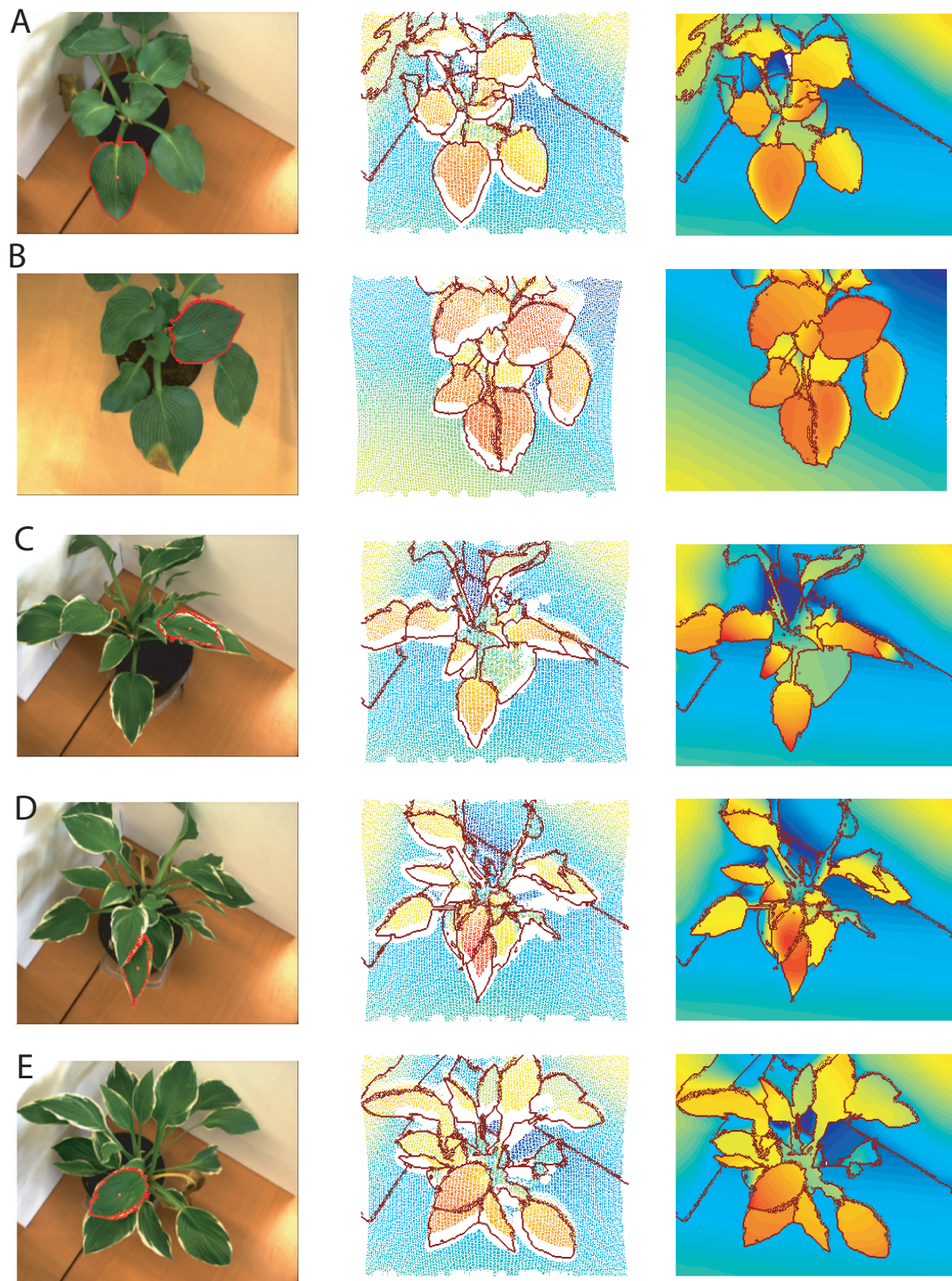


Figure 4. Segmentation results for plants. Left panels: Original color image together with an exemplary candidate segment boundary (marked in red) (see text). Middle panels: Initial PMD sparse depth plotted together with final segment boundaries. Right panels: Fitted depth using segment surface models plotted together with the final segment boundaries.

