# 3D OBJECT MODELING WITH A VOXELSET CARVING APPROACH

*Giovanni Dainese, Marco Marcon, Augusto Sarti, Stefano Tubaro*

Dipartimento di Elettronica e Informazione - Politecnico di Milano
Piazza Leonardo Da Vinci 32, 20133 Milano, Italy
phone: +39-0223999639, fax: +39-0223999611
email: dainese/marcon/sarti/tubaro@elet.polimi.it

## ABSTRACT

In the past few years several systems for object reconstruction based on the analysis of 2D images have been proposed. In order for such systems to be of practical use, the 3D data extraction process is expected to be fast and reliable. In this paper we propose a general approach for the reconstruction of complete 3D objects based on a mesh fusion algorithm. Every surface patch is obtained as a depth map using an algorithm based on graph cuts theory. Each depth map is then triangulated before using it in a fusion algorithm based on a voxel-set carving approach. The result of the process is a closed mesh representing the object surface with sub-voxel resolution.

## 1. INTRODUCTION

Reconstructing an object's tridimensional shape for a set of cameras is a classic vision problem. In the last few years, it has attracted a great deal of interest, partly due to the number of application both in vision and in graphics that require good reconstructions. In order to get the complete model of an object, we must extract 3D informations from a large set of cameras and this often leads to a time-expensive process. In this paper, we show how a divide and conquer approach can be used to speed up the entire process guaranteeing a good precision of the final model. In order to do this task, we chose to create a complete model of the interested object by linking together several surface patches reconstructed rapidly by a graph-cuts algorithm. The linking process is accomplished by a mesh fusion algorithm based on a volume of fluid approach, using a volumetric function.

## 2. DEPTH-MAP RECONSTRUCTION

In this section we show how to reconstruct accurately a portion of the surface of the object present in the analyzed scene. This is the crucial step of the reconstruction process. In fact we will link the surface patches resulting from this step to obtain the final complete object by the fusion algorithm described in the next section. We have chosen to use the known *graph cuts* approach [1, 2, 3, 5], adapting it to the problem of depth map reconstruction. In the next paragraph we propose a short description of the energy minimization approach; after that we will show how to formulate the problem of depth map reconstruction in term of energy minimization.

### 2.1. Energy minimization approach

Our approach to depth map reconstruction is similar to some recent work that give strong results for stereo matching and image restoration. It is well known that both problems can be elegantly stated in term of energy minimization [2, 3]. In the last few years powerful energy minimization algorithms have been developed based on graph cuts [2, 4, 6]. This methods are fast enough to be practical, but unlike simulated annealing, graph cuts methods cannot be applied to an arbitrary function. In this paper we will use some recent results [3] that give graph constructions for a quite general class of energy functions. The energy minimization formalism has several advantages. It allows a clean specification of the problem to be solved, as opposed to the algorithm used to solve it. In addiction, energy minimization naturally allows the use of soft constraints, such as spatial coherence. In an energy minimization framework, it is possible to cause ambiguities to be resolved in a manner that leads to a spatially smooth answer. Finally, energy minimization avoids being trapped by early hard decision.

### 2.2. Problem formulation

Suppose we are given $n$ calibrated images of the same scene taken from different viewpoints (or at different moments of time). Let assume a camera as the preferred one and let $\mathcal{P}$ be the set of pixels of the corresponding image. A pixel $p \in \mathcal{P}$ corresponds to a ray in 3D-space. Consider the first intersection of this ray with an object in the scene. Our goal is to find the depth of this point for all the pixel of the preferred image. So we want to find a labeling $f : \mathcal{P} \rightarrow \mathcal{L}$ where $\mathcal{L}$ is a discrete set of labels corresponding to increasing depths

from the preferred camera. Equivalently, we want to obtain the *depth map* of the pixels in the preferred image.

A pair $\langle p, l \rangle$ where $p \in \mathcal{P}$, $l \in \mathcal{L}$ corresponds to some point in 3D-space. We will refer to such points as *3D-points*. We define our energy function as consisting of two terms:

$$E(f) = E_{data}(f) + E_{smooth}(f)$$

In their work, Kolmogorov and Zabih [1] formulate the problem of scene reconstruction in a slightly different manner that permits to obtain a depth map for every image in the input set by an energy minimization approach. This leads to a computational expensive algorithm whose result is a unorganized clouds of point representing the surface of the visible part of the scene to reconstruct. Moveover, to have an effective reconstruction from the input set, cameras must respect some particular restrictive configuration, whereas with our definition we can treat a very large number of camera configurations without distinctions. It can be also noted that in our approach it is no long necessary the visibility term defined in [1]. In fact, assuming that the set of label corresponds to the increasing depths from the preferred camera, there cannot exists occluding pixels in the same image and consequently it is no long necessary for a visibility term. Moreover, also the other term are quite different. Our data term is defined as follow:

$$E_{data}(f) = \sum_{p \in \mathcal{P}} D(p)$$

where $D(p)$ is a non-positive value which results from the differences in intensity between corresponding pixels. $D(p)$ is computed for every pixel of the preferred image (we indicate this image with the index $j$) by this steps:

1. from $p$, we get the corresponding 3D-point by retroprojecting it from the preferred camera center of projection with the selected depth and then we project this 3D-point on each other calibrated image obtaining a set of $n - 1$ corresponding pixels $\{q_1, q_2, \ldots, q_i, \ldots, q_n | i \neq j\}$.
   The window dimensions usually range from $3 \times 3$ to $7 \times 7$: bigger windows are useful for poorly textured surfaces but become inefficient when cameras in the set present wide differences in viewpoint and rotations, their extension must then be chosen accordingly to the acquisition set.

2. on every non-preferred image we compute the SSD (Sum of Square Difference) using a square window centered on $q_i$ and the one centered on $p$, obtaining the set of values $\{d_1, d_2, \ldots, d_i, \ldots, d_n | i \neq j\}$.

3. finally we evaluate the energy data term for the $p$ point as follows:

$$D(p) = min(0, \sum_{\substack{i = 1 \\ i \neq j}}^{n} d_i - K) \qquad (1)$$

where $K$ is a positive constant large enough to capture significant variation of the SSD function (a typical value is $K = 30$).

The smoothness term is quite similar to the one used in [1] and its goal is to make neighboring pixels in the preferred image tend to have similar depths. The smoothness term is defined as follow:

$$E_{smooth}(f) = \sum_{\{p,q\} \in \mathcal{N}} V_{\{p,q\}}(f(p), f(q)) \qquad (2)$$

This term involves a notion of neighborhood: we assume that there is a neighborhood system on pixel

$$\mathcal{N} \subset \{\{p, q\} \mid p, q \in \mathcal{P}\}$$

This can be the usual 4-neighborhood system: pixels $p = (p_x, p_y)$ and $q = (q_x, q_y)$ are neighbors if they are in the same image and $|p_x - q_x| + |p_y - q_y| = 1$.

In [1], the function $V_{\{p,q\}}$ assumes the following form:

$$V_{\{p,q\}}(l_p, l_q) = \begin{cases} U_{\{p,q\}} & \text{if } l_p \neq l_q \\ 0 & \text{otherwise} \end{cases} \qquad (3)$$

where the $U_{\{p,q\}}$ is the following non-decreasing function:

$$U_{\{p,q\}} = \begin{cases} 3\lambda & \text{if } \Delta I(p, q) < 5 \\ \lambda & \text{otherwise} \end{cases} \qquad (4)$$

Where $\Delta I(p, r)$ is the average of values $|Intensity(p) - Intensity(r)|$ for all three bands(R,G,B). To make the reconstruction smooth while preserving discontinuities, we choose to follow a particular strategy in the use of the smoothness term. In fact, it is known that graph cuts techniques often yields flat and blocky results. This may not be important for disparity maps, but is crucial for shape reconstruction. To avoid this problem, we make a first cycle of the reconstruction algorithm with a limited set of labels, in order to reach rapidly a value of the energy near to the local minimum that could be got at convergence with the original algorithm. This corresponds to a good approximation of the position of the 3D-points, that can be improved with a second cycle at double resolution where we change the function $V_{\{p,q\}}$ defined in (3) with this new function:

$$\hat{V}_{\{p,q\}}(l_p, l_q) = \begin{cases} U_{\{p,q\}} & \text{if } |l_p - l_q| > z\_threshold \\ 0 & \text{otherwise} \end{cases}$$
$$(5)$$

In fact, this function relaxes the penalty mechanism of the smoothness term, giving a 0 penalty not only to the neighboring pixels that lie at the same depth but also to the ones that stay sufficiently near one another. The idea is supported by the fact that after the first cycle of the algorithm, only some of the pixels are approximatively well positioned in 3D-space by the consistency measure given by the data term, while the other are positioned only by the effect of the smoothness term which forces them to lie on the same level of neighboring pixel, resulting in flat blocks. Thus, relaxing the constraint imposed by the first smoothness term, neighboring pixels have greater chance to occupy adjacent depths correctly.

## 2.3. Graph cuts Algorithm

Thanks to our energy redefinition the results obtained from the standard graph cuts algorithm (as defined in [1]) are much more accurate. As shown in the next paragraph further depth map optimization guarantees an high fidelity to the reconstructed data.

## 2.4. Depth map optimization

Even though the graph cuts algorithm is able to reconstruct an accurate depth map, it works only with a limited set of depths and, thus, it introduce a considerable quantization error in the position of each 3D-point. To overcome this problem, it is necessary an optimization step which yields the depth map more regular. The output of this process is a new depth map, where the discontinuities are preserved while the other parts become smoothed. To do this work, we consider the depth map as a functions of two variables defined on the preferred image and we apply a sequence of bidimensional filters on it. In particular, we start with a median filter to eliminate possible outliers and then we apply a dithering technique: some white noise is added to the depth function and, then, a low pass filter is used to yields the depth map smooth. To preserve discontinuities, the bidimensional low pass filter keeps the information needed from the neighbors of a pixel only if the depth distance is below a certain threshold. The size of the filter windows and this threshold are empirically chosen on the basis of the current reconstruction.

## 3. MESH FUSION BY A VOLUMETRIC APPROACH

In order to create a complete model of the object to reconstruct we can think to melt together the several surface patches obtained with the previous graph cuts based method. We chose a volumetric representation of the scene that use a voxelset made of cubic voxels, with an approach similar to the already known volume-of-fluid technique. Each voxel

can assume a value in the $[-1, +1]$ interval. The entire voxelset can be seen as a volumetric function which represents the surface of the object as the zero levelset. Negative values of the function indicate the space inside the object while positive values stay for external space. Near the surface, each voxel assumes an intermediate value on the basis of its distance from the closer depth map. The algorithm starts initializing every voxel to the $-1$ value. By this way the volumetric function represents a solid block where subsequent steps will carve the surface of the object. At this point, we select a depht map and assign a value to each voxel of the voxelset with the steps explained in figure 1 using a bilinear interpolation between neighboor points in the depth map. The following criterion is followed: a voxel value can only be changed with a greater one.
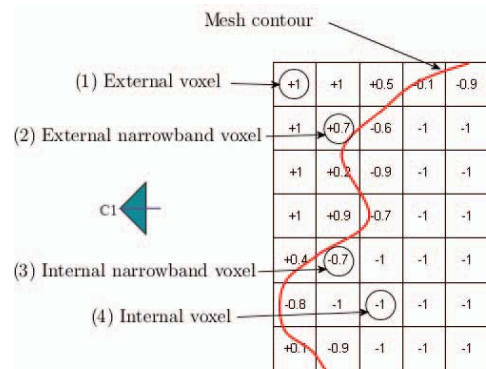


**Fig. 1**. Modelling of the volumetric function from a mesh.

Repeating this steps for every mesh will lead to a volumetric function whose zero leveset locates the object surface. The resulting object can be seen as a sort of convex hull obtained by linking together the meshes and taking only the part of the 3D space contained in their intersection.

## 4. EXPERIMENTAL RESULTS

The proposed algorithm has been applied to a set of images of a synthetic teapot and then to a set of images of a real object, a skull, acquired with a trinocular calibrated camera system. The teapot has been modeled by a 3D software and several snapshots have been rendered from it. The skull has been located on a turntable and a sequence of snapshots has been taken for every position of the turntable. For both the object, some images have been selected in triplets. From each triplet a depth map is reconstructed. Figure 2(a) shows a triplet of images of the teapot, while figure 2(b) shows the corresponding reconstructed surface patch; the complete object reconstructed by the volumetric algorithm is visualized in figure 2(c). Analogously, figure 3(a) shows a triplet of images of the skull, while figure 3(b) shows the corre-

sponding surface patch; the complete model of the skull is shown in figure 3(c). Both the final teapot and skull model have been obtained by melting together eight surface patches taken from different position. The parameters K of equation (1) and $\lambda$ of equation (4) are determined heuristically: optimal values depend on the images we are processing. The parameters can be varied to gain some insight about the algorithm: for big values of $\lambda$ the smoothness dominates the correlation, resulting in a map with many flat blocks of pixels, whereas little values of $\lambda$ yields to an irregular depth map with many wrong discontinuities. In our experiment, we chose the values $K = 30$ and $\lambda = 5$.

## 5. CONCLUSIONS

3D reconstruction from a set of images is a critical process. In order to perform this task we presented a reconstruction algorithm based on graph cuts theory. We have defined an energy function whose minimum represents the solution to our problem and we implemented a technique to raffinate the obtained depth maps. A virtue of this approach is the algorithm speed. In fact, we chose to build up a complete model of an object linking together several depth maps, reducing the computational effort either in the time needed and in the memory space required for reconstruct each of them. A volumetric approach has been used to do this task.

## 6. ACKNOLEDGEMENTS

## 7. REFERENCES

[1] V. Kolmogorov, R. Zabih. "Multi-camera Scene Reconstruction via Graph Cuts". *In European Conference on Computer Vision*, 2002.

[2] V. Kolmogorov, R. Zabih. "Computing Visual Corrispondence with Occlusion via Graph Cuts". *In International Conf. on Computer Vision*, 2001.

[3] V. Kolmogorov, R. Zabih. "What energy functions can be minimized via graph cuts?" *In European Conf. on Computer Vision*, 2002.

[4] Y. Boykov, O. Veksler, R. Zabih. "Fast Approximate Energy Minimization via Graph Cuts". *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 2001.

[5] D. Snow, P. Viola, R. Zabih. "Exact Voxel Occupancy with Graph Cuts". *In Proc. Computer Vision and Pattern Recognition Conf.*, 2000.

[6] Y. Boykov, O. Veksler, R. Zabih. "Markov Random Fields with efficient approximations". *IEEE Conf. on Computer Vision and Pattern Recognition*, 1998.
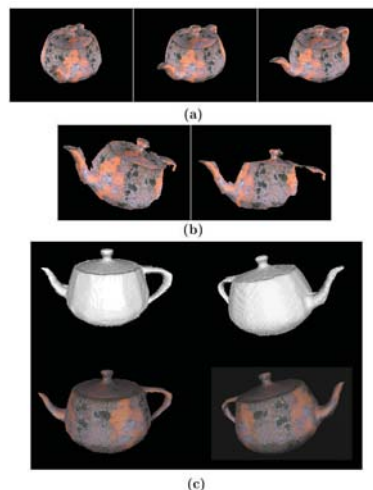
**Fig. 2**. (a) teapot image triplet. (b) corresponding surface patch (c) complete model of the teapot
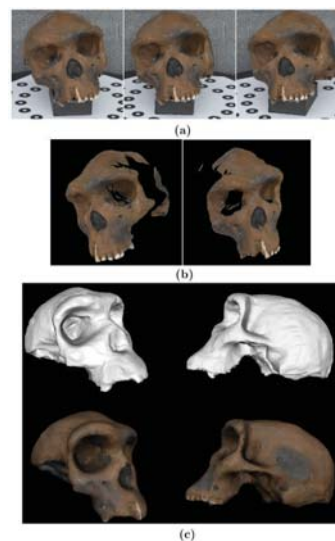


**Fig. 3**. ((a) skull image triplet. (b) corresponding surface patch (c) complete model of the skull