

Occlusions, Discontinuities, and Epipolar Lines in Stereo^{*}

Hiroshi Ishikawa and Davi Geiger

Department of Computer Science, Courant Institute of Mathematical Sciences
New York University, 251 Mercer Street, New York, NY 10012, U.S.A.
ishikawa@cs.nyu.edu, geiger@cs.nyu.edu

Abstract. Binocular stereo is the process of obtaining depth information from a pair of left and right views of a scene. We present a new approach to compute the disparity map by solving a global optimization problem that models occlusions, discontinuities, and epipolar-line interactions.

In the model, geometric constraints require every disparity discontinuity along the epipolar line in one eye to *always* correspond to an occluded region in the other eye, while at the same time encouraging smoothness across epipolar lines. Smoothing coefficients are adjusted according to the edge and junction information. For some well-defined set of optimization functions, we can map the optimization problem to a maximum-flow problem on a directed graph in a novel way, which enables us to obtain a global solution in a polynomial time. Experiments confirm the validity of this approach.

1 Introduction

Binocular stereo is the process of obtaining depth information from a pair of left and right images, which may be obtained biologically or via a pair of cameras. The fundamental issues in stereo are: (i) how the geometry and calibration of the stereo system are determined, (ii) what primitives are matched between the two images, (iii) what a priori assumptions are made about the scene to determine the disparity, (iv) how the disparity map is computed, and (v) how the depth is calculated from the disparity.

Here we assume that (i) is solved, and hence the correspondence between epipolar lines (see Fig.1) in the two images are known. Answering question (v) involves determining the camera parameters, triangulation between the cameras, and an error analysis, for which we refer the reader to [10]. We focus on problems (ii), (iii), and (iv) in this paper.

Main contributions of this paper to these problems are summarized as follows:

- (ii) A stereo algorithm solely based on matching pixels from the left and right views tends to have difficulties precisely locating the discontinuities. To remedy this problem, we use intensity edges and junctions as cues for the depth discontinuities. The significance of junctions to stereo has been pointed out in [1, 16]. Our new approach uses edges and junctions in a uniform manner where “ordinary” pixels, edge pixels, and junction pixels increasingly suggest discontinuities in this order. Although

^{*} This work was supported in part by NSF under contract IRI-9700446.

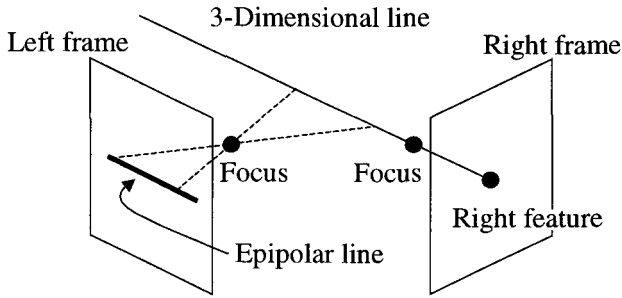


Fig. 1. A pair of frames (eyes) and an epipolar line in the left frame.

a use of window features can also accommodate for this problem, it requires another set of parameters to estimate their size and shape [15].

- (iii) Various algorithms, as in the cooperative stereo [17], have proposed a priori assumptions on the solution, including *smoothness* to bind nearby pixels and *uniqueness* to inhibit multiple matches. Occlusions and discontinuities must also be modeled to explain the geometry of the multiple-view image formation [3, 11]. Another aspect of stereo geometry is the interdependence between epipolar lines. This topic is often neglected because of a lack of optimal algorithms. We show that it is possible to account for all of these assumptions, including occlusions, discontinuities, and epipolar-line interactions, in computing the optimal solution. However, this result is under some restricted set of optimization functions and we discuss the scope of its applicability.
- (iv) Our method is a new use of the maximum-flow algorithm to find a globally optimal solution, which is represented by a cut of a *directed* graph that models a symmetric stereo (symmetric with respect to the left and right views.) We will describe the limitations imposed on the prior knowledge of surfaces to allow the use of the maximum-flow algorithm. More precisely, it can be shown that, to model interactions between epipolar lines, or discontinuities and occlusions along them, only convex functions can be used in the way described in this paper. Thus, for instance, the cost for discontinuities must be at least linear on the size of the discontinuities. Roy and Cox [22] introduced maximum-flow algorithms on *undirected* graphs for stereo, a more limited set than directed graphs, to compute disparity maps. They also claimed the approach to be a generalization of the dynamic programming algorithm. While our approach is also motivated by their work, we will show that their algorithm has several limitations in the ability to model stereo, e.g., no guarantee of ordering or uniqueness, no model for occlusions and discontinuities, and that it cannot be called a generalization of the dynamic programming algorithms in a strict sense.

1.1 Testing Stereo Algorithms

Evaluation of stereo theories and algorithms is a difficult task. Stereo must deliver accurate disparity maps. While obtaining real image results that “look” good is a definite

necessity, looking good is not enough to guarantee a high stereo quality. A careful examination of the disparity map is necessary.

One method of comparison would be to test all current stereo algorithms against each other. However, not only is it difficult to have access to most stereo algorithms, the criteria of comparisons are also unclear.

Alternatively, let us assume we have a ground-truth solution. Then we can compare the disparity map resulting from the algorithm against the ground truth by computing some measure of the error. Usually, ground truth is available only for synthetic examples. Synthetic examples can be created to capture geometrical scene properties, such as occlusions, discontinuities, and junctions, to maximize the information they provide about a stereo theory/algorithm. On the other hand, synthetic images generally lack realistic intensity information that accurately reflects such properties as textures, reflectivity, and illuminations. Thus, while it would take very sophisticated measures to study illumination and reflectivity-related issues with synthetic images, even the simplest synthetic examples using only black and white pixels can be quite instrumental to study the prior models of stereo geometry and the role of edges and junctions.

Two kinds of synthetic imagery are of special interest.

1. Random-dot stereograms (see Fig.7): They remove intensity considerations and essentially lack any features. The disambiguation and solution are derived by the prior model of surfaces. It is interesting to note that an occluded region in one of the pair of images is created by adding random dots to an empty region and, therefore, does not have a good correspondence in the other image.
2. Illusory Surfaces (see Fig.6): These are very important synthetic images to complement the random-dot stereograms, since they crucially require the study of edge and junction features as well as the way in which they relate to occlusions and discontinuities. The following issues deserve our attention: (a) the formation of discontinuities where no intensity edges are present, (b) the choice between front parallel panes with discontinuities and tilted planes, and (c) the role of epipolar-line interactions. In contrast to the random-dot stereograms, occluded regions are composed of pixels that in terms of feature match can have good correspondences, which however the geometrical constraints would not allow.

Finally, let us emphasize that the ultimate test of a stereo algorithm is its performance on real image pairs. We do address all these experimental issues in this paper.

1.2 Background and Comparison to Previous Work

A number of researchers, including Julesz [14]; Marr and Poggio [17]; Pollard, Mayhew and Frisby [21]; Grimson [13]; Kanade and Okutomi [15]; Ayache [2]; Roy and Cox [22], have addressed the problem of binocular stereo matching without explicitly modeling occlusions and its relation to discontinuities.

There is now abundant psychophysical evidence [1, 12, 19] that the human visual system does take advantage of the detection of occluded regions to obtain depth information. The earliest attempts to model occlusions and its relation to discontinuities, in Geiger, Ladendorf and Yuille [11], and independently in Belhumeur and Mumford [3],

had a limitation that they restrict the optimization function to account only for interactions along the epipolar lines. Malik [16] brought attention to the role of junctions in stereo, together with Anderson's experiments [1]. Our use of junctions in this paper is limited to the detection of discontinuities.

As for optimization, Roy and Cox [22] presented the study of epipolar-line interaction through the use of maximum-flow algorithm, where they claimed that the algorithm provides a generalization of the dynamic programming algorithm to the two-dimensional stereo. We point out that their algorithm does not model discontinuities and occlusions and that their algorithm obliges occlusion regions to have a match. Thus, for example, their algorithm would have difficulties with random-dot stereograms, where occluded regions do not have any good correspondence. Also, their method provides no guarantee of unique match, but almost always results in undesirable multiple matches, which it eliminates in an ad-hoc way. Moreover, contrary to their claim, the maximum-flow algorithms cannot be a generalization of the dynamic programming. We show that this way of using maximum-flow algorithm cannot model non-convex functions and, as a result, it cannot use sublinear penalties for discontinuities. Though our proposed maximum-flow algorithm is still limited in the same way, it guarantees the ordering constraint (or monotonicity of the solution) and uniqueness of the match, and can model occlusions, discontinuities, and epipolar-line interactions. When non-convex costs are needed, we propose the use of the dynamic programming algorithms, but at the expense of approximate solution when accounting for the epipolar-line interactions.

2 Matching and Surface Reconstruction

Here we describe our stereo model so that in the next section we can discuss the optimization issues.

2.1 Matching Features

We use gray-level pixels, edgeness, and cornerity as features. While we assume the error in the correspondence is solely based on gray-level values, the detection of discontinuities is helped by the intensity edges (edgeness) and, even more, by junctions (cornerity). That is, our model assumes that depth discontinuities are more likely to occur in the presence of edges, and of junctions even more.

If feature $I_{e,l}^L$ at pixel l on epipolar line e , or pixel (e,l) , in the left image matches feature $I_{e,r}^R$ at pixel (e,r) in the right image, then $\|I_{e,l}^L - I_{e,r}^R\|$ should be small, where $\|\cdot\|$ is some measure of feature distance. We assume a dense set of features, though it would be easy to extend the model to include sparse features like edges.

As in [11, 17], we use a matching process $M_{l,r}^e$ that is 1 if the feature at pixel (e,l) in the left eye matches the feature at pixel (e,r) in the right eye, and 0 otherwise. Given a pair of left and right image I^L and I^R , we define the input cost of the matching process M by

$$E_{\text{input}}(M | I^L, I^R) = \sum_{l,r,e} M_{l,r}^e \|I_{e,l}^L - I_{e,r}^R\|, \quad (1)$$

where $l = 0, \dots, N-1$ and $r = 0, \dots, N-1$ are indices that scan the left and right images along the epipolar lines $e = 0, \dots, N-1$ for $N \times N$ square images.

This model, in principle, can be derived from an image formation model. For example, in the case where $\|\cdot\|$ is the Euclidean norm, (1) assumes that the pixels $I_{e,l}^L$ and $I_{e,r}^R$ are related by the relation $I_{e,l}^L = I_{e,r}^R + x$ for corresponding points l and r , where x is a random variable distributed with $P(x) = e^{-x^2} / \sqrt{2\pi}$.

2.2 Uniqueness, Occlusion and Disparity

In order to prohibit multiple matches, we impose a *uniqueness* constraint:

$$\sum_{l=0}^{N-1} M_{l,r}^e \leq 1 \quad \text{for all } r \quad \text{and} \quad \sum_{r=0}^{N-1} M_{l,r}^e \leq 1 \quad \text{for all } l. \quad (2)$$

Notice that this guarantees that there can be at most one match per feature, while also allowing unmatched features to exist.

Remark: It is important to account for scenes composed of tilted planes. In the discretized setting, image pairs of tilted planes exhibit less pixels in one image than in the other. This property will force the pairing to break the uniqueness constraint. When thinking in a continuous setting, the concept of uniqueness is not broken by tilted planes, but in the discrete setting it is. When we describe the mapping of the model to the maximum-flow problem, we will account for multiple matching in the presence of tilted plane.

Occlusion and Disparity: For a stereoscopic image pair we define occlusions to be regions in one image that have no match in the other image. These may occur as a result of occlusions in the 3-D scene (see Fig.2.) We first define an occlusion field $O_{e,l}^L$ and $O_{e,r}^R$ as

$$O_{e,l}^L(M) = 1 - \sum_{r=0}^{N-1} M_{l,r}^e \quad \text{and} \quad O_{e,r}^R(M) = 1 - \sum_{l=0}^{N-1} M_{l,r}^e .$$

Due to uniqueness constraint (2), the occluded pixels are the ones with this field 1, when no matches occur.

Now, we define a disparity field $D_{e,l}^L$ and $D_{e,r}^R$ as

$$D_{e,l}^L(M) = \sum_{r=0}^{N-1} M_{l,r}^e (r-l) \quad \text{if} \quad O_{e,l}^L = 0 \quad \text{and}$$

$$D_{e,r}^R(M) = \sum_{l=0}^{N-1} M_{l,r}^e (r-l) \quad \text{if} \quad O_{e,r}^R = 0 ,$$

for unoccluded pixels and then linearly interpolate for occluded pixels, assuming the boundary condition $D_{e,-1}^L(M) = D_{e,N}^L(M) = D_{e,-1}^R(M) = D_{e,N}^R(M) = 0$. Then, it can be shown that the following holds:

$$\sum_{l=0}^{N-1} D_{e,l}^L(M) = \sum_{r=0}^{N-1} D_{e,r}^R(M) .$$

Since these two variables $O(M)$ and $D(M)$ are functions of the matching process M , our model is completely determined by M .

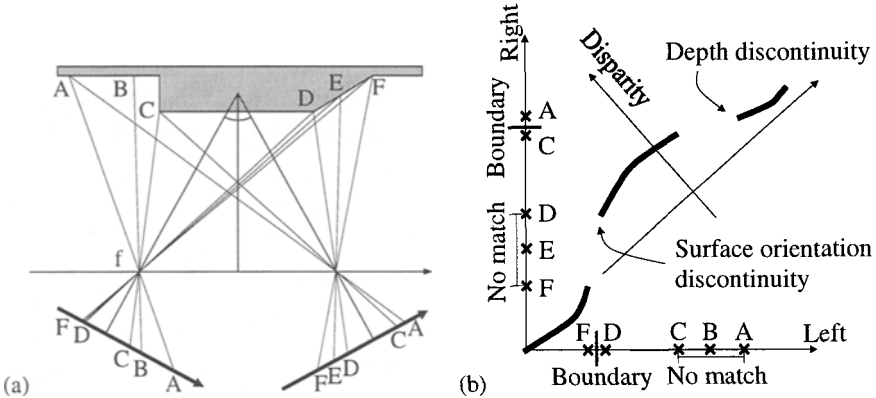


Fig. 2. (a) A polyhedron (shaded area) with self-occluding regions and with a discontinuity in the surface-orientation at feature D and a depth discontinuity at feature C. (b) A diagram of left and right images (1D slice) for the image of the ramp above. Notice that occlusions always correspond to discontinuities. Dark lines indicates where match occurs

2.3 The Monotonicity / Ordering constraint

The *monotonicity* constraint is a variant of the ordering constraint. It differs slightly from the standard ordering constraint because it requires neighboring points to match.

The monotonicity constraint is defined as follows:

Monotonicity Constraint: For every match (e, l, r) such that $M_{lr}^e = 1$,

$$M_{N_{e,l}^L, N_{e,r}^R}^e = 1 \tag{3}$$

holds, where

$$N_{e,l}^L = \max(\{l' \mid l' < l, O_{e,l'}^L = 0\})$$

$$N_{e,r}^R = \max(\{r' \mid r' < r, O_{e,r'}^R = 0\}) .$$

We call the match $N_{l,r}^e = (e, N_{e,l}^L, N_{e,r}^R)$ the *neighbor match* to match (e, l, r) .

The criteria to choose the optimal matching are based on the prior knowledge of surfaces as well as on the data matching term. We will formalize it next.

2.4 Surface Reconstruction

We now specify a prior model for surfaces in the scene. Depth changes are usually small compared to the viewer distance, except at depth discontinuities. Thus, we would like surfaces to be smooth, yet like to allow for large depth changes to occur. Since there is a simple trigonometric relation between disparity and depth, we consider the same constraint to hold for disparities. This can be modeled by some cost function E_{surface} depending upon the disparity change $\nabla_v D(M)$ and $\nabla_h D(M)$. We define the disparity change $\nabla_v D(M)$ along the epipolar line, which is defined only at unoccluded pixels, as follows:

$$\nabla_v D_{e,l,r}^e = D_{e,l}^L - D_{e,N_{e,l}^L}^L = (l - N_{e,l}^L) + (r - N_{e,r}^R) = D_{e,r}^R - D_{e,N_{e,r}^R}^R ,$$

and we define the terms as:

$$\nabla_v D_{e,l}^L = l - N_{e,l}^L \quad \nabla_v D_{e,r}^R = r - N_{e,r}^R .$$

Next, we define the disparity change $\nabla_h D(M)$ across the epipolar line as:

$$\begin{aligned} \nabla_h D_{e,l}^L &= D_{e,l}^L - D_{e-1,l}^L , \text{ and} \\ \nabla_h D_{e,r}^R &= D_{e,r}^R - D_{e-1,r}^R . \end{aligned}$$

We can now define the cost function E_{surface} as

$$\begin{aligned} E_{\text{surface}}(M | I^L, I^R) &= \sum_{e,l,r} M_{e,l,r}^e \left\{ \mu_{e,l,r}^{vL} F(\nabla_v D_{e,l}^L) + \mu_{e,l,r}^{vR} F(\nabla_v D_{e,r}^R) \right. \\ &\quad \left. + \mu_{e,l,r}^{hL} F(\nabla_h D_{e,l}^L) + \mu_{e,l,r}^{hR} F(\nabla_h D_{e,r}^R) \right\} , \quad (4) \end{aligned}$$

where the function $F(x)$ penalizes for the size of the disparity changes x . The parameters μ^{vL} and μ^{hL} control the smoothing in the left image along and across epipolar lines, respectively. Analogously, μ^{vR} and μ^{hR} control the smoothing in the right image along and across epipolar lines.

These smoothing parameters vary according to intensity edges and junction information. The edge information affects the parameters as

$$\begin{aligned} \mu_{e,l,r}^{vL} &\propto \frac{\gamma}{\gamma + (\Delta_v I_{e,r}^R)^2} , & \mu_{e,l,r}^{hL} &\propto \frac{\gamma}{\gamma + (\Delta_h I_{e,r}^R)^2} , \\ \mu_{e,l,r}^{vR} &\propto \frac{\gamma}{\gamma + (\Delta_v I_{e,l}^L)^2} , & \mu_{e,l,r}^{hR} &\propto \frac{\gamma}{\gamma + (\Delta_h I_{e,l}^L)^2} . \end{aligned}$$

where γ is a parameter to be estimated, and $\Delta_v I_{e,l}^L = I_{e,l}^L - I_{e,l-1}^L$, $\Delta_h I_{e,l}^L = I_{e,l}^L - I_{e-1,l}^L$, $\Delta_v I_{e,r}^R = I_{e,r}^R - I_{e,r-1}^R$, and $\Delta_h I_{e,r}^R = I_{e,r}^R - I_{e-1,r}^R$.

Large image gradients in one image reduce the smoothing coefficients and facilitate discontinuities in the other image to occur. To further reduce the cost of discontinuities occurring at junctions (without affecting the edge information,) we write

$$\begin{aligned} \mu_{e,l,r}^{vL} &= \mu \frac{\gamma}{\gamma + (\Delta_v I_{e,r}^R)^2} \left\{ \frac{\gamma}{\gamma + (\Delta_h I_{e,r}^R)^2} + \frac{\gamma}{\gamma + (\Delta_h I_{e,r-1}^R)^2} \right. \\ &\quad \left. + \frac{\gamma}{\gamma + (\Delta_h I_{e-1,r}^R)^2} + \frac{\gamma}{\gamma + (\Delta_h I_{e-1,r-1}^R)^2} \right\} , \end{aligned}$$

so that corners have even lower cost for discontinuities along either direction. We define the other smoothing parameters analogously.

The final optimization cost is then given by

$$E(M | I^L, I^R) = E_{\text{input}}(M | I^L, I^R) + E_{\text{surface}}(M | I^L, I^R) . \quad (5)$$

We can see this cost as the energy of a Gibbs probability distribution

$$P_{\text{stereo}}(M | I^L, I^R) = \frac{1}{Z} e^{-E(M | I^L, I^R)} ,$$

where Z is a normalization constant.

Observe that our theory, given by (1) and (4), is symmetric with respect to the two eyes. This is necessary to ensure that the full information can be extracted from occlusions. Moreover, the monotonicity constraint (3) will be applied as a hard constraint to simplify the optimization of the cost.

We will show in Sect.3 that the maximum-flow algorithm naturally lands to this theory but it restricts the function $F(x)$ to be a convex function. To account for this limitation, we will consider $F(x) = |x|$, since, among the convex functions, it penalizes least for large $|x|$.

3 Mapping the Optimization Problem to a Maximum-flow Problem

In this section, we explain the stereo-matching architecture that utilizes the maximum-flow algorithm to obtain the globally optimal matching, with respect to the energy (5), between left and right image.

3.1 The Directed Graph

We devise a directed graph and let a cut represent a matching so that the minimum cut corresponds to the optimal matching. The formulation explicitly handles the occlusion and is completely symmetric with respect to left and right, up to the reversal of all edges, under which the solution is invariant.

Let \mathcal{M} be the set of all possible matching between pixels, i.e., $\mathcal{M} = \{(e, l, r) \mid e, l, r \in [0, \dots, N - 1]\}$. We define a directed graph $G = (V, E)$ as follows:

$$\begin{aligned} V &= \{u_{lr}^e \mid (e, l, r) \in \mathcal{M}\} \cup \{v_{lr}^e \mid (e, l, r) \in \mathcal{M}\} \cup \{s, t\} \\ E &= E_M \cup E_C \cup E_P \cup E_E. \end{aligned}$$

In addition to the source s and the sink t , the graph has two vertices u_{lr}^e and v_{lr}^e for each possible matching $(e, l, r) \in \mathcal{M}$. The set E of edges is divided into subsets E_M , E_C , E_P , and E_E , each associated with a capacity with a precise meaning in terms of the model (5), which we will explain in the following subsections.

We denote a directed edge from vertex u to vertex v as (u, v) . Each edge (u, v) has a nonnegative capacity $c(u, v) \geq 0$. A *cut* of G is a partition of V into subsets S and

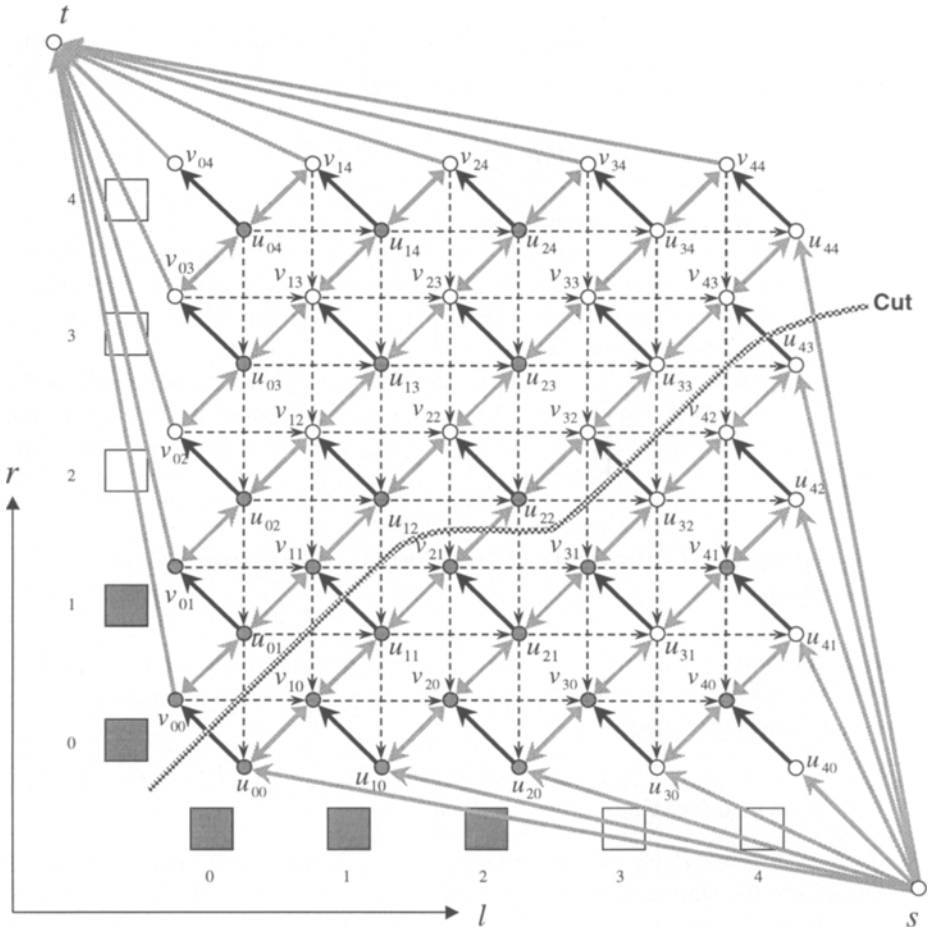


Fig. 3. An epipolar slice of the graph representing the stereo model, where we perform a maximum-flow algorithm. The full graph is naturally represented as three dimensional, with the third axis parametrizing epipolar lines. Edges are given by pairs of graph nodes (u, v) . A cut of the graph can be thought of as a surface that separates the two parts, and restrict to a curve (a path) in an epipolar slice. The optimal cut is the one that minimizes the sum of the capacities associated with the cut edges. In this example, we illustrate a cut (solution) that yields the matches $(l, r) = (0, 0), (1, 1), (3, 2)$, and $(4, 3)$ and assign an occlusion to gray (white) pixel 2 (4) in the left (right) image. In the following figures we explain each type of edge/capacity.

$T = V \setminus S$ such that $s \in S$ and $t \in T$ (see Fig.3). When two vertices of an edge (u, v) is separated by a cut with $u \in S$ and $v \in T$, we say that the edge is cut. This is the only case that the cost $c(u, v)$ of the edge contributes to the total cost, i.e., if the cut is through the edge (u, v) with $u \in T$ and $v \in S$, the cost is $c(v, u)$, which is in general different from $c(u, v)$. It is well known that by solving a maximum-flow problem one can obtain a *minimum cut*, a cut that minimizes the total cost $\sum_{u \in S, v \in T} c(u, v)$ over all cuts. (See [9] for details.)

Note that, in this formulation, we use a directed graph while in the method of [22] it suffices to use an undirected graph. The use of the directed graph is crucial to enforce the uniqueness constraint (2) and the ordering/monotonicity constraint (3).

Let us now explain each set of edges E_M , E_C , E_P , and E_E .

3.2 Matching Edges

Each pair of vertices are connected by a directed edge (u_{lr}^e, v_{lr}^e) with a capacity $\|I_{e,l}^L - I_{e,r}^R\|$. This edge is called a *matching edge* and we denote the set of matching edges by E_M :

$$E_M = \{(u_{lr}^e, v_{lr}^e) \mid (e, l, r) \in \mathcal{M}\} .$$

If a matching edge (u_{lr}^e, v_{lr}^e) is cut, we interpret this to represent a match between pixels (e, l) and (e, r) , i.e., $M_{l,r}^e = 1$. Thus, the sum of the capacities associated with the cut matching edges is exactly E_{input} in (5), which is defined in (1). Figure 3 shows the nodes and matching edges on an epipolar line. The cut shown represents a match $\{(l, r)\} = \{(0, 0), (1, 1), (3, 2), (4, 3)\}$.

Note that pixel 2 in the left image has no matching pixel in the right image, as well as pixel 4 in the right image, that is, these pixels are occluded. This is how the formulation represents occlusions and discontinuities, whose costs are accounted for by *penalty edges*.

3.3 Penalty Edges (Discontinuity, Occlusions, and Tilts)

Penalty edges are classified in four categories:

$$\begin{aligned} E_P &= E_L \cup E'_L \cup E_R \cup E'_R \\ E_L &= \{(v_{lr}^e, u_{l(r+1)}^e) \mid (e, l, r) \in \mathcal{M}, r < N - 1\} \cup \\ &\quad \{(s, u_{l0}^e) \mid e, l \in [0 \dots N - 1]\} \cup \\ &\quad \{(v_{l(N-1)}^e, t) \mid e, l \in [0 \dots N - 1]\} \\ E'_L &= \{(u_{l(r+1)}^e, v_{lr}^e) \mid (e, l, r) \in \mathcal{M}, r < N - 1\} \\ E_R &= \{(v_{lr}^e, u_{(l-1)r}^e) \mid (e, l, r) \in \mathcal{M}, l > 0\} \cup \\ &\quad \{(s, u_{(N-1)r}^e) \mid e, r \in [0 \dots N - 1]\} \cup \\ &\quad \{(v_{0r}^e, t) \mid e, r \in [0 \dots N - 1]\} \\ E'_R &= \{(u_{(l-1)r}^e, v_{lr}^e) \mid (e, l, r) \in \mathcal{M}, l > 0\} \end{aligned}$$

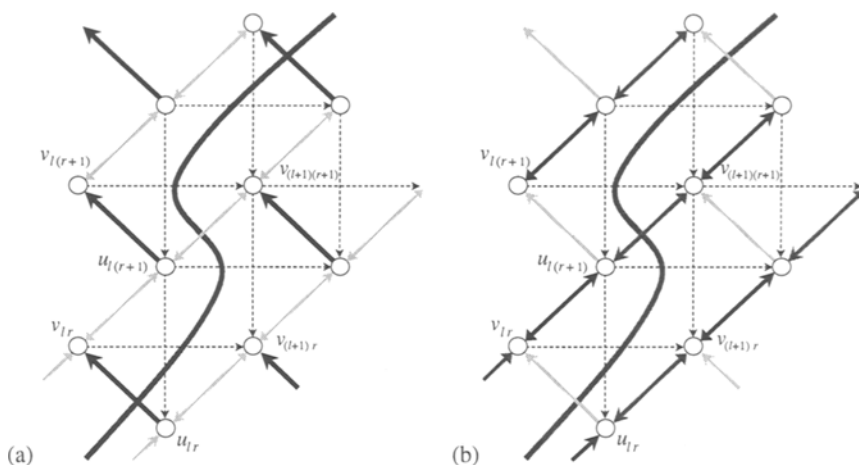


Fig. 4. (a) A *matching edge* is here represented by a black arrow, and we show a cut through a matching edge (u_{lr}^e, v_{lr}^e) . This capacity, given by $\|I_{e,l}^L - I_{e,r}^R\|$, provides the cost for the first term of the stereo-energy cost. (b) *Penalty edges* are represented by dark arrows. Every time a cut crosses these edges the cost is given by its capacity, which are given by either $\mu_{e,l,r}^{vL}$ or $\mu_{e,l,r}^{vR}$. By crossing consecutive penalty capacities, the cost is added linearly, accounting for the function $F(x) = |x|$. Here, the cut crosses the penalty edge $(v_{(l+1)(r+1)}^e, u_{l(r+1)}^e)$ with cost $\mu_{e,l+1,r+1}^{vR}$, accounting for the occlusion of pixel $r + 1$ on the right image.

These edges are for paying for discontinuities and occlusions. Edges in E_L are cut whenever a pixel in the left image has no matching pixel in the right image. For instance, if pixel (e, l) in the left image has no match, exactly one of the edges of the form $(v_{lr}^e, u_{l(r+1)}^e)$, (s, u_{l0}^e) , or $(v_{l(N-1)}^e, t)$ is cut (see Fig.4(b).) By setting the capacity for these edges, we control the smoothness along the epipolar line. Similarly, an edge in E_R corresponds to an occlusion in the right image. The penalty is given by the smoothing coefficients $\mu_{e,l,r}^{vL}$ for a edge $(v_{l(r-1)}^e, u_{lr}^e)$ in E_L and $\mu_{e,l,r}^{vR}$ for a edge $(v_{lr}^e, u_{l(r-1)}^e)$ in E_R . By crossing consecutive penalty capacities the cost is added linearly, accounting for the function $F(x) = |x|$.

Edges in E_L' are cut when a pixel in the left image matches two or more pixels in the right image. This corresponds to a tilted surface. Since every edge $(u_{l(r+1)}^e, v_{lr}^e)$ in E_L' has an opposite edge $(v_{lr}^e, u_{l(r+1)}^e)$ in E_L , by setting the capacity of $(u_{l(r+1)}^e, v_{lr}^e)$ higher or lower than the capacity of $(v_{lr}^e, u_{l(r+1)}^e)$, we can favor or disfavor tilted surface solution over occlusion/discontinuity solution. To strictly enforce the uniqueness constraint (2), we can make the capacity of the edge infinity, which cannot be done in an undirected graph.

F(x) must be convex: We briefly outline a proof that $F(x)$ has to be convex. We assume a uniformity of the cost, i.e., $F(x)$ is the same everywhere in the system. The cost of a cut is the sum of the capacities of the edges crossed by the cut. This sum can only increase the cost by the amount of each capacity, since the capacity is non-negative. The size of the discontinuity defines the minimum number of edges that is crossed,

and as the graph becomes more connected, the number of edges only increases, adding up the cost of a discontinuity more. The sum guarantees at least a linear cost, for the minimum connectivity used here, which represents the function $F(x) = |x|$. This proves that we cannot use non-convex functions as $F(x)$ and in particular cannot penalize for discontinuities or oclusions sublinearly.

In a graph with more connectivity, this cost can grow faster than linearity and so other convex functions can be created. Our interest is to have costs that increase the least with the discontinuity size and hence to have as few edges as possible, which also is helpful to improve the efficiency of the algorithm. This is a limitation of the use of the maximum-flow algorithm that is not present in dynamic programming algorithms (e.g., [11]). While there is this limitation, our experiments do yield good solution with the linear cost $F(x) = |x|$.

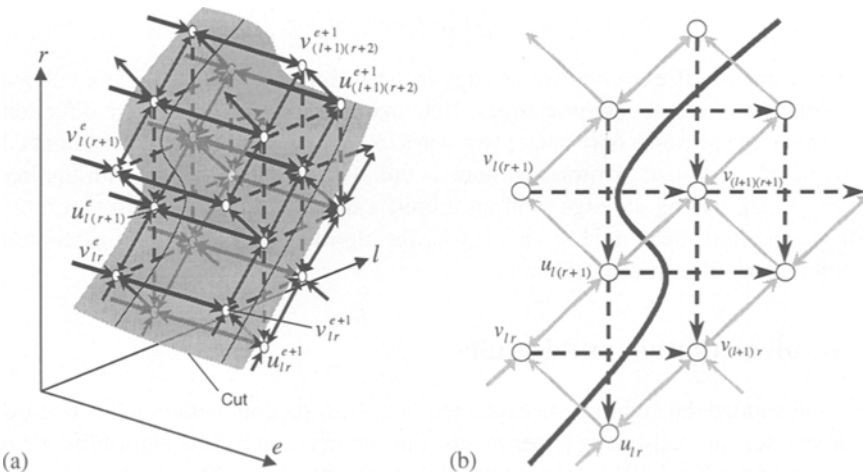


Fig. 5. (a) This figure shows a piece of the full three dimensional graph. Epipolar edges account for the epipolar-line interactions. They are represented by thick arrows. A cut is shown through two epipolar edges. (b) Constraint edges are depicted as dashed arrows. They enforce the monotonicity or ordering constraint (3), i.e., the edge capacity are set to infinity so that the optimal (minimum) cut does not cut them.

3.4 Epipolar edges

Epipolar edges are the only edges across epipolar lines. They simply connects vertices with the same (l, r) in both directions:

$$\begin{aligned}
 E_E = & \{(u_{lr}^e, u_{lr}^{e+1}) \mid (e, l, r) \in \mathcal{M}, e < N - 1\} \cup \\
 & \{(u_{lr}^{e+1}, u_{lr}^e) \mid (e, l, r) \in \mathcal{M}, e < N - 1\} \cup \\
 & \{(v_{lr}^e, v_{lr}^{e+1}) \mid (e, l, r) \in \mathcal{M}, e < N - 1\} \cup \\
 & \{(v_{lr}^{e+1}, v_{lr}^e) \mid (e, l, r) \in \mathcal{M}, e < N - 1\} .
 \end{aligned}$$

The capacity of an epipolar edge controls the smoothness of the solution across epipolar lines. In one extreme, if the cost is zero, the matching on epipolar lines have no influence to each other. In the other extreme, if the cost is infinity, all epipolar lines must have the same matching (see Fig.5(a).)

3.5 Constraint Edges

Constraint edges are for enforcing the monotonicity constraint (3) and defined as follows:

$$E_C = \{(u_{lr}^e, u_{(l+1)r}^e) \mid (e, l, r) \in \mathcal{M}, l < N - 1\} \cup \\ \{(u_{lr}^e, u_{l(r-1)}^e) \mid (e, l, r) \in \mathcal{M}, r > 0\} \cup \\ \{(v_{lr}^e, v_{(l+1)r}^e) \mid (e, l, r) \in \mathcal{M}, l < N - 1\} \cup \\ \{(v_{lr}^e, v_{l(r-1)}^e) \mid (e, l, r) \in \mathcal{M}, r > 0\} .$$

The capacity of each constraint edge is set to infinity. Therefore, any cut with a finite total flow cannot cut these edges. Note that, because the edges have directions, a constraint edge prevents only one of two ways to cut them. In Fig.5(b), constraint edges are depicted as dashed arrows, and none is cut. This cannot be done with undirected graphs, where having an edge with an infinite capacity is equivalent to merging two vertices, and thus meaningless. This is why the algorithm in [22] cannot guarantee the monotonicity/ordering constraint.

4 Implementation and Results

We implemented the architecture explained in the last section, with window features of small size for the real-image experiments. For the maximum-flow algorithm, we used the standard push-relabel method with global relabeling [8]. It took about 1 hour on 266Mhz Pentium II machine to compute the results on real images (512×512 pixels).

To see how the algorithm handles occlusions and discontinuities, we experimented with both random-dot stereogram and illusory surface images. These experiments were particularly important on the early phases of the work, when we needed to understand the role of each kind of edges and capacities in the graph. For instance, letting disparity discontinuity more likely to be present at intensity edges and junctions, as well as the epipolar-line interactions, was crucial for the performance of the algorithm on the illusory figures. The experiments on random-dot stereograms were important as they made it apparent that occluded regions do not match any region, and therefore having multiple matches at these points as the algorithm in [22] does cannot possibly be right. Results are shown in Fig.6 and Fig.7. In both cases, the program successfully found the occluded regions in both left and right images. The correspondence diagram in Fig.7 shows the symmetric nature of our result.

We also tested the algorithm on real images. Figure 8 shows the result for images Pentagon and Fruits. The results are in support of the model – they “look good”. Notice that the small variance of the height of the ground around the Pentagon can clearly be seen.

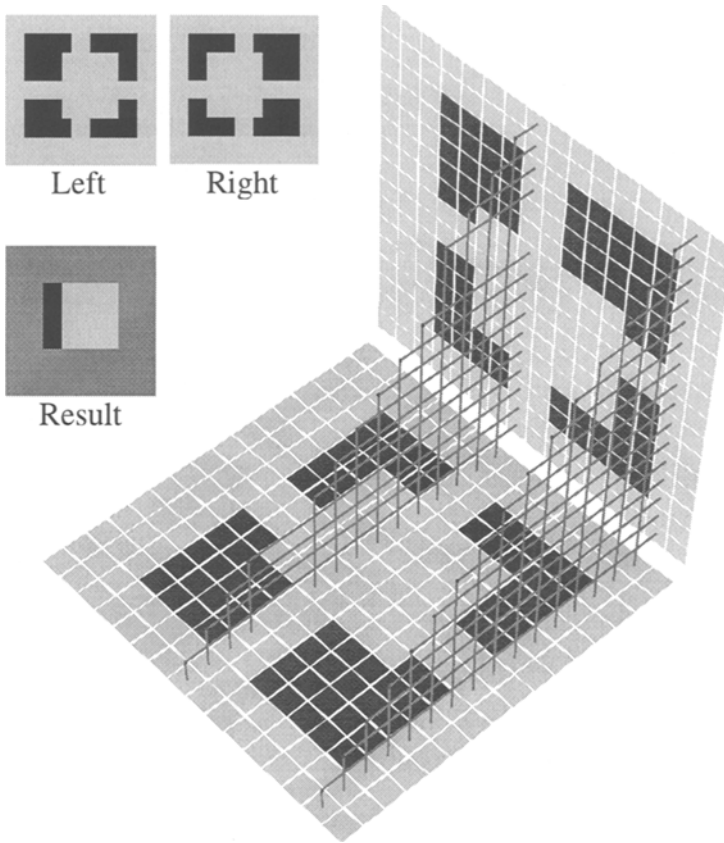


Fig. 6. Illusory surface. Shown as result is the disparity map for the left image. The diagram on the right shows the correspondence between left and right pixels for two epipolar lines. Note the symmetric nature of the result.

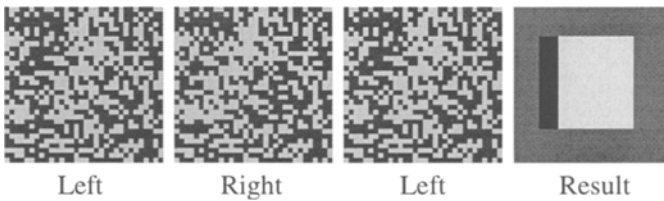


Fig. 7. Random-dot stereogram. Shown as result is the disparity map for the left image. Note the occluded region (black) at the left of the rectangle.

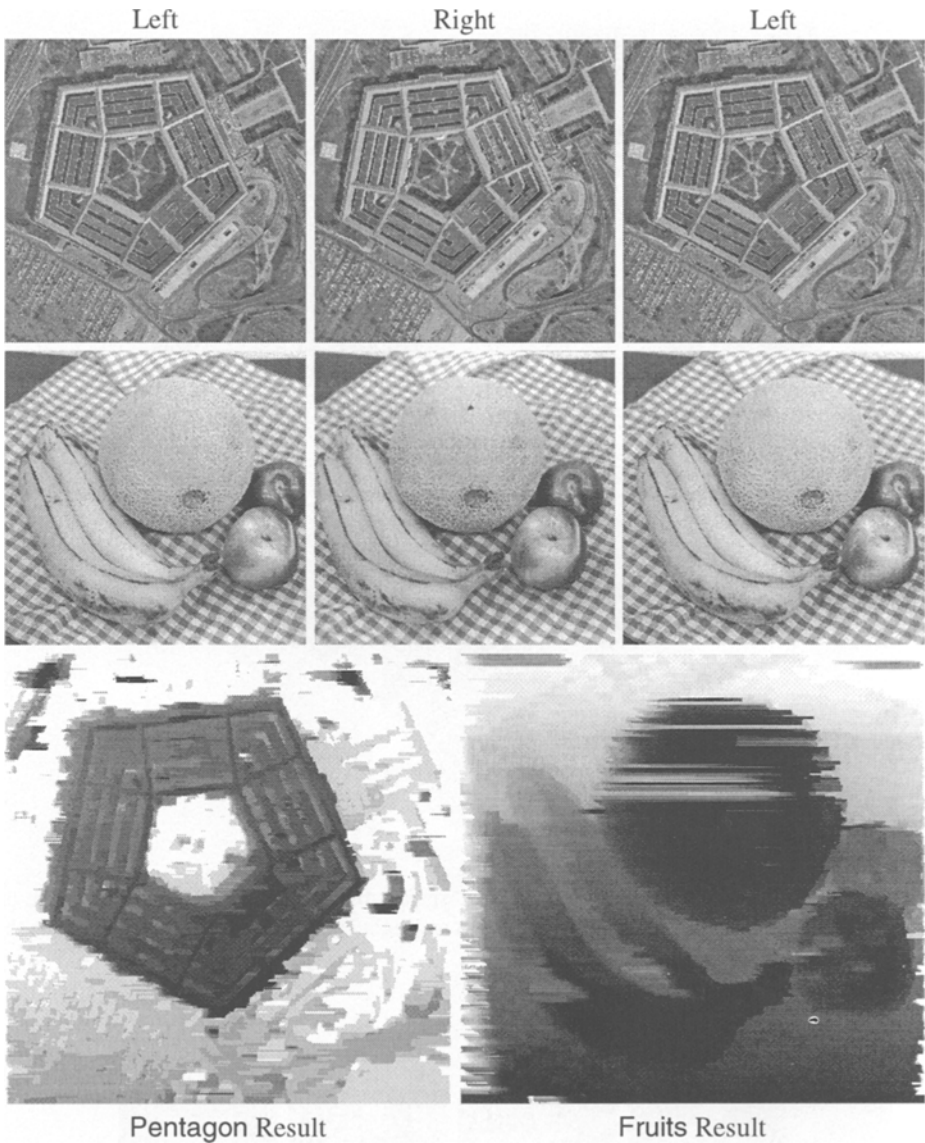


Fig. 8. Results (disparity maps for the left image) on real images **Pentagon** and **Fruits**. Notice the small variance of the height of the ground around the Pentagon can clearly be seen. Also, in **Pentagon**, the disparity map is slightly darker toward the bottom, suggesting that the pictures were not taken exactly over the site. In **Fruits**, one can check by doing stereo with own eyes that even the region in the middle of the cantaloupe that is giving a background disparity is correct, i.e., there is some correlating noise present in both images.

5 Conclusion

We have presented a new approach to compute the disparity map, by modeling occlusions, discontinuities, and epipolar-line interactions, then optimally solving the problem in a polynomial time. We have modeled binocular stereo, including the monotonicity/ordering constraint, mapped the model to a directed graph and used a maximum-flow algorithm to globally optimize the problem. We have improved on previous work, which either (i) models occlusion/discontinuities and ordering constraints, but does not incorporate epipolar-line interactions, with, e.g., as the dynamic programming does, or (ii) models epipolar-line interactions, but not ordering constraint or discontinuities, as in the case of maximum-flow approach on undirected graphs. We have shown that the discontinuities/occlusions cost have to be modeled by a convex function in order to map stereo to a directed graph, and thus used the linear function, since it least penalizes for large discontinuities and offers simplicity (fewer edges in the graph). The experiments support the model.

References

1. B. Anderson. The role of partial occlusion in stereopsis. *Nature*, 367:365–368, 1994.
2. N. Ayache. *Artificial Vision for Mobile Robots*. MIT Press, Cambridge, Mass., 1991.
3. P. N. Belhumeur and D. Mumford. A bayesian treatment of the stereo correspondence problem using half-occluded regions. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition*, 1992.
4. A. Blake and A. Zisserman. *Visual Reconstruction*. MIT Press, Cambridge, Mass., 1987.
5. P. Burt and B. Julesz. A disparity gradient limit for binocular fusion. *Science*, 208:615–617, 1980.
6. B. Cernushi-Frias, D. B. Cooper, Y. P. Hung, and P. Belhumeur. Towards a model-based bayesian theory for estimating and recognizing parameterized 3d objects using two or more images taken from different positions. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, PAMI-11:1028–1052, 1989.
7. A. Champolle, D. Geiger, and S. Mallat. Un algorithme multi-échelle de mise en correspondance stéréo basé sur les champs markoviens. In *13th GRETSI Conference on Signal and Image Processing*, Juan-les-Pins, France, Sept. 1991.
8. B. V. Cherkassky and A. V. Goldberg. On Implementing Push-Relabel Method for the Maximum Flow Problem. In *Proc. 4th Int. Programming and Combinatorial Optimization Conf.*, pages 157-171, 1995.
9. T. H. Cormen, C. E. Leiserson, and R. L. Rivest. *Introduction to Algorithms*. McGraw-Hill, New York, 1990.
10. O. Faugeras. *Three-Dimensional Computer Vision*. MIT Press, Cambridge, Mass., 1993.
11. D. Geiger and B. Ladendorf and A. Yuille. Occlusions and binocular stereo. *International Journal of Computer Vision*, 14, pp 211-226 March, 1995.
12. B. Gillam and E. Borsting. The role of monocular regions in stereoscopic displays. *Perception*, 17:603–608, 1988.
13. W. E. L. Grimson. *From Images to Surfaces*. MIT Press, Cambridge, Mass., 1981.
14. B. Julesz. *Foundations of Cyclopean Perception*. The University of Chicago Press, Chicago, 1971.
15. T. Kanade and M. Okutomi. A stereo matching algorithm with an adaptive window: theory and experiments. In *Proc. Image Understanding Workshop DARPA*, PA, September 1990.

16. J. Malik. On Binocularly viewed occlusion Junctions. In *Fourth European Conference on Computer Vision*, vol.1, pages 167–174, Cambridge, UK, 1996. Springer-Verlag.
17. D. Marr and T. Poggio. Cooperative computation of stereo disparity. *Science*, 194:283–287, 1976.
18. D. Marr and T. Poggio. A computational theory of human stereo vision. *Proceedings of the Royal Society of London B*, 204:301–328, 1979.
19. K. Nakayama and S. Shimojo. Da vinci stereopsis: depth and subjective occluding contours from unpaired image points. *Vision Research*, 30:1811–1825, 1990.
20. Y. Ohta and T. Kanade. Stereo by intra- and inter-scanline search using dynamic programming. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, PAMI-7(2):139–154, 1985.
21. S. B. Pollard, J. E. W. Mayhew, and J. P. Frisby. Disparity gradients and stereo correspondences. *Perception*, 1987.
22. S. Roy and I. Cox. A Maximum-Flow Formulation of the N-camera Stereo Correspondence Problem In *Proc. Int. Conf. on Computer Vision, ICCV'98*, Bombay, India 1998.