

Single-view modelling of free-form scenes

Li Zhang, Guillaume Dugas-Phocion, Jean-Sebastien Samson
and Steven M. Seitz*



This paper presents a novel approach for reconstructing free-form, texture-mapped, 3D scene models from a single painting or photograph. Given a sparse set of user-specified constraints on the local shape of the scene, a smooth 3D surface that satisfies the constraints is generated. This problem is formulated as a constrained variational optimization problem. In contrast to previous work in single-view reconstruction, our technique enables high-quality reconstructions of free-form curved surfaces with arbitrary reflectance properties. A key feature of the approach is a novel hierarchical transformation technique for accelerating convergence on a non-uniform, piecewise continuous grid. The technique is interactive and updates the model in real time as constraints are added, allowing fast reconstruction of photorealistic scene models. The approach is shown to yield high-quality results on a large variety of images. Copyright © 2002 John Wiley & Sons, Ltd.

Revised: 2 May 2002

KEY WORDS: shape reconstruction; hierarchical transformation; discontinuities; pictorial relief; free-form modelling; variational surfaces

Introduction

One of the most impressive features of the human visual system is our ability to infer 3D shape information from a single photograph or painting. A variety of strong single-image cues have been identified and used in computer vision algorithms (e.g. shading, texture, and focus) to model objects from a single image. However, existing techniques are not capable of robustly reconstructing free-form objects with general reflectance properties. This deficiency is not surprising given the ill-posed nature of the problem—from a single view it is not possible to differentiate an image of an object from an image of a flat photograph of the object. Obtaining good shape models from a single view therefore requires invoking domain knowledge.

In this paper, we argue that a reasonable amount of user interaction is sufficient to create high-quality 3D scene reconstructions from a single image, *without*

placing strong assumptions on either the shape or reflectance properties of the scene. To justify this argument, an algorithm is presented that takes as input a sparse set of user-specified constraints, including surface positions, normals, silhouettes, and creases, and generates a well-behaved 3D surface satisfying the constraints. As each constraint is specified, the system recalculates and displays the reconstruction in real time. The algorithm yields high-quality results on images with limited perspective distortion.

We cast the single-view modelling problem as a constrained variational optimization problem. Building upon previous work in hierarchical surface modelling,^{1–3} the scene is modeled as a piecewise continuous surface represented on a quad-tree-based adaptive grid and is computed using a novel hierarchical transformation technique. The advantages of our approach are:

- A general constraint mechanism: any combination of point, curve, and region constraints may be specified as image-based constraints on the reconstruction.
- Adaptive resolution: the grid adapts to the complexity of the scene, i.e., the quad-tree representation can be made more detailed around contours and regions of high curvature.
- Real-time performance: a hierarchical transformation technique is introduced that enables 3D reconstruction at interactive rates.

*Correspondence to: S. Seitz, Department of Computer Science and Engineering, University of Washington, Box 352350, Seattle, WA 98195-2350, USA. E-mail: seitz@cs.washington.edu

Contract/grant sponsor: Intel Corporation.
Contract/grant sponsor: Microsoft Corporation.
Contract/grant sponsor: National Science Foundation.
Contract/grant number: IIS-0049095.
Contract/grant number: IIS-9984672.



Figure 1. The 3D model at right is generated from a single image and user-specified constraints.

A technical contribution of our algorithm is the formulation of a hierarchical transformation technique that handles discontinuities. *Controlled-continuity stabilizers*⁴ have been proposed to model discontinuities in inverse visual reconstruction problems. Whereas hierarchical schemes^{1,2} and quad-tree splines³ enabled fast solutions for related variational problems, our approach for integrating discontinuity conditions into a hierarchical transformation framework is shown to yield significant performance improvements over these prior methods.

The remainder of the paper is structured as follows. We begin with a brief summary of related work. The section on ‘A Variational Framework for Single-View Modeling’ formulates single-view modelling as a constrained optimization problem in a high-dimensional space. In order to solve this large-scale optimization problem efficiently, a novel hierarchical transformation technique that operates on an adaptive grid is introduced in the section on ‘Hierarchical Transformation with Adaptive Resolution’. Experimental results will be presented followed by a conclusion to the paper.

Previous Work on Single-View Modeling

The topic of 3D reconstruction from a single image is a long-standing problem in the computer vision literature. Traditional approaches for solving this problem have isolated a particular cue, such as shading,⁵ texture,⁶ or focus.⁷ Because these techniques make strong assumptions on shape, reflectance, or exposure, they tend to produce acceptable results for only a restricted class of images. Of these, the topic of *shape from shading*, pioneered by Horn,⁵ is most related to our approach in its use of variational techniques and surface normal analysis.

More recent work by a number of researchers has shown that moderate user interaction is highly effective in creating 3D models from a single view. In particular, Horry *et al.*⁸ and Criminisi *et al.*⁹ reconstructed piecewise planar models based on user-specified vanishing points and geometric invariants. Shum *et al.*¹⁰ generated similar models from panoramas using a constraint system based on user input. The *Façade* system¹¹ modeled architectural scenes using collections of simple primitives from one or more images, also with the assistance of a user. A limitation of these approaches is that they are limited to scenes composed of planes or other simple primitives and do not permit modelling of free-form scenes. A different approach is to use domain knowledge; for example, Blanz and Vetter¹² have obtained remarkable reconstructions of human faces from a single view using a database of head models.

A primary source of inspiration for our work is a series of papers on the topic of *Pictorial Relief*.¹³ In this work, Koenderink and his colleagues explored the depth perception abilities of the human visual system by having several human subjects hand-annotate images with relative distance or surface normal information. They found that humans are quite proficient at specifying local surface orientation, i.e., normals, and that integrating a dense user-specified normal field leads to a well-formed surface that approximates the real object, up to a depth scale. Interestingly, the depth scale varies across individuals and is influenced by illumination conditions. We believe that the role of this depth scale is mitigated in our work, due to the fact that we allow the user to view the reconstruction from any viewpoint during the modelling process—the user will set the normals and other constraints so that the model appears correct from *all* viewpoints, rather than just the original view. The surface integration technique used by Koenderink *et al.* is not attractive as

a general-purpose modelling tool, due to the large amount of human labor needed to annotate every pixel or grid point in the image. Although it is also based on the principles put forth in the Pictorial Relief work, our modelling technique is more efficient, works from sparse constraints, and incorporates discontinuities and other types of constraints in a general-purpose optimization framework.

An interesting alternative to the approach advocated in this paper is to treat the scene as an intensity-coded depth image and use traditional image-editing techniques to sculpt the depth image.¹⁴⁻¹⁶ While our framework allows direct specification of depth values, we found that surface normals are easier to specify and provide more intuitive surface controls. This conclusion is consistent with Koenderink's findings¹³ that humans are more adept at perceiving local surface orientation than relative depth.

A Variational Framework for Single-View Modeling

The subset of a scene that is visible from a single image may be modeled as a piecewise continuous surface. In our approach, this surface is reconstructed from a set of user-specified constraints, such as point positions, normals, contours, and regions. The problem of computing the best surface that satisfies these constraints is cast as a constrained optimization problem.

Surface Representation

In this paper, the scene is represented as a piecewise continuous function, $f(x, y)$, referred to as the *depth map*. Samples of f are represented on a discrete grid, $g_{i,j} = f(id, jd)$, where the i and j samples correspond to pixel coordinates of the input image, and d is the distance between adjacent samples, assumed to be the same in x and y . Denote \mathbf{g} as the vector whose components are $g_{i,j}$.

A set of four adjacent samples, $A = (i, j)$, $B = (i + 1, j)$, $C = (i + 1, j + 1)$, and $D = (i, j + 1)$ define the corners of a grid cell. Note that a cell, written as A-B-C-D, is specified by its vertices listed in counter-clockwise order.

The technique presented in this paper reconstructs the smoothest surface that satisfies a set of user-specified constraints. A natural measure of surface smoothness is

the *thin plate* functional:⁴

$$Q_0(\mathbf{g}) = \frac{1}{d^2} \sum_{i,j} [\alpha_{i,j}(g_{i+1,j} - 2g_{i,j} + g_{i-1,j})^2 + 2\beta_{i,j}(g_{i+1,j+1} - g_{i,j+1} - g_{i+1,j} + g_{i,j})^2 + \gamma_{i,j}(g_{i,j+1} - 2g_{i,j} + g_{i,j-1})^2] \quad (1)$$

where $\alpha_{i,j}$, $\beta_{i,j}$, and $\gamma_{i,j}$ are weights that take on values of 0 or 1 and are used to define discontinuities, as will be described later.

Piecewise Continuous Surface Representation.

While it is convenient to represent a surface by a grid of samples, users should have the freedom to interact with a continuous surface by specifying constraints at any location with sub-grid accuracy. Given a sampled surface $g_{i,j}$, we represent the continuous surface $f(x, y)$ using a triangular mesh. Specifically, each grid cell is divided into four triangles by inserting a vertex at the centre with depth defined as the average of the depths of the four corner samples, and adding edges connecting the new vertex with the four corners. The resulting mesh defines a piecewise planar surface over the cell. The depth of each point in the cell can be expressed as a barycentric combination of the depth values of four corner samples. Grid cells that intersect discontinuity curves are omitted from the representation and appear as gaps in the reconstruction.

Constraints

Our technique supports five types of constraints: point constraints, depth discontinuities, creases, planar region constraints, and fairing curve constraints. Point constraints specify the position or the surface normal of any point on the surface. Surface discontinuity constraints identify tears in the surface, and crease constraints specify curves across which surface normals are not continuous. Planar region constraints determine surface patches that lie on the same plane. Fairing curve constraints allow users to control the smoothness of the surface along any curve in the image.

Point Constraints. A point constraint sets the depth and/or the surface normal of any point in the input image. A position constraint is specified by clicking at a point in the image to define the (sub-pixel) position (x_0, y_0) , and then dragging up or down to specify the depth value. A surface normal is specified by rendering

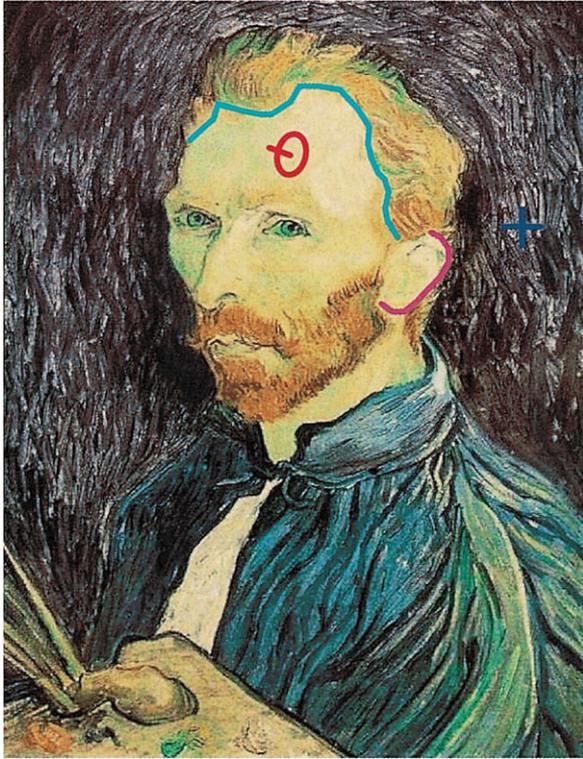


Figure 2. Specifying constraints. Surface orientation is specified by placing a disc with a needle pointing in the direction of the surface normal, shown in red. A position constraint is rendered as a blue cross. Depth discontinuity and crease curves are drawn in purple and cyan, respectively.

a figure representing the projection of a disk sitting on the surface with a short line pointing in the direction of the surface normal (Figures 2 and 3(a)). This figure is superimposed over the point in the image where the normal is to be specified and manually rotated until it appears to align with the surface in the manner

proposed by Koenderink.¹³ In order to uniquely determine the normal from its image plane projection, we assume orthographic projection.

A position constraint $f(x_0, y_0) = f_0$ defines the following constraint:

$$c_{00}g_{i,j} + c_{10}g_{i+1,j} + c_{01}g_{i,j+1} + c_{11}g_{i+1,j+1} = f_0 \quad (2)$$

where (x_0, y_0) is located in grid cell $[id, (i + 1)d] \times [jd, (j + 1)d]$ and c_{00}, c_{01}, c_{10} , and c_{11} are the barycentric coordinates of (x_0, y_0) . Specifying the normal of a point (x_0, y_0) to be $(N_x, N_y, N_z)^T$, defines the following pair of constraints:

$$\frac{f(x_0 + d, y_0) - f(x_0 - d, y_0)}{2d} = -\frac{N_x}{N_z} \quad (3)$$

$$\frac{f(x_0, y_0 + d) - f(x_0, y_0 - d)}{2d} = -\frac{N_y}{N_z} \quad (4)$$

Substituting Equation (2) for $f(x_0 \pm d, y_0)$ and $f(x_0, y_0 \pm d)$ yields two linear constraints on \mathbf{g} . An example of the effects of position and normal constraints is shown in Figure 3(a).

Depth Discontinuities and Creases. A depth discontinuity is a curve across which surface depth is not continuous, creating a tear in the surface. A crease is a curve across which the surface normal is not continuous, while the surface depth is continuous. Depth discontinuities and creases are introduced to model important features in real-world imagery. For example, mountain ridges can be modeled as creases and silhouettes of objects can be modeled as depth discontinuities. These features can be easily specified by users with a 2D graphics interface.

Depth discontinuities and creases are modeled by defining the weights $\alpha_{i,j}$, $\beta_{i,j}$, and $\gamma_{i,j}$ in the smoothness

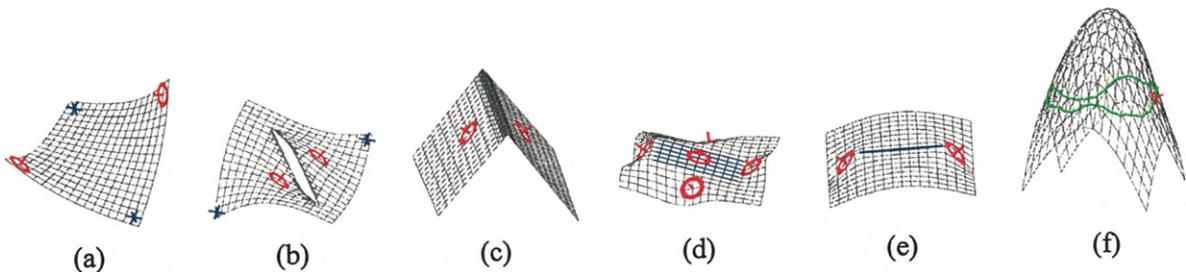


Figure 3. Modeling constraints. (a) The effects of position (blue crosses) and surface normal constraints (red discs with needles). (b) A depth discontinuity constraint creates a tear. (c) A crease constraint. (d) The blue region is a planar region constraint. (e) A fairing curve minimizing curvature. (f) A fairing curve minimizing torsion makes the surface bend smoothly given a single normal constraint—this type of constraint is useful for modeling silhouettes.

objective function of Equation (1). Given a depth discontinuity curve, let A-B-C be a set of three consecutive collinear grid points that cross the curve, and D-E-F-G a cell that the curve intersects. For each such tuple A-B-C, the term $(g_A - 2g_B + g_C)^2$ is dropped from Q_0 by setting α_B or γ_B to 0. For each such cell D-E-F-G, the term $(g_D - g_E - g_F + g_G)^2$ is also dropped by setting β_G to 0. Each crease curve is first *scan converted*¹⁷ to the sampling grid points. Then, all the terms $(g_A - 2g_B + g_C)^2$ are dropped if B is on the curve; all the terms $(g_D - g_E - g_F + g_G)^2$ are dropped if either edge D-E or edge F-G is on the curve. Otherwise, all the weights are 1 by default. Examples of depth discontinuity and crease constraints are shown in Figures 3(b) and (c) respectively.

Planar Region Constraints. The necessary and sufficient conditions for surface planarity over a region R are $f_{xx}(x, y) = f_{xy}(x, y) = f_{yy}(x, y) = 0, \forall(x, y) \in R$. For a discrete mesh \mathbf{g} , these constraints translate to

$$g_A - 2g_B + g_C = 0 \quad (5)$$

$$g_D - g_E - g_F + g_G = 0 \quad (6)$$

for all three consecutive collinear grid points A-B-C in R , and for all cells D-E-F-G in R . An example of a planar region constraint is shown in Figure 3(d).

Fairing Curve Constraints. It is often useful to control the smoothness of the surface along or across a specific curve. For example, surface depth is made to vary slowly *along* a curve in Figure 3(e), and the surface gradient is made to vary slowly *across* a curve in Figure 3(f). Fairing curves provide better control of the shape of the surface along salient contours such as silhouettes, and are achieved as follows.

Suppose that a user specifies a curve $\tau(l) = (x(l), y(l))^T$ in the image. To maximize the smoothness along the curve, the following integral is minimized:

$$Q_d(\tau) = \int_l \left(\frac{d^2}{dl^2} f(\tau(l)) \right)^2 dl \quad (7)$$

The gradient of the surface across τ is $(\nabla f)^T \mathbf{n}_\tau$, where $\nabla f = (f_x, f_y)^T$ is the gradient of the surface $f(x, y)$ at the point $\tau(l)$ and $\mathbf{n}_\tau(l) = (-\frac{d}{dl}y(l), \frac{d}{dl}x(l))^T$ is the normal of τ . To make the surface gradient across $\tau(l)$ have small variation, the integral

$$Q_s(\tau) = \int_l \left(\frac{d}{dl} ((\nabla f)^T \mathbf{n}_\tau) \right)^2 dl \quad (8)$$

is minimized. Note that $\frac{d}{dl} ((\nabla f)^T \mathbf{n}_\tau)$ is the derivative of the surface gradient *across* the curve with respect to the curve parameter.

The terms $\frac{d^2}{dl^2} f(\tau(l))$ and $\frac{d}{dl} ((\nabla f)^T \mathbf{n}_\tau)$ may be discretized as follows:

$$\frac{d^2}{dl^2} f(\tau(l_i)) = f(\tau(l_{i+1})) - 2f(\tau(l_{i-1})) + f(\tau(l_{i-1}))$$

$$((\nabla f)^T \mathbf{n}_\tau)(l_i) = f \left(\left(\tau + \frac{d}{2} \mathbf{n}_\tau \right) (l_i) \right)$$

$$- f \left(\left(\tau - \frac{d}{2} \mathbf{n}_\tau \right) (l_i) \right)$$

$$\frac{d}{dl} ((\nabla f)^T \mathbf{n}_\tau)(l_i) = ((\nabla f)^T \mathbf{n}_\tau)(l_{i+1}) - ((\nabla f)^T \mathbf{n}_\tau)(l_i)$$

where $\{\tau(l_i)\}$ are sampling points on the curve. Consequently, Equations (7) and (8) can be expressed as quadratic functions of \mathbf{g} . The resulting equations are added, with weights ζ_τ and η_τ , into Equation (1), resulting in a modified surface smoothness objective function $Q(\mathbf{g})$:

$$Q_c(\tau) = \zeta_\tau Q_d(\tau) + \eta_\tau Q_s(\tau) \quad (9)$$

$$Q(\mathbf{g}) = Q_0(\mathbf{g}) + \sum_\tau Q_c(\tau) \quad (10)$$

We call $\zeta_\tau Q_d(\tau)$ the *curvature* term and $\eta_\tau Q_s(\tau)$ the *torsion* term.

Linearly Constrained Quadratic Optimization

Based on the surface objective function and constraints presented above, finding the smoothest surface that satisfies these constraints may be formulated as a linearly constrained quadratic optimization. Point constraints and planar region constraints introduce a set of linear equations, Equations (2–6), for the depth map \mathbf{g} , expressed as $\mathbf{A}\mathbf{g} = \mathbf{b}$. Surface discontinuity and crease constraints define weights α, β , and γ and fairing curve constraints introduce the terms in Equation (9). $Q(\mathbf{g})$ is a quadratic form and can be expressed as $\mathbf{g}^T \mathbf{H}\mathbf{g}$. Consequently, our linearly constrained quadratic optimization problem is defined by

$$\begin{cases} \mathbf{g}^* = \arg \min_{\mathbf{g}} \{ Q(\mathbf{g}) = \mathbf{g}^T \mathbf{H}\mathbf{g} \} \\ \text{subject to } \mathbf{A}\mathbf{g} = \mathbf{b} \end{cases} \quad (11)$$

The Lagrange multiplier method is used to convert this problem into the following augmented linear system:

$$\begin{bmatrix} \mathbf{H} & \mathbf{A}^T \\ \mathbf{A} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{g} \\ \lambda \end{bmatrix} = \begin{bmatrix} \mathbf{0} \\ \mathbf{b} \end{bmatrix} \quad (12)$$

The Hessian matrix \mathbf{H} is a banded sparse matrix. For a grid of size N by N , \mathbf{H} is of size N^2 by N^2 , with band width of $O(N)$ and about 13 non-zero elements per row. Direct methods, such as LU decomposition, are of $O(N^4)$ time complexity, and therefore do not scale well for large grid sizes. Iterative methods are more applicable. We use the *minimum residue* method,¹⁸ designed for symmetric indefinite systems. However, the linear system arising from Equation (11) is often poorly conditioned, resulting in slow convergence of the iterative solver. To address this problem, a hierarchical-basis preconditioning approach with adaptive resolution is presented in the next section.

Hierarchical Transformation with Adaptive Resolution

The reason for the slow convergence of the minimum residue method is that it takes many iterations to propagate a constraint to its neighborhood, due to the sparseness of \mathbf{H} . The first row of Figure 6 shows an example of this constraint propagation process, where the two normal constraints generate only two small ripples after 200 iterations. Multigrid techniques¹⁹ have been applied to this type of problem; however, they are tricky to implement and require a fairly smooth solution to be effective.¹ Szeliski¹ and Gortler and Cohen² use hierarchical basis functions to accelerate the solution of linear systems like Equation (12). We review their approach next, to provide a foundation for our work which builds upon it. In the hierarchical approach, a regular grid is represented with a pyramid of coefficients,²⁰ where the number of coefficients is equal to the original number of grid points. The coarse-level coefficients in the pyramid determine a low-resolution surface sampling and fine-level coefficients determine surface details, represented as displacements relative to the interpolation of the low-resolution sampling. To convert from coefficients $\tilde{\mathbf{g}}$ to depth values \mathbf{g} , the algorithm starts from the coarsest level, doubles the resolution by linearly interpolating the values of the current level, adds in the displacement values defined by the coefficients in the next finer level,

moves to the next finer level, and repeats the procedure until the finest resolution is obtained. Using similar notation to Szeliski's,³ the process can be written

```

procedure CoefToDepth( $\tilde{\mathbf{g}}$ )
  for l=1 to L-1
    for every grid point P in level l
       $\mathbf{g}_P = \tilde{\mathbf{g}}_P + \sum_{Q \in N_P} w_{P,Q} * \mathbf{g}_Q$ 
    return  $\mathbf{g}$ 
end CoefToDepth

```

where L is the number of levels in the hierarchy, $\tilde{\mathbf{g}}_P$ is the hierarchical coefficient for P , N_P is the set of grid points in level $l-1$ used in interpolation for P in level l , and $w_{P,Q}$ is a weight that will be described later. Level 0 consists of a single cell, with coefficients defined to be the depth values at the corners of the cell.

In previous work,¹³ the weights $w_{P,Q}$ were implicitly defined to be constant, resulting in a simple averaging operation for computing P from N_P . This averaging approach implicitly assumes local smoothness within the region defined by N_P , resulting in poor convergence in the presence of discontinuities. In practice, this choice of weights causes the artifact that modifying the surface on one side of a discontinuity boundary disturbs the shape on the other side during the iterative convergence process. As a result, it takes longer to converge to a solution, and results in unnatural convergence behavior. The latter artifact is a problem in an incremental solver where the evolving surface is displayed to the user, as is done in our implementation. To address this problem, we next introduce a new interpolation rule to handle discontinuities between the grid points in N_P .

The basic unit in the 2D hierarchical transformation technique is the cell shown in Figure 4, where the depth for corners A, B, C, and D has already been computed and the task is to transform coefficients at E, F, G, H, and I to depth values at these points. With the same notation as in the procedure CoefToDepth, $N_E = \{A, B\}$, $N_F = \{B, C\}$, $N_G = \{C, D\}$, $N_H = \{D, A\}$, and $N_I = \{E, F, G, H\}$. g_E, g_F, g_G , and g_H are first interpolated from A, B, C, and D along edges, and then offset by their respective coefficients $\tilde{g}_{E'}$, $\tilde{g}_{F'}$, $\tilde{g}_{G'}$, and $\tilde{g}_{H'}$. Second, g_I is interpolated from g_E, g_F, g_G , and g_H and offset by its coefficient, \tilde{g}_I . The two interpolation steps above use *continuity-based interpolation* with weights defined as

$$w_{P,Q} = \begin{cases} \frac{e_{P,Q}}{\sum_{Q \in N_P} e_{P,Q}} & \text{if } \sum_{Q \in N_P} e_{P,Q} > 0; \\ 0 & \text{otherwise} \end{cases}$$

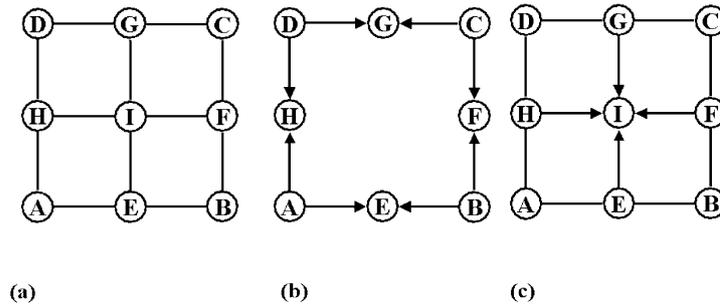


Figure 4. A cell is the primitive for 2D hierarchical transformation. (a) Interpolation on a cell; (b) edge interpolation; (c) cross interpolation. The depth at the center point I is interpolated from the midpoints E, F, G, and H, which are in turn interpolated along each edge of the cell.

where

$$e_{P,Q} = \begin{cases} 1 & \text{if edge } P - Q \text{ is continuous} \\ 0 & \text{otherwise} \end{cases}$$

In the absence of discontinuities, the proposed *continuity-based weighting* scheme is the same as simple averaging schemes used in previous work.^{1,3} In the presence of discontinuities, only locally continuous coarse-level grid points are used in the interpolation. The new scheme prevents interference across discontinuity boundaries and consequently accelerates the convergence of the Minimum Residue algorithm. The second and third rows of Figure 6 show a performance comparison between standard hierarchical transformation and our transformation with continuity-based weighting on a simple surface modelling problem with one discontinuity curve. The improvement of our algorithm is quite evident in this example. In the third row, our new transformation both accelerates the propagation of

constraints and removes the interference across the discontinuity boundary. We have found this kind of behavior very typical in practice and find that adding continuity-based weighting yields dramatic improvements in system performance.

To summarize our approach, instead of solving Equation (12) directly for \mathbf{g} , we instead compute the hierarchical coefficients $\tilde{\mathbf{g}}$. The conversion from $\tilde{\mathbf{g}}$ to \mathbf{g} is implemented by the procedure *CoefToDepth* with continuity-based weighting. *CoefToDepth* implements a concatenation of linear transformations and can be described by a matrix $\mathbf{S} = \mathbf{S}_1 \mathbf{S}_2 \dots \mathbf{S}_{L-1}$,¹ where \mathbf{S}_l performs the continuity-based interpolation and displacement operation for grid points in level l . Substituting $\mathbf{g} = \mathbf{S}\tilde{\mathbf{g}}$ into Equation (11) and applying the Lagrange Multiplier method yields the transformed linear system:²

$$\begin{bmatrix} \mathbf{S}^T \mathbf{H} \mathbf{S} & \mathbf{S}^T \mathbf{A}^T \\ \mathbf{A} \mathbf{S} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{g}} \\ \lambda \end{bmatrix} = \begin{bmatrix} \mathbf{0} \\ \mathbf{b} \end{bmatrix} \quad (13)$$

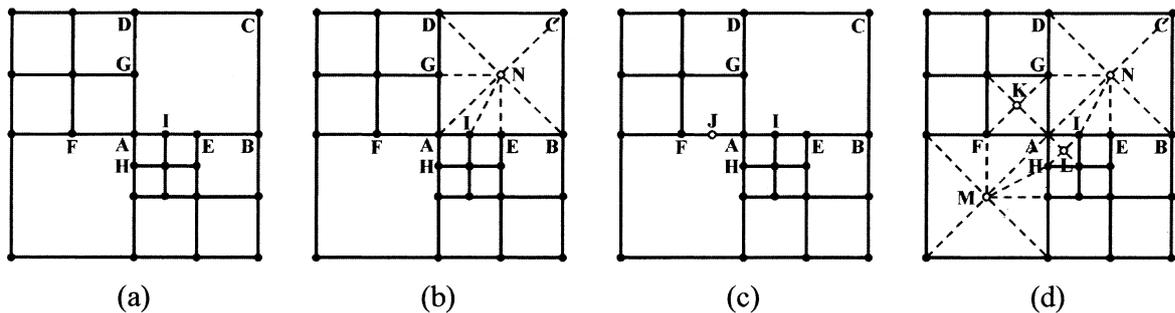


Figure 5. Depth is represented by a quad-tree grid. (a) An example of a quad-tree of three levels. (b) Triangulating the cell A-B-C-D makes it piecewise continuous. (c) The horizontal second-order derivative at grid point A is computed from the depths at I and J. J is a virtual vertex, whose depth is interpolated from grid points F and A. (d) Grid point A is weighted according to one-half of the area of all incident triangles, i.e., AHL, AIN, ANG, AGK, AKF, AFM, and AMH.

The matrix $S^T H S$ is shown to be better conditioned,¹ resulting in faster convergence. The number of floating point operations of the procedure CoefToDepth and its transpose¹ is approximately $2N^2$ for a grid size of $N \times N$. Considering that there are around 13 non-zero elements per row in H , the overhead introduced by S in multiplying $S^T H S$ with a vector is about 15%. Given the considerable reduction in number of iterations shown in Figure 6, the total run time is generally much lower using a hierarchical technique.

Adaptive Surface Resolution

As an alternative to solving for the surface on the full grid, it is advantageous to use an adaptive grid, with higher resolution used only in areas where it is needed. For example, the surface should be sampled densely along a silhouette and sparsely in areas where the geometry is nearly planar. We support adaptive resolution by allowing the user to specify the grid resolution for each region via a user interface. Subdivision may also occur automatically—in our implementation, discontinuity and crease curves are automatically subdivided to enable accurate boundaries. A quad-tree representation is used to represent the adaptive grid. By modifying our hierarchical transformation technique to operate on a quad-tree grid, as in Szeliski and Shum,³

the run time of the algorithm is proportional to the number of subdivided grid points, which is typically much smaller than the full grid.

Modifying the algorithm to operate on a quad-tree requires the following changes. First, at each level l , the CoefToDepth procedure should consider only grid points that are present in the quad-tree.

Second, the triangular mesh representation is modified so that each inserted vertex is connected to *all* neighboring grid points in the quad tree, not just to the four corners of the cell. The depth of each inserted vertex is defined as the average of the depths of the midpoints of the leaf cell's four edges. For example, a 3-level quad-tree pyramid is shown in Figure 5(a). To triangulate the leaf cell A-B-C-D, a vertex N is inserted at the center of the cell, with depth g_N defined as the average of $g_E, \frac{g_B+g_C}{2}, \frac{g_C+g_D}{2}$, and g_G . The cell A-B-C-D is triangulated by connecting N to A, I, E, B, C, D, and G, as shown in Figure 5(b).

Third, the second-order derivatives of f , in terms of g , should be derived from the quad-tree representation by interpolating a regular grid neighborhood around each point from the triangular mesh. For example, in Figure 5(c), $\frac{d^2}{dx^2} f(A)$ is approximated as $\frac{g_I+g_J-2g_A}{|A||AI|}$, where $|JA| = |AI|$ and g_J is interpolated as $\frac{|JA|}{|FA|} g_F + \frac{|FI|}{|FA|} g_A$. Similarly, $\frac{d^2}{dx dy} f(N)$ is approximated as $\frac{g_D-g_C-g_B+g_A}{|AB||AD|}$.

Finally, special care should be taken to approximate surface integrals by summations on the non-uniform

Methods	Iteration 0	Iteration 200	Iteration 1200	Iteration 2500	Iteration 9500
No hierarchical transformation					
Hierarchical transformation without continuity-based weighting					
Novel hierarchical transformation with continuity-based weighting					

Figure 6. Performance comparison of solving Equation (12) by using no hierarchical transformation, traditional transformation, and our novel transformation in terms of number of iterations. The model has approximately 1400 grid points, and four constraints.

grid. Specifically, the smoothness measure (Equation ((1)) is redefined as

$$Q_0(\mathbf{g}) = \sum_{\{A,B,C\}} \alpha_B \mu_B \left(\frac{g_A - 2g_B + g_C}{|AB||BC|} \right)^2 + \sum_{\{D,E,F,G\}} 2\beta_D \nu_D \left(\frac{g_G - g_F - g_E + g_D}{|DE||DF|} \right)^2 \quad (14)$$

where the first summation is over all consecutive collinear grid points A-B-C in the horizontal and vertical directions, and the second summation is over all leaf cells D-E-F-G. The binary coefficients α_B and β_D are determined by continuity conditions, as described previously. The coefficients μ_B and ν_D are the weights based on the size of the local quad-tree neighborhood. In our implementation, ν_D is simply the area of cell D-E-F-G, and μ_B is one half^a of the sum of the areas of triangles that are incident to B. For example, in Figure 5(d), the weight μ_A for the term $\left(\frac{g_I + g_J - 2g_A}{|JA||AI|} \right)^2$ is defined as

$$\mu_A = \frac{1}{2} [\text{area}(AHL) + \text{area}(ALI) + \text{area}(AIN) + \text{area}(ANG) + \text{area}(AGK) + \text{area}(AKF) + \text{area}(AFM) + \text{area}(AMH)] \quad (15)$$

Experimental Results

We have implemented the approach described in this paper and applied it to create reconstructions of a wide variety of objects. Three of these results are presented in this section but higher-resolution images and 3D VRML models can be found online.²¹ We encourage the reader to peruse these results online to better gauge the quality of the reconstructions.

Smooth objects without position discontinuities are especially easy to reconstruct using our approach. As a case in point, the jelly bean image in the first column of Figure 7 requires only isolated normals and creases to generate a compelling model, and can be created quite rapidly (about 20 minutes, including time to specify constraints) using our interactive system. The first column of Figure 7 shows the input image, quad-tree grid with constraints, a view of the quad-tree from a novel viewpoint, and a texture mapped rendering of the same

^aThe factor of one half arises because each triangle contains two quad-tree grid points, and therefore contributes twice to the integral.

view. For this example, the user worked with a 32×32 grid that was automatically subdivided as the crease curves were drawn. This model has 144 constraints in all, 3396 grid points, and required 25 seconds to converge completely on a 1.5GHz Pentium 4 processor, using our hierarchical transformation technique with continuity-based weighting. The system is designed so that new constraints may be added interactively at any time during the modelling process—the user does not have to wait until full convergence to specify more constraints. The second column of Figure 7 shows a single-view reconstruction of the Great Wall of China. This example is challenging due to the complex scene geometry and significant perspective distortions. Because our method assumes orthography, the scale of the tower and other foreground elements is exaggerated in the reconstruction, and the scene appears unrealistically compressed. Addressing perspective is an interesting topic for future work. The Great Wall of China model has 135 constraints, 2566 grid points, and required 40 seconds to converge completely.

An interesting application of single-view modelling techniques is to reconstruct 3D models from paintings. In contrast to other techniques,^{8,9,11,12} our approach does not make strong assumptions about geometry, making it amenable to impressionist and other non-photorealistic works. Here we show a reconstruction created from a self-portrait of Van Gogh. This model has 264 constraints, 3881 grid points, and required 45 seconds to converge. This was the most complex model we tried, requiring roughly 1.5 hours to design. For comparison, it takes 70 seconds to converge without using the hierarchical transformation and 3 minutes using the hierarchical transformation without continuity-based weighting, i.e., an inappropriately weighted hierarchical method can perform significantly worse than not using a hierarchy at all. Note, however, that there is significant room for optimization in our implementation; we expect that the timings for both hierarchical methods could be improved by a factor of 1.5 or 2.

Conclusions

In this paper, it was argued that a reasonable amount of user interaction is sufficient to create high-quality 3D scene reconstructions from a single image, *without* placing strong assumptions on either the shape or reflectance properties of the scene. To justify this argument, an algorithm was presented that takes as input a sparse set of user-specified constraints, including

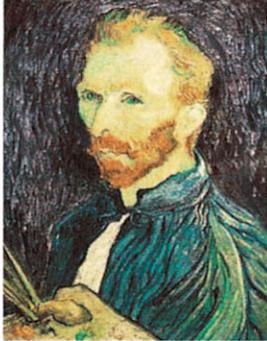
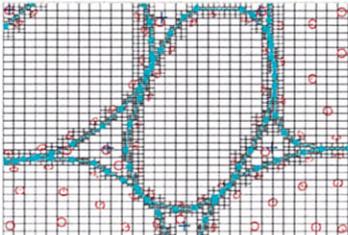
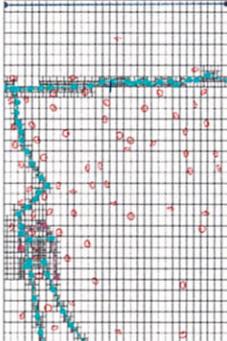
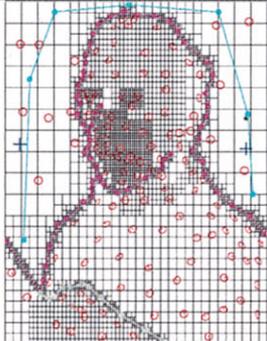
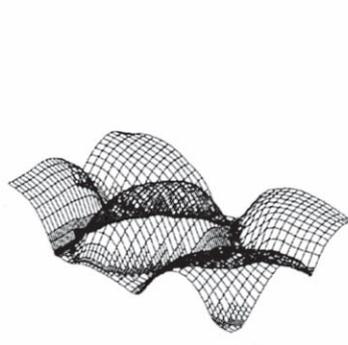
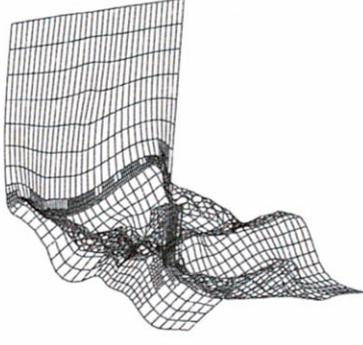
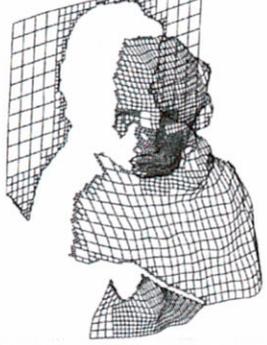
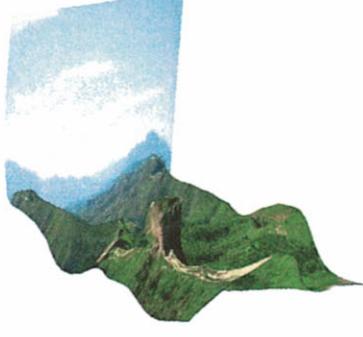
	Jelly Beans	The Great Wall of China	The Self Portrait of van Gogh
Original image			
Constraints			
Wireframe			
Novel view			

Figure 7. Examples of single-view modeling on different scenes. From top to bottom, the rows show the original images, user-specified constraints on adaptive grids, 3D wireframe rendering, and textured rendering.

surface positions, normals, silhouettes, and creases, and generates a well-behaved 3D surface satisfying the constraints. As each constraint is specified, the system recalculates and displays the reconstruction in real time. A technical contribution is a novel hierarchical transformation technique that explicitly models discontinuities and computes surfaces at interactive rates. The approach was shown to yield very good results on a variety of images.

There are a number of interesting avenues for future research in this area. In particular, single-view modelling has the inherent limitation that only visible surfaces in an image can be modeled, leading to distracting holes near occluding boundaries. Automatic hole-filling techniques could be developed that maintain the surface and textural attributes of the scene. Another important extension would be to generalize to perspective projection as well as other useful projection models like panoramas.

ACKNOWLEDGEMENTS

The support of Intel Corporation, Microsoft Corporation, and the National Science Foundation under grants IIS-0049095 and IIS-9984672 is gratefully acknowledged. The jelly bean photo is courtesy of Philip Greenspun, <http://philip.greenspun.com>.

References

1. Szeliski R. Fast surface interpolation using hierarchical basis functions. In *IEEE Transactions on PAMI* 1990; **12**(6): 513–528.
2. Gortler S, Cohen M. Variational modelling with wavelets. *TR-456-94*, Department of Computer Science, Princeton University, 1994.
3. Szeliski R, Shum H-Y. Motion estimation with quadtree splines. In *IEEE Transactions on PAMI* 1996; **18**(12): 1199–1209.
4. Terzopoulos D. Regularization of inverse visual problems involving discontinuities. 1996; In *IEEE Transactions on PAMI* 1986; **8**(4): 413–424.
5. Horn BKP. Height and gradient from shading. In *International Journal of Computer Vision* 1990; **5**(1): 37–75.
6. Super BJ, Bovik AC. Shape from texture using local spectral moments. In *IEEE Transactions on PAMI* 1995; **17**(4): 333–343.
7. Nayar SK, Nakagawa Y. Shape from focus. In *IEEE Transactions on PAMI* 1994; **16**(8): 824–831.
8. Horry Y, Anjyo K, Arai K. Tour into the picture: using a spidery mesh interface to make animation from a single image. In *ACM SIGGRAPH Proceedings* 1997; 225–232.
9. Criminisi A, Reid I, Zisserman A. Single view metrology. In *International Conference on Computer Vision* 1999; 434–442.
10. Shum H-Y, Han M, Szeliski R. Interactive construction of 3D models from panoramic mosaics. In *IEEE Conference on CVPR* 1998; 427–433.
11. Debevec P, Taylor C, Malik J. Façade: modelling and rendering architecture from photographs. In *ACM SIGGRAPH Proceedings* 1996; 11–20.
12. Blanz V, Vetter T. A morphable model for the synthesis of 3D faces. In *ACM SIGGRAPH Proceedings* 1999; 187–194.
13. Koenderink JJ. Pictorial Relief. *Philosophical Transactions of the Royal Society: Mathematics, Physics and Engineering Sciences* 1998; **356**(1740): 1071–1086.
14. Williams L. 3D paint. In *Proceedings of the Symposium on Interactive 3D Graphics Computer Graphics* 1990; 225–233.
15. Kang SB. Depth painting for image-based rendering applications. *Technical Report CRL*, Compaq Computer Corporation, Cambridge Research Laboratory, December 1998.
16. Oh BM, Chen M, Dorsey J, Durand F. Image-based modelling and photo editing. In *ACM SIGGRAPH Proceedings* 2001; 433–442.
17. Foley JD, van Dam A, Feiner SK, Hughes JF. *Computer Graphics: Principles and Practice*. Addison-Wesley: Reading, MA, 1990; 72–91.
18. Press WH, Flannery BP, Teukolsky SA, Vetterling WT. *Numerical Recipes in C*. Cambridge University Press: Cambridge, UK, 1988; 83–89.
19. Terzopoulos D. Image analysis using multigrid relaxation methods. In *IEEE Transactions on PAMI* **8**(2): 129–139.
20. Burt PJ, Adelson EH. The Laplacian pyramid as a compact image code. In *IEEE Transactions on Communication* 1983; **31**(4): 532–540.
21. Single view modelling project website. Online at <http://grail.cs.washington.edu/projects/svm/>.