

Multilevel modelling and rendering of architectural scenes

Akash M Kushal

Gaurav Chanda

Kanishka Shrivastava

Mohit Gupta

Subhajit Sanyal

T. V. N. Sriram

Prem Kalra

Subhashis Banerjee

Department of Computer Science and Engineering,
IIT Delhi, New Delhi 110016, India
email: {pkalra,suban}@cse.iitd.ernet.in

Abstract

We present a novel approach for multilevel modelling and rendering of architectural scenes using a small set of photographs. Our approach is based on interactive probing of intuitive measures like lengths, widths and heights using single view metrology. These measures can then be aggregated with additional input if need be, for defining high-level primitives such as planes, prismatic blocks, cuboids, spheres, and general surfaces of translation and revolution. The modelling approach is modular and incremental which enables ease and flexibility in building large architectural scenes. Further, our approach allows multilevel modelling and rendering. That is, the model can be enhanced/reduced both in terms of changing details at the primitive level and at the model level. This requires inter-registration of two or more images with common contents. Our approach is computationally simple, fast and robust. We demonstrate our results by building a model of a magnificent historical monument in Delhi - Humayun's tomb.

Categories and Subject Descriptors (according to ACM CCS): I.3.7 [Computer Graphics]: Three-Dimensional Graphics and Realism I.3.5 [Computer Graphics]: Computational Geometry and Object Modeling

1. Introduction

Recently the problem of interactive 3D modelling of architectures from a sparse set of photographs has attracted a lot of attention. First Debevec et al.² and then others^{9,8,4} have shown that such image based methods can provide high quality reconstruction and photo-realistic rendering without going through tedious and cumbersome CAD or digitizing systems.

In this paper we describe a system for interactive modelling and rendering of architectural scenes from one or more photographs with multiple levels of details. We present an incremental technique that is based on probing of lengths and widths along parallel directions. Our method can also deal with spheres, cones and frustums, surfaces of translation and surfaces of revolution.

Most methods for recovery of 3D structure from images rely on calibrated stereo or uncalibrated images using epipolar or trilinear constraints⁵. In contrast, methods based on single-view analysis^{2,9,8,4} exploit the structural constraints available in architectural photographs to directly estimate the location and the intrinsic parameters of the camera and reconstruct 3D models through user interactions from one or

more images. 3D reconstruction from a single image must necessarily be through an interactive process, where the user provides information and constraints about the scene structure. Such information may be in terms of vanishing points or lines^{6,1,8,9}, co-planarity¹¹, spatial inter-relationship of features^{2,4} and camera constraints¹².

The Facade system², one of the most successful for modelling and rendering of architectural scenes, consists of a hybrid image and geometry based approach which requires decomposing the scene into prismatic blocks and cuboids which are the basic primitives for modelling. Each primitive is defined in terms of six parameters describing its location and a small number of other parameters describing its internal details. The user needs to interactively define constraints on placements of the primitive blocks and mark corresponding edge segments in the image. In general a non-linear optimization method is required to solve for all the model and camera parameters (external) simultaneously. The internal parameters of the camera need to be calibrated a priori. In some special cases, especially if the blocks are aligned with respect to each other, or their rotation with respect to each other are known, linear methods can be used. While model fitting using simultaneous bundle adjustment⁵ appears to be an attractive approach, it may turn out to be problematic if the model requires a large number of primitive blocks oriented with respect to each other at unknown angles. In such a case the large number of unknown parameters in simultaneous bundle adjustment using non-linear optimization may make the method unsuitable. In addition, prismatic blocks and cuboids may not be the most convenient primitives to use for modelling in all situations (see for example Figure 1).

Other significant approaches to architecture modelling are based on estimation of vanishing lines in images⁹ and metric rectification of planes⁸. In the *PhotoBuilder* system^{9,10} a user needs to mark parallel and orthogonal edges in at least two images. These are used to compute the vanishing points in the images corresponding to the three orthogonal coordinate axes in the world. In addition, the user needs to mark the projection of a cuboid of known size in the images. The vanishing points and the images of the cuboid are sufficient to both localize the cameras in a common world frame and to estimate their intrinsic parameters independently. 3D reconstruction of points in the world is achieved by stereo triangulation from image correspondences. Since the basic method of 3D reconstruction and metrology is based on stereo triangulation, at least two images are required. Leibowitz et al.⁸ propose a method for metric rectification of planes and camera estimation based on computation of circular points in images from vanishing directions. Both these methods are limited to blocks world reconstruction of the scene.

We propose a method based on single view metrology^{1,7} for incremental model building. Our method is based on simple interactive probing of length ratios along directions parallel to the world coordinate frame. It is not necessary that

the three coordinate axes of the world frame are orthogonal and be specified by the user in the same scale. Often, accurate length measurements along the three axes directions in the world are not available, and, in such a case, since our probing technique is based on computing affine ratios along parallel directions in the world, we can still carry out an affine reconstruction of the world model. In case it is known that the three axes directions in the world are orthogonal, the affine reconstruction can be updated to Euclidean using rough knowledge of the internal parameters of the camera or experimenting with the scale factors of the coordinate axes to obtain visual alignment. In case the planes are orthogonal and the scales are known, we can recover the camera parameters and the model directly in Euclidean terms.

In the following sections we present the details of our method. In Section 2 we give an overview of our system. In Section 3 we present our basic method for single view metrology and describe how the basic probing technique can be used either incrementally or in aggregation. In Section 4 we extend the basic method to deal with curved objects. In Section 5 we describe our strategy and provide results for modelling and rendering of Humayun's tomb[†], and finally in Section 6 we conclude the paper.

2. Overview

In what follows we give a brief overview of our system. The method involves the following main steps:

Camera estimation: We first register two planes independently by identifying a rectangle (or parallelogram) on each in a coarse resolution photograph. The planes need not be orthogonal, though we refer to them as *vertical* and *horizontal* for convenience. See Figure 1(a) (also Figure 3). This sets up an affine coordinate system in the world and enables the recovery of corresponding camera parameters⁷. We obtain the vanishing points of the axes directions in the images directly as a by-product. In case the planes are orthogonal and the size of each rectangle is known accurately, then we can recover the camera in Euclidean terms and compute the internal parameters of the camera.

Probing heights and widths: Once the camera projection matrix is known we can compute the height (width) of any arbitrary 3D point by identifying the point (head) and its projection (foot) on the horizontal (vertical) plane in the image along the coordinate axis direction (See Figure 1(a)). In case the foot of the point is not visible on one of the reference plane but is clearly identifiable in a plane parallel to the reference plane, then the offset of the parallel plane from the reference plane must be computed a priori using a similar technique.

[†] <http://depts.washington.edu/uwch/silkroad/cities/india/delhi/humayun/humayun.html>

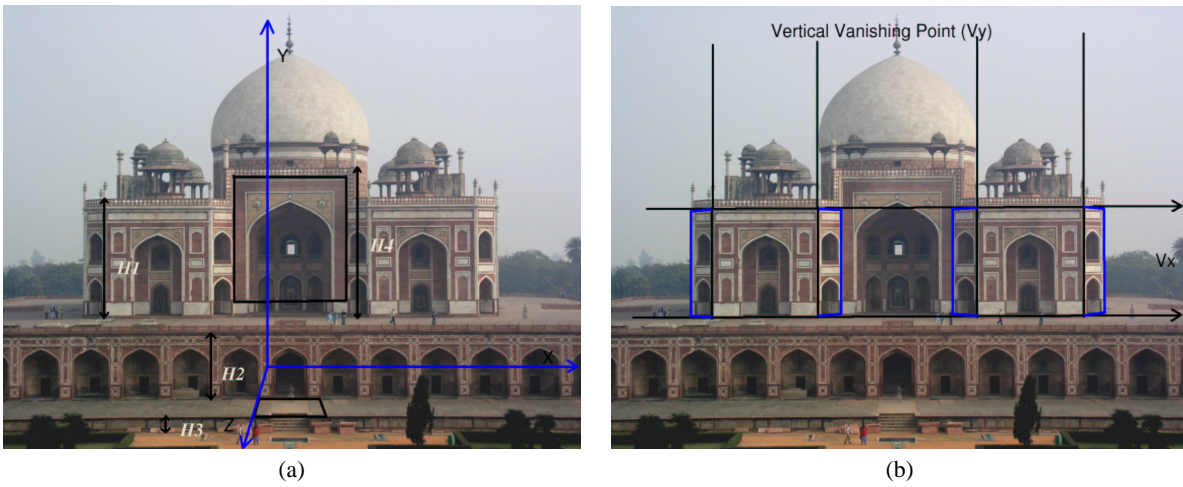


Figure 1: (a) *Camera determination and height measurements:* The user marks a rectangle each on a vertical and a horizontal plane from which the world coordinate frame and the camera is recovered. Head's and foot's can then be clicked to compute heights. (b) *Face Lifting:* The user may specify constraints for modelling of faces.

Modelling faces and blocks: The basic probing technique can be used directly to compute sizes of faces and blocks. Alternatively, several faces and blocks can be simultaneously modelled in a tightly constrained manner. For example, in the photograph of Figure 1(b), the user may choose to enforce the constraints that i) the entire front of the main monument consists of a sequence of planar faces adjacent to each other ii) several of them are parallel to each other iii) all the planar faces are orthogonal to the plane of the courtyard, hence the line joining the head and foot, sampled anywhere, must pass through the vanishing point of the Y direction iv) the top and bottom edges of several of the faces must pass through the vanishing point of the X direction v) the common height of the frontal faces can be estimated by considering a pencil of rays originating from the Y direction vanishing point and computing their intersections with the line through the top and bottom edges of the faces to mark of the *head's* and *foot's*. All the above constraints can be expressed as a single linear system which can be solved using *SVD*. Note that the entire model built as above is dependent on the height of the courtyard, and this dependence can even be expressed symbolically. This height can either be supplied, or computed separately by independent probing, or even be computed simultaneously using additional constraints. Thus, the user has the flexibility of modelling parts of the scene either incrementally or as an aggregate.

Modelling surfaces: The knowledge of the camera projection matrix and the axis is sufficient to compute an arbitrary surface of revolution from its silhouette edges. For the various surfaces of revolutions in the photograph of Figure 2(a) - the axis directions are known to be vertical, the *head* is directly visible in the image, and the foot posi-

tions are known from symmetry. In Section 4 we give the details of modelling arbitrary surfaces of revolutions.

Registering multiple levels of details: For faithful rendering of walk-throughs at close-ups we need to sub-divide a primitive modelled at a coarser resolution into finer primitives. For example, from the photograph in Figure 2(a) the big vertical face of the entrance to the main monument is modelled as a plane. However, at close-up, as in Figure 2(b), we need to carry out an independent metrology to model the inside in terms of several planes and surfaces of translation. For this we set up an independent coordinate system as depicted in Figure 2(b). The user needs to register the two independent metrologies by identifying correspondences as shown in Figure 2. The model with the internal detail of the face shown in wire-frames is depicted in Figure 7.

Rendering: For subsequent rendering of walk-throughs we employ a technique similar to view dependent texture mapping^{2,3}. The model is projected onto the original image(s) from which it is built. This establishes a correspondence of the model parameters with the images and provides the texture coordinates of the associated primitives. Often a single photograph will cover only a part of the model. This requires us to use, in most situations, multiple images at varying resolutions to be able to render the entire model. Compositing with α -channel blending is employed to decide the final pixel value. During the enhancement of the model (changing from one level of detail to another; see Figure 2) we do stenciling to separate the portions and pick appropriate textures.

In what follows we describe our basic probing technique.

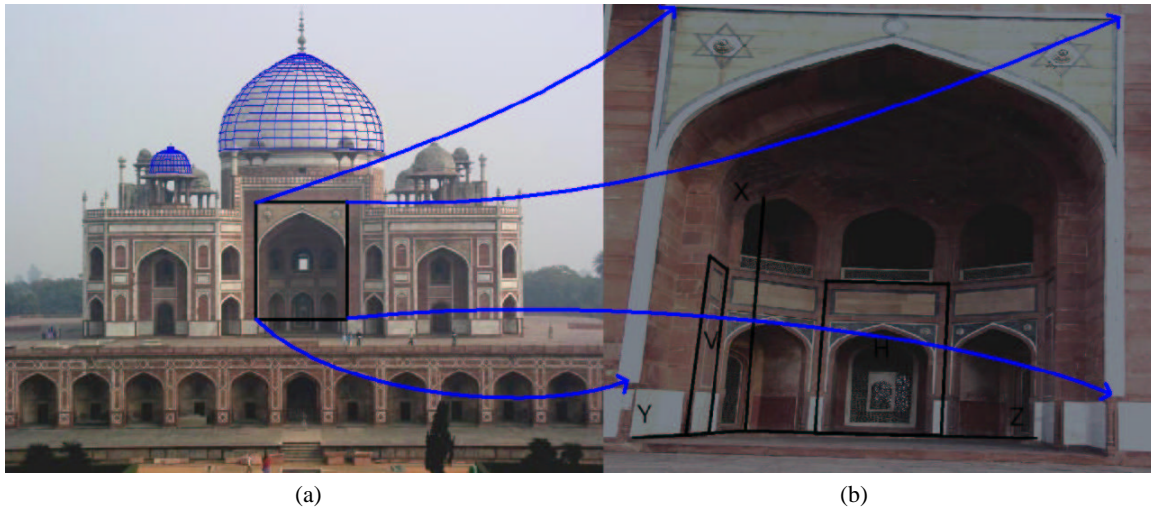


Figure 2: (a) Domes modelled as surfaces of revolution. (b) Independent metrology for multiple levels of details and tie-points.

3. Metrology

In a recent paper Criminisi et al.¹ show that under the customary assumption of the pin-hole camera model⁵, the minimal geometric information that a user needs to provide for computation of 3D measurements is the affine registration of i) a reference plane and ii) a reference direction away from the plane. Once the vanishing line of the reference plane and the vanishing point of the reference direction are established, the user can probe for length ratios on planes parallel to the reference plane and lines parallel to the reference direction and thereby recover 3D affine structure.

We follow a simple and practical method⁷ for implementation of the basic ideas of Criminisi et al.¹ for interactive 3D reconstruction from single images and also consider some extensions. The method is based on computation of homographies⁵ from two world planes to the image, which directly yield the necessary vanishing points and lines reliably. The user needs to identify the two planes by clicking four points on each (See Figures 1(a) and 3) so that the two planes can be independently calibrated in affine terms. Note that in our method we have to provide more information in comparison to the method described by Criminisi et al.¹ however the extra degrees of freedom as a result of the independent registration of the two planes give us constraints on (i) translation of the planes along the coordinate axis, and (ii) rotation of the coordinate system of one plane with respect to the other.

The interactive registration of two planes facilitates the recovery of the camera matrix. Furthermore the internal parameters of the camera can be recovered if registration of the two planes is done in a Euclidean coordinate system.

Let us assume the X dimension of the ground plane rectangle to be the unit length of the world. Let $1/\gamma$ be the Z

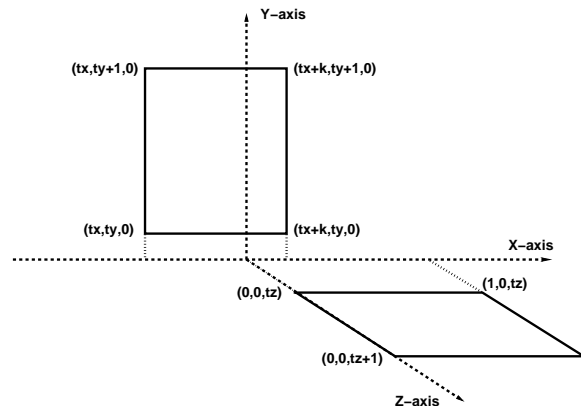


Figure 3: Determining the height.

dimension of the ground plane rectangle and $1/\alpha$ and $1/\beta$ be the X and Y dimensions of the vertical plane rectangle. Let the camera center be located at (C_x, C_y, C_z) . Also let t_x , t_y and t_z be the translations of the rectangles as shown in Figure 3.

1. A knowledge of one of α_x , β_{t_y} or γ_{t_z} enables us to determine the others. One can also then proceed to determine the quantities λC_z and $-\lambda \alpha C_z$. From these two values one can also figure out the value of α .
2. Instead, if α is known then the ratio λC_z can be obtained. Subsequently the unknowns (t_x , β_{t_y} and γ_{t_z}) can be solved in general.

Details of these methods could be accessed at Kushal et al.⁷.

3D affine reconstruction of structured scenes is possible by the probing technique as described below.

Once the camera center coordinates have been determined, in either affine or Euclidean terms, the coordinates of 3D points in the scene can be computed using simple geometry. The user clicks a point (call it the *head*) and its projection on the $X-Z$ plane (call it the *foot*) in the image along the $X-Y$ direction, for example the head and the foot of a person (see Figure 4). Now using the homography \mathbf{H}^{-1} (from the image plane to the $X-Z$ horizontal plane) we can transfer the head and the foot of the point on to the horizontal ($X-Z$ plane). We can then use simple similar triangles to calculate the height (Y coordinate) of the point.

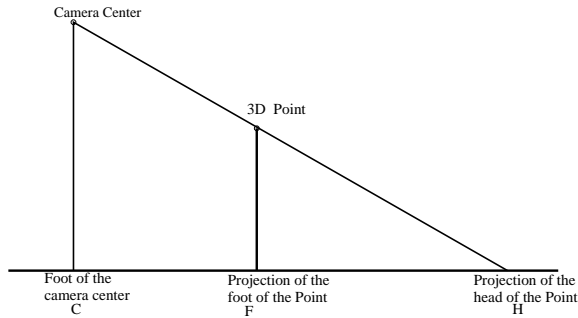


Figure 4: Determining the height.

The height can be computed as $h = \text{Camera height} * \frac{FH}{CH}$

Note that the height can also be obtained relative to the height of another object (with its head and foot given) by eliminating the camera height from the two similar triangle equations.

The actual degrees of freedom of the *head-foot* clicks provided by the user is only 3 since the vanishing point in the Y direction is constrained to lie on the line connecting the head and the foot. We follow a technique similar to Criminisi et al.¹, wherein a user specifies uncertainty areas in the image corresponding to the head and the foot, and the best head and foot estimates are calculated by minimizing the sum of the *Mahalanobis distances* between the estimates and the uncertainty areas specified by the user.

Apart from just the basic probing technique we also provide a method for extracting faces from the model. The user specifies points on the upper and lower edges of the face(s) to be extracted. We fit two lines through these sets of points. The user may specify the direction of these lines (see Figure 1(b)). In this case, the lines fitted are constrained to pass through the vanishing point in that direction. The heights are determined by picking up corresponding points (the line joining them passes through the vanishing point in the vertical direction) on the upper and lower lines. These corresponding points are taken as the *head-foot* pairs to calculate

the height of the face. Here the user may also choose to enforce the constraint that many of the faces have the same height. Other constraints like the fact that two or more of the lines for different faces are parallel or coincident may also be imposed. In this case the height will be solved for by minimizing the error for all the faces together.

4. Metrology of curved surfaces

In this section we extend the basic method to deal with metrology of curved surfaces. As with plane measurements, modelling surfaces of revolution and frustums from a single image also require the user to interactively specify some necessary constraints. In what follows we assume that the camera projection matrix \mathbf{P} is known and describe the constraints for some cases.

4.1. Modelling of general surfaces of revolution

Metrology of a general surface of revolution (SOR) involves computation of a space curve, which when rotated about a given axis produces that surface¹³. We find out the radii at different heights. As a result, we can now compute the generating curve and hence the surface of revolution.

The technique adopted requires the user to specify the axis of revolution in the world (by specifying a point on the axis \mathbf{O} and its direction \mathbf{dir} as shown in Figure 5). The axis direction can be computed using earlier methods of metrology. In addition, the user needs to mark the silhouette of the SOR. This method exploits the fact that the surface normal at any point on the SOR is co-planar with the axis of revolution.

A parametric curve is fitted to the silhouette and sampled uniformly. Consider the equation

$$\mathbf{C} = \mathbf{O} + \lambda_1 \mathbf{dir} + \lambda_2 \mathbf{n} + \lambda_3 \mathbf{r}$$

where \mathbf{n} is the surface normal at a silhouette point and \mathbf{r} is the direction vector from the silhouette point to the camera center. The tangent line at a point on the curve in the image gives us the equation of the plane tangent to the SOR at that point. And thus we can determine \mathbf{n} given a point on the curve. The direction vector \mathbf{r} can be determined by shooting a ray from the camera center through the point on the silhouette. Note that there exists a unique solution to the 3 variables (the λ 's).

Since the camera projection matrix is known, for a given point on the silhouette the corresponding point on the other silhouette at the same height can easily be computed. The radius for the height can be computed by enforcing the constraint that the corresponding points are at the same distance from the axis.

4.2. Frustums

The user specifies the base circle of the frustum in the image. The user also specifies the left and right silhouettes of the

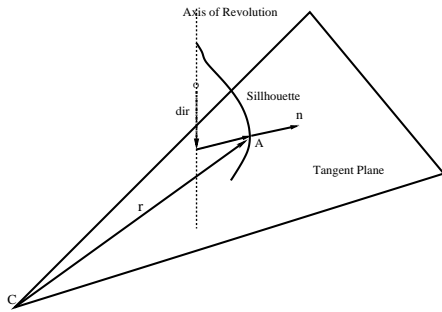


Figure 5: General surfaces of revolution.

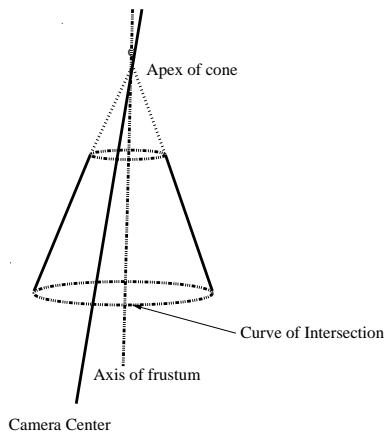


Figure 6: Frustum.

frustums which are lines in the image. The intersection of these two lines in the image is the image of the apex of the cone, of which the frustum is a part as shown in Figure 6. This localizes the apex to the line $\mathbf{P}^+ \mathbf{x}_{\text{apex}} + \lambda \mathbf{C}$ in the real world where \mathbf{P}^+ is the pseudo-inverse of \mathbf{P} (i.e. $\mathbf{P}^+ (\mathbf{P}\mathbf{P}^+)^{-1}$) and \mathbf{C} is the camera center⁵. The intersection of this line with the axis of the base circle determines the apex in the world. The cone is thus determined. To find the height of the frustum, the system projects the line (say ℓ_1) joining the apex to the point nearest to the camera center on the base circle, onto the image. The user then clicks the point of intersection of this projected line with the upper curve of intersection. We shoot a ray through the camera center and this intersection point. The required point has thus been localized as the point of intersection of this ray and ℓ_1 . The distance of this point from the base plane is the height of the frustum.

4.3. Surfaces of translation

To model arches in monuments one also requires surfaces of translation. These are parameterized by splines and the depth of translation. The user clicks points on one of the end planes of the surface of translation. We require that this plane

be known or computed by earlier metrology. Now using \mathbf{P} , we can project these points from the image to the plane to get their location in the world plane. Subsequently, we fit a spline on these points and translate the spline by the required depth to get the model of surface of translation. This depth must also be determined by earlier metrology.

5. Results

We have modelled the Humayun's Tomb, in New Delhi using the primitives described in the previous section. The gross modelling was done from the image in Figure 1.

First the camera center was localized by registering two planes in the image as shown in Figure 1. We specified the quantity βt_y mentioned in Section 3. Recall that this is just the height of the origin on the vertical plane from the horizontal plane in terms of the vertical plane's unit length (refer to Figure 3). This sets up a coordinate system in the image and determines the projection matrix \mathbf{P} of the camera in this coordinate system. We require the lengths of the units in the 3-axes to correct the \mathbf{P} to a Euclidean frame. The lengths in the X and Z directions were directly available from actual measurements, but the Y height was not possible to measure. This did not pose a serious difficulty because we could probe the height of the courtyard in terms of the unit of Y . We then used this height to relate the unit of Y to actual lengths in the world.

Once the height of the courtyard from the ground plane was known, all subsequent height measurements were carried out with respect to the plane of the courtyard. We solved for the common height of the planes together as explained in Section 3 and determined the height of the two side planes and the four angled planes visible in the image. The front plane's height was determined similarly and the rest of the sides of the monument which are not visible in the image were modelled using symmetry.

The error in measurement for the heights shown in Figure 1 is tabulated in Table 1.

Height	Measured Ht	Actual Ht	Error (in %)
H1	46.89	45.9	2.15
H2	21.96	22.0	0.18
H3	4.84	4.75	1.89
H4	58.12	60.0	3.13

Table 1: Error in reconstructed heights (Actual heights are obtained from Archaeological Survey of India)

For modelling the various surfaces of revolutions in the images, the axes and a point on each axis were determined from earlier metrology using symmetry information. Recall

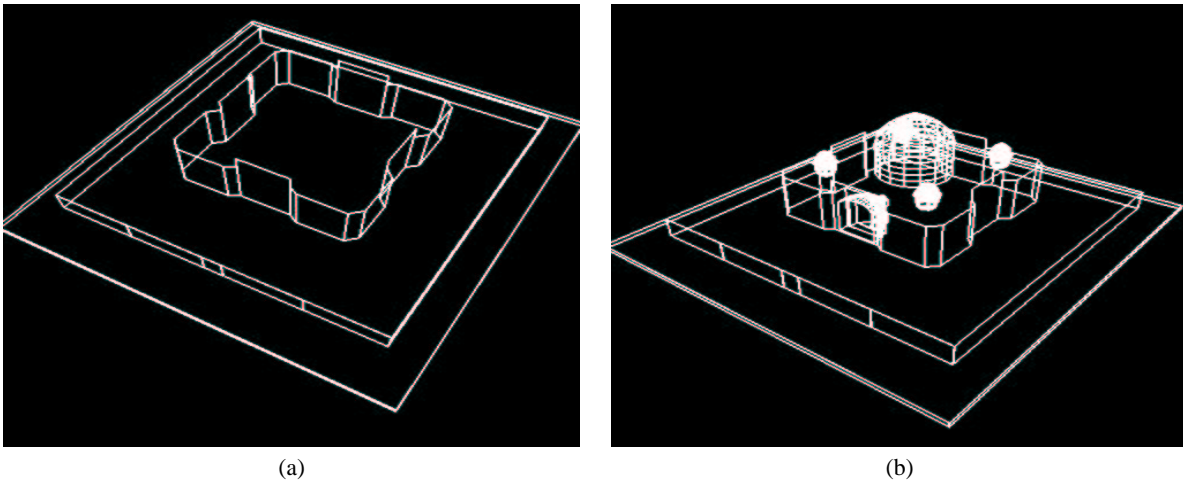


Figure 7: Wire-frame model of (a) the faces, (b) showing the complete structure with the internal details of the front face registered.

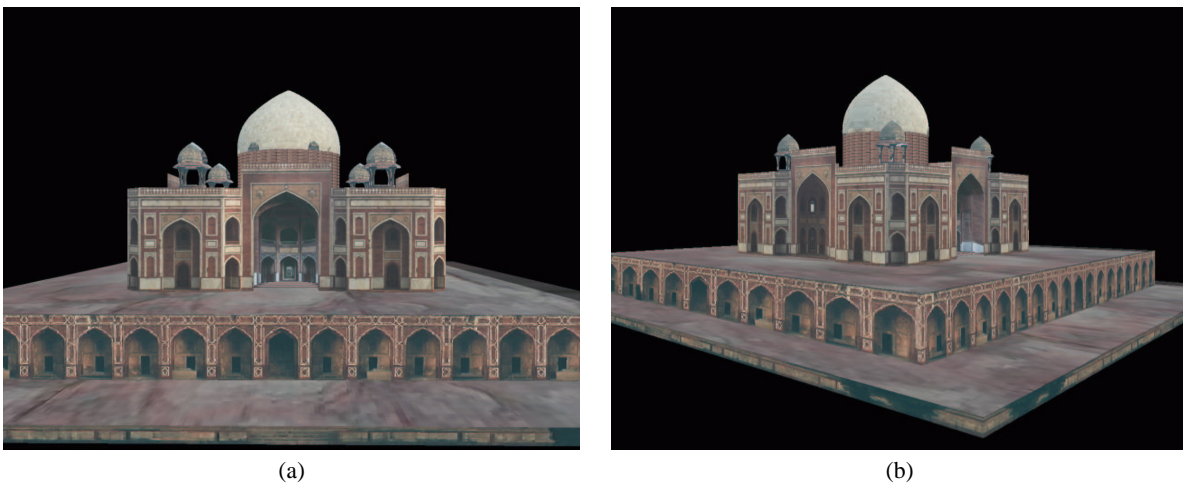


Figure 8: Two snapshots from the final walk-through

that these are necessary for fitting of SOR's (see Section 4). The surfaces of revolution were then modelled from their respective silhouette edges.

We used the image in Figure 2 to carry out a detailed metrology of the area inside the front facade. For this we marked the two rectangles (H and V) on orthogonal planes as shown in Figure 2. Here the scales along the X and Z axis on the new axis were determined from the tie points shown in Figure 2(b). The parallel plane technique was used to move up from the H plane to the front face plane in two steps. The tie point coordinates in this new frame were computed by the probing technique and were equated with the ones obtained from the far metrology to get the scaling between the far and the detailed coordinate systems.

The unit length along the Z axis of the detailed coordinate system could not be obtained from our previous metrology and was measured separately. Then the heights of the planes and other detailed measurements were made on Figure 2(b) and the inner details were included in the model. In Figure 7 we show the model building process in wireframe. In Figure 8 we show two snapshots from the final walk-through.

6. Conclusion

We have presented an intuitive and interactive method for multilevel modelling of architectural scenes. The method is simple and flexible and yields high quality interactive walk-throughs.

Our primitives are more atomic than prismatic blocks and

cuboids. At the same time the basic technique can be used to simultaneously compute parts or whole of the model in terms of higher level primitives like planes, sequence of planes, prismatic blocks, cuboids, spheres, cones and frustums, surfaces of translation and surfaces of revolution using simultaneous bundle adjustment. Consequently, a user has more flexibility to experiment and evolve a right strategy on a case to case basis and decide on whether to use the basic probing technique or to solve for several primitive parameters simultaneously in a tightly constrained manner.

We determine both the camera internal parameters and its location a priori, and hence all subsequent modelling steps are linear and computationally efficient. Further, modelling errors in a portion of the scene can be corrected without affecting the modelling of other parts.

The basic probing technique gives us added flexibility for incremental model building at multiple levels of detail. Since the atomic primitives are just length ratios and not more complex blocks, our method is more suitable for multilevel analysis and modelling.

The rendering techniques involved here can be improved by incorporating issues involved in seamless transitions (which will arise on account of rendering the model at different levels of detail). Furthermore, a more intuitive use of the methods presented here will be facilitated by the design of interfaces tailored to specific user needs. Another important direction of future work includes the incorporation of free-form surfaces into the current framework.

ACKNOWLEDGEMENTS

We are thankful to the Archaeological Survey of India for the help and cooperation they provided during this project. This project was in part supported by a grant issued by the Naval Research Board, India.

References

1. A. Criminisi, I. Reid, and A. Zisserman. Single view metrology. *IJCV*, **40**(2):123–148, 2000. 2, 4, 5
2. P. E. Debevec, C. J. Taylor, and J. Malik. Modeling and rendering architecture from photographs: a hybrid geometry and image based approach. In *Proc. SIGGRAPH*, 1996. 1, 2, 3
3. P. E. Debevec, Y. Yu, and G. D. Borshukov. Efficient view-dependent image-based rendering with projective texture-mapping. In *Eurographics Rendering Workshop*, pages 105–116, 1998. 3
4. E. Grossmann, D. Ortin, and J. Santos-Victor. Single and multi-view reconstruction of structured scenes. In *Proc. ACCV*, 2002. 1, 2
5. R. Hartley and A. Zisserman. *Multiple View Geometry in Computer Vision*. Cambridge University Press, 2000. 1, 2, 4, 6
6. Y. Horry, K. Anjyo, and K. Arai. Tour into the picture: using a spidery mesh interface to make animation from single image. In *Proc. SIGGRAPH*, pages 225–232, 1997. 2
7. A. M. Kushal, V. Bansal, and S. Banerjee. A simple method for interactive 3d reconstruction and camera calibration from a single view. In *Proc. Indian Conference in Computer Vision, Graphics and Image Processing*, 2002. (<http://www.cse.iitd.ernet.in/vglab/demo/single-view/2plane/>). 2, 4
8. D. Leibowitz, A. Criminisi, and A. Zisserman. Creating architectural models from images. In *Proc. Eurographics*, pages 39–50, 1999. 1, 2
9. R. Cipolla and D. Robertson. 3D models of architectural scenes from uncalibrated images and vanishing points. In *Proc. 10th IAPR*, pages 824–829, 1999. 1, 2
10. R. Cipolla, D. Robertson, and E. G. Boyer. Photobuilder - 3D models of architectural scenes from uncalibrated images. In *Proc. IEEE International Conference on Multimedia Computing and Systems*, pages 25–31, 1999. 2
11. P. F. Strum and S. J. Maybank. A method for interactive 3D reconstruction of piecewise planar objects from single views. In *Proc. BMVC*, pages 265–274, 1999. 2
12. M. Wilczkowiak, E. Boyer, and P. Strum. Camera calibration and 3D reconstruction from single images using parallelepipeds. In *Proc. ICCV*, pages 142–148, 2001. 2
13. K.-Y. K. Wong, P. R. S. Mendonça, and R. Cipolla. Reconstruction of surfaces of revolution from single uncalibrated views. In P. L. Rosin and D. Marshall, editors, *Proc. British Machine Vision Conference 2002*, volume 1, pages 93–102, Cardiff, UK, September 2002. British Machine Vision Association. 5