

View Synthesis from Uncalibrated Images Using Parallax

Andrea Fusiello, Stefano Calderer, Sara Ceglie, Nikolaus Mattern and Vittorio Murino

Dipartimento di Informatica, Università degli Studi di Verona

37134 Verona, Italy

{andrea.fusiello,vittorio.murino}@univr.it

Abstract

This work deals with the view synthesis problem, i.e., how to generate snapshots of a scene taken from a “virtual” viewpoint different from all the viewpoints of the real views. Starting from uncalibrated reference images, the geometry of the scene is recovered by means of the relative affine structure. This information is used to extrapolate novel views using planar warping plus parallax correction. The contributions of this paper are twofold. First we introduce an automatic method for specifying the virtual viewpoint based on the replication of the epipolar geometry linking two reference views. Second, we present a method for generating synthetic views of a soccer ground starting from a single uncalibrated image. Experimental results using real images are shown.

1. Introduction

Nowadays we see an increasing interest in the convergence of Computer Vision and Computer Graphics [10]. One of the most promising and fruitful topic is *Image-Based Rendering* (IBR) [12, 8]. While the traditional geometry-based systems use a 3-D model, in IBR views are generated by re-sampling one or more example images, using appropriate warping functions. The advantage is that photographs of real scenes can be used as a basis to create very realistic images.

The warping functions are based on the observation that certain relationships exist between the positions of pixels representing the same points in the scene observed from different viewpoints [3].

In the case of calibrated cameras, algorithms based on image interpolation yield satisfactory results [12, 14]. Uncalibrated techniques, that do not assume any knowledge on the imaging device, utilize image-to-image constraints such as the Fundamental matrix [9], trilinear tensors [1], plane+parallax [6], or homographies [2], to re-project pixels

from a small number of reference images to a given view.

Although uncalibrated point transfer algorithms are well understood, a “natural” way of specifying virtual viewpoints is missing. With an uncalibrated setting, one cannot specify the position and orientation of the virtual camera in the familiar Euclidean frame, because it is not accessible. Everything is specified in a projective frame that is linked to the Euclidean frame by an unknown projective transformation. This means that one has to specify some projective element, like the epipole.

In the first part of this work, we propose an automatic solution for specifying new viewpoints based on the replication of the (unknown) displacement that links two reference views.

In the second part of the paper we focus on the extraction of the information required for the view synthesis starting from a *single, uncalibrated* image.

Few works deal with view synthesis from a *single* image; in this case, additional constraints must be used. For example, the symmetry of human faces is exploited in [14]. We use the knowledge on the soccer ground structure and the fact that the players are in vertical position.

We follow the *relative affine structure* [16] approach, which will be reviewed in the next section. The rest of the paper is structured as follows. Section 4 introduces our first contribution. It is subdivided into two subsections. The first (Sec. 4.1) describe how the relative affine structure is recovered, the second (Sec. 4.2) deals the virtual viewpoint specification and the synthesis. Section 5 describes the synthesis from a *single* image. It is again subdivided in two parts. The recovery of the relative affine structure is discussed first (Sec. 5.1), then the generation of extrapolated views is described (Sec. 5.2).

2. Background

In this section we review some background notions needed to understand the paper. A complete discussion and formulation of the relative affine structure theory, and its

close relative plane+parallax, can be found in [16, 17]. A more general reference on the geometry of multiple views is [4].

Two views of a planar set of points are related via a homography, i.e. a non-singular linear transformation of the projective plane into itself. The most general homography is represented by a non-singular 3×3 matrix H .

If $m_i \in I_1$ and $m'_i \in I_2$ are projection in two different views I_1 and I_2 of the same 3-D point m_i belonging to some plane Π , we have

$$m'_i \sim H_{\Pi} m_i \quad (1)$$

where H_{Π} is the homography induced by plane Π , \sim means “equal up to a scale factor” and points are expressed in homogeneous coordinates.¹ The matrix H_{Π} has eight degrees of freedom, being defined up to a scale factor: four corresponding points in the two views define a homography.

For a general 3-D point m_i , we have

$$m'_i \sim H_{\Pi} m_i + k_i e' \quad (2)$$

where e' denotes the epipole in the second view, and k_i is the *relative affine structure*, which is proportional to the distance of the point P_i from the plane Π (denoted by “ a ” in Fig. 1).

This equation says that points are first transferred as if they were lying on the reference plane Π , and then their position gets corrected by a displacement $k_i e'$, called *parallax*, in the direction of the epipole e' , with magnitude proportional to the relative affine structure k_i . If $P_i \in \Pi$ then $k_i = 0$ and Eq. (2) reduces to Eq. (1).

Given the homography between two views and two off-plane conjugate pairs $(m_0; m'_0)$ and $(m_1; m'_1)$, following simple geometric consideration², the epipole is computed as the intersection between the line containing $H_{\Pi} m_0, m'_0$ and the line containing $H_{\Pi} m_1, m'_1$:

$$e' \sim (H_{\Pi} m_0 \times m'_0) \times (H_{\Pi} m_1 \times m'_1) \quad (3)$$

The geometry of two views in an uncalibrated framework is completely characterised by the Fundamental matrix [11, 4] F , defined by $m'^T F m = 0$. In terms of cameras' parameters it is given by:

$$F = A'^T [t]_{\times} R A^{-1} \quad (4)$$

The Fundamental matrix can be factored as

$$F \sim [e']_{\times} M \quad (5)$$

¹Points in the image plane are denoted as $p = (x_1, x_2, x_3) \sim (\frac{x_1}{x_3}, \frac{x_2}{x_3}, 1)$ with $(u, v) = (\frac{x_1}{x_3}, \frac{x_2}{x_3})$ being the corresponding Cartesian coordinates. The symbol \sim means equality up to a scale factor.

²In the projective plane, the line determined by two points is given by their cross product, as well as the point determined by two lines

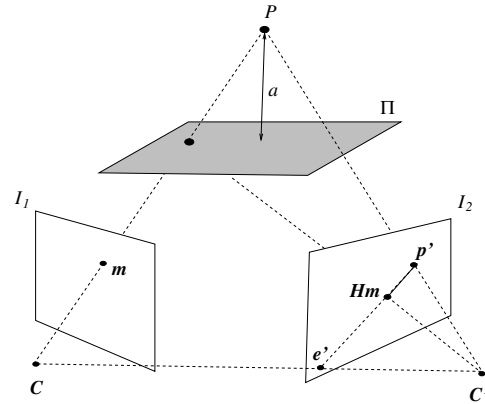


Figure 1. Relative affine structure. The segment joining p' and Hm is the parallax for point P .

Any matrix M that satisfies this factorization is said to be *compatible* (with F). It can be seen that any homography induced by a plane is compatible. Moreover, any compatible homography maps one epipole to its conjugate, namely:

$$e' \sim H_{\Pi} e \quad (6)$$

This relationship can be exploited to compute the epipoles given two compatible homographies. Since $e \sim H_{\Pi}^{-1} e'$, and $e' \sim H_{\Sigma} e$, it follows that:

$$e' \sim H_{\Sigma} H_{\Pi}^{-1} e' \quad (7)$$

The matrix $H_{\Sigma} H_{\Pi}^{-1}$ has three eigenvectors: one is the epipole, the other two belongs to the line $H_{\Sigma} \cap H_{\Pi}$ and are associated to two equal eigenvalues.

The previous observations lead to the conclusion that *any two compatible homographies determine the Fundamental matrix*.

3 View synthesis algorithm

A very important property is that the relative affine structure is independent of the choice of the second view. Therefore, arbitrary “second views” can be synthesized by specifying a plane homography and the epipole. This leads to the following view synthesis algorithm [17]:

1. given a set of conjugate pairs $(m'_i; m_i)$, $i = 0 \dots n$;
2. recover the epipole e' and the homography H_{Π} ;
3. choose a point m_0 and scale H_{Π} to satisfy³

$$m'_0 \sim H_{\Pi} m_0 + e';$$

³The scale factor is computed with a formula analogous to Eq. (8).

4. compute the relative affine structure k_i with⁴

$$k_i = \frac{(\mathbf{H}_{\Pi} \mathbf{m}_i \times \mathbf{m}'_i)^{\top} (\mathbf{m}'_i \times \mathbf{e}')}{\|\mathbf{m}'_i \times \mathbf{e}'\|^2}; \quad (8)$$

5. obtain a new epipole e'' and a new plane homography \mathbf{H}_{Σ} , properly scaled;
6. transfer points in the synthetic view with

$$\mathbf{m}''_i \sim \mathbf{H}_{\Sigma} \mathbf{m}_i + e'' k_i \quad (9)$$

Two problems are to be addressed here: i) how to compute correspondences, and ii) how to specify a new epipole e'' and a new homography \mathbf{H}_{Σ} , which fix the position and orientation of the virtual camera.

Two techniques will be presented: the first is innovative in the way e'' and \mathbf{H}_{Σ} are obtained, the second in the way the relative affine structure is recovered from just one image.

4. Two reference views

In this section we will consider the case of extrapolation from two (or more) uncalibrated reference views. We will present an automatic solution to the specification of the extrapolate viewpoint, based on the replication of the epipolar geometry that links two reference views, considered as an elementary displacement step.

4.1. Relative affine structure recovery

Let us now concentrate on the computation of the relative affine structure in the case that two reference views I_1 and I_2 are available; the extension to the case of more than two views is straightforward. We divide the problem in computing a dominant homography that caters for the motion of the majority of the pixels (usually the background), and a residual parallax.

The homography of the background plane \mathbf{H}_d is obtained as the one that explains the motion of the majority of the pixels in the image: the *dominant motion*. We are here implicitly assuming that the background is approximately planar, or that its depth variation is much smaller than its average distance from the camera. We use a feature-based technique: first, we extract and match corners obtaining a certain number of candidate conjugate pairs. Then, we compute the homography with Least Median of Squares [13], a robust parameter estimation technique that disregards wrong conjugate pairs (*outliers*), which are caused

⁴Eq. (8) can be derived from Eq. (2), given that $\mathbf{H}_{\Pi} \mathbf{m}_i, \mathbf{m}'_i$, and \mathbf{e}' are collinear, since they belong to the same epipolar line. See [17].

either by a wrong matching or by a correct matching of foreground points.

By warping I_1 with the dominant homography \mathbf{H}_d , we obtain another image I_w that (ideally) matches I_2 in those points that lie on the background plane. Therefore, the foreground segmentation can be determined by examining the difference between I_w and I_2 (change detection). To this end, we use the *likelihood ratio* [7] defined as

$$\lambda = \frac{\left[\frac{\sigma_1 + \sigma_2}{2} + \left(\frac{\mu_1 - \mu_2}{2} \right)^2 \right]^2}{\sigma_1 \sigma_2} \quad (10)$$

where μ and σ denote the mean gray value and the variance in a window around the pixel. Thresholding is then applied to the value of λ at each pixel, and the resulting binary image is processed using morphological filtering to remove isolated points and to fill small holes⁵.

From this segmentation we are able to build a mosaic of the background (possibly with holes) and to recover the relative affine structure for the foreground points. To this end, we distinguish two different cases. The first is when foreground is (approximately) planar. In this case we fit a homography \mathbf{H}_f to the foreground points as we did for the background.

When, on the contrary, the foreground is a free form surface, we refine the matching in the foreground region and compute the epipole with Eq. (3). As there are many epipole candidate, a voting technique is employed to select a winner. Finally, we compute the relative affine structure for the inlier foreground points. After scaling \mathbf{H}_d (step 3. of the algorithm), we obtain the value of k_i for each conjugate pair $(\mathbf{m}_i; \mathbf{m}'_i)$ from Eq. (8) (step 4. of the algorithm), and then we interpolate by fitting the set of k with a suitable function.

4.2. View extrapolation

Having extracted from the reference images all the information that are required, we can now use the synthesis equation (9) to construct a synthetic view, but first we need to specify \mathbf{H}_{Σ} and e'' (step 5. of the algorithm), that are projective elements.

Our idea is based on the observation that two compatible homographies completely determine the epipolar geometry of two views (see Sec. 2), and they can be inverted and chained consistently. Therefore, we encode the epipolar geometry of the reference views I_1 and I_2 , with the pair $(\mathbf{H}_d, \mathbf{H}_f)$, where \mathbf{H}_f is a compatible foreground homography. In the case of free-form foreground, \mathbf{H}_f is computed from three foreground points and the epipole.

Suppose we have two reference views I_1 and I_2 , and we want to extrapolate a synthetic view. In the case of planar foreground we simply warp separately the background of I_1

⁵We used the MATLAB “clean” and “fill” morphological operators.

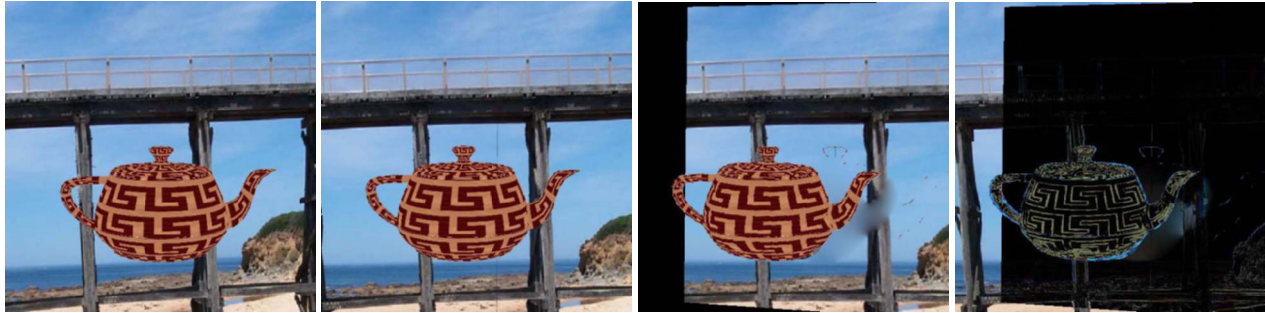


Figure 2. The first two images from left to right are the reference images (generated with OpenGL). The third is the image extrapolated by our algorithm. The last image is the difference with the ground truth.

with $\mathbf{H}_d \mathbf{H}_d$ and the foreground with $\mathbf{H}_f \mathbf{H}_f$. In the case of free-form foreground, we take $\mathbf{H}_\Sigma = \mathbf{H}_d \mathbf{H}_d$, and we compute $\mathbf{H}_\phi \triangleq \mathbf{H}_f \mathbf{H}_f$. The epipole e'' is computed from the eigen-decomposition of $\mathbf{H}_\Sigma \mathbf{H}_\phi^{-1}$ (Sec. 2). Then, assuming that m_0 is one of the points used in the computation of \mathbf{H}_f , we use the point $m_0'' = \mathbf{H}_\phi m_0$ to normalise \mathbf{H}_Σ , and proceed with step 6. of the view synthesis algorithm.

It is easy to be convinced that, in both cases, the new image I_3 is related to I_2 by the same fundamental matrix that links I_2 to I_1 . Indeed, this follows from the observation that we made at the end of Section 2, as \mathbf{H}_d and \mathbf{H}_f are two compatible homographies that transfers points from I_2 to I_3 , by construction.

Hence, provided that intrinsic parameters are constant, the virtual camera is displaced with respect to the second camera by the same rigid displacement that relates the second camera to the first. This follows from the fact that the matrix $\mathbf{E} \triangleq [t]_\times \mathbf{R}$ in Eq. (4), called Essential matrix, admits a unique factorization as a product of a nonzero skew-symmetric matrix (encoding translation) and a rotation matrix [5].

In formulae, if the pose (position and orientation) of the second camera with respect to the first camera is represented by a matrix \mathbf{G} (in homogeneous coordinates), then the pose if the virtual camera with respect to the first camera is given by $\mathbf{G}\mathbf{G}$. In the same way, if we take, for example, $\mathbf{H}_\Sigma = \mathbf{H}_d^{-1}$ and $\mathbf{H}_\phi = \mathbf{H}_f^{-1}$, we obtain a synthetic camera displaced by \mathbf{G}^{-1} from the first one.

Hence, we are able to move the synthetic camera in the space by steps \mathbf{G} and \mathbf{G}^{-1} . If the two reference view-points were displaced approximately horizontally, we could move along the horizontal direction. With a third reference view displaced vertically we could obtain synthetic images from above, below, left, right and any combination of them. The user just need to specify, in a graphical way, the direction toward which the virtual camera must move, and the

system computes automatically the required epipole and homography.

Image warping was performed using destination scan and bilinear interpolation for background and planar foreground, and source scan and pixel splatting [15] for free-form foreground.

5. One reference view

This section describes a method for generating extrapolated views of a soccer ground starting from a single uncalibrated reference image. The relative affine structure framework is particularly well suited for our case: the reference plane is the soccer ground, and the off-plane points are the players heads. The relative affine structure of the players is computed by exploiting the knowledge of the soccer ground geometry and the fact that the players are in vertical positions.

5.1. Relative affine structure recovery

Player silhouettes have to be extracted in order to find their position onto the soccer ground. We employ a simple color-based segmentation: non-green regions are labeled as *potential* players, then their size and shape descriptors are used to discard small areas (noise) and terraces in background. We assume that the players are in a green background (the soccer ground), and there are not intersections between the silhouettes (occlusions).

In order to synthesize geometrically correct arbitrary views we need to first estimate the relative affine structure of the players with respect to the ground plane. This could be easily done if one had two or more reference images (like in Sec. 4). Our key idea is to *synthesize* a second



Figure 3. From left to right. The two reference views, the extrapolated views (four steps) with free-form foreground and planar foreground.

reference view taken from the zenith⁶ of the soccer ground, by exploiting the knowledge of its layout.

The homography matrix H_z between the observed image and the zenithal view is estimated given four point and/or line correspondences. The GUI shows a layout of the soccer ground as seen from the zenith. The user specifies correspondences between the observed image and the schematic one. By applying H_z on the observed image, a synthetic zenithal view is obtained. Supposing that the players are oriented vertically and that viewing is approximately orthographic, the position of the players' head in the *correct* zenithal view coincides with their feet. Instead, in the synthetic zenithal view, as we disregarded the the players' 3D structure, they appear flattened onto the ground plane, and the segment joining the head and the feet is exactly the parallax.

The epipole e' is computed – as described in Section 2 – from the homography of the ground plane and the parallax of two players.

Then we scale H_z (step 3. of the algorithm) and, finally, we compute the relative affine structure k_i for all the players' heads (step 4. of the algorithm). This operation requires the user to specify head - feet correspondences in the synthetic zenithal view.

5.2. View extrapolation

The user specifies the new homography and epipole (step 5. of the algorithm) through GUI. A schematic representation of the soccer ground is shown to the user. He/she can rotate and translate a virtual camera in order to decide the new point of view . Then he/she specifies the correspondences between the observed image and the schematic view needed to compute the homography H_Σ . The positions of the players' feet are shown in the extrapolated view and the user is prompted to input the position of the corresponding heads of two players, which are used to estimate the new

⁶The zenithal view is an image of the soccer ground taken from above with the optical axis orthogonal to the ground plane.

epipole, as described previously. The scale factor for H_Σ is also estimated as before, setting m_0 to be one of the heads of the players used for the epipole.

Given the new epipole e'' and a new homography H_Σ , the position of the players' heads in the extrapolated view are computed from Eq. (9) (step 6. of the algorithm), using the relative affine structure k_i , whereas feet positions are recovered by setting $k_i = 0$. Finally, players' billboards are inserted in the extrapolated view with the correct height.

6. Results

We performed tests with synthetic and real images. The synthetic images are shown in Figure 2. We generated (with OpenGL) three images related by the same displacement of the viewing position. The first two are taken as reference images, the third as ground truth. Our extrapolated view should ideally coincide with it. The difference image confirm that the error is limited to few pixels.

As an example of view synthesis with real images, Figure 3 shows the extrapolated image after four displacement steps, and the corresponding extrapolated image obtained assuming a planar foreground. Please note how this extrapolation involves quite a dramatic shift in the viewing position. For this reason, the planar approximation for the foreground produces an appreciable distortion.

Figure 4 shows an example of synthetic view, in which the new point of view allows us to find the players in offside. Please note as the players appears correctly foreshortened in the synthetic view.

More examples can be found on the web at <http://www.sci.univr.it/~fusiello/demo/synth>.

7. Conclusion

In the first part of this paper we introduced a method for specifying the virtual camera position in an uncalibrated manner by replicating the epipolar geometry that links the model views, considered as an elementary displacement

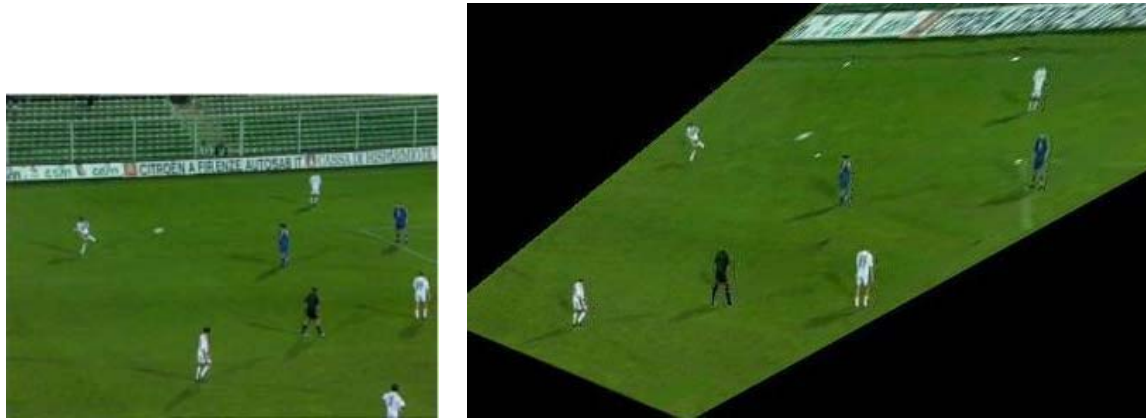


Figure 4. Actual image (left) and synthetic view (right)

step. The virtual viewpoint is not constrained to lie in between the positions of the real cameras. View synthesis examples are shown starting from uncalibrated reference views.

Future work will address the issue of compatible homography interpolation, will improve the segmentation and the forward warping in order to produce better quality images.

In the second part we introduced a method for generating extrapolated views of a soccer ground, starting from a single uncalibrated image. The trick is to obtain a virtual second reference image by generating a view from the zenith, using the knowledge of the ground geometry and the fact that players are vertical.

Further works will address the improvement of the segmentation using classification techniques and the reduction of user intervention.

Acknowledgements

Thanks to Andrea Colombari and Umberto Castellani for helping with the experiments.

References

- [1] S. Avidan and A. Shashua. Novel view synthesis in tensor space. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1034–1040, 1997.
- [2] B. S. Boufama. The use of homographies for view synthesis. In *Proceedings of the International Conference on Pattern Recognition*, pages 563–566, 2000.
- [3] O. D. Faugeras and L. Robert. What can two images tell us about a third one? In *Proceedings of the European Conference on Computer Vision*, pages 485–492, Stockholm, 1994.
- [4] R. Hartley and A. Zisserman. *Multiple view geometry in computer vision*. Cambridge University Press, 2000.
- [5] T. Huang and O. Faugeras. Some properties of the E matrix in two-view motion estimation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 11(12):1310–1312, December 1989.
- [6] M. Irani and P. Anandan. Parallax geometry of pairs of points for 3D scene analysis. In *Proceedings of the European Conference on Computer Vision*, pages 17–30, 1996.
- [7] R. Jain, R. Kasturi, and B. Schunk. *Machine Vision*. Computer Science Series. McGraw-Hill International Editions, 1995.
- [8] S. B. Kang. A survey of image-based rendering techniques. Technical Report CRL 97/4, Digital Cambridge Research Laboratories, August 1997.
- [9] S. Laveau and O. Faugeras. 3-D scene representation as a collection of images and fundamental matrices. Technical Report 2205, INRIA, Institut National de Recherche en Informatique et en Automatique, February 1994.
- [10] J. Lengyel. The convergence of graphics and vision. *IEEE Computer*, 31(7):46–53, July 1998.
- [11] Q.-T. Luong and O. D. Faugeras. The fundamental matrix: Theory, algorithms, and stability analysis. *International Journal of Computer Vision*, 17:43–75, 1996.
- [12] L. McMillan and G. Bishop. Plenoptic modeling: An image-based rendering system. In *SIGGRAPH 95 Conference Proceedings*, pages 39–46, Aug. 1995.
- [13] P. J. Rousseeuw and A. M. Leroy. *Robust regression & outlier detection*. John Wiley & sons, 1987.
- [14] S. M. Seitz and C. R. Dyer. View morphing: Synthesizing 3D metamorphoses using image transforms. In *SIGGRAPH 96 Conference Proceedings*, pages 21–30, Aug. 1996.
- [15] J. Shade, S. Gortler, L. He, and R. Szeliski. Layered depth images. In *SIGGRAPH 98 Conference Proceedings*, 1998.
- [16] A. Shashua and N. Navab. Relative affine structure: Theory and application to 3-D reconstruction from perspective views. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 483–489, 1994.
- [17] A. Shashua and N. Navab. Relative affine structure: Canonical model for 3D from 2D geometry and applications. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 18(9):873–883, September 1996.