

# Integrating Automated Range Registration with Multiview Geometry for the Photorealistic Modeling of Large-Scale Scenes

Ioannis Stamos · Lingyun Liu · Chao Chen ·  
George Wolberg · Gene Yu · Siavash Zokai

Received: 19 March 2007 / Accepted: 11 September 2007 / Published online: 10 November 2007  
© Springer Science+Business Media, LLC 2007

**Abstract** The photorealistic modeling of large-scale scenes, such as urban structures, requires a fusion of range sensing technology and traditional digital photography. This paper presents a system that integrates automated 3D-to-3D and 2D-to-3D registration techniques, with multiview geometry for the photorealistic modeling of urban scenes. The 3D range scans are registered using our automated 3D-to-3D registration method that matches 3D features (linear or circular) in the range images. A subset of the 2D photographs are then aligned with the 3D model using our automated 2D-to-3D registration algorithm that matches linear features between the range scans and the photographs. Finally, the 2D photographs are used to generate a second 3D model of the scene that consists of a sparse 3D point cloud, produced by applying a multiview geometry (structure-from-motion) algorithm directly on a sequence of 2D photographs. The last part of this paper introduces a novel algorithm for automatically recovering the rotation, scale, and translation that

best aligns the dense and sparse models. This alignment is necessary to enable the photographs to be optimally texture mapped onto the dense model. The contribution of this work is that it merges the benefits of multiview geometry with automated registration of 3D range scans to produce photorealistic models with minimal human interaction. We present results from experiments in large-scale urban scenes.

**Keywords** Range segmentation · Range-to-range registration · Range-to-image registration · Multiview geometry · Structure from motion · Photorealistic modeling

## 1 Introduction

The photorealistic modeling of large-scale scenes, such as urban structures, requires a combination of range sensing technology with traditional digital photography. A systematic way for registering 3D range scans and 2D images is thus essential. This paper presents a system that integrates multiview geometry and automated 3D registration techniques for texture mapping 2D images onto 3D range data. The novelty of our approach is that it exploits all possible relationships between 3D range scans and 2D images by performing 3D-to-3D range registration, 2D-to-3D image-to-range registration, and structure from motion from a sequence of 2D images. Several papers, including this one, provide frameworks for automated texture mapping onto 3D range scans (Ikeuchi 2003; Liu and Stamos 2005; Stamos and Allen 2001; Troccoli and Allen 2004; Zhao et al. 2005). These methods are based on extracting features (e.g., points, lines, edges, rectangles or rectangular parallelepipeds) and matching them between the 2D images and

---

I. Stamos (✉) · L. Liu · C. Chen  
Department of Computer Science, Hunter College/CUNY, New York, NY, USA  
e-mail: istamos@hunter.cuny.edu

L. Liu  
e-mail: lingyun77@gmail.com

C. Chen  
e-mail: cchen@gc.cuny.edu

G. Wolberg · G. Yu  
Department of Computer Science, City College of New York/CUNY, New York, NY, USA

G. Wolberg  
e-mail: wolberg@cs.ccny.cuny.edu

S. Zokai  
Brainstorm Technology LLC, New York, NY, USA

the 3D range scans. Our approach provides a solution of increased robustness, efficiency and generality with respect to previous methods. Our contribution is discussed in Sect. 2.

Despite the advantages of feature-based texture mapping solutions, most systems that attempt to recreate photorealistic models do so by requiring the manual selection of features among the 2D images and the 3D range scans, or by rigidly attaching a camera onto the range scanner and thereby fixing the relative position and orientation of the two sensors with respect to each other (Früh and Zakhor 2003; Pulli et al. 1998; Sequeira and Concalves 2002; Visual Information Technology Group 2005; Zhao and Shibasaki 2003). The fixed-relative position approach provides a solution that has the following major limitations:

1. The acquisition of the images and range scans occur at the same point in time and from the same location in space. This leads to a lack of 2D sensing flexibility since the limitations of 3D range sensor positioning, such as standoff distance and maximum distance, will cause constraints on the placement of the camera. Also, the images may need to be captured at different times, particularly if there were poor lighting conditions at the time that the range scans were acquired.
2. The static arrangement of 3D and 2D sensors prevents the camera from being dynamically adjusted to the requirements of each particular scene. As a result, the focal length and relative position must remain fixed.
3. The fixed-relative position approach cannot handle the case of mapping historical photographs on the models or of mapping images captured at different instances in time. These are capabilities that our method achieves.

In summary, fixing the relative position between the 3D range and 2D image sensors sacrifices the flexibility of 2D image capture. Alternatively, methods that require manual interaction for the selection of matching features among the 3D scans and the 2D images are error-prone, slow, and not scalable to large datasets. Laser range scanning is a laborious, tedious, time consuming, and costly operation that precludes easy and cost-effective recapturing of data. Therefore, it is best to separate the geometry capture mode from the image acquisition mode so that the latter can be done quickly with fairly constant lighting conditions. These limitations motivate the work described in this paper, making it essential for producing photorealistic models of large-scale urban scenes.

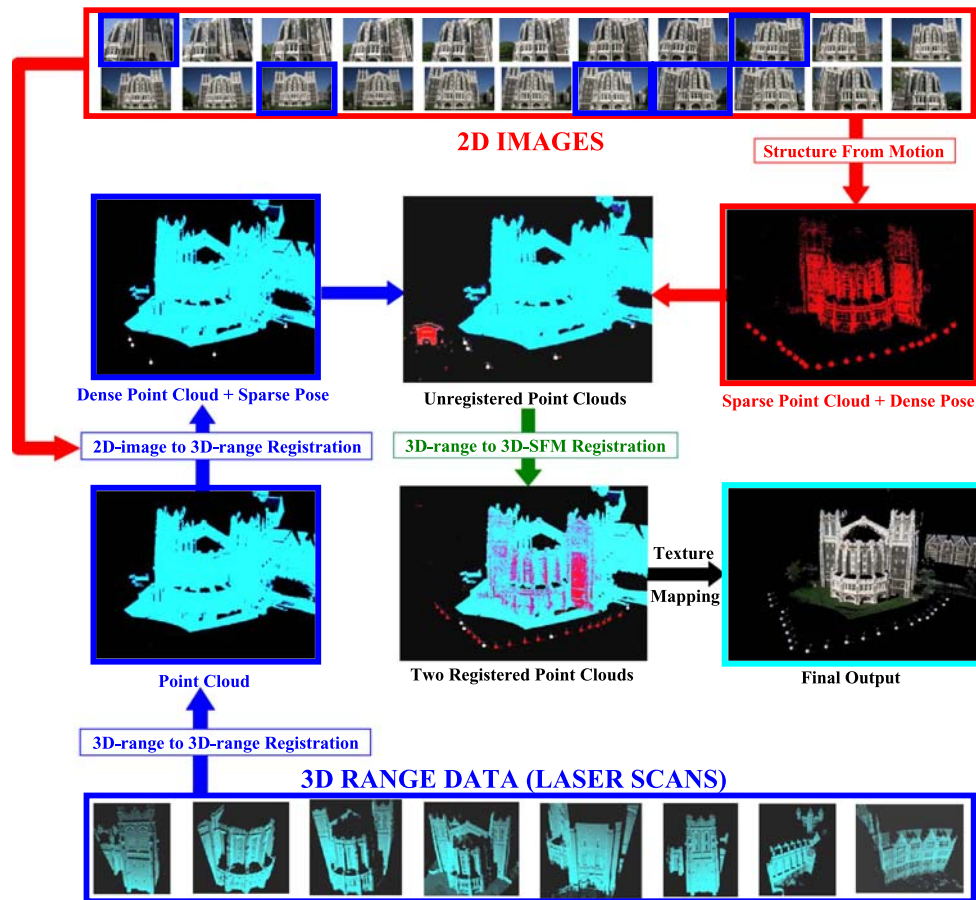
The texture mapping solution described in this paper merges the benefits of multiview geometry with automated 3D-to-3D range registration and 2D-to-3D image-to-range registration to produce photorealistic models with minimal human interaction. The 3D range scans and the 2D photographs are respectively used to generate a pair of 3D models of the scene. The first model consists of a dense 3D point

cloud, produced by using a 3D-to-3D registration method that matches 3D lines in the range images to bring them into a common reference frame. The second model consists of a sparse 3D point cloud, produced by applying a multiview geometry (structure-from-motion) algorithm directly on a sequence of 2D photographs to simultaneously recover the camera motion and the 3D positions of image features. This paper introduces a novel algorithm for automatically recovering the similarity transformation (rotation/scale/translation) that best aligns the sparse and dense models. This alignment is necessary to enable the photographs to be optimally texture mapped onto the dense model. No a priori knowledge about the camera poses relative to the 3D sensor's coordinate system is needed, other than the fact that one image frame should overlap the 3D structure (see Sect. 4). Given one sparse point cloud derived from the photographs and one dense point cloud produced by the range scanner, a similarity transformation between the two point clouds is computed in an automatic and efficient way. The framework of our system is shown in Fig. 1. Each of the framework elements listed below, is a distinct system module in Fig. 1.

- A set of 3D range scans of the scene is acquired and co-registered to produce a dense 3D point cloud in a common reference frame (Sect. 3).
- An independent sequence of 2D images is gathered, taken from various viewpoints that do not necessarily coincide with those of the range scanner. A sparse 3D point cloud is reconstructed from these images by using a structure-from-motion (SFM) algorithm (Sect. 5).
- A *subset* of the 2D images are automatically registered with the dense 3D point cloud acquired from the range scanner (Sect. 4).
- Finally, the *complete* set of 2D images is automatically aligned with the dense 3D point cloud (Sect. 6). This last step provides an integration of all the 2D and 3D data in the same frame of reference. It also provides the transformation that aligns the models gathered via range sensing and computed via structure from motion.

## 2 Related Work

A robust method that extracts distinguishable features from range images is very important for our method. Previous range image segmentation techniques include edge detection (Bellon and Silva 2002; Wami and Batchelor 1994), region growing (Besl and Jain 1988; Pulli and Pietikäinen 1993), and polynomial surface fitting (Besl and Jain 1988; Marshall and Martin 2001). Most of these methods provide edge maps and/or regions expressed as polynomial functions. This is useful for object modeling and reconstruction, but may not be suitable for feature matching. Our method



**Fig. 1** System framework. Several registered range scans of Shepard Hall (CCNY) constitute a dense 3D point cloud model  $M_{range}$  shown in the *leftmost column*. The *five white dots* correspond to the locations of five of the 26 color images (shown as thumbnails on top row) that are independently registered with the model  $M_{range}$  via a 2D-to-3D image-to-range registration algorithm. The *rightmost image of the second row* depicts the 3D model  $M_{sfm}$  produced by SFM. The points of  $M_{sfm}$  as well as *all* the recovered camera positions for the sequence of 2D images that produced  $M_{sfm}$  are shown as *red dots* in the figure.

Since SFM does not recover scale,  $M_{range}$  and  $M_{sfm}$  are not registered when brought to the same coordinate system, as shown in the *second row*. The 3D range model  $M_{range}$  overlaid with the 3D model  $M_{sfm}$  is shown in the *third row* of the figure after a 3D-range to 3D-SFM registration module aligns them together. The recovered camera positions from SFM can now be used to project the 26 color images onto  $M_{range}$ , which now properly sits in the  $M_{sfm}$  coordinate system, to produce the richly textured 3D model (Final Output) shown in the *right column*.

detects precise edges and extracts geometric features with concise descriptors that make them appropriate for feature matching.

Iterative Closest Point (ICP) is one of the most popular range registration algorithms (Besl and McKay 1992; Rusinkiewicz and Levoy 2001). ICP provides very accurate results but requires a good initial guess of the registration transformation. We, on the other hand, utilize ICP as a post-processing step after our automated method brings scans into alignment. A method that does not require knowledge of an initial registration transformation is presented in (Huber and Hebert 2003; Johnson and Hebert 1999) (spin images). The spin images approach does not rely on features of specific geometric type, but is sensitive to varying scan resolutions. Furthermore, the extracted point signatures have local support, the extent of which is specified by the user.

There are many approaches for the solution of the pose estimation problem from both point correspondences (Oberkampff et al. 1996; Quan and Lan 1999) and line correspondences (Christy and Horaud 1999; Horaud et al. 1997), when a set of matched 3D and 2D points or lines are known, respectively. In the early work of (Fischler and Bolles 1981), the probabilistic RANSAC method was introduced for automatically computing matching 3D and 2D points. RANSAC is a robust method and can handle a large number of outliers. It is possible that the method may fail if presented with degenerate data configurations. In such cases, RANSAC may fit the model to outliers instead of inliers. Solutions in automated matching of 3D with 2D features in the context of object recognition and localization include (Cass 1997; Hausler and Ritter 1999; Huttenlocher and Ullman 1990; Jacobs 1997; Jurie 1999; Wells 1997). Very few meth-

ods, though, attack the problem of automated alignment of images with dense point clouds derived from range scanners. This problem is of major importance for automated photorealistic reconstruction of large-scale scenes from range and image data. In (Liu and Stamos 2005; Stamos and Allen 2001) two methods that exploit orthogonality constraints (rectangular features and vanishing points) in man-made scenes are presented. The methods can provide excellent results, but will fail in the absence of a sufficient number of linear features. Ikeuchi (2003), on the other hand, presents an automated 2D-to-3D registration method that relies on the reflectance range image. However, the algorithm requires an initial estimate of the image-to-range alignment in order to converge. Finally, Troccoli and Allen (2004) presents a method that works under specific outdoor lighting situations.

A system whose goals are very similar to ours is described in (Zhao et al. 2005). In that work, continuous video is aligned onto a 3D point cloud obtained from a 3D sensor. First, an SFM/stereo algorithm produces a 3D point cloud from the video sequence. This point cloud is then registered to the 3D point cloud acquired from the range scanner by applying the ICP algorithm (Besl and McKay 1992). One limitation of this approach has to do with the shortcomings of the ICP algorithm. In particular, the 3D point clouds must be brought close to each other manually to yield a good initial estimate that is required for the ICP algorithm to work. The ICP may fail in scenes with few discontinuities, such as those replete with planar or cylindrical structures. Also, in order for the ICP algorithm to work, a very dense model from the video sequence must be generated. This means that the method of Zhao et al. (2005) is restricted to video sequences, which limits the resolution of the 2D imagery. Finally, that method does not automatically compute the difference in scale between the range model and the recovered SFM/stereo model.

Our contributions can be summarized as follows:

- We automatically register the 3D range scans by matching linear and circular features.
- Like Zhao et al. (2005), we compute a model from a collection of images via SFM. Our method for aligning the range and SFM models, described in Sect. 6, does not rely on ICP and thus does not suffer from its limitations.
- We are able to automatically compute the scale difference between the range and SFM models.
- We perform 2D-to-3D image-to-range registration (Liu 2007; Liu and Stamos 2005, 2007) for a few (at least one) images of our collection. This feature-based method provides excellent results in the presence of a sufficient number of linear features. Therefore, the images that contain enough linear features are registered using that method. The utilization of the SFM model allows us to align the remaining images with a method that involves robust point (and not line) correspondences.

- We generate an optimal texture mapping result by using contributions of all 2D images.

### 3 3D-to-3D Range Registration

The first step is to acquire a set of range scans  $S_i$  ( $i = 1, \dots, K$ ) that adequately covers the 3D scene. The laser range scanner used in our work is a Leica HDS 2500 (Leica Geosystems 2007), an active sensor that sweeps an eye-safe laser beam across the scene. It is capable of gathering one million 3D points at a maximum distance of 100 meters with an accuracy of 6 mm at 50 m. Each 3D point is associated with four values  $(x, y, z, l)^T$ , where  $(x, y, z)^T$  is its Cartesian coordinates in the scanner's local coordinate system, and  $l$  is the laser intensity of the returned laser beam.

Each range scan then passes through an automated segmentation algorithm (Stamos and Allen 2002) to extract a set of major 3D planes and a set of geometric 3D lines  $G_i$  from each scan  $i = 1, \dots, K$ . The geometric 3D lines are computed as the intersections of segmented planar regions and as the borders of the segmented planar regions. In addition to the geometric lines  $G_i$ , a set of reflectance 3D lines  $L_i$  are extracted from each 3D range scan. The range scans are registered in the same coordinate system via the automated 3D-to-3D feature-based range-scan registration method of (Chen and Stamos 2005; Stamos and Leordeanu 2003). The method is based on an automated matching procedure of linear features of overlapping scans. As a result, all range scans are registered with respect to one selected pivot scan. The set of registered 3D points from the  $K$  scans is called  $M_{\text{range}}$  (Fig. 1).

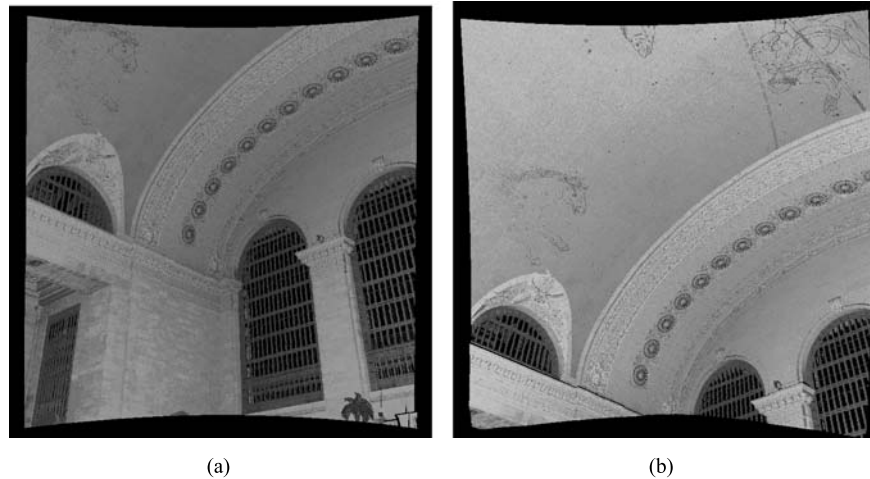
Despite the many linear features prevalent in large-scale urban scenes, there are some cases in which the extracted linear features are inadequate to register the range scans together. Nonlinear features such as circular arcs may be more appropriate in those cases. This is particularly true of range scans taken of curved ceilings with arches and round windows, as shown in Fig. 2. We have therefore developed a circle-based registration method (Chen and Stamos 2006; Chen 2007) that we describe in the following sections.

#### 3.1 3D Edge Detection

Each range scan  $S_i$  can be represented as a 2D array of 3D points  $\{\mathbf{P}(k, l), k = 1, \dots, W, l = 1, \dots, H\}$ .<sup>1</sup> Two such range images are shown in Fig. 2. Within each range image we consider 3D edges of the following two types: (a) edges caused by surface normal discontinuities (roof edges), and (b) edges caused by depth discontinuities (step edges). Step

<sup>1</sup>The indices  $k, l$  define the position and orientation of the laser-beam which produces the 3D point  $\mathbf{P}(k, l)$ .

**Fig. 2** Two range images of the interior of the Grand Central Terminal, NYC



edges are further divided into edges caused by one surface occluding another (occlusion edges), and edges caused by 3D surface boundaries (boundary edges).

We briefly summarize the algorithm for detecting edges of various types. First the surface orientation change at each point is decomposed into variations along four grid directions. This grid is the 2D structured grid on which each range image is defined. We thus obtain four values at every 3D point, that we call directional variation values. In the structured 2D grid we form four 2D images. The intensity value at each pixel is the surface variation (we define it properly in the next paragraphs) of the corresponding 3D point. We call the four 2D images directional variation images. 2D Canny-based edge detection is performed on each image. Finally the 2D edges are combined and projected<sup>2</sup> back to the 3D range image space, providing the final 3D edges due to surface normal discontinuities or depth discontinuities.

The directional variation images are obtained as follows: At each point  $P$ , let  $B_1$  and  $B_2$  be its two neighbors along one of the four grid directions (see Fig. 3(e)). The vector from  $P$  to  $B_1$  is  $\mathbf{V}_1$ , and from  $P$  to  $B_2$  is  $\mathbf{V}_2$ . The variation at each direction for point  $P$  is defined as  $Angle(\mathbf{V}_1, \mathbf{V}_2)/\pi$ . This provides a value in  $(0, 1]$  as the intensity value for this 2D directional variation image.

2D edge detection is performed on each of the four directional variation images. First, Gaussian smoothing is applied to suppress noise. Then, gradients along  $x$  and  $y$  direction,  $g_x$  and  $g_y$ , are computed at each pixel using Sobel operators. With  $g_x$  and  $g_y$  we compute the gradient magnitude  $g$ , and determine the edge direction at each point as one of the following: horizontal, vertical, positive diagonal and negative diagonal. Our algorithm then carries out hysteresis thresholding followed by non-maximum suppression to obtain thin continuous edges.

<sup>2</sup>Each pixel  $p(k, l)$  in the grid-point image corresponds to a 3D point  $\mathbf{P}(k, l)$ .

To this point, we have detected all roof and step edges. However, occlusion edges need to be identified and only the foreground edges should be kept in order to reflect the true geometry of the scene (similar to the shadows in 2D images). The earlier step of non-maximum suppression votes off edge points based on magnitude, regardless of whether it is on a foreground surface or a background surface. To find occlusion edges, we map the 2D edge points back to the 3D range scan and label a 3D point  $P$  if its corresponding pixel  $p$  is an edge point. For each edge point  $P$ , we check its two neighbors perpendicular to its edge direction. If one of these neighbors is much closer to the scanner and is not an edge point, we mark that neighbor to be a foreground edge point and mark  $P$  as non-edge.

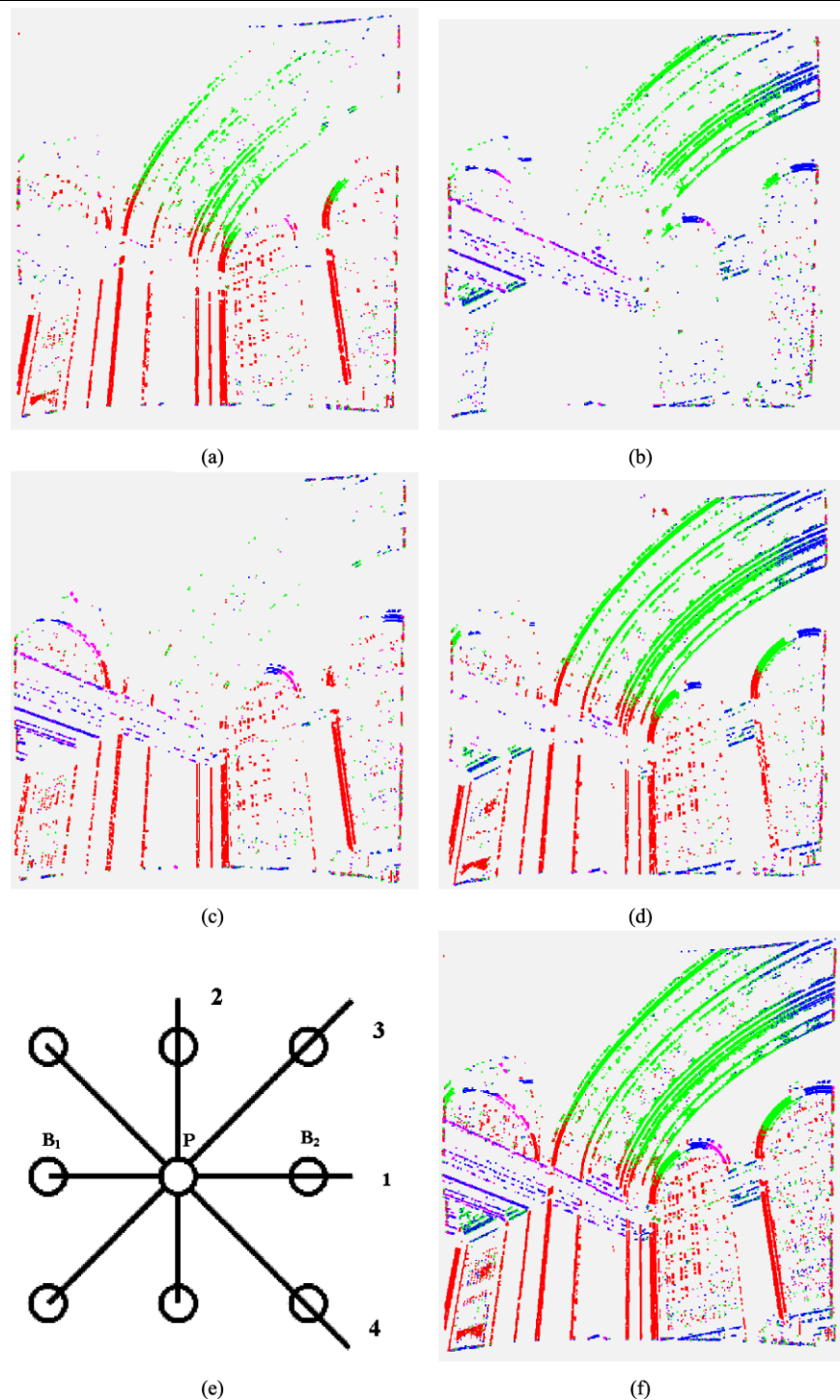
Next we remove corner points in order to break the connections between edges of different directions, thereby simplifying edge linking and fitting. Corner points are detected by applying the Harris corner detector to every edge point, and testing whether there are more than one principle directions formed by all edge points in its local neighborhood.

The last step is the combination of four edge maps by taking the union of all edge points. From the combined edge map, isolated edge points are deleted, and short gaps (1 or 2 pixels) are filled along the local edge direction. Then continuous edge points are linked by tracing along edge directions. The edge linking utilizes the structured grid on which the range image is represented for resolving neighbors. Only long edges (30 points or more) are kept for later processing. Figure 3(f) shows the final combined edge map, in which different colors indicate different edge directions, each detected by one directional variation image.

### 3.2 3D Circular Feature Extraction

Our linear feature extraction is described in (Chen and Stamos 2005; Stamos and Leordeanu 2003). In this section we

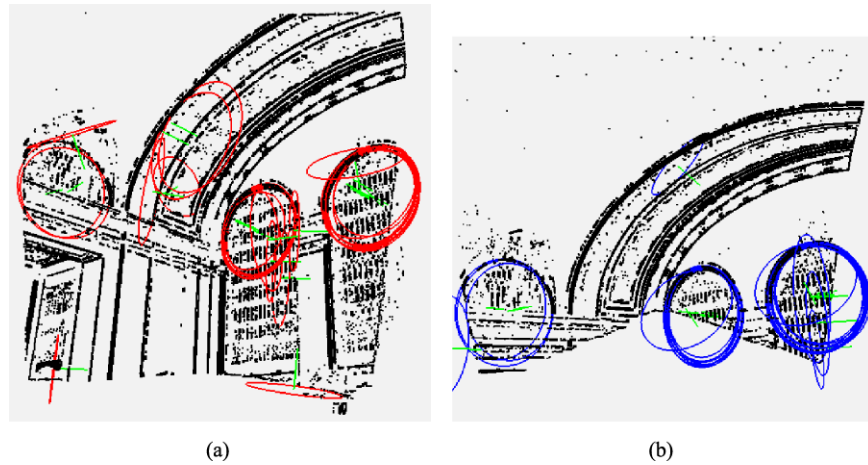
**Fig. 3** Edge points of range images of Fig. 2(a). Note that the color at each point (red/green/blue/magenta) indicates its edge direction (see text), hence the same point usually has the same color in the four edge images. **(a–d)** Edge points detected from Fig. 2(a). **(e)** Four grid directions: 1-Horizontal, 2-Vertical, 3-Positive Diagonal, 4-Negative Diagonal.  $B_1$  and  $B_2$  are  $P$ 's neighbors along direction 1. **(f)** Combined edge image from **(a)–(d)**



present our algorithms for detecting circular features in 3D space. This is necessary in those cases where linear features are inadequate to register the circular arcs that may be present in the scene. Although the ellipse also exists in our scenes, it is less robust to extract due to its feature of having two foci. Any error or noise in edge point extraction might greatly affect the parameters of the fitted ellipse, and further lead to incorrect computation on transformation.

Each linked set of edges describes a curve in 3D space. For each linked edge from Fig. 3(e), its best-fit line direction  $\mathbf{V}_{\max}$  and best-fit plane normal  $\mathbf{V}_{\min}$  are computed. A curve is considered linear if the line fitting error (average distance of all points to the fitted line) is less than a threshold 0.03 m (approximate distance between two neighboring 3D range points). For nonlinear curves, the average perpendicular distance of all points to the fitted plane is used to discard

**Fig. 4** Circles extracted from the range images of Figs. 2(a) and (b), respectively. All edge points are in black, and all fitted circular arcs are represented with colored full circles, with green lines indicating their normals. Note that the three circular windows are detected in both images. The images are rotated to the best angle to observe all circles. Therefore some of them appear as ellipses due to the viewing direction



3D curves that are nonplanar (a generous threshold of 0.5 m is used). For each of the remaining planar curves, all points are projected onto their fitted plane. After this process, the 3D curve becomes a set of 2D points in the 2D space  $\Pi$  of the fitted plane. Circle fitting is done in this space.

Taking the common approach of least square fitting, we compute the center  $(a, b)$  and radius  $r$  of the circle by finding an approximate null-vector of a  $n \times 4$  design matrix, where  $n$  is the number of points on the curve. Consider the circle function  $(x - a)^2 + (y - b)^2 - r^2 = 0$ . It can be written as  $x^2 + y^2 - 2ax - 2by + a^2 + b^2 - r^2 = 0$ . Let  $(x_i, y_i)$  be the 2D coordinates of all points  $p_i (i = 1, \dots, n)$  on the curve. Then the circle equation for all points can be expressed as a multiplication of the  $n \times 4$  matrix  $M = [M_1 \ M_2 \ \dots \ M_n]^T$  where  $M_i = [x_i^2 + y_i^2 \ -2x_i \ -2y_i \ 1]$  (for  $i = 1, \dots, n$ ), with unknown vector  $[1 \ a \ b \ a^2 + b^2 - r^2]^T$ . The null-vector of the design matrix, computed by SVD, provides the solution. Finally, the circle fitting error is computed as  $c_{err} = \sqrt{\frac{\sum_i^n (\text{distance}(p_i - \text{center}) - r)^2}{n}}$ . The ratio  $(\frac{c_{err}}{r})$  must fall below a threshold (0.02) to verify that the planar 3D curve is a circular arc. Finally, the center of the fitted circle is converted back from  $\Pi$  to the 3D space. We now have three parameters to represent each oriented 3D circle: 3D center point, radius, and plane normal. Figure 4 shows all the circles with radii between 3.0 m and 5.0 m. These are the ones most useful for matching in the next step. In the execution, we detect all circles with radii between 2.0 m and 20.0 m.

### 3.3 Feature Matching

Our linear feature matching is described in (Chen and Stamos 2005; Stamos and Leordeanu 2003). In this section we present our algorithms for matching circular features in 3D space.

After the oriented 3D circles are extracted from range images, possible matchings between them are hypothesized.

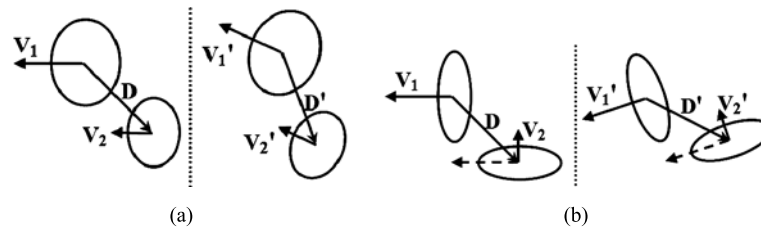
The computed transformations are graded using surface consistency (Huber and Hebert 2003) and average point distance in the overlapping area between the scans.

Similarity of radii, orientation and relative position between pairs of circles is utilized in the matching phase. In particular, consider a pair of circles  $(C_1, C_2)$  from scan  $S_1$  and another pair of circles  $(C'_1, C'_2)$  from scan  $S_2$ . The pairs would be considered as matching if

1. Circles  $C_1, C'_1$  have equal radii within a threshold (maximum difference of 0.1 m);
2. Circles  $C_2, C'_2$  have equal radii within a threshold (maximum difference of 0.1 m);
3. The distance between the centers of  $C_1, C'_1$  equals the distance between the centers of  $C_2, C'_2$  within a threshold (maximum difference of 0.2 m);
4. The angle between the normals of  $C_1, C'_1$  equals the angle between the normals of  $C_2, C'_2$  within a threshold (maximum difference of  $10^\circ$ ).

Furthermore, consider a pair of circles  $(C_1, C_2)$  from scan  $S_1$  and another pair of circles  $(C'_1, C'_2)$  from scan  $S_2$  that could be considered a valid match according to the previous definitions. A transformation (rotation  $R$  followed by a translation  $T$ ) can be computed by converting the correspondence of a pair of oriented circles to a pair of 3D oriented lines. This approach leads to a robust transformation computation, since it is based on relative position and orientation of the features rather than exact position and orientation of each feature. In particular, two cases are considered:

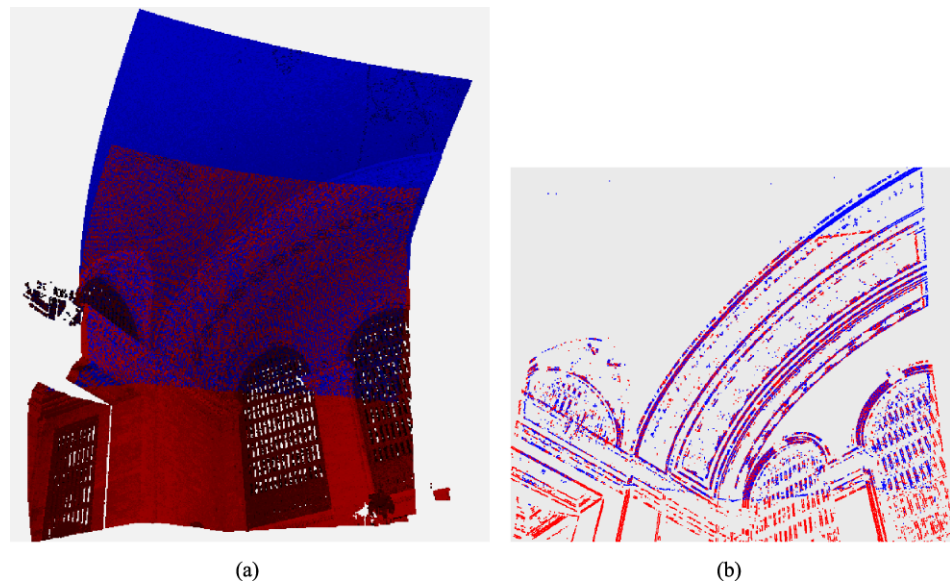
*Case 1:* Circles  $(C_1, C_2)$  have parallel normals  $\mathbf{V}_1$  and  $\mathbf{V}_2$  (the same is true for the normals  $\mathbf{V}'_1$  and  $\mathbf{V}'_2$  of circles  $(C'_1, C'_2)$ ) (Fig. 5(a)). Let us consider the oriented line  $\mathbf{D}$  that connects the centers of  $(C_1, C_2)$  and the oriented line  $\mathbf{D}'$  that connects the centers of  $(C'_1, C'_2)$ . If  $\mathbf{D}$  is not parallel to  $\mathbf{V}_1$  (that means that  $\mathbf{D}'$  is not parallel to  $\mathbf{V}'_1$ ), the match of the oriented line  $\mathbf{D}$  with  $\mathbf{D}'$  and  $\mathbf{V}_1$  with  $\mathbf{V}'_1$  can provide a reliable transformation (closed form formula, Stamos and Leordeanu 2003). Otherwise  $\mathbf{D}$  is parallel to  $\mathbf{V}_1$



**Fig. 5** Two cases of matching circle pairs. The *dash line* separates scan  $S_1$  from  $S_2$ . The radii and relative position of the two circles from  $S_1$  must be similar to those from  $S_2$ . (a) Case 1: two circles have par-

allel normals.  $V_1, D$  and  $V_1', D'$  are used to compute transformation. (b) Case 2: two circle normals are not parallel.  $V_1, V_2$  and  $V_1', V_2'$  are used to compute transformation

**Fig. 6** Registered images of Fig. 2(a) and (b). They are colored to highlight overlapping area. (a) All image points. (b) Edge points at overlapping area



and a reliable transformation cannot be computed. Note that if the length of  $\mathbf{D}$  (or  $\mathbf{D}'$ ) is below an empirically determined threshold (it is 5m for the Grand Central Terminal dataset) the above computation is not performed (that means that the candidate match is discarded). This effectively improves the performance by not considering pairs of spatially close features.<sup>3</sup>

*Case 2:* Circles ( $C_1, C_2$ ) do not have parallel normals  $\mathbf{V}_1$  and  $\mathbf{V}_2$  (the same is true for the normals  $\mathbf{V}'_1$  and  $\mathbf{V}'_2$  of circles ( $C'_1, C'_2$ )) (Fig. 5(b)). Then, the two pairs of oriented lines ( $\mathbf{V}_1, \mathbf{V}_2$ ) and ( $\mathbf{V}'_1, \mathbf{V}'_2$ ) are used for the computation of a reliable transformation (closed form formula (Stamos and Leordeanu 2003)).

From each valid matching circle pairs, a candidate transformation is computed as described above. Each transformation is verified for correctness as follows. Based on the fact that overlapping images are captured from nearby positions, we discard all rotation matrices with diagonal elements smaller than 0.7 (allowing  $45^\circ$  tilting of the range

scanner about each of its  $x/y/z$  axes). Note that this step reduces the number of possible transformations and thus speeds up the algorithm, but is not otherwise necessary. Then we test whether the transformation causes surface inconsistency (Huber and Hebert 2003). Finally, from all verified transformations, the one achieving the smallest average distance between overlapping range points is chosen as the best.<sup>4</sup>

### 3.4 Experimental Results

Our automated method is used for registration of the interior scans of Grand Central Terminal in NYC (a large-scale landmark urban structure). The best transformation of the two corner scans of Fig. 2 provides a registration error (average point distance in the 55.7% overlapping area) of 0.95 cm. Within a few iterations of ICP an optimal transformation with a registration error of 0.90 cm is obtained (Fig. 6).

Also, we registered other scans of this hall with the same technique. The entire hall is roughly rectangular in shape

<sup>3</sup>Our segmentation, as seen in Fig. 4, produces many similar spatially close circles. We decided not to average them in order to not degrade accuracy.

<sup>4</sup>Note that an approach similar to association graphs (Ambler et al. 1973) would generate a very large search space.



**Table 1** Experimental results of circle-based 3D-to-3D registration. Grand Central Terminal dataset. Columns: Number of circles in images ({left, right}); Number of candidate transformations; Average point distance assuming best transformation after circle-matching; Average point distance after ICP optimization; Percentage of overlap; Execution time, including circle extraction and matching (on a Linux-based 2 GHz Xeon-Processor with 2 GB of RAM)

Circles	RT	Dist_match	Dist_icp	Overlap	Time
{26, 24}	544	0.95 cm	0.90 cm	55.7%	6 min
{24, 39}	980	1.11 cm	1.00 cm	29.1%	9 min
{21, 39}	15	3.42 cm	1.28 cm	16.1%	1 min
{24, 20}	748	2.13 cm	0.77 cm	38.4%	7 min
{13, 30}	126	2.01 cm	0.84 cm	28.9%	2 min
{21, 26}	534	1.68 cm	0.90 cm	35.5%	6 min
{23, 11}	29	4.26 cm	0.87 cm	28.7%	1 min
{14, 31}	18	2.65 cm	0.93 cm	27.8%	2 min
{31, 13}	58	2.34 cm	0.98 cm	23.0%	2 min
{37, 26}	67	3.83 cm	0.87 cm	37.2%	2 min
{23, 35}	310	1.20 cm	0.84 cm	26.7%	7 min
{49, 41}	3054	2.81 cm	1.02 cm	38.7%	58 min
{50, 38}	931	1.83 cm	0.92 cm	44.6%	10 min

with an arched ceiling. Figure 7 shows a few typical scans on the front wall ((a), (c)) and the side wall ((e), (g)), together with circles extracted from them. Note that the lower parts of the walls (e.g. (c), (g)) contain lines and planes, and are therefore registered with our linear-feature based technique (Chen and Stamos 2005; Stamos and Leordeanu 2003). The upper regions with very few linear features, e.g. (a), (e), are registered with their lower neighboring scans (c), (g) respectively, by matching overlapping circular windows.

In Fig. 9, registered edge points from 23 scans are visualized. There are another 14 scans not shown for clarity of presentation. Among all 37 scans, 20 of them are lower parts registered with lines, 13 of them are the upper parts registered with their lower neighbor scans based on overlapping circular windows. Three scans are manually registered, because they are cylindrical ceiling patches without any distinguishing geometric shape information. In Table 1 we report the performance of 13 registrations based on circles. When registering the last two pairs, a long execution time is experienced due to a large number of valid transforms from the precisely extracted circles around the window frame (as in Fig. 7(b), (d)). To avoid unnecessary computations, we set the program to terminate when the average distance falls below 0.03 cm (approximate distance between two neighboring points). The values in columns *RT*, *Dist\_match*, *Time* are therefore recorded up to the point when an accurate enough result is reached. In Table 2 we report the performance of line-based registration (Chen and Stamos 2005; Stamos and Leordeanu 2003) in the lower part of the building. In Fig. 8, more results are shown, including registered

**Table 2** Experimental results of line-based 3D-to-3D registration (Chen and Stamos 2005; Stamos and Leordeanu 2003). Grand Central Terminal dataset. Columns: Number of lines in images ({left, right}); Number of matching line-pairs before ICP; Average point distance assuming best transformation from line-based matches; Number of matching line-pairs after ICP; Average point distance after ICP optimization; Execution time, for line matching (on a Linux-based 2 GHz Xeon-Processor with 2 GB of RAM)

Line Pairs	<i>N</i>	Dist_match	<i>N</i>	Dist_icp	Time
{318, 214}	5	2.31 cm	8	0.934 cm	1 sec
{318, 261}	19	2.16 cm	21	1.02 cm	10 sec
{261, 239}	13	1.51 cm	13	1.14 cm	10 sec
{239, 236}	22	1.05 cm	22	0.988 cm	10 sec
{241, 230}	3	2.26 cm	2	1.30 cm	9 sec
{241, 230}	3	1.53 cm	3	1.33 cm	10 sec
{230, 175}	1	2.66 cm	2	1.49 cm	11 sec
{175, 175}	16	1.05 cm	17	0.932 cm	10 sec
{175, 247}	2	3.32 cm	4	0.981 cm	9 sec
{262, 175}	2	3.37 cm	4	1.39 cm	10 sec
{247, 262}	5	2.97 cm	10	1.04 cm	9 sec
{262, 275}	16	2.00 cm	18	1.04 cm	10 sec
{275, 254}	11	1.15 cm	11	1.08 cm	10 sec
{194, 262}	2	3.79 cm	8	1.10 cm	11 sec
{198, 239}	10	3.13 cm	19	1.39 cm	12 sec
{253, 230}	2	2.26 cm	2	1.19 cm	12 sec

scenes, edge points, and part of a 3D model constructed using the registered range points.

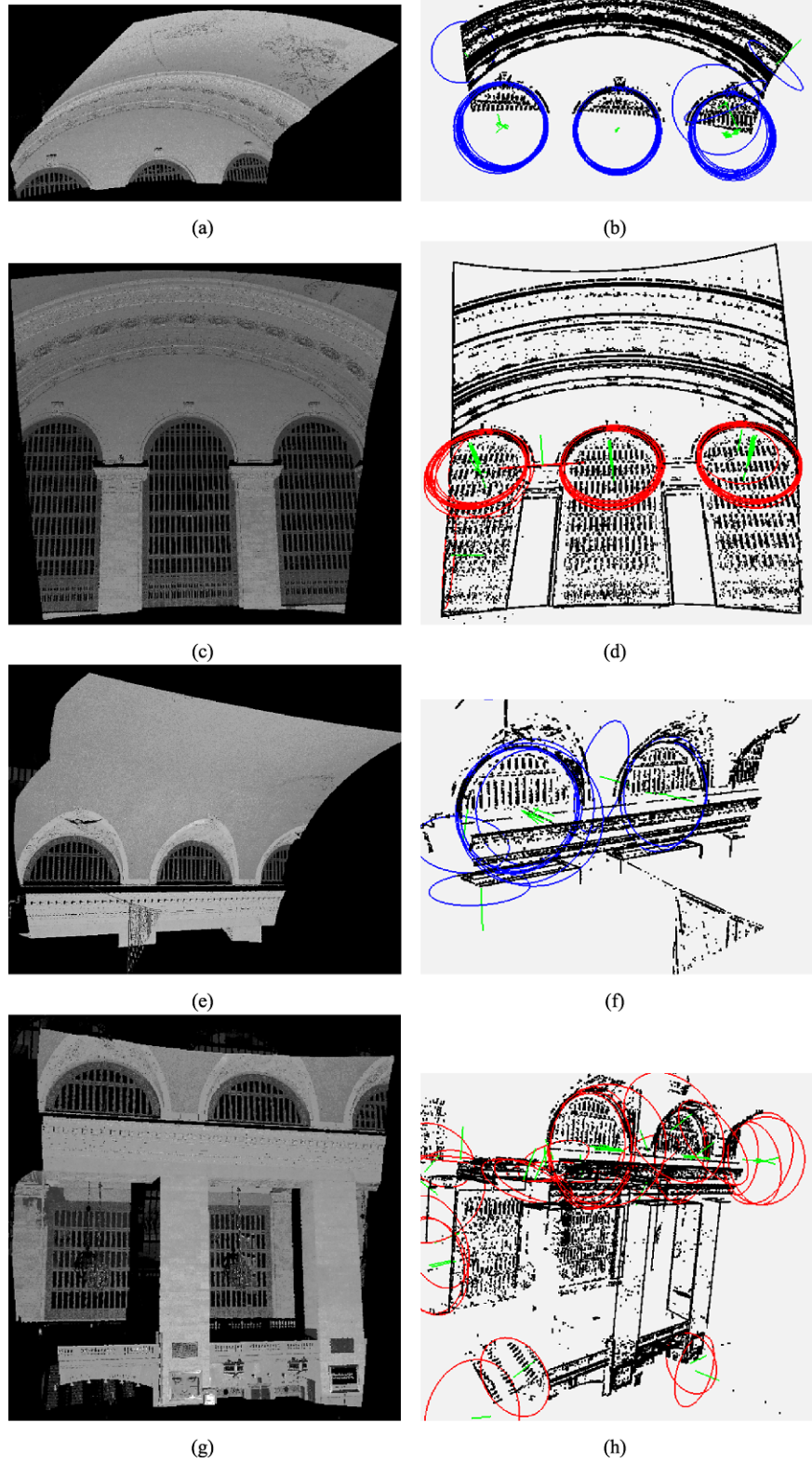
#### 4 2D-to-3D image-to-range registration

We present our automated 2D-to-3D image-to-range registration method used for the automated calibration and registration of a single 2D image  $I_n$  with the 3D range model  $M_{\text{range}}$ . The computation of the rotational transformation between  $I_n$  and  $M_{\text{range}}$  is achieved by matching at least two vanishing points computed from  $I_n$  with major scene directions computed from clustering the linear features extracted from  $M_{\text{range}}$ .

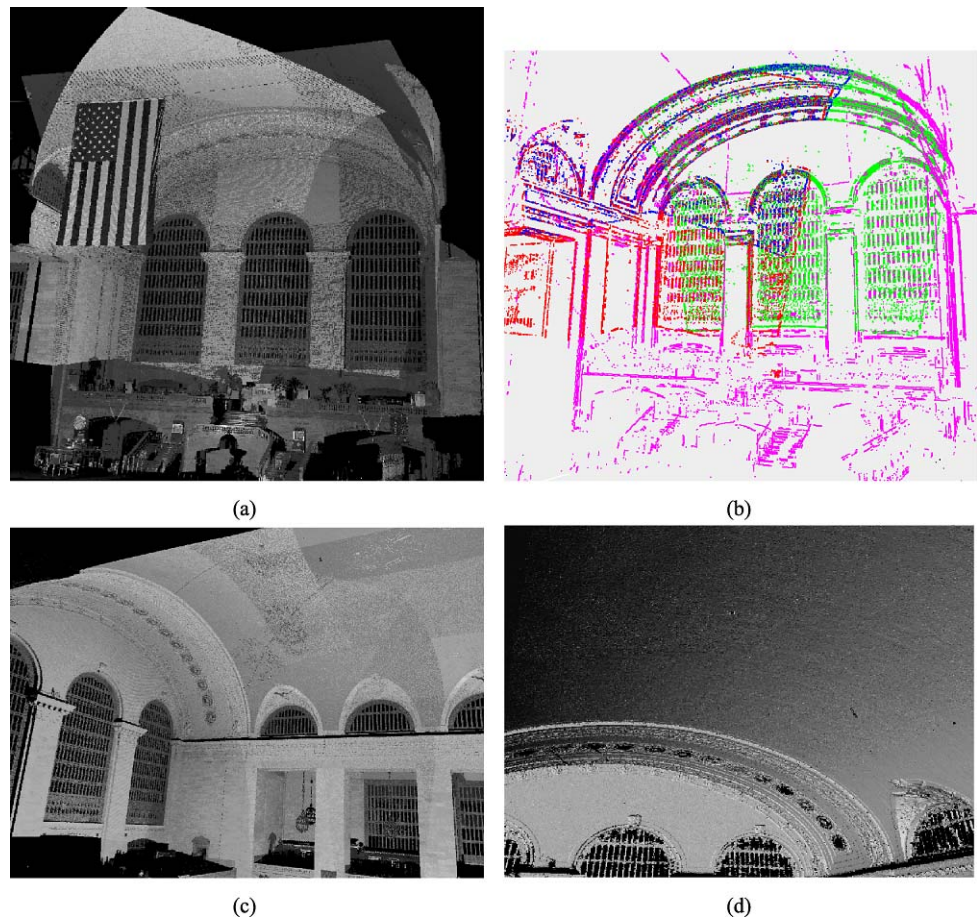
The internal camera parameters consist of focal length, principal point, and other parameters in the camera calibration matrix  $\mathcal{K}$  (Hartley and Zisserman 2003). They are derived from the scene's vanishing points, whereby the 2D images are assumed to be free of distortion after we remove lens deformations, as described in (Liu et al. 2006). Finally, the translation between  $I_n$  and  $M_{\text{range}}$  is computed after 2D and 3D lines from the 2D image and 3D model and are extracted and automatically matched.

With this method, a few 2D images can be independently registered with the model  $M_{\text{range}}$ . The algorithm will fail to produce satisfactory results in parts of the scene where

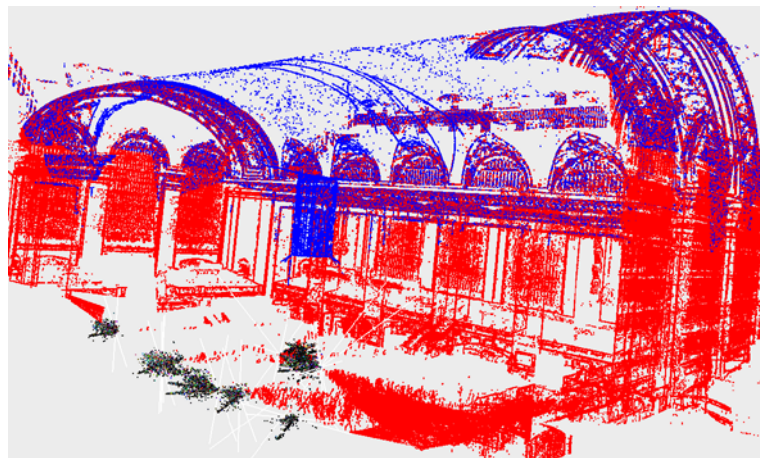
**Fig. 7** Four side wall scans (*left column*), extracted edges and fitted circles (*right column*)



**Fig. 8** Registration results. (a) Four out of the 37 automatically registered scans shown for clarity. (b) Edge points of (a). Four colors represent edge points from four scans. (c) Eight out of the 37 automatically registered scans shown for clarity. (d) 3D mesh model generated by the Ball Pivoting Algorithm (Bernardini and Rushmeier 2002). The smooth ceiling implies the registration is tight and seamless



**Fig. 9** Registered edges of 23 scans, with color indicating scans from upper parts (*blue*) and lower parts (*red*). The other 14 scans are not displayed for clarity of presentation; these scans compose a side wall closer to our point of view and symmetric to the wall being displayed. Tables 1 & 2 report the registration performance for the upper (circle-based registration) and lower (line-based registration) parts



there is a lack of 2D and 3D features for matching. Also, since each 2D image is independently registered with the 3D model, valuable information that can be extracted from relationships between the 2D images (SFM) is not utilized. In order to solve the aforementioned problems, an SFM module (Sect. 5) and final alignment module (Sect. 6) has been added into the system. These two modules increase the robustness of the reconstructed model, and improve the accuracy of the final texture mapping results. Therefore, the 2D-

to-3D image-to-range registration algorithm is used in order to register a few 2D images (five shown in Fig. 1) that produce results of high quality. The final registration of the 2D image sequence with the range model  $M_{\text{range}}$  is performed after SFM is utilized (Sect. 5).

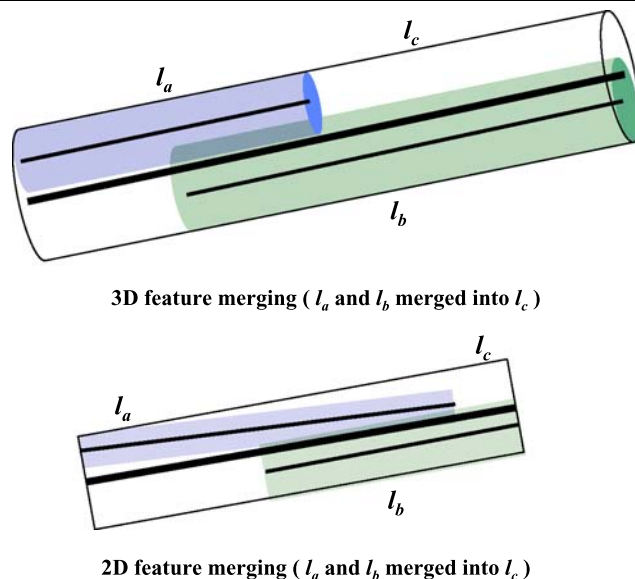
In this section, we present a system that can automatically register 2D images with 3D range data at interactive rates. Our contributions with respect to 2D-to-3D registration can be summarized as follows:

- We have developed a working system that is able to independently register 2D images to 3D models at interactive rates. This system requires minimal user interaction. Note that after a few 2D images are registered to the 3D model the multiview geometry approach (Sect. 5) is utilized for registering all images with the 3D range model.
- The whole space of possible matches between 3D and 2D linear features is explored efficiently (unlike probabilistic methods like (Stamos and Allen 2001)). That improves the possibility of convergence of our algorithm.
- Earlier systems (Liu and Stamos 2005; Stamos and Allen 2001) require the extraction of major facades, rectangles, or other higher-order structures from the 2D and 3D datasets. Our new method, on the other hand, utilizes 3D and 2D linear features for matching without significant grouping. This increases the generality of our algorithm since we make fewer assumptions about the 3D scene. Scenes with various layers of planar facades, or without clear major facades can thus be handled.
- This paper's method utilizes vanishing points and major 3D directions, but it does not require them to be orthogonal as most earlier methods assume.

#### 4.1 3D Feature Extraction for 2D-to-3D Registration

The range-image acquisition process and segmentation is described in Sect. 3. As we described in that section, all range scans are registered with respect to one selected pivot scan, in the scene's coordinate system.

The 3D line extraction step is based on the segmentation method of Stamos and Allen (2002), whereas the major directions clustering is based on (Liu and Stamos 2005). (Note that if 3D information is provided in terms of a CAD model, then the 3D line extraction step is trivial.) The result of this process is a set of line clusters  $\mathcal{L}^{3D}$ . Each line in a cluster has similar orientation as every other line in the same cluster. The set of line clusters are then sorted based on the number of lines in each cluster. We do not assume knowledge of vertical or horizontal directions for the line clusters as in our previous method (Liu and Stamos 2005). Each 3D line is thus associated with a cluster id, e.g. for the 3D lines in cluster  $\mathcal{L}_i^{3D}$ , their cluster id is  $i$ . In the next step, 3D features are extracted. First, an initial user-defined radius (e.g. 0.1 m) is assigned to each 3D line. Then, a line merging step generates the final 3D features. This reduces the number of features, and thus increases the efficiency of the matching stage (Sect. 4.3). In this step, each pair of 3D lines ( $l_a, l_b$ ) with the same cluster id is merged into a new line  $l_c$  (Fig. 10) if (a) the distance between them is smaller than the sum of their radii, and (b) their projections on  $l_c$  overlap. The merging procedure is continued until there are no two remaining 3D lines that can be merged. The final result is a set of 3D lines, each of which is associated with a cluster id and radius.



**Fig. 10** Example of new type of 3D and 2D features and their merging steps

#### 4.2 2D Feature Extraction, Internal Camera Calibration, & Rotation Computation

The internal camera parameters (focal length and principal point) of the camera sensor can be calculated from one 2D image if the image contains at least two vanishing points (i.e. the 3D scene that the image is viewing has at least two major scene directions). We use our previously developed robust methods to generate and cluster 2D lines from a 2D image (Stamos and Allen 2001). The result is a set of major vanishing points  $\mathbf{VP} = \{\mathbf{VP}_1, \mathbf{VP}_2, \dots, \mathbf{VP}_n\}$ . Using the methods described in (Stamos and Allen 2001) we can compute the center of projection  $\mathbf{COP} = [C_x, C_y, C_z]^T$  (effective focal length and principal point expressed in pixels) by utilizing **three** orthogonal vanishing points. In (Liu and Stamos 2005) we described an iterative method for estimating the internal calibration parameters from two orthogonal vanishing points.

In this section we present an additional method (Liu 2007; Liu and Stamos 2007) for the calculation of focal length  $f$  and rotation  $\mathbf{R}$ . We use two vanishing points and two major 3D directions, but we do not assume that these directions are orthogonal to each other. Orthogonality is prominent in urban scenes, but is not always present. Our method starts with an initial estimate  $f_{\text{init}}$  of the effective focal length, and of the principal point  $\mathbf{P}_{\text{init}}$ .  $f_{\text{init}}$  is included in the Exif meta-data, information that is now provided by most digital cameras.  $\mathbf{P}_{\text{init}}$  is estimated by the center of the image. Based on these estimates an initial center of projection  $\mathbf{C}_{\text{init}}$  is determined. This is the origin of the camera coordinate system (Fig. 11).

Let us consider a vanishing point  $\mathbf{V}_i$  extracted from a 2D image. The 3D coordinates of  $\mathbf{V}_i$  in the camera coordinate

system are  $[(V_i)_x - (P_{init})_x, (V_i)_y - (P_{init})_y, f_{init}]^T$ . Thus, the normalized vector  $\mathbf{D}_i^{2D} = u(\mathbf{C}_{init}\mathbf{V}_i)^5$  represents the 3D direction that generates the vanishing point  $\mathbf{V}_i$ . This direction is expressed in the camera coordinate system. Our goal is to match each vanishing point with its corresponding 3D direction extracted by the 3D range model (see Sect. 4.1). This correspondence leads to the calculation of the focal length and the rotation  $\mathbf{R}$ . Let us represent each 3D line cluster in  $\mathcal{L}^{3D}$  (Sect. 4.1) by its 3D direction  $\mathbf{D}_j^{3D}$ ,  $j = 1, \dots, n$  (where  $n$  is the number of extracted 3D clusters).

The next step is to find the matching pairs of directions  $\langle \mathbf{D}_i^{2D}, \mathbf{D}_j^{3D} \rangle$ . Consider for the moment that we know the correspondence between vanishing points (expressed in the camera coordinate system) and 3D directions (expressed in the world coordinate system). It is known that with the principal point fixed at the center of the image, two pairs  $(\langle \mathbf{D}_a^{2D}, \mathbf{D}_a^{3D} \rangle, \langle \mathbf{D}_b^{2D}, \mathbf{D}_b^{3D} \rangle)$  of matching vanishing point/3D directions are enough for the computation of the focal length  $f$ . The focal length  $f$  (which is  $|\mathbf{CP}|$  in Fig. 11) can be computed via the following equations (triangles  $\mathbf{CV}_a\mathbf{P}$ ,  $\mathbf{CV}_b\mathbf{P}$  and  $\mathbf{CV}_a\mathbf{V}_b$ ):

$$\begin{aligned} |\mathbf{CV}_a|^2 &= |\mathbf{PV}_a|^2 + f^2, \\ |\mathbf{CV}_b|^2 &= |\mathbf{PV}_b|^2 + f^2, \\ |\mathbf{V}_a\mathbf{V}_b|^2 &= |\mathbf{CV}_a|^2 + |\mathbf{CV}_b|^2 - 2 \cdot |\mathbf{CV}_a| \cdot |\mathbf{CV}_b| \cdot \cos \alpha \end{aligned}$$

where  $\alpha$  is the angle between  $\mathbf{D}_a^{3D}$  and  $\mathbf{D}_b^{3D}$ . (Note that the vanishing points  $\mathbf{V}_a$  and  $\mathbf{V}_b$  have been computed by using the initial estimates  $f_{init}$  and  $\mathbf{P}_{init}$ . The above computation leads to the calculation of a focal length that conforms to the 3D directions  $\mathbf{D}_a^{3D}$  and  $\mathbf{D}_b^{3D}$ .) From the above equations, we can get a quartic equation:

$$a \cdot f^4 + b \cdot f^2 + c = 0$$

where  $a = \sin^2 \alpha$ ,  $b = \sin^2 \alpha (|\mathbf{PV}_a|^2 + |\mathbf{PV}_b|^2) - |\mathbf{V}_a\mathbf{V}_b|^2$ ,  $c = (\frac{|\mathbf{V}_a\mathbf{V}_b|^2 - |\mathbf{PV}_a|^2 - |\mathbf{PV}_b|^2}{2})^2 - \cos^2 \alpha |\mathbf{PV}_a|^2 |\mathbf{PV}_b|^2$ . Solving this equation, one obtains the refined focal length:  $f = \sqrt{\frac{\sqrt{b^2 - 4ac} - b}{2a}}$ . Since  $\mathbf{D}_a^{3D} \neq \mathbf{D}_b^{3D}$ ,  $\sin \alpha$  will never be equal to 0. Finally, the rotation  $\mathbf{R}$  is computed based on these two pairs of matching directions (Faugeras 1996).

Based on the above analysis, the task of our system is to find two matching pairs of vanishing point/3D directions. Intuitively, pairs  $(\langle \mathbf{D}_a^{2D}, \mathbf{D}_a^{3D} \rangle, \langle \mathbf{D}_b^{2D}, \mathbf{D}_b^{3D} \rangle)$  for which the angle between  $\mathbf{D}_a^{2D}$  and  $\mathbf{D}_b^{2D}$  is not similar to the angle between  $\mathbf{D}_a^{3D}$  and  $\mathbf{D}_b^{3D}$  can be rejected. As a result, we have a list of matching candidates, each of which contains two pairs of matching vanishing points and 3D directions, a refined focal length and a rotation. For each one of these candidates we can apply the algorithm described in the next section for

calculating the camera position, and finally keep the result that provides the maximal alignment between the 2D image and 3D model.

In the worst case scenario though all pairs of directions have similar angles (this scenario is easily realizable in urban scenes where most angles between major directions is 90 degrees). In this case there are  $\binom{n}{2} \binom{m}{2}$  candidate matching pairs of directions (where  $n$  is the number of 3D and  $m$  the number of vanishing points). Even though this is not a large search space ( $n$  and  $m$  are small in most urban scenes), testing all hypotheses involves the computation of the translation (see Sect. 4.3). This is computationally inefficient for the purposes of an interactive system, where a response time of up to 10 seconds per image is appropriate. For these reasons we let the user provide the correct pair of matching directions, by rotating the 3D model to an orientation that produces a rendering that is similar (but not exactly the same) to the real 2D image. As shown in Figs. 15(b) and 16(c), the rotated 3D view (left) is similar (but not exactly the same) to the 2D image (right). This user-assisted rotation can approximately align the corresponding 2D and 3D directions.

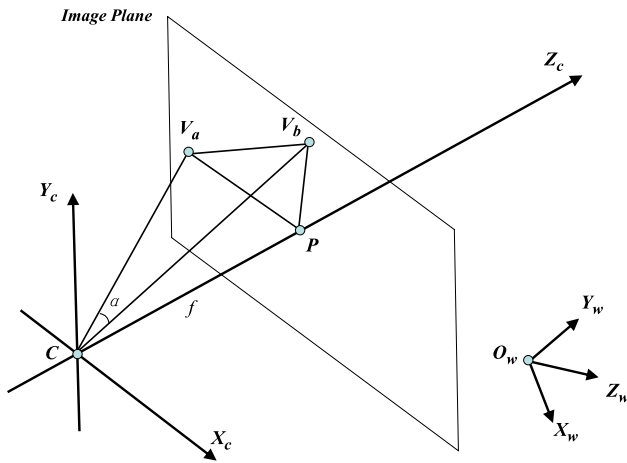
The aforementioned user interaction not only increases the computational efficiency of the whole system, but also makes the registration problem tractable. In general, without constraining the possible locations of 2D cameras wrt the 3D model, the 2D-to-3D registration problem becomes intractable. This is due to the existence of a possible large set of solutions. For example, a photograph of one of the columns of the 3D structure of Fig. 16 can be matched with any of the symmetric 3D columns of the real scene. By selecting a synthetic view that is similar, but not exactly the same as the 2D image, the user can provide an approximate field of view to help the matching algorithm. In particular, only 3D features that are viewable in the synthetic 3D view are used for matching 2D image features. Note here that all earlier approaches still require implicit user interaction in order to assist in that direction. For example in our earlier work (Liu and Stamos 2005) the user needs to explicitly provide the match between vanishing points and 3D directions. In that earlier system, the user also needs to match facades between the 2D image and 3D model. The approach presented in this section is more natural and leads to faster interaction time.

### 4.3 Translation Computation

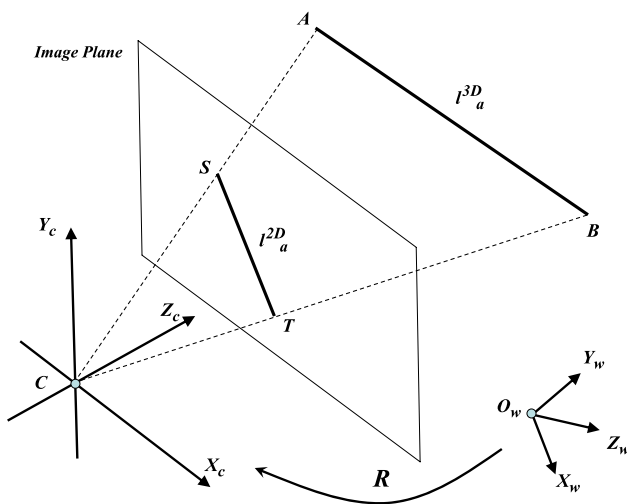
In this section we present the algorithm that automatically computes the translation between one 3D range scan and one 2D image (Fig. 13).

As described in Sect. 4.2, a list of matching candidates, named  $\mathcal{M}$ , is obtained. Each element in  $\mathcal{M}$  contains a

<sup>5</sup>We use the notation  $u(\mathbf{v})$  for describing the unit vector derived from  $\mathbf{v}$ .



**Fig. 11** Rotation and focal length computation based on two vanishing points and their corresponding 3D directions (not shown in this image)



**Fig. 12** Camera position computation based on a match between 3D feature  $AB$  with image feature  $ST$

matching pair of two vanishing points and two 3D directions, a refined focal length and a rotation. In this section, a 2D camera position will be computed for each candidate in  $\mathcal{M}$ . Our method of finding the camera position follows a hypothesis-and-test scheme by matching the extracted 3D and 2D features based on our original framework (Liu and Stamos 2005). A number of major differences with the aforementioned method make our algorithm more general and more robust. In particular, our current algorithm does not require the extraction of planar facades, and does not require the grouping of low-level features in higher-order structures. Scenes that do not contain clear major facades (such as the example of Figs. 16, where various layers of planar facades exist) can now be successfully handled. Also since all low-level features are used without significant grouping, more robust results are achieved.

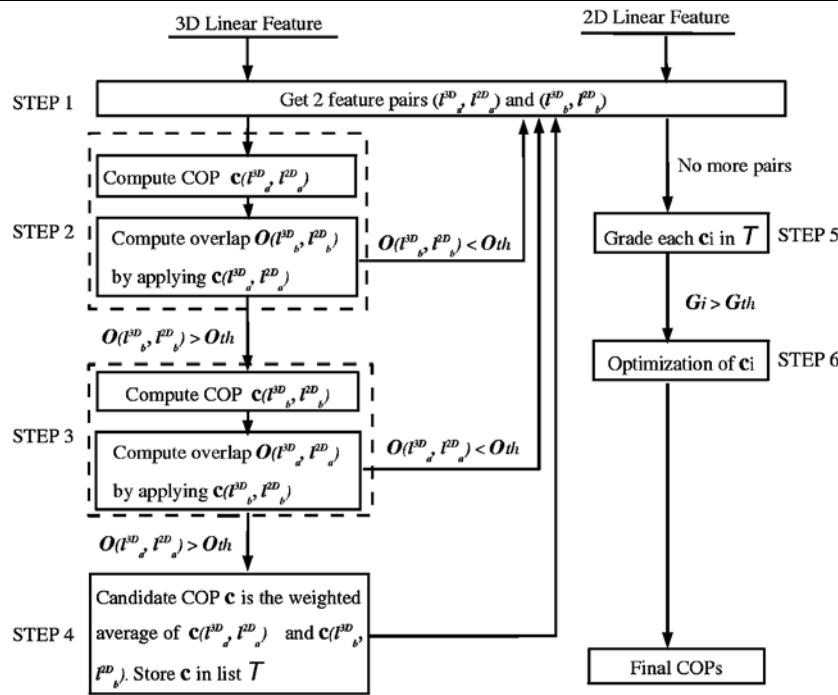
We now present a detailed description of our algorithm. First, a candidate from  $\mathcal{M}_i$  is selected, i.e. the matching pair of vanishing points and 3D directions are  $\langle \mathbf{V}_a, \mathbf{V}_b \rangle$  and  $\langle \mathbf{D}_a^{3D}, \mathbf{D}_b^{3D} \rangle$ ; the refined focal length is  $f_i$  and the rotation is  $\mathbf{R}_i$ . The camera position (translation) is then computed in the following six steps (Fig. 13):

**Step 1** A hypothetical match between two pairs of 3D and 2D lines is selected (the algorithm will go over all possible such selections). Let us call these pairs  $(l_a^{3D}, l_a^{2D})$  and  $(l_b^{3D}, l_b^{2D})$  ( $l_a^{3D}$  and  $l_b^{3D}$  are 3D lines extracted from the 3D model, and  $l_a^{2D}$  and  $l_b^{2D}$  2D lines extracted from the 2D image).

**Step 2** (Computation of camera position in world coordinate system (translation) based on the match of  $l_a^{3D}$  with  $l_a^{2D}$ ) As shown in Fig. 12,  $\mathbf{A}$  and  $\mathbf{B}$  are the endpoints of  $l_a^{3D}$  and  $\mathbf{S}$  and  $\mathbf{T}$  are the endpoints of  $l_a^{2D}$ .  $\mathbf{C}$  is the center of projection. If  $l_a^{3D}$  matches exactly with  $l_a^{2D}$ , then in the camera coordinate system,  $\mathbf{C}$ ,  $\mathbf{S}$  and  $\mathbf{A}$  should be collinear. The same applies for  $\mathbf{C}$ ,  $\mathbf{T}$  and  $\mathbf{B}$ . We thus consider  $\mathbf{C}$  as the intersection point of the following two lines: (a) one that passes through  $\mathbf{A}$  having the orientation of line  $\mathbf{CS}$  and (b) one that passes through  $\mathbf{B}$  having the orientation of line  $\mathbf{CT}$ . To compute the world coordinates of  $\mathbf{C}$ , we need to know the orientations of  $\mathbf{CS}$  and  $\mathbf{CT}$  in the world coordinate system. We know, however, the orientations of  $\mathbf{CS}$  and  $\mathbf{CT}$  in the camera coordinate system, say  $\mathbf{n}_a$  and  $\mathbf{n}_b$ . Furthermore, we have also computed the rotation  $R$  that brings the camera and world coordinate systems into alignment (see previous section). We can thus compute the orientations of  $\mathbf{CS}$  and  $\mathbf{CT}$  in the world coordinate system as:  $R \cdot \mathbf{n}_a$  and  $R \cdot \mathbf{n}_b$ . Then, the camera position is obtained by finding the intersection of two 3D lines: (a) one of which passes through  $\mathbf{A}$  with the orientation of  $R \cdot \mathbf{n}_a$  and (b) one which passes through  $\mathbf{B}$  with the orientation of  $R \cdot \mathbf{n}_b$ .<sup>6</sup> Finally, this computed center of projection is used to project  $l_b^{3D}$  onto the image plane. If the projection of  $l_b^{3D}$  overlaps with  $l_b^{2D}$  (within a threshold of 80%), then the camera position computed using  $(l_a^{3D}, l_a^{2D})$  is verified by the pair  $(l_b^{3D}, l_b^{2D})$ . We therefore move to the next step. Otherwise, we return to step 1 (i.e. the match is discarded) to pick another set of hypothetical matching lines.

**Step 3** Step 2 is repeated assuming as hypothesis the match between  $l_b^{3D}$  and  $l_b^{2D}$ . The newly computed center of projection is used to compute the overlap between  $l_a^{2D}$  and the projection of  $l_a^{3D}$ . If this overlap is less than a threshold (i.e. the computed  $\mathbf{C}$  is not verified by  $(l_a^{3D}, l_a^{2D})$ ), we return to step 1 (i.e. the match is discarded). Otherwise, we proceed to the next step.

<sup>6</sup> $\mathbf{A}$  and  $\mathbf{B}$  are both expressed in the world coordinate system.



**Fig. 13** Camera position (translation) computation flowchart [COP stands for camera position (or center of projection) in this flowchart]. Through step 1 all possible pairs of matched 3D and 2D lines  $(l_a^{3D}, l_a^{2D})$  and  $(l_b^{3D}, l_b^{2D})$  are selected ( $l_a^{3D}$  and  $l_b^{3D}$  are 3D lines extracted from the 3D model, and  $l_a^{2D}$  and  $l_b^{2D}$  2D lines extracted from the 2D image). Step 2 computes a camera position based on  $(l_a^{3D}, l_a^{2D})$ . The pair  $(l_b^{3D}, l_b^{2D})$  is used for the verification of this position. If the overlap between  $l_b^{2D}$  and the projection of  $l_b^{3D}$  on the image is smaller than  $O_{th}$

(20%) (i.e. the position is not verified) a new pair is selected (step 1). Otherwise a similar computation is carried out for the pair  $(l_b^{3D}, l_b^{2D})$  (step 3). If steps 2 and 3 produce two verifiable camera positions, a weighted average is computed (step 4). This average represents the position that is generated by the hypothesis  $(l_a^{3D}, l_a^{2D})$  and  $(l_b^{3D}, l_b^{2D})$ . All verified camera positions are stored in a list  $T$ . After all pairs have been explored, each position in  $T$  is graded by projecting all 3D lines on the 2D image space (step 5). Positions with high grade (greater than  $G_{th}$  number of matches) survive to the final optimization step 6

**Step 4** Step 2 has thus computed a camera position  $C_1$  by the hypothesis  $(l_a^{3D}, l_a^{2D})$  [verified by  $(l_b^{3D}, l_b^{2D})$ ], while step 3 has computed a camera position  $C_2$  by the hypothesis  $(l_b^{3D}, l_b^{2D})$  [verified by  $(l_a^{3D}, l_a^{2D})$ ]. In this step, the weighted average (based on the amount of overlap) of these two camera positions is computed and saved in a list  $T$ .

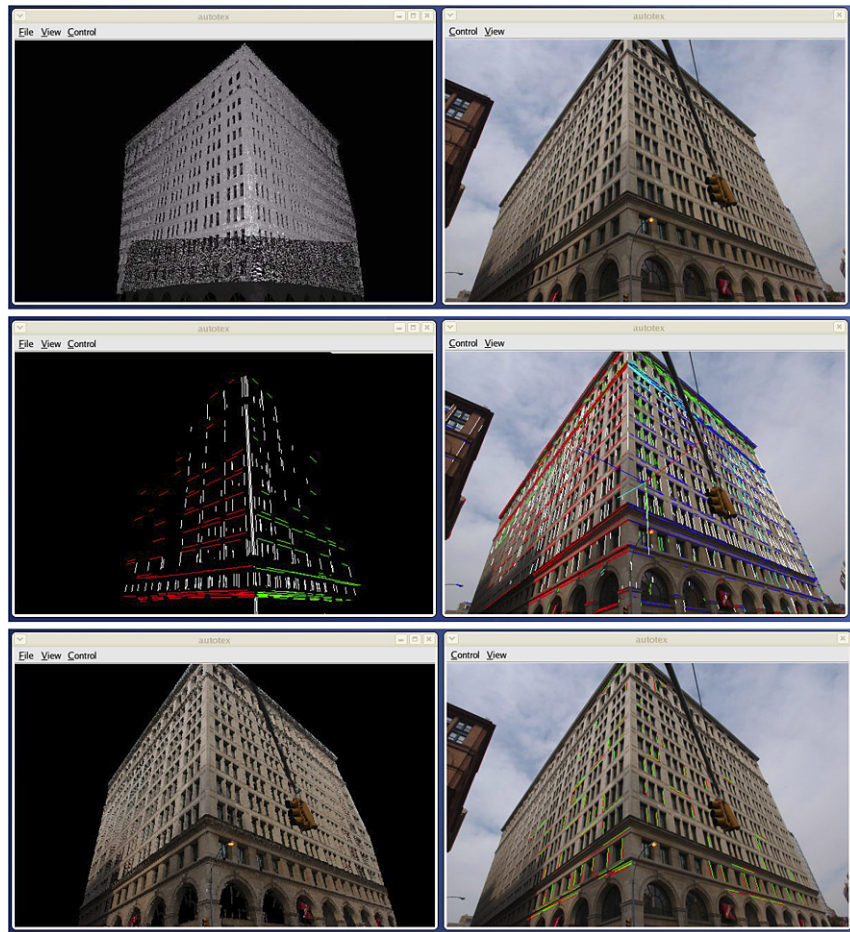
**Step 5** Steps 1 to 4 are repeated for all possible pairs of pairs of 3D and 2D lines  $(l_a^{3D}, l_a^{2D}), (l_b^{3D}, l_b^{2D})$ . All verified camera positions (see Step 4) are stored in a list  $T$ . Then, for each position in  $T$ , all 3D lines are projected onto the image plane. For each of the projected 3D lines, a possible matching 2D line is found by searching around its projection. This region is bounded by the radius of the 3D and 2D lines. The number of found matches grades this camera position. If the grade of a camera position is less than a threshold, it is removed from the list  $T$ .

**Step 6** The remaining camera positions in  $T$  are optimized by two steps. First, for each camera position  $C_i$  a refined position is found. This is achieved by searching around a small neighborhood of  $C_i$  in order to maximize the overlap

between the matching 3D and 2D lines. Then this refined position is further optimized by an iterative algorithm. In each iteration, the current camera position is used to generate a list of matching 2D and 3D lines from the whole 2D and 3D feature space. A new camera position is found by minimizing the error distance between the 2D lines and the projections of their matching 3D lines. The algorithm converges at the point when the error distance remains constant. The camera position computed after the two optimization steps is the final result.

The camera position in  $T$  with the maximum grade is picked as the best one for the matching candidate  $\mathcal{M}_i$ . This is normally correct, but the list is still kept as well in case the one with the maximum grade is not the best. Then, the user can select other positions in the list. This maximum grade is also used as the grade for  $\mathcal{M}_i$ . For each matching candidate in  $\mathcal{M}$ , a list of camera positions is computed by these 6 steps and a grade is assigned. Then, the list  $\mathcal{M}$  is sorted based on the grade and the one with the maximum grade is selected as the best one, but the user also can select other results in  $\mathcal{M}$ .

**Fig. 14** Registration result of Building 2. *Top row:* Initial state (before registration). The 3D range model (*left column*) and 2D image (*right column*) are loaded and displayed in the interface. Note that two overlapping range scans are shown. *Middle row:* The state of the system after the feature extraction. The 3D viewer (*left column*) shows the clustered 3D lines while the 2D viewer (*right column*) shows the clustered 2D lines that are drawn on the original 2D image. Different clusters are represented by different colors for clarity. *Bottom row:* The final registration. The 2D image is automatically registered with the 3D range data. The 3D viewer (*left*) shows the texture mapped 3D range data. The 2D viewer (*right*) shows the matching 2D and 3D line features (2D lines are displayed as red, while projected 3D lines are highlighted in green). Note that objects that are not part of the 3D model cannot be texture-mapped (corner of other building shown in the 2D image)



#### 4.4 Results & Conclusions

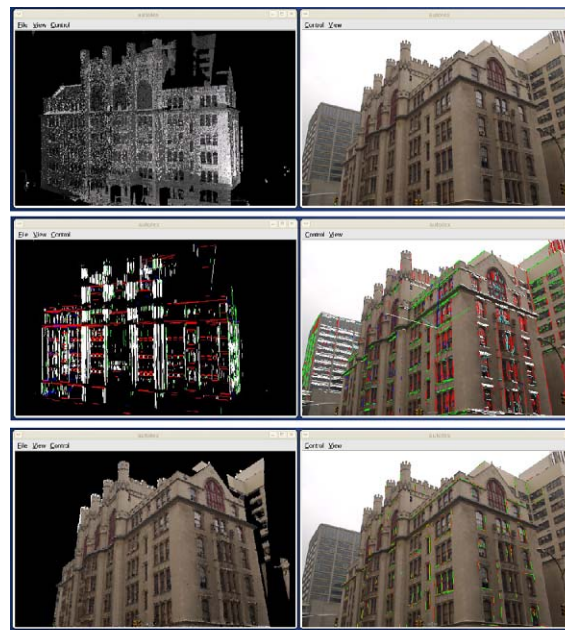
We are presenting results from real experiments in three urban settings that we name 1 (Fig. 15), 2 (Fig. 14), and 3 (Fig. 16). Buildings 1 and 2 are the exteriors of regular urban structures. Building 3 is the interior of Grand Central Station, a scene of architectural complexity and beauty. The 3D range model for all buildings was generated by our group. First a number of 3D range scans of each structure was acquired using a Leica HDS 2500 time-of-flight laser range scanner (Leica Geosystems 2007). This scanner provides absolute 3D range measurements up to a distance of 100 meters, and at an accuracy of 6 mm. Each 3D point is associated with reflectance information, that corresponds to the amount of laser-intensity getting back to the range sensor.<sup>7</sup> We then segment each range scan, extract linear 3D features, and register the scans in a common coordinate system.

<sup>7</sup>Note that the reflectance depends on various parameters (distance, orientation and surface material) and is not the actual color of the object as captured by a 2D digital camera.

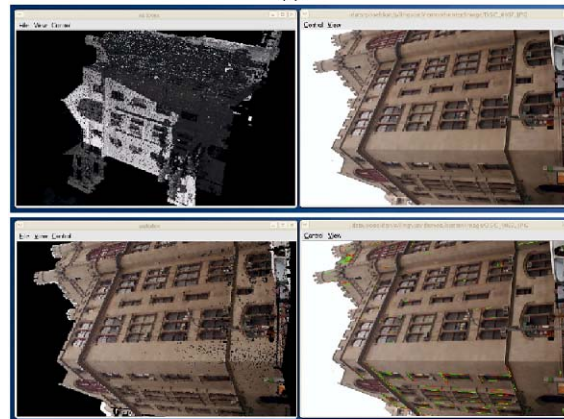
Figures 14, 15, and 16 provide individual registration results, as described in our technical sections. Note that in the case of Figs. 15(b) and 16(c) the user needs to orient the 3D range model in a position that simulates the 2D color image. As you can see from these figures this simulation need not be exact. It is necessary for assistance in matching vanishing points with 3D directions (Sect. 4.2) in order for our system to perform at near-interactive rates (5–10 seconds for matching per image). Table 3 presents quantitative results for successful automated registrations (see caption for more details). A 3–5 pixel distance between the matched 2D and projected 3D lines is due to noise in the line extraction process. Nevertheless, our texture mapping results are of extremely high quality though. Out of 18 total images tried for building 1, 13 were registered successfully, whereas 5 have slight errors. Out of 8 total images tried for building 2, 7 were registered successfully, whereas 1 has slight errors. Finally, out of 10 total images tried for building 3, 6 were registered successfully, whereas 4 have slight errors. In all cases the first step (Sect. 4.2) never fails since these scenes contain at least two vanishing points. The second step (Sect. 4.3), however, depends on the quality of



**Fig. 15** Registration results from building 1. **(a)** For description see caption of Fig. 14. **(b) (Top row)**: The user rotates the 3D model so that it is orientated similarly (note that it does not have to be exactly matched) to the 2D image. **(Bottom row)**: The *right image* shows the 2D image along with the matched 2D and projected 3D features (see caption of Fig. 14). The *left image* shows the texture-mapped 3D range model after successful registration



(a)



(b)

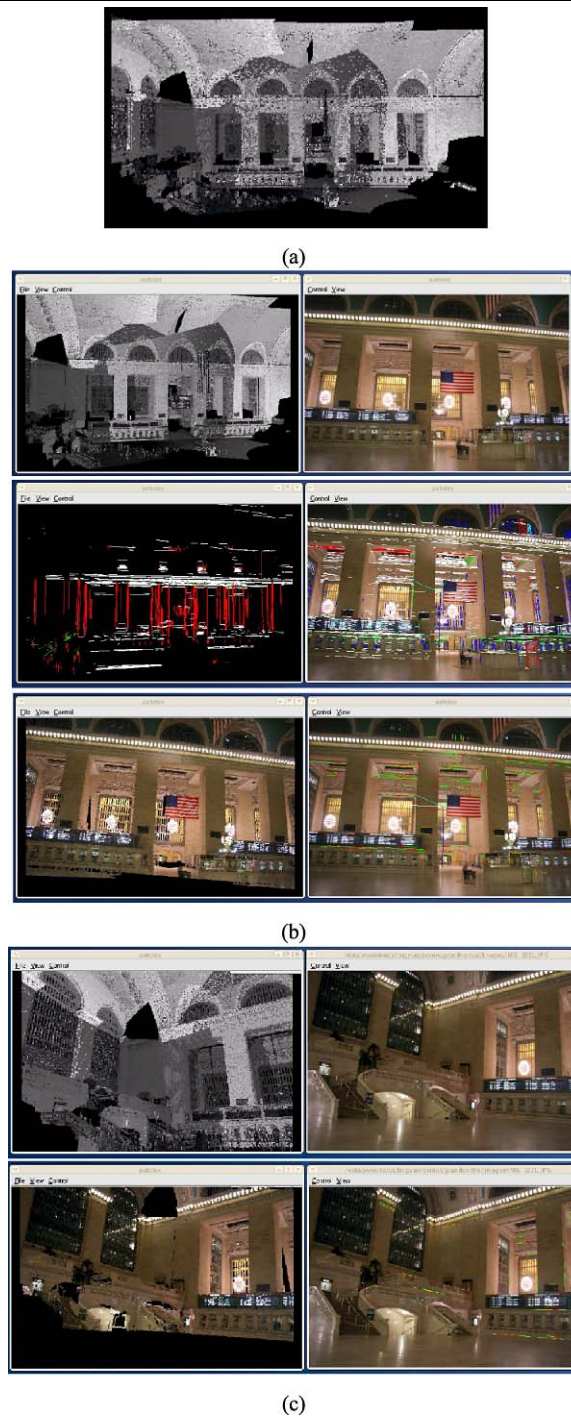
the extracted low-level 2D and 3D linear features. In cases for which we cannot extract features of high quality (due to low contrast in 2D images), this method will not be able to perform correctly. On the other hand few correct 2D-to-3D registrations can be enhanced with multiview-geometry solutions to bring sequences in alignment with a model (see Sect. 5).

We have presented a systematic way for registering individual 2D images with a 3D range model. Our methods assume the existence of at least two vanishing points in the scene (not necessarily orthogonal). No higher-order grouping of features is necessary. Our system allow us to register 2D images with a 3D model at interactive rates. In our future work we would like to be able to handle scenes of general configuration not containing any major vanishing points. This would allow the exploration of registration algorithms in non-urban scenes.

## 5 Multiview Pose Estimation and 3D Structure Reconstruction

In this section, we describe our algorithms for the generation of a sparse 3D model (Msfm) from a sequence of images. The input to our system is a sequence  $\mathbf{I} = \{I_n \mid n = 1, \dots, N\}$  of high resolution still images that capture the 3D scene. This is necessary to produce photorealistic scene representations. Therefore we have to attack the problem of finding correspondences in a sequence of wide-baseline high resolution images, a problem that is much harder than feature tracking from a video sequence. Fortunately, there are several recent approaches that attack the wide-baseline matching problem (Schaffalitzky and Zisserman 2001; Tuytelaars and Gool 2004; Lowe 2004). For the purposes of our system, we have adopted the scale-invariant feature transform (SIFT) method (Lowe 2004) for pairwise

**Fig. 16** Registration results from the interior of building 3 (Grand Central Terminal, NYC). **(a)** Registered point cloud of building 3. **(b)** For description see caption of Fig. 14. **(c)** (*Top row*): The user rotates the 3D model so that it is orientated similarly (note that it does not have to be exactly matched) to the 2D image. (*Bottom row*) The *right image* shows the 2D image along with the matched 2D and projected 3D features (see caption of Fig. 14). The *left image* shows the texture-mapped 3D range model after successful registration. Note that surfaces that are not part of the 3D model cannot be texture-mapped and appear as black holes. For example the floor is missing from our range model



feature extraction and matching. In general, structure from motion (SfM) from a set of images has been rigorously studied (Faugeras et al. 2001; Hartley and Zisserman 2003; Ma et al. 2003). Our method for pose estimation and partial structure recovery is based on sequential updating. The method is similar to work explained in (Beardsley et al. 1997; Pollefeys et al. 2004). In order to get very accurate pose estimation, we assume that the camera(s) are pre-calibrated. It

is, of course, possible to recover unknown and varying focal length by first recovering pose and structure up to an unknown projective transform and then upgrading to Euclidean space as shown in (Heyden and Astrom 1996; Pollefeys and Gool 1997; Triggs 1996). However, some of the assumptions that these methods make (e.g., no skew, approximate knowledge of the aspect ratio and principal point) may produce visible mismatches in a high resolution texture map. Thus,

for the sake of accuracy we are utilizing the camera calibration method of Zhang (2000). More details of our structure-from-motion method can be found in (Liu et al. 2006).

## 6 Alignment of 2D Image Sequences onto 3D-Range Point Clouds

The 2D-to-3D registration module described in Sect. 4 facilitates the texturing of the 3D range model  $M_{\text{range}}$  with photographs taken of the scene. A drawback to this approach is that only a subset of the images are successfully registered with  $M_{\text{range}}$  and so the texturing is restricted to these images for which the camera pose information is recovered. In order to more fully exploit the use of all images, we opt to use multiview geometry to derive a dense set of camera pose information and a sparse 3D model directly from the images. Merging the results of Sect. 4 with a structure-from-motion algorithm will produce a superior phototextured result that captures the full detail of a complete collection of photographs. We therefore create a sparse 3D model,  $M_{\text{sfm}}$  (see Sect. 5), and register it to  $M_{\text{range}}$ . The rotation, scale, and translation necessary to align these two models also applies to the dense set of camera poses derived from the images. After applying this transformation to the camera poses derived from structure-from-motion, the cameras will sit in the coordinate frame of  $M_{\text{range}}$  where images can then be projected onto the dense model. The method is described more fully below.

The set of dense range scans  $\{S_i \mid i = 1, \dots, K\}$  are registered in the same reference frame (Sect. 3), producing a 3D range model called  $M_{\text{range}}$ . On the other hand, the sequence of 2D images  $\mathbf{I} = \{I_n \mid n = 1, \dots, N\}$  produces a sparser 3D model of the scene (Sect. 5) called  $M_{\text{sfm}}$ . Both of these models are represented as clouds of 3D points. The distance between any two points in  $M_{\text{range}}$  corresponds to the actual distance of the points in 3D space, whereas the distance of any two points in  $M_{\text{sfm}}$  is the actual distance multiplied by an unknown scale factor  $s$ . In order to align the two models a similarity transformation that includes the scale factor  $s$ , a rotation  $R$  and a translation  $T$  needs to be computed. In this section, a novel algorithm (Liu et al. 2006) that automatically computes this transformation is presented. The transformation allows for the optimal texture mapping of all images onto the dense  $M_{\text{range}}$  model, and thus provides photorealistic results of high quality.

Every point  $X$  from  $M_{\text{sfm}}$  can be projected onto a 2D image  $I_n \in \mathbf{I}$  by the following transformation:

$$\mathbf{x} = \mathcal{K}_n[R_n \mid T_n]\mathbf{X} \quad (1)$$

where  $\mathbf{x} = (x, y, 1)$  is a pixel on image  $I_n$ ,  $\mathbf{X} = (X, Y, Z, 1)$  is a point of  $M_{\text{sfm}}$ ,  $\mathcal{K}_n$  is the projection matrix,  $R_n$  is the rotation transformation and  $T_n$  is the translation vector. These

matrices and points  $\mathbf{X}$  are computed by the SFM method (Sect. 5).

Some of the 2D images  $\mathbf{I}' \subset \mathbf{I}$  are also automatically registered with the 3D range model  $M_{\text{range}}$  (Sect. 4). Thus, each point of  $M_{\text{range}}$  can be projected onto each 2D image  $I_n \in \mathbf{I}'$  by the following transformation:

$$\mathbf{y} = \mathcal{K}_n[R'_n \mid T'_n]\mathbf{Y} \quad (2)$$

where  $\mathbf{y} = (x, y, 1)$  is a pixel in image  $I_n$ ,  $\mathbf{Y} = (X, Y, Z, 1)$  is a point of model  $M_{\text{range}}$ ,  $\mathcal{K}_n$  is the projection matrix of  $I_n$ ,  $R'_n$  is the rotation and  $T'_n$  is the translation. These transformations are computed by the 2D-to-3D registration method (Sect. 4).

The key idea is to use the images in  $I_n \in \mathbf{I}'$  as references in order to find the corresponding points between  $M_{\text{range}}$  and  $M_{\text{sfm}}$ . The similarity transformation between  $M_{\text{range}}$  and  $M_{\text{sfm}}$  is then computed based on these correspondences. In summary, the algorithm works as follows:

1. Each point of  $M_{\text{sfm}}$  is projected onto  $I_n \in \mathbf{I}'$  using (1). Each pixel  $p_{(i,j)}$  of  $I_n$  is associated with the closest projected point  $\mathbf{X} \in M_{\text{sfm}}$  in an  $L \times L$  neighborhood on the image. Each point of  $M_{\text{range}}$  is also projected onto  $I_n$  using (2). Similarly, each pixel  $p_{(i,j)}$  is associated with the projected point  $\mathbf{Y} \in M_{\text{range}}$  in an  $L \times L$  neighborhood (Fig. 17). Z-buffering is used to handle occlusions.
2. If a pixel  $p_{(i,j)}$  of image  $I_n$  is associated with a pair of 3D points  $(\mathbf{X}, \mathbf{Y})$ , one from  $M_{\text{sfm}}$  and the other from  $M_{\text{range}}$ , then these two 3D points are considered as candidate matches. Thus, for each 2D-image in  $\mathbf{I}'$  a set of matches is computed, producing a collection of candidate matches named  $\mathbf{L}$ . These 3D-3D correspondences between points of  $M_{\text{range}}$  and points of  $M_{\text{sfm}}$  could be potentially used for the computation of the similarity transformation between the two models. The set  $\mathbf{L}$  contains many outliers, due to the very simple closest-point algorithm utilized. However,  $\mathbf{L}$  can be further refined (Sect. 6.1) into a set of robust 3D point correspondences  $\mathcal{C} \subset \mathbf{L}$ .
3. Finally, the transformation between  $M_{\text{range}}$  and  $M_{\text{sfm}}$  is computed by minimizing a weighted error function  $E$  (Sect. 6.1) based on the final robust set of correspondences  $\mathcal{C}$ .

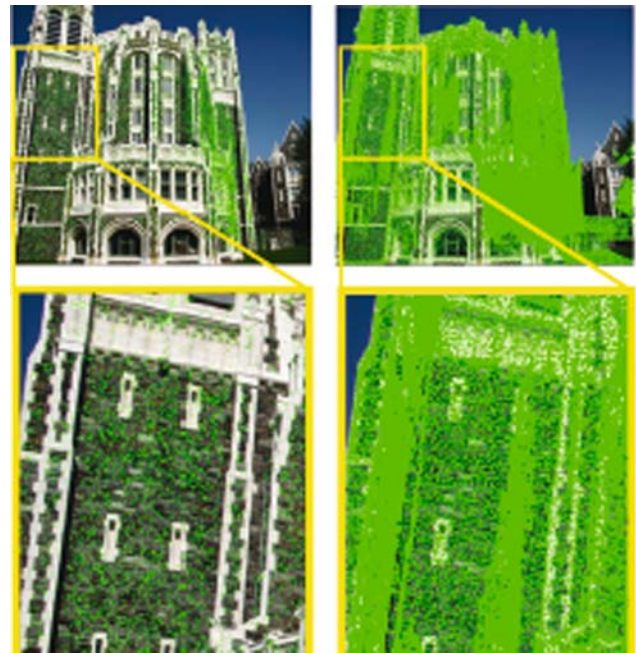
Even though the two models ( $M_{\text{sfm}}$  and  $M_{\text{range}}$ ) have different noise characteristics, they have similar errors in the operating range of our experiments. We ultimately exploit principles from close-range photogrammetry and therefore we do not use photos that are too far from the object. In that case, points that are recovered that are deemed to be too far away (which also may be due to small baseline) will be dismissed. The remaining projected points from the SFM model can therefore be matched against the projections of the 3D points from the range data.

**Table 3** Building 1 (13 images) - Building 2 (7 images) - Building 3 (6 images). Each row presents results from successful registration of a different 2D image with the 3D range model. The upper part of the table presents results of the registration of 13 images with a 3D range model of building 1. The middle part shows results from registering 7 images with a 3D range model of building 2. Finally, the lower part describes results from the registration of 6 images with the 3D range model of building 3. The registration (matching phase) of each image requires on average 5 to 10 seconds (2 GHz Xeon Intel processor, 2 GB of RAM). The first two columns show the numbers of 3D and 2D features used for matching. “Fi” is the initial focal length extracted from the Exif meta-data of the image, while “Fr” is the refined focal length. “M” is the number of matched features of the best transformation. Finally, “E” is the average line-to-line distance (in pixels) after the optimization (Step 6)

F3D	F2D	Fi	Fr	M	E
672	412	3065.83	3072.42	119	4.4492
583	345	3065.83	3075.34	103	4.9394
409	390	3065.83	3071.90	112	4.8973
392	230	3065.83	3069.45	93	4.2109
321	312	3065.83	3073.23	187	4.9021
456	387	3065.83	3072.12	134	4.3902
402	390	3065.83	3071.29	94	3.9827
390	219	3065.83	3069.22	87	4.2023
592	539	3065.83	3071.90	212	4.3003
390	416	3065.83	3061.39	145	3.9203
271	392	3065.83	3073.38	123	3.2900
430	456	3065.83	3076.19	209	4.1293
390	549	3065.83	3063.56	115	4.5902
<hr/>					
438	789	1185.03	1165.65	114	4.3215
421	654	1185.03	1175.89	83	4.2142
389	520	1185.03	1172.90	88	3.8992
402	432	1185.03	1179.34	101	4.2390
389	598	1185.03	1172.90	91	4.5009
435	621	1185.03	1169.39	156	4.1290
419	535	1185.03	1178.17	182	4.4923
<hr/>					
543	245	2805.81	2833.45	63	4.4439
569	312	2805.81	2831.32	45	3.9082
389	245	2805.81	2829.39	42	4.2312
390	190	2805.81	2839.93	50	4.9821
493	231	2805.81	2812.24	63	3.9023
301	189	2805.81	2829.39	58	3.8910

### 6.1 Correspondence Refinement and Optimization

The set of candidate matches  $\mathbf{L}$  computed in the second step of the previous algorithm contains outliers due to errors introduced from the various modules of the system (SFM, 2D-to-3D registration, range sensing). It is thus important to filter out as many outliers as possible through verification procedures. A natural verification procedure involves the difference in scale between the two models. Consider two pairs of plausible matched 3D-points  $(\mathbf{X}_1, \mathbf{Y}_1)$  and  $(\mathbf{X}_2, \mathbf{Y}_2)$  ( $\mathbf{X}_i$  de-



**Fig. 17** *Left column:* The points of model  $M_{\text{sfm}}$  projected onto one 2D image  $I_n$  (Sect. 5). The projected points are shown in green. *Right column:* The points of model  $M_{\text{range}}$  projected onto the same 2D image  $I_n$  (projected points shown in green) after the automatic 2D-to-3D registration (Sect. 4). Note that the density of 3D range points is much higher than the density of the SFM points, due to the different nature of the two reconstruction processes. Finding corresponding points between  $M_{\text{range}}$  and  $M_{\text{sfm}}$  is possible on the 2D image space of  $I_n$ . This yields the transformation between the two models (Sect. 6)

notes points from the  $M_{\text{sfm}}$  model, while  $\mathbf{Y}_j$  points from the  $M_{\text{range}}$  model). If these were indeed correct correspondences, then the scale factor between the two models would be  $s = \|\mathbf{X}_1 - \mathbf{X}_2\| / \|\mathbf{Y}_1 - \mathbf{Y}_2\|$ . Since the computed scale factor should be the same no matter which correct matching pair is used, then a robust set of correspondences from  $\mathbf{L}$  should contain only these pairs that produce the same scale factor  $s$ . The constant scale factor among correctly picked pairs is thus an invariant feature that we exploit. We now explain how we achieve this robust set of correspondences.

For each image  $I_n \in \mathbf{I}'$ , let us call the camera’s center of projection as  $\mathbf{C}_n^{\text{sfm}}$  in the local coordinate system of  $M_{\text{sfm}}$  and  $\mathbf{C}_n^{\text{mg}}$  in the coordinate system of  $M_{\text{range}}$ . These two centers have been computed from two independent processes: SFM (Sect. 5) and 2D-to-3D registration (Sect. 4). Then for any candidate match,  $(\mathbf{X}, \mathbf{Y}) \in \mathbf{L}$ , a candidate scale factor  $s_1(\mathbf{X}, \mathbf{Y})$  can be computed as:

$$s_1(\mathbf{X}, \mathbf{Y}) = \frac{\|\mathbf{X} - \mathbf{C}_n^{\text{sfm}}\|}{\|\mathbf{Y} - \mathbf{C}_n^{\text{mg}}\|}.$$

If we keep the match  $(\mathbf{X}, \mathbf{Y})$  fixed and we consider every other match  $(\mathbf{X}', \mathbf{Y}') \in \mathbf{L}$ ,  $L - 1$  candidate scale factors

$s_2(\mathbf{X}', \mathbf{Y}')$  and  $L - 1$  candidate scale factors  $s_3(\mathbf{X}', \mathbf{Y}')$  ( $L$  is the number of matches in  $\mathbf{L}$ ) are computed as:

$$s_2(\mathbf{X}', \mathbf{Y}') = \frac{\|\mathbf{X}' - \mathbf{C}_n^{\text{sfm}}\|}{\|\mathbf{Y}' - \mathbf{C}_n^{\text{rng}}\|}, \quad s_3(\mathbf{X}', \mathbf{Y}') = \frac{\|\mathbf{X} - \mathbf{X}'\|}{\|\mathbf{Y} - \mathbf{Y}'\|}.$$

That means that if we keep the match  $(\mathbf{X}, \mathbf{Y})$  fixed, and consider all other matches  $(\mathbf{X}', \mathbf{Y}')$  we can compute a triple of candidate scale factors:  $s_1(\mathbf{X}, \mathbf{Y}), s_2(\mathbf{X}', \mathbf{Y}')$ , and  $s_3(\mathbf{X}', \mathbf{Y}')$ . We then consider the two pairs of matches  $(\mathbf{X}, \mathbf{Y})$  and  $(\mathbf{X}', \mathbf{Y}')$  as *compatible* if the scale factors in the above triple are close with respect to each other. By fixing  $(\mathbf{X}, \mathbf{Y})$ , all matches that are compatible with it are found. The confidence in the match  $(\mathbf{X}, \mathbf{Y})$  is the number of compatible matches it has. By going through all matches in  $\mathbf{L}$ , their confidence is computed via the above procedure. Out of these matches the one with the highest confidence is selected as the most prominent:  $(\mathbf{X}_p, \mathbf{Y}_p)$ . Let us call  $\mathbf{L}_n$  the set that contains  $(\mathbf{X}_p, \mathbf{Y}_p)$  and all other matches that are compatible with it. Note that this set is based on the centers of projection of image  $I_n$  as computed by SFM and 2D-to-3D registration. Let us also call  $s_n$  the scale factor that corresponds to the set  $\mathbf{L}_n$ . This scale factor can be computed by averaging the triples of scale factors of the elements in  $\mathbf{L}_n$ . Finally a different set  $\mathbf{L}_n$  and scale factor  $s_n$  is computed for every image  $I_n \in \mathbf{I}'$ .

From the previous discussion it is clear that each  $\mathbf{L}_n$  is a set of matches that is based on the center of projection of each image  $I_n$  independently. A set of matches that will provide a globally optimal solution should consider all images of  $\mathbf{I}'$  simultaneously. Out of the scale factors computed from each set  $\mathbf{L}_n$ , the one that corresponds to the largest number of matches is the one more robustly extracted by the above procedure. That computed scale factor,  $s_{\text{opt}}$ , is used as the final filtration for the production of the robust set of matches  $\mathcal{C}$  out of  $\mathbf{L}$ . In particular, for each candidate match  $(\mathbf{X}, \mathbf{Y}) \in \mathbf{L}$ , a set of scale factors are computed as

$$s'_n = \frac{\|\mathbf{X} - \mathbf{C}_n^{\text{sfm}}\|}{\|\mathbf{Y} - \mathbf{C}_n^{\text{rng}}\|}$$

where  $n = 1, 2, \dots, J$ , and  $J$  is the number of images in  $\mathbf{I}'$ . The standard deviation of those scale factors with respect to  $s_{\text{opt}}$  is computed and if it is smaller than a user-defined threshold,  $(\mathbf{X}, \mathbf{Y})$  is considered as a *robust* match and is added to the final list of correspondences  $\mathcal{C}$ . The robustness of the match stems from the fact that it verifies the robustly extracted scale factor  $s_{\text{opt}}$  with respect to most (or all) images  $I_n \in \mathbf{I}'$ . The pairs of center of projections  $(\mathbf{C}_n^{\text{sfm}}, \mathbf{C}_n^{\text{rng}})$  of images in  $\mathbf{I}'$  are also added to  $\mathcal{C}$ .

The list  $\mathcal{C}$  contains robust 3D point correspondences that are used for the accurate computation of the similarity transformation (scale factor  $s$ , rotation  $R$ , and translation  $T$ ) be-

tween the models  $M_{\text{range}}$  and  $M_{\text{sfm}}$ . The following weighted error function is minimized with respect to  $sR$  and  $T$ :

$$E = \sum_{(\mathbf{X}, \mathbf{Y}) \in \mathcal{C}} w \|sR \cdot \mathbf{Y} + T - \mathbf{X}\|^2$$

where the weight  $w = 1$  for all  $(\mathbf{X}, \mathbf{Y}) \in \mathcal{C}$  that are not the centers of projection of the cameras, and  $w > 1$  (user-defined) when  $(\mathbf{X}, \mathbf{Y}) = (\mathbf{C}_n^{\text{sfm}}, \mathbf{C}_n^{\text{rng}})$ . By associating higher weights to the centers we exploit the fact that we are confident in the original pose produced by SFM and 2D-to-3D registration. The unknown  $sR$  and  $T$  are estimated by computing the least square solution from this error function. Note that  $s$  can be easily extracted from  $sR$  since the determinant of  $R$  is 1.

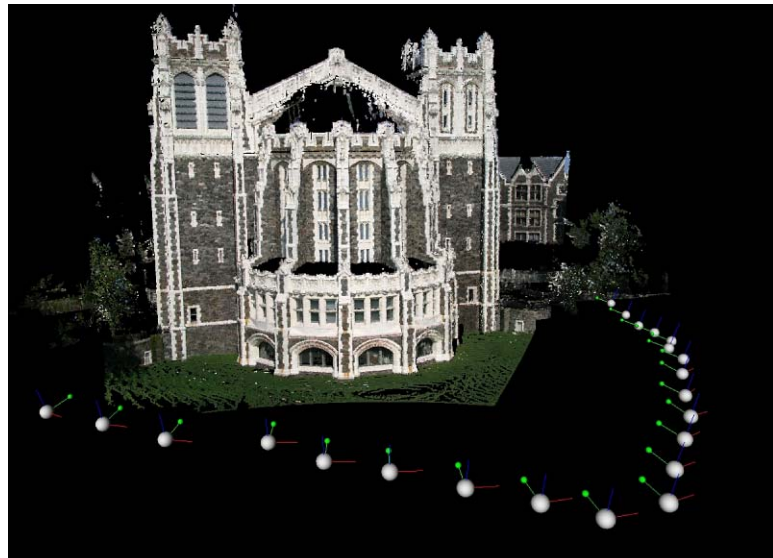
In summary, by utilizing the invariance of the scale factor between corresponding points in  $M_{\text{range}}$  and  $M_{\text{sfm}}$ , a set of robust 3D point correspondences  $\mathcal{C}$  is computed. These 3D point correspondences are then used for an optimal calculation of the similarity transformation between the two point clouds. This provides a very accurate texture mapping result of the high resolution images onto the dense range model  $M_{\text{range}}$ .

### 7 Results & Conclusions

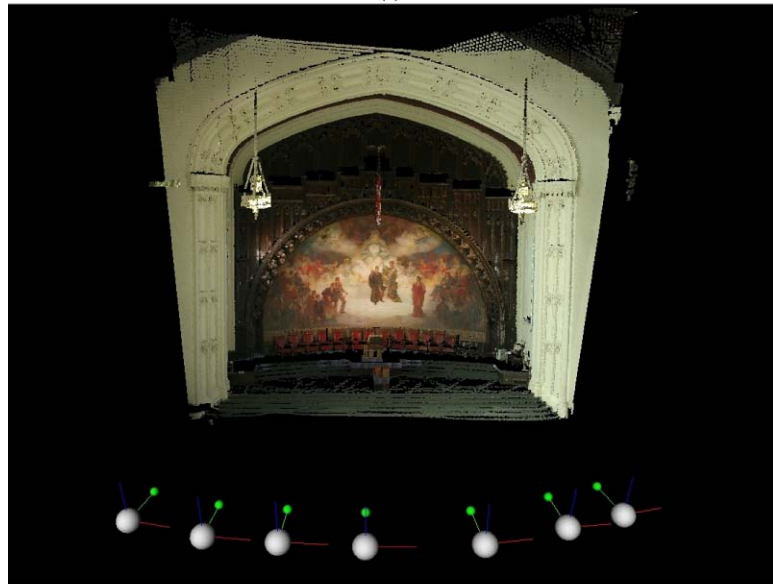
We tested our algorithms using range scans and 2D images acquired from a large-scale urban structure (Shepard Hall/CCNY) and from an interior scene (Great Hall/CCNY). 22 range scans of the exterior of Shepard Hall were automatically registered (Fig. 1) to produce a dense model  $M_{\text{range}}$ . In one experiment, ten images were gathered under the same lighting conditions. All ten of them were independently registered (2D-to-3D registration Sect. 4) with the model  $M_{\text{range}}$ . The registration was optimized with the incorporation of the SFM model (Sect. 5) and the final optimization method (Sect. 6). In a second experiment, 22 images of Shepard Hall that covered a wider area were acquired. Although the automated 2D-to-3D registration method was applied to all the images, only five of them were manually selected for the final transformation (Sect. 6) on the basis of visual accuracy. For some of the 22 images the automated 2D-to-3D method could not be applied due to lack of linear features. However, all 22 images were optimally registered using our novel registration method (Sect. 6) after the SFM computation (Sect. 5). Figure 1 shows the alignment of the range and SFM models achieved through the use of the 2D images. In Fig. 18(a) the accuracy of the texture mapping method is visible. Figure 18(b) displays a similar result of an interior 3D scene.<sup>8</sup> Table 4 provides some quantitative

<sup>8</sup><http://www.cs.hunter.cuny.edu/~ioannis/IjcvSI> contains images and videos of our results.

**Fig. 18** (a) Range model of Shepard Hall (CCNY) with 22 automatically texture mapped high resolution images. (b) Range model of interior scene (Great Hall) with seven automatically texture mapped images. The set of camera poses solved via multiview geometry are shown. Notice the accuracy of the photorealistic result



(a)



(b)

results of our experiments. Notice the density of the range models versus the sparsity of the SFM models. Also notice the number of robust matches in  $\mathcal{C}$  (Sect. 6) with respect to the possible number of matches (i.e., number of points in SFM). The final row in Table 4 displays the elapsed time for the final optimization on a Dell PC running Linux on an Intel Xeon-2 GHz, 2 GB-RAM machine.

We have presented a system that integrates multiview geometry and automated 3D registration techniques for texture mapping high resolution 2D images onto dense 3D range data. The benefits of multiview geometry (SFM) and automated 2D-to-3D/3D-to-3D registration are merged for the production of photorealistic models with minimal human interaction. Our approach provides a solution of increased

**Table 4** Quantitative results

	Shepard Hall		Great Hall
Number of points ( $M_{\text{range}}$ )	12,483,568		13,234,532
Number of points ( $M_{\text{sfm}}$ )	2,034	45,392	1,655
2D-images used	10	22	7
2D-to-3D registrations (Sect. 4)	10	5	3
No. of matches in $\mathcal{C}$ (Sect. 6)	258	1632	156
Final optimization (Sect. 6)	8.65 s	19.20 s	3.18 s

robustness, efficiency and generality with respect to previous methods.

**Acknowledgements** This work was supported in part by the following grants: NSF CAREER IIS-02-37878, NSF MRI/RUI EIA-0215962, and DOE DE-FG52-06NA27503.

## References

- Ambler, A. P., Barrow, H. G., Brown, C. M., Burstall, R. M., & Poplestone, R. J. (1973). A versatile computer-controlled assembly system. In *The 3rd international joint conference on artificial intelligence*.
- Beardsley, P. A., Zisserman, A. P., & Murray, D. W. (1997). Sequential updating of projective and affine structure from motion. *International Journal of Computer Vision* 23(3), 235–259.
- Bellon, O. R. P., & Silva, L. (2002). New improvements to range image segmentation by edge detection. *IEEE Signal Processing Letters* 9(2), 23–45.
- Bernardini, F., & Rushmeier, H. (2002). The 3D model acquisition pipeline. *Computer Graphics Forum* 21(2), 149–172.
- Besl, P. J., & Jain, R. C. (1988). Segmentation through variable-order surface fitting. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 10(2), 167–192.
- Besl, P. J., & McKay, N. D. (1992). A method for registration of 3D shapes. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 14(2) (1992).
- Cass, T. (1997). Polynomial-time geometric matching for object recognition. *International Journal of Computer Vision* 21(1–2), 37–61.
- Chen, C. (2007). *Range segmentation and registration for 3D modeling of large-scale urban scenes*. PhD thesis, City University of New York.
- Chen, C., & Stamos, I. (2005). Semi-automatic range to range registration: a feature-based method. In *The 5th international conference on 3-D digital imaging and modeling* (pp. 254–261), Ottawa, June 2005.
- Chen, C., & Stamos, I. (2006). Range image registration based on circular features. In *International symposium on 3D data processing, visualization and transmission*, Chapel Hill, June 2006.
- Christy, S., & Horaud, R. (1999). Iterative pose computation from line correspondences. *Journal of Computer Vision and Image Understanding* 73(1), 137–144.
- Faugeras, O. (1996). *Three-dimensional computer vision*. Cambridge: MIT Press.
- Faugeras, O., Luong, Q. T., & Papadopoulos, T. (2001). *The geometry of multiple images*. Cambridge: MIT Press.
- Fischler, M. A., & Bolles, R. C. (1981). Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Graphics and Image Processing* 24(6), 381–395.
- Früh, C., & Zakhor, A. (2003). Constructing 3D city models by merging aerial and ground views. *Computer Graphics and Applications* 23(6), 52–11.
- Hartley, R., & Zisserman, A. (2003). *Multiple view geometry in computer vision*. (2nd ed.). Cambridge: Cambridge University Press.
- Hausler, G., & Ritter, D. (1999). Feature-based object recognition and localization in 3D-space, using a single video image. *Journal of Computer Vision and Image Understanding* 73(1), 64–81.
- Heyden, A., & Astrom, K. (1996) Euclidean reconstruction from constant intrinsic parameters. In *International conference on pattern recognition* (Vol. 1, pp. 339–343) (1992).
- Horaud, R., Dornaika, F., Lamiro, B., & Christy, S. (1997). Object pose: the link between weak perspective, paraperspective, and full perspective. *International Journal of Computer Vision*, 22(2).
- Huber, D. F., & Hebert, M. (2003). Fully automatic registration of multiple 3D data sets. *Image and Vision Computing* 21(7), 637–650.
- Huttenlocher, D., & Ullman, S. (1990). Recognizing solid objects by alignment with an image. *International Journal of Computer Vision* 5(7), 195–212.
- Ikeuchi, K. (2003). The great Buddha project, In *IEEE ISMAR03*, Tokyo, Japan, November 2003.
- Jacobs, D. W. (1997). Matching 3-D models to 2-D images. *International Journal of Computer Vision* 21(1–2), 123–153.
- Johnson, A. E., & Hebert, M. (1999). Using spin images for efficient object recognition in cluttered 3D scenes. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 21(5), 433–449.
- Jurie, F. (1999). Solution of the simultaneous pose and correspondence problem using Gaussian error model. *Journal of Computer Vision and Image Understanding* 73(3), 357–373.
- Leica Geosystems. <http://hds.leica-geosystems.com/>.
- Liu, L. (2007). *Automated registration of 2D images with 3D range data in a photorealistic modeling system of urban scenes*. PhD thesis, City University of New York.
- Liu, L., & Stamos, I. (2005). Automatic 3D to 2D registration for the photorealistic rendering of urban scenes. In *IEEE conference on computer vision and pattern recognition* (Vol. II, pp. 137–143), San Diego, CA, June 2005.
- Liu, L., & Stamos, I. (2007). A systematic approach for 2D-image to 3D-range registration in urban environments. In *Virtual Representations and modeling of large-scale environments (VRML), in conjunction with ICCV*. Brazil: Rio De Janeiro.
- Liu, L., Stamos, I., Yu, G., Wolberg, G., & Zokai, S. (2006). Multi-view geometry for texture mapping 2D images onto 3D range data. In *IEEE conference on computer vision and pattern recognition* (Vol. II, pp. 2293–2300), New York City, June 2006.
- Lowe, D. G. (2004). Distinctive image features from scale-invariant keypoint. *International Journal of Computer Vision* 60(2), 91–110.
- Ma, Y., Soatto, S., Kosecka, J., & Sastry, S. (2003). *An invitation to 3-D vision: from images to geometric models*. Berlin: Springer.
- Marshall, G. L. D., & Martin, R. (2001). Robust segmentation of primitives from range data in the presence of geometric degeneracy. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 23(3), 304–314.
- Oberkampff, D., DeMenthon, D., & Davis, L. (1996) Iterative pose estimation using coplanar feature points. *Computer Vision Graphics and Image Processing*, 63(3) (1996).
- Pollefeys, M., & Gool, L. V. (1997). A stratified approach to metric self-calibration. In *IEEE conference on computer vision and pattern recognition* (pp. 407–412) (1997).
- Pollefeys, M., Gool, L. V., Vergauwen, M., Verbiest, F., Cornelis, K., Tops, J., & Koch, R. (2004). Visual modeling with a hand-held camera. *International Journal of Computer Vision* 59(3), 207–232.
- Pulli, K., & Pietikäinen, M. (1993). Range image segmentation based on decomposition of surface normals. In *Proceedings of the Scandinavian conference on image analysis*.
- Pulli, K., Abi-Rached, H., Duchamp, T., Shapiro, L. G., & Stuetzle, W. (1998). Acquisition and visualization of colored 3D objects. In *International conference on pattern recognition* (Vol. 1, p. 11), Australia, 1998.
- Quan, L., & Lan, Z. (1999). Linear N-point camera pose determination. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 21(7).
- Rusinkiewicz, S., & Levoy, M. (2001). Efficient variants of the ICP algorithm. In *The 3rd international conference on 3-D digital imaging and modeling*, June 2001.
- Schaffalitzky, F., & Zisserman, A. (2001). Viewpoint invariant texture matching and wide baseline stereo. In *IEEE international conference on computer vision* (pp. 636–643), July 2001.
- Sequeira, V., & Concalves, J. (2002). 3D reality modeling: photorealistic 3D models of real world scenes. In *International symposium on 3D data processing, visualization and transmission* (pp. 776–783).

- Stamos, I., & Allen, P. K. (2001). Automatic registration of 3-D with 2-D imagery in urban environments. In *IEEE International conference on computer vision* (pp. 731–736).
- Stamos, I., & Allen, P. K. (2002). Geometry and texture recovery of scenes of large scale. *Journal of Computer Vision and Image Understanding* 88(2), 94–118.
- Stamos, I., & Leordeanu, M. (2003). Automated feature-based range registration of urban scenes of large scale. In *IEEE conference on computer vision and pattern recognition* (Vol. 2, pp. 555–561).
- Triggs, B. (1996). Factorization methods for projective structure and motion. In *IEEE conference on computer vision and pattern recognition* (pp. 845–851).
- Troccoli, A., & Allen, P. K. (2004). A shadow based method for image to model registration. In *2nd IEEE workshop on video and image registration*, July 2004.
- Tuytelaars, T., & Gool, L. J. V. (2004). Matching widely separated views based on affine invariant regions. *International Journal of Computer Vision* 59(1), 61–85.
- Visual Information Technology Group, Canada. (2005). [http://iit-iti.nrc-cnrc.gc.ca/about-sujet/vit-tiv\\_e.html](http://iit-iti.nrc-cnrc.gc.ca/about-sujet/vit-tiv_e.html).
- Wami, M. A., & Batchelor, B. G. (1994). Edge-region-based segmentation of range images. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 16(3), 314–319.
- Wells, W. (1997). Statistical approaches to feature-based object recognition. *International Journal of Computer Vision* 21(1–2), 63–98.
- Zhang, Z. (2000). A flexible new technique for camera calibration. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 22(11), 1330–1334.
- Zhao, H., & Shibasaki, R. (2003). Reconstructing a textured CAD model of an urban environment using vehicle-borne laser range scanners and line cameras. *Machine Vision and Applications* 14(1), 35–41.
- Zhao, W., Nister, D., & Hsu, S. (2005). Alignment of continuous video onto 3D point clouds. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 27(8), 1305–1318.