# Data Processing Algorithms for Generating Textured 3D Building Facade Meshes from Laser Scans and Camera Images

CHRISTIAN FRUEH, SIDDHARTH JAIN AND AVIDEH ZAKHOR

*Video and Image Processing Laboratory, Department of Electrical Engineering and Computer Sciences,*
*University of California, Berkeley*

frueh@eecs.berkeley.edu
morpheus@eecs.berkeley.edu
avz@eecs.berkeley.edu

**Abstract.**    In this paper, we develop a set of data processing algorithms for generating textured facade meshes of cities from a series of vertical 2D surface scans and camera images, obtained by a laser scanner and digital camera while driving on public roads under normal traffic conditions. These processing steps are needed to cope with imperfections and non-idealities inherent in laser scanning systems such as occlusions and reflections from glass surfaces. The data is divided into easy-to-handle quasi-linear segments corresponding to approximately straight driving direction and sequential topological order of vertical laser scans; each segment is then transformed into a depth image. Dominant building structures are detected in the depth images, and points are classified into foreground and background layers. Large holes in the background layer, caused by occlusion from foreground layer objects, are filled in by planar or horizontal interpolation. The depth image is further processed by removing isolated points and filling remaining small holes. The foreground objects also leave holes in the texture of building facades, which are filled by horizontal and vertical interpolation in low frequency regions, or by a copy-paste method otherwise. We apply the above steps to a large set of data of downtown Berkeley with several million 3D points, in order to obtain texture-mapped 3D models.

## 1.  Introduction

Three-dimensional models of urban environments are useful in a variety of applications such as urban planning, training and simulation for urban terrorism scenarios, and virtual reality. Currently, the standard technique for creating large-scale city models in an automated or semi-automated way is to use stereo vision approaches on aerial or satellite images (Frere et al., 1998; Kim et al., 2001). In recent years, advances in resolution and accuracy of airborne laser scanners have also rendered them suitable for the generation of reasonable models (Haala and Brenner, 1997;

Maas, 2001). Both approaches have the disadvantage that their resolution is only in the range of 1 to 2 feet, and more importantly, they can only capture the roofs of the buildings but not the facades. This essential disadvantage prohibits their use in photo realistic walk or drive-through applications.

There exist a number of approaches to acquire the complementary ground-level data and to reconstruct building facades; however, these approaches are typically limited to one or few buildings. Debevec et al. (1996) propose to reconstruct buildings based on few camera images in a semi-automated way. Dick et al. (2001), Koch et al. (1999), and Wang et al.

160    *Frueh, Jain and Zakhor*

(2002) apply automated vision-based techniques for localization and model reconstruction, but varying lighting conditions, the scale of the environment, and the complexity of outdoor scenes with many trees and glass surfaces generally pose enormous challenges to purely vision-based methods.

Stamos and Allen (2002) use a 3D laser scanner and Thrun et al. (2000) use 2D laser scanners mounted on a mobile robot to achieve complete automation, but the time required for data acquisition of an entire city is prohibitively large; in addition, the reliability of autonomous mobile robots in outdoor environments is a critical issue. In Zhao and Shibasaki (1999), use a vertical laser scanner mounted on a van, which is localized by using odometry, an inertial navigation system, and the Global Positioning System (GPS), and thus with limited accuracy. While GPS is by far the most common source of global position estimates in outdoor environments, even expensive high-end Differential GPS systems become inaccurate or erroneous in urban canyons where there are not enough satellites in a direct line of sight.
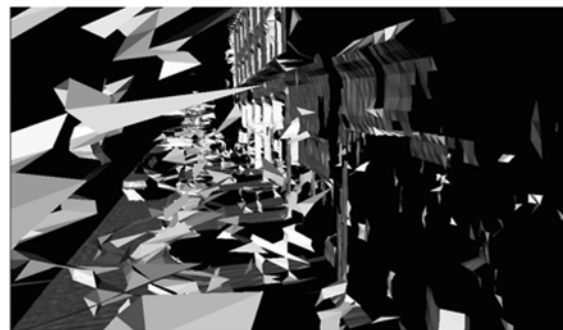
In previous work, we have developed a fast, automated data acquisition system capable of acquiring 3D geometry and texture data for an entire city at the ground level by using a combination of a horizontal and a vertical 2D laser scanners and a digital camera (Frueh et al., 2001; Frueh and Zakhor, 2001a). This system is mounted on a truck, moving at normal speeds on public roads, collecting data to be processed offline. It is similar to the one independently proposed by Zhao and Shibasaki (2001), which also use 2D laser scanners in horizontal and vertical configuration; however, our system differs from that of Zhao and Shibasaki (2001) in that we use a normal camera instead of a line camera. Both approaches have the advantage that data can be acquired continuously, rather than in a stop-and-go fashion, and are thus extremely fast; relative position changes are computed with centimeter accuracy by matching successive horizontal laser scans against each other. In Frueh and Zakhor (2001b), we proposed to use the particle-filtering-based Monte-Carlo Localization (Fox et al., 2000) to correct accumulating pose uncertainty by using airborne data such as an aerial photo or a digital surface model (DSM) as a map. An advantage of our approach is that both scan points and camera images are registered with airborne data, facilitating a subsequent fusion with models derived from this data (Frueh and Zakhor, 2003).

In this paper, we describe our approach to processing the globally registered scan points and camera images obtained in our ground-based data acquisition, and to creating detailed, textured 3D facade models. As there are many erroneous scan points, e.g. due to glass surfaces, and foreground objects partially occluding the desired buildings, the generation of a facade mesh is not straightforward. A simple triangulation of the raw scan points by connecting neighboring points whose distance is below a threshold value does not result in an acceptable reconstruction of the street scenery, as shown in Figs. 1(a) and (b). Even though the 3D structure can be easily recognized when viewed from a viewpoint near the original acquisition position as in Fig. 1(a), the mesh appears cluttered due to several reasons; first, there are holes and erroneous vertices due to reflections off the glass on windows; second, there are pieces of geometry "floating in the air", corresponding to partially captured objects or measurement errors. The mesh appears to be even more problematic when viewed from other viewpoints such as the one shown in Fig. 1(b); this is because in this case the large holes in the building facades caused by occluding foreground



(a)



(b)

*Figure 1.* Triangulated raw points: (a) front view; (b) side view.

objects, such as cars and trees, become entirely visible. Furthermore, since the laser scan only captures the frontal view of foreground objects, they become almost unrecognizable when viewed sideways. As we drive by a street only once, it is not possible to use additional scans from other viewpoints to fill in gaps caused by occlusions, as is done in Curless and Levoy (1996) and Stamos and Allen (2002). Rather, we have to reconstruct occluded areas by using cues from neighboring scan points; as such, there has been little work to solve this problem (Stulp et al., 2001).

In this paper, we propose a class of data processing techniques to create visually appealing facade meshes by removing noisy foreground objects and filling holes in the geometry and texture of building facades. Our objectives are robustness and efficiency with regards to processing time, in order to ensure scalability to the enormous amount of data resulting from a city scan. The outline of this paper is as follows: In Section 2, we introduce our data acquisition system and position estimation; Section 3 discusses data subdivision and depth image generation schemes. We describe our strategy to transform the raw scans into a visually appealing facade mesh in Sections 4 through 6; Section 7 discusses foreground and background segmentation of images, automatic texture atlas generation, and texture synthesis. The experimental results are presented in Section 8.

## 2. Data acquisition and Position Estimation

As described in Frueh et al. (2001) and Frueh and Zakhor (2001a), we have developed a data acquisition system consisting of two Sick LMS 2D laser scanners, and a digital color camera with a wide-angle lens. As



*Figure 2.* Truck with data acquisition equipment.

seen in Fig. 2, this system is mounted on a rack approximately 3.6 meters high on top of a truck, in order to obtain measurements that are not obstructed by pedestrians and cars. The scanners have a $180°$ field of view with a resolution of $1°$, a range of 80 meters and an accuracy of $\pm 3.5$ centimeters. Both 2D scanners face the same side of the street and are mounted at a 90-degree angle. The first scanner is mounted vertically with the scanning plane orthogonal to the driving direction, and scans the buildings and street scenery as the truck drives by. The data captured by this scanner is used for reconstructing 3D geometry as described in this paper. The second scanner is mounted horizontally and is used for determining the position of the truck for each vertical scan. Finally, the digital camera is used to acquire the appearance of the scanned building facades. It is oriented in the same direction as the scanners, with its center of projection approximately in the intersection line of the two scanning planes. All three devices are synchronized with each other using hardware-generated signals, and their coordinate systems are calibrated with respect to each other prior to the acquisition. Thus, we obtain long series of vertical scans, horizontal scans and camera images that are all associated with each other.

We introduce a Cartesian world coordinate system $[x, y, z]$ where $x, y$ is the ground plane and $z$ points into the sky. While our truck performs a 6 degree-of-freedom motion, its primary motion components are $x$, $y$, and $\theta$ (yaw), i.e. its two-dimensional (2D) motion. As described in detail in Frueh and Zakhor (2001a), we reconstruct the driven path and determine the global pose for each scan by using the horizontal laser scanner: First, an estimate of the 2D relative pose $(\Delta x, \Delta y, \Delta\theta)$ between each pair of subsequent scans is obtained via scan-to-scan matching; these relative estimates are concatenated to form a preliminary estimate for the driven path. Then, in order to correct the global pose error resulting from accumulation of error due to relative estimates, we utilize an aerial image or a DSM as a global map, and apply Monte-Carlo-Localization (Frueh and Zakhor, 2001b). Matching ground-based horizontal laser scans with edges in the global map, we track the vehicle and correct the preliminary path accordingly to obtain a globally registered 2D trajectory as shown in Fig. 3. As described in Frueh and Zakhor (2003), we obtain the secondary motion components $z$ and pitch by utilizing the altitude information provided by the DSM, and the roll motion by correlating subsequent camera images, respectively.
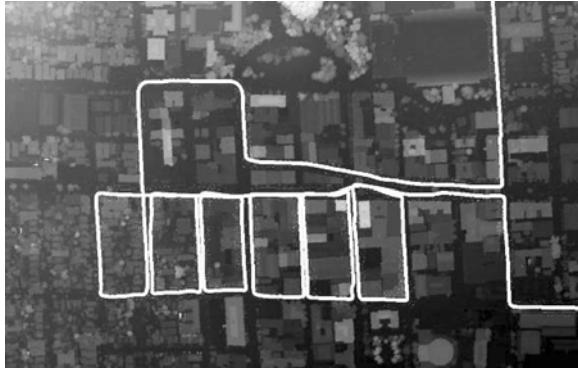
162    *Frueh, Jain and Zakhor*



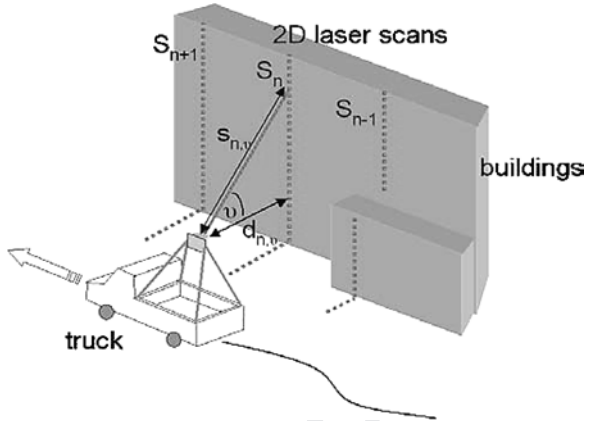*Figure 3.* Driven path superimposed on top of a DSM.



*Figure 4.* Scanning setup.

**203** While we use the full 6 degree-of-freedom pose to
**204** compute the final $x$, $y$, $z$ coordinates of each scan point
**205** in the final model, we can for convenience and sim-
**206** plicity neglect the 3 secondary motion components for
**207** most of the intermediate processing steps described in
**208** the following sections of this paper. Furthermore, to re-
**209** duce the amount of required processing and to partially
**210** compensate for the unpredictable, non-uniform speed
**211** of the truck, we do not utilize all the scans captured
**212** during slow motion; rather, we subsample the series of
**213** vertical scans such that the spacing between succes-
**214** sive scans is roughly equidistant. Thus, in our process-
**215** ing steps described in this paper, we assume the scan
**216** data to be given as a series of roughly equally spaced
**217** vertical scans $S_n$ with an associated tuple $(x_n, y_n, \theta_n)$
**218** describing 2D position and orientation of the scanner
**219** in the world coordinate system during acquisition. Fur-
**220** thermore, we use $s_{n,\upsilon}$ to denote the distance measure-
**221** ment on a point in scan $S_n$ with azimuth angle $\upsilon$, and
**222** $d_{n,\upsilon} = \cos(\upsilon) \cdot s_{n,\upsilon}$ to denote the depth value of this
**223** point with respect to the scanner, i.e. its orthogonal
**224** projection into the ground plane, as shown in Fig. 4.

**225** **3.    Data Subdivision and Depth Image Generation**

**226** *3.1.    Segmentation of the Driving Path into Quasi*
**227** *Linear Segments*

**228** The captured data during a 20-minute drive consists
**229** of tens of thousands of vertical scan columns. Since
**230** successive scans in time correspond to spatially close
**231** points, e.g. a building or a side of a street block, it is
**232** computationally advantageous not to process the entire
**233** data as one block, rather to split it into smaller segments
**234** to be processed separately. We impose the constraints

that (a) path segments have low curvature, and (b) scan **235**
columns have a regular grid structure. This allows us **236**
to readily identify the neighbors to right, left, above **237**
and below for each point, and, as seen later, is essential **238**
for the generation of a depth image and segmentation **239**
operations. **240**

Scan points for each truck position are obtained as **241**
we drive by the streets. During straight segments, the **242**
spatial order of the 2D scan rows is identical to the **243**
temporal order of the scans, forming a regular topol- **244**
ogy. Unfortunately, this order of scan points can be **245**
reversed during turns towards the scanner side of the **246**
car. Figure 5(a) and (b) show the scanning setup dur- **247**
ing such a turn, with scan planes indicated by the two **248**
dotted rays. During the two vertical scans, the truck per- **249**
forms not only a translation but also a rotation, making **250**
the scanner look slightly backwards during the second **251**
scan. If the targeted object is close enough, as shown in **252**
Fig. 5(a), the spatial order of scan points 1 and 2 is still **253**
the same as the temporal order of the scans; however, if **254**
the object is further away than a critical distance $d_{\text{crit}}$, **255**
the spatial order of the two scan points is reversed, as **256**
shown in Fig. 5(b). **257**

For a given truck translation of $\Delta s$, and a rotation **258**
$\Delta\theta$ between successive scans, the critical distance can **259**
be computed as **260**

$$d_{\text{crit}} = \frac{\Delta s}{\sin(\Delta\theta)}.$$

Thus, $d_{\text{crit}}$ is the distance at which the second scan- **261**
ning plane intersects with the first scanning plane. For **262**
a particular scan point, the order with its predecessors **263**
should be reversed if its depth $d_{n,\upsilon}$ exceeds $d_{\text{crit}}$; this **264**
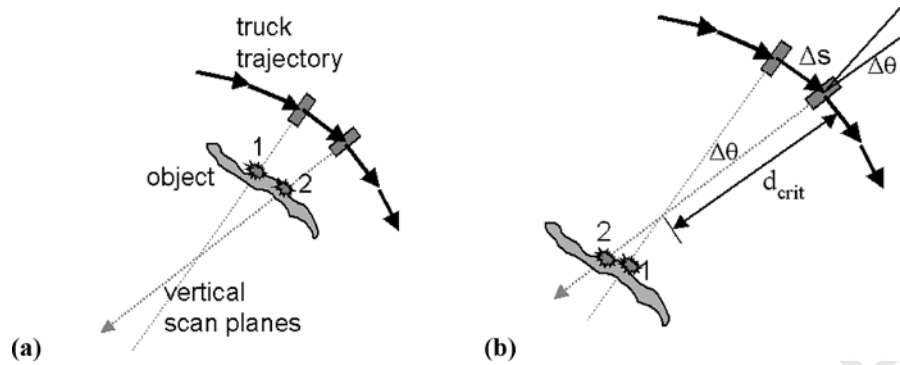
*Figure 5.* Scan geometry during a turn: (a) normal scan order for closer objects; (b) reversed scan order for farther objects.
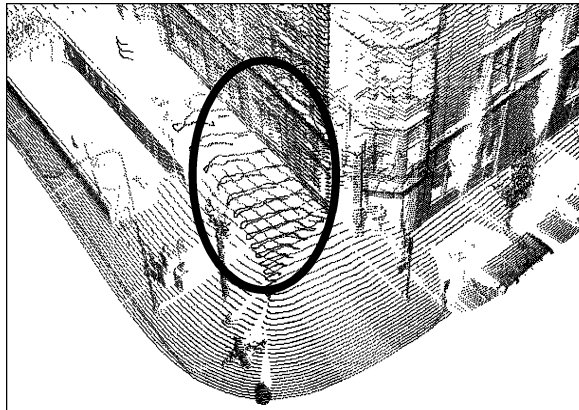


*Figure 6.* Scan points with reversed order.



*Figure 7.* Driven path: (a) before segmentation; (b) after segmentation into quasi-linear segments.

**265** means that its geometric location is somewhere in be-
**266** tween points of previous scans. The effect of such order
**267** reversal can be seen in the marked area in Fig. 6. At the
**268** corner, the ground and the building walls are scanned
**269** twice, first from a direct view and then from an oblique
**270** angle, and hence with significantly lower accuracy. For
**271** the oblique points, the scans are out of order, destroy-
**272** ing the regular topology between neighboring scan
**273** points.
**274** Since the "out of order" scans obtained in these sce-
**275** narios correspond to points that have already been cap-
**276** tured by "in order" scans, and are therefore redundant,
**277** our approach is to discard them and use only "in or-
**278** der" scans. For typical values of displacement, turn-
**279** ing angle, and distance of structures from our driving
**280** path, this occurs only in scans of turns with significant
**281** angular changes. By removing these "turn" scans and
**282** splitting the path at the "turning points", we obtain path
**283** segments with low curvature that can be considered as
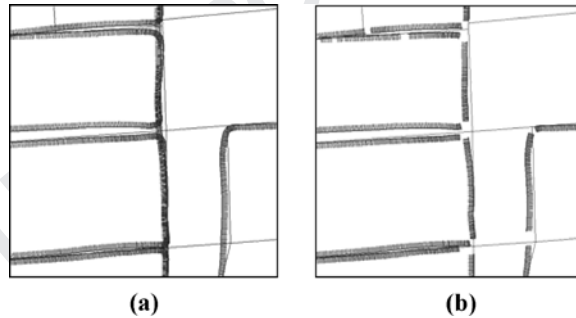**284** locally quasi-linear, and can therefore be conveniently

**285** processed as depth images, as described later in this
**286** section. In addition, to ensure that these segments are
**287** not too large for further processing, we subdivide them
**288** if they are larger then a certain size; specifically, in
**289** segments that are longer than 100 meters, we identify
**290** vertical scans that have the fewest scan points above
**291** street level, corresponding to gaps between buildings,
**292** and segment at these locations. Furthermore, we detect
**293** redundant path segments for areas captured multiple
**294** times due to multiple drive-bys, and use only one of
**295** them for reconstruction purposes. Figures 7(a) and (b)
**296** show an example of an original path, and the resulting
**297** path segments overlaid on a road map, respectively.
**298** The small lines perpendicular to the driving path indi-
**299** cate the scanning plane of the vertical scanner for each
**300** position.

### 3.2. Converting Path Segments into Depth Images    **301**

In the previous subsection, we described how to create **302**
path segments that are guaranteed to contain no scan **303**

164    *Frueh, Jain and Zakhor*

pairs with permuted horizontal order. As the vertical order is inherent to the scan itself, all scan points of a segment form a 3D scan grid with regular, quadrilateral topology. This 3D scan grid allows us to transform the scan points into a depth image, i.e. a 2.5D representation where each pixel represents a scan point, and the gray value for each pixel is proportional to the depth of the scan point. The advantage of a depth image is its intuitively easy interpretation, and the increased processing speed the 2D domain provides. However, most operations that are performed on the depth image can be done just as well on the 3D point grid directly, only not as conveniently.

A depth image is typically used for representing data from 3D scanners. Even though the way the depth value is assigned to each pixel is dependent on the specific scanner, in most cases it is the distance between scan point and scanner origin, or its cosine with respect to the ground plane. As we expect mainly vertical structures, we choose the latter option and use the depth $d_{n,v} = \cos(v) \cdot s_{n,v}$ rather than the distance $s_{n,v}$, so that the depth image is basically a tilted height field. The advantage is that in this case points that lie on a vertical line, e.g. a building wall, have the same depth value, and are hence easy to detect and group. Note that our depth image differs from one that would be obtained from a normal 3D scanner, as it does not have a single center from which the scan points are measured; instead, there are different centers for each individual vertical column along the path segment. The obtained depth image is neither a polar nor a parallel projection; it resembles most to a cylindrical projection. Due to non-uniform driving speed and non-linear driving direction, these centers are in general not on a line, but on an arbitrary shaped, though low-curvature curve, and the spacing between them is not exactly uniform. Because of this, strictly speaking the grid position only specifies the topological order of the depth pixels, and not the exact 3D point coordinates. However, as topology and depth value are a good approximation for the exact 3D coordinates, especially within a small neighborhood, we choose to apply our data processing algorithms to the depth image, thereby facilitating use of standard image processing techniques such as region growing. Moreover, the actual 3D vertex coordinates are still kept and used for 3D operations such as plane fitting. Figure 8(a) shows an example of the 3D vertices of a scan grid, and Fig. 8(b) shows its corresponding depth image, with a gray scale proportional to $d_{n,v}$.
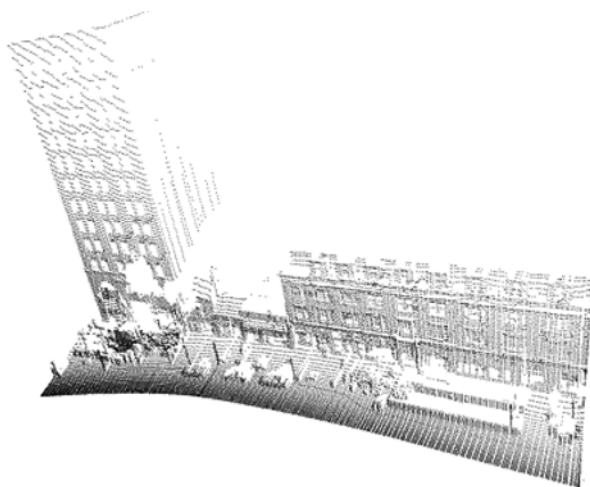
## 4. Properties of City Laser Scans

In this section, we briefly describe properties of scans taken in a city environment, resulting from the physics of a laser scanner as an active device measuring time-of-flight of light rays. It is essential to understand these properties and the resulting imperfections in distance measurement, since at times they lead to scan points that appear to be in contradiction with human eye perception or a camera. As the goal of our modeling approach is to generate a photo realistic model, we are interested in reconstructing what the human eye or a camera would observe while moving around in the city. As such, we discuss the discrepancies between these two different sensing modalities in this section.

### 4.1. Discrepancies Due to Different Resolution

The beam divergence of the laser scanner is about 15 milliradians (mrad) and the spacing, hence the angular resolution, is about 17 mrad. As such, this is much lower than the resolution of the camera image with about 2.1 mrad in the center and 1.4 mrad at the image borders. Therefore, small or thin objects, such as cables, fences, street signs, light posts and tree branches, are clearly visible in the camera image, but only partially captured in the scan. Hence they appear as "floating" vertices, as seen in the depth image in Fig. 9.

### 4.2. Discrepancies Due to the Measurement Physics

Camera and eye are passive sensors, capturing light from an external source; this is in contrast with a laser scanner, which is an active sensor, and uses light that it emits itself. This results in substantial differences in measurement of reflecting and semitransparent surfaces, which are in form of windows and glass fronts frequently present in urban environments. Typically, there is at least 4% of the light reflected at a single glass/air transition, so a total of at least 8% per window; if the window has a reflective coating, this can be larger. The camera typically sees a reflection of the sky or a nearby building on the window, often distorted or merged with objects behind the glass. Although most image processing algorithms would fail in this situation, the human brain is quite capable of identifying windows. In contrast, depending on the window reflectance, the laser beam is either entirely reflected, most times in a different direction from the laser itself,

Data Processing Algorithms for Generating Textured 3D Building Facade Meshes    165



**(a)**



**(b)**

*Figure 8.*    Scan grid representations: (a) 3D vertices; (b) depth image.



*Figure 9.*    "Floating" vertices.



*Figure 10.*    Laser measurement in case of a glass window.

**398** resulting in no distance value, or is transmitted through
**399** the glass. In the latter case, if it hits a surface as shown
**400** in Fig. 10, the backscattered light travels again through
**401** the glass. The resulting surface reflections on the glass
**402** only weaken the laser beam intensity, eventually below

the detection limit, but do not otherwise necessarily af- **403**
fect the distance measurement. To the laser, the window **404**
is quasi non-existent, and the measurement point is gen- **405**
erally not on the window surface, unless the surface is **406**
orthogonal to the beam. In case of multi-reflections, the **407**

166    *Frueh, Jain and Zakhor*

408  situation becomes even worse as the measured distance
409  is almost random.

*4.3. Discrepancies Due to Different Scan
410
411      and Viewpoints*

412  Laser and camera are both limited in that they can only
413  detect the first visible/backscattering object along a
414  measurement direction and as such cannot deal with
415  occlusions. If there is an object in the foreground, such
416  as a tree in front of a building, the laser cannot cap-
417  ture what is behind it; hence, generating a mesh from
418  the obtained scan points results in a hole in the build-
419  ing. We refer to this type of mesh hole as *occlusion*
420  *hole*. As the laser scan points resemble a cylindrical
421  projection, but rendering is parallel or perspective, in
422  presence of occlusions, it is impossible to reconstruct
423  the original view without any holes, even for the view-
424  points from which data was acquired. This is a special
425  property of our fast 2D data acquisition method. An
426  interesting fact is that the wide-angle camera images
427  captured simultaneously with the scans often contain
428  parts of the background invisible to the laser. These
429  could be potentially used either to fill in geometry us-
430  ing stereo techniques, or to verify the validity of the
431  filled in geometry obtained from using interpolation
432  techniques.
433     For a photo realistic model, we need to devise
434  techniques for detecting discrepancies between the
435  two modalities, removing invalid scan points, and
436  filling in holes, either due to occlusion or due to
437  unpredictable surface properties; we will describe
438  our approaches to these problems in the following
439  sections.

440  **5.  Multi-Layer Representation**

441  To ensure that the facade model looks reasonable from
442  every viewpoint, it is necessary to complete the geom-
443  etry for the building facades. Typically, our facades are
444  2 1/2 D objects rather than full 3D objects, and hence
445  we introduce a representation based of multiple depth
446  layers for the street scenery, similar to the one pro-
447  posed in Chang and Zakhor (1999). Each depth layer
448  is a scan grid, and the scan points of the original grid
449  are assigned to exactly one of the layers. If at a certain
450  grid location there is a point in a foreground layer, this
451  location is empty in all layers behind it and needs to be
452  filled in.

453  Even though the concept can be applied to an arbi-
454  trary number of layers, we found that it is in our case
455  sufficient to generate only two, namely a foreground
456  and a background layer. To assign a scan point to ei-
457  ther one of the two layers we make the following as-
458  sumptions about our environment: Main structures, i.e.
459  buildings, are usually (a) vertical, and (b) extend over
460  several feet in horizontal dimension. Furthermore, we
461  assume that (c) building facades are roughly perpen-
462  dicular to the driving direction and that (d) most scan
463  points correspond to facades rather than to foreground
464  objects, as it can occur in residential areas with houses
465  hidden behind trees. Under these conditions, we can ap-
466  ply the following steps to identify foreground objects:
467     For each vertical scan *n* corresponding to a column in
468  the depth image, we define the main depth as the depth
469  value that occurs most frequently, as shown in Fig. 11.
470  The scan vertices corresponding to the main depth lie
471  on a vertical line, and the first assumption suggests that
472  this is a main structure, such as a building, or perhaps
473  other vertical objects, such as a street light or a tree
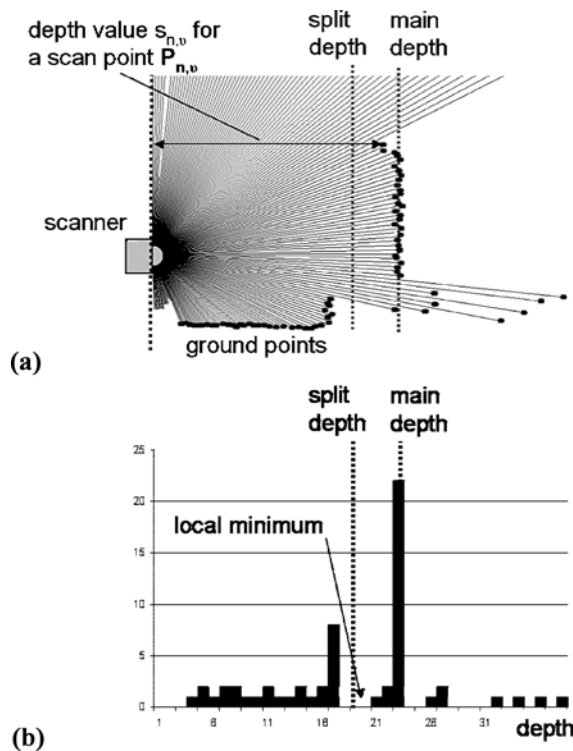474  trunk. With the second assumption, we filter out the



*Figure 11.* Main depth computation for a single scan n: (a) laser scan with rays indicating the laser beams and dots at the end the corresponding scan points; (b) computed depth histogram.

Data Processing Algorithms for Generating Textured 3D Building Facade Meshes    167
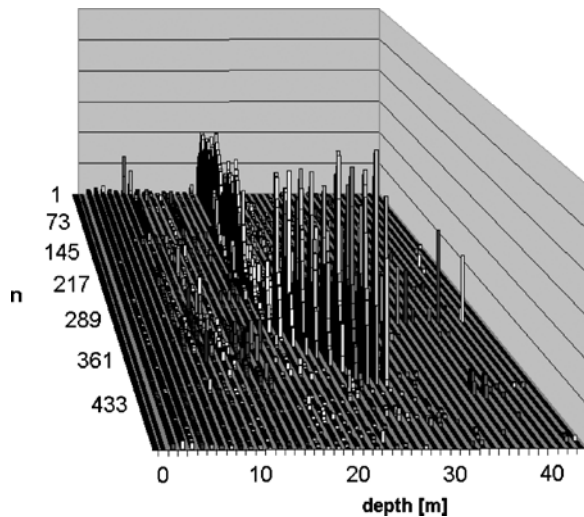


*Figure 12.*    Two-dimensional histogram for all scans.

475 latter class of vertical objects. More specifically, our
476 processing steps can be described as follows:

477     We sort all depth values $s_{n,v}$ for each column $n$ of
478 the depth image into a histogram as shown in Fig. 11(a)
479 and (b), and detect the peak value and its correspond-
480 ing depth. Applying this to all scans results in a 2D
481 histogram as shown in Fig. 12, and an individual main

482 depth value estimate for each scan. Based on the second
483 assumption, isolated outliers are removed by applying
484 a median filter on these main depth values across the
485 scans, and a final depth value is assigned to each col-
486 umn $n$. We define a "split" depth, $\gamma_n$, for each column
487 $n$, and set it to the first local minimum of the histogram
488 occurring immediately before main depth, i.e. with a
489 depth value smaller than the main depth. Taking the first
490 minimum in the distribution instead of the main value
491 itself has the advantage that points clearly belonging
492 to foreground layers are splits off, whereas overhang-
493 ing parts of buildings, for which the depth is slightly
494 smaller than the main depth, are kept in the main layer
495 where they logically belong to, as shown in Fig. 11.

496     A point can be identified as a ground point if its $z$ co-
497 ordinate has a small value and its neighbors in the same
498 scan column have a similarly low $z$ value. We prefer
499 to include the ground in our models, and as such, as-
500 sign ground points also to the background layer. There-
501 fore, we split layers by assigning a scan point $P_{n,v}$ to
502 the background layer, if $s_{n,v} > \gamma_n$ or $P_{n,v}$ is a ground
503 point, and to the foreground layer otherwise. Figure 13
504 shows an example for the resulting foreground and
505 background layers.

506     Since the steps described in this section assume the
presence of vertical buildings, they cannot be expected



*Figure 13.*    (a) Foreground layer; (b) background layer.

**507** to work for segments that are dominated by trees; this
**508** also applies to the processing steps we introduce in
**509** the following sections. As our goal is to reconstruct
**510** buildings, path segments can be left unprocessed and
**511** included "as is" in the city model, if they do not contain
**512** any structure. A characteristic of a tree area is its fractal-
**513** like geometry, resulting in a large variance among ad-
**514** jacent depth values, or even more characteristically,
**515** many significant vector direction changes for the edges
**516** between connected mesh vertices. We define a coeffi-
**517** cient for the fractal nature of a segment by counting
**518** vertices with direction changes greater than a specific
**519** angle, e.g. twenty degrees, and dividing them by the
**520** total number of vertices. If this coefficient is large, the
**521** segment is most likely a tree area and should not be
**522** made subject to the processing steps described in this
**523** section. This is for example the case for the segment
**524** shown in Fig. 9.

**525** After splitting layers, all grid locations occupied in
**526** the foreground layer are missing in the background
**527** layer as the vertical laser does not capture any oc-
**528** cluded geometry; in the next section we will describe
**529** an approach for filling these missing grid locations
**530** based on neighboring pixels. However, in our data ac-
**531** quisition system there are 3D vertices available from
**532** other sources, such as stereo vision and the horizon-
**533** tal scanner used for navigation. Thus, it is conceiv-
**534** able to use this additional information to fill some
**535** in the depth layers. Our approach to doing so is as
**536** follows:

**537** Given a set of 3D vertices $V_i$ obtained from a dif-
**538** ferent modality, determine the closest scan direction
**539** for each vertex and hence the grid location $(n, \upsilon)$ it
**540** should be assigned to. As shown in Fig. 14, each $V_i$
**541** is assigned to the vertical scanning plane, $S_n$, with the
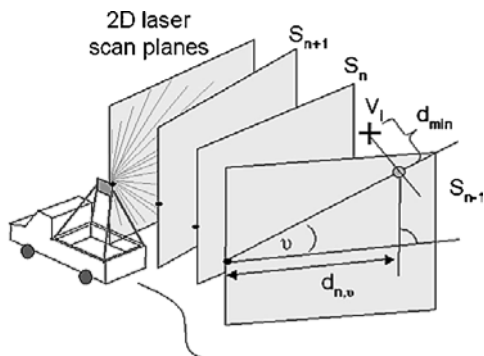**542** smallest Euclidean distance, corresponding to column



*Figure 14.*    Sorting additional points into the layers.



*Figure 15.*    Background layer after sorting in additional points from other modalities.

$n$ in the depth image. Using simple trigonometry, the **543**
scanning angle under which this vertex appears in the **544**
scanning plane, and hence the depth image row $\upsilon$, can **545**
be computed, as well as the depth $d_{n,\upsilon}$ of the pixel. **546**

We can now use these additional vertices to fill in **547**
the holes. To begin with, all vertices that do not belong **548**
to background holes are discarded. If there is exactly **549**
one vertex falling onto a grid location, its depth is di- **550**
rectly assigned to that grid location; for situations with **551**
multiple vertices, median depth value for this location **552**
is chosen. Figure 15 shows the background layer from **553**
Fig. 13(b) after sorting in 3D vertices from stereo vi- **554**
sion and horizontal laser scans. As seen, some holes **555**
can be entirely filled in, and the size of others becomes **556**
smaller, e.g. the holes due to trees in the tall building on **557**
the left side. Note that this intermediate step is optional **558**
and depends on the availability of additional 3D data. **559**

## 6.    Background Layer Postprocessing **560**
and Mesh Generation **561**

In this section, we will describe a strategy to remove **562**
erroneous scan points, and to fill in holes in the back- **563**
ground layer. There exists a variety of successful hole **564**
filling approaches, for example based on fusing mul- **565**
tiple scans taken from different positions (Curless and **566**
Levoy, 1996; Stamos and Allen, 2002). Most previ- **567**
ous work on hole filling in the literature has been fo- **568**
cused on reverse engineering applications, in which a **569**
3D model of an object is obtained from multiple laser **570**
scans taken from different locations and orientations. **571**
Since these existing hole filling approaches are not ap- **572**
plicable to our experimental setup, our approach is to **573**
estimate the actual geometry based on the surrounding **574**
environment and reasonable heuristics. One cannot ex- **575**
pect this estimate to be accurate in *all* possible cases, **576**
rather to lead to an acceptable result in *most* cases, thus **577**

578 reducing the amount of further manual interventions
579 and postprocessing drastically. Additionally, the esti-
580 mated geometry could be made subject to further veri-
581 fication steps, such as consistency checks by applying
582 stereo vision techniques to the intensity images cap-
583 tured by the camera.
584    Our data typically exhibits the following character-
585 istics:

586  • Occlusion holes, such as those caused by a tree,
587    are large and can extend over substantial parts of a
588    building.
589  • A significant number of scan points surrounding a
590    hole may be erroneous due to glass surfaces.
591  • In general, a spline surface filling is unsuitable, as
592    building structures are usually piecewise planar with
593    sharp discontinuities.
594  • The size of data set resulting from a city scan is huge,
595    and therefore the processing time per hole should be
596    kept to a minimum.

597    Based on the above observations, we propose the
598 following steps for data completion.

### 6.1. Detecting and Removing Erroneous Scan Points in the Background layer

601 We assume that erroneous scan points are due to
602 glass surfaces, i.e. the laser measured either an in-
603 ternal wall/object, or a completely random distance
604 due to multi-reflections. Either way, the depth of the
605 scan points measured through the glass is substantially
606 greater than the depth of the building wall, and hence
607 these points are candidates for removal. Since glass
608 windows are usually framed by the wall, we remove
609 the candidate points only if they are embedded among
610 a number of scan points at main depth. An example
611 of the effect of this step can be seen by comparing the
612 windows of the original image in Fig. 16(a) with the
613 processed background layer in Fig. 16(b).

### 6.2. Segmenting the Occluding Foreground Layer into Objects

616 In order to determine holes in the background layer
617 caused by occlusion, we segment the occluding fore-
618 ground layer into objects and project segmentation onto
619 the background layer. This way, holes can be filled in
620 one "object" at a time, rather than all at the same time;
621 this approach has the advantage that more localized



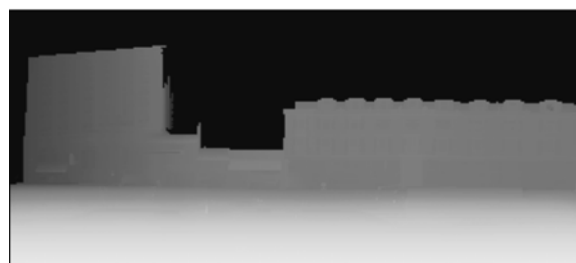*Figure 16.* Processing steps of depth image. (a) Initial depth im-
age. (b) Background layer after removing invalid scan points. (c)
Foreground layer segmented. (d) Occlusion holes filled. (e) Final
background layer after filling remaining holes.

170    *Frueh, Jain and Zakhor*

hole filling algorithms are more likely to result in visually pleasing models than global ones. We segment the foreground layer by taking a random seed point that does not yet belong to a region, and applying a region growing algorithm that iteratively adds neighboring pixels if their depth discontinuity or their local curvature is small enough. This is repeated until all pixels are assigned to a region, and the result is a region map as shown in Fig. 16(c). For each foreground region, we determine boundary points on the background layer; these are all the valid pixels in the background layer that are close to hole pixels caused by the occluding object.

### 6.3. Filling Occlusion Holes in the Background Layer for Each Region

As the foreground objects are located in front of main structures and in most cases stand on the ground, they occlude not only parts of a building, but also parts of the ground. Specifically, an occlusion hole caused by a low object, such as a car, with a large distance to the main structure behind it, is typically located only in the ground and not in the main structure. This is because the laser scanner is mounted on top of a rack, and as such has a top down view of the car. As a plane is a good approximation to the ground, we fill in the ground section of an occlusion hole by the ground plane. Therefore, for each depth image column, i.e. each scan, we compute the intersection point between the line through the main depth scan points and the line through ground scan points. The angle $v'_n$ at which this point appears in the scan marks the virtual boundary between ground part and structure part of the scan; we fill in structure points above and ground points below this boundary differently.

Applying a RANSAC algorithm, we find the plane with the maximum consensus, i.e. maximum number of ground boundary points on it, as the optimal ground plane for that local neighborhood. Each hole pixel with $v < v'_n$ is then filled in with a depth value according to this plane. It is possible to apply the same technique for the structure hole pixels, i.e. the pixels with $v > v'_n$, by finding the optimal plane through the structure boundary points and filling in the hole pixels accordingly. However, we have found that in contrast to the ground, surrounding building pixels do not often lie on a plane. Instead, there are discontinuities due to occluded boundaries and building features such as marquees or lintels, in most cases extending horizontally

across the building. Therefore, rather than filling holes with a plane, we fill in structure holes line by line horizontally, in such a way that the depth value at each pixel is the linear interpolation between the closest right and left structure boundary point, if they both exist; otherwise no value is filled in. In a second phase, a similar interpolation is done vertically, using the already filled in points as valid boundary points. This method is not only simple and therefore computationally efficient, it also takes into account the surrounding horizontal features of the building in the interpolation. The resulting background layer is shown in Fig. 16(d).

### 6.4. Postprocessing the Background Layer

The resulting depth image and the corresponding 3D vertices can be improved by removing scan points that remain isolated, and by filling small holes surrounded by geometry using linear interpolation between neighboring depth pixels. The final background layer after applying all processing steps is shown in Fig. 16(e).

In order to create a mesh, each depth pixel can be transformed back into a 3D vertex, and each vertex $P_{n,v}$ is connected to a depth image neighbor $P_{n+\Delta n, v+\Delta v}$ if

$$|s_{n+\Delta n, v+\Delta v} - s_{n,v}| < s_{max} \quad \text{or if}$$
$$\cos \varphi > \cos \varphi_{max}$$

with

$$\cos \varphi = \frac{(\vec{P}_{n-\Delta n, v-\Delta v} - \vec{P}_{n,v}) \cdot (\vec{P}_{n,v} - \vec{P}_{n+\Delta n, v+\Delta v})}{|\vec{P}_{n-\Delta n, v-\Delta v} - \vec{P}_{n,v}| \cdot |\vec{P}_{n,v} - \vec{P}_{n+\Delta n, v+\Delta v}|}$$

Intuitively, neighbors are connected if their depth difference does not exceed a threshold $s_{max}$ or the local angle between neighboring points is smaller than threshold angle $\varphi_{max}$. The second criteria is intended to connect neighboring points that are on a line, even if their depth difference exceeds $s_{max}$. The resulting quadrilateral mesh is split into triangles, and mesh simplification tools such as Qslim (Garland and Heckbert, 1997) can be applied to reduce the number of triangles.

## 7.    Atlas Generation for Texture Mapping

As photorealism cannot be achieved by using geometry alone, we need to enhance our model with texture data. To achieve this, we equip our data acquisition system with a digital color camera with a wide-angle lens. The
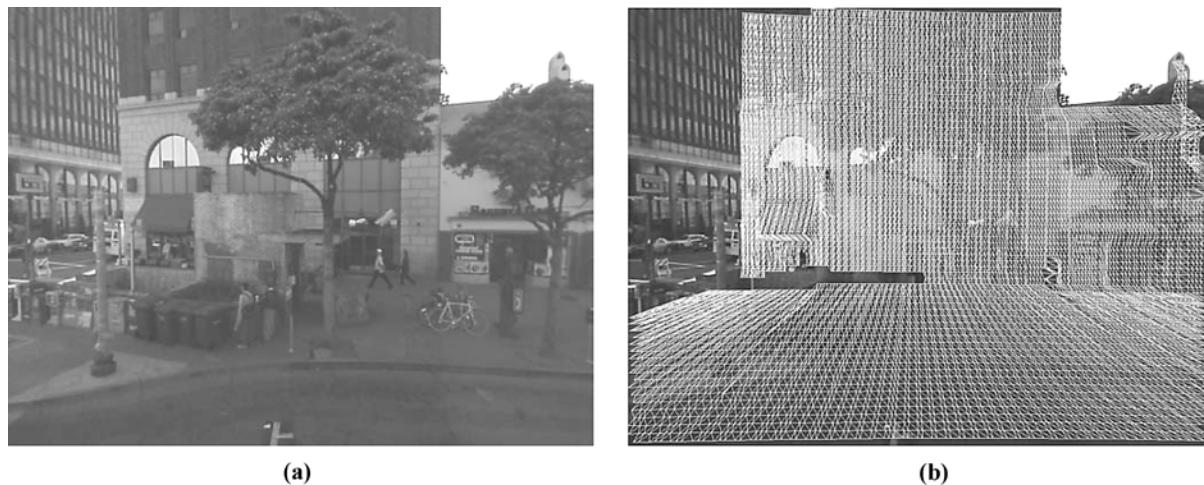
Data Processing Algorithms for Generating Textured 3D Building Facade Meshes     171



*Figure 17.*    Background mesh triangles projected onto camera images. (a) Camera image. (b) Hole filled background mesh projected onto the image and shown as white triangles; occluded background triangles project onto foreground objects. The texture of foreground objects such as the trees should not be used for texturing background triangles corresponding to the building facade.

708 camera is synchronized with the two laser scanners,
709 and is calibrated against the laser scanners' coordinate
710 system; hence, the camera position can be computed for
711 all images. After calibrating the camera and removing
712 lens distortion in the images, each 3D vertex can be
713 mapped to its corresponding pixel in an intensity image
714 by a simple projective transformation. As the 3D mesh
715 triangles are small compared to their distance to the
716 camera, perspective distortions within a triangle can
717 be neglected, and each mesh triangle can be mapped
718 to a triangle in the picture by applying the projective
719 transformation to its vertices.

720     As described in Section 4, camera and laser scanners
721 have different viewpoints during data acquisition, and
722 in most camera pictures, at least some mesh triangles
723 of the background layer are occluded by foreground
724 objects; this is particularly true for triangles that con-
725 sist of filled-in points. An example of this is shown in
726 Fig. 17 where occluded background triangles project
727 onto foreground objects such as the tree. The back-
728 ground triangles are marked in white in Fig. 17. Al-
729 though the pixel location of the projected background
730 triangles is correct, some of the corresponding texture
731 triangles merely correspond to the foreground objects,
732 and thus should not be used for texture mapping the
733 background triangles.

734     In this section, we address the problem of segment-
735 ing out the foreground regions in the images so that their
736 texture is not used for the background mesh triangles.
737 After segmentation, multiple images are combined into

738 a single texture atlas; we then propose a number of tech-
739 niques to fill in the texture holes in the atlas resulting
740 from foreground occlusion. The resulting hole filled at-
741 las is finally used for texture mapping the background
742 mesh.

### 7.1.  Foreground/Background Segmentation in the Images

745 A simple way of segmenting out the foreground objects
746 is to project the foreground mesh onto the camera im-
747 ages and mark out the projected triangles and vertices.
748 While this process works adequately in most cases, it
749 could miss out some parts of the foreground objects
750 such as those shown in Fig. 18, where projected fore-
751 ground geometry is marked in white. As seen in the
752 figure, some small portions of the foreground tree are
753 incorrectly considered as background. This is due to
754 following reasons:

755 1. The foreground scan points are not dense enough
756    for segmenting the image with pixel accuracy, es-
757    pecially at the boundaries of foreground objects.
758 2. The camera captures side views of foreground
759    objects whereas the laser scanner captures a di-
760    rect view, as illustrated in Fig. 19. Hence, some
761    foreground geometry does not appear in the
762    laser scans and as such cannot be marked as
763    foreground.
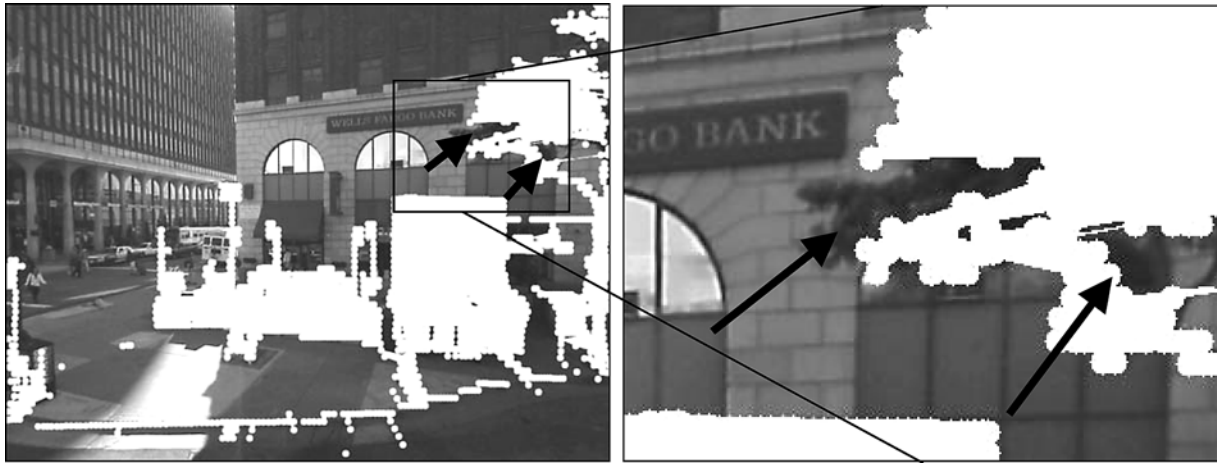
172    *Frueh, Jain and Zakhor*



*Figure 18.*    Identifying foreground in images by projection of the foreground mesh. White denotes the projected foreground and thus image areas not to be used for texture mapping of facades.
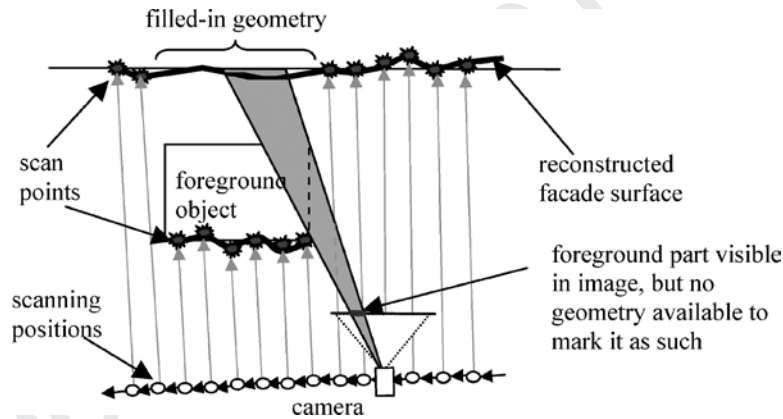


*Figure 19.*    Some foreground objects at oblique viewing angle are not entirely marked in camera images.

To overcome this problem, we have developed a second, more sophisticated method for pixel-accurate foreground segmentation based on the use of correspondence error. The overview of our approach is as follows:

After splitting the scan points into the foreground and background layers, the foreground scan points are projected onto the images. A flood-filling algorithm is applied to all the pixels within a window centered at each of the projected foreground pixels using cues of color constancy and correspondence error. The color at every pixel in the window is compared to that of the center pixel. If the colors are in agreement, and the correspondence error value at the test pixel is close or higher than the value at the center pixel, the test pixel is assigned to the foreground.

In what follows we describe the notion of correspondence error in more detail. Let $I = \{I_1, I_2, \ldots, I_n\}$ denote the set of camera images available for a quasi-linear path segment. Consider two consecutive images $I_{c-1}$ and $I_c$. Consider a 3D point $\mathbf{x}$ belonging to the background mesh obtained after geometry hole filling described in Section 7. $\mathbf{x}$ is projected to the images $I_{c-1}$ and $I_c$ using the available camera position. Assuming that the projected point is within the clip region of both images, let its coordinates in $I_{c-1}$ and $I_c$ be denoted by $u_{c-1}$ and $u_c$ respectively. If $\mathbf{x}$ is not occluded by any foreground object in an image, then its pixel coordinates in the image belong to the background and represent $\mathbf{x}$; otherwise its pixel coordinates correspond to the occluding foreground object. This leads to three cases described below, and illustrated in Fig. 20:

Data Processing Algorithms for Generating Textured 3D Building Facade Meshes    173
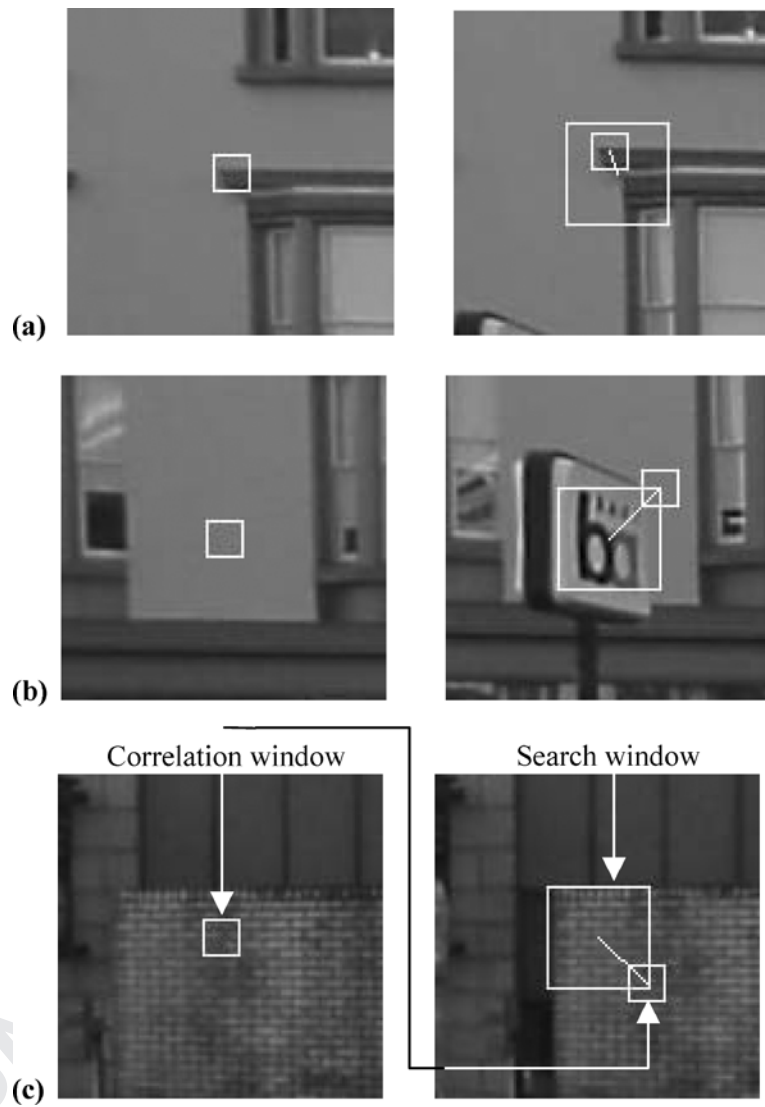


*Figure 20.* Illustration of correspondence error. (a) background scan point is unoccluded in both images. (b) background scan point occluded in one of the images. (c) background scan point occluded in both images. The search window and correlation window are marked for clarity. The line represents the correspondence error vector. The correlation window slides in the search window in order to find the best matching window.

796  1. **x** is occluded in neither images as shown in
797   Fig. 20(a); $u_{c-1}$, and $u_c$ both belong to the back-
798   ground. If the camera position is known precisely,
799   $u_c$ would be the correspondence point for $u_{c-1}$. In
800   practice, the camera position is known only approx-
801   imately, and taking $u_{c-1}$ as a reference, its corre-
802   spondence point in $I_c$ can be located close to $u_c$.
803  2. **x** is occluded only in one of the images as shown in
804   Fig. 20(b); one of $u_{c-1}$ or $u_c$ belongs to a foreground
805   object due to occlusion of point **x**, and the other
806   belongs to the background.

807  3. Point **x** is occluded in both images as shown in
808   Fig. 20(c), and both $u_{c-1}$ and $u_c$ belong to fore-
809   ground objects.

810  In all three cases the best matching pixel to $u_{c-1}$
811  in $I_c$, denoted by $u_{c-1,c}$, is found by searching in a
812  window centered around $u_c$, and performing color cor-
813  relation as illustrated in Fig. 20. The length of vec-
814  tor $\mathbf{v}(u_c, u_{c-1,c})$ then denotes the correspondence error
815  between $u_{c-1}$ and $u_c$. If $|\mathbf{v}(u_c, u_{c-1,c})|$ is large, one
 or both of $u_{c-1}$ and $u_c$ belong to a foreground object

174    *Frueh, Jain and Zakhor*

**816** resulting in cases 2 or 3. In the next step when im-
**817** ages $I_c$ and $I_{c+1}$ are considered, $\mathbf{v}(u_{c+1}, u_{c,c+1})$ is com-
**818** puted and we define the correspondence error at pixel
**819** $u_c$ as:

$$\varepsilon(u_c) = \max(|\mathbf{v}(u_c, u_{c-1,c})|, |\mathbf{v}(u_{c+1}, u_{c,c+1})|)$$

**820** Intuitively, if the correspondence error at a pixel is large
**821** the pixel likely belongs to a foreground object. The
**822** above equation is used to compute the correspondence
**823** error at all the pixels corresponding to projected back-
**824** ground scan points. To compute the correspondence
**825** error at all other pixels within the window centered at
**826** each of the projected foreground scan points, we apply
**827** nearest neighbor interpolation. Each pixel in the win-
**828** dow is declared to be foreground if (a) its color is in
**829** agreement with the center pixel, and (b) its correspon-
**830** dence error value is close or higher than the value at
**831** the center pixel.

**832** The max operation in the above equation has the ef-
**833** fect of not missing out any foreground pixels. Even
though this approach results in large values of cor-

respondence error at some background pixels corre- **834**
sponding to case 2 above, we choose to adopt it for **835**
following reasons: **836**

1. The flood filling algorithm is applied to projected **837**
   foreground scan points only within a square win- **838**
   dow $w$, the size of which is $61 \times 61$ pixels in our **839**
   case; so if a background pixel has a high value of $\varepsilon$ **840**
   but has no projected foreground scan point within a **841**
   neighborhood equal to size of $w$, it is never sub- **842**
   jected to flood filling and thus never marked as **843**
   foreground. **844**
2. Marking non-foreground pixels as foreground is **845**
   not as problematic as leaving foreground pixels un- **846**
   marked. This is because the same 3D point is ob- **847**
   served in multiple camera images, and even though **848**
   it may be incorrectly classified as foreground in **849**
   some images, it is likely to be correctly classified as **850**
   background in others. On the other hand incorrect **851**
   assignment of foreground pixels to the background **852**
   and using then for texturing, results in a erroneous **853**
   texture as discussed before.



*Figure 21.* (a), (b), (c) sequence of three camera images $I_{c-1}$, $I_c$, $I_{c+1}$. (d) correspondence error for $I_c$ shown as gray values. White corresponds to low value and black corresponds to high value of $\varepsilon$. Red pixels are pixels where no background scan points projected. $\varepsilon$ is not computed at these pixels. (e) Foreground scan points marked as white pixels. (f) Foreground regions of $I_c$ marked as white, using color constancy and correspondence error. The green triangles are the triangles used for texture mapping/atlas generation from this image.

**854** Figures 21(a)–(c) show a sequence of three cam-
**855** era images, and Fig. 21(d) shows the correspondence
**856** error for the center image shown as gray values; the
**857** gray values have been scaled so that 0 or black corre-
**858** sponds to maximum value of $\varepsilon$, and 255 or white cor-
**859** responds to minimum value of $\varepsilon$. The correspondence
**860** error has been computed for each projected background
**861** scan point. A $7 \times 7$ window is centered at each pro-
**862** jected background scan point, and $\varepsilon$ at all pixels in the
**863** window has been determined using nearest neighbor
**864** interpolation. The red pixels denote those for which
**865** $\varepsilon$ has not been computed or interpolated in the im-
**866** age. The image looks like a roughly segmented fore-
**867** ground and background. Figure 21(e) shows the pro-
**868** jected foreground scan points marked as white pixels.[1]
**869** Figure 21(f) shows the foreground segmentation using
**870** flood-filling with color and correspondence error com-
**871** parisons as explained in this section. The foreground
**872** has been marked in white color. The green triangles
**873** are the triangles used for texture mapping/atlas gener-
**874** ation from this image. As seen, there are some back-
**875** ground pixels that have been incorrectly assigned to the
**876** foreground. This can be attributed to the fact that our
**877** algorithm has been purposely biased to maximize the
**878** size of foreground region in order to avoid erroneously
**879** assigning background pixels to foreground.

**880** *7.2. Texture Atlas Generation*

**881** Since most parts of a camera image correspond to ei-
**882** ther foreground objects, or facade areas visible in other
**883** images at a more direct view, we can reduce the amount
**884** of texture imagery by extracting only the parts actually
**885** used. The vertical laser scanner results in a vertical col-
**886** umn of scan points, and triangulation of the scan points
**887** thus results in a mesh with a row-column structure as
**888** can be seen in Fig. 17(b). The inherent row-column
**889** structure of the triangular mesh permits to assemble a
**890** new artificial image with a corresponding row-column
**891** structure, and reserved spaces for each texture triangle.
**892** This so-called *texture atlas* is created by performing
**893** the following steps: (a) Determining the inter-column
**894** and inter-row spacing for each consecutive column and
**895** row pair in the mesh and using this to reserve space in
**896** the atlas. (b) Warping each texture triangle to fit to the
**897** corresponding reserved space in the atlas and copying
**898** it into the atlas. (c) Setting texture coordinates of the
**899** mesh triangles to the location in the atlas.
**900** Since in this manner the mesh topology of the tri-
**901** angles is preserved and adjacent triangles align auto-

**902** matically due to the warping process, the resulting tex-
**903** ture atlas resembles a mosaic image. While the atlas
**904** image might not visually look precisely proportionate
**905** due to slightly non-uniform spacing between vertical
**906** scans, these distortions are inverted by the graphics
**907** card hardware during the rendering process, and are
**908** thus negligible.
**909** Figures 22(a) and (b) illustrate the atlas generation:
**910** From the acquired stream of images, the utilized texture
**911** triangles are copied into the texture atlas as symbolized
**912** by the arrows. In this illustration, only five original im-
**913** ages are shown; in this example we have actually com-
**914** bined 58 images of $1024 \times 768$ pixels size to create
**915** a texture atlas of $3180 \times 540$ pixels. Thus, the texture
**916** size is reduced from 45.6 million pixels to 1.7 mil-
**917** lion pixels, while the resolution remains the same. If
**918** occluding foreground objects and building facade are
**919** too close, some facade triangles might not be visible
**920** in any of the captured imagery, and hence cannot be
**921** texture mapped at all. This leaves visually unpleasant
**922** *holes* in the texture atlas, and hence in final rendering
**923** of the 3D models. In the following, we propose ways of
**924** synthesizing plausible artificial texture for these holes.

**925** *7.3. Hole Filling of the Atlas*

**926** Early work relating to disocclusion in images was done
**927** by Nitzberg et al. (1993). Significant improvements
**928** to this were made in Masnou and Morel (1998) and
**929** Ballester et al. (2000, 2001). These methods are capable
**930** of filling in small holes in non-textured regions and
**931** essentially deal with *local Inpainting*; they thus cannot
**932** be used for filling in large holes or holes in textured
**933** regions (Chan and Shen, 2001). We propose a simple
**934** and efficient method of hole filling that first completes
**935** regions of low spatial frequency by interpolating the
**936** values of surrounding pixels, and then uses a copy-paste
**937** method to synthesize artificial texture for the holes.
**938** In what follows, we explain the above steps in more
**939** detail.

**940** ***Horizontal and Vertical Interpolation.*** Our pro-
**941** posed algorithm first fills in holes in regions of low
**942** variance using linear interpolation of surrounding pixel
**943** values. A generalized two-dimensional (2D) linear in-
**944** terpolation is not advantageous over a one-dimensional
**945** (1D) interpolation in a man-made environment where
**946** features are usually either horizontal or vertical e.g.
**947** curbs run across the streets horizontally, edges of fa-
**948** cades are vertical, banners on buildings are horizontal.

176     *Frueh, Jain and Zakhor*



*Figure 22.*   (a) Images obtained after foreground segmentation are combined to create a texture atlas. In this illustration only five images are shown, whereas in this particular example 58 images were combined to create the texture atlas. (b) Atlas with texture holes for the facade portions that were not visible in any image. (c) Artificial texture is synthesized in the texture holes to result in a filled in atlas that is finally used for texturing the background mesh.

949 One-dimensional interpolation is simple, and is able to
950 recover most sharp discontinuities and gradients. We
951 perform 1D horizontal interpolation in the following
952 way: for each row, pairs of pixels between which RGB
953 information is missing are detected. The missing values
954 are filled in by a linear interpolation of the boundary
955 pixels if (a) the boundary pixels at the two ends have
956 similar values, and (b) the variances around the bound-
957 aries are low at both ends. We follow this by vertical
958 interpolation in which for each column the missing val-
959 ues are interpolated vertically.
960     Figure 23(a) shows part of a texture atlas with holes
961 marked in red. Figure 23(b) shows the image after a
962 pass of 1D horizontal interpolation. As seen, horizontal

963 edges such as the blue curb are completed. Figure 23(c)
964 shows the image after horizontal and vertical interpo-
965 lation. We find the interpolation process to be simple,
966 fast, and to complete the low frequency regions well.

967 *The Copy-Paste Method.*    Assuming that building fa-
968 cades are highly repetitive, we fill holes that could not
969 be filled by horizontal and vertical interpolation, by
970 copying and pasting blocks from other parts of the im-
971 age. This approach is similar to the one proposed in
972 Efros and Freeman (2001) where a large image is cre-
973 ated with a texture similar to a given template. In our
974 copy-paste method the image is scanned pixel by pixel
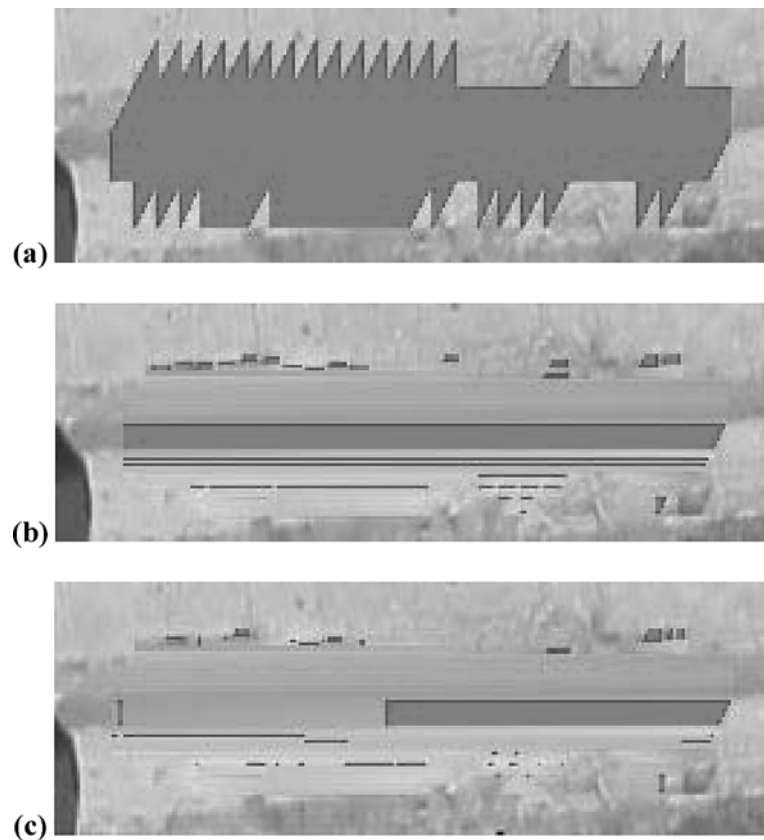975 in raster scan order, and pixels at the boundary of holes

*Figure 23*.  (a) part of a texture atlas with holes marked in red (b) after horizontal interpolation (c) after horizontal and vertical interpolation.

976  are stored in an array to be processed. A square win-
977  dow $w$ of size $(2M + 1) \times (2M + 1)$ pixels is centered
978  at a hole pixel $p$, and the atlas is searched for a win-
979  dow denoted by *bestmatch*($w$) which (a) has the same
980  size as $w$, (b) does not contain more than 10% hole
981  pixels, and (c) matches best with $w$. If the difference
982  between $w$ and *bestmatch*($w$) is below a threshold, the
983  *bestmatch* is classified as a good match to $w$ and hole
984  pixels of $w$ are replaced with corresponding pixels in
985  *bestmatch*($w$). The method is illustrated in Fig. 24.
986       For the method to work well, we need a suitable met-
987  ric that accurately measures the perceptual difference
988  between two windows, an efficient search process that
989  finds the *bestmatch* of a window $w$, a decision rule that
990  classifies whether the *bestmatch* found is good enough,
991  and a strategy to deal with cases when the *bestmatch*
992  of a window $w$ is not a good match. In our proposed
993  scheme, the difference between two windows consists
     of two components: (a) the sum of color differences

994  of corresponding pixels in the two windows, and (b)
995  the number of outliers for the pair of windows. These
996  components are weighted appropriately to compute the
997  resulting difference. An efficient search is performed
998  by constructing a hierarchy of Gaussian pyramids, and
999  performing an exhaustive search at a coarse level to
1000  find a few good matches, which are then successively
1001  refined at finer levels of the hierarchy. In cases when
1002  no good match is found the window size is changed
1003  adaptively. If a window of size $(2M + 1) \times (2M + 1)$
1004  does not result in a good match, the algorithm finds
1005  the *bestmatch* for a smaller window of size $(M + 1) \times$
1006  $(M + 1)$ and this process continues until the window
1007  size becomes too small, in our case $9 \times 9$ pixels. If no
1008  good match is found even after reducing the window
1009  size, the hole pixels are filled by averaging the known
1010  neighbors provided the pixel variance of the neighbors
1011  is low; otherwise the colors of hole pixels are set to the
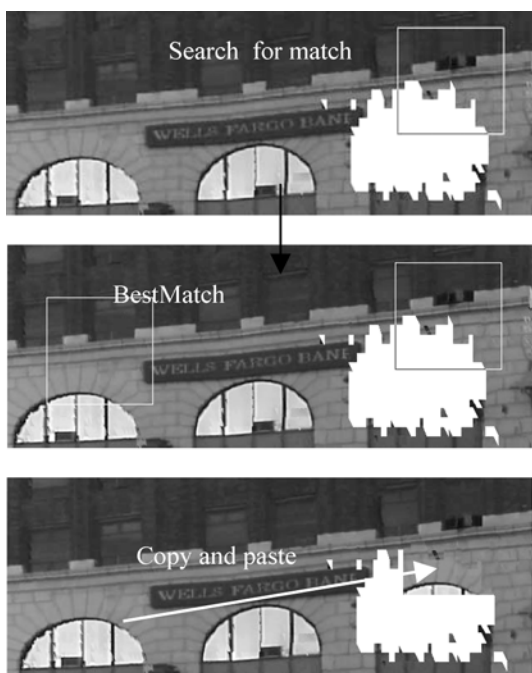1012  value of randomly chosen neighbors.

178  *Frueh, Jain and Zakhor*



*Figure 24*. Illustrating the copy-paste method.

**1013** **8. Results**

**1014** We drove our equipped truck on a 6769 meters
**1015** long path in downtown Berkeley, starting from Blake
**1016** street through Telegraph avenue, and in loops around
**1017** the downtown blocks. During this 24-minute-drive,
**1018** we captured 107,082 vertical scans, consisting of
**1019** 14,973,064 scan points. For 11 minutes of driving time
**1020** in the downtown area, we also recorded a total of 7,200
**1021** camera images. Applying the described path splitting
**1022** techniques, we divide the driven path into 73 segments,
**1023** as shown in Fig. 25 overlaid with a road map. There is
**1024** no need for further manual subdivision, even at Shat-
**1025** tuck Avenue, where Berkeley's street grid structure is
**1026** not preserved.

**1027** *8.1. Geometry Reconstruction*

**1028** For each of the 73 segments, we generate two meshes
**1029** for comparison: the first mesh is obtained directly from
**1030** the raw scans, and the second one from the depth im-
**1031** age to which we have applied the postprocessing steps
**1032** described in previous sections. For 12 out of the 73
**1033** segments, additional 3D vertices derived from stereo
**1034** vision techniques are available, and hence, sorting in



*Figure 25*. Entire path after split in quasi-linear segments.

these 3D points into the layers based on Section 5 **1035**
does fill some of the holes. For these specific holes, **1036**
we have compared the results based on stereo vision **1037**
vertices with those based on interpolation alone as de- **1038**
scribed in Section 6, and have found no substantial dif- **1039**
ference; often the interpolated mesh vertices appear to **1040**
be more visually appealing, as they are less noisy than **1041**
the stereo vision based vertices. Figure 26(a) shows **1042**
an example before processing, and Fig. 26(b) shows **1043**
the tree holes completely filled in by stereo vision ver- **1044**
tices. As seen, the outline of the original holes can **1045**
still be recognized in Fig. 26(b), whereas the points **1046**
generated by interpolation alone are almost indistin- **1047**
guishable from the surrounding geometry, as seen in **1048**
Fig. 26(c). **1049**

We have found our approach to work well in the **1050**
downtown areas, where there are clear building struc- **1051**
tures and few trees. However, in residential areas, **1052**
where the buildings are often almost completely hid- **1053**
den behind trees, it is difficult to accurately estimate **1054**
the geometry. As we do not have the ground truth **1055**
to compare with, and as our main concern is the vi- **1056**
sual quality of the generated model, we have manu- **1057**
ally inspected the results and subjectively determined **1058**
the degree to which the proposed postprocessing pro- **1059**
cedures have improved the visual appearance. The **1060**
evaluation results for all 73 segments before and af- **1061**
ter postprocessing techniques described in this paper **1062**
are shown in Table 1; the postprocessing does not uti- **1063**
lize auxiliary 3D vertices from horizontal laser scan- **1064**
ner or the camera. Even though 8% of all processed **1065**
segments appear visually worse than the original, the **1066**
overall quality of the facade models is significantly im- **1067**
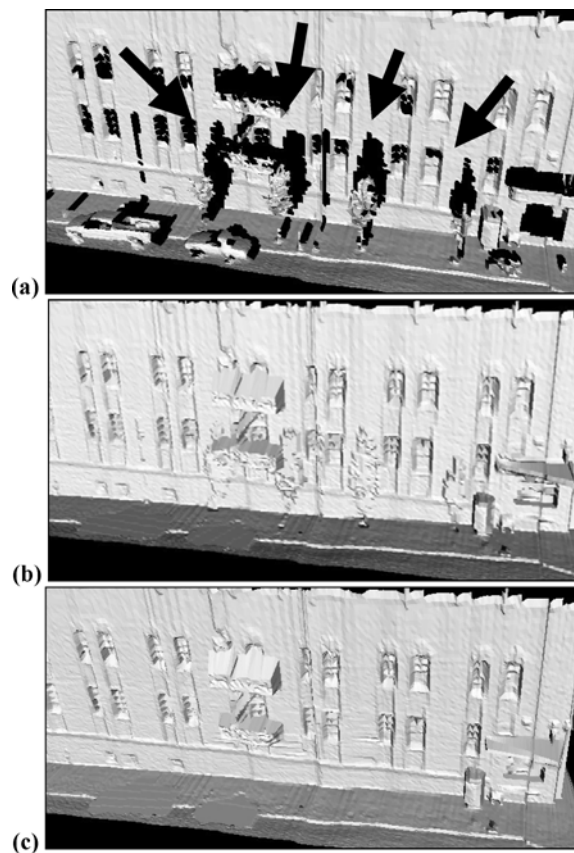proved. The important downtown segments are in most **1068**

*Figure 26.* Hole filling. (a) Original mesh with holes behind occluding trees; (b) filled by sorting in additional 3D points using stereo vision; (c) filled by using the interpolation techniques of Section 6.

**1069** cases ready to use and do not require further manual
**1070** intervention.
**1071**　　The few problematic segments all occur in residen-
**1072** tial areas, consisting mainly of trees. The tree detection
**1073** algorithm described in Section 5 classifies ten segments
**1074** as "critical" in that too many trees are present; all six
**1075** problematic segments corresponding to "worse" and
**1076** "significantly worse" rows in Table 1 are among them,
**1077** yet none of the improved segments in rows 1 and 2 are

*Table 1.* Visual comparison of the processed mesh vs. the original mesh for all 73 segments.

| | | |
|---|---|---|
| Significantly better | 35 | 48% |
| Better | 17 | 23% |
| Same | 15 | 21% |
| Worse | 5 | 7% |
| Significantly worse | 1 | 1% |
| Total | 73 | 100% |

*Table 2.* Visual comparison of the processed mesh vs. the original mesh for the segments automatically classified as non-tree-areas.

| | | |
|---|---|---|
| Significantly better | 35 | 56% |
| Better | 17 | 27% |
| Same | 11 | 17% |
| Worse | 0 | 0% |
| Significantly worse | 0 | 0% |
| Total | 63 | 100% |

detected as critical. This is significant because it shows **1078** that (a) all problematic segments correspond to regions **1079** with a large number of trees, and (b) they can be suc- **1080** cessfully detected and hence not be subjected to the **1081** proposed steps. Table 2 shows the evaluation results if **1082** only non-critical segments are processed. As seen, the **1083** postprocessing steps described in this paper together **1084** with the tree detection algorithm improve over 80% of **1085** the segments, and never result in degradations for any **1086** of the segments. **1087**

　　In Fig. 27 we show before and after examples, and **1088** the corresponding classifications according to Tables 1 **1089** and 2. As seen, except for pair "f", the proposed post- **1090** processing steps result in visually pleasing models. Pair **1091** f in Fig. 27 is classified by our tree detection algorithm **1092** as critical, and hence, should be left "as is" rather than **1093** processed. **1094**

*8.2. Texture Reconstruction* **1095**

For 29 path segments or $3\frac{1}{2}$ city blocks, we recorded **1096** camera images for texture mapping, and hence we re- **1097** construct texture atlases as described in Section 7. Most **1098** facade triangles which were occluded in the direct view **1099** could be texture mapped from some other image with **1100** an oblique view. Only 1.7% of the triangles were not **1101** visible in any image, and therefore required texture **1102** synthesis. **1103**

　　Figure 28 demonstrates our texture synthesis algo- **1104** rithm. Figure 28(a) shows a closer view of the facade to- **1105** gether with holes caused by occlusion from foreground **1106** objects. The holes are marked in white. Figure 28(b) **1107** shows the result using the hole filling technique de- **1108** scribed in Section 7. As seen, the synthesized texture **1109** improves the visual appearance of the model. For com- **1110** parison purposes, Fig. 28 (c) shows the image resulting **1111** from the inpainting algorithm described in Bertalmio **1112** et al. (2000). A local algorithm such as inpainting only **1113** uses the information contained in a thin band around the **1114**
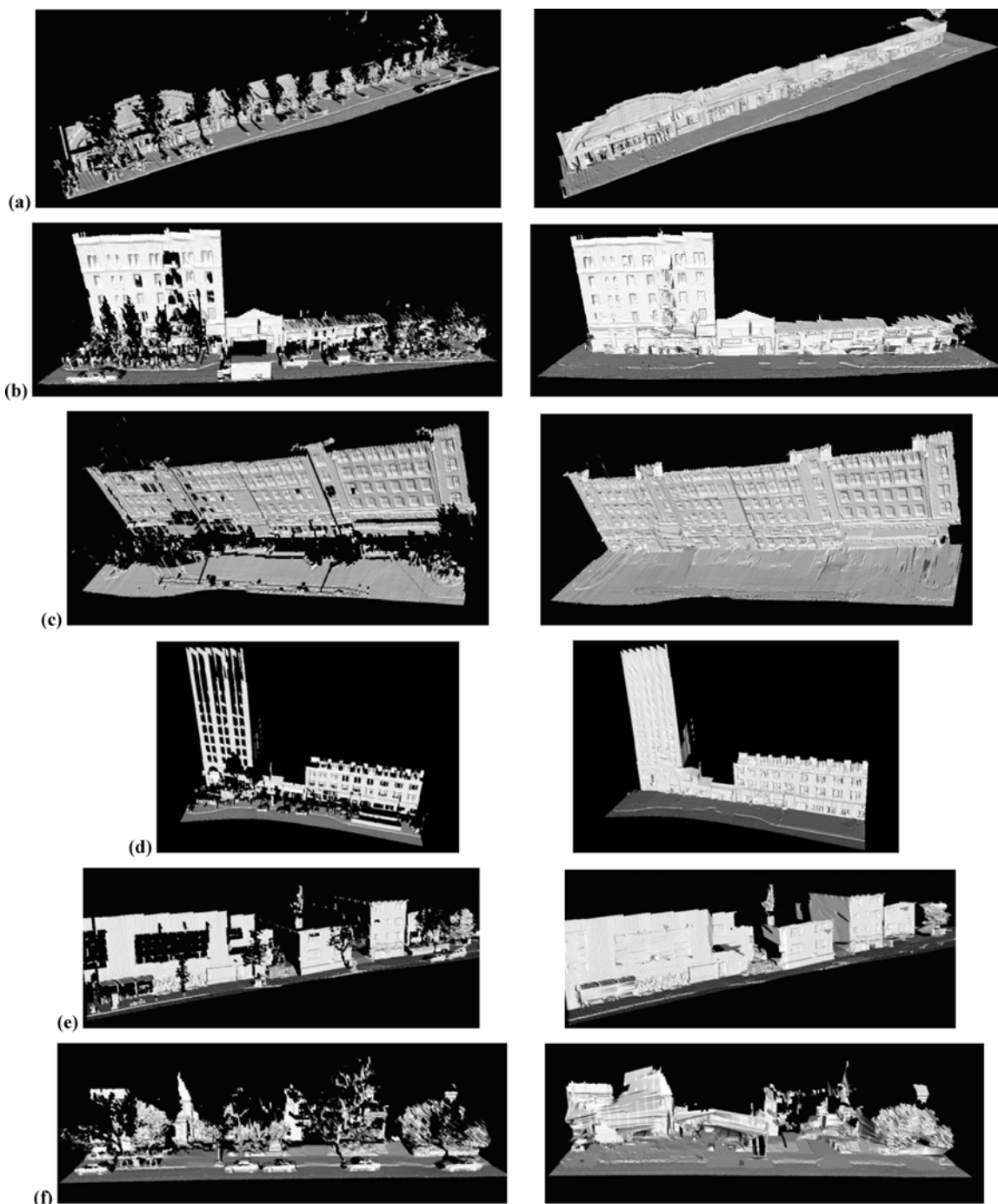
180    *Frueh, Jain and Zakhor*



*Figure 27*.    Generated meshes, left side original, right side after the proposed foreground removal and hole filling procedure. The classification for the visual impression is "significantly better" for the first four image pairs, "better" for pair e and "worse" for pair f.

1115 hole, and hence interpolation of surrounding boundary
1116 values cannot possibly reconstruct the window arch or
1117 the brick pattern on the wall. The copy-paste method on
1118 the other hand, is able to reconstruct the window arch

and brick pattern by copying and pasting from other 1119
parts of the image. 1120
  In Fig. 29 we apply the texture atlas of Fig. 28 to the 1121
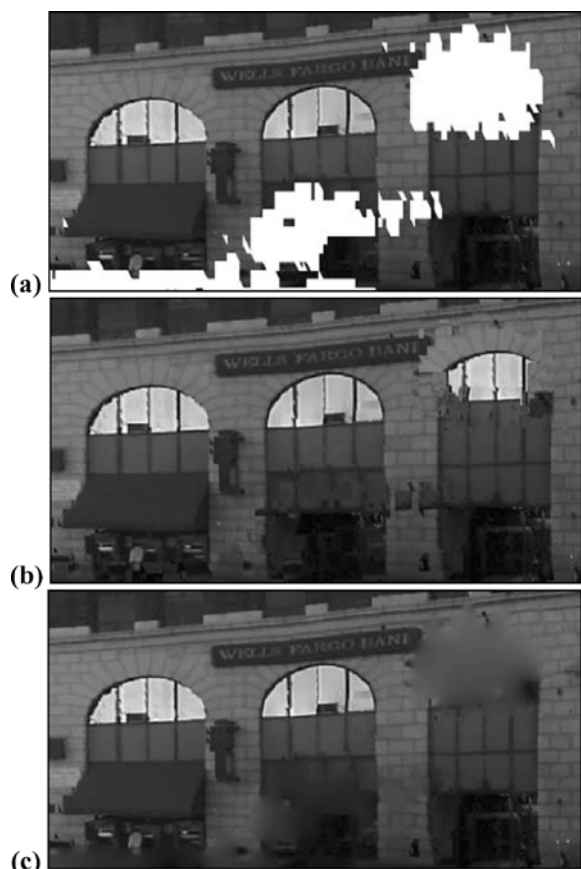geometry shown in Fig. 27(d) and compare the model 1122

Data Processing Algorithms for Generating Textured 3D Building Facade Meshes     181



*Figure 28.* (a) part of texture atlas with holes marked in white; (b) hole filled atlas using the copy-paste method described in Section 7; (c) result of Inpainting.

1123 with and without the data processing algorithms de-
1124 scribed in this paper. Figure 29(a) shows the model
1125 without any processing, Fig. 29(b) the same model af-
ter our proposed geometry processing, and Fig. 29(c)

the model after both geometry processing and texture 1126
synthesis. Note that in the large facade area occluded 1127
by the two trees on the left part of the original mesh, 1128
geometry has been filled in; while most of it could 1129
be texture mapped using oblique camera views, a few 1130
remaining triangles could only be textured via synthe- 1131
sis. As seen, the visual difference between the original 1132
mesh and the processed mesh is striking and appears 1133
to be even larger than in Fig. 27(d). This is because 1134
texture distracts the human eye from missing details 1135
and geometry imperfections introduced by hole filling 1136
algorithms. Finally, Fig. 30 shows the facade model for 1137
the entire $3\frac{1}{2}$ city blocks area. 1138

### 8.3. Complexity and Processing Time     1139

Table 3 shows the processing time measured on a 2 1140
GHz Pentium 4 PC for the automated reconstruction of 1141
the $3\frac{1}{2}$ complete street blocks of downtown Berkeley 1142
shown in Fig. 30. Without the texture synthesis tech- 1143
nique of Section 7, thus leaving 1.7% of the triangles 1144
untextured, the processing time for the model recon- 1145
struction is 2 hours and 17 minutes. Due to the size 1146
of the texture, our texture synthesis algorithm is much 1147
slower, with processing time varying between <1 min 1148
and 8 hours per segment, depending on the number and 1149
the size of the holes. If quality is more important than 1150
processing speed, the entire model can be reconstructed 1151
with texture synthesis in about 23 hours. 1152

Our approach is not only fast, but also automated: 1153
Besides the driving, which took 11 minutes for the 1154
model shown, the only manual step in our modeling 1155
approach is one mouse click needed to enter the ap- 1156
proximate starting position in the digital surface map 1157
for Monte-Carlo Localization, which is needed once

*Table 3.* Processing times for $3\frac{1}{2}$ downtown Berkeley blocks.

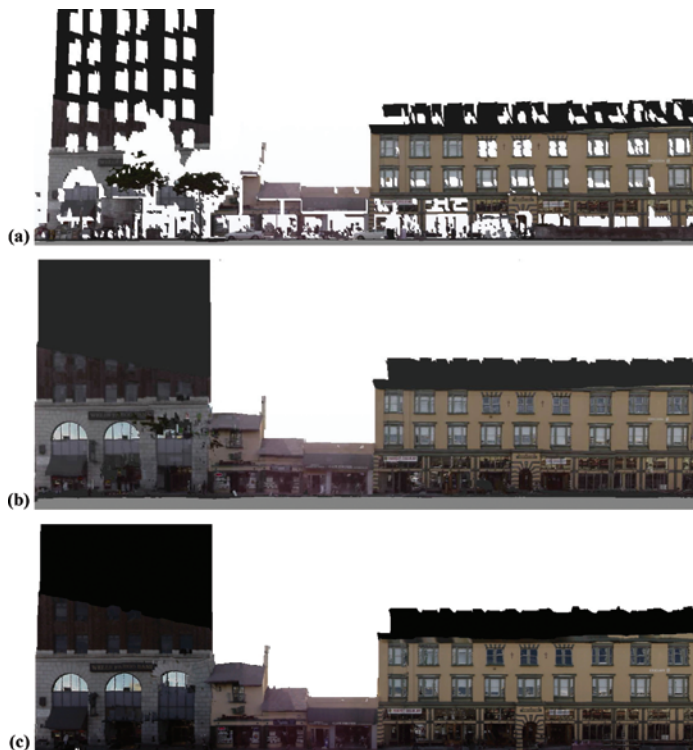| Processing Times for Automated Reconstruction on 2 GHz Pentium 4 | |
| --- | --- |
| Data conversion | 14 min |
| Path reconstruction based on scan matching and global correction with Monte Carlo Localization (with DSM and 5,000 particles) | 70 min |
| Path segmentation | 1 min |
| Geometry reconstruction | 6 min |
| Texture mapping and atlas generation | 27 min |
| Texture synthesis for atlas holes (including pixel-accurate image foreground removal) | 20 h 51 min |
| Model optimization for rendering | 19 min |
| Total model generation time without texture synthesis | 2 h 17 min |
| Total model generation time with texture synthesis | 23 h 08 min |

182      *Frueh, Jain and Zakhor*



*Figure 29.* Textured facade mesh: (a) without any processing; (b) with geometry processing; and (c) with geometry processing, pixel-accurate foreground removal and texture synthesis.
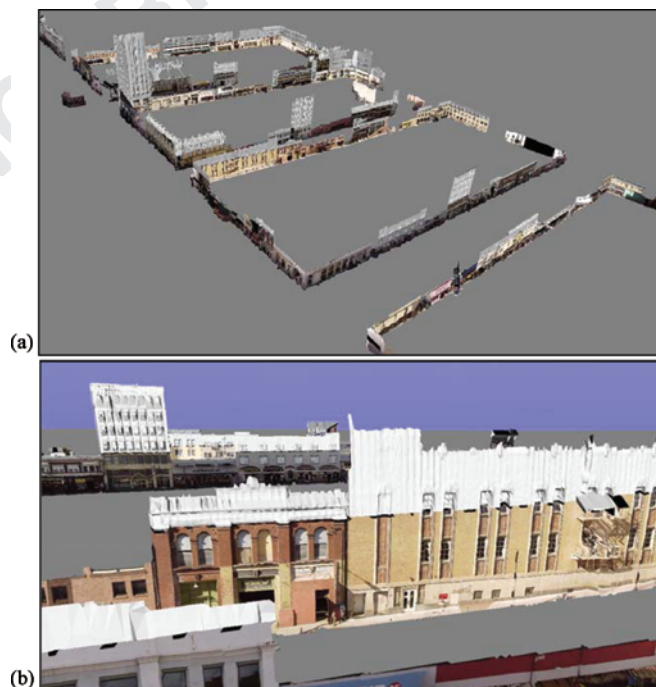


*Figure 30.* Reconstructed facade models: (a) overview; (b) close-up view.

**1158** at the beginning of a model acquisition, and could be
**1159** automated by using a low-cost GPS.

**1160** *8.4.  Accuracy, Limitations, and Possible*
**1161**      *Failure Scenarios*

**1162** We have demonstrated that our approach is capable
**1163** of reconstructing facade models for a large-scale ur-
**1164** ban area. Since we do not have access to ground-truth
**1165** geometry or texture data, it is difficult, if not impossi-
**1166** ble, to assess the accuracy of the reconstructed models.
**1167** However, the following observations can be made:
**1168**     The accuracy of the reconstructed model depends on
**1169** (a) the accuracy of the raw scan points and (b) errors
**1170** made during hole filling and mesh reconstruction. The
**1171** vertical scan points have a basic random error of $\sigma_s =$
**1172** $\pm 3.5$ centimeters due to the scanner's measurement
**1173** noise. As determined in Frueh (2002), the horizontal
**1174** scan matching is accurate to within $\sigma_x = \sigma_y = 1$ cm
**1175** for successive horizontal scans, which are on average
**1176** about 1 meter apart. Thus, the relative position accu-
**1177** racy for a path corresponding to N matched horizon-
**1178** tal scans, or about N meters, is $\sigma_N = \sqrt{N \cdot (\sigma_x^2 + \sigma_y^2)}$.
**1179** Therefore, the total uncertainty between 2 scan points
**1180** p1, p2 recorded within $\underline{N \text{ meters of driving}}$ can be esti-
**1181** mated to $\sigma_{p1,p2} = \sqrt{N \cdot (\sigma_x^2 + \sigma_y^2) + 2\sigma_s^2}$. For exam-
**1182** ple for a 10 meter wide facade, $\sigma_{p1,p2}$ is 6.67 centimeter.
**1183** Additionally, our Monte-Carlo-Localization-based ap-
**1184** proach utilizes a DSM to correct drift-like global pose
**1185** offsets in the vehicle's path by redistributing correction
**1186** vectors among the relative motion estimates. By virtue
**1187** of the parameters chosen in our Monte-Carlo localiza-
**1188** tion, these correction vectors are designed to be of the
**1189** same order of magnitude as $\sigma_x$. While the correction
**1190** vectors are intended to compensate for errors made
**1191** during the horizontal scan matching, they can add to
**1192** the uncertainty due to inaccuracies in the DSM itself.
**1193** Thus, our models are accurate, locally to about $\sigma_{p1,p2}$,
**1194** e.g. few centimeters, and globally to the accuracy of
**1195** the DSM as a global map, e.g. one meter.
**1196**     Errors made during hole filling and mesh reconstruc-
**1197** tion can be severe, depending on the scene and the
**1198** amount of geometry that needs to be "invented". First,
**1199** facades perpendicular to the driving direction or en-
**1200** tirely occluded by large foreground objects are invisi-
**1201** ble to the laser scanner and hence not even result in a
**1202** hole to be filled in—such structures do not appear in
**1203** the model at all. Similarly, facades that are nearly all
**1204** glass without surrounding solid walls would not pro-
**1205** vide enough vertical scan points to be recognized as a

**1206** facade and would therefore not be reconstructed. Sec-
**1207** ond, complicated facade objects such as fences, fire es-
**1208** capes, or wires cannot be adequately reconstructed; due
**1209** to their non-contiguous structure, corresponding scan
**1210** points are classified as outliers and removed. Third, it
**1211** is obvious that even a human operator can be wrong
**1212** in filling a hole, since clues at the boundaries might be
**1213** misleading; this is more so for an automated hole filling
**1214** algorithm such as ours, which is based on interpolation
**1215** and hence implicitly assumes a rather simple geometric
**1216** structure. Forth, and most importantly, there are scenes
**1217** for which a simple foreground/facade layer concept is
**1218** not sufficient. Examples of these are more complex
**1219** staged building structures with porches, pillars, oriels,
**1220** or non-vertical walls, and residential areas with many
**1221** trees. In these cases, our assumptions of Section 5 do
**1222** not hold true any longer; using histogram analysis to
**1223** separate the scan points into either foreground or fa-
**1224** cades is inadequate and results in oddly reconstructed
**1225** models as seen in Fig. 27(f).
**1226**     As a matter of fact, for complicated structures which
**1227** differ substantially from a foreground/background sce-
**1228** nario, our drive-by approach with one single vertical
**1229** scanner does not provide enough data to successfully
**1230** reconstruct a satisfactory model and hence is inappli-
**1231** cable. Fortunately however, as demonstrated for down-
**1232** town Berkeley, the street scenery in most downtown
**1233** areas consists of a foreground/background composi-
**1234** tion. As a solution to more complicated structures, mul-
**1235** tiple vertical laser scanners could be mounted at dif-
**1236** ferent orientations; similar to merging 3D scans taken
**1237** from multiple viewpoints, these oblique scans could be
**1238** used if direct scans are not sufficient.

**9.  Conclusions**    **1239**

**1240** We have proposed a method to reconstruct building fa-
**1241** cade meshes from large laser surface scans and camera
**1242** images, even in presence of occlusion. Future work will
**1243** focus on using color and texture cues to verify filled-
**1244** in geometry. Additionally, foreground objects could be
**1245** classified and replaced by appropriate generics.

184     *Frueh, Jain and Zakhor*

## 1251  Note

1251  **Note**

1252  1. The original image is more than 4 times larger in each dimension.
1253     This image is produced by subsampling the original image in
1254     a special way. Each white pixel corresponding to a foreground
1255     scan point in the original image is retained as a white pixel in the
1256     subsampled image. This gives a false impression that the density
1257     of foreground scan points is very high. On the other hand if the
1258     image is subsampled in the normal fashion, there would almost
      be no white pixels left in the subsampled image.

## 1259  References

Ballester, C., Bertalmio, M., Caselles, V., Sapiro, G., and Verdera, J. 2001. Filling in by joint interpolation of vector fields and gray levels. In *IEEE Transactions on Image Processing*, pp. 1200–1211.

Ballester, C., Caselles, V., Verdera, J., Bertalmio, M., and Sapiro, G. A variational model for filling-in gray level and color images. In *Proc. 8th IEEE Int'l Conference on Computer Vision*, vol. 1, pp. 10–16.

Bertalmio, M., Sapiro, G., Ballester, C., and Caselles, V. 2000. Image inpainting. In *Proc. SIGGRAPH* 2000, pp. 417–424.

Chan, T. and Shen, J. 2001. Mathematical models for local nontexture inpaintings. *SIAM Journal on Applied Mathematics* 62(3):1019–1043.

Chang, N.L. and Zakhor, A. 1999. A multivalued representation for view synthesis. In *Proc. Int'l Conference on Image Processing*, Kobe, Japan, vol. 2, pp. 505–509.

Curless, B. and Levoy, M. A volumetric method for building complex models from range images. In *SIGGRAPH*, New Orleans, pp. 303–312.

Debevec, P.E., Taylor, C.J., and Malik, J. 1996. Modeling and rendering architecture from photographs. In *Proc. ACM SIGGRAPH*.

Dick, A., Torr, P., Ruffle, S., and Cipolla, R. 2001. Combining single view recognition and multiple view stereo for architectural scenes. In *International Conference on Computer Vision*, Vancouver, Canada, pp. 268–274.

Efros, A. and Freeman, W. 2001. Image quilting for texture synthesis and transfer. In *Proc. SIGGRAPH*, pp. 341–346.

Fox, D., Thrun, S., Dellaert, F., and Burgard, W. 2000. Particle filters for mobile robot localization. In *Sequential Monte Carlo Methods in Practice*, A. Doucet, N. de Freitas, and N. Gordon (Eds.), Springer Verlag, New York.

Frere, D., Vandekerckhove, J., Moons, T., and Van Gool, L. 1998. Automatic modelling and 3D reconstruction of urban buildings from aerial imagery. In *IEEE International Geoscience and Remote Sensing Symposium Proceedings*, Seattle, pp. 2593–2596.

Frueh, C. 2002. Automated 3D Model Generation for Urban Environments. Ph.D. Thesis, University of Karlsruhe.

Frueh, C., Flynn, J., Foroosh, H., and Zakhor, A. 2001. Fast 3D model generation in urban environments. In *Workshop on the Convergence of Graphics, Vision, and Video (CGVV'01)*, Berkeley, USA.

Frueh, C. and Zakhor, A. 2001a. Fast 3D model generation in urban environments. In *IEEE Conf. on Multisensor Fusion and Integration for Intelligent Systems*, Baden-Baden, Germany, pp. 165–170.

Frueh, C. and Zakhor, A. 2001b. 3D model generation of cities using aerial photographs and ground level laser scans. In *Computer Vision and Pattern Recognition*, Hawaii, USA, vol. 2.2. pp. II-31-8.

Frueh, C. and Zakhor, A. 2003. Constructing 3D city models by merging ground-based and airborne views. *IEEE Computer Graphics and Applications*, Special Issue Nov/Dec. pp. 52–61.

Garland, M. and Heckbert, P. 1997. Surface *Simplification Using Quadric Error Metrics*. In *SIGGRAPH '97*, Los Angeles, pp. 209–216.

Haala, N. and Brenner, C. 1997. Generation of 3D city models from airborne laser scanning data. In *Proc. EARSEL Workshop on LIDAR Remote Sensing on Land and Sea*, Tallin, Esonia, pp. 105–112.

Kim, Z., Huertas, A., and Nevatia, R. 2001. Automatic description of Buildings with complex rooftops from multiple images. In *Computer Vision and Pattern Recognition*, Kauai, pp. 272–279.

Koch, R., Pollefeys, M., and van Gool, L. 1999. Realistic 3D Scene modeling from uncalibrated image sequences. In *ICIP'99*, Kobe, Japan, Oct., pp. II: 500–504.

Maas, H.-G. 2001. The suitability of airborne laser scanner data for automatic 3D object reconstruction. 3. In *Int'l Workshop on Automatic Extraction of Man-Made Objects*, Ascona, Switzerland.

Masnou, S. and Morel, J. 1998. Level-lines based disocclusion. In *5th IEEE Int'l Conference on Image Processing*, pp. 259–263.

Nitzberg, M., Mumford, D., and Shiota, T. 1993. *Filtering, Segmentation and Depth*. Springer-Verlag: Berlin.

Stamos, I. and Allen, P.K. 2002. Geometry and texture recovery of scenes of large scale. In *Computer Vision and Image Understanding (CVIU)*, 88(2):94–118.

Stulp, F., Dell'Acqua, F., and Fisher, R.B. 2001. Reconstruction of surfaces behind occlusions in range images. In *Proc. 3rd Int. Conf. on 3-D Digital Imaging and Modeling*, Montreal, Canada, pp. 232–239.

Thrun, S., Burgard, W., and Fox, D. 2000. A real-time algorithm for mobile robot mapping with applications to multi-robot and 3D mapping. In *Proc. of International Conference on Robotics and Automation*, San Francisco, vol. 1. pp. 321–328.

Wang, X., Totaro, S., Taillandier, F., Hanson, A., and Teller, S. 2002. Recovering facade texture and microstructure from real-world images. In *Proc. 2nd International Workshop on Texture Analysis and Synthesis in Conjunction with European Conference on Computer Vision*, pp. 145–149.

Zhao, H. and Shibasaki, R. 1999. A system for reconstructing urban 3D objects using ground-based range and CCD images. In *Proc. of International Workshop on Urban Multi-Media/3D Mapping*, Tokyo.

Zhao, H. and Shibasaki, R. 2001. Reconstructing urban 3D model using vehicle-borne laser range scanners. In *Proc. of the 3rd International Conference on 3D Digital Imaging and Modeling*, Quebec, Canada.

Au: Pls. provide year