DataProcessingAlgorithmsforGeneratingTextured 3DBuildingFaçadeMeshesFromLaserScansand CameraImages

ChristianFrühandAvidehZakhor

VideoandImageProcessingLab Dept.ofElectricalEngineeringandComputerScienc es UniversityofCalifornia,Berkeley

algorithms In this paper, we develop a set of data processing forgeneratingtexturedfacademeshesofcitiesfro maseriesof vertical 2D surface scans, obtained by a laser scan ner while drivingonpublicroadsundernormaltrafficcondit ions.These processing steps are needed to cope with imperfecti ons and non-idealities inherent in laser scanning systems s uch as occlusionsandreflectionsfromglasssurfaces. The drivenpath is cut into easy-to-handle quasi-linear segments wi th approximatelystraightdrivingdirection. The verti calscansin each segment are recorded in a sequential topologic al order. and are transformed into a depth image; if availabl e. additional3Ddatafromothermodalitiessuchasst ereovision is also taken into consideration. Dominant building structures are detected in the depth images, and points are cl assifiedinto foreground and background layers. Large holes in th e backgroundlayer, caused by occlusion from foregrou ndlayer objects, are filled in by planar or horizontal inte rpolation.The depth image is further processed by removing isolat ed points and filling remaining small holes, to obtain a text uredsurface mesh. We apply the above steps to a large set of dat a with several million 3D points, and show photorealistic texture mapped3Dmodels.

Keywords: 3D city model, depth image, occlusion, image restauration, urbansimulation

I. I NTRODUCTION

Three-dimensional models of urban environments are useful in a variety of applications such as urban p lanning, training and simulation for urban terrorism scenari os. and virtual reality. Currently, the standard technique for creating large scale city models in an automated or semiautomated way, is to use stereovision approaches o naerial or satellite images [5,9,14]. In recent years, ad vances in resolutionandaccuracyofairbornelaserscanners havealso rendered them suitable for the generation of reason able models[8,10].Bothapproacheshavethedisadvanta gethat theirresolutionisonlyintherangeof1to2fee t,andmore importantly, they can only capture theroofs of the buildings but not the facades. This essential disadvantage p rohibits their use in photo realistic walk or drive-through applications. There have been various attempts to a cquire mplete the complementary ground-level datanecessary to co existing airborne models, either using stereo visio n[3] or

3D laser scanners [11]. But as data has to be acqui red in a slow stop-and-go fashion, these approaches do nots cale to more than few buildings.

In previous work, we have developed a method capabl eof rapidly acquiring 3D geometry and texture data for an entire city at the ground level by using fast 2D la ser scanners and a digital camera [6, 7]. The data acqu isition system is mounted on a truck, moving at normal spee dson public roads, collecting data to be processed offli ne. This approach has the advantage that data can be acquire d continuously, rather than in a stop-and-go fashion, and is therefore extremely fast. Relative position changes are computed with centimeter accuracy by matching succe ssive horizontal laser scans against each other; global p ositionis determined by additional use of aerial photos and d igital roadmaps[7]. As a result, façade scan points arer egistered with the aerial photos or airborne laser scans, fac ilitating subsequentfusionwithmodelsderivedfromairborne data.

In this paper, our goal is to create a detailed, te xtured 3D façade mesh to represent the building walls at the highest level of detail. As there are many erroneous scan p oints, 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 wh ose distance is below a threshold value, does not resul t in an acceptable reconstruction of the street scenery, as shownin Figures 1(a) and 1(b). Even though the 3D structure canbe easily recognized when viewed from a viewpoint near the original acquisition position as in Figure 1(a), th e mesh appears noisy due to several reasons; first, there are holes and erroneous vertices due to reflections off the g lass on windows: second, there are many pieces of geometry "floating in the air", corresponding to partially c aptured objects and measurement errors. The mesh appears to be evenmoreproblematic when viewed from other viewpo ints suchastheoneshowninFigure1(b);thisisbecau seinthis case the large holes in the building facades caused by occluding foreground objects, such as cars and tree s, become visible. Furthermore, since the laser scan o nly captures the frontal view of foreground objects, th ey become almost unrecognizable when viewed sideways. As we drive by a street only once, it is not possible to use additionalscansfromotherviewpointstofilling apscaused by occlusions, as is done in [2, 11]. Rather, we ha ve to reconstruct occluded areas by only using cues from neighboring scanpoints; as such, there has been little work to solve this problem [12].



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

In this paper, we propose a class of data processin g techniques to create visually appealing façade mesh es by removing noisy foreground objects and filling holes inthe building facades. Our objectives are robustness and efficiency with regards to processing time, in orde r to ensure scalability to the enormous amount of datar esulting fromacityscan. Theoutline of this paperisas follows:In section II we introduce our data acquisition system and position estimation; sections III and IV discuss pa th splitting and depthimage generation schemes. We de scribe our strategy to transform the raw scans into a visu ally appealing facademeshinsections Vthrough VII, au tomatic texture mapping in section VIII, and the experiment al resultsinsectionIX.

II. D ATAACQUISITIONAND POSITION ESTIMATION

As described in [6], 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 seen in Figure 2, this is mounted on a rack of approximatel y 3.6 meters height on top of a truck, in order to obtain measurements that are not obstructed by pedestrians and cars. The scanners have a 180 ⁰ field of view with a resolution of 1 ⁰, arange of 80 meters and an accuracy of ± 6

centimeters. Both 2D scanners are facing the same s ide of the street. One is mounted horizontally and is used for positionestimation. The estimation is based on sca n-to-scan matching and global correction algorithms using aer ial photos, and is described in detail in [7]. The othe rscanner ismountedvertically with the scanning plane or tho gonalto the driving direction, and scans the buildings and street sceneryasthetruckdrivesby.Thecameraisorien tedinthe same direction as the scanners, with its center of projection approximately in the intersection line of the two s canning planes. All three devices are synchronized with eac hother using hardware-generated signals, and their coordin ate systemsareregistered with respect to each other.



Figure2:Truckwithdataacquisitionequipment.



Figure3:Drivenpathoverlaidwithroadmap.

We introduce a Cartesian world coordinate system [x ,y,z] where x,y is the ground plane and z points into the sky. Assuming that the city streets are flat, the positi on of the truck can be described by the two coordinates x, y a nd an orientation angle θ of the truck coordinate system. Using the localization methods in [6, 7], current speed, position andorientationofthetruckcanbeestimatedfore achscan. Thus, the entire "capture" path of the acquisition truckcan bereconstructed, as shown in Figure 3, togetherwi thalong series of vertical 2D scans, associated with scanne

position. To partially compensate for the unpredict able, non-uniform motion of the truck, the scan series is subsampled such that the spacing between successivescans _nusedforthe isroughlyequidistant.ForeachverticalscanS 3D reconstruction, there is a tuple (x $_{n}$, y $_{n}$, θ_{n}) which describes position and orientation of the scanner i n the world coordinate system during acquisition. Further more, let s _{n,v} be the distance measurement on a point in scan S with azimuth angle υ . Then, $d_{n,\upsilon} = \cos(\upsilon) \cdot s_{n,\upsilon}$ is the depth value of this point with respect to the scanner, i. e. its orthogonal projection into the ground plane, as sho wn in Figure4.



III. S EGMENTATIONOFTHEDRIVINGPATHINTOQUASI LINEARSEGMENTS

The captured data during at wenty minuted rive cons istsof tens of thousands of scan columns. Since successive scans intime correspond to spatially close points, e.g. abuilding or a side of a street block, it is computationally advantageous not to process the entire data as one block, rather to split it into smaller segments to be proc essed separately. We impose the constraints that (a) path segments have low curvature, and (b) scan columns h avea regular grid structure. This allows us to readily i dentifythe neighbors to right, left, above and below for each point, and, as seen later, is essential for the generation ofadepth imageandsegmentationoperations.

Scanpointsforeachtruckpositionareobtainedas wedrive bythestreets.Duringstraightsegments,thespati alorderof the 2D scan rows is identical to the temporal order ofthe scans, forming aregular topology. Unfortunately, t hisorder of scan points can be reversed during turns towards the scanner side of the car. Figures 5(a) and (b) show the scanning setup during such a turn, with scan planes indicated by the two dotted rays. During the two ve rtical scans, the truck performs not only a translation bu t also a

rotation, making thes canner look slightly backward sduring the second scan. If the targeted object is close en ough, as shown in Figure 5(a), the spatial order of scanpoints is however, if the object is further away than a critical distance d_{crit}, the spatial order of the two scanpoints is reversed, as shown in Figure 5(b).



Figure5:Scangeometryduringaturn, (a)normalscanorderforcloserobjects; (b)reversedscanorderforfurtherobjects.

For a given truck translation of s, and a rotation θ between successive scans, the critical distance can be computed as

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

Thus, d crit is the distance at which the second scanning plane intersects with the first scanning plane. For а particularscanpoint, the order with its predecess orsshould bereverse difits depthd _{n,v} exceedsd _{crit}; this means that its geometric location is somewhere in between points o f previous scans. The effect of such order reversal c an be seen in the marked area in Figure 6. At the corner, the ground and the building walls are scanned twice, fi rstfrom a direct view and then from an oblique angle, and h ence with significantly lower accuracy. For the oblique points, the scans are out of order, destroying the regular topology betweenneighboringscanpoints.



Figure6:Scanpointswithreversedorder.

Since the "out of order" scans obtained in these sc enarios correspond to points that have already been capture dby"in order" scans, and are therefore redundant, our appr oachis to discard them and use only "in order" scans. For typical values of displacement, turning angle, and distance of structures from our driving path, this occurs only in scans of turns with significant angular changes. By remov ing these "turn" scans and splitting the path at the "t urning points", we obtain path segments with low curvature that canbe considered as locally quasi-linear, and can therefore be conveniently processed as depthimages, as descr ibedin the following section. In addition, to ensure that these segments are not too large for further processing, we subdivide them if they are larger then a certain si ze: specifically, in segments that are longer than 100 meters, we identify vertical scans that have the fewest sca npoints abovestreetlevel, corresponding to emptyregions inspace, and segment at these locations. Furthermore, we det ect redundant path segments for areas captured multiple times due to multiple drive bys, and only use one of them for reconstruction purposes. Figures 7(a) and 7(b) show an exampleofanoriginalpath, and the resulting path segments overlaid on a road map, respectively. The small lin es perpendicular to the driving path indicate the scan ning planeoftheverticalscannerforeachposition.



Figure7:Drivenpath,(a)beforesegmentation;(b) after segmentationintoquasi-linearsegments.

IV. C ONVERTINGPATHSEGMENTSINTODEPTHIMAGES

In the previous section, we create path segments th at are guaranteed to contain no scan pairs with permuted horizontalorder.Astheverticalorderisinherent tothescan itself,allscanpointsofasegmentforma3Dscan gridwith regular,quadrilateraltopology.This3Dscangrid allowsus to transform the scan points into a 2.5D representa tion, i.e. a depth image where each pixel represents a scan p oint, and the gray value for each pixel is proportional t o the depth of the scan point. The advantage of a depth i mageis its intuitively easy interpretation, and the increa sed processing speed the 2D domain provides. However, m ost operations that are performed on the depth image ca n be done just as well on the 3D point grid directly, on lynotas conveniently.

A depth image is typically used for representing th e data from 3D scanners. Even though the way the depth val ueis 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 g round 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 thanthedistances _{n,v}, so that the depth image is basically a tilted height field. The advantage is that in this case points thatlieonaverticalline, e.g. abuilding wall, havethesame 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 h ave a single center from which the scan points are measur ed: instead, there are different centers for each indiv idual vertical columnalong the path segment. The obtaine ddepth image is neither a polar nor a parallel projection; it resembles most to a cylindrical projection. Due to nonuniform driving speed and non-linear driving direct ion, these centers are ingeneral not on a line, but on anarbitrary shaped, though low-curvature curve, and the spacing between them is not exactly uniform. Because of thi s, strictly speaking the grid position only specifies the topological order of the depth pixels, and not the exact3D point coordinates. However, as topology and depth v alue are a good approximation for the exact 3D coordinat es. especially within a small neighborhood, we choose t oapply our data processing algorithms to the depth image, thereby facilitating use of standard image processing techn iques such as region growing. Moreover, the actual 3D ver tex coordinatesarestillkeptandusedfor3Doperatio nssuchas plane fitting. Figure 8(a) shows an example of the 3D vertices of a scan grid, and Figure 8(b) shows its corresponding depth image, with a gray scale propor tional tod n.v.





(U)

Figure 8: Scan grid representations; (a) 3D vertice s; (b)depthimage.

V. P ROPERTIESOF CITY LASER SCANS

Inthissection, webriefly describe properties of scanstaken inacity environment, resulting from the physics o falaser scannerasanactivedevicemeasuringtime-of-fligh toflight rays. It is essential to understand these propertie s and the resulting imperfections in distance measurement, si nce 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 reconstructin gwhatthe human eye or a camera would observe while moving around in the city. As such, we discuss the discrep ancies between these two different sensing modalities in t his section.

a) Discrepancies due to different resolution

The beam divergence of the laser scanner is about 1 5 milliradians (mrad) and the spacing, hence the angu lar 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 border s. Therefore, small or thin objects, such as cables, f ences. streetsigns, lightposts and tree branches, are cl earlyvisible inthecameraimage, but only partially captured in thescan.

Hence they appear as "floating" vertices, as seen i n the depthimageinFigure9.



Figure9:"Floating"vertices.

b) Discrepancies due to the measurement physics

Cameraandevearepassivesensors, capturing light froman external source; this is in contrast with a laser s canner. which is an active sensor, and uses light that ite mitsitself. This results in substantial differences in measurem ent of reflecting and semitransparent surfaces, which are in form of windows and glass fronts frequently present in u rban environments. Typically, there is at least 4% of th e light reflectedatasingleglass/airtransition, soato talofatleast 8% perwindow; if the window has a reflective coat ing,this canbelarger. The camera typically sees are flecti onofthe skyor anearby building on the window, often disto rtedor mergedwithobjectsbehindtheglass.Althoughmost 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, resulting innodi stancevalue, or is transmitted through the glass. In the latter case, if it hits a lambertian surface as shown in Figure 10, th e backscattered light travels again through the glass . The resulting surface reflections on the glass only wea ken the laser beam intensity, eventually below the detectio n limit, but do not otherwise necessarily affect the distanc e measurement. To the laser, the window is quasi nonexistent, and the measurement point is generally no tonthe window surface, unless the surface is orthogonal to the beam. In case of multi-reflections, the situation b ecomes evenworseasthemeasureddistanceisalmostrando m.



Figure10:Lasermeasurementincaseofaglasswin dow

c)Discrepanciesduetodifferentscanandviewpoin ts

Laser and camera are both limited in that they can only detect the first visible/backscattering object alon g a measurement direction and as such cannot deal with occlusions. If there is an object in the foreground .suchasa tree in front of a building, the laser cannot captu re what is behindit; hence, generating a mesh from the obtain edscan pointsresultsinaholeinthebuilding.Werefer tothistype of mesh hole as occlusion hole. As the laser scan points resemble a cylindrical projection, but rendering is parallel orperspective, in presence of occlusions, it is im possibleto reconstruct theoriginal view without any holes, ev enforthe viewpoints from which data was acquired. This is a special property of our fast 2D data acquisition method. An interesting fact is that the wide-angle camera imag es captured simultaneously with the scans often contai n parts of the background invisible to the laser. These cou ld in potentially be either used to fill in geometry usin g stereo techniques, or toverify the validity of the filled ingeometry obtained from using interpolation techniques.

For a photo realistic model, we need to devise tech niques for detecting discrepancies between the two modalit ies, removing invalids canpoints, and filling inholes, either due to occlusion or due to unpredictable surface proper ties; we will describe our approaches to these problems in t he following sections.

VI. MULTI-LAYER REPRESENTATION

To ensure that the facade model looks reasonable fr om every viewpoint, it is necessary to complete the ge ometry for the building facades. As our facades are not on 1y manifolds, butalso resemble a heightfield, it is possibleto introducearepresentationbasedofmultipledepth layersfor the street scenery, similar to the one proposed in [1]. Each depthlayerisascangrid, and the scanpoints of theoriginal gridareassignedtoexactlyoneofthelayers.If atacertain grid location there is a point in a foreground laye r, this location is empty in all layers behind it and needs to be filledin.

Even though the concept can be applied to an arbitr ary numberoflayers, for our problem it is sufficient togenerate onlytwolayers, a foreground and a background. To assign a scan point to either one of the two layers we mak e the following assumptions about our environment: Main structures, i.e. buildings, are usually (a) vertica l, and (b) extend over several feet in horizontal dimension. F oreach vertical scan n corresponding to a column in the de pth image, we define the main depth as the depth value that occurs most frequently, as shown in Figure 11. The scan vertices corresponding to the main depth lie on a v ertical line, and the first assumption suggests that this i s a main structure, such as a building, or perhaps other ver tical objects, such as a street light or a tree trunk. With the second assumption, we filter out the latter class of the fvertical objects. More specifically, our processing steps can be described as follows:

Wesortalldepthvaluess _{n.v}foreachcolumnofthedepth image into a histogram as shown in Figures 11(a) an d(b), and detect the peak value and its corresponding dep th. Applying this to all scans results in a 2D histogra m as shown in Figure 12, and an individual main depth va lue estimate for each scan. According to the second assumption, isolated outliers are removed by applyi ng a median filter on these main depth values across the scans, n. We and a final depth value is assigned to each column definea" split" depth, γ_n , for each column, and set it to the firstlocalminimumofthehistogramoccurringimme diately before main depth, i.e. with a depth value smaller than the main depth. Taking the first minimum in the distrib ution insteadofthemainvalueitselfhastheadvantage thatpoints clearly belonging to foreground layers are splits o ff, whereasoverhangingpartsofbuildings, for which t hedepth isslightlysmallerthanthemaindepth, arekepti nthemain layer where the ylogically belong to, as shown in Figure11.



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.



Figure12:Two-dimensionalhistogramforallscans.

A point can be identified as a ground point if its Ζ coordinate has a small value and its neighbors int hesame scan column have a similarly low z value. We prefer to include the ground in our models, and as such, assi gn ground points also to the background layer. Therefo re, we $_{n,\upsilon}$ to the perform the layer split by assigning a scan point P background layer, if s $_{n,\nu} > \gamma_n$ or P $_{n,\nu}$ is a ground point, and otherwise, to the foreground layer. Figure 13 shows an example for the resulting foreground and background layers.



Figure13:(a)Foregroundlayer;(b)backgroundlay er.

Since the steps described in this section assume th e presence of vertical buildings, they cannot be expe cted to work for segments that are dominated by trees; this also applies to the processing steps we introduce in the following sections. As our goal is to reconstruct b uildings, pathsegments can be left unprocessed and included "asis" in the city model, if they do not contain any struc ture. A characteristicofatreeareaisitsfractalgeomet ry, resulting in a large variance among adjacent depth values, or even more characteristic, many significant vector direct ion changes for the edges between connected mesh vertic es. We define a coefficient for the fractal nature of a segment by counting vertices with direction changes greater than a specific angle, e.g. twenty degrees, and dividing t hem by the total number of vertices. If this coefficient i slarge, the segment is most likely a tree area and should not b e made subject to the processing steps described in this s ection. This is for example the case for the segment shown in Figure9.

After layer splitting, all grid locations occupied in the foreground layer are empty in the background layer. The vertical laser does not capture any occluded geomet ry, and in the next section we will describe an approach fo rfilling these empty grid locations based on neighboring pix els. However, in our data acquisition system there are 3 D vertices available from other sources, such as ster eo vision and the horizontal scanner used for navigation. Thu s, it is conceivable to use the sead ditional information tofillsome inthedepthlayers.Ourapproachtodoingsoisas follows:

Given a set of 3D vertices V i obtained from a different modality, determine the closest scan direction for each vertex and hence the grid location (n, v) it should be assigned to. As shown in Figure 14, each V is assigned to the vertical scanning plane, S_n , with the smallest Euclidean distance, corresponding to column n in the depth im age. Usingsimpletrigonometry, the scanning angle under which this vertex appears in the scanning plane, and henc e the depthimagerow v, can be computed, as well as the depth d_{n.v}ofthepixel.



Figure14:Sortingadditionalpointsintothelayer s.

Wecannowusetheseadditionalverticestofillin theholes. Tobegin with, all vertices that do not belong to b ackground holes are discarded. If there is exactly one vertex falling onto a grid location, its depth is directly assigne d to that grid location; for situations with multiple vertice s, median depthvalueforthislocationischosen.Figure15 showsthe background layer from Figure 13(b) after sorting in 3D vertices from stereo vision and horizontal laser sc ans. As seen, some holes can be entirely filled in, and the size of othersbecomessmaller.e.g.theholesduetotrees inthetall building on the left side. Note that this intermedi atestepis optional and depends on the availability of additio nal 3D data.



Figure 15: Background layer after sorting in additi onal points from other modalities.

VII. B ACKGROUND LAYERPOSTPROCESSINGAND MESH GENERATION

In this section, we will describe a strategy to rem ove erroneousscanpoints, and to fill inholes in the background layer. There exists a variety of successful hole fi lling approaches, for example based on fusing multiple sc ans taken from different positions [2]. Most previous w ork on hole filling in the literature has been focused on reverse engineering applications, in which a 3D model of an object is obtained from multiple laser scans taken from di fferent locations and orientations. Since these existing ho lefilling approaches are not applicable to our experimental s etup, our approach is to estimate the actual geometry bas ed on thesurroundingenvironmentandreasonableheuristi cs.One cannot expect this estimate to be accurate in all possible cases, rather to lead to an acceptable result in most cases. thus reducing the amount of further manual interven tions and postprocessing drastically. Additionally, the e stimated geometry could be made subject to further verificat ion steps, such as consistency checks, by applying ster eovision techniquestotheintensityimagescapturedbythe camera.

Ourdatatypicallyexhibitsthefollowingcharacter istics:

- Occlusion holes, such as those caused by a tree, arelargeandcanextendoversubstantialparts of building.
- Asignificant number of scan points surrounding a holemay beer roneous due toglass surfaces.

- Ingeneral, aspline surface filling is unsuitable, as building structures are usually piecewise planar withsharpdiscontinuities.
- The size of data set resulting from a city scan is huge, and therefore the processing time per hole shouldbekepttoaminimum.

Basedontheaboveobservations, we propose the fol lowing steps for data completion:

1. Detecting and removing erroneous scan points in the backgroundlayer

We assume that erroneous scan points are due to gla SS surfaces, i.e. the laser measured either an interna wall/object, or a completely random distance due to multireflections. Either way, the depth of the scan poin ts measured through the glassis substantially greater thanthe depth of the building wall, and hence these points are candidates for removal. Since glass windows are usu ally framedbythewall, weremove the candidate points onlyif they are embedded among a number of scanpoints atmain depth. An example of the effect of this step can be seenby comparing the windows of the original image in Figu re 16(a)withtheprocessedbackgroundlayerinFigure 16(b).

2.Segmentingtheoccludingforegroundlayerintoo **bjects** Inordertodetermineholesinthebackgroundlayer caused by occlusion, we segment the occluding foreground l ayer into objects and project segmentation onto the back ground " at a layer. This way, holes can be filled in one "object time, rather than all at the same time; this approa chhasthe advantage that more localized hole filling algorith ms are morelikelytoresultinvisuallypleasingmodelst hanglobal ones.Wesegmenttheforegroundlayerbytakingar andom seedpointthatdoesnotyetbelongtoaregion,an dapplying a region growing algorithm that iteratively adds neighboringpixelsiftheirdepthdiscontinuityor theirlocal curvature is small enough. This is repeated until a llpixels are assigned to a region, and the result is a regio n map as shown in Figure 16(c). For each foreground region, we determine boundary points on the background layer; these areallthevalidpixelsinthebackgroundlayerth atareclose toholepixelscausedbytheoccludingobject.

3. Filling occlusion holes in the background layer for each region

As the foreground objects are located in front of m ain structures and in most cases stand on the ground, t hev occlude not only parts of a building, but also part s 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 t he ground and not in the main structure. This is because the laser scanneris mounted on top of a rack, and as such ha satop down view of the car. As a plane is a good approxim ation 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 betweenthelinethrough the main depths can points and the



(a)Initialdepthimage.



(b)Backgroundlayerafterremovinginvalidscanpo ints.



(c)Foregroundlayersegmented.



(d)Occlusionholesfilled.



(e)Finalbackgroundlayerafterfillingremaining holes. Figure16:Processingstepsofdepthimage.

linethroughgroundscanpoints. The angle υ'_n at which this point appears in the scan marks the virtual boundar between ground part and structure part of the scan; we fill in structure points above and ground points below t boundary differently.

ApplyingaRANSAC algorithm, we find the plane with the maximum consensus, i.e. maximum number of ground boundary points on it, as the optimal ground plane forthat local neighborhood. Each hole pixel with $\upsilon < \upsilon'_n$ is then filled in with a depth value according to this plan e. It is possible to apply the same technique for the struct urehole pixels, i.e. the pixels with $\upsilon > \upsilon'_n$, by finding the optimal plane through the structure boundary points and fil ling in the hole pixels accordingly. However, we have found that incontrasttotheground, surrounding building pix elsdonot oftenlieonaplane.Instead,therearediscontinu itiesdueto occluded boundaries and building features such as marquees or lintels, in most cases extending horizo ntally acrossthebuilding. Therefore, rather than filling holeswith a plane, we fill in structure holes line by line ho rizontally, insuchawaythatthedepthvalueateachpixelis thelinear interpolation between the closest right and left st ructure boundary point, if they both exist; otherwise no va lue is filled in. In a second phase, a similar interpolati onisdone vertically, using the already filled in points as v alid 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 lave r is showninFigure16(d).

4. Postprocessing the background layer

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

 $\begin{array}{ll} \mbox{In order to create a mesh, each depth pixel can be} \\ \mbox{transformed back into a 3D vertex, and each vertex} \\ \mbox{connected to a depth imageneighbor P} & P_{n,\nu \, \nu \, \nu} \mbox{if} \end{array} \end{array} \\ \begin{array}{ll} P_{n,\nu \, is} \\ P_{n,\nu \, is} \end{array}$

 $|s_{n+\Delta n, \upsilon+\Delta \upsilon} - s_{n,\upsilon}| < s_{\max}$ orif

 $\cos \varphi > \cos \varphi_{max}$

with

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

 s_{max} . The resulting quadrilateral meshis split into triangles, and mesh simplification tools such as Qslim [4] or VTK decimation can be applied to reduce the number of triangles.



(a)CameraImage.



(b)Meshtrianglesprojectedintotheimage,withs ome foregroundandbackgroundtrianglesprojectingtot he sameimagearea(arrow).



(c)Foregroundobjectsmarkedwhite.

Figure17:Meshtrianglesprojectedintocameraima ges.

VIII. A UTOMATIC TEXTURE MAPPING

As photorealism cannot be achieved by using geometr y alone, and requires color texture as well, our data acquisition system includes a digital color camera with a wide angle lens. The camera is synchronized with th etwo laser scanners, and is calibrated against the laser scanners' coordinate system, and hence, the camera positions canbe computed for all images. After calibrating the came ra and removing the lens distortion in the images, each 3D vertex can be mapped to its corresponding pixel in an inte nsity image by a simple projective transformation. As the 3D mesh triangles are small compared to their distance to the camera, perspective distortions within a triangle c an be neglected, and each mesh triangle can be mapped to а triangle in the picture by applying the projective transformationtothethreecornerpoints.

As described in section V, camera and laser scanner have different viewpoints during data acquisition, and i n most camera pictures at least some mesh triangles of the background layer are occluded by foreground objects ; this is particularly true for triangles that consist of filled-in points. An example of this is shown in Figure 17(b) ,where triangles of the tree and the building project to t he same areas in the image. Although the pixel location of the projectedbackgroundtrianglesiscorrect, the corr esponding texturetrianglesareincorrect, and merely corresp ondtothe foregroundobjects.

However, after splitting the scan points to the two layers, the foreground geometry is readily identified, and both foregroundscanpointsandtrianglescanbemarked ineach camerapicture, as shown in Fig 17(c) with white co lor.In ordertoselectspecificcameraimagestobeusedf ortexture mapping of a specific mesh triangle of the backgrou nd layer, for each picture containing a triangle in it s field of view, we determine whether any of the triangle's co rner points maps onto a white marked pixel; if this is t he case, thepicturecannotbeusedfortexturingthetriang le.Forthe particular wide angle lens we use and typical build ing topologies, a background layer point is usually in the field of view of about 10 to 20 pictures. In most situati ons, we can use multiple images, and as such, we choose the most direct view to texture map the triangle. However, i f foreground object and building façade are too close , some façade triangles might not be visible in any pictur e and hence cannot be texture mapped at all. A possible s olution for this case is to apply a texture synthesis algor ithm in order to create artificial texture. We are current ly investigatingthisapproach.

IX. R ESULTS

Wedroveourequippedtruckona6769meterslongp athin downtown Berkeley, starting from Blake street throu gh Telegraph avenue, and in loops around the downtown blocks. During this 24-minute-drive, we captured 10 7,082 vertical scans, consisting of 14,973,064 scan point S. Applying the described path splitting techniques, w e cut this pathinto 73 segments, as shown in Figure 180 verlaid with a road map. There is no need for further manua cutting, even at Shattuck Avenue, where the "Manhat tan geometry" of Berkeley is not preserved. For each of the73 segments, we generate two meshes for comparison: th efirst mesh is obtained directly from the raw scans, and t he secondonefromthedepthimagetowhichwehaveap plied thepostprocessingstepsdescribedinprevioussect ions.The processing time for the millions of scanpoints of theentire pathisabout2hoursona2GHzPentium4PC.



Figure18:Entirepathaftersplitinquasi-linear segments.

For 12 out of the 73 segments, additional 3D vertic es derived from stereo vision techniques are available , and hence, sorting in these 3D points into the layers b ased on section VI does fill some of the holes. For these s pecific holes, we have compared the results based on stereo vision vertices with those based on interpolation alone as described in section VII, and have found no substan tial difference; often the interpolated meshvertices ap peartobe slightlymoreaccurate, as they are less noisy than thestereo vision based vertices. Figure 20(a) shows an exampl e before processing, and Figure 20(b) shows the tree holes completely filled in by stereo vision vertices. As seen, the outline of the original holes can still be recogniz ed in Figure 20(b), whereas the points generated by inter polation alone are almost indistinguishable from the surroun ding geometry, asseen in Figure 20(c).

Significantlybetter	35	4	8%
Better	17	2	3%
Same	15	2	1%
Worse	5	79	6
Significantlyworse	1	1%	
Total	73	1	00%

Table1:Subjectivecomparisonoftheprocessedmeshvs.theoriginalmeshforall73segments.



Figure 19: Hole filling (a) original mesh with hole s behindoccludingtrees;(b)filledbysortinginad ditional 3D points using stereo vision; (c) filled by using the interpolationtechniquesofsectionVII.

Wehavefoundourapproachtoworkwellinthedown town areas, where there are clear building structures an d few trees.However,inresidentialareas,wherethebui ldingsare often almost completely hidden behind trees, it is difficult to accurately estimate the geometry. As we do not h avethe ground truth to compare with, and as our main conce rnis thevisual quality of the generated model, we have manually inspectedtheresultsandsubjectivelydeterminedt hedegree to which the proposed postprocessing procedures hav e improved the visual appearance. The evaluation resu lts for all 73 segments before and after postprocessing tec hniques described in this paper are shown in Table 1; the postprocessing does not utilize auxiliary 3D vertic es from horizontallaserscannerorthecamera.Eventhough 8% of all processed segments appear visually worse than t he original, the overall quality of the facade models is significantly improved. The important downtown segm ents are in most cases ready to use and do not require f urther manualintervention.

Thefewproblematicsegmentsalloccurinresidenti alareas, consisting mainly of trees. The tree detection algo rithm "critical" described in section VI classifiesten segments as in that too many trees are present; all six problem atic segments corresponding to "worse" and "significantl У worse" rows in Table 1 are among them, yet none of the improved segments in rows 1 and 2 are detected as c ritical. This is significant because it shows that (a) all p roblematic segments correspond to regions with a large number of trees, and (b) they can be successfully detected an dhence not be subjected to the proposed steps. Table 2 sho wsthe evaluation results if only non-critical segments ar e processed. As seen, the postprocessing steps descri bed in this paper together with the tree detection algorit hm improve over 80% of the segments, and never result in degradationsforanyofthesegments.

Significantlybetter	35	5	5%
Better	17	2	7%
Same	11	1	7%
Worse	0	09	6
Significantlyworse	0	0%	
Total	63	1	00%

Table 2: Subjective comparison of the processed meshvs. the original mesh for the segments automaticallyclassifiedasnon-tree-areas.y

In Figure 20 we show before and after examples, and the corresponding classifications according to Tables 1 and 2. As seen, except for pair "f", the proposed postproc essing stepsresultinvisuallypleasing models. Pairfin Figure 20 is classified by our tree detection algorithm as cr itical, and hence, should be left "asis" rather than processed .Wehave further reduced all 73 meshes using the Qslim mesh simplification tool [4] to create multiple levels o f details. This enables us to render the façades for the entir e path with a standard VRML viewer as shown in Figure 21. For bettervisualization, these facades are superimpose doveran aerialimage.

Note that the evaluation of our technique is based on comparing the non-textured geometry; in a compariso n in whichbothmodelsaretexturemapped, the processed mesh is more likely to be visually superior to the origi nal. Texture distracts the human eye from geometry imperfections such as those introduced by hole fill ing algorithms. In Figure 22 we compare a textured mesh ofa cityblock without and with processing. As seen, th evisual difference between the two meshes is striking, for the reasons described above. Note the façade area occlu dedby the two trees on the left side of the original mesh hasbeen filled in, and texture mapped using camera views as visiblein describedinsectionVIII. A few triangles are not any camera image and are therefore left untextured. We

plan to subject these to a texture synthesis algori thm in a nearfuture.

X. C ONCLUSIONSAND FUTURE WORK

We have proposed a strategy to create building faca de meshes from large laser surface scans. Future work will focus on using color and texture cues to verify fil led-in geometry. Additionally, foreground objects could be classified and replaced by appropriate generics; th egroundlevel façade meshes can also be merged with rooftop modelsobtained from aerial view.

XI. A CKNOWLEDGEMENT

This work was sponsored by Army Research Office und er contract DAAD19-00-1-0352. We wish to thank Sick, I nc. for their support. We also wish to thank John Flynn for providing the stereovision results.

XII. R EFERENCES

- N.L. Chang and A. Zakhor, "A Multivalued Representation for ViewSynthesis", Proc. Int'lConferenceon ImageProcessing, Kobe, Japan, 1999, vol.2, pp. 505-509
- [2] B. Curless and M. Levoy, "A volumetric method f or building complex models from range images", SIGGRAPH, New Or leans, 1996, p. 303-312
- [3] A. Dick, P. Torr, S. Ruffle, and R. Cipolla, "C ombining Single View Recognition and Multiple View Stereo for Archi tectural Scenes", International Conference on Computer Visio n, Vancouver, Canada, 2001, p. 268-74
- M. Garland and P. Heckbert, "Surface Simplifica tion Using QuadricErrorMetrics",SIGGRAPH'97,LosAngeles, 216
- [5] D. Frere, J. Vandekerckhove, T. Moons, and L. V an Gool, "Automaticmodellingand3Dreconstructionofurban aerialimagery", IEEE International Geoscience and SymposiumProceedings, Seattle, 1998, p.2593-6
- [6] C. Frueh and A. Zakhor, "Fast 3D model generati environments", IEEE Conf. on Multisensor Fusion and forIntelligentSystems, Baden-Baden, Germany, 2001, p.165-170
- [7] C. Frueh and A. Zakhor, "3D model generation of aerial photographs and ground level laser scans", C andPatternRecognition, Hawaii, USA, 2001, p. II-3
 [7] C. Frueh and A. Zakhor, "3D model generation of omputer Vision
- [8] N. Haala and C. Brenner, "Generation of 3D city models from airborne laser scanning data", Proc. EARSEL worksho remotesensingonlandandsea, Tallin, Esonia, 199 7, p. 105-112
- Z.Kim, A. Huertas, and R. Nevatia, "Automatic d escription of Buildings with complex rooftops from multiple image s", Computer VisionandPatternRecognition, Kauai, 2001, p.272 -279
- [10] H.-G. Maas, "The suitability of airborne laser automatic3Dobjectreconstruction", 3.Int'lWorks ExtractionofMan-MadeObjects, Ascona, Switzerland ,2001
- [11] I. Stamos and P.E. Allen, "3-D model construct and image data." Computer Vision and Pattern Recogn HeadIsland, 2000, p.531-6
- F. Stulp, F. Dell'Acqua, and R. B. Fisher, "Re surfaces behind occlusions in range images", Proc.
 3rd Int. Conf. on 3-DDigital Imaging and Modeling, Montreal, Canada, 2001, p.232-239
- [13] S. Thrun, W. Burgard, and D. Fox, "A real-time algorithm for mobile robot mapping with applications to multi-rob ot and 3D mapping", Proc. of International Conference on Robotics and Automation, San Francisco, 2000, p. . 321-8, vol. 1.
- [14] C. Vestri and F. Devernay, "Using Robust Metho ds for Automatic extraction of buildings", Computer Vision and Pattern Recognition, Hawaii, USA, 2001, p.I-133-8, vol. 1.



Figure20:Generatedmeshes,leftsideoriginal,ri Theclassificationforthevisualimpressionis"si "worse"forpairf.

ghtsideaftertheproposedforegroundremovaland holefillingprocedure. gnificantlybetter"forthefirstfourimagepairs, "better"forpaireand



Figure 21: Non-textured facade model for the entire

path, overlaid ontopofanaerial photo.



Figure22:Texturedfacademeshwithout(top)andw ith(bottom)processing.