

# Reconstructing Building Mass Models from UAV images

Minglei Li<sup>a,b</sup>, Liangliang Nan<sup>a,\*</sup>, Neil Smith<sup>a</sup>, Peter Wonka<sup>a</sup>

<sup>a</sup>Visual Computing Center, KAUST, KSA

<sup>b</sup>Institute of Remote Sensing and Digital Earth, CAS, P.R.China

---

## Abstract

We present an automatic reconstruction pipeline for large scale urban scenes from aerial images captured by a camera mounted on an unmanned aerial vehicle. Using state-of-the-art Structure from Motion and Multi-View Stereo algorithms, we first generate a dense point cloud from the aerial images. Based on the statistical analysis of the footprint grid of the buildings, the point cloud is classified into different categories (i.e., buildings, ground, trees, and others). Roof structures are extracted for each individual building using Markov random field optimization. Then, a contour refinement algorithm based on pivot point detection is utilized to refine the contour of patches. Finally, polygonal mesh models are extracted from the refined contours. Experiments on various scenes as well as comparisons with state-of-the-art reconstruction methods demonstrate the effectiveness and robustness of the proposed method.

*Keywords:* urban reconstruction, aerial images, point cloud, Markov random field, graph cut

---

## 1. Introduction

Digital 3D models of urban scenes are important for a variety of applications such as urban planning, navigation, simulation, virtual reality, and entertainment. However, the digitization of urban scenes with complex architectural structures still remains a challenge [1, 2, 3]. Most of traditional surface reconstruction techniques reconstruct objects with smooth surfaces by exploiting either increasingly sophisticated solvers or better formulation of prior knowledge [4]. For urban scenes, since automatic semantic segmentation is very hard to achieve, the reconstruction process (especially for complex architectural structures) requires tedious manual effort.

In the last two decades, a considerable amount of reconstruction approaches have been developed, aiming at automatically modeling large scale urban scenes. Most of these approaches, however, are designed to deal with Light Detection and Ranging (LiDAR) point clouds obtained from airborne planes or ground level vehicles, which usually face expensive device cost and unavoidable severe occlusions. Most recently, state-of-the-art Structure from Motion (SfM) and Multi-View Stereo (MVS) methods [5, 6, 7] have produced extremely compelling results on a wide variety of scenes. A typical SfM and MVS pipeline starts by automatically matching features among the input image sequences, then it recovers the internal and external camera parameters, and produces a sparse and finally a dense 3D point cloud of the scene. To further enhance the data acquisition, in this work we exploit an Unmanned Aerial Vehicle (UAV) mounted with a camera, which provides

more flexibility and significantly improves the efficiency for capturing large scale urban scenes.

Although UAV imagery is more effective in capturing all sides of urban buildings and robust against occlusion, the point clouds computed from SfM and MVS are still noisy and sparse, which hinders automatic processing and reconstruction. To overcome these problems, statistical information from different resolution are extracted to enhance the segmentation and reconstruction. The proposed method manages to classify the point cloud and reconstruct architectural models automatically, and it is robust to a wide range of data qualities.

The contributions of our work include:

- a novel framework for automatic reconstruction of large scale urban scenes from UAV images, which provides realistic reconstruction with semantic information.
- an object level point cloud segmentation algorithm and a roof extraction algorithm based on a regularized MRF formulation, which significantly speeds up the whole reconstruction pipeline.
- an effective contour refinement method based on pivot point detection, which ensures compact final reconstruction.

## 2. Related Work

The reconstruction of urban scenes has been a hot topic in computer graphics and computer vision in the last two decades with large number of approaches recently developed [2, 4]. In this section, we review the work that are most related to the proposed method. We divide these work into three categories according the data sources they use.

---

\*Corresponding author

Email address: liangliang.nan@gmail.com (Liangliang Nan)

58 **Image-based reconstruction.** Using street level ortho-  
59 rectified photographs, Müller et al. [8] devised a procedural  
60 modeling strategy that identifies repeated elements in the facade  
61 image using mutual information. Sinha et al. [9] proposed  
62 an interactive system to recover the 3D structure of buildings  
63 by manually drawing outlines overlaid on 2D photographs  
64 and calculating their intersections. Enhanced by 3D depth  
65 information recovered from SfM, Xiao et al. [10] proposed a  
66 semi-automatic image-based approach for facade modeling.  
67 Garcia-Dorado et al. [11] first calibrated aerial images and  
68 fused them with GIS meta-data to compute a per-building 2.5D  
69 volumetric reconstruction using graph cut.

70 **Laser scan-based reconstruction.** In the last decades, laser  
71 scanners have provided a new type of data source for urban  
72 reconstruction. Lin et al. [12] first classified the point clouds  
73 of a large scale residential area into different categories, and  
74 then performed reconstruction based on segmentation of each  
75 building into basic symmetric and convex blocks. Lafarge  
76 and Mallet [13] proposed a non-supervised approach for point  
77 cloud classification. Then, regular roof sections are represented  
78 by basic geometric primitives and irregular roof components  
79 are represented by the combination of a set of geometric  
80 primitives. By assuming piecewise planar structures, Lafarge  
81 and Alliez [14] reconstruct surfaces using a point consolidation  
82 strategy that preserves of the buildings' structure at a given  
83 scale.

84 Aiming at 2.5D reconstruction, Zhou and Neumann [15]  
85 proposed a data-driven approach to detect a set of principal  
86 directions to align roof boundaries. They used these roof  
87 boundaries to produce a footprint for the reconstruction. Poullis  
88 and You [16] created compact city models from high elevation  
89 LiDAR data by simplifying boundaries of fitted planes. Lafarge  
90 et al. [17] employed Bayesian decision to assemble simple  
91 urban structures as the reconstruction from a single Digital  
92 Surface Model (DSM). By extending the traditional dual  
93 contouring algorithm into 2.5D, Zhou and Neumann [18]  
94 optimized the 2D boundaries for the roofs, which enables the  
95 reconstruction of buildings with arbitrarily shaped roofs. In  
96 their following work [19, 20], the authors further incorporated  
97 topology control and global regularity to improve the dual  
98 contouring results, yielding impressive performance.

99 To reconstruct facade details, Nan et al. [21] proposed an  
100 interactive reconstruction method that exploits the repetitive  
101 structure of the facades. During the drag-and-drop operation,  
102 each facade element is snapped to its proper location based on  
103 discrete optimization that balances between a regularity term  
104 and a data fitting term. By using Manhattan World assumption,  
105 Venegas et al. [22] first segmented the point cloud into walls,  
106 edges, corners, and edge-corners. They then organized the  
107 classified points into clusters to extract a volumetric description  
108 of the buildings.

109 **MVS-based reconstruction.** As images of urban scenes  
110 becomes easier to acquire from both the internet (e.g., flicker)  
111 and cameras (e.g., smart phones), more and more recent  
112 research interests have focused on reconstructing urban scenes  
113 from a set of images or videos.

114 Pollefeys et al. [23] designed a real-time system to generate

115 street level city models from video frames captured by onboard  
116 cameras. Given a dense point cloud reconstructed from a set  
117 of images using SfM and MVS techniques, Arikan et al. [24]  
118 proposed O-Snap, an interactive reconstruction system that  
119 fits planar primitives along with boundary polygons, and then  
120 snaps polygons together to obtain a mesh model of a building  
121 through non-linear optimization. Using the same data source,  
122 Nan et al. [25] proposed to reconstruct detailed urban models  
123 by assembling facade details onto a set of manually extruded  
124 coarse models based on linear integer programming. Some  
125 other approaches [26, 27] are also proposed for reconstruction  
126 of large scale scenes based on MVS. These methods can  
127 generate high resolution results, but semantic information are  
128 ignored during the reconstruction and usually suffer from data  
129 storage difficulties.

130 To obtain a level-of-detail representation of urban scenes,  
131 Verdie et al. [28] introduced an abstraction step between  
132 the classification and reconstruction steps to regularize planar  
133 structures from a large set of plane candidates. Finally, a  
134 surface model is extracted from a set of 3D arrangements  
135 based on a min-cut formulation. Compared with methods  
136 using ground level images, airborne-based data sources cover  
137 larger area of the scene. Most of existing airborne-based  
138 methods [15, 16, 17, 18] describe data in 2.5D due to the data  
139 acquisition strategy. This strategy makes quality reconstruction  
140 of building faces not possible, since only the roof information  
141 is available in the data. In this paper, we focus on the automatic  
142 generation of lightweight urban models from airborne point sets  
143 reconstructed from UAV images.

### 144 3. Overview

145 Our method takes as input a sequence of images of a scene  
146 and outputs 3D polygonal mesh models of the scene. The  
147 images are captured by a camera mounted on a UAV. In a pre-  
148 processing step, we extract a point cloud from these images  
149 using SfM and MVS [29]. Then there are two core steps  
150 for automatic generation of the urban models: point cloud  
151 classification and roof extraction. The main idea for automating  
152 these processes relies on a regularized MRF labeling strategy.  
153 An overview of our method is shown in Figure 1.

154 We first classify the point cloud of a large scene into four  
155 different categories, i.e., buildings, ground, trees, and others.  
156 We define a set of point features based on a 2D supporting  
157 grid by projecting the point set onto the ground. Then the  
158 classification is achieved using a regularized MRF formulation.  
159 Graph cut [30, 31] is used to solve the labeling problem (see  
160 Section 4).

161 After point cloud classification, the data for each building  
162 is processed independently. By projecting the points onto the  
163 ground plane, a depth map is first generated to represent the  
164 2.5D structure of a building. Based on the depth map, a higher  
165 resolution regularized MRF formulation is used to extract the  
166 roof structure of the building, followed by a regularization step  
167 for the roof contours. Finally, polygonal mesh models are  
168 generated by extruding the roof patches onto the ground (see  
169 Section 5).

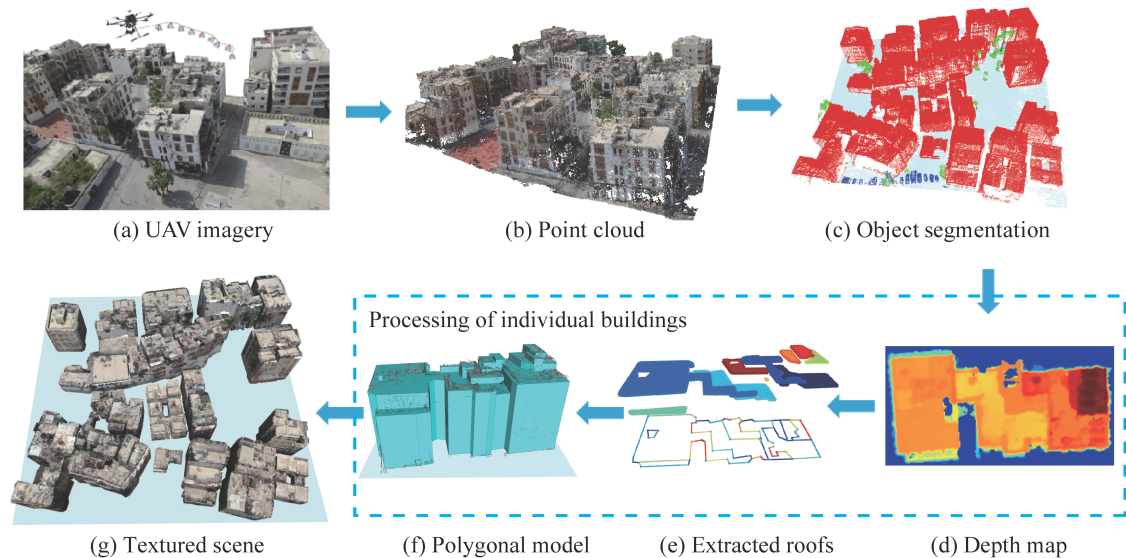


Figure 1: An overview of the proposed reconstruction pipeline. From a sequence of images captured by the camera mounted on an UAV (a), a point cloud (b) is generated using SfM and MVS. Then an object level segmentation is performed to decompose the entire scene into buildings and other objects (c). For each individual building, we extract the roofs (e) from its depth map (d) defined on a grid representation. Then a polygonal model (f) is extracted from the roofs (e). Finally, the entire scene can be textured (g) for various applications.

#### 170 4. Object Level Segmentation

171 The goal of the object level segmentation step is to separate  
 172 each individual building from others. By doing so, the point  
 173 cloud of each building in the large scene can be processed  
 174 independently. In our work, this procedure focuses on three  
 175 categories, i.e., buildings, ground, and trees. We first describe  
 176 how the statistical information is obtained, then we exploit  
 177 graph cut to segment the points into the above three categories.

##### 178 4.1. Point features

179 Inspired by previous work [13, 32] that exploits geometric  
 180 features defined on single points to perform classification, our  
 181 approach relies on a statistical analysis of the neighborhoods of  
 182 the points.

183 In order to classify the points into the aforementioned  
 184 different categories, we first compute the statistical information  
 185 of the data based on a 2D supporting grid. Specifically, the  
 186 entire region of the the scene is discretized into a grid defined  
 187 on the ground plane using predefined grid resolution  $r_g$ . We  
 188 project all the points onto the grid and within each grid cell  
 189 we compute attributes for the points projected into this cell.  
 190 In the grid, each cell has the standard 4-connected neighbors.  
 191 An illustration of a 2D supporting grid is shown in Figure 2.  
 192 Note in the classification step, our goal is to extract individual  
 193 buildings and it is not necessary to extract precise contours  
 194 for buildings, thus we choose to use a larger (compared with  
 195 the one used for roof extraction described in Section 5.1) grid  
 196 resolution for the classification. Empirically, the grid resolution  
 197 is set to 0.35 m.

198 To extract discriminative features for classification, we  
 199 analyze the spatial distribution and structure of the points for

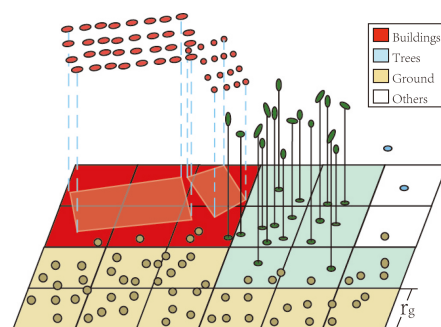


Figure 2: An illustration of the 2D supporting grid for object level segmentation.

200 each category based on the 2D supporting grid. Considering  
 201 different objects in an urban area often exhibit strong structural  
 202 regularities, e.g., buildings often exhibit planar regions, sharp  
 203 corners, and axis aligned dominant planes; points of trees  
 204 have more random distribution for both positions and normal  
 205 directions; the ground plane is usually regarded as a single large  
 206 segment that is relatively planar and low in height, allowing  
 207 for it to be identified separately. Our point cloud classification  
 208 algorithm incorporates features defined by these observations.

209 We introduce an identification function  $F(\cdot)$  that measures  
 210 the probability of a grid cell  $c_i \in C$  belonging to one of these  
 211 three categories. Similar to [13], our identification function  
 212 is defined on a set of features extracted from the point set, as  
 213 below:

- 214 • the maximum height of the cell from the ground:  $h_i =$   
 215  $\max\{p_i \rightarrow z\} - z_{ground}$
- 216 • the standard deviation of the absolute value of z compo-

- 217 ment of the normal vectors in a cell:  $\sigma_{Nz}$
- 218 • the standard deviation of the height of the points in a cell:
  - 219  $\sigma_H$

Then, the normalized identification function  $F(\cdot)$  is defined as

$$\begin{aligned} F_{ground} &= \max(1 - h_i/\bar{h}, 0) \\ F_{building} &= \min(\max(h_i/\bar{h} - \gamma \cdot \sigma_{Nz}, 0), 1) , \\ F_{tree} &= \min(\alpha \cdot \sigma_H + \beta \cdot \sigma_{Nz}, 1) \end{aligned} \quad (1)$$

220 where  $z_{ground}$  denotes the elevation of the ground plane;  $\bar{h}$   
221 is a threshold such that a point is considered belonging to a  
222 building if its height from ground is higher than  $\bar{h}$ ;  $\alpha$  and  $\beta$  are  
223 weights that balance between the elevation feature and the point  
224 distribution. Since the heights of trees in the experimented areas  
225 are less than  $6m$  and  $\sigma_{Nz}$  ranges from 0.01 to 0.4, we set  $\bar{h} = 6m$ ,  
226  $\gamma = 3$ ,  $\alpha = 0.03$ , and  $\beta = 3$  through all our experiments.  
227 Intuitively, a value of  $F(f_c)$  closer to 1 means the cell in the  
228 grid has higher possibility to be assigned the label  $f_c$ , and vice  
229 versa.

#### 230 4.2. Point classification

231 As the point cloud is discretized and embedded into a  
232 uniform 2D grid, the goal of the classification is to classify the  
233 corresponding cells into different categories. This classification  
234 is a typical labeling problem. We compute an assignment of  
235 labels  $f_c$  to elements  $c \in C$  such that the joint labeling  $f$   
236 minimizes an objective function  $E(f)$ . Our energy function  
237 consists of two terms: data and smoothness costs.

238 **Data cost.** The data cost  $D(c, f_c)$  measures how well the label  
239 assignment fits to the cells  $C$ . The normalized identification  
240 functions  $F(\cdot)$  provide the initial labeling estimation for all the  
241 cells. We define the data cost for each category as follows

$$D(c, f_c) = \begin{cases} 1 - F_{ground} & \text{if } f_c = \textit{ground} \\ 1 - F_{building} & \text{if } f_c = \textit{building} \\ 1 - F_{tree} & \text{if } f_c = \textit{tree} \end{cases} . \quad (2)$$

242 **Smoothness cost.** The smoothness term measures the spatial  
243 correlation of neighboring cells. Given two adjacent elements  
244  $p$  and  $q$ , the smoothness energy term is defined by

$$V_{p,q} = \frac{1}{\gamma \cdot |h_p - h_q| + 1} \cdot \mathbb{1}(p, q), \quad (3)$$

245 where  $\mathbb{1}(p, q)$  is an indicator function that has value 0 if  $p$  and  $q$   
246 are signed the same label, otherwise it has value 1. Intuitively,  
247 the smoothness term penalizes assigning different labels to a  
248 pair of adjacent cells  $(p, q)$  that have smaller difference in their  
249 heights, i.e.,  $|h_p - h_q|$ . For all the examples shown in the paper,  
250  $\gamma$  is set to 10.

**Optimization.** Thus the overall energy function is

$$E(f) = \sum_{c \in C} D(c, f_c) + \lambda \sum_{p,q \in N} V_{p,q}. \quad (4)$$

251 Finding a solution to this labeling problem is equivalent  
252 to the minimization of the above energy function. In our

253 implementation, we use graph cut [30, 31] to find the optimal  
254 labeling assignment. Compared with previous point cloud  
255 classification methods [32, 13] that use features defined on local  
256 neighborhood of the points, our statistic based classification  
257 can obtain more reliable results especially for complex scenes  
258 with higher level of noise and is more consistent with human  
259 perception.

#### 260 4.3. Object segmentation

261 Since our final goal is to reconstruct buildings exhibited in  
262 the scene, we perform a segmentation step that aggregates and  
263 extracts individual buildings using a simple label based region  
264 growing algorithm.

265 We first extract buildings, trees, and ground by querying  
266 and combining neighboring cells that have the same label  
267 assigned in the previous classification step. The remaining  
268 points are more likely distributed in small regions with irregular  
269 geometries, thus are classified into the fourth category (i.e.,  
270 *others*). Using the features defined in Section 4.1, some tall  
271 objects (e.g., wire poles) may be misclassified as *building*.  
272 We filter out these false positives using a simple thresholding  
273 mechanism. In our implementation, if the 2D area of a  
274 projected object labeled as building contains less than 200 grid  
275 cells (i.e.,  $24.5 m^2$ ), the object is then assigned as *others*. A  
276 point set of building may still contain some points that may  
277 belong to ground or other categories, but these outliers are only  
278 restricted within no more than one cell outward of the contour  
279 of the building structure. So these outliers will have little effect  
280 on the final reconstruction.

### 281 5. Polygonal Mesh Extraction

282 Given the point clouds of individual buildings separated from  
283 the scene, our next goal is to reconstruct mesh models from  
284 these point clouds. Automatic reconstruction is challenging due  
285 to the following two reasons. First, point clouds reconstructed  
286 from images using SfM and SVM are usually nonuniform and  
287 contain a higher level of noise compared with laser scans.  
288 Second, missing data is an unavoidable problem during the data  
289 acquisition process due to occlusions, lighting conditions, and  
290 the trajectory planing of the UAV. We observe that in the point  
291 clouds generated from aerial images the walls of the buildings  
292 are extremely sparse and incomplete if the trajectory for the  
293 UAV are not carefully designed, while roofs are relatively  
294 denser and more complete than the walls. Thus, quite a  
295 few previous work mainly utilize only roof information for  
296 reconstruction [15, 16, 17]. These methods are either based on  
297 region growing for roof extraction [16, 17], or detection of roof  
298 contours by measuring certain point features (e.g., [15]), thus  
299 they suffer difficulties caused by noise and missing data. In this  
300 work, we propose a regularized MRF formulation to extract the  
301 roof structure of the building, followed by a refinement step for  
302 the roof contours. Finally, building models are extruded from  
303 the roof patches.

304 Compared with previous graph-cut based approaches for  
305 surface reconstruction [33, 13, 14, 11], where their formulations

are based on either the irregular graph of the Delaunay tetrahedron, or points, or triangulated meshes, our formulation makes use of a graph with a four-neighbor grid structure in 2D space. Thus, our strategy significantly simplifies the roof extraction process, resulting in better stability and efficiency.

### 5.1. Roof extraction

In order to reliably extract roof structures, we employ another MRF-based segmentation algorithm on a grid with higher resolution defined on the point set of the building. This grid is similar to the one used in the previous classification stage, but with smaller cells that ensure more details of the roof structures can be recovered. Another difference is that in the new grid each cell stores an elevation value of the local points projected into the cell. Thus this grid can also be regarded as a depth map that provides us an effective way to process the data. With the depth map representation, processing can be conducted efficiently and effectively on the depth map despite the imperfections of the data.

Before extracting the roof structures from the point set, it would be helpful to reduce the noise in the data. To this end, we run a median filter on the depth map, since median filters are well known for reducing noise and outliers, and meanwhile preserve features (i.e., edges) of the data.

After the 2D filtering preprocess, we generate a set of plane hypotheses from the depth map using RANSAC [34]. Thus each cell in the depth map is assigned with an initial hypothesis label. Performing RANSAC on the depth map is extraordinarily efficient as the depth map significantly reduces the amount of data and maintains sufficient 2.5D information of the point cloud. We now perform a global optimization over all the cells in the depth map, to consistently segment the depth map into a set of planar regions (including roofs and the ground). This optimization is formulated as a cell-wise labeling problem, which is similar to the one used in the previous classification step (see Section 4.2). Our objective function still has a data cost term and a smoothness cost term.

**Data cost.** The data cost term  $D(p, f_p)$  encodes the likelihood of assigning a label  $f_p$  to a cell  $p \in P$ . It is defined as the distance measured from  $p$  to the corresponding plane with label  $f_p$

$$\begin{aligned} D(p, f_p) &= \text{dist}(p, f_p) \\ &= \mathbf{x}_p \cdot \mathbf{n}_f + D. \end{aligned} \quad (5)$$

where  $\mathbf{x}_p$  is the position of cell  $p$  in the grid,  $\mathbf{n}_f$  is the normal vector of the plane, and  $D$  is the constant coefficient in the plane equation denoted as  $Ax + By + Cz + D = 0$ .

**Smoothness cost.** Smoothness cost term  $V_{p,q}$  penalizes the assignment of two different labels to adjacent cells  $p$  and  $q$ , and thus encourages the coherence between neighboring cell pairs:

$$V_{p,q} = \begin{cases} 0 & \text{if } l_p = l_q \\ \delta_1 & \text{if } l_p \neq l_q, l_p = l_{\text{ground}} \text{ or } l_q = l_{\text{ground}}, \\ \delta_2 & \text{otherwise} \end{cases} \quad (6)$$

where  $l_{\text{ground}}$  denotes the ground plane. The penalty term  $\delta_1$  is a constant term that makes the penalty robust to region

boundaries.  $\delta_2$  is another penalty term defined as the distance between the projected points on different planes

$$\delta_2 = \left\| \text{proj}_i(p), \text{proj}_j(p) \right\|_2,$$

where  $\text{proj}(p)$  is the projection of the point on the corresponding plane.

**Optimization.** By combining the above two terms, the overall energy is defined similar to that in Equation 4:

$$E(f) = \sum_{p \in P} D(p, f_p) + \mu \sum_{p,q \in N} V_{p,q}, \quad (7)$$

where  $P$  is the cells set and  $N$  represents the standard 4-neighborhood. Parameter  $\mu$  is a weight that balances between the two terms. To optimize the above energy, we use the same graph cut algorithm used in Section 4.2.

### 5.2. Contour refinement and model extraction

Minimizing the above energy defined in Equation 7 will decompose the depth map of a building into a set of roof patches (see Figure 3). In our experiments, we observed that directly optimizing Equation 7 tends to generate zigzag artifacts in the roof contours.

Considering planar and orthogonal structures are common in architecture (i.e., most building structures are aligned with three dominant principal directions), we add a rotation step before the grid structure is built. Specifically, we detect two dominant directions by analyzing the normals of the original point set, and then transform the point cloud such that these directions are aligned with the  $X$  and  $Y$  axes of the 2D coordinate system. Experiments show that the alignment of the grid with the 2D coordinate system significantly helps to eliminate the zigzag artifacts in the extracted roof patches. Figure 3 shows the extracted roofs after the rotation step.

To extrude polygonal models from the roof patches, we first extract straight line segments from the roof contours obtained in the previous process. Specifically, we divide each roof contour into the following two categories of small contours according to the contents linked by the contours.

- Building boundaries: a *roof* patch is on one side of the contour and a *ground* patch is on another side;
- Roof boundaries: patches on both sides of the contour are of *roof*.

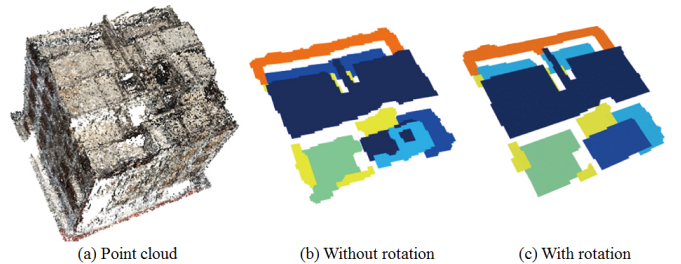


Figure 3: Roof extraction without and with the rotation step. Color denotes different roof patches.

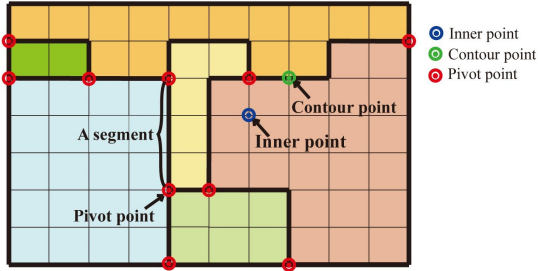


Figure 4: An illustration of pivot points and contour segments. The colors represent different roof patches.

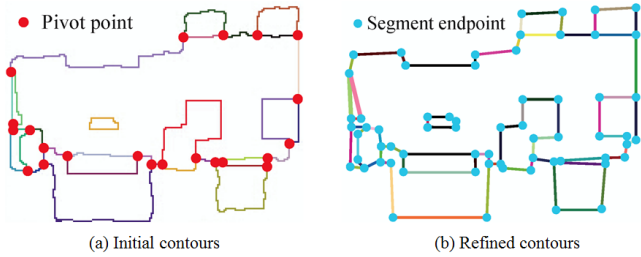


Figure 5: Contour simplification using the Douglas-Peucker polygonal approximation algorithm [35].

377 Since each cell in the grid has been assigned with a roof  
 378 patch, pivot points are detected by checking the number of roof  
 379 patches associated with the junctions in the grid. Specifically,  
 380 a junction in the grid is considered as a pivot point if the cells  
 381 associated with this junction belongs to at least 3 different roof  
 382 patches (see Figure 4).

383 We linearize, and thus simplify the contours of roof  
 384 patches using the Douglas-Peucker polygonal approximation  
 385 algorithm [35]. This algorithm decomposes the contours into a  
 386 sequence of straight line segments by recursively finding a point  
 387 that has the maximum distance to the simplified segments, and  
 388 this point is discarded if it is closer than a threshold  $\epsilon$  to the  
 389 approximating segments. The recursion is continuing until no  
 390 more points can be found that have distances greater than  $\epsilon$  to  
 391 the simplified segments. In our experiment, we set  $\epsilon$  to 0.2 m.  
 392 Figure 5 shows an example of contour simplification results.

393 In the end, we finish the whole pipeline by constructing a  
 394 polygonal mesh from the refined contours. Specifically, we  
 395 construct a 2D polygon for each boundary loop in the contours,  
 396 and further extrude them to the ground plane by adding vertical  
 397 walls that are orthogonal to the roofs and the ground plane. The  
 398 result is 2.5D reconstruction of the building in the scene. By  
 399 performing the same processing on each individual buildings,  
 400 then entire scene is reconstructed.

## 401 6. Results and Discussion

402 We have tested our approach on several datasets of large  
 403 scenes acquired by a high resolution camera mounted on an  
 404 UAV. After the images are obtained, we generate colored point  
 405 clouds from these images using SfM and MVS. Since SfM

406 and MVS are based on local image features, the computed  
 407 point clouds usually suffer from serious noise, occlusions, and  
 408 nonuniform densities. We then reconstruct polygonal models  
 409 from the point clouds using our proposed method. Figures 6  
 410 and 7 show the reconstruction of these scenes.

411 Figure 6 shows a portion of the United Nations Educational  
 412 Scientific and Culture Organization cultural heritage site of  
 413 Al-Balad, Jeddah, Saudi Arabia. The area scanned by the  
 414 UAV consists of many unique 100-300 year old buildings  
 415 with complex architectural features, cluttered rooftops, lattice  
 416 shuttered windows, and balconies. In the last fifty years, the  
 417 cultural heritage site has lost over 600 historical buildings  
 418 and within even the last several months homes have been  
 419 destroyed by accidental fire. We were given special permission  
 420 to scan the area due to its endangerment with the intent to  
 421 document the remaining buildings and generate a master plan  
 422 of the area. This study will help in digitizing the remaining  
 423 375 buildings that would be too time-consuming to do using  
 424 manual methods. The dataset consists of 1,518 images captured  
 425 during three 10-minute autonomous flights with a Sony QX100  
 426 camera (20M pixels) and 24mm (equivalent lens) achieving a  
 427 ground sampling density of 2.5 – 3.0cm per pixel. The three  
 428 flights were repeated over the same area at an elevation of 50m  
 429 (oblique), 75m (oblique), and 75m (nadir). The total generated  
 430 point cloud contains 20 million colored points. Although a  
 431 dense point cloud was generated a higher frequency of noise  
 432 especially in low-feature surface areas (e.g., windows, white  
 433 walls, metallic surfaces, etc.) was created. Our method benefits  
 434 from the statistical analysis of the imperfect point cloud, which  
 435 compensates the low quality of the data in an excellent way.  
 436 As can be seen from this figure, although the roofs of the  
 437 buildings are noisy and have missing regions, our method  
 438 successfully detected and reconstructed all buildings in these  
 439 regions, resulting in crack-free models.

440 Figure 7 shows a portion of a large modern residential area  
 441 consisting of a mix of two story homes with garage ports  
 442 and multiple balconies, multi-story apartment buildings, and  
 443 a residential park. Two 15-minute flights were conducted to  
 444 capture the entire area (approximately  $125 \times 100 m^2$ ) using a  
 445 larger UAV with a Sony Nex-7 camera (24M pixels) mounted  
 446 on a gimbal angled at  $50^\circ$  achieving a ground sampling  
 447 density of 1cm per pixel. The dataset consists of 924 images  
 448 and the point cloud generated from these images contains  
 449 approximately 80 million colored points. The reconstructed  
 450 polygonal models fit the initial point cloud in a precise manner,  
 451 and significantly reduce the storage. By converting the point  
 452 cloud representation into polygonal models, the storage of the  
 453 scene is reduced from 1.2 GB (initial point cloud with color, in  
 454 binary format) to 530 KB (polygonal model).

455 **Robustness to parameters.** Our MRF formulations for the  
 456 object level point cloud segmentation and roof extraction relies  
 457 on two key parameters:  $\lambda$  and  $\mu$ . In our experiments, we  
 458 found the final reconstruction results are not sensitive to these  
 459 parameters.

460 Figure 8 demonstrates the object level segmentation results  
 461 for a historical downtown scene with increasing value of  
 462 parameter  $\lambda$ . As can be seen from this figure, smaller values of  $\lambda$

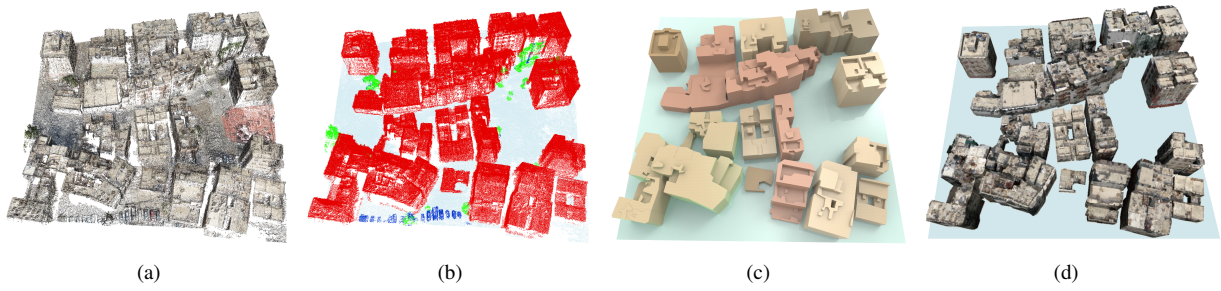


Figure 6: Segmentation and reconstruction of an old downtown area. (a) Initial point cloud; (b) Object level segmentation result; (c) Polygonal models reconstructed by the proposed method; (d) Textured polygonal models.

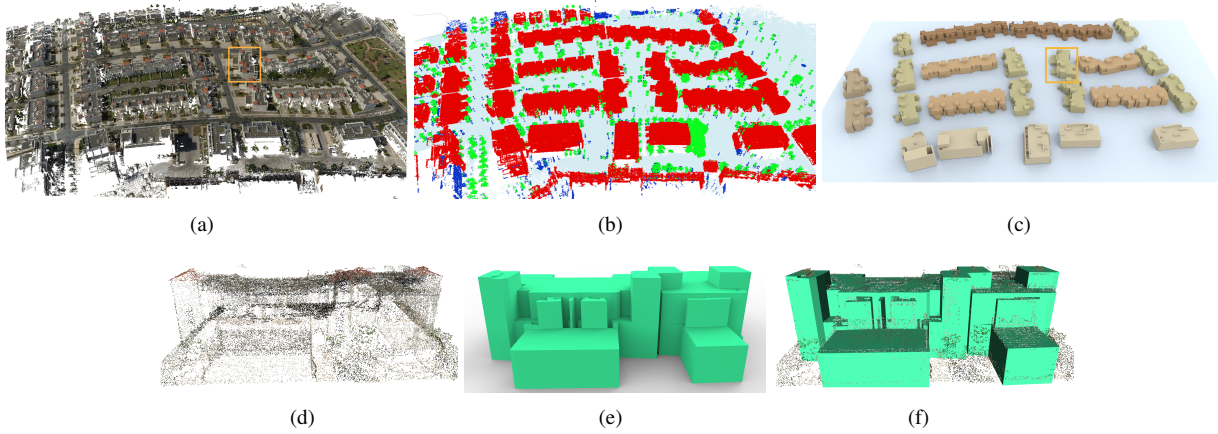


Figure 7: Segmentation and reconstruction of a large modern residential area. (a) Initial point cloud; (b) Object level segmentation result; (c) Polygonal models reconstructed by the proposed method; (d), (e), and (f) are the zoomins of the marked building in the scene.

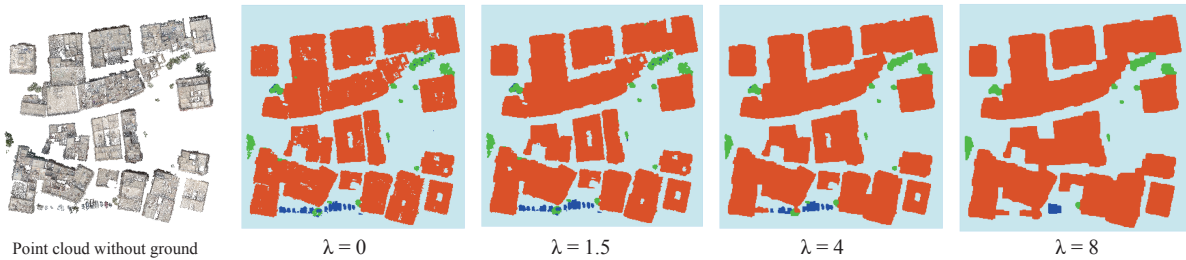


Figure 8: The effect of varying parameter  $\lambda$  (in Equation 4) on the segmentation results. Red, green, and blue colors represent *building*, *tree*, and *others* respectively.

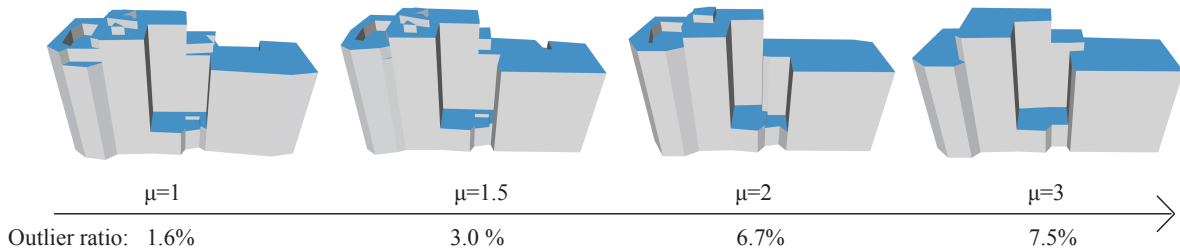


Figure 9: The effect of varying parameter  $\mu$  (in Equation 7) on the final reconstruction results. Here outlier ratio is defined as the percentage of points whose distances to the 3D model are larger than  $0.8m$ . Blue color represents the building roofs.

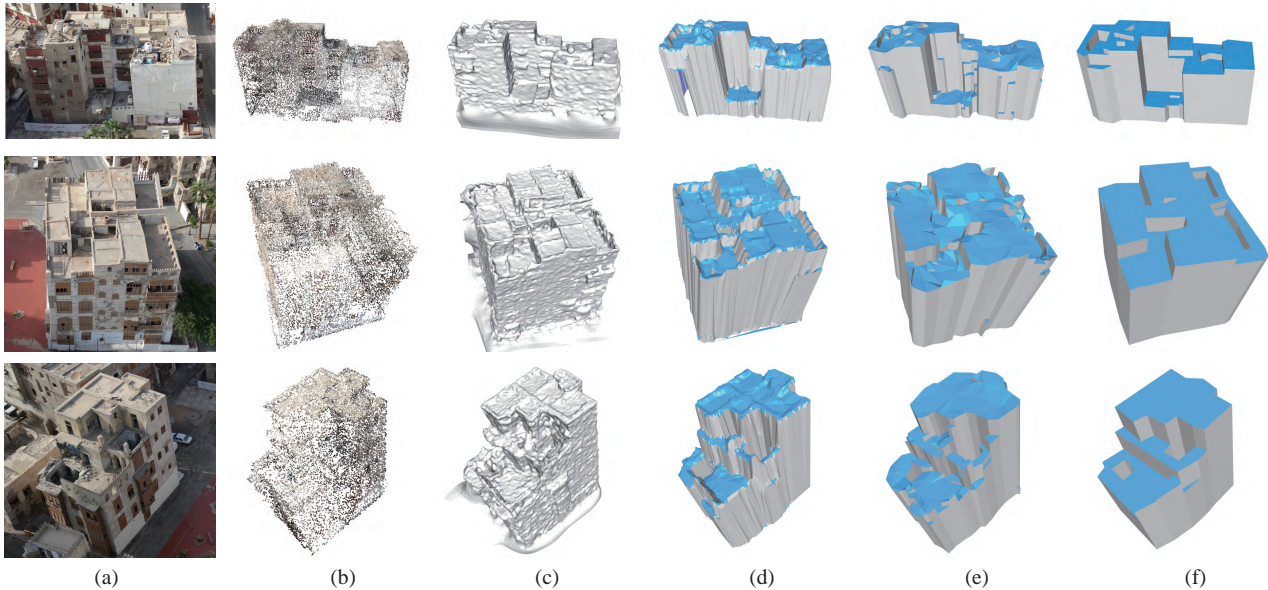


Figure 10: Comparison of the reconstruction results of our approach with other methods on three individual buildings (in different rows). (a) Photo of the building; (b) Point cloud; (c) Surface model reconstructed by Screened Poisson Reconstruction algorithm [36]; (d) DEM simplification result [37]; (e) Result of 2.5D Dual Contouring method [18]; (f) Our result.

463 ignore more smoothness constraints, which results in more gaps  
 464 and holes in the segmentation results due to noise and missing  
 465 data. On the contrary, increasing the value of  $\lambda$  will encourage  
 466 close segments to be merged. However, as our experiments  
 467 demonstrate, the value of  $\lambda$  in the range [1.4, 4.3] can guarantee  
 468 similar satisfactory segmentation results.

469 In Figure 9, we demonstrate the robustness of our roof  
 470 extraction algorithm on the final reconstruction result in terms  
 471 of varying parameter  $\mu$ . Similar to the effect of varying  $\lambda$   
 472 in the object level segmentation step,  $\mu$  controls how much  
 473 smoothness constraints are preferred in the energy function.  
 474 Intuitively, increasing the value of  $\mu$  encourages larger planar  
 475 roofs in the final reconstruction. Our experiments reveal that  
 476 the value of  $\mu$  in a range of [1.5, 3.0] usually generates similar  
 477 compact 3D models.

478 **Comparison.** We also conduct comparisons with three  
 479 methods: Surface simplification from Digital Elevation Model  
 480 (DEM) [37], Screened Poisson Reconstruction [36], and 2.5D  
 481 Dual Contouring [18].

482 Figure 10 shows the reconstruction results of three individual  
 483 buildings. The results of the DEM simplification method are  
 484 competitive in terms of fitting quality to the point clouds.  
 485 However, it can not produce straight roof boundaries. The  
 486 Screened Poisson Reconstruction method [36] can generate  
 487 an isotropic dense mesh surface from the point clouds. This  
 488 method, however, can not handle local incompleteness (i.e.,  
 489 holes in the point clouds) caused by occlusions. Besides,  
 490 since the result is represented as a single surface approximating  
 491 the entire scene, it is rather difficult to differentiate individual  
 492 buildings in the reconstruction. The results from the 2.5D  
 493 Dual Contouring [18] method contain large areas of small  
 494 bumps. This is because the 2.5D Dual Contouring algorithm

495 is initially designed to deal with airborne LiDAR point clouds  
 496 that mainly consist of points of building roofs with uniform  
 497 density and higher accuracy. Thus, it is sensitive to our noisy  
 498 point clouds computed from images using SfM and MVS.  
 499 Compared with these approaches, our method can generate  
 500 a simplified polygonal model that is visually pleasing and  
 501 satisfactory for various applications or can be used as input for  
 502 further processing.

503 In Table 1, we show a quantitative comparison with the  
 504 aforementioned methods on the buildings shown in Figure 10.  
 505 As can be seen from this table, the Screened Poisson  
 506 Reconstruction method wins in terms of precision, but the  
 507 final surfaces are more fluctuating. Our method has similar  
 508 accuracy as the 2.5D Dual Contouring method, but it has a  
 509 more compelling performance and our results have the simplest  
 510 geometric structure. Our approach is seeking a tradeoff between  
 511 accuracy and automatic reconstruction.

512 Furthermore, we also run our method on LiDAR point cloud  
 513 data provided by [18]. As shown in Figure 11, our method also  
 514 can deal with LiDAR data and can obtain a similar compact  
 515 reconstruction results as the primitive-based method proposed  
 516 in [13].

517 **Accuracy and scalability.** To intuitively evaluate the  
 518 accuracy of the reconstructed models, we show the overlay of  
 519 the point clouds onto the polygonal models in Figure 12, where  
 520 color coding indicates the error magnitude. Our method has an  
 521 average fitting error less than 0.2 m for the scene.

522 Besides the individual buildings, the experiments also  
 523 demonstrate that our reconstruction framework has satisfactory  
 524 performance on large scenes (see Figures 6 and 7). We record  
 525 the running times for these scenes, which can be seen in Table 2.  
 526 Both the object level segmentation and roof extraction for the



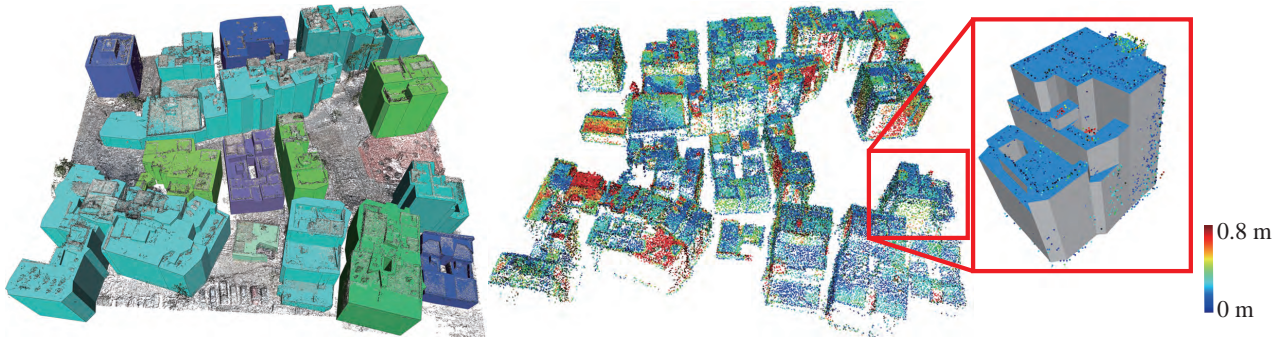


Figure 12: Point cloud overlaid on the reconstructed models. Color indicates the distances from points to their nearest faces in the model.

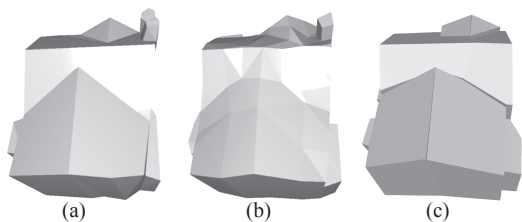


Figure 11: A comparison of our method with the primitive-based method proposed in [13] on a LiDAR point cloud. (a) The model obtained by [13]; (b) 2.5D Dual Contouring result [18]; (c) Our result.

Table 1: Statistical comparison of running times (in seconds), mesh sizes (face number), and mean errors (in meters, defined as the average distance of the points to the model) of our method with 2.5D Dual Contouring [18] (2.5D for short) and Screened Poisson reconstruction [36] (SPR) methods on the buildings shown in Figure 10.

		2.5D [18]	SPR [36]	Ours
Figure 10 (top) 20.6 k points	Time	0.31	7.66	0.40
	# Faces	2,250	56,344	675
	Error	0.076	0.068	0.096
Figure 10 (middle) 29.6 k points	Time	0.28	8.42	0.38
	# Faces	3,599	110,287	387
	Error	0.086	0.044	0.053
Figure 10 (bottom) 13.9 k points	Time	0.17	1.97	0.19
	# Faces	1,382	66,968	207
	Error	0.093	0.103	0.106

two scenes take only a few seconds. Thus, our method is quite suitable for processing large scale urban environments.

**Limitations.** During the reconstruction, we mainly rely on the roof information of the buildings. We assume that there is only one flat ground in each scene and the roofs are parallel to the ground plane. Since the models are obtained by extruding prisms from the ground plane to the roofs, the reconstructed buildings always lie in the same ground plane, and they are actually 2.5D reconstructions. Given point clouds with vertical facades, our current formulation simply ignores these vertical facade information and only uses the information given by the roof points.

Another limitation is that the piecewise planar roof structure

Table 2: Running times (in seconds) of the two core steps (object level segmentation and roof extraction) for the two large scenes shown in Figure 6 and Figure 7.

	Segmentation	Roof extraction
Figure 6	2.36	3.36
Figure 7	15.01	6.92

assumption becomes too restrictive when dealing with atypical architectures, e.g., buildings with curved roofs or facades. Currently our method can not handel these types of buildings.

## 7. Conclusions and Future Work

This paper presented an automatic framework for reconstructing large scale urban scenes from UAV images. We introduce an effective segmentation algorithm which segments the data based on statistical analysis of their geometric properties using a low resolution grid structure. Roofs are extracted and their contours are simplified and refined using a similar grid structure of higher resolution. By using the proposed MRF formulations on the statistical information, our method is able to handle a higher level of noise and outliers. Experiments on various scenes show the reconstructed polygonal models are more compact and regular compared with state-of-the-art methods.

Currently, we only use the roof information for the reconstruction. Although the walls are sparse, they do provide extra constraints on the geometry of the buildings. As a future work, we would like to exploit the wall information to regularize the roof extraction algorithm. Another interesting problem could be approximating the trees in the scenes using template models from a database.

## Acknowledgements

We thank the anonymous reviewers for their valuable comments and suggestions. This work was supported by the KAUST Visual Computing Center.

## 567 References

- 568 [1] N. Haala, M. Kada, An update on automatic 3d building reconstruction,  
569 ISPRS Journal of Photogrammetry and Remote Sensing 65 (2010) 570–  
570 580.
- 571 [2] P. Musialski, P. Wonka, D. G. Aliaga, M. Wimmer, L. van Gool,  
572 W. Purgathofer, A survey of urban reconstruction, in: Computer Graphics  
573 Forum (STAR Proceedings of Eurographics), 2013.
- 574 [3] F. Rottensteiner, G. Sohn, M. Gerke, J. Wegner, U. Breitkopf,  
575 J. Jung, Results of the isprs benchmark on urban object detection and 3d  
576 building reconstruction, ISPRS Journal of Photogrammetry and Remote  
577 Sensing 93.
- 578 [4] M. Berger, A. Tagliasacchi, L. M. Seversky, P. Alliez, J. A. Levine,  
579 A. Sharf, C. Silva, State of the art in surface reconstruction from point  
580 clouds, in: S. Lefebvre, M. Spagnuolo (Eds.), Eurographics 2014 - State  
581 of the Art Reports, 2014.
- 582 [5] S. Agarwal, Y. Furukawa, N. Snavely, I. Simon, B. Curless, S. M. Seitz,  
583 R. Szeliski, Building rome in a day, Communications of the ACM 54 (10)  
584 (2011) 105–112.
- 585 [6] Y. Furukawa, J. Ponce, Accurate, dense, and robust multiview stereopsis,  
586 IEEE Transactions on Pattern Analysis and Machine Intelligence 32 (8)  
587 (2010) 1362–1376.
- 588 [7] C. Wu, S. Agarwal, B. Curless, S. M. Seitz, Multicore bundle adjustment,  
589 in: Proceedings of the 2011 IEEE Conference on Computer Vision and  
590 Pattern Recognition, CVPR '11, IEEE Computer Society, Washington,  
591 DC, USA, 2011, pp. 3057–3064. doi:10.1109/CVPR.2011.5995552.  
592 URL <http://dx.doi.org/10.1109/CVPR.2011.5995552>
- 593 [8] P. Müller, G. Zeng, P. Wonka, L. Van Gool, Image-based procedural  
594 modeling of facades, in: ACM SIGGRAPH 2007 Papers, SIGGRAPH  
595 '07, 2007.
- 596 [9] S. N. Sinha, D. Steedly, R. Szeliski, M. Agrawala, M. Pollefeys,  
597 Interactive 3d architectural modeling from unordered photo collections,  
598 ACM Transactions on Graphics (TOG) 27 (5) (2008) 159:1–159:10.
- 599 [10] J. Xiao, T. Fang, P. Tan, P. Zhao, E. Ofek, L. Quan, Image-based façade  
600 modeling, in: ACM Transactions on Graphics (TOG), Vol. 27, ACM, New  
601 York, NY, USA, 2008, pp. 161:1–161:10. doi:10.1145/1409060.1409114.  
602 URL <http://doi.acm.org/10.1145/1409060.1409114>
- 603 [11] I. Garcia-Dorado, I. Demir, D. G. Aliaga, Automatic urban modeling  
604 using volumetric reconstruction with surface graph cuts, Computers &  
605 Graphics 37 (7) (2013) 896–910.
- 606 [12] H. Lin, J. Gao, Y. Zhou, G. Lu, M. Ye, C. Zhang, L. Liu, R. Yang,  
607 Semantic decomposition and reconstruction of residential scenes from  
608 lidar data, ACM Transactions on Graphics (TOG) 32 (4) (2013) 66:1–  
609 66:10.
- 610 [13] F. Lafarge, C. Mallet, Building large urban environments from  
611 unstructured point data, in: Computer Vision (ICCV), 2011 IEEE  
612 International Conference on, Barcelona, Spain, 2011, pp. 1068–1075.  
613 doi:10.1109/ICCV.2011.6126353.
- 614 [14] F. Lafarge, P. Alliez, Surface reconstruction through point set structuring,  
615 Computer Graphics Forum 32 (2) (2013) 225–234.
- 616 [15] Q.-Y. Zhou, U. Neumann, Fast and extensible building modeling from  
617 airborne lidar data, in: Proceedings of the 16th ACM SIGSPATIAL  
618 international conference on Advances in geographic information systems,  
619 ACM, 2008, p. 7.
- 620 [16] C. Poullis, S. You, Automatic reconstruction of cities from remote  
621 sensor data, in: Computer Vision and Pattern Recognition, 2009.  
622 CVPR 2009. IEEE Conference on, IEEE, 2009, pp. 2775–2782.  
623 doi:10.1109/CVPR.2009.5206562.
- 624 [17] F. Lafarge, X. Descombes, J. Zerubia, M. Pierrot-Deseilligny, Structural  
625 approach for building reconstruction from a single dsm, IEEE  
626 Transactions on Pattern Analysis and Machine Intelligence 32 (1) (2010)  
627 135–147. doi:10.1109/TPAMI.2008.281.
- 628 [18] Q.-Y. Zhou, U. Neumann, 2.5d dual contouring: A robust approach to  
629 creating building models from aerial lidar point clouds, in: Proceedings  
630 of the 11th European Conference on Computer Vision Conference on  
631 Computer Vision: Part III, ECCV'10, 2010, pp. 115–128.
- 632 [19] Q.-Y. Zhou, U. Neumann, 2.5d building modeling with topology control,  
633 in: CVPR, IEEE Computer Society, 2011, pp. 2489–2496.
- 634 [20] Q.-Y. Zhou, U. Neumann, 2.5d building modeling by discovering global  
635 regularities, in: CVPR, IEEE Computer Society, 2012, pp. 326–333.
- 636 [21] L. Nan, A. Sharf, H. Zhang, D. Cohen-Or, B. Chen, Smartboxes for  
637 interactive urban reconstruction, ACM Transactions on Graphics (TOG)  
638 29 (4) (2010) 93:1–93:10. doi:10.1145/1778765.1778830.  
639 URL <http://doi.acm.org/10.1145/1778765.1778830>
- 640 [22] C. A. Vanegas, D. G. Aliaga, B. Benes, Automatic extraction of  
641 manhattan-world building masses from 3d laser range scans, IEEE  
642 Transactions on Visualization & Computer Graphics 18 (10) (2012)  
643 1627–1637.
- 644 [23] M. Pollefeys, D. Nistér, J. M. Frahm, A. Akbarzadeh, P. Mordohai,  
645 B. Clipp, C. Engels, D. Gallup, S. J. Kim, P. Merrell, C. Salmi,  
646 S. Sinha, B. Talton, L. Wang, Q. Yang, H. Stewenius, R. Yang, G. Welch,  
647 H. Towles, Detailed real-time urban 3d reconstruction from video, Int. J.  
648 Comput. Vision 78 (2-3) (2008) 143–167.
- 649 [24] M. Arikan, M. Schwärzler, S. Flöry, M. Wimmer, S. Maierhofer, O-  
650 snap: Optimization-based snapping for modeling architecture, ACM  
651 Transactions on Graphics (TOG) 32 (1) (2013) 6:1–6:15.
- 652 [25] L. Nan, C. Jiang, B. Ghanem, P. Wonka, Template assembly for detailed  
653 urban reconstruction, Eurographics 2015, Computer Graphics Forum  
654 34 (2).
- 655 [26] A. Akbarzadeh, J.-M. Frahm, P. Mordohai, B. Clipp, C. Engels,  
656 D. Gallup, P. Merrell, M. Phelps, S. N. Sinha, B. Talton, L. W.  
657 0002, Q. Yang, H. Stewenius, R. Yang, G. Welch, H. Towles, D. Nistr,  
658 M. Pollefeys, Towards urban 3d reconstruction from video, in: 3DPVT,  
659 IEEE Computer Society, 2006, pp. 1–8.
- 660 [27] V. Hiep, R. Keriven, J. P. P. Labatut, Towards high resolution large-scale  
661 multi-view stereo, Miami, US, 2009, pp. 1430–1437.
- 662 [28] Y. Verdier, F. Lafarge, P. Alliez, Lod generation for urban scenes, ACM  
663 Transactions on Graphics (TOG) 34 (3) (2015) 15.
- 664 [29] C. Wu, Visualsfm: A visual structure from motion system, URL:  
665 <http://homes.cs.washington.edu/~ccwu/vsfm> 9.
- 666 [30] Y. Boykov, V. Kolmogorov, Computing geodesics and minimal surfaces  
667 via graph cuts, in: Proceedings of the Ninth IEEE International  
668 Conference on Computer Vision, Vol. 2 of ICCV '03, 2003, p. 26.
- 669 [31] Y. Boykov, V. Kolmogorov, An experimental comparison of min-cut/max-  
670 flow algorithms for energy minimization in vision, IEEE Transactions on  
671 Pattern Analysis and Machine Intelligence 26 (9) (2004) 1124–1137.
- 672 [32] M. Carlberg, P. Gao, G. Chen, A. Zakhor, Classifying urban landscape  
673 in aerial lidar using 3d shape analysis, in: Image Processing (ICIP), 16th  
674 IEEE International Conference on, IEEE, 2009, pp. 1701–1704.
- 675 [33] P. Labatut, J.-P. Pons, R. Keriven, Robust and efficient surface  
676 reconstruction from range data, Computer Graphics Forum 28 (8) (2009)  
677 2275C2290. doi:10.1111/j.1467-8659.2009.01530.x.
- 678 [34] R. Schnabel, R. Wahl, R. Klein, Efficient ransac for point-cloud shape  
679 detection, Computer Graphics Forum 26 (2) (2007) 214–226.
- 680 [35] D. H. Douglas, T. K. Peucker, Algorithms for the reduction of the  
681 number of points required to represent a digitized line or its caricature,  
682 Cartographica: The International Journal for Geographic Information and  
683 Geovisualization 10 (2) (1973) 112–122.
- 684 [36] M. Kazhdan, H. Hoppe, Screened poisson surface reconstruction, ACM  
685 Transactions on Graphics (TOG) 32 (3) (2013) 29:1–29:13.
- 686 [37] G. Priestnall, J. Jaafar, A. Duncan, Extracting urban features from lidar  
687 digital surface models, Computers, Environment and Urban Systems 24  
688 (2000) 65–78.