Block World Reconstruction from Spherical Stereo Image Pairs

Hansung Kim and Adrian Hilton

Centre for Vision, Speech and Signal Processing University of Surrey, Guildford, Surrey, GU2 7XH, UK h.kim@surrey.ac.uk

Abstract

We propose a block-based scene reconstruction method using multiple stereo pairs of spherical images. We assume that the urban scene consists of axis-aligned planar structures (Manhattan world). Captured spherical stereo images are converted into six central-point perspective images by cubic projection and façade alignment. Depth information is recovered by stereo matching between images. Semantic regions are segmented based on colour, edge and normal information. Independent 3D rectangular planes are constructed by fitting planes aligned with the principal axes of the segmented 3D points. Finally cuboid-based scene structure is recovered from multiple viewpoints by merging and refining planes based on connectivity and visibility. The reconstructed model efficiently shows the structure of the scene with a small amount of data.

Keywords:

3D reconstruction, Scene modelling, Spherical imaging, Block world interpretation

1 1. Introduction

² 3D scene reconstruction from photographic images has been an important research topic for various domains. Applications include visual sets in film and game production, 3D map generation, virtual tourism and urban planning. There have been many studies into outdoor scene reconstruction from multi-view images [1, 2, 3]. Strecha et al. created a benchmarking site for the quantitative evaluation of algorithms against ground-truth by LIDAR scanning [4]. However, the quality of pure image-based reconstruction largely
depends on the capture environment.

Firstly, real environments include complex appearance causing errors in reconstruction from images. Textureless and non-Lambertian surfaces often result in errors in matching and reconstruction. Scenes reflected on glass or water induce false depth. Moving pedestrians and cars in the scene can be occluders in urban scene modelling.

Another problem is that normal cameras with a limited field-of-view 15 (FOV) capture only a partial observation of the surrounding environment. 16 Reconstruction of a complete model of the 3D environment requires addi-17 tional views to capture the scene and occluded regions. Reconstruction of 18 scene models from multiple images or video acquired with a standard cam-19 era has been the focus of considerable research. However, the limited FOV 20 presents a challenging problem to ensure complete scene coverage for recon-21 struction. Agarwal st al. [5] reconstructed full 3D street models from 150,000 22 photos from the internet using grid computing. Pollefeys et al. [6] used 3,000 23 video frames to reconstruct one building and 170,000 frames for a small town. 24 The relatively narrow FOV and low resolution of normal cameras require ac-25 quisition and processing of large image sets for scene model reconstruction. 26

Finally, conventional dense reconstruction methods such as LIDAR scans 27 or image-based reconstruction result in millions of points with a high-level 28 redundancy which do not efficiently represent the scene structure. The task 29 of extracting a structured representation for subsequent visualisation is typ-30 ically performed manually. When we applied our previous dense reconstruc-31 tion algorithm [7] for datasets covering areas of 30m diameter surrounded by 32 buildings, it produced more than 100 million faces with 60 million vertices. 33 This occupies huge amount of system memory and may require out of core 34 techniques [8] to visualise and render. Applications such as 3D structure rep-35 resentation and pre-visualisation require scene models in a structured form 36 for efficient storage, transmission and rendering. 37

Piecewise-planar, plane-based and block-based scene modelling methods provide a good solution for the above problems. These approaches start from the assumption that man-made environments such as urban areas or building interiors are composed of piecewise planar surfaces. Furukawa et al.[9] and Gupta et al.[10] used the strong assumption of a piecewise-axis-aligned-planar world (Manhattan world).

We previously presented a dense environment model reconstruction [7] and a plane-based reconstruction [11] using a line-scan camera and manual segmentation. In this paper, we propose an automatic block-based environment model reconstruction method based on the same input data. This produces a more complete scene model with a compact representation for storage and transmission. The geometry can be refined for higher resolution mesh models if dense depth information is available. The approach provides a compact scene model for hierarchical geometry representation of the detailed scene structure.

⁵³ The main contributions of this paper are:

- We propose a 3D block-based scene reconstruction system. This is a
 simple and efficient way to represent the structure of a scene with high
 completeness for transmission and interactive visualization.
- Spherical stereo imaging enables full scene reconstruction with a small
 number of input images. This saves considerable time in scene capture
 and reconstruction.

We propose a façade alignment algorithm to find regions in the scene
 for optimal alignment and cubic projection. Cubic projection decompose the spherical image into six central-point perspective images. The
 central-point perspective image is advantageous in feature matching
 and 3D plane reconstruction because it is distortion-free and has a
 vanishing point at the centre of the image aligned with the principal
 axes for a Manhattan world.

- We propose an automatic extraction of plane and cuboid structure
 from colour and depth images. Optimal block-based representation of
 the scene is recovered based on visibility, occupancy, point density and
 physical stability.
- We provide an optional user interaction to constrain primitive recon struction to keep specific geometrical details or refine erroneous regions.
- High resolution texture mapping from the original images to the block
 based representation gives a quick rendering of the scene.

The rest of this paper is organised as follows: Section 2 introduces related previous works and Section 3 outlines overview of the proposed method. Section 4 presents capture method and cubic projection with façade alignment.



Figure 1: Categories of Simplified scene modelling methods

Depth reconstruction and region segmentation methods are proposed in Section 5. In Section 6, we introduce plane primitives reconstruction and structured block reconstruction methods. Experimental results and discussion are
given in Section 7, and Section 8 makes conclusions of this work. Supplemental video is also available at: http://www.cvssp.org/hkim/BlockWorld/
BlockRecon-CVIU.mov showing results of reconstruction for various scenes.

⁸⁴ 2. Related Work

Simplified scene modelling has been a long-standing area of research. Previous approaches can be separated into two categories: interactive and fully automatic methods. The automatic method are divided into grammarbased and matching-based approaches according to the registration strategy and the matching-based approach uses various input modalities as illustrated in Fig. 1

The FAÇADE system introduced by Debevec et al.[12] pioneered interactive environment modelling from images. In this approach, a simplified geometric model of the architecture is recovered interactively with manual correspondence using multiple view geometry. Novel views are rendered using view-dependent texture mapping, and additional geometric detail is recovered automatically through stereo correspondence. Their research was ⁹⁷ commercialised as ImageModeler¹.

Hengel et al. [13] proposed an interactive 3D modelling method from video 98 frames by tracing the shape of objects in the scene. They used structure from gc motion (SfM), feature point tracking and superpixel segmentation to get 3D 100 information from 2D video frames. If users draw 2D primitives such as lines 101 and circles on frames, then the system automatically builds 3D primitives 102 from the user's input and reconstructed 3D information. This concept was 103 extended by Sinha et al. [14] using feature-matching and SfM methods with 104 line detection and vanishing point detection algorithms for interactive 3D 105 architectural modelling from photo collections. SketchUp² provides a simple 106 3D reconstruction tool from multiple photos. This is similar to the Sinha's 107 method but it does not use any matching method, just manual vanishing 108 point alignment for photo registration to 3D coordinates. This tool is useful 109 to build very simple scenes but has limitations in building complex scenes 110 because it requires manual matchings for each primitive. 111

Automatic scene reconstruction can be divided into two categories: grammar-112 based and matching-based reconstruction. Grammar-based reconstruction 113 uses semantic region detection and recognition to compose the world ac-114 cording to pre-defined rules. Gupta et al.[10] proposed block world recon-115 struction from a single outdoor image, inspired by the "Blocks World" work 116 in the 1960's and Hoiem et al.'s "pop up 3D" [15]. They assume that the 117 world is composed of blocks and match 2D image regions into 3D block view 118 classes. They also estimate the density of each block using visual cues and 119 use it to generate 3D parse graphs which describe geometric and mechanical 120 relationships between objects within an image. Muller et al. [16] proposed 121 a rule-base city modelling method using shape grammar rules from façade 122 images, now commercialised as CityEngine³. Xiao et al.[17] proposed an 123 automatic approach to generate street-side 3D photo-realistic models from 124 images captured along streets at ground level with an assumption that build-125 ing façades have two principal directions. They use a SfM method for the 126 initial point cloud reconstruction and apply a multi-view semantic segmenta-127 tion method for classifying regions into semantic models in the hand-labelled 128 image database. Then, independent blocks are reconstructed using major 120

¹ImageModeler, http://usa.autodesk.com/adsk/servlet/pc/index?id= 11390028&siteID=123112

²SketchUp, http://www.sketchup.com/

³City Engine, http://www.esri.com/software/cityengine/

line structures and the final facade scene is modelled by inverse patch-based 130 orthographic composition and structure analysis. This approach generates 131 clean facade scenes but is highly computationally expensive, taking 23 hours 132 on a cluster of 15 computers for semantic segmentation of 202 building blocks. 133 Bellotti et al. [18] proposed an Architectonic Style Area (ASA) algorithm for 134 procedural generation of buildings in an urban area, based on the concept 135 of "architectonic likelihood". The algorithm accepts facade pictures from 136 sample buildings and statistical description of the elements and styles as in-137 put, and composes facade models by statistically assembling sample images 138 of architectonic components. Components are classified in an ontology based 139 on the classic principles of architecture. The algorithm relies on rules that 140 encode the semantics of the ontology. Simon et al. [19] proposed a grammar-141 based modelling method with basic shapes (roof, wall, window, balcony, floor, 142 door, shop, etc.) and deviation tree for the procedural geometry. Mathias 143 et al. [20] proposed a similar grammar-driven approach for reconstruction 144 of buildings and landmarks, but they used an inverse procedural modelling 145 strategy for SfM and image-based analysis. Satkin et. al [21] present a data-146 driven approach using repositories of 3D models to find the identities, poses 147 and styles of objects in a scene. However, the grammar-based approach has 148 a serious problem because semantic segmentation is not always stable, and 149 this approach works only within the given rule and categories. Any object 150 or building out of the given categories induces errors in reconstruction. 151

Multi-View Stereo (MVS) and SfM reconstruction is the most popular ap-152 proach, not only in full geometry reconstruction, but also in piece-wise planar 153 reconstruction. Schindler et al.[22] proposed a novel method for recovering 154 the 3D-line structure of a scene from multiple widely separated views. 2D 155 lines aligned to major axes are detected by EM-based vanishing point es-156 timation. Those 2D lines are reconstructed as 3D lines to provide guide 157 lines for 3D structure reconstruction. Hane et al. [23] proposed a piece-wise 158 planar depth map fusion, which formulates an energy term in stereo match-159 ing using patch-based priors to reconstruct piece-wise planar scenes. Sinha 160 et al. [24] suggested extracting vanishing directions and fitting point clouds 161 into 3D planes reconstructed based on the vanishing directions. Gallup et 162 al.[25] proposed a stereo method handling scenes containing both planar and 163 non-planar regions by segmentation and planar region detection. The planar 164 regions are represented by planes and the non-planar regions are modelled 165 by the results of a standard multi-view stereo algorithm. One problem of 166 this approach is the lack of completeness due to small independent planes. 167

Toldo et al. [26] proposed planar patch extraction based on photo consis-168 tency from point clouds using the J-linkage algorithm [27] and reconstructed 169 the scene with a view clustering tree and hierarchical reconstruction. Some 170 research has invoked the stronger Manhattan-world assumption [28] which 171 states that the world is piecewise planar and aligned to orthogonal axes. 172 Micusik et al. [29] proposed a super-pixel stereo on a Markov Random Field 173 (MRF) and aligned surfaces to three dominant directions based on the grav-174 ity vector and vertical vanishing point. Furukawa et al. [9] also built indoor 175 and outdoor scenes by axis aligned depth map integration relying on the 176 Manhattan world assumption. The approach starts from point clouds gener-177 ated by their Patch-based Multi-view Stereo (PMVS) algorithm [1] and finds 178 an optimal minimum volume solution with plane hypotheses. Compensat-179 ing or concealing the occlusion part of the scene is an important problem 180 in 3D reconstruction. Chauve et al.[30] used additional ghost primitives to 181 fill gaps between detected basic primitives by inducing cell complex. We use 182 a similar plane extension technique in this paper to detect intersections of 183 reconstructed partial planes. Kowdle et al.[31] proposed an active learning 184 technique. They used an energy minimization framework for piecewise pla-185 nar reconstruction but allowed simple user interaction to provide support for 186 the uncertain regions. 187

Some approaches reconstruct geometry from point cloud datasets gener-188 ated by an active sensor such as LIDAR without the help of image data. City 189 modelling from aerial scans is one typical example. Zhouet al.[32] proposed 190 a method to produce crack-free models composed of complex roofs and verti-191 cal walls from aerial LIDAR point clouds. Poullis et al. [33] also developed a 192 fully automatic method for extracting high-fidelity geometric models directly 193 from aerial LIDAR scans using 2D roof boundaries extraction based on GMM 194 and camera pose estimation using Levenberg-Marquardt optimisation. Li et 195 al.[34] introduced an idea for modelling algorithm from range data that ex-196 ploits a priori knowledge that buildings can be modelled from cross-sectional 197 contours using extrusion and tapering operations. Nguatem et al. [35] pro-198 posed an automatic cuboid fitting algorithm using a line sweep to recon-199 struct cuboid-based building model from point clouds. Xiao et al. [36] also 200 developed a virtual walkthrough system with regularized texture-mapped 3D 201 model using an inverse constructive solid geometry for large indoor scenes 202 from ground-level photographs and 3D laser points. 203

The use of spherical imaging provides a simple approach to overcome the limited FOV of conventional cameras. Sturm[37] suggested a method for 3D

plane-based scene reconstruction from a single panoramic image. He used a 206 priori constraints on the 3D structure such as: co-planarity of points, per-207 pendicularity of planes and lines, and parallelism of planes and lines. Kang 208 et al. [38] also proposed a similar 3D plane reconstruction method using the 209 normal vector of plane and vanishing points from a single panoramic image. 210 Point Grey developed an omnidirectional multi-camera system, the Lady-211 bug⁴, consisting of six XGA color CCDs to provide high resolution spherical 212 images. Micusik et al. [39] used this camera for piecewise planar city mod-213 elling. They back-projected images to quadrangular planes and applied MRF 214 superpixel stereo and depth sweeping algorithms for depth map reconstruc-215 tion. The reconstructed depth maps were fused into surfaces aligned to three 216 dominant directions. They assumed that the cameras are pre-calibrated and 217 that reference images are also pre-segmented. Google also developed their 218 own omnidirectional multi-camera system to reconstruct and render street 210 models [40]. They simultaneously utilised range sensor to obtain a base-220 structure of the street scene and refined the model with optical flow esti-221 mation from captured images. In their approach, accurate registration of 222 photometric and geometric information is important. Simultaneous sensing 223 from different locations requires calibration and registration to align depth 224 and image information. 225

Instead of omnidirectional or panoramic images, Feldman et al.[41] used the Cross Slits (X-Slits) projection with a rotating fisheye camera to generate a high quality spherical image and to reduce the dimension of the plenoptic function. In this research we use a similar line-scan camera to capture latitude-longitude image which has advantage in stereo matching and 3D reconstruction.

232 3. Overview of the Proposed System

In this research, we propose a simple and efficient method to reconstruct a simplified structured environment model from spherical image pairs. Figure 2shows a block diagram for the whole process.

A linescan camera captures a full surrounding scene at multiple locations as vertical stereo pairs. The captured images are latitude-longitude images. They are projected into a unit cube with a novel façade alignment algorithm

⁴Pointgrey, http://ww2.ptgrey.com/spherical-vision



Figure 2: Block diagram of the system

based on the Hough transform. Each face image of the cubic projection is a 239 distortion-free central-point perspective image whose three principal axes are 240 aligned to vertical, horizontal, and the image centre directions, respectively. 241 To reconstruct depth information from stereo pairs, disparity estima-242 tion is performed. For automatic initial region segmentation, we propose a 243 graph-based region segmentation extending Felzenszwalb and Huttenlocher's 244 algorithm[42] to perform segmentation based on colour, normal direction and 245 detected Hough lines. From the segmentation and disparity maps, inde-246 pendent 3D rectangular planes are constructed by plane fitting. The plane 247 structure is refined by merging, expanding, cropping and eliminating planes 248 validated against the reliability, visibility and occupancy. 249

Finally connected planes are extruded in the counter normal direction to construct block models. An optimal block structure is recovered based on the point density in each cuboid. The result represents the scene structure as a set of cuboids which can be used to render the scene with texture mapping.

²⁵⁴ 4. Line-scan Capture and Cubic Projection

255 4.1. Spherical stereo acquisition

In this work, we use a commercial off-the-shelf line-scan camera⁵ with 256 a fisheye lens in order to capture the full environment as a high resolution 257 spherical image. This camera samples rays on a hemisphere about the centre 258 of projection and stitches together from the rotating slits together to form 259 a new image. The camera rotates about axis through its optical centre. As 260 a result the imaging geometry of the line-scan capture can be regarded as 261 conventional perspective projection, and the result is a latitude-longitude 262 image like a world map. 263

 $^{^5}$ Spheron, https://www.spheron.com/products.html

In order to recover depth information from the images, the scene is captured with the camera at two different heights. This vertical stereo line-scan camera capture has the following advantages:

(1) Relatively simple calibration is required. Depth reconstruction only requires knowledge of the baseline distance between the stereo image pair and correction of radial distortion in the vertical direction. Radial distortion is rectified using a 1D lookup table to evenly map pixels on the vertical central line to the $[0, \pi]$ range. Lens distortion parameters are fixed so that this mapping can be calculated for the lens in advance.

273

(2) Stereo matching can be simplified to a 1D search along the vertical scan 274 line as discussed above, while normal spherical images require a complex 275 search along conic curves or rectification of the images. In the latitude-276 longitude geometry, the great circles intersecting at the epipoles of the 277 spherical geometry become parallel straight lines. Therefore, the con-278 ventional correlation-based matching on an 1D search range can be used 279 to compute the disparity of spherical stereo images if they are vertically 280 aligned. Error in the alignment can be corrected by rectification using 281 the method proposed by Banno and Ikeuchi [43]. 282

283

(3) High resolution images can be captured by a line-scan camera because
the sensor array is 1D and the resolution about the axis depends on the
step size. High resolution images provide more accurate depth estimation
and high quality texture mapping.

288 4.2. Cubic projection and Façade alignment

The latitude-longitude images can be directly used for 3D reconstruction. 289 However, we propose to convert the image into distortion-free perspective 290 images via projection of the spherical image to a cube, referred to here as 291 cubic projection. The cubic projection projects all pixels on the unit sphere 292 to the unit cube in the range [-1, 1] in each axis. The converted image is 293 decomposed into six perspective images as illustrated Fig. 3 (a) and (b). 294 We set 0° of longitude as the x-axis in the cubic projection of Fig. 3 (b) 295 and each side face image of the cubic projection has two vanishing points. 296 These images have a single-vanishing point if we set the axes of the cubic 297 projection to be aligned to the Manhattan world axes, which we refer to as 298 façade alignment. 299



(d) Facade aligned projection image (Carpark1)



Façade alignment is the process of matching the main façades in the 300 scene to be perpendicular to the principal axes of the cubic projection by 301 rotating the spherical image around the vertical axis. Fig. 3 (c) shows the 302 projection result when 126° of the longitude is set on the x-axis. Fig. 3 (d) is 303 another example of the façade aligned cubic projection. We can observe that 304 the horizontal and vertical lines in the scene were aligned to horizontal and 305 vertical directions in each images, respectively, and that the lines aligned to 306 the depth direction converges to the image centre. Therefore we can consider 307 these images as formed by central-point perspective projection. Central-308 point perspective projection has significant advantage for axis-aligned plane 309 reconstruction. Most of the current plane-based reconstruction algorithms 310



Figure 4: Block world reconstruction from façade-aligned cubic projection image

³¹¹ use vanishing point and principal directions detection in 3D space [9, 24, 29].
³¹² Cubic projection with the façade alignment can detect 3D principal directions
³¹³ in 2D images. Therefore, aligned 3D planes or 3D blocks can be built by
³¹⁴ extruding detected 2D planes in the depth direction as shown in Fig. 4.

In order to find the most reliable angle shift t_{opt} to set the x-axis, we considers the number, sparseness, average length and average angle errors of image lines resulting from the probabilistic Hough transform [44]. The angle shift in longitude is equivalent to the horizontal pixel shift in the line-scan image. We detect the following three kinds of Hough lines as aligned among all detected lines H: Horizontal Hough lines H_h , Vertical Hough lines H_v and Perspective Hough lines (to the depth direction) H_p .

322 -
$$H_h = \{h_h | h_h \in H, |\theta(h_h)| < 1^\circ$$

323 324

-
$$H_v = \{h_v | h_v \in H, |\theta(h_v) - 90^\circ| < 1^\circ\}$$

• -
$$H_p = \{h_p | h_p \in H, D(I_c, h_p) <$$

where θ is the angle of the line to the horizontal direction, and D(I, L) is the distance between the image centre point I_c and the line L. The distance threshold r_c varies depending on the image resolution. In the above aligned Hough lines, H_h are most important to detect the façades of the scene and H_v are almost the same over any shift. We estimate the optimal rotational shift t_{opt} to find the façade direction of the scene by maximising the following energy term:

 r_c



Figure 5: Facade energy according to angle shift

$$t_{opt} = \underset{0^{\circ} \le t < 360^{\circ}}{\operatorname{argmax}} E_F(t) \tag{1}$$

$$E_F(t) = \lambda_R E_R(t) + \lambda_S E_S(t) + \lambda_L E_L(t)$$
(2)

$$E_R = \frac{|H_h \cup H_v \cup H_p|}{|H|}$$
$$E_S = \sigma_{H_h}$$
$$E_L = \frac{1}{|H_h|} \sum_{H_h} \log(l(h_h)/\theta(h_h) + \varepsilon))$$

In Eq.(2), E_R represents the ratio of the number of aligned Hough lines 332 to all Hough lines. E_S is the standard deviation of average y-position of H_h 333 which relates to sparseness of the horizontal Hough lines. We give higher 334 priority to sparse features in the scene because dense Hough lines can be 335 detected from small areas which have complicate patterns and bias the op-336 timisation. E_L represents the magnitude and accuracy of the detected H_h 337 where l(h) is the length of the line. The weighting factors λ_R , λ_S and λ_L 338 can be adjusted according to the scene characteristics, but we fix $\lambda_R=1.0$, 339

 $\lambda_S = 1.0, \lambda_L = 0.3$ throughout our experiments to show that those parameters are applicable to general scenes.

Figure 5 (a) shows $E_F(t)$ against all angle shifts from 0° to 360°. The 342 Cathedral scene shows four distinctive peaks at around 90° intervals, while 343 the CarPark scene has an ambiguous peak around 100° because the ground 344 is slightly slanted and the high frequency texture of the tree and brick walls 345 induce many outliers in the Hough transform. Using the assumption that 346 façades in each side face are perpendicular to each other based on the Man-347 hattan world assumption, we detect the optimal shift by maximising the 348 following energy sum which results in Fig. 5 (b) with a more distinct peak 349 point. 350

$$t_{opt} = \underset{0^{\circ} \le t < 90^{\circ}}{\operatorname{argmax}} \sum_{k=0}^{3} E_F(t + k * 90^{\circ})$$
(3)

The façade-aligned cubic projection image is a distortion-free central-351 point perspective image. It has several advantages over alternative projec-352 tions. First, it is easy to extract axis-aligned planes from the image because 353 it does not require any vanishing point detection. Second, it is easy to find 354 matched features between multi-view images because they do not have dis-355 tortion of the appearance while the spherical images have serious radial dis-356 tortion according to the angle. Finally, multi-view registration is a simple 357 3 DOF problem (only translation) because the façades direction is already 358 aligned for all views. 359

³⁶⁰ 5. Depth Reconstruction and Region Segmentation

³⁶¹ 5.1. Depth reconstruction from spherical stereo

One of the most important problems in depth estimation is locating cor-362 responding points in the images, a process referred to as disparity estima-363 tion. The estimated disparity fields can be converted into depth information 364 by camera geometry. Depth reconstruction from images captured by con-365 ventional cameras require a calibration step to extract camera parameters. 366 However, the spherical stereo pair and cubic projection pair used in this re-367 search do not require a complex calibration step because pixel positions in 368 each image directly correspond to 3D spherical coordinates as described in 369 Section 4.1. 370



Figure 6: Spherical and cubic stereo geometry



(a) Disparity map

(b) Depth map with cubic projection

Figure 7: Depth reconstruction result for Cath2 image set

The angle disparity d between two image pairs is defined as illustrated 371 in Fig. 6. If we assume the angles of the projection of the point P onto 372 the spherical or cubic projection image pair displaced along the y-axis are 373 θ_t and θ_b respectively, then the angle disparity d of point $p_t(x_t, y_t)$ can be 374 calculated as $d(p_t) = \theta_t - \theta_b$. The distance of the scene point P from the two 375 cameras is calculated by triangulation as Eq. (4), where B is the baseline 376 distance between the camera's center of projection and r_t and r_b represent 377 the distance from P to the top and bottom cameras. 378

$$r_{t} = B / \left(\frac{\sin \theta_{t}}{\tan(\theta_{t} + d)} - \cos \theta_{t} \right)$$

$$r_{b} = B / \left(\cos \theta_{b} - \frac{\sin \theta_{b}}{\tan(\theta_{b} - d)} \right)$$
(4)

Stereo matching can be carried our in either the spherical image pair or cubic projection image pairs. In the latitude-longitude image, the epipolar line for correspondence search is a vertical scan line. In cubic projection images, epipolar lines are vertical lines for side faces, and radial lines from

the centre for the top and bottom face images. Both image types have a 383 trade-off according to the disparity estimation method. Disparity estima-384 tion on the latitude-longitude images is good for pixel-base approaches or 385 global optimisation, but contains errors in area-based approaches like block 386 matching because of the distortion of the image. On the other hand, cubic 387 projection images show better results in area-based matching but correspon-388 dence should be independently estimated for each face image and requires 389 boundary processing between face image. In any case, disparity estimation 390 results can easily be converted between formats by the projection geometry 391 in Fig. 6. 392

Any disparity estimation algorithm can be used for the proposed system as long as it does not produce too many outliers. We use latitude-longitude images and a PDE-based variational disparity estimation method previously proposed to generate accurate disparity fields with sharp depth discontinuities for surface reconstruction [7].

Figure 7 shows the result of the estimated angle disparity field and its depth map followed by cubic projection. Depth is mapped to grey scale according to their disparity or depth range.

401 5.2. Region segmentation for plane reconstruction

Felzenszwalb and Huttenlocher [42] proposed a simple and intuitive segmentation concept that: "The intensity differences across the boundary of two regions are perceptually important if they are large relative to the intensity difference inside at least one of the regions". We modify Felzenszwalb's segmentation method [42] to embrace 3D features such as surface normal direction and aligned Hough lines.

⁴⁰⁸ A graph G = (V, E) is constructed for each face image domain, where ⁴⁰⁹ $v_i \in V$ is the set of pixels and $(e_{ij}) \in E$ is the edge between neighbouring ⁴¹⁰ elements (v_i, v_j) with a weight $w(e_{ij})$. We set the affinity weights according ⁴¹¹ to the colour difference W_c , face normal angle difference W_o and edge penalty ⁴¹² W_e in Eq. (5).



(c) Felzenszwalb

Figure 8: Region segmentation results(Cathedral2 and Carpark1)

$$W = W_c + \lambda_o W_o + W_e$$

$$W_c = dist(I(v_i) - I(v_j))$$

$$W_o = |cos^{-1}(\overline{O(v_i)} \cdot \overline{O(v_j)})|$$

$$W_e = \begin{cases} a & \text{if } v_i \text{ is on Aligned Hough lines,} \\ b & \text{else if } v_i \text{ is on Canny edge lines,} \\ 0 & \text{otherwise.} \end{cases}$$
(5)

I(v) is a colour value in (R, G, B) coordinates and O(v) is a surface normal 413 vector calculated from the depth map generated in section in 5.1. We experi-414 mentally set λ_o as 40 and edge penalties a and b as 400 and 200, respectively, 415 in all of our experiments. We also set the region size preference parameter k416 in the Felzenszwalb's algorithm as 1200. 417

Figure 8 shows results of region segmentation for the main façade of 418 the Cathedral and Carpark scenes. Figure 8 (a) shows the surface normal 419 map projected into the (R, G, B) domain and Fig. 8 (b) shows Canny edge 420 and Aligned Hough lines. Fig. 8(c) and (d) are results of Felzenszwalb's 421 segmentation algorithm and the proposed algorithm respectively. We can 422

observe that the side walls and objects are clearly segmented by the proposed 423 algorithm owing to the integration of surface normal and edge information. 424 However, the results are still over-segmented for plane reconstruction. Mi-425 cusik et al. [29] used Felzenszwalb's segmentation for generating super-pixels 426 and refined them iteratively using a 3D MRF. This method is computation-427 ally expensive to refine the super-pixel segmentation into meaningful copla-428 nar regions and is unstable in many cases. Therefore we introduce a method 429 to merge segmented regions in the plane reconstruction stage by considering 430 their reliability and spatial relationship. 431

432 6. Block World Reconstruction

3D structured scene is reconstructed from the 2D images, region segments 433 and disparity information. First, 3D rectangular plane elements are con-434 structed by projecting segmented regions to 3D with the depth information. 435 The resulting 3D planes are merged, eliminated and connected to generate 436 the 3D plane structure of the scene. If multiple stereo reconstructions are 437 available, they are registered and refined into one complete structure. Finally 438 a block world model is generated from the plane structure by fitting cuboids. 439 We introduce a scene scale parameter S_c which represents the level of 440 detail in the scene reconstruction. S_c defines the minimum size of objects to 441 be reconstructed, and also to merge or eliminate less reliable planes. Small 442 S_c can reconstruct details of the scene but large coherent area can be divided 443 into small planes including erroneous pieces. Large S_c produces rough scene 444 structure with less pieces but may lose scene details. We set S_c as 0.8m for 445 outdoor scenes and smaller value of 0.2-0.4m for indoor scene according to the 446 preference for scene details (S_c values for indoor scenes are given in Section 447 7). However, applying a single scale parameter to the whole scene can miss 448 important details. We allow user interaction as an option to introduce hard 449 constraint to specific regions to keep their properties in reconstruction. 450

451 6.1. Plane reconstruction

All 2D points $V_p \subset V$ in each segment can be projected into 3D space with the depth information to form a 3D point cloud. Rectangular planes are constructed from the segments and point clouds. Two different approaches can be used for plane detection in a point cloud segment: total least squares [45] and RANSAC-based [46]. The total least squares fitting is a form of linear regression and provides a solution to the problem of finding the best

Table 1: Plane classification				
Class	Constraint			
X plane	$n_x \sim N(0, \sigma_{n_x}^2)$			
Y plane	$n_y \sim N(0, \sigma_{n_y}^2)$			
Z plane	$n_z \sim N(0, \sigma_{n_z}^2) \& t(z) \geq T_g$			
Ground plane	$n_z \sim N(0, \sigma_{n_z}^2) \& t(z) < T_g$			
Arbitrary plane	otherwise (eliminated)			

fitting 3D plane through a set of points, while the RANSAC-based plane 458 detection iteratively selects a small subset of points at random to fit a model 459 to that subset and remove outliers. The total least square method is fast and 460 converges to a single solution, but the result can be biased by outliers. The 461 RANSAC-based method can be more accurate if there are many outliers, but 462 it is computationally expensive. In our case, the plane reconstruction is a 463 large set of small problems, and the disparity estimation algorithm provides 464 smooth and accurate depth fields over the surface except near region bound-465 aries. Therefore, we exclude 10% of points close to region boundaries in the 466 segmentation and apply the total least squares (orthogonal regression) fitting 467 algorithm [45] and bounding box extraction. If the 3D point cloud is noisy, 468 the RANSAC-based approach can be applied as an alternative. 469

Once all regions are fitted to planes, they are categorised into five classes 470 (X, Y, Z, Ground and Arbitrary) according to the constraints with their 471 normal vectors **n** and centre point $\mathbf{t}(x, y, z)$ as shown in Table 1. T_g , set as 472 $S_c/2$, is a threshold to define *Ground* planes among Z-planes and $\sigma_{n_i}^2$ is the 473 variation from the ideal normal vector for a particular plane orientation. We 474 set $\sigma_{n_i}^2$ as 0.15 for the block world reconstructions that do not use arbitrary 475 planes in our experiments. All reconstructed planes are saved as a vector 476 list: 477

$$\mathbf{P} = \{p_i\} = \{[\mathbf{n}_i \, \mathbf{t}_i \, w_i \, h_i]\} \tag{6}$$

478 where w and h are the height and width of the plane.

Plane elements reconstructed from the region segments may include false planes and partial planes which can be merged into a larger plane. In order to refine those planes, we measure the following reliability factors for each plane: reconstruction confidence R_c , plane size R_s , distance from the camera R_d , distance between planes R_b and angle to the camera view direction R_{θ} . $R_c(p_i) = MSE(p_i) \tag{7}$

$$R_s(p_i) = w_i \times h_i \tag{8}$$

$$R_d(p_i) = \|O_c - \mathbf{t}_i\| \tag{9}$$

$$R_b(p_i, p_j) = \|p_i - p_j\|_d \tag{10}$$

$$R_{\theta}(p_i) = |\cos^{-1}(\overline{O_c \mathbf{t}_i} \cdot \mathbf{n}_i)| \tag{11}$$

 $MSE(p_i)$ is the mean squared error calculated in the plane fitting, O_c is the location of the camera that the plane belongs to in the unified coordinate system and $\|\cdot\|_d$ is the minimum distance between two planes.

⁴⁸⁷ Merging Similar Planes: Two neighbouring planes p_i and p_j are merged ⁴⁸⁸ into one plane and the bounding box is newly set to cover both regions if ⁴⁸⁹ they satisfy the following conditions.

490 - Two planes are in the same category 491 - $R_b(p_i, p_j) < S_c$ 492 - $\{R_s(p_i)_{new} < R_s(p_i)_{old}\} \cap \{R_s(p_j)_{new} < R_s(p_j)_{old}\}$

⁴⁹³ $R_s(p)_{old}$ is the original area of p and $R_s(p)_{new}$ is the new area to be extended ⁴⁹⁴ to build the merged plane. The third condition keeps structures with a big ⁴⁹⁵ hole such as bridge or door. The position of the new plane is set to the ⁴⁹⁶ position of the old plane with lower $R_c(p)$.

Elimination of Unreliable Planes: According to the observations in [7],
we assume that a plane is unreliable if it is too distant from the camera or
its angle to the camera is too big. Therefore, the plane is eliminated if it
satisfies the following conditions.

501 - {
$$R_s(p_i) < S_c^2$$
} \cup { $R_d(p_i) > d_{max}$ } \cup { $R_\theta(p_i) < \theta_{min}$ }

 d_{max} is set to 20m for outdoor scenes and 5m for indoor scenes. θ_{min} is set to 15°. Figure 9 shows the original planes reconstructed from the segments and their refinement results by merging similar planes and eliminating unreliable ones.

⁵⁰⁶ Plane Intersection Refinement: All plane-to-plane intersections are checked ⁵⁰⁷ if they have any intersection with each other in the extension range of $S_c/2$. If



Figure 9: Plane refinement (Cathedral2, Randomly coloured planes)

any intersection is found, the length of intersection and visibility are checked 508 to determine the type of intersection. If the intersection is larger than half 509 of the bigger plane, two planes are welded at the intersection to generate 510 a corner (Fig. 10 (a)). Otherwise, only the smaller plane stops growing at 511 the intersection to generate a T-junction (Fig. 10 (b)). If two planes al-512 ready have an intersection, residual parts are eliminated based on visibility 513 constraints [47] (Fig. 10 (c)). If the plane does not meet any intersection 514 during the extension in any direction, we keep the original boundary. Figure 515 10 (d) illustrates examples of the plane intersection refinement observed in 516 the cathedral dataset. 517

Filling gaps from self-occlusion: The scene captured from a single fixed 518 location inevitably has self-occlusions in the scene as illustrated in Fig. 11 519 (a). We adopt the minimum volume solution proposed by Furukawa et al. [48] 520 based on the Manhattan World assumption. The plane occluded at the rear 521 of the front plane is extended to the boundary of the orthogonal line-of-522 sight (LOS) from the surface normal direction as shown in Fig. 11 (b). 523 The occluded region perpendicular to the orthogonal LOS is compensated 524 by other viewpoints or the block reconstruction presented in the following 525 subsections as demonstrated in Fig. 11 (c). 526



(d) Examples (Cathedral2, Randomly coloured planes)Figure 10: Intersection refinement

527 6.2. Multiple-view plane structure reconstruction and texture mapping

The spherical camera is more advantageous in environment capture than 528 normal camera with a limited FOV as mentioned in Section 1. However, 529 spherical imaging does not capture the complete scene due to self-occlusion. 530 Simple cases can be compensated by the minimum volume solution, but there 531 is no way to get information for occluded regions behind any object from a 532 single viewpoint. Another problem of the single-view spherical capture is the 533 fact that the accuracy of the depth estimation is inversely-proportional to 534 the distance and the angle of surface normal. These problems can be over-535 come by captures from multiple viewpoints. Merging reconstructions from 536 multiple stereo pairs can be integrated into a common 3D scene structure. 537 Kim and Hilton^[7] proposed a mesh-fusion algorithm for dense meshes by 538 mesh registration and reliable surface selection by considering surface visi-539



Figure 11: Occlusion filling

⁵⁴⁰ bility, orientation and distance. They calculate 3D rigid transforms between
⁵⁴¹ viewpoints using 2D feature matching.

We propose a similar but much simpler and faster method based on the observations that: 1) the mesh registration is not optimised for rotation and translation (r, s), but only for translation s because the façade direction is already aligned for all viewpoints; 2) the surface reliability test is not applied for each vertex on the surface, but for the whole plane.

⁵⁴⁷ We use SURF feature matching [49] between captured images for different ⁵⁴⁸ stereo pairs. The resulting 2D matches are projected into 3D space with the ⁵⁴⁹ estimated depth field. However, these points are not reliable enough to be ⁵⁵⁰ used in registration because of SURF matching and depth estimation errors. ⁵⁵¹ We use a RANSAC-based least square minimisation for the following error, ⁵⁵² where *i* and *j* are corresponding matching points in the model set *m* and ⁵⁵³ reference set *M*, respectively.

$$E_t(s) = \sum_{(i,j)} \|m_i - (M_j + s)\|^2$$
(12)

554

Once all viewpoints are registered into a unified coordinate system, plane



(a) Multiple capture (b) Reconstructed planes

(c) Texture mapping

Figure 12: Multiple-view plane reconstruction

primitives are reconstructed independently for each viewpoint. All planes in the X, Y, and Z classes are then refined by the same method as the single view reconstruction.

UV mapping [50] is used for texture representation. If texture mapping 558 is required, all planes are subdivided into small regular triangles with their 559 corresponding vertices in the texture image so that the mapped texture is 560 not distorted. If the plane is merged from multiple viewpoints, the dominant 561 viewpoint is decided by comparing their camera view direction R_{θ} in Eq. (11) 562 and the texture is obtained from the dominant view image. Multiple blending 563 is not used because the blending result can be blurred or result in ghosting 564 artefacts due to the simplified geometry. 565

Figure 12 shows the result of plane reconstruction from multiple pairs of spherical stereo. The Cathedral scene was captured at three different locations and each reconstructions is merged into the central viewpoint. All major objects in the scene are reconstructed.



(a) Primitives (b) Min volume (c) Max volume

Figure 13: Cuboid reconstruction from two planes



Figure 14: Cuboid reconstruction from three planes

570 6.3. Block-world reconstruction

Plane-based reconstruction describes simplified scene structure. However, block-based visualisation can provide better perception of the scene with surface normal orientations motivated by Gupta et al.[10]. This provides a model with higher completeness and an efficient representation of the scene because each block has only six degrees of freedom (3D location and dimensions).

Here we propose a cuboid fitting method starting from plane primitives reconstructed in Section 6.2. As mentioned in Section 4.2, cuboid reconstruction can be considered as an outward extrusion process (counter surface normal direction) of each cubic projection face.

If a plane primitive is connected to other perpendicular plane primitives whose extrusion directions overlap and they have different boundary lengths in the weld junction, the volume of the cuboid is decided by the original 3D point density in the primitive regions. We define an discrete objective function $D_o(P)$ as the density of 3D points belonging to the region as in Eq. (13).



Figure 15: Examples of cuboid reconstruction (Cathedral, partial view)

$$D_o(P) = \frac{Number \ of \ 3D \ points \ inP}{Area \ of \ region \ P}$$
(13)

⁵⁸⁷ We start from the minimum volume for connected planes and check all ⁵⁸⁸ possible volumes made up from planes up to the maximum volume. The ⁵⁸⁹ objective function $D_o(P)$ is calculated for each volume and the volume with ⁵⁹⁰ the maximum value is the optimal cuboid.

Let us consider a simple example as illustrated in Fig. 13. From the 591 two plane primitives, we can consider two cases of cuboid reconstruction: 592 Minimum volume in Fig. 13 (b) based on the regions B and C, and Maximum 593 volume in Fig. 13 (b) based on the regions A, B and C. We compare $D_o(B \cup$ 594 C) and $D_o(A \cup B \cup C)$ to choose the volume with the higher density. In 595 the more complex case with three planes in Fig. 14, there are 18 different 596 cases between minimum and maximum volumes. We calculate $D_o(P)$ for 597 all possible volumes in the same way, and choose the case with the highest 598 $D_o(P)$. 599

If the planes are isolated, they are extruded to an initial depth d_{init} . In case of multi-view reconstruction, planes are extruded in the counter normal direction. In the extrusion process, cuboids can intersect each other. If any intersection is detected, the original plane primitives and larger objects take



Figure 16: Cathdral main building reconstructed with different scene scale parameters

⁶⁰⁴ priority over limiting the extrusion of the smaller object.

Finally the block structures are refined based on their physical stability [10]. There may be floating blocks which do not meet the ground plane due to occlusion or disparity errors in the automatic reconstruction. These blocks are physically unstable, violating the law of gravity. If any block is not supported by other stable blocks and is close to the ground plane, the block is extended to the ground to retain the physical stability.

Figure 15 shows examples of the cuboid reconstruction from plane primitives in the Cathedral scene.

613 6.4. Optional user interaction to constrain primitive reconstruction

The proposed pipeline is a fully automatic method from the image input to the block world reconstruction. However, there are two possible problems in applying this automatic pipeline to various environments. First, automatic plane primitives reconstruction can fail to build meaningful coplanar regions due to the errors in disparity estimation or region segmentation. Second, applying a single scene scale parameter S_C can miss important details in the scene geometry.

Geometrical details can be preserved by adjusting the scene scale parameter S_C in reconstruction, but this may result in an over segmented scene with cluttered geometry. Figure 16 shows the Cathedral main building reconstructed for different scale parameters S_C . The smaller S_C in Fig. 16 (a) includes geometric details such as the sculptures and window regions in the main façade, and the eaves of the side wings. However, it can not fully resolve the steps and there are holes in the façade.

In order to overcome these problems, we implemented a simple user interface as an option to constrain the segmentation step. The user can merge or split regions by scribbling and assign X,Y and Z-plane class to specific regions as a hard constraint so that these regions are not affected by the



(b) Highstreet

Figure 17: Hard-constrained primitive reconstruction introduced by user interaction. In both the Cathedral and Highstreet scene user-interaction required < 1 minute to introduce constraints on the reconstruction (Left: Fully-automatic, Right: With hard constraints)

region refinement step. It takes less than a minute for each image to define semantic constraints which are kept as independent clusters in the scene
reconstruction.

Figure 17 illustrates results of introducing hard constraints on primitive reconstruction with user interaction. We observe that the eaves and steps are reconstructed regardless of the large scene scale parameter in the Cathedral scene. The cluttered background is also simplified by adding constraints. In the Hightstreet scene, the cluttered fence region and erroneous shop window regions are reconstructed by similar user interaction.

The inclusion of simple user-interaction to constrain the reconstruction is left as an option according to the application requirements for full or semiautomatic scene modelling. All experimental results in the following sections are produced by the fully automatic process without user interaction.

645 7. Experimental Results

All scenes presented in this section were captured with a Spheron commercial line scan camera introduced in Section 4.1. We attached a Nikon 16mm f/2.8 AF fisheye lens to the system and captured vertical stereo pairs with a baseline of 60cm for outdoor scenes and 20cm for indoor scenes, respectively. The resolution of spherical images is 3143×1414 .

651 7.1. Evaluation against LIDAR ground-truth

The goal of our proposed approach is to reconstruct an approximate rep-652 resentation of the scene structures. Evaluation of geometric accuracy against 653 ground-truth scene geometry therefore only provides a partial measure. In 654 this section, we evaluated reconstruction results from test scenes against 655 ground-truth models from LIDAR scans to show how close the proposed ap-656 proach can represent the scenes. We compared accuracy and completeness of 657 representation with a dense reconstruction method represented in Kim and 658 Hilton^[7]. 659

Figure 18 shows the ground-truth from multiple LIDAR scans and the 660 reconstructed models from three viewpoints using the proposed algorithm. 661 The "Gate" scene has a width of 9m and a height of 6m. Stereo pairs are 662 captured with a baseline of 60cm and the scene scale parameter S_c is set to 663 0.2m. The reconstructed plane primitives represent the approximate struc-664 ture of the scene. Figure 19 also shows the ground-truth model from seven 665 LIDAR scans and reconstructions of the main building for the "Cathedral" 666 outdoor scene in Fig. 12. The main building has a width of 30m and a height 667 of 20m. 668

Accuracy (how close the reconstruction is to the ground-truth) and com-669 pleteness (how much of the ground-truth is modelled by the reconstruction) 670 are measured based on the evaluation methodology proposed in Seitz et 671 al. [51]. Both ground-truth and reconstructions are incomplete, therefore 672 we considered only subset regions of the target model in measuring error 673 distance. The reconstruction and ground-truth are registered in the same 674 coordinate frame. Then Hausdorff distance from each vertex in the source 675 model to the closest point in the target model is calculated. Accuracy is 676 measured by the RMS error from the reconstruction to the ground-truth, 677 and completeness is measure by the ratio of vertices in the ground-truth 678



whose closest points to the reconstruction exist within an allowable distance d_c . The plane and block reconstructions have vertices only at the corners of each plane. In order to measure the accuracy and completeness, all planes are regularly sampled on a $2 \text{cm} \times 2 \text{cm}$ grid. The block reconstruction includes redundant planes to complete cuboid structure. Therefore we use plane reconstruction results for the accuracy test and block reconstruction results for the completeness test.

Overlapped models and accuracy maps of the plane primitives against the ground-truth are illustrated in Fig. 20. Table 2 shows comparison of accuracy and completeness with dense reconstruction results from the same data sets using Kim and Hilton[7]. In measuring completeness, we set the allowable distance d_c as 0.25m for the Gate scene and 1.0m for the Cathedral scene



Figure 19: Cathedral scene (main building only)

considering the scene scales and the range of errors. The dense reconstruction 691 shows better results in the accuracy test, but the proposed block reconstruc-692 tion also shows competitive accuracy especially with the Gate scene whose 693 scene scale factor is set small. The proposed method also shows higher com-694 pleteness with the Cathedral scene because occluded regions are covered in 695 block reconstruction. Although the proposed algorithm cannot reconstruct 696 geometric details in the scene, the reconstructed plane primitives are reliable 697 and provide an efficient approximation of the scene structure. 698

699 7.2. Scene reconstruction results

We evaluated the proposed algorithm on four outdoor and one indoor scenes. The capture points and spherical images are shown in the first and second columns of Fig. 21. The Cathedral scene has a complex structure with many self-occlusions, there is sufficient overlap between views to reconstruct the complete cathedral façade. The Carpark scene was captured in a relatively small but complex area of $20m \times 25m$ including occlusions by



(b) Cathedral

Figure 20: Accuracy maps of plane reconstruction results (Left: Overlaid reconstruction with ground-truth, Right: Accuracy map)

cars. There are relatively few overlapping regions between view 1 and 3. The Highstreet scene covers a street of 80m with four image pairs and includes many small and non-planar objects such as benches and trees. The Plaza scene covers a large and relatively complex area of $60m \times 80m$ with five image pairs. The Reception is an indoor scene captured in three locations. It covers an area of $20m \times 7m$ and the main area is connected to other corridors and rooms. The scene scale parameter S_c for the reception scene is set as 0.4m.

In the captured spherical images, most horizontal straight lines are distorted and it is hard to understand the structure of the scenes from the images. The columns 3-5 in Fig. 21 show automatically façade aligned cubic projection images for the spherical captures. We can observe that all horizontal and vertical lines in the scenes are aligned to x and y axes in the image planes and the vanishing point is located at the centre of each image.

Detect	Accur	acy (RMSE)	Completeness (%)		
Dataset	Gate	Cathedral	Gate	Cathedral	
Dense	0.11	0.32	90.18	74.43	
Proposed	0.17	0.57	82.28	88.37	

Table 2: Accuracy and completeness evaluation of the dense reconstruction [7] and proposed reconstruction methods against ground-truth

Following automatic façade alignment and cubic projection the direction of
the principal axis for Manhattan world plane reconstruction is easily found
from the image axis.

Figure 22 and 23 show the reconstructed block-based structure of the 722 scenes and their texture mapped results. The Cathedral set in Fig. 22723 (a) clearly shows main structure of the scene though details such as narrow 724 steps and awnings in the main building are not reconstructed. The Carpark 725 scene in Fig. 22 (b) is more complicated but is efficiently represented with 726 cuboids. This shows some errors in structure around cars and the wrong 727 texture in some regions because of occlusions between objects and walls. The 728 Highstreet scene in Fig. 22 (c) consists of a long street with many windows 729 on buildings. Some buildings or parts of buildings are missing because they 730 have large windows where reconstructed depth information is unreliable due 731 to the reflection and transparency of the windows. Small windows can be 732 reconstructed with the proposed algorithm because the plane location for the 733 surrounding region is estimated in the refinement process. The Plaza scene in 734 Fig. 23 (a) also includes large reflective regions and produces a few erroneous 735 planes dominated by the scenes reflected on the glass, but they are removed 736 in the refinement process. In the Plaza scene, viewpoints 3 and 4 do not have 737 sufficient overlap in the images. Manual feature matching is performed for 738 multiple view registration. The Reception scene in Fig. 23 (b) demonstrates 739 how the proposed system works for an indoor environment. Textureless walls 740 in the scene may cause serious distortion and errors in dense reconstruction. 741 However, the majority of the walls are reconstructed in correct positions by 742 the proposed system. Free-viewpoint video rendering of the scenes is available 743 from: http://www.cvssp.org/hkim/BlockWorld/BlockRecon-CVIU.mov. 744

From the examples above, we can see that the proposed method generates
a coarse approximation of the scene structure. Texture mapping produces
natural rendering results.

Table 3 shows an analysis of runtime for the Cathedral dataset. We assume that we already have scanned spherical stereo image pairs and their

Table 5: Running time analysis (Cati	lieurai dataset)
Step	Time (sec)
Data loading	0.29
Façade alignment	19.56
Region segmentation	3.91
Plane primitives reconstruction	4.12
Multi-view plane refinement	6.51
Cuboid reconstruction	2.58
Total	36.97

Table 3: Running time analysis (Cathedral dataset)

disparity maps. The proposed algorithms were run on a Intel Core i7 3.40GHz
Windows machine with 32GB RAM. In the table, steps from data loading
to plane primitives reconstruction were performed for individual pairs in
parallel. It takes approximately 40 seconds per datasets.



(a) Cathedral



(b) Carpark



(c) Highstreet



(d) Plaza



(e) Reception

Figure 21: Test datasets (First column: Capture points on maps (from http://maps.google.com), Second column: Spherical capture from three selected points (top image of captured stereo pairs), Three right columns: Façade aligned cubic projection of the selected images



(a) Cathedral



(b) Carpark



(c) Highstreet

Figure 22: Reconstructed structure and texture mapping result 1



(b) Reception

Figure 23: Reconstructed structure and texture mapping result 2

	Dense reconstruction		Proposed method			
Dataset	# of	Data file	# of	Plane data	# of	Cuboid data
	triangles	size	planes	file size	cuboids	file size
Cathedral	792,512	81.9MB	34	4.52KB	25	2.65KB
Carpark	480,644	$52.7 \mathrm{MB}$	43	$5.69 \mathrm{KB}$	35	$3.55 \mathrm{KB}$
Highstreet	$987,\!220$	108.2MB	45	$5.95 \mathrm{KB}$	38	$3.82 \mathrm{KB}$
Plaza	$1,\!254,\!356$	$132.5 \mathrm{MB}$	48	$6.34 \mathrm{KB}$	34	$3.46 \mathrm{KB}$
Reception	325,780	40.5MB	50	$6.60 \mathrm{KB}$	43	4.27KB

Table 4: Comparison with dense reconstruction method [7]

754 7.3. Comparison with other methods

We compared the amount of data produced for the scene reconstruction 755 with dense geometry reconstruction using the approach of Kim and Hilton^[7]. 756 We used the same input images and disparity estimation methods for both re-757 constructions. The dense reconstruction results were saved in the obj format 758 with vertex positions, vertex normals, UV texture and triangle information. 759 The plane primitive information was saved as an ASCII file with the format 760 in Eq. (6) with headers. The cube information was also saved as an ASCII 761 file with the position and length in each direction. Texture index numbers 762 are also included to identify the correct texture for rendering. In Table 4, 763 we see that the size of data required to represent the scenes is reduced by 764 three to four orders of magnitude. The block world provides a compact rep-765 resentation of the scene as a set of 3D cuboid proxies for rendering. Detailed 766 geometry is not represented, but texture mapping enables rendering of the 767 appearance of detailed geometry suitable for scene visualisation. 768

We also compared visualisation quality, processing time and data file size 769 including texture information for the Cathedral scene with other methods 770 in Fig. 24 and Table 5. The dense reconstruction method [7] produced a 771 huge amount of data. It recovered fine details of the scene but shows geo-772 metrical errors in the occluded or ambiguous regions such as the ceiling and 773 windows. The LIDAR model was created from 7 LIDAR scans and dozens 774 of reference stills using MAYA software⁶ by a professional CG designer. It is 775 a clean model with high accuracy but it took one full day even by the spe-776 cialist to build the mesh model and generate textures from the raw sources. 777 The SketchUp result was modelled from nine photographs using the Google 778

⁶Autodest MAYA, http://www.autodesk.co.uk/products/maya/



(a) Final model (From Left: Dense recon., CG from LIDAR, SketchUp and Proposed)



(b) Rendering results (From Left: Photograph, Dense recon., CG from LIDAR, SketchUp and Proposed)

Figure 24: Comparison with other methods (Cathedral)

Table 5. 1 rocessing time and volume comparison with other methods (Cathedrar)				
Method	Dense recon.[7]	CG from LIDAR	SketchUp	Proposed
Processing Time	$12 \mathrm{~mins}$	1 day	3 hours	$37 \mathrm{secs}$
File size (inc. texture)	$84 \mathrm{MB}$	12.9 MB	12.5 MB	$1.33 \ \mathrm{MB}$

Table 5: Processing time and volume comparison with other methods (Cathedral)

SketchUp tool. It took about 3 hours to align vanishing points and geometrical primitives to the original photographs. Blending of multiple photographs
for texture mapping resulted in incorrect or blurred textures in some regions.
The proposed method is the fastest in building geometry and shows relatively
clear structure and texture with the minimum amount of data.

784 8. Conclusions

In this paper, we propose a block-based simplified 3D scene reconstruction method from spherical stereo image pairs. Vertical spherical stereo pairs are captured at multiple locations in the scene and converted into cubic projection images which are aligned to principal axes. A façade alignment algorithm is proposed which automatically generates central point perspective images

aligned with the principal building faces. This is advantageous in 3D struc-790 ture reconstruction as it is free from spherical distortion and has a vanishing 791 point at the centre of the image aligned with the principal axes. From the 792 captured images and estimated disparity maps, planar regions are segmented 793 and reconstructed. Reconstructed planes from multiple capture locations are 794 merged and refined to obtain a more complete scene reconstruction. Finally, 795 optimal cuboid structures are reconstructed based on the density of plane 796 primitives. 797

Results show that the proposed algorithm produces a simplified struc-798 tured representation of the scene requiring several orders of magnitude less 790 storage compared with dense scene reconstruction. The resulting scene rep-800 resentation provides a compact 3D proxy for visualisation of the scene. Po-801 tential future extensions of this research include: 1) Simplified scene recon-802 struction with arbitrary planes not aligned to the principal axes and various 803 type of 3D structure primitives; 2) Bundle adjustment for large scale loop 804 closure and precise registration; 3) Texture blending and occlusion mapping 805 from multiple viewpoints. 806

807 Acknowledgement

This research was supported by the European Commission, FP7 IMPART project (grant agreement No 316564).

810 References

- [1] Y. Furukawa, J. Ponce, Accurate, dense, and robust multiview stere opsis, IEEE Trans. Pattern Analysis and Machine Intelligence 32 (8)
 (2010) 1362–1376.
- [2] H. Vu, R. Keriven, P. Labatut, J. Pons, Towards high-resolution large scale multi-view stereo, in: Proceedings of CVPR, 2009, pp. 1430–1437.
- [3] N. Salman, M. Yvinec, Surface reconstruction from multi-view stereo,
 in: Proceedings of ACCV, 2009.
- [4] C. Strecha, W. Hansen, L. Gool, P. Fua, U. Thoennessen, On benchmarking camera calibration and multi-view stereo for high resolution imagery, in: Proceedings of CVPR, 2008, pp. 1–8.

- [5] S. Agarwal, N. Snavely, I. Simon, S. Seitz, R. Szeliski, Building rome in a day, in: Proceedings of ICCV, 2009, pp. 72–79.
- [6] M. Pollefeys, D. Nistér, J. Frahm, A. Akbarzadeh, P. Mordohai,
 B. Clipp, C. Engels, D. Gallup, S. Kim, P. Merrell, C. Salmi, S. Sinha,
 B. Talton, L. Wang, Q. Yang, H. Stewénius, R. Yang, G. Welch,
 H. Towles, Detailed real-time urban 3d reconstruction from video, International Journal of Computer Vision 78 (2) (2008) 143–167.
- [7] H. Kim, A. Hilton, 3d scene reconstruction from multiple spherical stereo pairs, International Journal of Computer Vision 104 (1) (2013) 94–116.
- [8] G. Sellers, J. Obert, P. Cozzi, K. Ring, E. Persson, J. de Vahl, J. M. P.
 van Waveren, Rendering massive virtual worlds, in: Proceedings of SIGGRAPH 2013 Courses, 2013, pp. 23:1–23:88.
- [9] Y. Furukawa, B. Curless, S. Seitz, R. Szeliski, Manhattan-world stereo,
 in: Proceedings of CVPR, 2009.
- [10] A. Gupta, A. A. Efros, M. Hebert, Blocks world revisited: Image un derstanding using qualitative geometry and mechanics, in: Proceedings
 of ECCV, 2010.
- [11] H. Kim, A. Hilton, Planar urban scene reconstruction from spherical
 images using facade alignment, in: Proceedings of IVMSP, 2013.
- [12] P. Debevec, C. Taylor, J. Malik, Modeling and rendering architecture
 from photographs: A hybrid geometry- and image-based approach, in:
 Proceedings of SIGGRAPH, 1996, pp. 11–20.
- [13] A. V. Hengel, A. Dick, T. Thormählen, B. Ward, P. H. S. Torr, Videotrace: Rapid interactive scene modelling from video, in: Proceedings of
 SIGGRAPH, 2007.
- [14] S. N. Sinha, D. Steedly, R. Szeliski, M. Agrawala, M. Pollefeys, Interactive 3d architectural modeling from unordered photo collections, in: Proceedings of SIGGRAPH ASIA, 2008.
- ⁸⁴⁹ [15] D. Hoiem, A. Efros, M. Hebert, Recovering surface layout from an image, International Journal of Computer Vision 75 (1).

- [16] P. Müller, G. Zeng, P. Wonka, L. Gool, Image-based procedural modeling of facades, in: Proceedings of SIGGRAPH, 2007.
- In J. Xiao, T. Fang, P. Zhao, M. Lhuillier, L. Quan, Image-based street-side
 city modeling, in: Proceedings of SIGGRAPH ASIA, 2009.
- [18] F. Bellotti, R. Berta, R. Cardona, A. D. Gloria, An architectural approach to efficient 3d urban modeling, Computers & Graphics 35 (5) (2011) 1001–1012.
- L. Simon, O. Teboul, P. Koutsourakis, N. Paragios, Random exploration
 of the procedural space for single-view 3d modeling of buildings, Inter national Journal of Computer Vision 93 (2011) 253–271.
- [20] M. Mathias, A. Martinovic, J. Weissenberg, L. J. V. Gool, Procedural 3d building reconstruction using shape grammars and detectors, in: Proceedings of 3DIMPVT, 2011, pp. 304–311.
- ⁸⁶⁴ [21] S. Satkin, M. Rashid, J. Lin, M. Hebert, 3dnn: 3d nearest neighbor,
 ⁸⁶⁵ International Journal of Computer Vision 111 (1) (2015) 69–97.
- [22] G. Schindler, P. Krishnamurthy, F. Dellaert, Line-based structure from
 motion for urban environments, in: Proceedings of 3DPVT, 2006, pp.
 846–853.
- [23] C. Hane, C. Zach, B. Zeisl, M. Pollefeys, A patch prior for dense 3d reconstruction in man-made environments., in: Proceedings of 3DIMPVT, 2012, pp. 563–570.
- ⁸⁷² [24] S. Sinha, D. Steedly, R. Szeliski, Piecewise planar stereo for image-based
 ⁸⁷³ rendering, in: Proceedings of ICCV, 2009.
- [25] D. Gallup, J.-M. Frahm, M. Pollefeys, Piecewise planar and non-planar
 stereo for urban scene reconstruction, in: Proceedings of CVPR, 2010,
 pp. 1418–1425.
- R. Toldo, A. Fusiello, Photo-consistent planar patches from unstructured cloud of points, in: Proceedings of ECCV, 2010.
- R. Toldo, A. Fusiello, Robust multiple structures estimation with j linkage, in: Proceedings of ECCV, 2008.

- [28] J. Coughlan, A. Yuille, Manhattan world: orientation and outlier detection by bayesian inference, Neural Computation 15 (5) (2003) 1063–
 1088.
- ⁸⁸⁴ [29] B. Micusik, J. Kosecka, Multi-view superpixel stereo in urban environments, International Journal of Computer Vision 89 (1) (2010) 106–119.
- [30] A.-L. Chauve, P. Labatut, J.-P. Pons, Robust piecewise-planar 3d reconstruction and completion from large-scale unstructured point data,
 in: Proceedings of CVPR, 2010, pp. 1261–1268.
- [31] A. Kowdle, Y.-J. Chang, A. C. Gallagher, T. Chen, Active learning for
 piecewise planar 3d reconstruction, in: Proceedings of CVPR, 2011, pp.
 929–936.
- [32] Q.-Y. Zhou, U. Neumann, 2.5d dual contouring: A robust approach to
 creating building models from aerial lidar point clouds, in: Proceedings
 of ECCV, 2010, pp. 115–128.
- ⁸⁹⁵ [33] C. Poullis, S. You, 3d reconstruction of urban areas, in: 3DIMPVT, ⁸⁹⁶ 2011, pp. 33–40.
- ⁸⁹⁷ [34] W. Li, G. Wolberg, S. Zokai, Lightweight 3d modeling of urban buildings
 ⁸⁹⁸ from range data, in: Proceedings of 3DIMPVT, 2011, pp. 124–131.
- [35] W. Nguatem, M. Drauschke, H. Mayer, Finding cuboid-based building
 models in point clouds, in: Proceedings of ISPRS, 2012, pp. 149–154.
- ⁹⁰¹ [36] J. Xiao, Y. Furukawa, Reconstructing the worlds museums, Interna-⁹⁰² tional Journal of Computer Vision 110 (3) (2014) 243–258.
- ⁹⁰³ [37] P. Sturm, A method for 3d reconstruction of piecewise planar objects
 ⁹⁰⁴ from single panoramic images, in: Proceedings of IEEE Workshop on
 ⁹⁰⁵ Omnidirectional Vision, 2000, pp. 119–126.
- [38] H.-D. Kang, K.-H. Jo, 3d reconstruction of planar objects using the
 properties of plane and vanishing points from a single panoramic image,
 in: Proceedings of World Multi-Conference on Systemics, Cybernetics
 and Informatic, 2006, pp. 342–346.

- [39] B. Micusik, J. Kosecka, Piecewise planar city 3d modeling from street
 view panoramic sequences, in: Proceedings of CVPR, 2009, pp. 2906– 2912.
- [40] D. Anguelov, C. Dulong, D. Filip, C. Frueh, S. Lafon, R. Lyon, A. Ogale,
 L. Vincent, J. Weaver, Google street view: Capturing the world at street
 level, IEEE Computer 43 (6) (2010) 32–38.
- [41] D. Feldman, D. Weinshall, Realtime ibr with omnidirectional crossedslits projection, in: Proceedings of ICCV, 2005, pp. 839–845.
- [42] P. F. Felzenszwalb, D. P. Huttenlocher, Efficient graph-based image seg mentation, International Journal of Computer Vision 59 (2).
- [43] A. Banno, K. Ikeuchi, Omnidirectional texturing based on robust 3d
 registration through euclidean reconstruction from two spherical images,
 Computer Vision and Image Understanding 114 (4) (2010) 491–499.
- [44] J. Matas, C. Galambos, J. Kittler, Robust detection of lines using the
 progressive probabilistic hough transform, Computer Vision and Image
 Understanding (2000) 119–137.
- ⁹²⁶ [45] D. Eberly, Least Squares Fitting of Data, http://www.
 ⁹²⁷ geometrictools.com/Documentation/LeastSquaresFitting.pdf,
 ⁹²⁸ [Online; accessed 22-Aug-2013] (2008).
- ⁹²⁹ [46] R. Schnabel, R. Wahl, R. Klein, Efficient ransac for point-cloud shape
 ⁹³⁰ detection, Computer Graphics Forum 26 (2) (2007) 214–226.
- [47] A. Hilton, Scene modelling from sparse 3d data, Image and Vision Computing 23 (10) (2005) 900–920.
- [48] Y. Furukawa, B. Curless, S. M. Seitz, R. Szeliski, Reconstructing building interiors from images, in: Proceedings of ICCV, 2009.
- [49] H. Bay, A. Ess, T. Tuytelaars, L. Gool, Surf: Speeded up robust features, Computer Vision and Image Understanding 110 (2008) 346–359.
- ⁹³⁷ [50] T. Mullen, Mastering Blender, 1st Edition, Wiley Publishing, Inc., 2009.
- ⁹³⁸ [51] S. Seitz, B. Curless, J. Diebel, D. Scharstein, R. Szeliski, A compari⁹³⁹ son and evaluation of multi-view stereo reconstruction algorithms, in:
 ⁹⁴⁰ Proceedings of CVPR, 2006, pp. 519–528.