

000  
001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

# Growing semantically meaningful models for visual SLAM.

Anonymous CVPR submission

Paper ID \*\*\*\*

## Abstract

Though modern Visual Simultaneous Localisation and Mapping (vSLAM) systems are capable of localising robustly and efficiently even in the case of a monocular camera, the maps produced are typically sparse point-clouds that are difficult to interpret and of little use for higher-level reasoning such as obstacle avoidance or navigation. In this paper we begin to address this deficiency, presenting progress on expanding the competency of visual SLAM systems to build richer maps. Specifically, we concentrate on modelling indoor scenes using semantically meaningful surfaces and accompanying labels, such as “floor”, “wall”, and “ceiling” — an important step towards a representation that can support higher-level reasoning and planning.

We leverage the Manhattan world assumption and show how to extract vanishing directions jointly across a video stream. We then propose a guided line detector that utilises known vanishing points to extract extremely subtle axis-aligned edges. We utilise recent advances in single view structure recovery to building geometric scene models and demonstrate our system operating in real-time.

## 1. Introduction

The simultaneous localisation and mapping problem has received considerable attention over past decades, which is unsurprising given its centrality in fields from mobile robotics to augmented reality. Considerable progress has been made over this period and modern SLAM systems perform efficiently and robustly even in the case of a monocular video stream [8].

Many high-level reasoning and planning problems can benefit from an accurate underlying SLAM system, but SLAM point clouds alone provide a poor basis upon which to reason about scene semantics as they represent just a fraction of the information present in the original images. Photometric cues for edges, surfaces, occlusion boundaries, and texture information, among others, are lost entirely.

An important step towards higher-level reasoning tasks is to represent the structure of the scene at a semantic

level, using meaningful concepts such as “floor” and “wall”. These entities assist in reasoning as they correlate with, for example, the locations in which objects might appear or the actions that might be feasible in various locations.

The present work investigates the extraction of such a semantic scene model from video sequences, utilising an underlying SLAM system together with rich photometric cues. We focus on indoor scenes as they exhibit a rich set of regularities that assist in model building. One such regularity is the prevalence of three mutually orthogonal orientations around which man-made environments are often built up. The use of this observation to restrict the space of possible scene interpretation has come to be known in the literature as the Manhattan-world assumption. Typically, the orientation of the camera with respect to these dominant directions is *a priori* unknown and must be extracted explicitly, for example by identifying vanishing points. We propose a new method that leverages the camera poses provided by SLAM to estimate vanishing directions jointly across a video sequence, allowing frames with salient edge information to inform the system about vanishing directions in frames with poor or non-existent edge information.

We then return to the images and re-identify line segments using the known vanishing points to inform our search. We propose a novel line detector that takes vanishing point locations into account to identify important structural edges that are aligned with a dominant direction but which may exhibit weak gradient magnitudes, whilst ignoring stronger gradients generated by surface texture or occlusion boundaries.

The final component of our system joins the line segments into Manhattan building structures, inspired by recent work in single view reconstruction [3]. We enumerate possible building structures that the observed line segments could generate under the Manhattan world assumption, evaluating each for consistency with surface orientations estimated from photometric cues. The key contribution of this section is the extension to multiple views of the hypothesis testing framework, which we show gives significantly improved results.

The remainder of this paper is organised as follows. In

108 the next section we overview prior work in this field. Following this we describe the three primary components of  
109 our system in order: the joint vanishing point estimator, the  
110 guided line search algorithm, and the reconstruction sys-  
111 tem, with results given in each section. Finally we discuss  
112 the results and present closing remarks.  
113  
114

## 115 2. Background

116 In recent years the need for semantics to be connected  
117 with SLAM maps has been recognised by several re-  
118 searchers. Stachniss *et al.* [14] have taken an image-centric  
119 approach wherein 3D features are projected into frames  
120 and used together with photometric cues to classify en-  
121 vironments into semantic categories such as “corridor” or  
122 “room”. Posner *et al.* [10] take this a step further by seg-  
123 menting incoming frames based jointly on 3D and photo-  
124 metric cues, after which image segments are assigned la-  
125 bels such as “grass” or “vehicle”. Xiao and Quan [17] have  
126 approached this problem by solving a multiple label MRF  
127 over superpixels from two or more views, with 2D and 3D  
128 cues combined to form node potentials. Brostow *et al.* [1]  
129 also combined photoemtric and geometric cues but used  
130 simpler cues and did not force consistency across views.  
131

132 Buschka and Saffiotti [2] opt instead to reason directly in  
133 the map. They build an occupancy grid and identify room  
134 boundaries by repeated application of morphological filters.  
135 More recently, Golovinskiy *et al.* [5] learn to segment and  
136 identify objects in city-wide reconstructions using machine  
137 learning techniques. These approaches discard the original  
138 images after building a map, which we believe discards  
139 many useful cues not captured in the map representation.  
140

141 Furukawa *et al.* [4] have shown how to produce recon-  
142 structions of indoor scenes under the Manhattan-world as-  
143 sumption. Their reconstructions are of very high quality but  
144 they quoted computation times of between one minute and  
145 one hour forty minutes, which makes their approach unsuit-  
146 able to our on-line situation. Their system uses multiple-  
147 view stereo to first reconstruct a dense point cloud, whereas  
148 we are interested in working with a sparse point cloud and  
149 leveraging photometric cues for real-time performance.  
150

151 Several authors have recently demonstrated impressive  
152 single view reconstruction systems. Hoiem *et al.* [7] pose  
153 the problem as a multi-class segmentation problem, with  
154 labels corresponding to 3D geometry. Saxena *et al.* [12] ob-  
155 obtain reconstructions by estimating surface normals of image  
156 patchlets. Gould *et al.* [6] label with both geometry and  
157 object classes for improved accuracy.  
158

159 Lee *et al.* [3] take a geometric approach in which de-  
160 tected line segments are connected to form hypotheses  
161 about 3D structure. The authors show that the Manhattan  
162 worlds highly constrain the set of possible building struc-  
163 tures, so an exhaustive search is possible. This work forms  
164 the basis for Section 5 in which we extend this approach to  
165

166 multiple views.  
167

168 Many researchers have proposed methods for comple-  
169 ting partial lines or identifying subtle-yet-important line  
170 segments that humans see easily. For example, Sarti *et al.*  
171 [11] iteratively fill missing boundaries starting from a refer-  
172 ence point, Tuytelaars *et al.* [16] iterate between detecting  
173 lines and identifying intersections using the Hough trans-  
174 form, and Shufelt [13] models line detection errors explic-  
175 itely. This body of literature deals with the single image case,  
176 whereas we use known vanishing points obtained from mul-  
177 tiple views to improve the accuracy and speed of the search.  
178

## 179 3. Extracting a canonical coordinate frame

180 In order to make use of the Manhattan world assump-  
181 tion we must first discover the orientation of the Manhattan  
182 world with respect to the camera. Equivalent to the Man-  
183 hattan world assumption is the statement that there exists a  
184 “canonical” coordinate frame in which world surfaces are  
185 axis-aligned. The problem is therefore to discover the ro-  
186 tation  $R_w$  mapping canonical coordinate to SLAM coordi-  
187 nates. Since our SLAM system has already determined the  
188 relative poses between successive frames, this rotation is  
189 fixed for all frames.  
190

191 In the past researchers have discovered  $R_w$  for single  
192 images by identifying vanishing points (cite), or in the  
193 structure-from-motion setting by clustering surface nor-  
194 mals estimated from local neighbourhoods of point clouds  
195 (cite Furukawa). The former approach fails for frames con-  
196 taining edges in only one or none of the three dominant di-  
197 rections — a common occurrence for video sequences of  
198 indoor environments — and for frames containing many  
199 edges not aligned with one of the dominant directions.  
200 Since  $R_w$  is fixed for all frames it makes sense to leverage  
201 all available data during estimation rather than to estimate  
202 vanishing points separately for each frame. The surface nor-  
203 mal approach requires a dense scene reconstruction, which  
204 is currently prohibitive in for real-time applications.  
205

206 We begin by sampling a set of key frames from the in-  
207 put sequence with the number of key frames chosen for  
208 tractability. On each key frame we run the Canny edge  
209 detector followed by a standard edge linking algorithm to  
210 identify a set of straight line segments  $L_i = \{x : l_i^T x =$   
211  $0\}$  and associated confidences  $c_i$ , computed from the ratio  
212 of eigenvalues of each line segment’s constituent pixels.  
213

214 The SLAM system provides a pose  $P_i$  for each frame,  
215 which contains a rotation  $R_i$  and translation  $t_i$ . These are  
216 measured with respect to some arbitrary coordinate frame  
217 determined during initialization, which we will henceforth  
218 refer to as the “SLAM” coordinate frame.  
219

220 In the canonical coordinate system the vanishing points  
221 have homogeneous coordinates  $(1, 0, 0, 0)$ ,  $(0, 1, 0, 0, 0)$ ,  
222 and  $(0, 0, 1, 0)$ . Their projections into camera  $i$  are given  
223

216 by  
 217  
 218  
 219  
 220  
 221  
 222  
 223

$$\mathbf{v}_x = R_i R_w \mathbf{e}_x \quad (1)$$

$$\mathbf{v}_y = R_i R_w \mathbf{e}_y \quad (2)$$

$$\mathbf{v}_z = R_i R_w \mathbf{e}_z \quad (3)$$

(4)

224 where  $\mathbf{e}_x$ ,  $\mathbf{e}_y$ , and  $\mathbf{e}_z$  are unit vectors in the  $x$ ,  $y$ , and  $z$   
 225 directions respectively, and  $R_w$  is the mapping from canonical  
 226 coordinate to SLAM coordinates.

227 We can now write down an error function to be min-  
 228 imised in terms of the world rotation  $R_w$ :

$$E(R_w) = \sum_i \sum_j \sum_{k=1}^3 r_{jk} (l_j^T R_i R_w e_k)^2 \quad (5)$$

233 where  $r_{jk}$  is the responsibility of the  $k^{\text{th}}$  vanishing point  
 234 for the  $j^{\text{th}}$  line segment. The term  $R_i R_w e_k$  represents the  
 235 projection of the  $k^{\text{th}}$  vanishing point into the  $i^{\text{th}}$  frame, and  
 236  $(l_j^T R_i R_w e_k)^2$  is the deviation of the  $j^{\text{th}}$  line segment from  
 237 it. The full error function (5) is then a sum over all line seg-  
 238 ments and vanishing points, with the deviations weighted  
 239 by the estimated responsibilities of each vanishing point for  
 240 the line segments.

241 While other authors search for vanishing points by clus-  
 242 tering on the Gaussian sphere, enforcing orthogonality con-  
 243 straints afterwards, we prefer to optimise in terms of  $R_w$   
 244 directly, which builds the orthogonality constraints into the  
 245 core of the estimation process.

246 We now proceed to describe our EM algorithm. Our ex-  
 247 pectation step computes the responsibilities of the vanishing  
 248 points for the line segments. We assume a Gaussian likeli-  
 249 hood

$$p(l_j | v_k) = G(l_j^T v_k; \mu) \quad (6)$$

250 as well as a fixed prior on observing a spurious line segment

$$p(s) = \rho \quad (7)$$

$$p(l_j | s) = 1 \quad (8)$$

256 The responsibilities are then computed as the normalized  
 257 likelihoods

$$r_{jk} = \frac{p(l_j | v_k)}{\rho + \sum_k p(l_j | v_k)} \quad (9)$$

(10)

263 The M step consists of optimising  $R_w$  with respect to  
 264 the error function (5). There is no closed form solution for  
 265 the optimal  $R_w$  so we instead perform gradient descent. We  
 266 represent  $R_w$  in the Lie algebra as a member of the special  
 267 orthogonal group  $SO(3)$ .

$$R_w = \exp(\sum m_i G_i) \quad (11)$$

270 where the  $G_i$  are the generator matrices for  $SO(3)$  and  $m_i$   
 271 are the respective coefficients. The advantage of using the  
 272 Lie algebra is that at each step we are guaranteed that  $R_w$   
 273 remains a pure rotation, whereas in other representations,  
 274 such as simply optimizing the elements of the  $3 \times 3$  rotation  
 275 matrix, this is not the case.

276 We now differentiate (5) with respect to the coefficient  
 277 vector  $\mathbf{m}$ :

$$\begin{aligned} \nabla f(x) &= \sum_i \sum_j \sum_{k=1}^3 (2r_{jk} l_j^T R_i R_w e_k - r_{jk} l_j^T R_i \nabla R_w) \\ \nabla R_w &= [G_1 \mathbf{e}_1, G_2 \mathbf{e}_2, G_3 \mathbf{e}_3] \end{aligned} \quad (12)$$

283 In accordance with the standard gradient descent ap-  
 284 proach our update rule is then

$$\mathbf{m}^{t+1} = \mathbf{m}^t - \frac{f(\mathbf{m}^t)}{\|\nabla f(\mathbf{m}^t)\|_2} \nabla f(\mathbf{m}^t) \quad (14)$$

288 In summary, to obtain  $R_w$  we iterate between assigning  
 289 responsibilities (the E step) and optimising the rotation be-  
 290 tween world and SLAM coordinates (the M step). Each M  
 291 step consists of a gradient descent in the Lie algebra. In  
 292 practice we found that our system converges in around 25  
 293 iterations of the EM algorithm, and around 20 steps are re-  
 294 quired for each gradient descent.

295 Figure 1 shows the vanishing points identified in one of  
 296 our sequences. Since each frame is informed by the entire  
 297 sequence we are able to identify a globally consistent co-  
 298 ordinate frame where single-image vanishing point detec-  
 299 tion fails. Figure 2 shows a side-by-side comparison with  
 300 the single-image vanishing point detector of [9]. Recently  
 301 proposed improvements to single-image approach [15] may  
 302 improve slightly on these but we found that in cases where  
 303 the single-image approach fails there is often simply not  
 304 enough information available in individual frames to iden-  
 305 tify the appropriate coordinate frame, so any single-image  
 306 approach will necessarily fail.

### 3.1. Identifying the Vertical Direction

310 Of the three axes in the canonical coordinate frame, the  
 311 one corresponding to the vertical direction is semantically  
 312 distinct from the others since it defines the orientation of the  
 313 ground and ceiling planes, and also the direction in which  
 314 gravity operates. It is easy to identify this axis since humans  
 315 necessarily move over the ground plane when capturing a  
 316 sequence, and have limited scope for moving the camera in  
 317 the up-down direction. We therefore set the vertical direc-  
 318 tion to that in which the range of camera poses is smallest.  
 319 Having identified  $R_w$  there are only three possible choices  
 320 for the vertical direction and we found this heuristic to work  
 321 effectively in all of our sequences.

322 An alternative approach would be to make use of the ten-  
 323 dency for humans to orient the camera so that the vertical

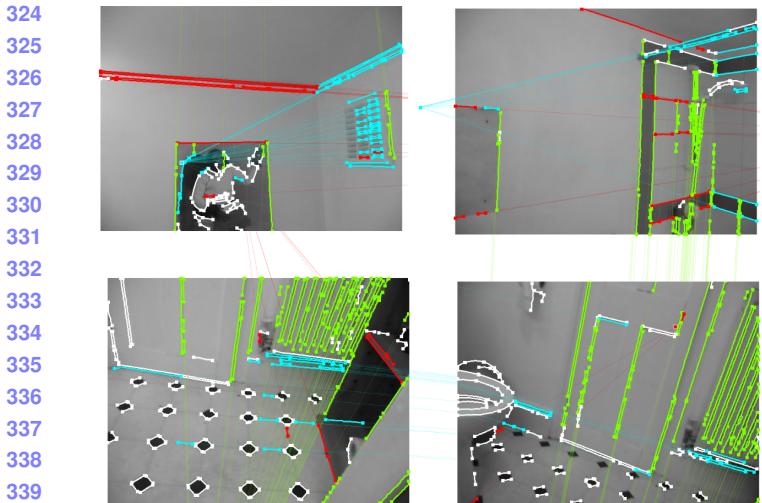


Figure 1. Four frames from the “bathroom” sequence and the detected vanishing points. The vanishing points are correctly identified despite the strong distractor gradients generated by the floor tiles, which is possible only by integrating information from multiple views into the estimation process. Also note that the up–down, east–west, and north–south directions are labelled consistently in each frame.

direction in the world corresponds to the up–down direction in the image. The major drawback of this approach arises in the common case that an entire sequence is filmed with the camera rotated by 90° so as to capture a high, narrow environment, in which case this test fails.

### 3.2. Relaxing the Manhattan world assumption

The strong Manhattan assumption states that any pair of surfaces of interest are either parallel or orthogonal to one another. Amongst indoor scenes, one common deviation from this is scenes with vertical walls that are not pair-wise orthogonal. We define the weak Manhattan assumption as “the environment consists of a horizontal ground plane and corresponding ceiling plane, and a set of vertical wall segments extending continuously between them.” This formulation permits pairs of wall segments that are not orthogonal to one another.

Weakly Manhattan scenes contain much of the regularity of strongly Manhattan scenes. We deal with the weak Manhattan assumption as follows. First, we run the rotation recovery algorithm described above. Next, for each line marked as spurious by the EM algorithm we find its intersection with the horizon,

$$\mathbf{v}_+ = R_w^{-T} R_i^{-T} \mathbf{l}_j \times \mathbf{e}_z \quad (15)$$

which would be its vanishing point if it were horizontal. The horizontal lines on all vertical wall segments in with a given orientation will coincide in their intersection with the horizon. We parameterize each such intersection according

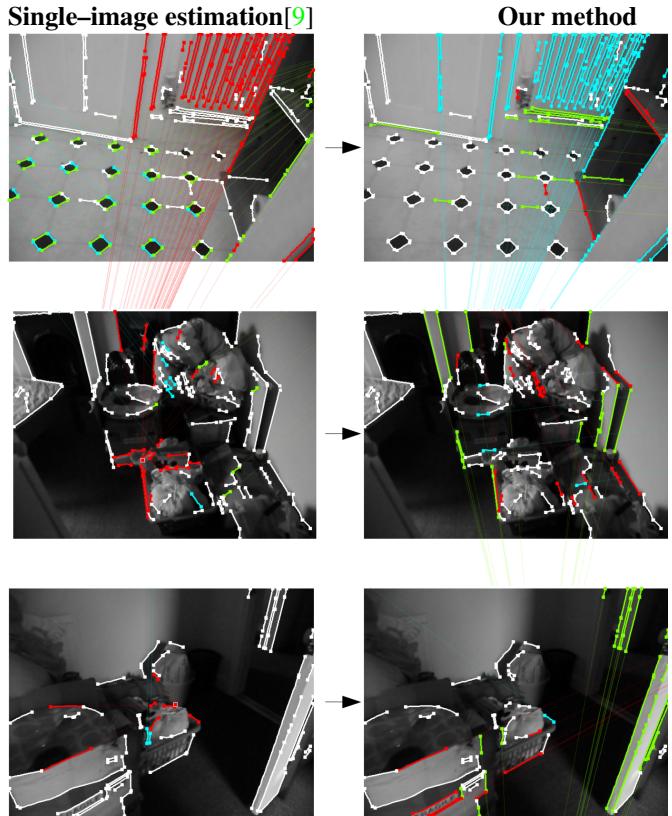


Figure 2. Comparison between vanishing points estimated for single views (left column) and joint estimates from 20 frames in a video sequence (right column). Our method is able to identify vanishing directions correctly in these difficult cases, whereas the single view estimator is confused by non–Manhattan line segments.

to the angle it makes with the  $x$  and  $y$  axes

$$\theta_j = \text{atan}(\mathbf{e}_y^T \mathbf{v}_+, \mathbf{e}_x^T \mathbf{v}_+) \quad (16)$$

We accumulate the  $\theta_j$  into histogram bins and identify any bins above a fixed threshold  $k$  that are local maxima. Such bins represent the orientation of the additional vertical wall segments we seek. Finally, we re–estimate the vanishing point for each orientation by minimising the likelihood described in the previous section.

## 4. Directed Line Search

Many of the structurally important edges in indoor scenes generate gradients far weaker than those generated by surface texture or occlusion boundaries. For example, in the common scenario that the walls in a room are painted the same colour, the image gradients at wall intersections will be generated only from subtle lambertian lighting differences.

We have found that the standard line detection approach (Canny followed by edge linking) is incapable of detecting many important structural edges unless the thresholds are

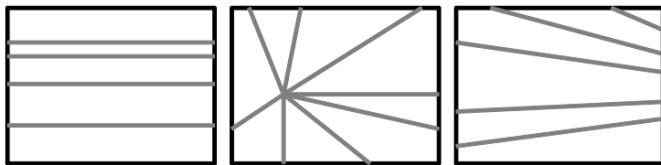
432 lowered to such a point that any textured surface generates  
 433 many thousands of spurious line segments. However, if the  
 434 Canny thresholds are set to some level high enough to re-  
 435 duce spurious detections to a manageable level then the true  
 436 positive line segments are plentiful enough to robustly de-  
 437 tect vanishing directions but not to recover building struc-  
 438 tures (as described in the following section).  
 439

440 To overcome this, we return to the images after deter-  
 441 mining the vanishing directions and search again for line  
 442 segments, this time with the known vanishing directions in-  
 443 forming the search.  
 444

445 Each vanishing point is associated with a one-parameter  
 446 family of lines extending from it. Parameterizing this fam-  
 447 ily is non-trivial since there exists a continuum of scenarios  
 448 from parallel lines meeting at infinity to lines radiating out-  
 449 wards from a vanishing point inside the image bounds (see  
 450 Figure 3). We choose to parameterize the lines by the angle  
 451  $\theta$  they form in the plane generated by the other two vanish-  
 452 ing points  
 453

$$\theta(i, \mathbf{x}) = \text{atan}(\mathbf{x}^T \mathbf{v}_j, \mathbf{x}^T \mathbf{v}_k) \quad (17)$$

454 where all vectors are in homogeneous coordinates. This  
 455 representation is free of singularities and sampling uni-  
 456 formly in  $\theta$  space produces roughly uniformly distributed  
 457 lines in the image.  
 458



459 Figure 3. Lines meeting at a vanishing point.  
 460  
 461

462 For each pixel  $\mathbf{x}$  we begin by estimating its vanishing  
 463 point association. Since vanishing points are always well-  
 464 separated (due to orthogonality) we obtain a reasonable esti-  
 465 mate of the vanishing point a line containing the pixel would  
 466 be associated with from the image gradient  $\mathbf{g}_x$ .  
 467

$$\text{assoc}(\mathbf{x}) = \text{argmin}_i ((\mathbf{x} - \mathbf{v}_i)^T \mathbf{g}_x) \quad (18)$$

470 where  $i$  ranges over the three possible vanishing points.  
 471

472 Next we build a histogram over  $\theta$  for each of the three  
 473 vanishing points. Pixels vote for the lines they generate ac-  
 474 cording to the strength of the gradient at the pixel and the  
 475 agreement between the pixel's gradient orientation and the  
 476 direction to its associated vanishing point. Formally, the  
 477 weight with which a pixel at  $\mathbf{x}$  votes for the line  $\mathbf{l}_{i,x}$  it gen-  
 478 erates with the  $i^{\text{th}}$  vanishing point is given by  
 479

$$w_{i,x} = (\|\mathbf{g}_x\|_2)^\alpha \cos^\beta ((\mathbf{x} + \mathbf{g}_x)^T \mathbf{l}_{i,x}) \quad (19)$$

480 where  $\alpha$  and  $\beta$  are parameters that determine the relative  
 481 importance of the gradient orientation and magnitude.  
 482

483 The bin width for the histograms is set to the minimum  
 484 size such that no bin spans more than two pixels anywhere  
 485 in the image. Histogram peaks are identified by applying  
 486 non-maximum suppression followed by thresholding. The  
 487 final line segments are identified by walking along lines  
 488 corresponding to the identified histogram peaks and linking  
 489 edge pixels using hysteresis. A line segment is started each  
 490 time a gradient magnitude above the high threshold  $k_{\text{high}}$   
 491 is detected, and is ended when the gradient magnitude drops  
 492 below the low threshold  $k_{\text{low}}$ .  
 493

494 This algorithm is highly efficient since only one pass  
 495 over the image is needed to build all three histograms. Iden-  
 496 tifying histogram peaks and walking along the correspond-  
 497 ing lines is then computationally trivial.  
 498

499 Figure 4 shows four example frames and the lines corre-  
 500 sponding to the histogram peaks. The detector fires at subtle  
 501 axis-aligned gradients while ignoring strong but non-axis-  
 502 aligned segments. Figure 5 shows a side-by-side compari-  
 503 son with the line detector of [9], which employs the stan-  
 504 dard Canny edge detector followed by an edge linking al-  
 505 gorithm. In each example our detector is able to identify  
 506 important structural edges that the Canny detector does not  
 507 respond to. We found that lowering the Canny thresholds  
 508 sufficiently that these edges were detected generated many  
 509 thousands of spurious line segments on the textured carpet  
 510 and other areas.  
 511



512 Figure 4. Four frames from the “lab” sequence, with peaks in  
 513 the histograms over  $\theta$  highlighted. The rays capture the  
 514 important geometric structure extremely accurately, with almost no false  
 515 positives.  
 516

## 5. Recovering Building Structures

517 The final component of our system estimates building  
 518 structures from the set of axis-aligned lines produced by  
 519 our line detector. Lee *et al.* [3] have shown that under as-  
 520 sumptions of perspective projection and a planar scene  
 521 the lines intersect at a vanishing point. We can use this  
 522 fact to estimate the camera parameters and the 3D struc-  
 523 ture of the scene. We use a bundle adjustment approach  
 524 to refine the camera parameters and the 3D structure.  
 525

526  
 527  
 528  
 529  
 530  
 531  
 532  
 533  
 534  
 535  
 536  
 537  
 538  
 539

The figure consists of three vertically stacked panels, each showing a grayscale depth map of an indoor environment. The panels are separated by horizontal black lines. Each panel contains several thick, multi-colored lines (purple, blue, green, red) that outline specific objects or regions of interest. In the top panel, the lines form a U-shape around a dark object on a floor. In the middle panel, the lines define a rectangular region around a person sitting on a bench. In the bottom panel, the lines outline a doorway and a small shelf. The background in all panels is a plain, light-colored wall.

Figure 5. Comparison between our guided line search (right-hand column) and the Canny/edge linking detector (left-hand column). Our system is able to recover several subtle yet structurally significant line segments that Canny misses.

umption of a Manhattan world containing infinite floor and ceiling planes between which vertical wall segments extend, there are relatively few possible intersections to consider. Floor/wall intersections and ceiling/wall intersections generate line segments with vanishing directions on the horizon. Vertical line segments arise either when two wall segments meet or when one wall segment occludes another. In the latter case the possible vanishing directions of the wall segments on the left and right of the occlusion are constrained by the position of the intersection in the image. With these constraints, the entire set of building structures can be enumerated by branch-and-bound. For further details, the reader is referred to [3].

Lee *et al.* evaluate hypotheses by checking for consistency with pixelwise orientation estimates in the image. We extend this to use multiple images for hypothesis evaluation. A building structure  $B$  defines a unique 3D model up to an unknown scale factor  $s^*$ . Lee *et al.* ignore this parameter since they build reconstructions for single-images only. However, in our case  $s^*$  is required to transfer the model between frames.

To determine  $s^*$ , we leverage the observation that some map points will fall on the surfaces we are trying to reconstruct. We therefore devise the following voting scheme. For each map point visible in the current frame we identify the surface it falls upon within the hypothesized building

structure. Next we compute the scaling  $s^-$  such that the reconstructed surface contains that 3D point. Each map point then votes for the scale it induces on the building structure, where the votes are discretized to allow efficient vote counting. To avoid patches of detailed texture dominating other cues, each surface is allowed to vote at most once for any  $s^*$ .

Knowledge of  $s^*$  permits construction of a full 3D model and hence allows transfer of the building structure between frames. We test each building hypothesis according to its consistency with surface orientation estimates in the original frame and the  $K$  preceding frames. The surface orientation estimates are obtained on a per-frame basis using the same method as Lee *et al.*

Enforcing consistency with multiple views improves hypothesis evaluation considerably, particularly since the surface orientation estimates of any one frame are often noisy and incomplete.

## 6. Concluding Remarks

## References

- [1] G. J. Brostow, J. Shotton, J. Fauqueur, and R. Cipolla. Segmentation and recognition using structure from motion point clouds. In *ECCV '08: Proceedings of the 10th European Conference on Computer Vision*, pages 44–57, Berlin, Heidelberg, 2008. Springer-Verlag. 2 617  
618  
619  
620  
621

[2] P. Buschka and A. Saffiotti. A virtual sensor for room detection. In *Intelligent Robots and System, 2002. IEEE/RSJ International Conference on*, volume 1, pages 637–642 vol.1, 2002. 2 622  
623  
624  
625

[3] M. H. David Changsoo Lee and T. Kanade. Geometric reasoning for single image structure recovery. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2009. 1, 2, 5, 6 626  
627  
628

[4] Y. Furukawa, B. Curless, S. Seitz, and R. Szeliski. Manhattan-world stereo. *Computer Vision and Pattern Recognition, IEEE Computer Society Conference on*, 0:1422–1429, 2009. 2 629  
630  
631  
632

[5] A. Golovinskiy, V. G. Kim, and T. Funkhouser. Shape-based recognition of 3d point clouds in urban environments. In *Proc 12th IEEE Int Conf on Computer Vision*, volume 2, 2009. 2 633  
634  
635  
636

[6] S. Gould, R. Fulton, and D. Koller. Decomposing a scene into geometric and semantically consistent regions. In *Proc 12th IEEE Int Conf on Computer Vision*, volume 2, 2009. 2 637  
638  
639

[7] D. Hoiem, A. A. Efros, and M. Hébert. Geometric context from a single image. In *Proc 10th IEEE Int Conf on Computer Vision*, pages 654–661, 2005. 2 640  
641  
642

[8] G. Klein and D. Murray. Parallel tracking and mapping for small AR workspaces. In *Proc. Sixth IEEE and ACM International Symposium on Mixed and Augmented Reality (ISMAR'07)*, Nara, Japan, November 2007. 1 643  
644  
645

[9] J. Koseckà and W. Zhang. Video compass. In *Proc 7th European Conf on Computer Vision*, volume 2353 of *Lec-* 646  
647

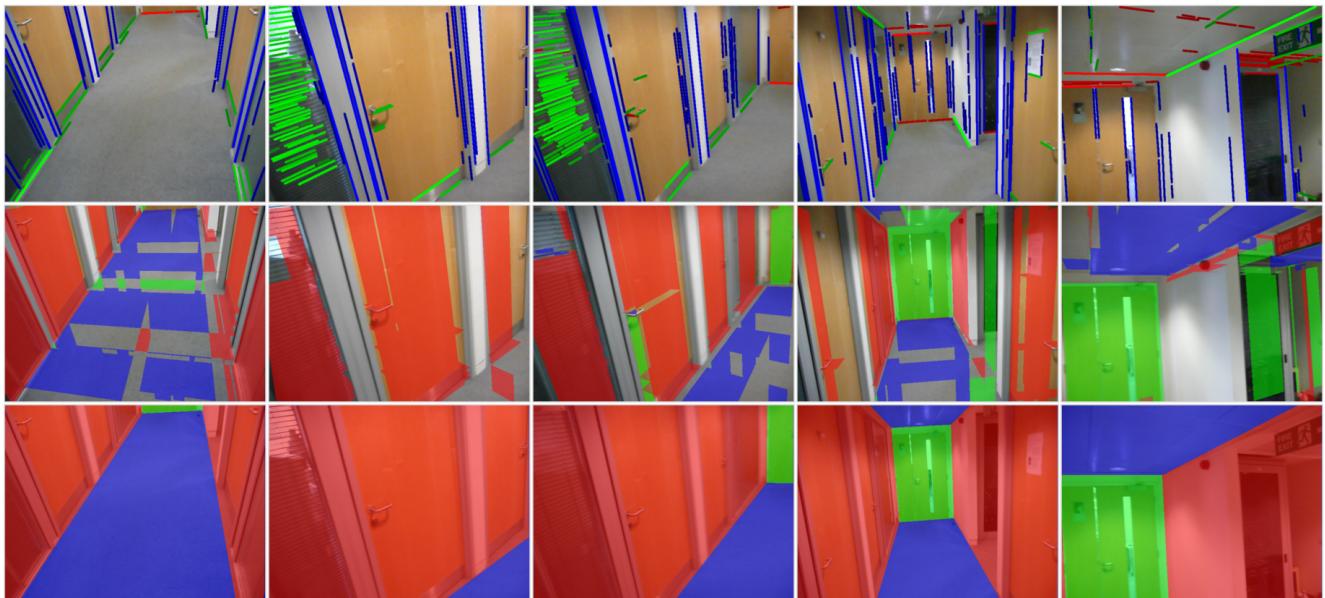


Figure 6. The final reconstruction for the “lab” sequence. The top row shows the line segments identified by the guided line search; the second row shows the surface orientation estimates from the individual frames; the bottom row shows the final model projected into five frames. The orientation estimates shown in the second row are noisy and incomplete but we are able to obtain an accurate model by combining information from all the views. The reconstruction accurately captures the primary surfaces within the scene.

ture Notes in Computer Science, pages 4: 476–490. Springer, 2002. 3, 4, 5

- [10] I. Posner, D. Schroeter, and P. Newman. Online generation of scene descriptions in urban environments. *Robot. Auton. Syst.*, 56(11):901–914, 2008. 2
- [11] A. Sarti, R. Malladi, and J. A. Sethian. Subjective surfaces: A geometric model for boundary completion. *International Journal of Computer Vision*, 46(3):201–221, 2002. 2
- [12] A. Saxena, M. Sun, and A. Y. Ng. Make3d: Learning 3d scene structure from a single still image. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31(5):824–840, 2009. 2
- [13] J. Shufelt. Performance evaluation and analysis of vanishing point detection techniques. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 21(3):282–288, Mar 1999. 2
- [14] C. Stachniss, O. Martnez-mozos, A. Rottmann, and W. Burgard. Semantic labeling of places. In *in Proceedings of the International Symposium on Robotics Research*, 2005. 2
- [15] J.-P. Tardif. Non-iterative approach for fast and accurate vanishing point detection. In *Proc 12th IEEE Int Conf on Computer Vision*, volume 2, 2009. 3
- [16] T. Tuytelaars, M. Proesmans, and L. J. V. Gool. The cascaded hough transform as support for grouping and finding vanishing points and lines. In *AFPAC '97: Proceedings of the International Workshop on Algebraic Frames for the Perception-Action Cycle*, pages 278–289, London, UK, 1997. Springer-Verlag. 2
- [17] J. Xiao and L. Quan. Multiple view semantic segmentation for street-view images. In *Proc 12th IEEE Int Conf on Computer Vision*, volume 2, 2009. 2