# Manhattan World: Compass Direction from a Single Image by Bayesian Inference

James M. Coughlan                    A.L. Yuille

Smith-Kettlewell Eye Research Institute
2318 Fillmore St.
San Francisco, CA 94115

## Abstract

*When designing computer vision systems for the blind and visually impaired it is important to determine the orientation of the user relative to the scene. We observe that most indoor and outdoor (city) scenes are designed on a Manhattan three-dimensional grid. This Manhattan grid structure puts strong constraints on the intensity gradients in the image. We demonstrate an algorithm for detecting the orientation of the user in such scenes based on Bayesian inference using statistics which we have learnt in this domain. Our algorithm requires a single input image and does not involve pre-processing stages such as edge detection and Hough grouping. We demonstrate strong experimental results on a range of indoor and outdoor images. We also show that estimating the grid structure makes it significantly easier to detect target objects which are not aligned with the grid.*

## 1    Introduction

Recently there has been growing interest in building computer vision navigational systems for the blind [9], [10]. These systems can be used, for example, for navigation and for the detection and reading of informational signs. The goal of this paper is to determine the orientation of the viewer in the scene (indoor or outdoor) from a single image. A useful spin-off is the ability to detect target objects which are not aligned with the Manhattan grid.

Most indoor and outdoor city scenes are based on a cartesian coordinate system [3, 6] which we can refer to as a Manhattan grid. This grid defines an $\vec{i}, \vec{j}, \vec{k}$ coordinate system. This gives a natural reference frame for the viewer. If the viewer can determine his/her position relative to this frame – in other words, estimate the $\vec{i}, \vec{j}$ or $\vec{k}$ directions – then it becomes significantly easier to interpret the scene. In particular, it will be a lot easier to determine the most important lines in the scene (corridor boundaries and doors, street boundaries and traffic lights) because they will typically lie in either the $\vec{i}, \vec{j}$ or $\vec{k}$ directions. Knowledge of this reference frame will make it significantly easier and faster to detect informational signs. We will assume that the camera direction lies approximately in the horizontal plane and so lines in the $\vec{k}$ direction map to approximately vertical lines in the image. There is, of course, an ambiguity in the orientations of $\vec{i}$ and $\vec{j}$ so the compass heading can only be obtained modulo 90°.

## 2    Previous Work and Three- Dimensional Geometry

There has been an enormous amount of work in projective geometry [3, 6]. Techniques from projective geometry have been applied to finding the vanishing points [1], [5]. For a recent application to vision systems for the blind see [9] for the detection of pedestrian crossings using
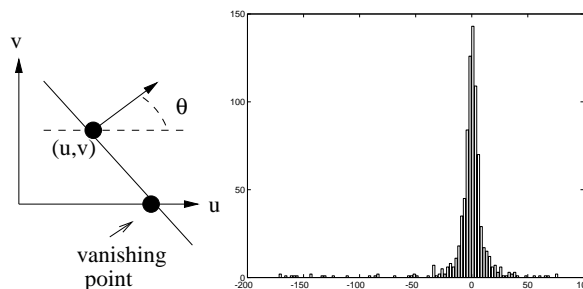
Figure 1: (Left). Geometry of an $\vec{i}$ line projected onto $(u,v)$ image plane. $\theta$ is the normal orientation of the line in the image. Because our camera is assumed to point in a horizontal direction, the vanishing point lies on the $u$ axis. (Right) Histogram of edge orientation error (displayed modulo 180°). Observe the strong peak at 0°, indicating that the image gradient direction at an edge is usually very close to the true normal orientation of the edge. We modelled this distribution using a simple box function.

projection geometry. This work, however, has typically proceeded through the stages of edge detection, Hough transforms, and finally the calculation of the geometry. Alternatively, a sequence of images over time can be used to estimate the geometry, see for example [8]. In this paper, we demonstrate that accurate results can be obtained from a single image directly without the need for techniques such as edge detection and Hough transforms.

For completeness, we give the basic geometry. We assume that the camera is oriented in the horizontal plane. This is a reasonable assumption and it turned out to be approximately correct for the images in our datasets (all of which were photographed without taking this into account). (In our current work we are relaxing this constraint to allow for any camera configuration.)

We define $\Psi$ to be the compass angle. This defines the orientation of the camera with respect to the Manhattan grid: the camera points in direction $\cos\Psi\vec{i}-\sin\Psi\vec{j}$. Camera coordinates $\vec{u}=(u,v)$ are related to the Cartesian scene coordinates $(x,y,z)$ by the equations:

$$u = \frac{f\{-x\sin\Psi - y\cos\Psi\}}{x\cos\Psi - y\sin\Psi}, \quad v = \frac{fz}{x\cos\Psi - y\sin\Psi}, \tag{1}$$

where $f$ is the focal length of the camera (which we determined to be 797 pixel units for our images).

By standard geometry, the vanishing points of lines in the $\vec{i}$ and $\vec{j}$ directions lie at $(-f\tan\Psi, 0)$ and $(f\cot\Psi, 0)$ respectively in the $(u,v)$ plane. (Lines in the $\vec{k}$ direction are all vertical in the image given our compass-world assumption.)

It is a straightforward calculation to show that a point in the image at $\vec{u}=(u,v)$ with intensity gradient at $(\cos\theta,\sin\theta)$ is *consistent with an $\vec{i}$ line in the sense that it points to the vanishing point* if $-v\tan\theta = u + f\tan\Psi$ (observe that this equation is unaffected by adding $\pm\pi$ to $\theta$ and so it does not depend on the polarity of the edge). We get a similar expression $v\tan\theta = -u + f\cot\Psi$ for lines in the $\vec{j}$ direction. (See Figure 1 (left) for an illustration of this geometry.)

## 3   $P_{on}$ and $P_{off}$: Characterizing Edges Statistically

A key element of our approach is that we do not use a binary edge map. Such edge maps make premature decisions based on too little information. (The poor quality of some of the images – underexposed and overexposed – makes edge detection particularly difficult).

Instead we use the power of statistics. Following work by Konishi *et al.* [4], we determine probabilities $P_{on}(E_{\vec{u}})$ and $P_{off}(E_{\vec{u}})$ for the probabilities of the response $E_{\vec{u}}$ of an edge filter at position $\vec{u}$ in the image *conditioned on whether we are on or off an edge*. These distri-
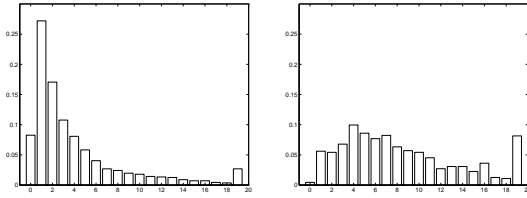
2

Figure 2: $P_{off}(y)$ (left) and $P_{on}(y)$(right), the empirical histograms of edge responses off and on edges, respectively. Here the response $y = \left|\vec{\nabla}I\right|$ is quantized to take 20 values and is shown on the horizontal axis. Note that the peak of $P_{off}(y)$ occurs at a lower edge response than the peak of $P_{on}(y)$. These distributions were very consistent for a range of images.

butions were learnt by Konishi *et al* for the Sowerby image database which contain one hundred presegmented images. The more different $P_{on}$ is from $P_{off}$ then the easier edge detection becomes, see Figure 2. A suitable measure of difference is the Chernoff Information [2] $C(P_{on}, P_{off}) = -\min_{0 \leq \lambda \leq 1} \log \sum_y P_{on}^{\lambda}(y) P_{off}^{1-\lambda}(y)$. Konishi *et al* tested a variety of different edge filters and ranked them by their effectiveness based on their Chernoff information. For this project, we chose a very simple edge detector $\left|\vec{\nabla}G_{\sigma=1} * I\right|$ – the magnitude of the gradient of the grayscale image $I$ filtered by a Gaussian $G_{\sigma=1}$ with standard deviation $\sigma = 1$ pixel units – which has a Chernoff of 0.26 nats. More effective edge detectors are available – for example, the gradient at multiple scales using colour has a Chernoff of 0.51 nats. But we do not need these more sophisticated detectors.

We extend the work of Konishi *et al* by putting probability distributions on how accurately the edge filter gradient estimates the true perpendicular direction of the edge. These were learnt for this dataset by measuring the true orientations of the edges and comparing them to those estimated from the image gradients.

This gives us distributions on the magnitude and direction of the intensity gradient $P_{on}(\vec{E}_{\vec{u}}|\beta), P_{off}(\vec{E}_{\vec{u}})$, where $\vec{E}_{\vec{u}} = (E_{\vec{u}}, \phi_{\vec{u}})$, $\beta$ is the true normal orientation of the edge, and $\phi_{\vec{u}}$ is the gradient direction measured at point $\vec{u}$. We make a *factorization assumption* that $P_{on}(\vec{E}_{\vec{u}}|\beta) = P_{on}(E_{\vec{u}})P_{ang}(\phi_{\vec{u}} - \beta)$ and $P_{off}(\vec{E}_{\vec{u}}) = P_{off}(E_{\vec{u}})U(\phi_{\vec{u}})$. $P_{ang}(.)$ (with argument evaluated modulo $2\pi$ and normalized to 1 over the range 0 to $2\pi$) is based on experimental data, see Figure 1 (right), and is peaked about 0 and $\pi$. In practice, we use a simple box function model: $P_{ang}(\delta\theta) = (1 - \epsilon)/4\tau$ if $\delta\theta$ is within angle $\tau$ of 0 or $\pi$, and $\epsilon/(2\pi - 4\tau)$ otherwise (i.e. the chance of an angular error greater than $\pm\tau$ is $\epsilon$ ). In our experiments $\epsilon = 0.1$ and $\tau = 4°$ for indoors and $6°$ outdoors. By contrast, $U(.) = 1/2\pi$ is the uniform distribution.

## 4 Bayesian Model

We devised a Bayesian model which combines knowledge of the three-dimensional geometry of the Manhattan world with statistical knowledge of edges in images. The model assumes that, while the majority of pixels in the image convey no information about camera orientation, most of the pixels with high edge responses arise from the presence of $\vec{i}, \vec{j}, \vec{k}$ lines in the three-dimensional scene. The edge orientations measured at these pixels provide constraints on the camera angle, and although the constraining evidence from any single pixel is weak, the Bayesian model allows us to pool the evidence over all pixels (both on and off edges), yielding a sharp posterior distribution on the camera angle. An important feature of the Bayesian model is that *it does not force us to decide prematurely which pixels are on and off* (or whether an on pixel is due to $\vec{i}, \vec{j}$, or $\vec{k}$), *but allows us to sum over all possible interpretations of each pixel.*

3

## 4.1  Evidence at one pixel

The image data $\vec{E}_{\vec{u}}$ at pixel $\vec{u}$ is explained by one of five models $m_{\vec{u}}$: $m_{\vec{u}} = 1, 2, 3$ mean the data is generated by an edge due to an $\vec{i}, \vec{j}, \vec{k}$ line, respectively, in the scene; $m_{\vec{u}} = 4$ means the data is generated by a random edge (not due to an $\vec{i}, \vec{j}, \vec{k}$ line); and $m_{\vec{u}} = 5$ means the pixel is off-edge. The prior probability $P(m_{\vec{u}})$ of each of the edge models was estimated empirically to be $0.02, 0.02, 0.02, 0.04, 0.9$ for $m_{\vec{u}} = 1, 2, \dots, 5$.

Using the factorization assumption mentioned before, we assume the probability of the image data $\vec{E}_{\vec{u}}$ has two factors, one for the magnitude of the edge strength and another for the edge direction:

$$P(\vec{E}_{\vec{u}}|m_{\vec{u}}, \Psi, \vec{u}) = P(E_{\vec{u}}|m_{\vec{u}})P(\phi_{\vec{u}}|m_{\vec{u}}, \Psi, \vec{u}) \qquad (2)$$

where $P(E_{\vec{u}}|m_{\vec{u}})$ equals $P_{off}(E_{\vec{u}})$ if $m_{\vec{u}} = 5$ or $P_{on}(E_{\vec{u}})$ if $m_{\vec{u}} \neq 5$. Also, $P(\phi_{\vec{u}}|m_{\vec{u}}, \Psi, \vec{u})$ equals $P_{ang}(\phi_{\vec{u}} - \theta(\Psi, m_{\vec{u}}, \vec{u}))$ if $m_{\vec{u}} = 1, 2, 3$ or $U(\phi_{\vec{u}})$ if $m_{\vec{u}} = 4, 5$. Here $\theta(\Psi, m_{\vec{u}}, \vec{u}))$ is the predicted normal orientation of lines determined by the equation $-v \tan \theta = u + f \tan \Psi$ for $\vec{i}$ lines, $v \tan \theta = -u + f \cot \Psi$ for $\vec{j}$ lines, and $\theta = 0$ for $\vec{k}$ lines.

In summary, the edge strength probability is modeled by $P_{on}$ for models 1 through 4 and by $P_{off}$ for model 5. For models 1,2 and 3 the edge orientation is modeled by a distribution which is peaked about the appropriate orientation of an $\vec{i}, \vec{j}, \vec{k}$ line predicted by the compass angle at pixel location $\vec{u}$; for models 4 and 5 the edge orientation is assumed to be uniformly distributed from 0 through $2\pi$.

Rather than decide on a particular model at each pixel, we marginalize over all five possible models (i.e. creating a mixture model):

$$P(\vec{E}_{\vec{u}}|\Psi, \vec{u}) = \sum_{m_{\vec{u}}=1}^{5} P(\vec{E}_{\vec{u}}|m_{\vec{u}}, \Psi, \vec{u})P(m_{\vec{u}}) \qquad (3)$$

In this way we can determine evidence about the camera angle $\Psi$ at each pixel without knowing which of the five model categories the pixel belongs to.

## 4.2  Evidence over all pixels

To combine evidence over all pixels in the image, denoted by $\{\vec{E}_{\vec{u}}\}$, we assume that the image data is conditionally independent across all pixels, given the compass direction $\Psi$:

$$P(\{\vec{E}_{\vec{u}}\}|\Psi) = \prod_{\vec{u}} P(\vec{E}_{\vec{u}}|\Psi, \vec{u}) \qquad (4)$$

Thus the posterior distribution on the compass direction is given by $\prod_{\vec{u}} P(\vec{E}_{\vec{u}}|\Psi, \vec{u})P(\Psi)/Z$ where $Z$ is a normalization factor and $P(\Psi)$ is a uniform prior on the compass angle.

To find the MAP (maximum a posterior) estimate, we need to maximize the log posterior term (ignoring $Z$, which is independent of $\Psi$) $\log[P(\{\vec{E}_{\vec{u}}\}|\Psi)P(\Psi)] = \log P(\Psi) + \sum_{\vec{u}} \log[\sum_{m_{\vec{u}}} P(\vec{E}_{\vec{u}}|m_{\vec{u}}, \Psi, \vec{u})P(m_{\vec{u}})]$. Our algorithm evaluates the log posterior numerically for the compass direction $\Psi$ in the range $-45°$ to $+45°$, in increments of $1°$.

# 5  Experimental Results

We tested our model on two datasets of indoor and outdoor scenes. These images were taken by an unskilled photographer unfamiliar with the goals of the study. No special attempt was made to hold the camera horizontal. The camera was set on automatic so some images are over- or under- exposed. Experiments performed by a blind user (W. Gerrey) at the Smith-Kettlewell Institute demonstrate that similar quality images can be attained by a camera mounted on the chest of a blind user (personal communication – Dr. J. Brabyn, Director of the Rehabilitation, Engineering, and Research Center, Smith-Kettlewell Eye Research Institute, San Francisco, CA 94115. 1998).

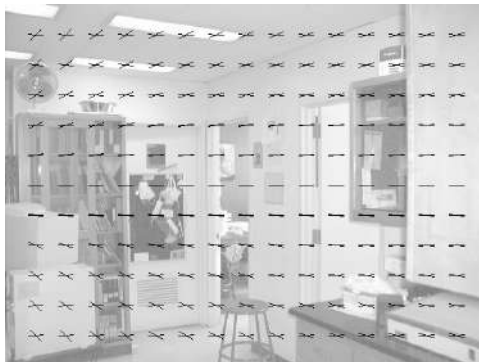Our results show strong success of our approach in both domains.

Figure 3: Estimates of the compass angle and geometry obtained by our algorithm. The estimated orientations of the $\vec{i}, \vec{j}$ lines are indicated by the black line segments drawn on the input image. At each point on a subgrid two such segments are drawn – one for $\vec{i}$ and one for $\vec{j}$. Observe how the $\vec{i}$ directions align with the wall on the right hand side and with features parallel to this wall. The $\vec{j}$ lines align with the wall on the left (and objects parallel to it). (Indoor 17).

## 5.1  Indoor Scenes

A total of twenty-five images were tested. On twenty-three images, the estimated angle was accurate to within $5°$. On two images, the orientation of the camera was far from horizontal and the estimation was poor. Examples of successes, demonstrating the range of images used, are shown in Figures 3,4,5, 6. The log posteriors for typical images, plotted as a function of $\Psi$, are shown in Figure 7.

## 5.2  Outdoor Scenes

We next tested the accuracy of estimation on outdoor scenes. Again we used twenty-five test images (taken by a naive photographer). In these scenes the vast majority of the results (twenty-two) were accurate up to $10°$. On three of the images the angles were worse than $10°$, see Figure 8. Inspection of these images showed that the log posterior had multiple peaks, see Figure 10. There was always a peak corresponding to the true compass angle (to within $10°$), however, there were false peaks which were higher in these cases. What causes these errors? Observe in Figure (8) that the vanishing point of the $\vec{i}$ lines occurs near a car *whose edges are aligned only approximately to the Manhattan grid.* The car's edges may therefore cause a small distortion in the vanishing point estimate. The correct alignment for this image can be obtained, see Figure (9), by ignoring the image data within a circle of radius 100 pixels centered around the vanishing point for each compass angle considered (this means the car will no longer contribute when evaluating the likelihood of the compass angle corresponding to the false vanishing point). Observe the difference, see Figure (10), between the log posteriors for the compass angle with and without this procedure (i.e. ignoring, or not ignoring, the circle). This new procedure, however, is intended only to show proof of concept and a thorough stability analysis is required (this is current work).

On twenty-two of the twenty-five images, however, the algorithm gave estimates accurate to $10°$ which is sufficient for the task (observe that a blind user will typically have access to a sequence of images which can be used to improve the compass estimate). See Figure 11 for a representative set of images on which the algorithm was successful.

## 6  Detecting Objects in Manhattan world

We now consider applying the Manhattan assumption to the alternative problem of detecting target objects in background clutter. To perform such a task effectively requires modelling the properties of the background clutter in addition to those of the target object. It has recently

Figure 4: Another indoor scene. Standard conventions for display of $\vec{i}, \vec{j}$ directions. Observe that the $\vec{i}, \vec{j}$ directions align with the appropriate walls despite the poor quality of the image (i.e. under-exposed). (Indoor 15).

been appreciated [7] that simple models of background clutter based on Gaussian probability distributions are often inadequate and that better performance can be obtained using alternative probability models [11].

The Manhattan world assumption gives an alternative way of probabilistically modelling background clutter. The background clutter will correspond to the regular structure of buildings and roads and its edges will be aligned to the Manhattan grid. The target object, however, is assumed to be unaligned (at least, in part) to this grid. *Therefore many of the edges of the target object will be assigned to model 4 by the algorithm.* (Note the algorithm first finds the MAP estimate $\Psi^*$ of the compass angle, see section (4), and then estimates the model by doing MAP of $P(m_{\vec{u}}|\vec{E}_{\vec{u}}, \Psi^*, \vec{u})$ to estimate $m_{\vec{u}}$ for each pixel $\vec{u}$.) This enables us to significantly simplify the detection task by removing all edges in the images except those assigned to model 4.

This idea is demonstrated in Figure (12) where the target is a bike and a robot respectively. Observe how most of the edges in the image are eliminated as target object candidates because of their alignment to the Manhattan grid. The bike and the robot stand out as outliers to the grid.

This simple example illustrates a method of modelling background clutter which we refer to as *scene clutter* because it is effectively the same as defining a probability model for the entire scene. Observe that scene clutter models require external variables – in this case the $\Psi$ angle – to determine the orientation of the viewer relative to the scene axes. These variables must be estimated to help distinguish between target and clutter. This differs from standard models used for background clutter [7],[11] where no external variable is used.

## 7  Summary and Conclusions

Our work has demonstrated proof of concept and shows the potential of our approach. The system, however, needs to be tested more extensively before it will be suited for blind users.

One obvious limitation is that we have assumed that the only unknown variable is the compass angle. This is only correct if the camera is held approximately horizontal although our results have shown robustness to this condition. It is straightforward to adjust our theory to extend the theory to estimate all three orientation angles simultaneously.

Other improvements would come from using better filters. As demonstrated by Konishi *et al* [4] the use of colour and multi-scale can give quantifiably better measures of edgeness (improving
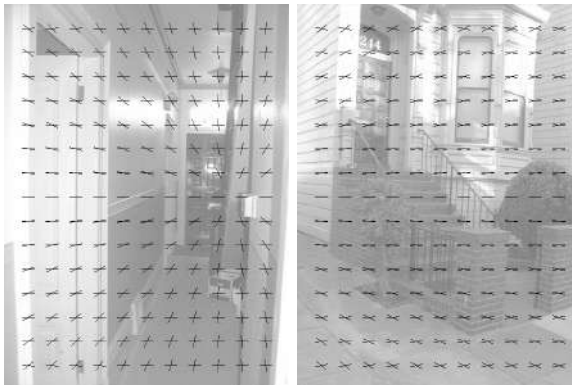
6

Figure 5: Another indoor scene and its exterior. Same conventions as above. The vanishing points are estimated to within 5° (perfectly adequate for our purposes). Note poor quality of the indoor image (i.e. over-exposed). (Indoor 23 and Outdoor 12).



Figure 6: Another indoor scene. Same conventions as above. (Indoor 8).

the Chernoff information from 0.26 nats to 0.51 nats). We anticipate that such filters will give more accurate angle estimates.

Further statistical analysis of the domains is also required. We should quantify the amount of outliers, particularly in the outdoor scenes. In particular, we should investigate the number of structured outliers and determine techniques to detect them. In addition, we should use error analysis to improve our estimates of the probability distributions and, in particular, to see how the angle errors change as a function of distance from a vanishing point. This will enable us to do performance analysis such as estimating Cramer-Rao lower bounds for the accuracy of the estimates.

We should mention the issues of algorithmic speed. At present the algorithm takes a minute which is too slow for practical use. However, this is for unoptimized code when it is run on images of size $640 \times 480$. Optimizing the code (e.g. by using look-up tables to pre-compute trigonometric functions) and subsampling the image will allow the algorithm to work significantly faster. Other techniques involve rejecting image pixels where the edge detector response is so low that there is no realistic chance of an edge being present. This would mean that at least 70% of the image pixels could be removed from the computation. We observe that the algorithm is entirely parallelizable. Overall, there seems little difficulty in getting this algorithm to work in a few seconds –which is perfectly adequate for blind users.
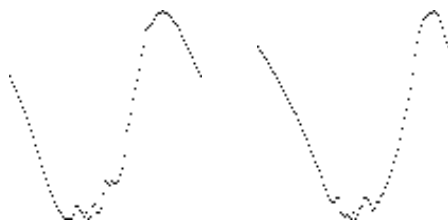
Figure 7: The log posteriors as a function of compass angle (from $-45°$ to $45°$ along the horizontal axis) for images Indoor 17 (left) and Indoor 15 (right). These results are typical for both the indoor and outdoor dataset. See Figure 10 for an exception where there are multiple peaks.
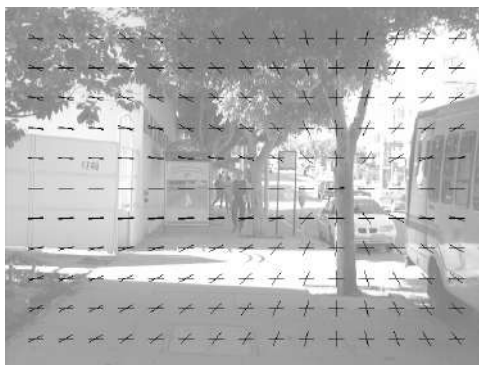


Figure 8: Incorrect estimation of compass angle for outdoor scene. The algorithm computes the vanishing point to be more than $10°$ to the right of the true vanishing point. (Outdoor 35).

## Acknowledgments

## References

[1] B. Brillault-O'Mahony. "New Method for Vanishing Point Detection". *Computer Vision, Graphics, and Image Processing.* 54(2). pp 289-300. 1991.

[2] T. M. Cover and J. A. Thomas. *Elements of Information Theory.* Wiley Interscience Press. New York. 1991.

[3] O.D. Faugeras. **Three-Dimensional Computer Vision.** MIT Press. 1993.

[4] S. Konishi, A. L. Yuille, J. M. Coughlan, and S. C. Zhu. "Fundamental Bounds on Edge Detection: An Information Theoretic Evaluation of Different Edge Cues." *Proc. Int'l conf. on Computer Vision and Pattern Recognition*, 1999.

[5] E. Lutton, H. Maître, and J. Lopez-Krahe. "Contribution to the determination of vanishing points using Hough transform". *IEEE Trans. on Pattern Analysis and Machine Intelligence.* 16(4). pp 430-438. 1994.

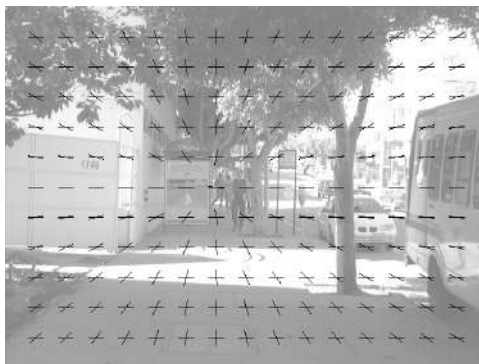[6] J.L. Mundy and A. Zisserman. (Eds). **Geometric Invariants in Computer Vision**. MIT Press. 1992.

Figure 9: A correct estimate of the compass angle for the previous figure can be obtained by ignoring data from image points within a circle of radius 100 pixels centered about the vanishing point for each compass angle considered.



Figure 10: Log posterior as a function of compass angle for the previous two figures. Observe that for these images the log posterior has multiple peaks. For the original algorithm, the false peak had higher probability (left). For the modified algorithm which ignores the central circle of data (right) the true peak is higher.

[7] J. A. Ratches, C. P. Walters, R. G. Buser and B. D. Guenther. "Aided and Automatic Target Recognition Based upon Sensory Inputs from Image Forming Systems". *IEEE Trans. on PAMI*, vol. 19, No. 9, Sept. 1997.

[8] P. Torr and A. Zisserman. "Robust Computation and Parameterization of Multiple View Relations". In *Proceedings of the International Conference on Computer Vision.* ICCV'98. Bombay, India. pp 727-732. 1998.

[9] S. Utcke. "Grouping based on Projective Geometry Constraints and Uncertainty". In *Proceedings of the International Conference on Computer Vision.* ICCV'98. Bombay, India. pp 739-746. 1998.

[10] A.L. Yuille, D. Snow, and M. Nitzberg. "Signfinder". In *Proceedings of the International Conference of Computer Vision.* (ICCV'98). Bombay. India. January 1998.

[11] S. C. Zhu, A. Lanterman, and M. I. Miller. "Clutter Modeling and Performance Analysis in Automatic Target Recognition". In *Proceedings Workshop on Detection and Classification of Difficult Targets.* Redstone Arsenal, Alabama. 1998.

Figure 11: Results on four outdoor images. Same conventions as before. Observe the accuracy of the $\vec{i}, \vec{j}$ projections in these varied scenes despite the poor quality of some of the images.
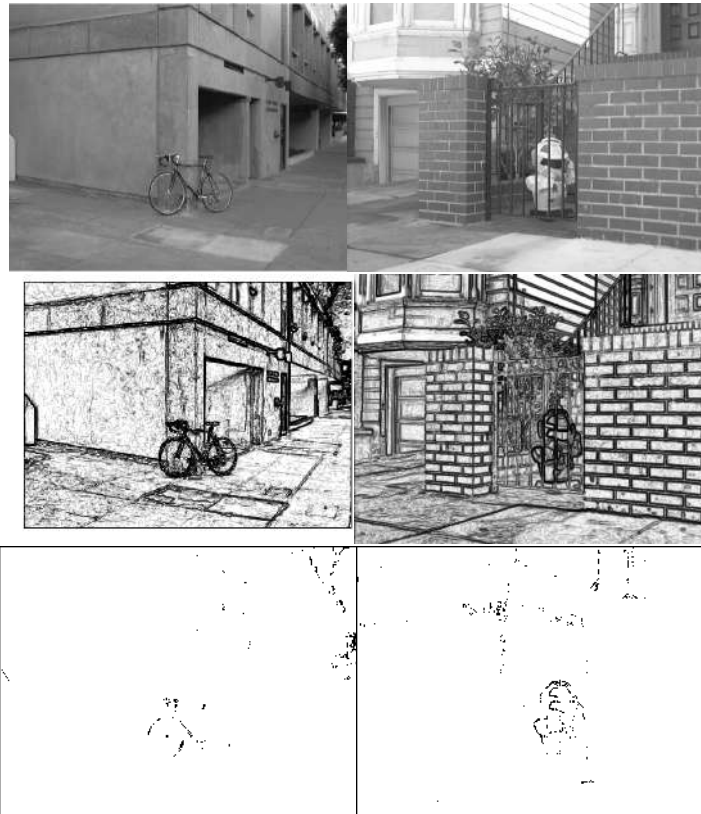


Figure 12: Detecting bikes (left column) and robots (right column) in urban scenes. The original image (top row) and the edge maps (centre row) computed as $\log P_{on}(E_{\vec{u}})/P_{off}(E_{\vec{u}})$ – see Konishi *et al* 1999 – displayed as a grey-scale image where black is high and white is low. In the bottom row we show the edges assigned to model 4 (i.e. the outliers) in black. Observe that the edges of the bike and the robot are now highly salient (and make detection straightforward) because most of them are unaligned to the Manhattan grid.

10