

# Putting Objects in Perspective

Derek Hoiem

Alexei A. Efros

Martial Hebert

Carnegie Mellon University, Robotics Institute  
{dhoiem,efros,hebert}@cs.cmu.edu

## Abstract

*Image understanding requires not only individually estimating elements of the visual world but also capturing the interplay among them. In this paper, we provide a framework for placing local object detection in the context of the overall 3D scene by modeling the interdependence of objects, surface orientations, and camera viewpoint.*

*Most object detection methods consider all scales and locations in the image as equally likely. We show that with probabilistic estimates of 3D geometry, both in terms of surfaces and world coordinates, we can put objects into perspective and model the scale and location variance in the image. Our approach reflects the cyclical nature of the problem by allowing probabilistic object hypotheses to refine geometry and vice-versa. Our framework allows painless substitution of almost any object detector and is easily extended to include other aspects of image understanding. Our results confirm the benefits of our integrated approach.*

## 1. Introduction

Consider the street scene depicted on Figure 1. Most people will have little trouble seeing that the green box in the middle contains a car. This is despite the fact that, shown in isolation, these same pixels can just as easily be interpreted as a person's shoulder, a mouse, a stack of books, a balcony, or a million other things! Yet, when we look at the entire scene, all ambiguity is resolved – the car is unmistakably a car. How do we do this?

There is strong psychophysical evidence (e.g. [3, 25]) that context plays a crucial role in scene understanding. In our example, the car-like blob is recognized as a car because: 1) it's sitting on the road, and 2) it's the "right" size, relative to other objects in the scene (cars, buildings, pedestrians, etc). Of course, the trouble is that everything is tightly interconnected – a visual object that uses others as its context will, in turn, be used as context by these other objects. We recognize a car because it's on the road. But how do we recognize a road? – because there are cars! How does one attack this chicken-and-egg problem? What is the right framework for connecting all these pieces of the recognition puzzle in a coherent and tractable manner?

In this paper we will propose a unified approach for modeling the contextual symbiosis between three crucial ele-



Figure 1. General object recognition cannot be solved locally, but requires the interpretation of the entire image. In the above image, it's virtually impossible to recognize the car, the person and the road in isolation, but taken together they form a coherent visual story. Our paper tells this story.

ments required for scene understanding: low-level object detectors, rough 3D scene geometry, and approximate camera position/orientation. Our main insight is to model the contextual relationships between the visual elements, *not in the 2D image plane* where they have been projected by the camera, but *within the 3D world* where they actually reside. Perspective projection obscures the relationships that are present in the actual scene: a nearby car will appear much bigger than a car far away, even though in reality they are the same height. We “undo” the perspective projection and analyze the objects in the space of the 3D scene.

### 1.1. Background

In its early days, computer vision had but a single grand goal: to provide a complete semantic interpretation of an input image by reasoning about the 3D scene that generated it. Indeed, by the late 1970s there were several image understanding systems being developed, including such

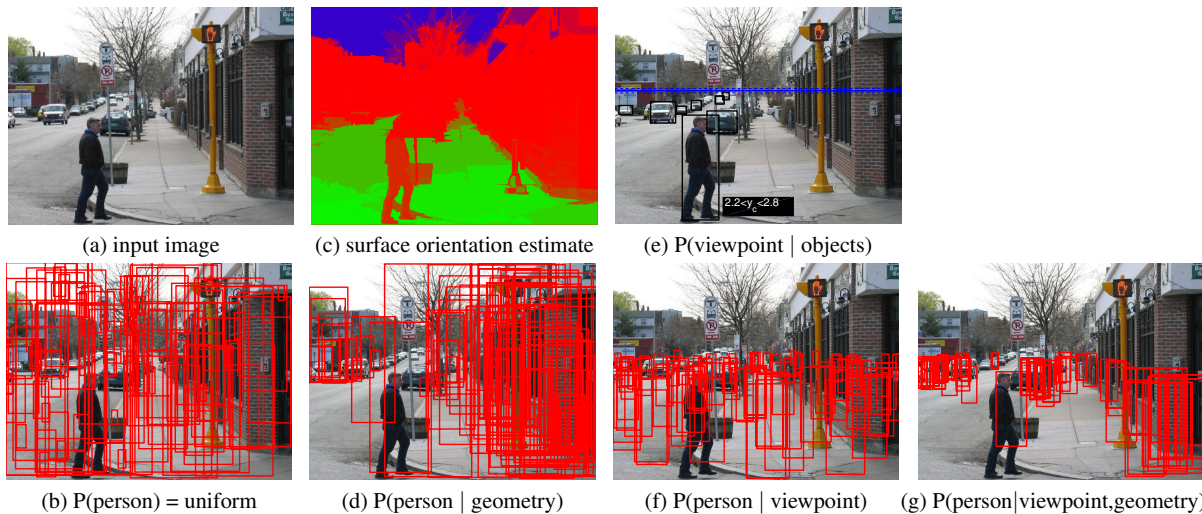


Figure 2. Watch for pedestrians! In (b,d,f,g), we show 100 boxes sampled according to the available information. Given an input image (a), a local object detector will expect to find a pedestrian at any location/scale (b). However, given an estimate of rough surface orientations (c), we can better predict where a pedestrian is likely to be (d). We can estimate the camera viewpoint (e) from a few known objects in the image. Conversely, knowing the camera viewpoint can help in predict the likely scale of a pedestrian (f). The combined evidence from surface geometry and camera viewpoint provides a powerful predictor of where a pedestrian might be (g), before we even run a pedestrian detector! Red, green, and blue channels of (c) indicate confidence in vertical, ground, and sky, respectively. Best viewed in color.

pioneering work as Brooks' *ACRONYM* [4], Hanson and Riseman's *VISIONS* [9], Ohta and Kanade's outdoor scene understanding system [19], Barrow and Tenenbaum's intrinsic images [2], etc. For example, *VISIONS* was an extremely ambitious system that analyzed a scene on many interrelated levels including segments, 3D surfaces and volumes, objects, and scene categories. However, because of the heavy use of heuristics, none of these early systems were particularly successful, which led people to doubt the very goal of complete image understanding.

We believe that the vision pioneers were simply ahead of their time. They had no choice but to rely on heuristics because they lacked the computational resources to *learn* the relationships governing the structure of our visual world. The advancement of learning methods in the last decade brings renewed hope for a complete image understanding solution. However, the currently popular learning approaches are based on looking at small image windows at all locations and scales to find specific objects. This works wonderfully for face detection [23, 29] (since the inside of a face is much more important than the boundary) but is quite unreliable for other types of objects, such as cars and pedestrians, especially at the smaller scales.

As a result, several researchers have recently begun to consider the use of contextual information for object detection. The main focus has been on modeling direct relationships between objects and other objects [15, 18], regions [10, 16, 28] or scene categories [18, 24], all within the 2D image plane. Going beyond the 2D image plane, Hoiem *et al.* [11] propose a mechanism for estimating rough 3D scene

geometry from a single image and use this information as additional features to improve object detection. From low-level image cues, Torralba and Oliva [26] get a sense of the viewpoint and mean scene depth, which provides a useful prior for object detection [27]. Forsyth *et al.* [7] describe a method for geometric consistency of object hypotheses in simple scenes using hard algebraic constraints. Others have also modeled the relationship between the camera parameters and objects, requiring either a well-calibrated camera (e.g. [12]), a stationary surveillance camera (e.g. [14]), or both [8].

In this work, we draw on several of the previous techniques: local object detection (based on Murphy *et al.* [18]), 3D scene geometry estimation [11], and camera viewpoint estimation. Our contribution is a statistical framework that allows *simultaneous* inference of object identities, surface orientations, and camera viewpoint using a *single image* taken from an uncalibrated camera.

## 1.2. Overview

To evaluate our approach, we have chosen a very challenging dataset of outdoor images [22] that contain cars and people, often partly occluded, over an extremely wide range of scales and in accidental poses (unlike, for example, the framed photographs in Corel or CalTech datasets). Our goal is to demonstrate that substantial improvement over standard low-level detectors can be obtained by reasoning about the underlying 3D scene structure.

One way to think about what we are trying to achieve is to consider the likely places in an image where an ob-

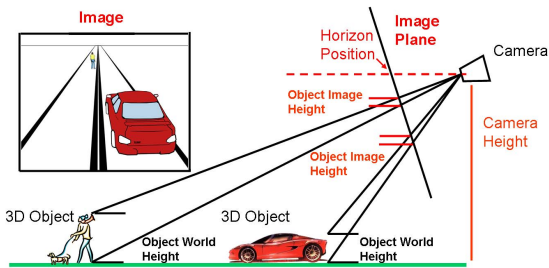


Figure 3. An object's height in the image can be determined from its height in the world and the viewpoint.

ject (e.g. a pedestrian) could be found (Figure 2). Without considering the 3D structure of the scene, all image positions and scales are equally likely (Figure 2b) – this is what most object detectors assume. But if we can estimate the rough surface geometry in the scene, this information can be used to adjust the probability of finding a pedestrian at a given image location (Figure 2d). Likewise, having an estimate of the camera viewpoint (height and horizon position) supplies the likely scale of an object in the image (Figure 2f). Combining these two geometric cues together gives us a rather tight prior likelihood for the location and scale of a pedestrian, as in Figure 2g. This example is particularly interesting because this is still only a prior – we have not applied a pedestrian detector yet. Notice, as well, that the pattern of expected pedestrian detections is very reminiscent of typical human eye-tracking experiments, where subjects are asked to search for a person in an image.

Of course, just as scene and camera geometry can influence object detection, so can the detected objects alter the geometry estimation. For example, if we know the locations/scales of some of the objects in the image, we can use this to better estimate the camera viewpoint parameters (see the 90% confidence bounds in Figure 2e). In general, our aim is to combine all these pieces of evidence into a single coherent image interpretation framework.

The rest of the paper will be devoted to exploring our two primary conjectures: 1) 3D reasoning improves object detection, even when using a single image from an uncalibrated camera, and 2) the more fully the scene is modeled (more properties, more objects), the better the estimates will be. We will first describe the mathematics of projective geometry as it relates to our problem (Section 2). We will then define the probabilistic model used for describing the relationships within the 3D scene (Section 3) and how it can be learned (Section 4). Finally, we present quantitative and qualitative results demonstrating the performance of our system on a difficult dataset (Section 5).

## 2. Scene Projection

Under a zero-skew, unit aspect ratio perspective camera model, we can compute a grounded object's height in the

scene, given only the camera height and horizon line (see Figure 3). First, let's rotate and translate our image coordinates  $(u, v)$  to the coordinates  $(\hat{u}, \hat{v})$  so that  $\hat{v} = 0$  for every point on the horizon and  $\hat{v} > 0$  for every point below the horizon. The world height  $y$  of a point can be recovered from  $\hat{v} = (y_c - y) \frac{f}{z}$  where  $y_c$  is the camera height,  $z$  is the depth, and  $f$  is the camera focal length. Without loss of generality, we define the object to rest on the plane  $y = 0$ . The object's height can be recovered from  $\frac{\hat{v}_1}{\hat{v}_2 - \hat{v}_1} = \frac{y_c}{y}$ , where  $\hat{v}_1$  is the bottom and  $\hat{v}_2$  is the top of the object. To get  $\hat{v}$  from pixel coordinates, we simply compute the distance of the horizon line to the point. In this paper, since photographs typically have little roll, we define the horizon line by the image row  $v_0$ . Letting  $v_i$  and  $h_i$  denote the bottom position and height of an object in the image, we have the following relationship:

$$y_i = \frac{h_i y_c}{v_i - v_0}. \quad (1)$$

From equation 1, we can compute the image height  $h_i$  of an object given its image position  $v_i$ , 3D height  $y_i$ , and the viewpoint  $(v_0, y_c)$ . Of course, for an uncalibrated camera, we do not know the viewpoint or the object's true height *a priori*. However, since people do not take photos in a completely random manner and since objects have a small range of possible 3D sizes, we can estimate an informative distribution of viewpoint and object size and, from it, derive a distribution for  $h_i$  given  $v_i$ .

In our paper, we assume that all objects of interest rest on the ground plane. While this assumption may seem restrictive (cannot find people on the rooftops), humans seem to make the same assumption (we don't notice the security standing on the rooftops at political rallies unless we specifically look for them). If the ground is sloped, as in Figure 2, the coordinates and parameters are computed with respect to that slope, and the relationship between viewpoint and objects in the image still holds.

## 3. Modeling the Scene

We want to determine the viewpoint, object identities, and surface geometry of the scene from an image. We could estimate each independently, but our estimates will be much more accurate if we take advantage of the interactions between the scene elements. We consider the objects (e.g., cars, pedestrians, background) and geometric surfaces to each produce image evidence. The viewpoint, defined by the horizon position in the image and the camera height, directly affects the position and size of the objects in the image. In turn, the objects directly affect nearby geometric surfaces. We assume that local geometric surfaces are independent given their corresponding object identities and that the object identities are independent given the viewpoint. In Figure 4, we represent these conditional independence

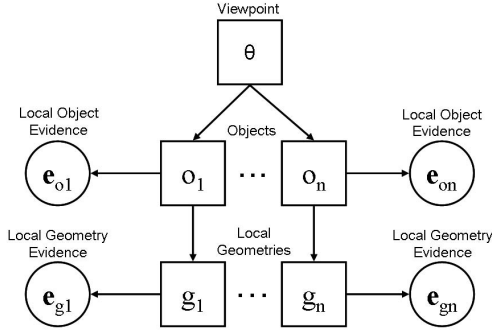


Figure 4. Graphical model of conditional independence for viewpoint  $\theta$ , object identities  $\mathbf{o}$ , and the 3D geometry of surfaces  $\mathbf{g}$  surrounding the objects. Viewpoint describes the horizon position in the image and the height of the camera in the 3D scene (in relation to the objects of interest). Each image has  $n$  object hypotheses, where  $n$  varies by image. The object hypothesis  $o_i$  involves assigning an identity (e.g., pedestrian or background) and a bounding box. The surface geometry  $g_i$  describes the 3D orientations of the  $i^{th}$  object surface and nearby surfaces in the scene.

assumptions in a graphical model, denoting objects as  $\mathbf{o}$ , surface geometries as  $\mathbf{g}$ , object evidence as  $\mathbf{e}_o$ , geometry evidence as  $\mathbf{e}_g$ , and the viewpoint as  $\theta$ .

Our model implies the following decomposition:

$$P(\theta, \mathbf{o}, \mathbf{g} | \mathbf{e}) \propto P(\theta) \prod_i P(o_i | \theta) \frac{P(o_i | \mathbf{e}_o)}{P(o_i)} P(g_i | o_i) \frac{P(g_i | \mathbf{e}_g)}{P(g_i)} \quad (2)$$

The proportionality Equation 2 is with respect to terms of the observed evidence ( $\mathbf{e} = \{\mathbf{e}_o, \mathbf{e}_g\}$ ) that are constant within an image.

Our approach allows other researchers to easily integrate their own detectors into our framework. When interactions among elements of the scene are defined, each addition to the framework adds evidence that can be used to improve estimation of the other elements.

### 3.1. Viewpoint

The viewpoint  $\theta$  involves two variables: the horizon position in the image  $v_0$  and the camera height (in meters)  $y_c$ . We consider camera height and horizon position to be independent *a priori* so that  $P(\theta) = P(v_0)P(y_c)$ . We investigated using image statistics, including vanishing points [13] and surface geometry [11], as evidence for the horizon position but found that a simple Gaussian prior performed just as well. Similarly, for the camera height  $y_c$ , we estimate a prior distribution using kernel density estimation over the  $y_c$  values (computed based on objects of known height in the scene) in a set of training images.

Figure 5 displays the viewpoint prior (e) and an example of the revised likelihood (f) when object and surface geometry evidences are considered. *A priori*, the most likely camera height is 1.67m, which happens to be eye level for a

typical adult male, and the most likely horizon position is 0.50. While the viewpoint prior does have high variance, it is much more informative than the uniform distribution that is implicitly assumed when scale is considered irrelevant.

### 3.2. Objects

An object candidate  $o_i$  consists of a type  $t_i \in \{object, background\}$  (e.g. “pedestrian”) and a bounding box  $bbox_i = \{u_i, v_i, w_i, h_i\}$  (lower-left coordinate, width, and height, respectively). The object term of our scene model is composed as follows:

$$P(o_i | \mathbf{e}_o, \theta) \propto \frac{P(o_i | \mathbf{e}_o)}{P(o_i)} P(o_i | \theta) \quad (3)$$

Our window-based object detector outputs the class-conditional log-likelihood ratio at each position and scale (with discrete steps) in the image. From these ratios and a prior  $P(o_i)$ , we can compute the probability of an object occurring at a particular location/scale:

$$P(t_i = obj, bbox_i | \mathbf{I}_i) = \frac{1}{1 + \exp[-c_i - \log \frac{P(o_i)}{1 - P(o_i)}]} \quad (4)$$

where  $c_i$  is the log-likelihood ratio estimated by the detector, based on local image information  $\mathbf{I}_i$  at the  $i^{th}$  bounding box. Typically, researchers perform non-maxima suppression, assuming that high detection responses at neighboring positions could be due to an object at either of those positions (but not both). Making the same assumption, we also apply non-maxima suppression, but we form a point distribution out of the non-maxima, rather than discarding them. An object candidate is formed out of a group of closely overlapping bounding boxes.<sup>1</sup> The candidate’s likelihood  $P(t_i = obj | \mathbf{e}_o)$  is equal to the likelihood of the highest-confidence bounding box, and the likelihoods of the locations given the object identity  $P(bbox_i | t_i = obj, \mathbf{e}_o)$  are directly proportional to  $P(t_i = obj, bbox_i | \mathbf{I})$ . After thresholding to remove detections with very low confidences from consideration, a typical image will contain several dozen object candidates (determining  $n$ ), each of which has tens to hundreds of possible position/shapes.

An object’s height depends on its position when given the viewpoint. Formally,  $P(o_i | \theta) \propto p(h_i | t_i, v_i, \theta)$  (the proportionality is due to the uniformity of  $P(t_i, v_i, w_i | \theta)$ ). From Equation 1, if  $y_i$  is normal, with parameters  $\{\mu_i, \sigma_i\}$ , then  $h_i$  conditioned on  $\{t_i, v_i, \theta\}$  is also normal, with parameters  $\frac{\mu_i y_i (v_o - v_i)}{y_c}$  and  $\frac{\sigma_i y_i (v_o - v_i)}{y_c}$ .

### 3.3. Surface Geometry

Most objects of interest can be considered as vertical surfaces supported by the ground plane. Estimates of the local

<sup>1</sup>Each detector distinguishes between one object type and background in our implementation. Separate candidates are created for each type of object.



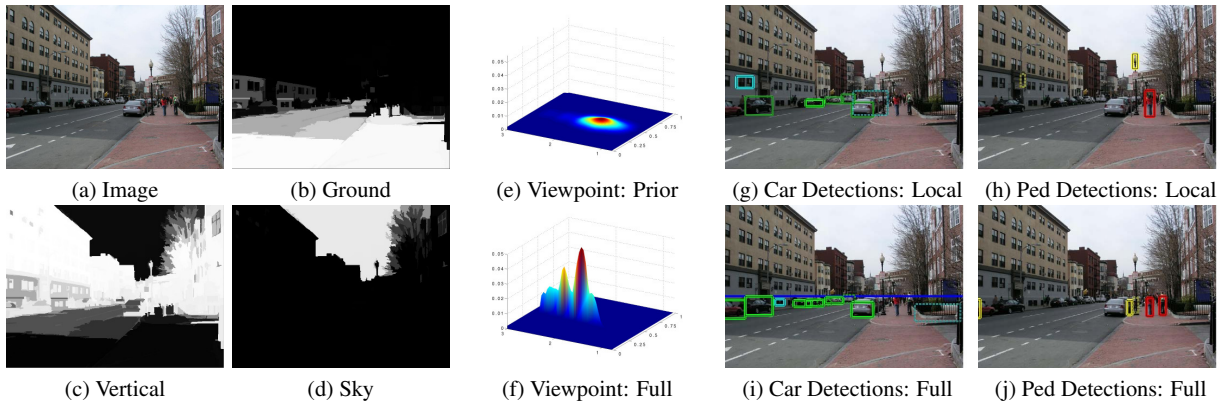


Figure 5. We begin with geometry estimates (b,c,d), local object detection confidences (g,h), and a prior (e) on the viewpoint. Using our model, we improve our estimates of the viewpoint (f) and objects (i,j). In the viewpoint plots, the left axis is camera height (meters), and the right axis is horizon position (measured from the image bottom). The viewpoint peak likelihood increases from 0.0037 *a priori* to 0.0503 after inference. At roughly the same false positive (cars:cyan, peds:yellow) rate, the true detection (cars:green, peds:red) rate doubles when the scene is coherently modeled.

surface geometry could, therefore, provide additional evidence for objects. To obtain the rough 3D surface orientations in the image, we apply the method of [11] (we use the publicly available executable), which produces confidence maps for three main classes: “ground”, “vertical”, and “sky”, and five subclasses of “vertical”: planar, facing “left”, “center”, and “right”, and non-planar “solid” and “porous”. Figure 5(b,c,d) displays the confidence maps for the three main surface labels.

We define  $g_i$  to have three values corresponding to whether the object surface is visible in the detection window and, if so, whether the ground is visible just below the detection window. For example, we consider a car’s geometric surface to be planar or non-planar solid and a pedestrian’s surface to be non-planar solid. We can compute  $P(g_i|o_i)$  and  $P(g_i)$  by counting occurrences of each value of  $g_i$  in a training set. If  $o_i$  is background, we consider  $P(g_i|o_i) \approx P(g_i)$ . We estimate  $P(g_i|e_g)$  based on the confidence maps of the geometric surfaces. In experiments, we found that the average geometric confidence in a window is a well-calibrated probability for the geometric value.

### 3.4. Inference

Inference is well-understood for tree-structured graphs like our model (Figure 4). We use Pearl’s belief propagation<sup>2</sup> algorithm [20] from the Bayes Net Toolbox [17]. Once the model is defined and its parameters estimated, as described above, it can answer queries, such as “What is the expected height of this object?” or “What are the marginal probabilities for cars?” or “What is the most probable

<sup>2</sup>To simplify the BP algorithm, we quantize all continuous variables ( $v_0$  and  $y_c$  into 50 and 100 evenly-spaced bins);  $o_i$  is already discrete due to sliding window detection.

explanation of the scene?”. In this paper, we report results based on marginal probabilities from the sum-product algorithm. Figure 5 shows how local detections (g,h) improve when viewpoint and surface geometry are considered (i,j).

## 4. Training

**Viewpoint.** To estimate the priors for  $\theta$ , we manually labeled the horizon in 60 outdoor images from the LabelMe database [22]. In each image, we labeled cars (including vans and trucks) and pedestrians (defined as an upright person) and computed the maximum likelihood estimate of the camera height based on the labeled horizon and the height distributions of cars and people in the world. We then estimated the prior for camera height using kernel density estimation (`ksdensity` in Matlab).

**Objects.** Our baseline car and pedestrian detector uses a method similar to the local detector of Murphy, Torralba, and Freeman [18]. We used the same local patch template features but added six color features that encode the average  $L^*a^*b$  color of the detection window and the difference between the detection window and the surrounding area. The classifier uses a logistic regression version of Adaboost [5] to boost eight-node decision tree classifiers. For cars, we trained two views (front/back: 32x24 pixels and side: 40x16 pixels), and for pedestrians, we trained one view (16x40 pixels). Each were trained using the full PASCAL dataset [1].

To verify that our baseline detector has reasonable performance, we trained a car detector on the PASCAL challenge training/validation set, and evaluated the images in test set 1 using the criteria prescribed for the official competition. For the sake of comparison in this validation experiment, we did not search for cars shorter than 10% of

the image height, since most of the official entries could not detect small cars. We obtain an average precision of 0.423 which is comparable to the best scores reported by the top 3 groups: 0.613, 0.489, and 0.353.

To estimate the height distribution of cars (in the 3D world), we used Consumer Reports ([www.consumerreports.org](http://www.consumerreports.org)) and, for pedestrians, used data from the National Center for Health Statistics ([www.cdc.gov/nchs/](http://www.cdc.gov/nchs/)). For cars, we estimated a mean of 1.59m and a standard deviation of 0.21m. For adult humans, the mean height is 1.7m with a standard deviation of 0.085m. Alternatively, the distribution of (relative) object heights and camera heights could be learned simultaneously using the EM algorithm if the training set includes images that contain multiple objects.

**Surface Geometry.**  $P(g_i|o_i)$  was found by counting the occurrences of the values of  $g_i$  for both people and cars in the 60 training images from LabelMe. We set  $P(g_i)$  to be uniform, because we found experimentally that learned values for  $P(g_i)$  resulted in the system over-relying on geometry. This over-reliance may be due to our labeled images (general outdoor) being drawn from a different distribution than our test set (streets of Boston) or to the lack of a modeled direct dependence between surface geometries. Further investigation is required.

## 5. Evaluation

Our test set consists of 422 random outdoor images from the LabelMe dataset [22]. The busy city streets, sidewalks, parking lots, and roads provide realistic environments for testing car and pedestrian detectors, and the wide variety of object pose and size and the frequency of occlusions make detection extremely challenging. In the dataset, 60 images have no cars or pedestrians, 44 have only pedestrians, 94 have only cars, and 224 have both cars and pedestrians. In total, the images contain 923 cars and 720 pedestrians.

We detect cars with heights as small as 14 pixels and pedestrians as small as 36 pixels tall. To get detection confidences for each window, we reverse the process described in Section 3.2. We then determine the bounding boxes of objects in the standard way, by thresholding the confidences and performing non-maxima suppression.

Our goal in these experiments is to show that, by modeling the interactions among several aspects of the scene and inferring their likelihoods together, we can do much better than if we estimate each one individually.

**Object Detection Results.** Figure 6 plots the ROC curves for car and pedestrian detection on our test set when different subsets of the model are considered. Figure 7 displays and discusses several examples. To provide an estimate of how much other detectors may improve under

	Cars			Pedestrians		
	1FP	5FP	10FP	1FP	5FP	10FP
+Geom	6.6%	5.6%	7.0%	7.5%	8.5%	17%
+View	8.2%	16%	22%	3.2%	14%	23%
+GeomView	<b>12%</b>	<b>22%</b>	<b>35%</b>	<b>7.2%</b>	<b>23%</b>	<b>40%</b>

Table 1. Modeling viewpoint and surface geometry aids object detection. Shown are percentage reductions in the missed detection rate while fixing the number of false positives per image.

	Mean	Median
Prior	10.0%	8.5%
+Obj	7.5%	4.5%
+ObjGeom	7.0%	3.8%

Table 2. Object and geometry evidence improve horizon estimation. Mean/median absolute error (as percentage of image height) are shown for horizon estimates.

	Horizon	Cars (FP)		Ped (FP)	
Car	7.3%	5.6	7.4	—	—
Ped	5.0%	—	—	12.4	13.7
Car+Ped	3.8%	5.0	6.6	11.0	10.7

Table 3. Horizon estimation and object detection are more accurate when more object models are known. Results shown are using the full model in three cases: detecting only cars, only pedestrians, and both. The horizon column shows the median absolute error. For object detection we include the number of false positives per image at the 50% detection rate computed over all images (first number) and the subset of images that contain both cars and people (second number).

our framework, we report the percent reduction in false negatives for varying false positive rates in Table 1. When the viewpoint and surface geometry are considered, about 20% of cars and pedestrians missed by the baseline are detected for the same false positive rate! The improvement due to considering the viewpoint is especially amazing, since the viewpoint uses no direct image evidence. Also note that, while individual use of surface geometry estimates and the viewpoint provides improvement, using both together improves results further.

**Horizon Estimation Results.** By performing inference over our model, the object and geometry evidence can also be used to improve the horizon estimates. We manually labeled the horizon in 100 of our images that contained both types of objects. Table 2 gives the mean and median absolute error over these images. Our prior of 0.50 results in a median error of 0.085% of the image height, but when objects and surface geometry are considered, the median error reduces to 0.038%. Notice how the geometry evidence provides a substantial improvement in horizon estimation, even though it is separated from the viewpoint by two variables in our model.

**More is Better.** Intuitively, the more types of objects that we can identify, the better our horizon estimates will

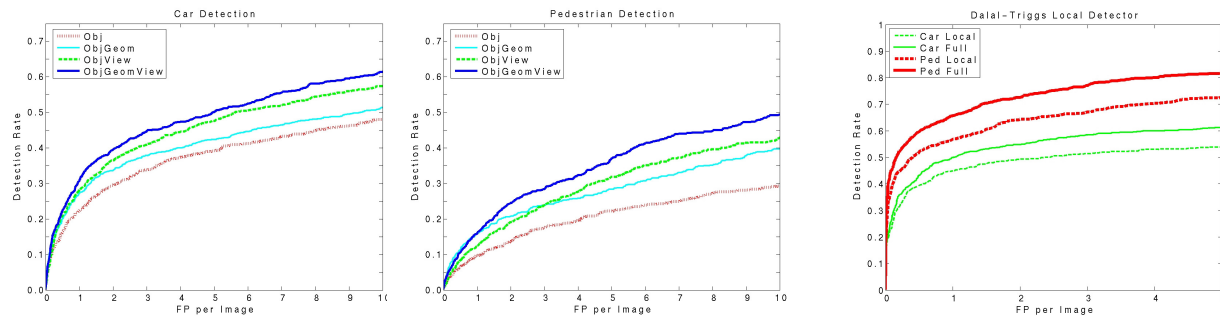


Figure 6. Considering viewpoint and surface geometry improves results over purely local object detection. The left two plots show object detection results using only local object evidence (Obj), object and geometry evidence (ObjGeom), objects related through the viewpoint (ObjView), and the full model (ObjViewGeom). On the right, we plot results using the Dalal-Triggs local detector [6].

be, leading to improved object detection. We verify this experimentally, performing the inference with only car detection, only pedestrian detection, and both. Table 3 gives the accuracy for horizon estimation and object detection when only cars are detected, when only pedestrians are detected, and when both are detected. As predicted, detecting two objects provides better horizon estimation and object detection than detecting one.

**Dalal-Triggs Detector.** To support our claim that any local object detector can be easily improved by plugging it into our framework, we performed experiments using the Dalal-Triggs detector [6] after converting the SVM outputs to probabilities using the method of [21]. We used code, data, and parameters provided by the authors, training an 80x24 car detector and 32x96 and 16x48 (for big and small) pedestrian detectors. The Dalal-Triggs local detector is currently among the most accurate for pedestrians, but its accuracy (Figure 6) improves considerably with our framework, from 57% to 66% detections at 1 FP per image.

## 6. Discussion

In this paper, we have provided a “skeleton” model of a scene – a tree structure of camera, objects, and surface geometry. Our model-based approach has two main advantages over the more direct “bag of features/black box” classification method: 1) subtle relationships (such as that object sizes relate through the viewpoint) can be easily represented; and 2) additions and extensions to the model are easy (the direct method requires complete retraining whenever anything changes).

To add a new object to our model, one needs only to train a detector for that object and supply the distribution of the object’s height in the 3D scene. Our framework could also be extended by modeling other scene properties, such as scene category. By modeling the direct relationships of objects and geometry (which can be done in 3D, since perspective is already part of our framework) further improvement is possible.

As more types of objects can be identified and more aspects of the scene can be estimated, we hope that our framework will eventually grow into a vision system that would fulfill the ambitions of the early computer vision researchers – a system capable of complete image understanding.

**Acknowledgements.** We thank Bill Freeman for useful suggestions about the inference, Navneet Dalal for providing code and data, Moshe Mahler for his illustration in Figure 2, and Takeo Kanade for his car-road illustrative example. This research was funded in part by NSF CAREER award IIS-0546547.

## References

- [1] The PASCAL object recognition database collection. Website, 2005. <http://www.pascal-network.org/challenges/VOC/>.
- [2] H. Barrow and J. Tenenbaum. Recovering intrinsic scene characteristics from images. In *Computer Vision Systems*, 1978.
- [3] I. Biederman. On the semantics of a glance at a scene. In M. Kubovy and J. R. Pomerantz, editors, *Perceptual Organization*, chapter 8. Lawrence Erlbaum, 1981.
- [4] R. Brooks, R. Greiner, and T. Binford. Model-based three-dimensional interpretation of two-dimensional images. In *Proc. Int. Joint Conf. on Art. Intell.*, 1979.
- [5] M. Collins, R. Schapire, and Y. Singer. Logistic regression, adaboost and bregman distances. *Machine Learning*, 48(1-3), 2002.
- [6] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *Proc. CVPR*, 2005.
- [7] D. A. Forsyth, J. L. Mundy, A. Zisserman, and C. A. Rothwell. Using global consistency to recognise euclidean objects with an uncalibrated camera. In *Proc. CVPR*, 1994.
- [8] M. Greienhagen, V. Ramesh, D. Comaniciu, and H. Niemann. Statistical modeling and performance characterization of a real-time dual camera surveillance system. In *Proc. CVPR*, 2000.
- [9] A. Hanson and E. Riseman. VISIONS: A computer system for interpreting scenes. In *Computer Vision Systems*, 1978.
- [10] X. He, R. S. Zemel, and M. Á. Carreira-Perpiñán. Multiscale conditional random fields for image labeling. In *Proc. CVPR*, 2004.



Figure 7. We show car and pedestrian results from our baseline local detector and after inference using our model. The blue line shows the horizon estimate (always 0.5 initially). The boxes show detection estimates (green=true car, cyan=false car, red=true ped, yellow=false ped), with the solid lines being high confidence detections (0.5 FP/Image) and the dotted lines being lower confidence detections (2 FP/Image). In most cases, the horizon line is correctly recovered, and the object detection improves considerably. In particular, boxes that make no sense from a geometric standpoint (e.g. wrong scale (d), above horizon (b), in the middle of the ground (e)) usually are removed and objects not initially detected are found. Of course, improvement is not guaranteed. Pedestrians are often hallucinated (c,e) in places where they could be (but are not). In (f), a bad geometry estimate and repeated false detections along the building windows causes the horizon estimate to become worse and the car to be missed.

- [11] D. Hoiem, A. A. Efros, and M. Hebert. Geometric context from a single image. In *Proc. ICCV*, 2005.
- [12] S. G. Jeong, C. S. Kim, D. Y. Lee, S. K. Ha, D. H. Lee, M. H. Lee, and H. Hashimoto. Real-time lane detection for autonomous vehicle. In *ISIE*, 2001.
- [13] J. Kosecka and W. Zhang. Video compass. In *Proc. ECCV*. Springer-Verlag, 2002.
- [14] N. Krahnstoeber and P. R. S. Mendona. Bayesian autocalibration for surveillance. In *Proc. ICCV*, 2005.
- [15] S. Kumar and M. Hebert. Discriminative random fields: A discriminative framework for contextual interaction in classification. In *Proc. ICCV*, 2003.
- [16] S. Kumar and M. Hebert. A hierarchical field framework for unified context-based classification. In *Proc. ICCV*, 2005.
- [17] K. Murphy. The bayes net toolbox for matlab. In *Computing Science and Statistics*, volume 33. 2001.
- [18] K. Murphy, A. Torralba, and W. T. Freeman. Graphical model for recognizing scenes and objects. In *Proc. NIPS*. 2003.
- [19] Y. Ohta. *Knowledge-Based Interpretation Of Outdoor Natural Color Scenes*. Pitman, 1985.
- [20] J. Pearl. *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. Morgan Kaufmann Publishers Inc., 1988.
- [21] J. C. Platt. Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. In *Advances in Large Margin Classifiers*, 2000.
- [22] B. C. Russell, A. Torralba, K. P. Murphy, and W. T. Freeman. LabelMe: a database and web-based tool for image annotation. Technical report, MIT, 2005.
- [23] H. Schneiderman. Learning a restricted bayesian network for object detection. In *Proc. CVPR*, 2004.
- [24] E. Sudderth, A. Torralba, W. T. Freeman, and A. Wilsky. Learning hierarchical models of scenes, objects, and parts. In *ICCV*, 2005.
- [25] A. Torralba. *Contextual Influences on Saliency*, pages 586–593. Academic Press / Elsevier, 2005.
- [26] A. Torralba and A. Oliva. Depth estimation from image structure. *PAMI*, 24(9), 2002.
- [27] A. Torralba and P. Sinha. Statistical context priming for object detection. In *Proc. ICCV*, 2001.
- [28] Z. Tu, X. Chen, A. L. Yuille, and S. C. Zhu. Image parsing: Unifying segmentation, detection, and recognition. *IJCV*, 63(2):113–140, 2005.
- [29] P. Viola and M. J. Jones. Robust real-time face detection. *IJCV*, 57(2), 2004.