

Unconstrained 1D Range and 2D Image Based Human Detection and Segmentation

Mehmet Kocamaz¹ and Fatih Porikli²

Abstract—In this paper, an accurate and computationally very efficient multi-modal human detector and a novel technique to refine a detection window are presented. The human detector fuses 1D range scan and 2D image information via an effective geometric descriptor and a silhouette based visual representation within a radial basis function kernel support vector machine learning framework. Unlike the existing approaches, the proposed 1D+2D detector does not make any restrictive assumptions on the range scan positions, thus it is applicable to a wide range of real-life detection tasks. It works robustly under challenging imaging conditions and achieves several orders of magnitude performance improvement while reducing the computational load drastically.

The proposed multi-modal refinement technique takes the detection window as the input and produces the point-wise mask of the human which could be more beneficial for pose estimation and activity recognition tasks. It combines shape information encoded in a point-wise descriptor and color model of the human in graph cut framework. The shape descriptor is formed by utilizing the projected 1D range scan points in the image space. Our experiments demonstrate that the presented refinement method achieves significantly better performance than single-modal algorithms.

I. INTRODUCTION

Detecting humans by the intelligent driving systems is an important task for several reasons. First, it can prevent traffic accidents, and save lives of thousands of pedestrians [1]. Also, it provides understanding of their behaviors, recognition of their activities and forecasting of their actions. In order to serve for all of these purposes, it is more desirable to have a system which accurately detects humans and outputs their point-wise representations.

There are two sets of challenges that make this task complicated. The first one is the external factors. These factors are not object depended and often caused by environmental elements. Illumination variations, insufficient street lighting, saturation due to headlights, cast shadows, reflections, weather conditions, existence of human-like objects and clutter, and imaging noise fit into this category. External factors have absolute effects to the performance of the detection process.

The second set of challenges are due to the human itself, thus may be called as the internal factors. Humans have articulated body parts that move, rotate, and deform. They stand up, walk, run, bend and make body gestures. The appearance, height, weight, and clothing might differ significantly from one to another. Therefore, their bodies appear in different shapes and silhouettes. In addition, human body has various

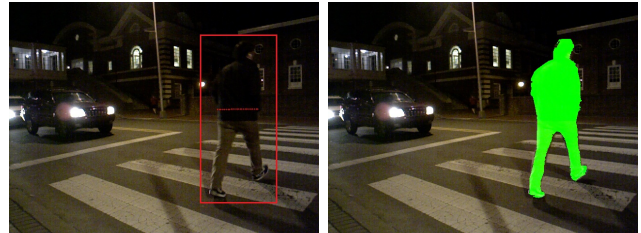


Fig. 1. Left image shows the detected pedestrian by the proposed multi-modal human classifier under severe illumination conditions. Right image displays the result of the proposed refinement technique if the detection window is provided as the input.

poses from distinct view points. All these factors make the objective of the human detection considerably more difficult than detection of rigid objects.

In this paper, we present two solutions for the human detection and segmentation problems. Firstly, a novel multi-modal human detector that fuses 1D range scans from a LIDAR (Laser Imaging Detection And Ranging) sensor and 2D monocular images from an optical camera is presented. The proposed algorithm integrates the photometric and depth features obtained from both data modalities in a joint classifier. It is robust under difficult environmental conditions. Unlike the existing approaches, it can detect humans even if the range scan beams hit upper torso and head of the body without making any assumptions about the visibility of the legs. This is critical for real world applications. For instance, the scan beam may easily miss the legs when the road climbs over a hill or there is a short subject, e.g. a child, in the detection range. The legs can be occluded due to skirts, bags, strollers, etc. When the pedestrian stands up sideways, only one leg is visible. Fusing multiple modalities not only increases the detection accuracy for such examples but also improves the computational time. Since it efficiently narrows down the search region in the image, our detector runs very fast.

In addition, a novel multi-modal *refinement* technique which takes a detection window, W , provided by the human detector and outputs the segmented human body points is proposed. This method uses both LIDAR and color image data in the *refinement* process. It combines the color and shape related cues in the graph cut framework to segment the human. A point-wise descriptor which employs the relative vectorial distance information to the projected LIDAR points on the image is constructed. These descriptors are used to train a Random Decision Forests (*RDF*) classifier.

¹ Carnegie Mellon University

² Australian National University

Based on the point-wise confidence scores obtained from *RDF* classifier, the human and background color models are constructed. A final graph cut step produces the mask of the human by fusing the color information and the *RDF* confidence scores. The proposed *refinement* method shows superior performance for the experimented dataset.

This paper makes several improvements to the human detection and segmentation problem in the following ways:

1: A highly accurate and computationally fast multi-modal human detector that fuses 1D range scans and 2D images is presented. This 1D+2D detector does not make any restrictive assumptions about the range scan positions.

2: A simple yet effective geometric descriptor is introduced for LIDAR data. A single-modal human detector, 1D+, using this descriptor is developed. This detector achieves higher accuracy than the state-of-the-art human classifiers based on 1D range scans.

3: An accurate multi-modal technique which uses LIDAR and color image data is proposed to refine a detection window. It presents a novel point-wise shape descriptor.

4: It is shown that the multi-modal human detector and the *refinement* classifier can be trained with less precise range information, for instance using Kinect sensor depth data, to eliminate the need for expensive and cumbersome manual labeling.

A review of the related work is summarized in the next section. Descriptors, fusion modules, training algorithm and the proposed human classifiers are described in Section 3. The details of the refinement method are explained in Section 4. Human dataset collection process, and the ground truth generation are described in Section 5. In the same section, experimental results of both single- and multi-modal human classifiers and the refinement technique are analyzed. Finally, the future directions of this work are drawn in the last section.

II. RELATED WORK

Computer vision and robotics communities have been conducting extensive research on human detection and segmentation for years. In both fields, selected sensor types has a direct impact on the fundamentals of the developed method.

A. Single Modal Human Detection

There are two essential sensor types used for single modal detection. First group includes the visual sensors, such as monocular cameras. Sensors that provide 3D geometric cues, such as one or multi layer LIDAR and Kinect, form the second group.

Visual human detectors take an input image, compute descriptors within all possible subwindows and ask a classifier to determine whether there is a human inside the subwindows or not. In earlier image based human detection works [2], [3], Haar wavelets are used to construct descriptors and train multiple linear Support Vector Machines (SVMs).

A seminal human detection technique that uses the Histogram of Oriented Gradients (HOG) features is proposed in [4]. For speed improvement, a rejection cascade of AdaBoost

classifiers using the HOG features is described in [5]. The region covariance features (COV) are first introduced in [6] and a classifier based on the underlying Riemannian manifold is deployed in [7]. These holistic methods achieve remarkable results, but they may suffer from occlusions.

Alternatively, human detection can be done by identifying body parts and their common shapes [8], [9], [10], [11], [12]. In these methods, local features for body parts are determined and combined to form human models. In [13] and [14] human silhouette information is also taken into account. These methods are more robust to occlusion. However, their performance highly depends on the image resolution of the human body parts.

Detectors that rely only on geometric cues often extract features from 3D or range scan data. For example, [15] applies a set of oriented filters to the spatial depth histograms. Instead of a classifier, a simple threshold operation is performed to find the humans. Depth images are converted to 3D point clouds in [16], [17]. A dictionary is constructed from the geodesic local interest points by [16]. This method has a high detection rate as long as humans are not occluded and touch other objects. A large feature vector that employs the histograms of the local depth information is used to represent humans [17]. This approach is robust to occlusions yet it is computationally very demanding and not suitable for real time applications. Only a single LIDAR range scan is processed to form a leg descriptor in [18] and [19]. These approaches extract a number of predefined features from the segmented line parts and train classifiers. They can detect humans if the legs are visible and the LIDAR beam hits at the lower torso level.

Integrated human detection and tracking solutions that use 3D data from Velodyne LIDAR are described in [20], [21]. These methods are more accurate than [18] and [19], yet the comparably expensive cost of the sensor limits their applicability.

B. Multi-Modal Human Detection

The underlying idea of using multiple modalities is to combine their complementary advantages.

Multi-modal detection algorithms can be centralized and cascaded. Centralized approaches combine the features obtained from different sensors in a single feature vector [22] and train a single classifier.

Cascaded approaches construct multiple descriptors and train separate classifiers for each modality. They compute classifier confidences [23], [24] or impose one of the classifier to reduce the search space of the other classifier [25]. The classifiers explained in [22], [23], [24] use 1D range scans and color images to construct features and extract features. [23] and [22] use the HOG and COV features, whereas [24] uses Haar-like features to form their visual descriptors. Similar to [18], they cannot handle the situations where the range scans hit human body other than the legs.

The method described in [26] utilizes 3D features obtained from image and 3D point cloud. However, it is computationally expensive to retrieve the 3D geometric features

for realtime applications. [25] shares a similar concept and focuses on reducing the computational load. This method uses 3D information retrieved from stereo images to limit the search. Since it has no geometric feature extraction or information fusion mechanism, it still suffers from abrupt changes in the illumination conditions.

C. Low-dimensional Human Shape Refinement

To best of our knowledge, a multi-modal low-dimensional human shape refinement method which utilize 1D LIDAR data and 2D color image does not exist. However, the algorithms that serve for this purpose can be categorized into two groups. First group of techniques [27] [28] [29] [30] [31] do not carry any features or information from the internal steps of the human detector to the refinement process. They can be independently run providing only the detection window as the input. *GrabCut* [27] can be considered as one of the most well-known methods in this category. It is designed as a semi-automatic and iterative segmentation algorithm which takes a box surrounding the object as the input.

Second group of the refinement methods leverage features or confidence scores taken from the internal steps of the methods which provide the low-dimensional human shape [32] [33] [34] [35] [36]. These methods do not follow a generic way to refine any given low-dimensional human shape. Hence, they depend on their rough shape estimators. The human silhouette feature retrieved from the HOG detector [4] builds the essential parts of the human model used in [34] [35] [36]. The faces of the humans are detected by Haar-like features [37] and help to initialize the seed points of *GrabCut* in [35] [36]. The method explained in [32] uses a set of features obtained from a human body part detector. A pre-processing step which utilizes Edgelet features defines the region of the interest in [33]. Then, the points of the human body are segmented in this region [33]. *Humanising GrabCut* [38] is a specialized version of *GrabCut* to refine the low-dimensional human shape. The predictions of the HOG detector are used to build the appearance models to initialize *GrabCut*.

III. 1D+2D DETECTOR

To take the advantages of the geometric and visual information, our 1D+2D multi-modal human detector combines the range scan and image descriptors into a single representation. It works in the joint higher-dimensional feature space. A diagram of the classifier is given in Fig. 2.

For a training image window W_i , the corresponding 1D range scan segment $L_i = (d_1, \dots, d_{m_i})$ within the window can be obtained either from the LIDAR or from the depth camera. In the case of the LIDAR sensor, there is a single, horizontal, synchronously acquired range scan segment within the window. On the other hand, the depth camera can provide multiple horizontal range scan segments, which are particularly valuable for training. Here, d is the depth, i.e. the distance of the sensor to a scene point in camera normal direction.

A. Geometric Descriptor

In contrast to [18] that assumes the geometric descriptor corresponds to leg region, our geometric descriptor f^{1D} applies to every part of the human body. It is obtained by the following procedure:

1) Depending on the size and depth of the human objects, range scans L_i for positive samples form arbitrary length vectors

$$f_i^{1D} = [d_1, \dots, d_{m_i}]_i^T, \quad 1 \leq m_i \leq \max(\|w_W\|) \quad (1)$$

where $\|w_W\|$ is the width of the window. In order to map the arbitrary length feature vectors onto a uniform, fixed dimensional feature space R^m , an m -point bilinear interpolation, B_m , is performed on f_i^{1D} . After the interpolation, the dimension of f_i^{1D} , that is m_i , becomes m

$$f_i^{1D} = [d_1, \dots, d_m]_i^T \leftarrow B_m(f_i^{1D}). \quad (2)$$

2) The distance between the sensor setup and a human differs significantly in the scene. To compensate for this distance, the closest point depth, d_C , in f_i^{1D} to the sensor setup is determined. Then, d_C is subtracted from f_i^{1D} .

$$d_C = \min(d_1, \dots, d_m), \quad d_C \neq 0 \quad (3)$$

$$f_i^{1D} \leftarrow f_i^{1D} - d_C = [d_1 - d_C, \dots, d_m - d_C]^T. \quad (4)$$

3) Human objects stand in front of all kinds of backgrounds. Background clutter, as well as other objects in the scene, may be positioned at different distances from the human objects. This causes considerable geometric feature variation around the silhouette of the human body. Capturing all this variation in the training data would be one approach. Yet, this requires a huge amount of training data, which would be impractical. Besides, it may cause the classifier to fail because of the weakened discriminative power of the descriptors.

Therefore, the depth values of the feature vector elements that are above a human shape threshold are upper bracketed. The threshold, d_H , is set to the maximum possible radius of a human. If a point in the feature vector f^{1D} has a depth value larger than the threshold, it is set to the maximum radius. As a result, the variation due to the other objects and background clutter are eliminated effectively:

$$d_k = \begin{cases} d_H & \text{if } d_k \geq d_H \\ d_k & \text{otherwise} \end{cases}. \quad (5)$$

B. Visual Descriptor

Due to its shape representation ability, computational simplicity, and robustness to illumination changes up to a certain degree, the HOG feature is used to form the visual part of our human descriptor, $f^{2D} = [v_1, \dots, v_n]^T$ in the classifier. The HOG features can represent efficiently the local appearance by a distribution of the edge gradients in a cell within an image region. These cells, either overlapping or on a regular grid, are smaller components of an image window. A histogram is obtained within a cell. These local cell histograms are concatenated into a larger window

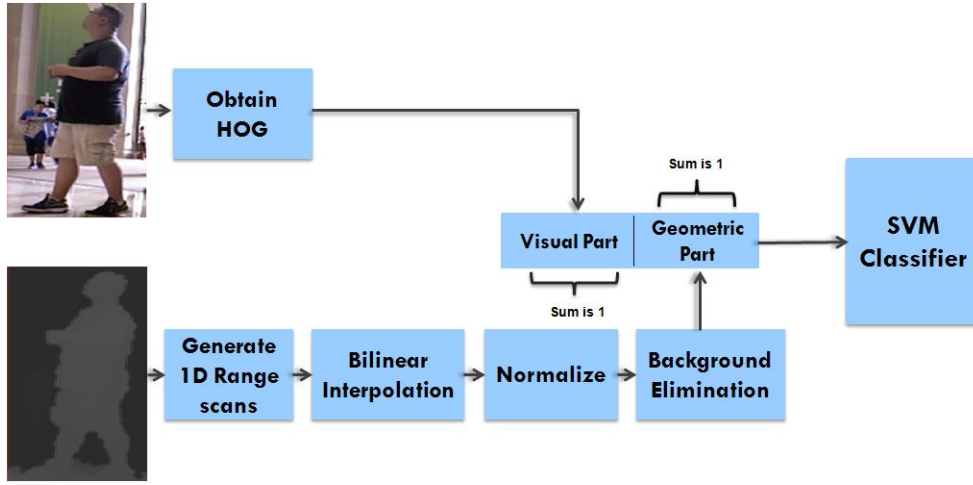


Fig. 2. Training process of the 1D+2D detector.

descriptor. All cell histograms of the window descriptor are normalized using the accumulated energy within the window for additional illumination robustness.

C. Combined Descriptor & Classifier Training

The geometric f^{1D} and visual f^{2D} features are concatenated in the same feature vector to form the final multi-modal human descriptor, f .

The raw geometric and visual feature vectors have different dimensions, thus their individual contributions in the combined multi-modal descriptor are not balanced. To overcome this issue, individual vectors are normalized to unit norm:

$$f^{1D} \leftarrow \frac{f^{1D}}{\sum_{k=1}^m d_k} \quad (6)$$

and

$$f^{2D} \leftarrow \frac{f^{2D}}{\sum_{k=1}^n v_k}. \quad (7)$$

The combined descriptor in R^{m+n} is then $f = [f^{1D} f^{2D}]^T$.

In training, the negative samples are chosen from the windows where there are no human objects. Since the window size changes according to the depth value of the window center, the size variation of the negative samples comes naturally. Even though in practice only LIDAR sensor data is available with the image, our training process still benefits from the additional depth camera data.

We use Support Vector Machines (SVMs) as our base classifiers. SVMs fits a hyperplane between the positive and negative training samples in the feature space. The decision boundary is defined by a set of support vectors that separate the positive and negative samples in a maximum margin. The decision function of SVM is

$$h(f) = \sum_{i=1}^m \alpha_i [\phi(f) \cdot \phi(f_i^*)] \quad (8)$$

where α_i are the weight of the corresponding m support vectors f_i^* and ϕ is a mapping function to a space \mathcal{H} . The

dot products in the decision function can be replaced by a kernel function:

$$k(f, f_i^*) = \phi(f) \cdot \phi(f_i^*) \quad (9)$$

By using a kernel function the classifier becomes a hyperplane in \mathcal{H} , yet it may be non-linear in the original input space. For given a set of labeled samples (x_i, y_i) where the labels $y_i = \{-1, 1\}$, the learning problem of SVM can be formulated as the minimization of

$$\min_{w, \varepsilon, b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \varepsilon_i \quad (10)$$

subject to

$$y_i(w \cdot f_i - b) \geq 1 - \varepsilon_i, \quad \varepsilon_i \geq 0 \quad (11)$$

where ε_i a penalty for the misclassified samples. The above optimization tries to classify as many training sample as possible correctly. Also, the minimization of $\|w\|$ makes the margin as large as possible. C is a variable term to set the relative influence.

We use the Radial Basis Function (RBF) as the kernel function of SVM:

$$\phi(f) \cdot \phi(f_i^*) = \exp(-\gamma \|f - f_i^*\|^2) \quad (12)$$

where γ is the width of Gaussian kernel width. By using RBF, it is always possible to find a decision function that perfectly represents a shape in a higher, possibly infinite, dimensional space. By incorporating RBF, SVM decision function takes the final form of

$$h(f) = \sum_{i=1}^m \alpha_i \exp(-\gamma \|f - f_i^*\|^2) \quad (13)$$

the result of final classification is the sign of $h(f)$. This decision function depends on the distance between the support vectors and the data, thus normalizing the geometric f^{1D} and visual f^{2D} feature vectors to unit norm, as formulated in Eqns. 6 and 7, is necessary. Otherwise, higher dimensional features would be favored by the SVM decision function.

In addition to the above 1D+2D detector, a single-modal classifier, called as 1D+ detector, is also trained by SVM using only the 1D range scans to assess the discriminative power of the proposed geometric descriptor.

D. Fast Detection

Since the speed of the human detection is an important factor, the 1D+2D detector is employed in a joint fashion that takes advantage of the depth information to eliminate the unnecessary window evaluations.

To determine whether a test window depicts a human, the corresponding 1D and 2D features are computed on the registered data. The range scan line L is aligned with the 2D image I by a perspective transformation $L_I : T(L)$ to obtain a set of image pixel coordinates $L_I = (p_1, \dots, p_n)$ in the image.

A search window $W(x, y, \delta x, \delta y)$ centered around p_k is slid on the coordinates of L_I . The size (width δx and height δy) of W is set according to the original depth value d_k of the point p_k such that for smaller depth values (objects closer to the sensor setup) the window size becomes larger. The window size is also proportional to the average human size at the corresponding original depth value.

There is no guarantee that the LIDAR beam always hits a specific level of the human body in a real application, thus the vertical position y of the image window W is not fixed. Instead, multiple windows at different vertical positions $y \pm \Delta y_j$ are tested for each p_k . Similar to the selection of the window size, the number of the vertical windows and their separation are determined by the original depth value of the center point. In this case, if d_k has a large value, a smaller vertical jumps Δy_j between multiple windows is desirable.

Within each window, the geometric descriptor f^{1D} and visual descriptors f^{2D} are computed, normalized, and concatenated into f . If the sign of the $h(f)$ in the SVM classifier is positive, a human is detected by the multi-modal classifier. Algorithm 1 outlines the testing procedure.

In contrast to the conventional visual-only human detectors that need to search entire image at different scales, our 1D+2D classifier reduces drastically the search space. It eliminates completely the image scaling step. Using L_I help to prune most of the image areas, which decreases the computational load greatly.

In practice, window evaluations can be ordered from nearest to far based on the LIDAR sensor depth values to determine the most critical object first.

IV. POINT-WISE REFINEMENT OF LOW DIMENSIONAL HUMAN SHAPE

In order to obtain the point-wise mask of human body points in the image, the low dimensional human shape, $W_H(x, y, \delta x, \delta y)$, which contains the detected human can be *refined*. The *refinement* process utilizes learned point-wise shape prior and the color of the points in the image. These cues are jointly combined in the graph cut framework to have the final human body points in $W_H(x, y, \delta x, \delta y)$.

Algorithm 1 Detection Algorithm

- Inputs:** $L = (d_1, \dots, d_n)$ range scan points, I, T, h
- 1: * Compute L_I , by $L_I : T(L)$
 - 2: **for** $k=1, \dots, n$ (all points in L_I)
 - 3: * Scale search window W by $1/d_k$
 - 4: * Compute geometric descriptor $f^{1D} = [d_1, \dots, d_m]^T$ inside W using Eqs. 1-5
 - 5: * Determine, Δy_j , vertical jump offsets from d_k
 - 6: **for each** Δy_j for W
 - 7: * Compute HOG $f^{2D} = [v_1, \dots, v_n]^T$
 - 8: * Normalize f^{1D} and f^{2D}
 - 9: * Concatenate f^{1D} and f^{2D} to $f = [f^{1D} f^{2D}]^T$
 - 10: * Compute $h(f) = \sum_{i=1}^m \alpha_i \exp(-\gamma \|f - f_i^*\|^2)$
 - 11: * if $h(f) > 0$ detect human, remove underlying points from L_I
-

A. Point-wise Shape Descriptor

A point-wise shape descriptor, f_s^{2D} , is formed for each point, p_w , in the detection window, W_H , in the following way:

1) All points, p_w , in the detection window, W_H , are normalized according to the top left point, x and y , of W_H .

$$p_w^x \leftarrow \frac{p_w^x - x}{\delta x} \quad (14)$$

$$p_w^y \leftarrow \frac{p_w^y - y}{\delta y} \quad (15)$$

2) The range scan line, L^W , which is in W_H is aligned with the 2D image I by a perspective transformation $L_I^W : T(L^W)$ to obtain a set of normalized image pixel coordinates $L_I^W = (p_1, \dots, p_n)$ in the image.

3) Two feature vectors, f_s^x and f_s^y , are formed from each coordinates, that are x and y , of L_I^W .

$$f_s^x = [p_1^x, \dots, p_i^x, \dots, p_n^x]^T, \quad p_i \in L_I^W \quad (16)$$

$$f_s^y = [p_1^y, \dots, p_i^y, \dots, p_n^y]^T, \quad p_i \in L_I^W \quad (17)$$

3) The dimensions of f_s^x and f_s^y are fixed by performing separate m -point bilinear interpolations, B_m . After the interpolation, the size of the dimensions become m .

$$f_s^x \leftarrow B_m(f_s^x) \quad (18)$$

$$f_s^y \leftarrow B_m(f_s^y) \quad (19)$$

4) The vectorial spatial distance of a point, p_w , to the range scan points in W_H is computed by subtracting each dimension of p_w from f_s^x and f_s^y .

$$f_s^x \leftarrow f_s^x - p_w^x = [p_1^x - p_w^x, \dots, p_m^x - p_w^x]^T \quad (20)$$

$$f_s^y \leftarrow f_s^y - p_w^y = [p_1^y - p_w^y, \dots, p_m^y - p_w^y]^T \quad (21)$$

5) We already know that the size of f_s^x , f_s^y and the geometric descriptor, f^{1D} , which is computed in the previous

section are m . Since the background points are set to d_H in the geometric descriptor, f^{1D} , the same indices of f_s^x and f_s^y are labeled as background by setting them to 1.

6) f_s^x and f_s^y are concatenated in f_s^{2D} to have the final point-wise shape descriptor, $f_s^{2D} = [f_s^x f_s^y]^T$.

B. Training the Refinement Classifier

The refinement of the detection window, W_H , can be considered as a 2-label classification problem of the points in W_H . The label of the first class is ‘‘human’’, while the other label is for the non-human points and is called ‘‘background’’. In order to train a point-wise refinement classifier, *Refinement-Classifier*, positive human descriptors, f_s^{2D+} , and negative human descriptors, f_s^{2D-} , are necessary. The samples of f_s^{2D-} are chosen from the non-human body points inside W_H .

Randomized Decision Forests (*RDF*) are a fast and effective machine learning technique [39] [40] [41] [42] which are suitable and applicable for wide range of different tasks and problems [43] [44] [45]. Therefore, it is used to train *Refinement-Classifier*. A Decision Forest consists of some number, T , of decision trees. A tree includes split and leaf nodes. Each split node consists of an axis, $f_s^{2D}(x)$, of f_s^{2D} , and a threshold τ . To classify the given descriptor of a point, f_s^{2D} , the split nodes of the decision tree are evaluated by starting from the root of the tree. Whenever a leaf node is hit in a tree, t , a decision distribution, $P_t(d|f_s^{2D})$, is obtained.

In our problem, $P_t(d|f_s^{2D})$ can be considered as a 2-bin histogram. The labels of this histogram are the human and background. The result label of the randomized decision forest classifier can be the average of all distributions given by the trees in the forest:

$$P(d|f_s^{2D}) = \frac{1}{T} \sum_{t=1}^T P_t(d|f_s^{2D}) \quad (22)$$

Or the result can be the label with the maximum number of votes by each decision tree, t , in the forest as formulated in the following equations:

$$L_t(P_t(d|f_s)) = \begin{cases} 1 & \text{if } P_t(d_H|f_s^{2D}) \geq P_t(d_B|f_s) \\ -1 & \text{otherwise} \end{cases} \quad (23)$$

where $L_t(x)$ is the decision label function of a given decision distribution of a tree, and t . d_H and d_B are the bin values of the human and background labels in the distribution. The normalized confidence score of a point which belongs to the human region becomes:

$$C_i = 0.5 + \frac{1}{T} \sum_{t=1}^T L_t(P_t(d|f_s^{2D})) \quad (24)$$

Then, the final decision label of the forest, $L(P(d|f_s^{2D}))$ becomes:

$$L(P(d|f_s^{2D})) = \begin{cases} 1 & \text{if } C_i \geq 0.5 \\ -1 & \text{otherwise} \end{cases} \quad (25)$$

Each tree is trained on a different set of randomly selected positive and negative samples using the following algorithm [44]:

1) Randomly obtain a set of splitting candidates for a tree node, $\Phi = (f_s^{2D}(x), \tau)$. $f_s^{2D}(x)$ is an axis of a point-wise descriptor, and τ is the split threshold.

2) The set of training points, $S = \{p_i\}$, are divided into two sets, S_l and S_r , for left and right leaves of the node by each Φ :

$$S_l(\Phi) = \{p_i \mid f_s^{2D}(x) \leq \tau\} \quad (26)$$

$$S_r(\Phi) = S - S_l(\Phi) \quad (27)$$

3) Find the best splitting candidate, Φ^* , which produces the largest information gain:

$$\Phi^* = \underset{\Phi}{\operatorname{argmax}} G(\Phi) \quad (28)$$

$$G(\Phi) = H(S) - \sum_{\psi \in (l, r)} \frac{|S_w(\Phi)|}{|S|} H(S_w(\Phi)) \quad (29)$$

where $H(S)$ is the Shannon Entropy. It is computed on the normalized distribution of the labels of the points in the set of S as in the following equation:

$$H(S) = - \sum_{i=1}^n Pr(l_i|P_L) \log_2 Pr(l_i|P_L) \quad (30)$$

where P_L is the label distribution in the set S , and l_i is the label name.

4) If the current depth of the tree is under a maximum threshold, create left and right children of the current node by using left and right subsets, $S_l(\Phi^*)$ and $S_r(\Phi^*)$.

C. Combining Shape Prior and Color in Graph Cut

Graph cut [46] [47] [48] [49] provides a powerful framework to produce globally optimal segmentation results. Its graph structure enables the combination of multiple different kinds of features in one joint framework. In our approach, graph cut is chosen as the infrastructure to incorporate the cues for a joint final solution.

It is difficult to generalize the color models of the human and background for all possible scenes. The point-wise descriptor, f_s^{2D} , described in the previous section does not include the color cue of the human body. However, the refinement procedure can utilize the color of the points. This can be achieved by employing these cues jointly in the graph cut framework.

Graph Cut Energy Functions: The energy function of the graph cut consists of two terms, namely the regional term, R , and the boundary term, B , as in the standard graph cut energy equation:

$$E(\tilde{L}) = \sum_{i \in V} R(l_i) + \sum_{\{i,j\} \in V} B_{i,j}(l_i, l_j) \quad (31)$$

where i and j are the nodes of any edge in the graph.

The regional term of the graph cut employs two cues. They are (a) the confidence score obtained from *Refinement – Classifier*, C_i , as defined in Equation 24 and (b) the color of the point. More precisely, the regional term of the graph cut energy function becomes:

$$\sum_{i \in V} R(l_i) = \sum_{i \in V} -\ln(p_s(l_i)) - \beta \ln(p_c(l_i)) \quad (32)$$

where β sets relative influence. The shape likelihood of being in the object region of a point, p_s , is formulated using the confidence score, C_i , produced by *Refinement – Classifier* as following:

$$p_s(l_i) = \begin{cases} C_i & \text{if } l = \text{"human"} \\ 1 - C_i & \text{if } l = \text{"background"} \end{cases} \quad (33)$$

In order to compute the color likelihood of a point, p_c , the color histogram models of background, M_B , and human, M_H , are built. The points whose confidence scores, C_i , produced by *Refinement – Classifier* are larger than a threshold, T_H , are used for M_H . To construct the background model, M_B , the points which are outside of the detection window, W_H , and whose confidence scores are less than a threshold, T_B , are considered. Then, p_c becomes as following:

$$p_c(l_i) = \begin{cases} Pr(c_i|M_H) & \text{if } l = \text{"human"} \\ Pr(c_i|M_B) & \text{if } l = \text{"background"} \end{cases} \quad (34)$$

where c_i is the color of the point p_i . The color discontinuity is defined in the boundary term of the energy function, $E(\tilde{L})$. As in [50], the color discontinuity between neighbor points in the graph is formed by the following equation:

$$B_{i,j \in V_L} = \lambda_1 \frac{1}{dist(i,j)} e^{-\frac{||c_i - c_j||^2}{2\sigma^2}} \quad (35)$$

where c_i and c_j are the colors of the points i and j , $dist(i,j)$ is the standard L_2 Euclidean norm yielding point distance, and σ^2 is the average squared norm in the image.

The final labeling, \tilde{L}_F , can be achieved by minimizing the energy function in Equation 31:

$$\tilde{L}_F = \underset{\tilde{L}}{\operatorname{argmin}} E(\tilde{L}) \quad (36)$$

The graph cut algorithm in [49] is used to minimize this equation. The proposed refinement process which utilizes the graph cuts and *Refinement – Classifier* is called *Refinement – Classifier_{GC}* after this point.

V. DATASET AND EXPERIMENTS

A. DontHitMe Dataset & Sensor Setup

In supervised learning, the quality and quantity of training data are very critical for the final performance of the classifier. More training data prevents from overfitting, improves generality, and enables trained models to capture possible variations of target class samples. Since our purpose is to construct an inclusive and unconstrained classifier that performs accurately without making any assumption about the range scan position on the human body, a large number

of training samples is required for training. However, it is cumbersome to collect such a large number of registered LIDAR and camera data where range scans hit humans on different parts of their bodies. To capture different pose, appearance variations and scan positions, the height and position of the LIDAR must be modified excessively. This is definitely a tedious and inefficient task with no guarantee of capturing sufficient amount and quality of data.

To our advantage, it is possible to generate a high number of diverse range scans for positive and negative samples by using a depth camera that provides the 3D structure of the scene. Any number of scans can be obtained from a depth image by converting the geometric information into LIDAR-like readings synthetically.

Towards this goal, a sensor setup composed of an Asus Xtion Pro Live IR and color camera, and a Hokuyo URG-04LX LIDAR was used. Three sensors, IR camera, color camera and LIDAR were registered in the same coordinate system. A multi-modal human data set, called as *DontHitMe*, was collected in outdoors (parking lots, streets, etc.) and indoors (campus, etc.) buildings. Since the IR camera is sensitive to the sunlight, outdoor data was recorded when there was no direct sunlight in the scene. In addition to the color and depth images, this dataset also includes registered 1D LIDAR range scans. It contains 40,000 images of 450 different humans in different poses, appearance variations, lighting conditions, and shadow artifacts. Several human shapes that present a challenge to existing human classifiers, such as women in skirts and small children were recorded. To capture the variance of the human poses, images are recorded sequentially at 8 fps. The location and height of the setup was changed during the collection process to collect samples in different backgrounds. Modifying the height of the sensor setup was diversified the recorded 1D range scans.

The original LIDAR range scans hit human body on different parts from the legs to the head. A total of 3,600 manual ground truth positions in images, depth camera data, and range scans were annotated. Each human in the dataset was labeled with a bounding box, $W(x, y, \delta x, \delta y)$. *DontHitMe* dataset is divided into two different categories. The first dataset, called as *DontHitMe-Indoor*, includes 30K frames and 3,000 ground truths which are recorded indoor campus buildings. The second dataset is collected in outdoors at night. It contains more challenging cases for human detectors, such as insufficient lighting and severe illumination changes because of car headlights. This dataset contains 10K frames and 600 ground truths, called as *DontHitMe-Night*.

To complement the original LIDAR data, the depth camera data in *DontHitMe-Indoor* were processed to obtain additional synthetic range scans as shown in Fig. 3. These horizontal scans were produced by uniformly sampling multiple positions vertically along the labeled human window W for the positive samples. In this way, multiple scans were generated from each part of human body, from the legs to the head. A depth scan $L_i = (d_1, \dots, d_{m_i})_i$ was discarded if it contained points where the depth camera does not provide a valid distance.



Fig. 3. 1D range scans are generated from the depth camera data for each positive window.

In addition to the low-dimensional ground truths, point-wise ground truths positions of 1016 humans were labeled in LIDAR, color and depth camera data. This dataset, called as *DontHitMe-Refine*, was used to train and test *Refinement-Classifier_{GC}*.

B. Analyzing Performance of 1D+2D Human Classifier

Several experiments were conducted to quantify the performance of the proposed multi-modal human classifier, 1D+2D, and its range scan only version, 1D+.

In the first experiment, we analyzed the performance of 1D+ detector. We obtained 46,000 positive and 376,000 negative samples from the LIDAR sensor scans and depth images of *DontHitMe-Indoor* dataset. A total of 43,000 positive samples are generated synthetically from the depth images by uniform sampling and additional 3,000 positive samples were obtained from the recorded 1D range scans.

In order to reduce the variability in the testing scores, we performed multiple rounds of 10-fold cross-validation. We aimed to see the performance of the 1D+ detector at the different parts of the human body. Therefore, the positive samples in *DontHitMe-Indoor* dataset are divided into 3 categories, as upper body, torso and lower body.

The outcomes of the proposed and the existing state-of-the-art classifiers for the separate human body parts and for negative samples can be seen in Table I. The results are compared to [18], which is a 1D range scan based human classifier. As visible, our 1D+ detector outperforms [18] at least by 20.2% for each part of the human body. The result of this experiment shows that assumptions on the visibility of the legs is not valid for real-life scenarios. The 1D+ is more robust and achieves remarkable accuracy at each level of the human body as can be seen in the detection performance curves of the classifiers in Fig. 4. As expected, the method explained in [18] shows its best performance if the range scans hit the lower part of the human body. Whereas the performance of our detector is almost same at different parts of the body. Proposed 1D+ does not miss any human at 89% false detection level. One of the main reasons of the consistent performance of our classifier at each part is that the positive samples are provided to our detector uniformly from different body parts in the training phase. Also, it learns more diverse geometric cues from every different part of the the body from head to the feet.

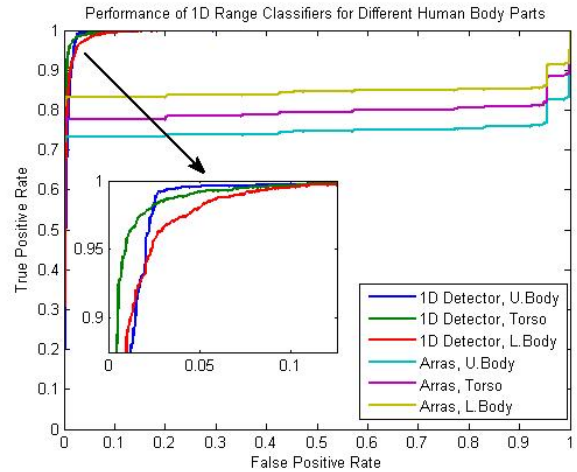


Fig. 4. ROC curves of the 1D+ Detector and Arras's classifier [18] at different parts of the human body.

TABLE I
COMPARISONS OF 1D RANGE SCAN BASED HUMAN DETECTORS FOR DIFFERENT HUMAN BODY PARTS

Test Set	1D+ Detector	Arras et al. [18]
Upper Body	97.5%	78.6%
Torso	97.9%	82.7%
Lower Body	96.8%	86.6%
Negative Samples	96.5%	5.6%

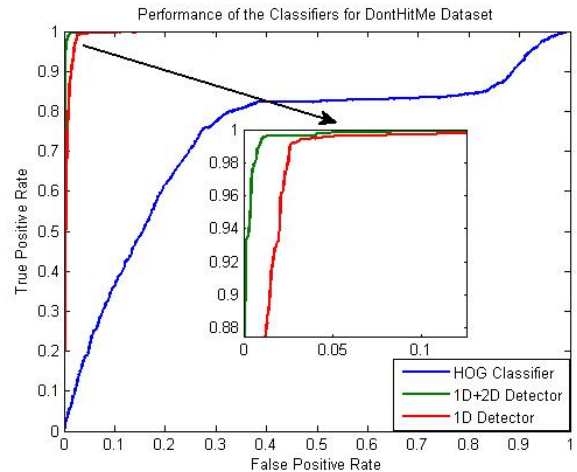


Fig. 5. Performance of the benchmark HOG [4] and the proposed 1D+2D and 1D+ human classifiers tested on *DontHitMe* dataset.

Another experiment was conducted to measure the performance of the proposed 1D+2D detector. A total of 1,000 positive and 10,000 negative visual descriptors were obtained from *DontHitMe-Indoor* dataset. For each visual descriptor, 20 different geometric descriptors were generated synthetically from different parts of the body by uniformly sampling in their corresponding depth images. In this way, total of 20,000 positive and 200,000 negative multi-modal samples which merge visual and geometric descriptors were gen-

erated from *DontHitMe-Indoor* dataset to train the 1D+2D detector. Also, for comparison purposes, 1D+ detector was trained only with the geometric descriptors and the HOG human classifier [4] was trained with the visual descriptors of this set. The accuracy of the proposed 1D+2D detector and 1D+ detector were compared to the HOG human classifier. As in the previous test, multiple 10-fold cross-validations were performed. During this experiment, it was ensured that the test fold and training folds include the samples obtained from different humans. In this way, testing of the geometric and visual descriptors obtained from the same positive samples used in training are prevented. The ROC curves of this experiment can be seen in Fig. 5. The 1D+2D detector and 1D+ detector perform significantly better than the visual only detector.

The proposed classifiers were tested with 600 labeled ground truth images of *DontHitMe-Night* dataset to quantify the performance of the classifiers under severe illumination conditions in outdoor. In this experiment, the classifiers trained in the previous experiment were applied on the night dataset. No new classifier was trained by using *DontHitMe-Night* and no synthetic range scans were generated from the depth images of this dataset. The tested geometric human descriptors were obtained only from the recorded LIDAR scans. The ROC curves of the 1D+2D, 1D+, and [4] detectors are displayed in Fig. 6. It can be seen that the HOG descriptor is not enough to represent the human under insufficient lighting and at night times. Our single-modal human descriptor achieved better accuracy than the HOG descriptor. Fusing the visual and geometric cues in a joint feature vector helped to improve the performance; 1D+2D detector outperforms consistently the other alternatives.

Since our geometric descriptor is obtained from LIDAR scans, our 1D+2D detector is more capable of handling image motion blur than the HOG classifier under low-light imaging conditions. Such motion blur examples (at the foot level of the pedestrians) can be seen in Fig. 5.

Note that, since it is accurate and computationally feasible at the same time, we compare against the HOG detector that uses SVM-RBF [4] for the most objective evaluations. There are other visual features that can generate higher detection results. Yet, such methods have prohibitively high computational loads for most practical applications.

C. Analyzing Performance of Proposed Refinement Method

The performance of *Refinement – Classifier_{GC}* was analyzed using *DontHitMe-Refine* dataset. As it was described in the previous section, a set of synthetic range scans were generated from the depth image data of *DontHitMe-Refine* by uniformly sampling multiple positions, $S_L = \{L_1, \dots, L_i, \dots, L_k\}$, in the window, W_H . Positive samples, f_s^{2D+} , are collected from the point-wise ground truth region of W_H . Negative samples, f_s^{2D-} , were taken from the non-human points in W_H . For each sampling position of, $L_i \in S_L$, a different *Refinement – Classifier_i* which includes 5 decision trees was trained. In the testing process, depending

TABLE II
MEDIAN, MEAN, MIN AND MAX OVERLAP SCORES OF PROPOSED *Refinement – Classifier_{GC}* AND *GrabCut* [27] FOR *DontHitMe – Refine* DATASET.

Overlap Scores	<i>GrabCut</i> [27]	<i>Refinement – Classifier_{GC}</i>
Median	0.61	0.94
Mean	0.59	0.935
Max	0.96	0.998
Min	0	0.83

on the location of Lidar scan in W_H , the corresponding *Refinement – Classifier_i* is chosen.

We performed multiple rounds of 10-fold cross-validation test to generate the refinement results of the proposed *Refinement – Classifier_{GC}* for *DontHitMe-Refine* dataset. Also, the results of *GrabCut* [27] were produced by providing W_H as the input to it. The following polygon area overlap formula is used to measure the overlap between the ground-truth and the result of the methods suggested by [51]:

$$O(\mathcal{R}_1, \mathcal{R}_2) = A(\mathcal{R}_1 \cap \mathcal{R}_2)^2 / (A(\mathcal{R}_1)A(\mathcal{R}_2)) \quad (37)$$

where \mathcal{R}_1 and \mathcal{R}_2 are the two regions to calculate the overlap between.

Overlap scores of the methods can be seen in Table II. The proposed *Refinement – Classifier_{GC}* significantly outperformed *GrabCut* by achieving median overlap score of 0.94, whereas the score of *GrabCut* is 0.61. Combining the shape and color cues in the proposed multi-modal refinement technique yields better results than *GrabCut* which utilize only color image data. The learned shape prior of the human body by *Refinement – Classifier* provides two main advantages that are (a) helps to build better human color models, M_H , and (b) carries a shape confidence score to the graph cuts. However, *GrabCut* uses all points in W_H to construct the object model at its first iteration. Hence, this superior performance of the proposed method can be expected. Some overlaid outputs of the methods are displayed in Figure 8.

D. Computational Load

A 64x128 detection window size was chosen for both the HOG and the proposed 1D+2D detector in the experiments. The dimension, m , of geometric feature f^{1D} is set to 40. The visual feature, f^{2D} , has the dimension of 3780. We used a machine which has 32GB RAM and Intel i7-2760QM quad processor to train and test the classifiers. The classifiers are implemented in native C++ language of Visual Studio 2010 Pro. The training phase of the 1D+2D detector consumed the largest memory among the classifiers in the second experiment since it requires 220,000 descriptors to fit into 29GB RAM, which took ~ 5 hours.

We compare the computational time and accuracies of the classifiers for *DontHitMe-Night* dataset experiment as can be seen in Table III. The average processing time of a 640×480 scale-space image (10,000 detection windows) by

TABLE III

AVERAGE RUNNING TIME AND FALSE ALARM RATE (AT 95% TRUE DETECTION RATE) OF DIFFERENT CLASSIFIERS FOR *DontHitMe-Night* DATASET.

Classifier	Time (in sec)	FAR at 0.95 TDR
HOG [4]	0.6	86%
1D+ Detector	0.0002	0.5%
1D+2D Detector	0.05	0%

TABLE IV

AVERAGE RUNNING TIME OF THE REFINEMENT METHODS AND THEIR MEDIAN OVERLAP SCORES FOR *DontHitMe-Refine* DATASET.

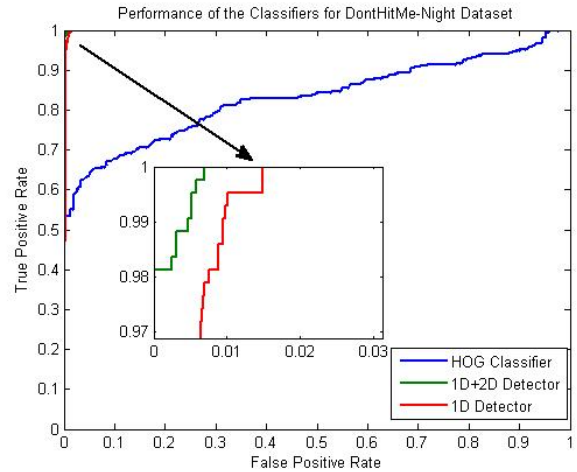
Method	Time (in sec)	Median O. Score
<i>GrabCut</i> [27]	0.74	0.61
<i>Refinement – Classifier_{GC}</i>	0.07	0.94

the benchmark HOG classifier is about 0.6 second. At 95% true detection rate, false alarm rate of it is 86%, whereas the false alarm rate of the 1D+2D detector is 0 on the tested dataset. Since the search space of the 1D+2D detector is reduced efficiently by the factors explained above, its average processing time is just 0.05 second. The proposed geometric descriptor has much less dimensions in comparison to other descriptors and it is easy to compute. Thus, 1D+ detector was able to run at 0.0002 second per scan in the same experiment.

In the training and testing steps of *Refinement – Classifier*, the size of the shape descriptor, f_s^{2D} , is set to 80. Each decision tree of the *RDF* was trained with 15,000 positive and 15,000 negative samples. The maximum depth of the tree was set to 20. In order to reduce the refinement time, the segmentation graph was built only for the pixels inside of the detection window, W_H . Table IV shows the average running time and median overlap scores of *GrabCut* and *Refinement – Classifier_{GC}*. Since the size of our segmentation graph and the shape descriptor, f_s^{2D} , are small and easy to build, our proposed method took only average of 0.07 seconds, whereas *GrabCut* took 0.74 seconds to refine a W_H .

VI. CONCLUSION

In this paper, we proposed a novel multi-modal human detector and a refinement technique to obtain the point-wise representation of the human, if a detection window is provided as the input. Our multi-modal human detector is accurate and computationally very fast. It combines 1D range scan and 2D image information within a SVM-RBF framework. Unlike the existing approaches, the proposed 1D+2D detector does not make any restrictive assumptions on the range scan positions, thus this unconstrained detector is applicable to a wide range of real-life detection tasks. We also discuss a range scan only version 1D+. Our extensive experiments demonstrate that the 1D+2D detector works robustly under challenging imaging conditions and achieves several orders of magnitude performance improvement (99%

Fig. 6. ROC curves of the classifiers for *DontHitMe-Night* dataset.

true detection at 0.005% false alarm rate in comparison to 54% true detection at 0.005% same false alarm rate on the benchmark) while reducing the computational load drastically (from 0.6 sec to 0.05 sec).

The presented multi-modal refinement technique trains a *RDF* classifier to provide shape confidence scores for the human body. These confidence scores are used to form color models of the background and human in the detection window for the graph cut. In order to segment the human region, graph cut is applied only inside of the detection window by fusing classifier confidence scores and the color information. Our method significantly outperforms *GrabCut* which is naturally suitable algorithm for this purpose.

As future work, the presented human detector can be extended to multi-class problems. Also, both of the proposed human detector and refinement method can be extended using multi-layer LIDAR data.

REFERENCES

- [1] National Center for Statistics and Analysis, “Traffic safety facts annual reports,” October 2012.
- [2] Constantine Papageorgiou and Tomaso Poggio, “A trainable system for object detection,” *International Journal of Computer Vision*, vol. 38, no. 1, pp. 15–33, June 2000.
- [3] Anuj Mohan, Constantine Papageorgiou, and Tomaso Poggio, “Example based object detection in images by components,” *IEEE Transactions Pattern Analysis and Machine Intelligence*, vol. 23, pp. 349–361, 2001.
- [4] Navneet Dalal and Bill Triggs, “Histograms of oriented gradients for human detection,” in *Proceedings of the 2005 IEEE Conference on Computer Vision and Pattern Recognition, (CVPR)*, 2005.
- [5] Qiang Zhu, Mei-Chen Yeh, Kwang-Ting Cheng, and Shai Avidan, “Fast human detection using a cascade of histograms of oriented gradients,” in *Proceedings of the 2006 IEEE Conference on Computer Vision and Pattern Recognition, (CVPR)*, 2006.
- [6] Oncel Tuzel, Fatih Porikli, and Peter Meer, “Region covariance: a fast descriptor for detection and classification,” in *Proceedings of the 9th European Conference on Computer Vision, (ECCV)*, 2006.
- [7] Oncel Tuzel, Fatih Porikli, and Peter Meer, “Human detection via classification on Riemannian manifolds,” in *Proceedings of the 2007 IEEE Conference on Computer Vision and Pattern Recognition, (CVPR)*, 2007.
- [8] Pedro F. Felzenszwalb and Daniel P. Huttenlocher, “Pictorial structures for object recognition,” *International Journal of Computer Vision*, vol. 61, no. 1, pp. 55–79, jan 2005.

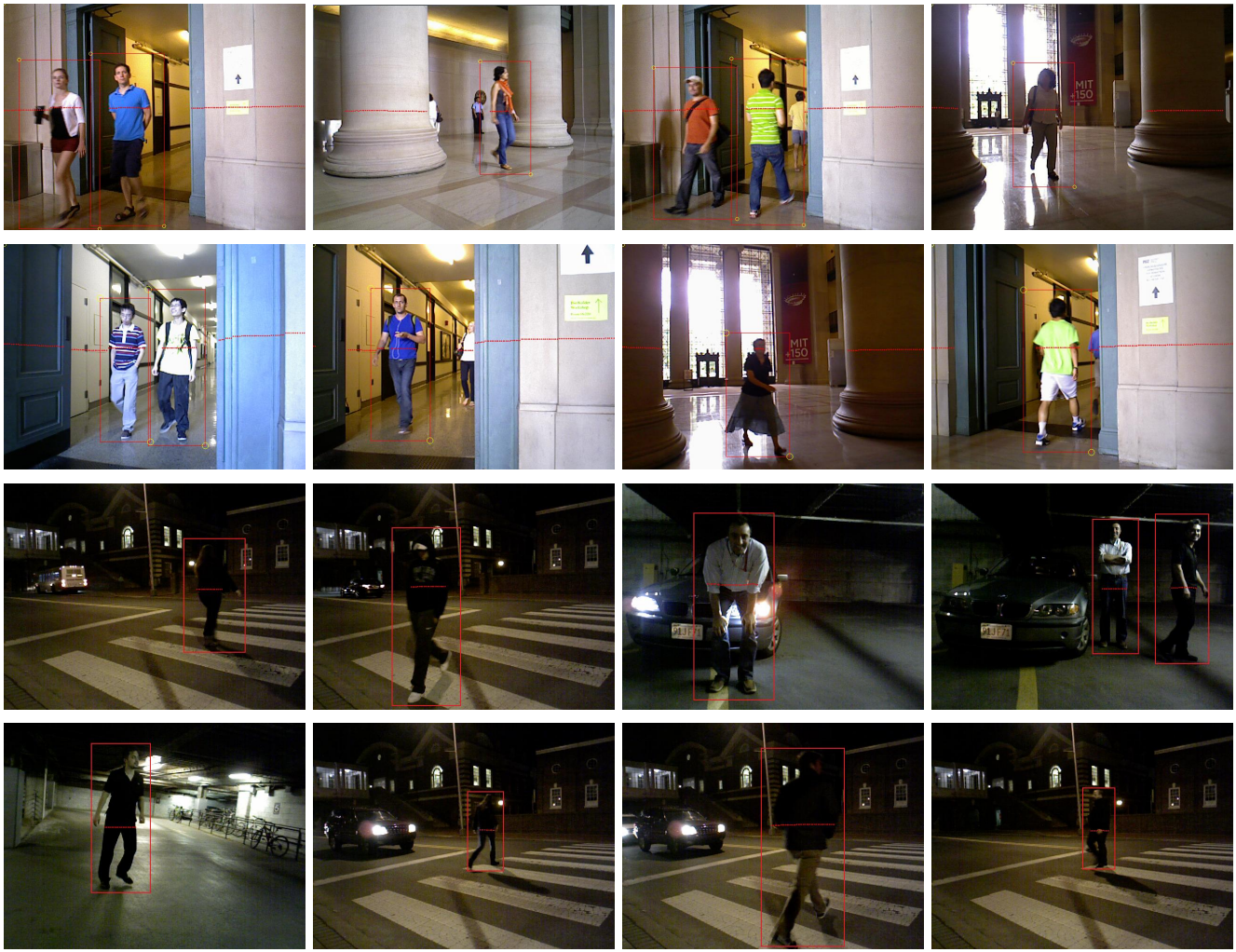


Fig. 7. Sample results of the 1D+2D multi-modal human detector. First two rows display detections in sample images from DontHitMe-Indoor dataset. Last two rows show sample results that were missed by the HOG based SVM-RBF [2] but accurately detected by the 1D+2D in DontHitMe-Night dataset.

- [9] S. Ioffe and D. A. Forsyth, "Probabilistic methods for finding people," *International Journal of Computer Vision*, vol. 43, no. 1, pp. 45–68, jun 2001.
- [10] Remi Ronfard, Cordelia Schmid, and Bill Triggs, "Learning to parse pictures of people," in *Proceedings of the 7th European Conference on Computer Vision, (ECCV)*, 2002.
- [11] Krystian Mikolajczyk, Bastian Leibe, and Bernt Schiele, "Multiple object class detection with a generative model," in *Proceedings of the 2006 IEEE Conference on Computer Vision and Pattern Recognition, (CVPR)*, 2006.
- [12] Krystian Mikolajczyk, Cordelia Schmid, and Andrew Zisserman, "Human detection based on a probabilistic assembly of robust part detectors," in *Proceedings of the 8th European Conference on Computer Vision, (ECCV)*, 2004.
- [13] Darius Gavrila and Vasanth Philomin, "Real-time object detection for smart vehicles," in *Proceedings of the Seventh IEEE International Conference on Computer Vision (ICCV)*, 1999.
- [14] Andreas Opelt, Axel Pinz, and Andrew Zisserman, "Incremental learning of object detectors using a visual shape alphabet," in *Proceedings of the 2006 IEEE Conference on Computer Vision and Pattern Recognition, (CVPR)*, 2006.
- [15] Taewan Kim, Sangho Cho, Jongmin Yoon, and Daijin Kim, "Pose robust human detection in depth image using four directional 2d elliptical filters," in *Proceedings of 11th IEEE International Symposium on Multimedia, (ISM)*, 2009.
- [16] Christian Plagemann, Varun Ganapathi, Daphne Koller, and Sebastian Thrun, "Real-time identification and localization of body parts from depth images," in *IEEE Int. Conf. on Rob. and Autom. (ICRA)*, 2010.
- [17] Sho Ikemura and Hironobu Fujiyoshi, "Real-time human detection using relational depth similarity features," in *Proceedings of the 10th Asian Conference on Computer vision, (ACCV)*, 2010.
- [18] Kai O. Arras, Oscar Martinez Mozos, and Wolfram Burgard, "Using boosted features for the detection of people in 2D range data," in *IEEE Int. Conf. on Rob. and Autom. (ICRA)*, 2007.
- [19] Christiano Pretebida, Oswaldo Ludwig, and Urbano Nunes, "Exploiting LIDAR-based features on pedestrian detection in urban scenarios," in *IEEE Conference on Intelligent Transportation Systems, (ITSC)*, 2009.
- [20] Luis Ernesto Navarro-Serment, Christoph Mertz, Nicolas Vandapel, and Martial Hebert, "LADAR-based pedestrian detection and tracking," in *Proc. 1st. Workshop on Human Detection from Mobile Robot Platforms, (ICRA)*, 2008.
- [21] Bharath Kalyan, K. W. Lee, W. Sardha Wijesoma, D. Moratuwage, and Nicholas M. Patrikalakis, "A random finite set based detection and tracking using 3d LIDAR in dynamic environments," in *Proceedings of the IEEE International Conference on Systems, Man and Cybernetics, (SMC)*, 2010.
- [22] Cristiano Pretebida, Oswaldo Ludwig, and Urbano Nunes, "LIDAR and vision-based pedestrian detection system," *Journal of Field Robotics*, vol. 26, no. 9, pp. 696–711, September 2009.
- [23] Luciano Spinello and Roland Siegwart, "Human detection using multimodal and multidimensional features," in *IEEE Int. Conf. on Rob. and Autom. (ICRA)*, 2008.

- [24] Z. Zivkovic and B. Krose, "Part based people detection using 2D range data and images," in *Proceedings of the IEEE International Conference on Intelligent Robots and Systems, (IROS)*, 2007.
- [25] Rodrigo Benenson, Markus Mathias, Radu Timofte, and Luc J. Van Gool, "Pedestrian detection at 100 frames per second," in *Proceedings of the 2012 IEEE Conference on Computer Vision and Pattern Recognition, (CVPR)*, 2012.
- [26] Stephen Gould, Paul Baumstarck, Morgan Quigley, Andrew Y. Ng, and Daphne Koller, "Integrating visual and range data for robotic object detection," in *ECCV workshop on Multi-camera and Multimodal Sensor Fusion Algorithms and Applications (M2SFA2)*, 2008.
- [27] C. Rother, V. Kolmogorov, and A. Blake, "Grabcut: Interactive foreground extraction using iterated graph cuts," *ACM Transactions on Graphics*, vol. 23, pp. 309–314, 2004.
- [28] Vinay Sharma and James W. Davis, "Integrating appearance and motion cues for simultaneous detection and segmentation of pedestrians," in *Proceedings of the International Conference on Computer Vision (ICCV)*, 2007, pp. 1–8.
- [29] Jian Zhao and Sen-ching Cheung, "Human segmentation by geometrically fusing visible-light and thermal imageries," *Multimedia Tools and Applications*, vol. 1, no. 1, pp. 1–29, 2012.
- [30] Jian Zhao and S Cheung Sen-ching, "Human segmentation by geometrically fusing visible-light and thermal imageries," in *International Conference on Computer Vision (ICCV) - Workshops*, 2009, pp. 1185 – 1192.
- [31] R. H. Luke, D. Anderson, J. M. Keller, and Skubic M., "Human segmentation from video in indoor environments using fused color and texture features," Tech. Rep., Electrical and Computer Engineering Department, University of Missouri, 2008.
- [32] Vibhav Vineet, Jonathan Warrell, Lubor Ladicky, and Philip Torr, "Human instance segmentation from video using detector-based conditional random fields," in *Proceedings of British Machine Vision Conference (BMVC)*, 2011, pp. 80.1–80.11.
- [33] Bo Wu and Ram Nevatia, "Simultaneous object detection and segmentation by boosting local shape feature based classifier," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2007.
- [34] Cyrille Migniot, Pascal Bertolino, and Jean-Marc Chassery, "Iterative human segmentation from detection windows using contour segment analysis," in *International Conference on Computer Vision Theory and Applications (VISAPP)*, 2013, pp. 405 – 412.
- [35] Antonio Hernandez-Vela, Miguel Reyes, Vctor Ponce, and Sergio Escalera, "Grabcut-based human segmentation in video sequences," *Sensors*, vol. 12, no. 11, pp. 15376–15393, 2012.
- [36] Antonio Hernandez-Vela, Miguel Reyes, Sergio Escalera, and Radeva Petia, "Spatio-temporal grabcut human segmentation for face and pose recovery," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR) - Workshops*, 2010, pp. 33 – 40.
- [37] Paul Viola and Michael Jones, "Robust real-time face detection," *International Journal of Computer Vision*, vol. 57, pp. 137–154, 2004.
- [38] Varun Gulshan, Victor Lempitsky, and Andrew Zisserman, "Humanising grabcut: Learning to segment humans using the kinect," in *Proceedings of the International Conference on Computer Vision (ICCV) - Workshops*, 2011.
- [39] J. R. Quinlan, "Induction of decision trees," *Machine Learning*, vol. 1, no. 1, pp. 81–106, 1986.
- [40] Leo Breiman, "Random forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, Oct 2001.
- [41] Amit Yali and Geman Donald, "Shape quantization and recognition with randomized trees," *Neural Computation*, vol. 9, pp. 1545–1588, 1997.
- [42] B. A. Shepherd, "An appraisal of a decision tree approach to image classification," in *Proceedings of International Joint Conference on Artificial Intelligence*, 1983, pp. 473–475.
- [43] Frank Moosmann, Bill Triggs, and Frederic Jurie, "Fast discriminative visual codebooks using randomized clustering forests," in *Proceeding of Advances in Neural Information Processing Systems*, 2007.
- [44] Vincent Lepetit, Pascal Lagger, and Pascal Fua, "Randomized trees for real-time keypoint recognition," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2005, pp. 775–781.
- [45] Jamie Shotton, Matthew Johnson, and Roberto Cipolla, "Semantic texton forests for image categorization and segmentation," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2008, pp. 1–8.
- [46] Yuri Boykov, Olga Veksler, and Ramin Zabih, "A new algorithm for energy minimization with discontinuities," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR) - Workshops*, 1999, pp. 26–29.
- [47] Y. Boykov, O. Veksler, and R. Zabih, "Fast approximate energy minimization via graph cuts," in *Proceedings of the International Conference on Computer Vision (ICCV)*, 1999.
- [48] Y. Boykov, O. Veksler, and R. Zabih, "Fast approximate energy minimization via graph cuts," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 23, no. 11, pp. 1222–1239, 2001.
- [49] Y. Boykov and V. Kolmogorov, "An experimental comparison of min-cut/max-flow algorithms for energy minimization in vision," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 26, no. 9, pp. 1124–1137, 2004.
- [50] Y. Boykov and M. Jolly, "Interactive graph cuts for optimal boundary and region segmentation of objects in N-D images," in *Proceedings of the International Conference on Computer Vision (ICCV)*, 2001.
- [51] S. Sclaroff and L. Liu, "Deformable shape detection and description via model-based region grouping," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 23, no. 5, pp. 475–489, 2001.

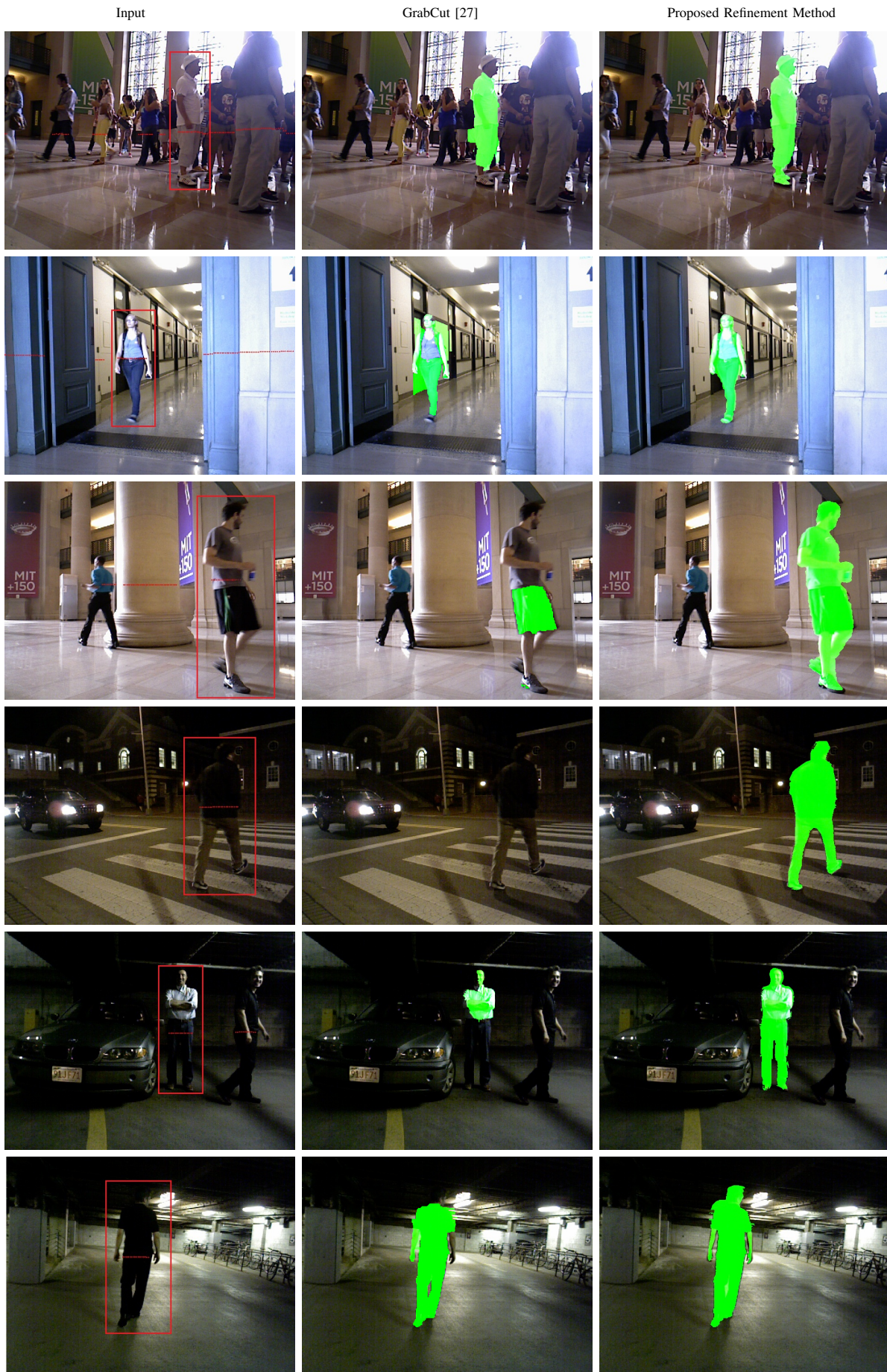


Fig. 8. Refinement results of *GrabCut* [27] and proposed *Refinement – Classifier_{GC}* for *DontHitMe-Refine* dataset.