

Fast and Accurate Pedestrian Detection in a Truck’s Blind Spot Camera

Kristof Van Beeck^{1,2(✉)} and Toon Goedemé^{1,2}

¹ EAVISE, KU Leuven - Campus De Nayer, J. De Nayerlaan 5, 2860
Sint-katelijne-waver, Belgium

² ESAT/PSI - VISICS, KU Leuven, Kasteelpark Arenberg 10, 3001 Leuven, Belgium
{kristof.vanbeeck,toon.goedeme}@kuleuven.be

Abstract. We propose a multi-pedestrian detection and tracking framework targeting a specific application: detecting vulnerable road users in a truck’s blind spot zone. Existing safety solutions are not able to handle this problem completely. Therefore we aim to develop an active safety system which warns the truck driver if pedestrians are present in the truck’s blind spot zone, using solely the vision input from the truck’s blind spot camera. This is not a trivial task, since—aside from the large distortion induced by such cameras—the application inherently requires real-time operation while at the same time attaining very high accuracy. To achieve this, we propose a fast and efficient pedestrian detection and tracking framework based on our novel *perspective warping window* approach. Experiments on real-life data show that our approach achieves excellent accuracy results at real-time performance, using a single core CPU implementation only.

Keywords: Computer vision · Pedestrian tracking · Real-time · Active safety systems

1 Introduction

Fast and meanwhile accurate pedestrian detection is necessary for many applications. Unfortunately these two demands are contradictory, and thus very difficult to unite. Even with today’s cheaply available computational power it remains very challenging to achieve both goals. Indeed, recent state-of-the-art pedestrian detectors achieving real-time performance heavily rely on the use of parallel computing devices (e.g. multicore CPUs or GPUs) to perform this task. This often makes it unfeasible to use these algorithms in real-life applications, especially if these applications rely on embedded systems to perform their tasks.

In this paper we propose an efficient multi-pedestrian detection and tracking framework for a specific application: detection of pedestrians in a truck’s blind spot zone. Statistics indicate that in the European Union alone, these blindspot accidents cause each year an estimated 1300 casualties [12]. Several commercial systems were developed to cope with this problem, both *active* and



Fig. 1. Example frame from our blind spot camera.

passive systems. Active safety systems automatically generate an alarm if pedestrians enter dangerous zones around the truck (e.g. ultrasonic distance sensors), whereas passive safety systems still rely on the focus of the truck driver (e.g. blind spot mirrors). However, none of these systems seem to adequately cope with this problem since each of these systems have their specific disadvantages. Active safety systems are unable to interpret the scene and are thus not able to distinguish static objects from actual pedestrians. Therefore they tend to generate many false alarms (e.g. with traffic signs). In practice the truck driver will find this annoying and often disables these type of systems. Existing passive safety systems are far from the perfect solution either. In fact, although blind spot mirrors are obliged by law in the European Union since 2003, the number of casualties did not decrease [18]. This is mainly due to the fact that these mirrors are not adjusted correctly; research indicates that truck drivers often use these mirrors to facilitate maneuvering. A passive blind-spot camera system with a monitor in the truck’s cabin is always adjusted correctly, however it still relies on the attentiveness of the driver.

To overcome these problems we aim to develop an active safety system based on the truck’s blind spot camera. Our final goal is to automatically detect vulnerable road users in the blind spot camera images, and warn the truck driver about their presence. Such an active safety system has multiple advantages over existing systems: it is independent of the truck driver, it is always adjusted correctly and it is easily implemented in existing passive blind spot camera systems. Due to the specific nature of this problem, this is a challenging task. Vulnerable road users are a very diverse class: besides pedestrians also bicyclists, mopeds, children and wheelchair users are included. Furthermore the specific position and type of the blind spot camera induces several constraints on the captured images. These wide-angle blind spot cameras introduce severe distortion while the sideways-looking view implies a highly dynamical background. See Fig. 1 for an example frame from our blind spot dataset.

However, the most challenging part is undoubtedly the hard real-time constraint, combined with the need for high accuracy. In this paper we present part of such a total safety solution: we propose an efficient multi-pedestrian



Fig. 2. Similarity vs perspective transformation model.

tracking- and detection framework based on blind spot camera images. Our algorithm achieves both high accuracy and high detection speeds. Using a single-core CPU implementation we reach an average of 13 FPS on our datasets.

In previous work [24,25] we proposed our initial *warping window approach*. However, this initial approach was based solely on a naive similarity warp, running up against its limit (e.g. w.r.t. accuracy for our application). In this paper we propose our *perspective warping window approach*: we extensively redesigned and improved our previous work making it more elegant and accurate, without significantly increasing the algorithmic complexity. Moreover, we even obtain higher computation speeds. Figure 2 concisely compares our previous and our improved novel approach presented here.

Our proposed algorithm briefly works as follows. Traditional state-of-the-art pedestrian detectors use a *sliding window* paradigm: each possible position and scale in the image is evaluated. This however is unfeasible in real-time applications. Instead, we proposed our *warping window* approach: we eliminate the need to perform a full scale-space search using the exploitation of scene constraints. That is, at each position in the input image we locally model the transformation induced by the specific camera viewpoint and the lens distortion. During detection, we can then warp the regions of interest (ROIs) in the image and use a standard pedestrian detector at a single scale on each ROI. This approach is integrated in a tracking-by-detection framework and combined with temporal information, making it more robust while reducing the detection time. We performed extensive experiments to evaluate our algorithm concerning both speed and accuracy. For this we recorded several realistically simulated dangerous blind spot scenarios with a real truck.

The remainder of this paper is organised as follows. In the next section we describe related work concerning this topic. Section 3 describes our algorithm in more detail, while in Sect. 4 we propose our experiments and evaluation results. We then conclude our work in Sect. 5.

2 Related Work

In the past few years the accuracy of pedestrian detectors has been significantly improved. Currently, even on challenging datasets excellent accuracy results are presented [9]. Initially, Dalal and Triggs [5] proposed a pedestrian detection

framework based on the Histograms of Oriented Gradients (HOG) combined with an SVM (Support Vector Machine) for classification. This idea was further refined in Felzenszwalb et al. [14] where the authors extended the concept with a part-based HOG model rather than a single rigid template. Evidently, this increases calculation time. To partially cope with this problem they proposed a more efficient cascaded framework [13]. Apart from increasing the model complexity, one can opt to increase the number of features to improve detection accuracy. Indeed, such a detector is presented in [7], called *Integral Channel Features*. However, each of these detectors still uses a sliding window approach. Across the entire image the features are calculated at all scales. To avoid such an exhaustive full scale-space search several optimisation techniques were proposed; e.g. Lampert et al. [17] proposed an efficient subwindow search. Dollár et al. [6] introduced the *Fastest Pedestrian Detector in the West (FPDW)* approach, in which they approximate feature responses from scales nearby thus eliminating the need to fully construct the scale-space pyramid. Extensive comparative works have been published [8, 10] to determine the most accurate approach. Both conclude that the HOG-based approach outperforms existing methods.

More recently, a benchmark between sixteen state-of-the-art pedestrian detectors was presented [9]. The authors conclude that part-based HOG detectors still achieve the highest accuracy, while the FPDW is one order of magnitude faster with only small loss in accuracy. Based on these conclusions, we performed extensive benchmark experiments with both pedestrian detectors to determine the most optimal one for our framework. These results, and more in-depth information on how both pedestrian detectors work are given in Sect. 3.1.

Concerning speed, several GPU optimisations were proposed. Prisacariu and Reid [22] proposed a fast GPU implementation of the standard HOG model. In [21], Pedersoli et al. presented a pedestrian detection system using a GPU implementation of the part-based HOG model. Benenson et al. [3] proposed work in which they perform model rescaling instead of image rescaling, and combined with their stixel world approximation [2] they achieve fast pedestrian detection. Recently the authors proposed their *Roerei* detector [1]. Based on a single rigid model they achieve excellent accuracy results. However, in real-life applications using embedded systems such high-end GPU computing devices are often not available. Therefore our algorithm focuses on real-time performance, while maintaining high accuracy, on standard hardware.

Speed optimisation is also achieved using pedestrian tracking algorithms, of which several are proposed in the literature. They often rely on a fixed camera, and use a form of background modelling to achieve tracking [23, 26]. Since in our application we have to work with moving camera images, this cannot be used. Pedestrian tracking algorithms based on moving cameras mostly use a forward-looking view [11] or employ disparity information [15]. Cho et al. [4] proposed a pedestrian tracking framework related to our work, exploiting scene constraints to achieve real-time detection. However, they use a basic ground-plane assumption whereas our approach is much more flexible and generic. Moreover, our specific datasets are much more challenging due to the severe distortion.

We significantly differ from all of the previously mentioned approaches. We aim to develop a monocular multi-pedestrian tracking framework with a challenging backwards/sideways looking view, targeting high accuracy at real-time performance. Furthermore, most of these classic sliding window approaches assume only object scale variation. Other geometrical variations (e.g. rotation [16] and aspect ratio [19]) are usually covered by an exhaustive search approach. Our proposed warping approach offers a solution that can even cope with perspective distortion. In fact, without our warping window paradigm it would be unfeasible in practice to perform such an exhaustive search in a perspective distortion space.

3 Algorithm Overview

As mentioned above, existing pedestrian detectors employ a sliding window approach. Across all positions and scales in the image the features are calculated and evaluated, making it almost impossible to meet the stringent real-time demands needed in most safety applications. To achieve real-time detection speeds with high accuracy we propose our novel *perspective warping window* approach.

Our idea is mainly based on the following observation. Looking at an example frame from our dataset (see Fig. 1) one clearly notices that the pedestrians appear rotated, scaled and perspectively transformed. This is due to the specific position and the wide-angle lens of our blind spot camera. The crux of the matter is that this transformation only depends on the position in the image. Thus each pixel coordinate $\mathbf{x} = [x, y]$ uniquely defines the transformation at that specific position. If at each pixel position this transformation is known, we can dramatically speed-up pedestrian detection. Based on this transformation we can locally warp each region of interest to upright pedestrians at a fixed height, and run a single-scale pedestrian detector on each warped ROI image patch. This approach effectively eliminates the need to construct a scale-rotation-transformation-space pyramid, and thus is very fast. Moreover, this approach is easily generalisable to other applications where such distortion occurs due to non-standard camera viewpoints and/or wide-angle lens distortions (e.g. surveillance cameras). To determine this transformation at each pixel coordinate a one-time calibration step is needed. To further increase both accuracy and speed, we integrate this warping window approach into an efficient tracking-by-detection framework. We use temporal information to predict future positions of pedestrians, thus further reducing the search space. Below we describe these steps in more detail. In Subsect. 3.1 we describe how our new perspective warping approach models the transformation, and motivate important algorithmic design choices such as the pedestrian detector, and the optimal scale parameter. In Subsect. 3.2 we then show how we integrate each of these steps into our total framework, thus describing our complete algorithm.

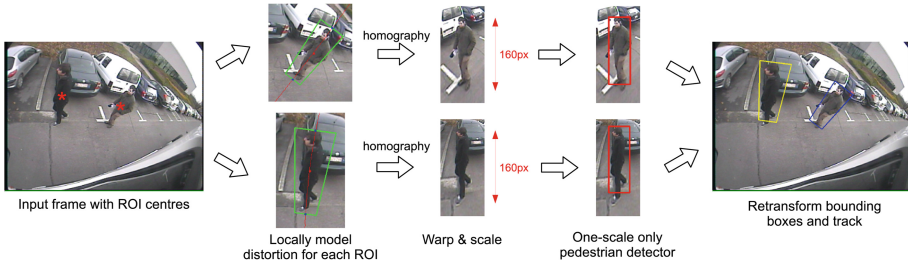


Fig. 3. Illustration of our novel perspective warping window approach. At each position in the image we locally model the distortion, warp the ROIs to a standard scale and use a one-scale only pedestrian detector.

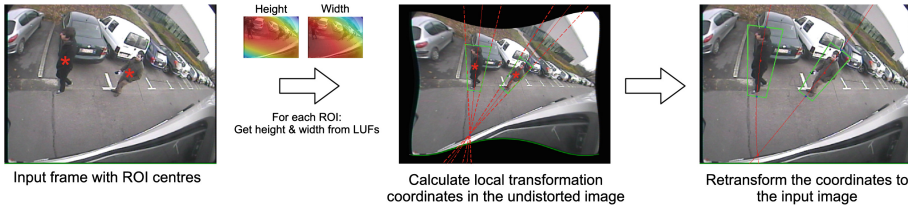


Fig. 4. The transformation is modeled as a perspective transformation, calculated in the undistorted image.

3.1 Warp Approach

Figure 3 illustrates our perspective warping window approach. Starting from input images as given in Fig. 1, pedestrians appear rotated, scaled and perspective distorted. If we assume a flat groundplane, these transformation parameters only depend on the specific position in the image. If we know the transformation we can model the perspective distortion for that ROI, extract and warp the ROI image patch to a fixed-scale (160 pixels - motivated further in this work) and perform pedestrian detection on a single scale only. We thus eliminate the need to construct a scale-space pyramid. Note that, although we perform detection on a single scale only, the pedestrian model still provides some invariance with respect to the pedestrian height. However, if large deviations from the standard height (e.g. children) need to be detected, an extra scale needs to be evaluated. After detection, the coordinates of the detected bounding boxes are retransformed and fed into our tracking framework. Next we describe further details of our algorithm: how this position-specific transformation is mathematically modeled and how the calibration is performed. We further motivate the choice of our baseline pedestrian detector and determine the optimal fixed-height parameter.

Transformation Modelling. Figure 4 illustrates how the transformation is locally modeled. We use a perspective distortion model in the lens-distortion-corrected image. At each position, the height and width (at the ground) are

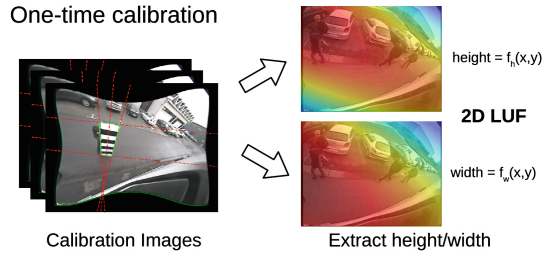


Fig. 5. A one-time calibration is needed to determine the local perspective distortion.

known after a one-time calibration step (see further). These are visualised as two heat maps (the so-called look-up-functions or LUFs) in Fig. 4. The transformation coordinates are determined as follows. Each ROI centre coordinate (indicated with the red asterisk in the leftmost image) is first transformed into the undistorted image. This lens undistortion is simply based on the traditionally used radial lens distortion model:

$$\mathbf{x}' = \mathbf{x}(1 + k_1r^2 + k_2r^4) \tag{1}$$

$$r^2 = x^2 + y^2 \tag{2}$$

Here, \mathbf{x}' denotes the corrected pixel coordinate, \mathbf{x} the input coordinate and k_1 and k_2 indicate the radial distortion coefficients.

Next we calculate the vantage line through this ROI centre in the undistorted image, and determine the height and width (at the bottom) from the two LUFs. Based on these data we construct the perspective model in the undistorted image. The rotation of the image patch is determined from the angle of the vantage line, and the length ratio between the top and bottom is calculated based on the distance to the vantage point (visualised in the middle of Fig. 4). We thus locally model the pedestrians as if they are planar objects standing upright, faced towards the camera (that is, perpendicular to the optical axis of our blind spot camera). Our experiments show that this is a valid approximation for pedestrians. These coordinates are then retransformed to the distorted input image. Note that evidently only the coordinates are transformed, the middle image displayed here is only used for visualisation purposes. Based on the coordinates in the distorted image, and the known calibration data we apply a homography on the ROI image patch, thereby effectively undoing the local perspective distortion (visualised in Fig. 3).

Calibration. To obtain these two LUFs, a one-time calibration step is needed. To achieve this, we manually annotated about 200 calibration images. We utilised a planar calibration board of 0.5×1.80 m, and captured calibration positions homogeneously spread over the entire image (Fig. 5). The labeling was performed in the undistorted image. These images yield the vantage point, and the height and width of a pedestrian (at the ground) at each position for that image.

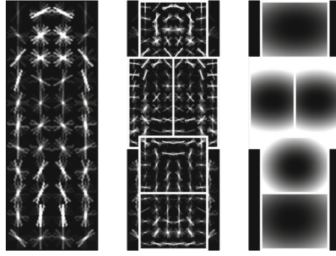


Fig. 6. The part-based pedestrian model [14]. (L) root filter (M) part filters (R) score distribution over parts.

Next we interpolated these datapoints using two-dimensional second order polynomial functions for both the height and the width: $f_h(x, y)$ and $f_w(x, y)$ with:

$$f_i(x, y) = p_0 + p_1x + p_2y + p_3x^2 + p_4xy + p_5y^2 \quad (3)$$

Both functions are displayed as heat maps in Fig. 5: for each pixel coordinate they effectively give the height and width of the calibration pattern at that location. If for some reason the position of the camera w.r.t. the ground plane changes, a recalibration needs to be performed. This is highly unlikely though, due to the robust camera mounting on the truck. Thus to summarise, detection is composed of four different steps: calculate the local perspective distortion model at each ROI centre, perform a homography and transform the pedestrians to an undistorted, upright position at a fixed height of 160 pixels, run a pedestrian detector at one scale, and finally retransform the coordinates of the detected bounding boxes to the original input image.

Pedestrian Detector. Based on the comparative works given in Sect. 2 we conclude that, since we aim for high accuracy, two approaches towards pedestrian detection are most suited for our application: the deformable part-based HOG models (DPM), and the rigid model-based FPWD. The FPWD has only slightly lower accuracy on established datasets [9] and is much faster. However, since we need to evaluate only one scale, no feature pyramid is constructed, thus this speed advantage is here not relevant. Selecting the most appropriate pedestrian detector for integration in our framework thus boils down to the selection of the most accurate detector on our dataset. For this we performed accuracy measurements for both detectors. To perform a fair comparison, let us briefly discuss how both pedestrian detectors work.

The DPM detector uses a pretrained model, consisting of HOG features (see Fig. 6). It consists of a root filter and a number of part filters representing the head and limbs of the pedestrian. To use this model, first a scale-space pyramid is constructed, using repeated smoothing and subsampling. For each pyramid layer the HOG features are computed. Then, for a specific scale the response of the root filter and the feature map is combined with the response of the part filters to calculate a final detection score. On our 640×480 resolution images this

detector off-the-shelf needs an average of 2.3 s per frame, while their cascaded version needs on average 0.67 s per frame.

As opposed to the deformable part-based detector, the FPDW detector utilises a rigid model, not making use of deformable parts. Again, a scale-space pyramid is constructed. Next, features are calculated on transformed versions of the original input image, called *channels*. Examples are color channels and gradient channels. The most basic features are a sum over a rectangle region in one of the channels. These features are combined into higher-order features, and classification is performed using a depth two decision tree boosted classifier (AdaBoost). Fast rejection of candidate windows is possible through the use of a *soft-cascade* approach. Essentially, this detector uses ICF as a baseline detector, and achieves a speed-up through a more efficient construction of this feature pyramid; intermediate feature scales are approximated using scales nearby, avoiding the need to compute all feature scales. Out-of-the-box calculation time for ICF on average equals 451 ms per frame (on 640×480 images), while using the FPDW approach the calculation time drops to 147 ms per frame.

We altered both detectors into single-scale detectors to utilise in our framework. We need to determine the most optimal scale for each detector; the rescale height at which the maximal accuracy is reached.

Determining the Optimal Scale Factor. As mentioned, we rescale the pedestrians to a fixed height in order to reduce the calculation time. For this, an optimal value needs to be determined. To achieve this, we labeled and extracted 6000 pedestrian images at different locations in the images from our dataset, and performed the warp operation as given above. These pedestrians were warped to fixed heights, and we then performed accuracy measurements with both single-scale pedestrian detectors to determine the optimal height for each of them. Besides our *perspective transformation model* presented in this paper, we also warped the pedestrians using the *similarity transformation model* as explained in [25], simply consisting of a rotation and scaling operation (see Fig. 2 for a qualitative comparison). This was done to analyse the benefit of our more complex

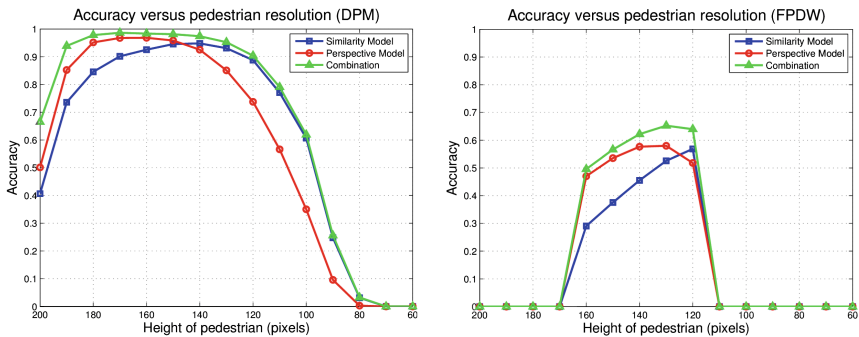


Fig. 7. Determining the optimal scale parameter. **Left:** results for the deformable-part detector. **Right:** results for the FPDW detector.

perspective model. Figure 7 displays our results for both pedestrian detectors (left: DPM, right: FPDW). Besides the individual transformations, we also give the combined accuracy. Evidently, optimal accuracy is reached when pedestrians in the rescaled image patches approximate the height of the detection models. Note that the one-scale DPM detector achieves much better accuracy results for all transformation models as compared to the FPDW detector. The reason for this significant difference is found in the design methodology of both detectors. Due to the part-based approach, the DPM detection model is much more flexible, making this detector invariant to slight deviations in height between the pedestrians that need to be detected, and the actual pedestrian model. The FPDW is much more sensitive for this due to the rigidity of their detection model. Since in our image patches slight differences between the actual and estimated pedestrian height exist (due to small calibration errors and the inherent height differences between pedestrians), more search scales would be needed to obtain higher accuracy with the FPDW approach. Further note that for both detectors the optimal resolution of the perspective and similarity transformation model differs. Concerning DPM, for the first transformation model the optimal height lies at 160 pixels, whereas the latter reaches its optimum at 140 pixels. As can be seen, the perspective model has a clear accuracy advantage over the similarity model. If both models were combined, an even higher accuracy is achieved. This, however, would double the calculation time. Although at lower accuracy, a similar trend is noticeable for the FPDW detector.

These insights favor the use of the deformable part-based model over the rigid FPDW approach. However, these experiments exclude the influence of the position in the image w.r.t. the detection accuracy. Specifically for our images, the position inherently defines the scale and the amount of transformation. For example, specific positions in the image require significant upscaling to transform these patches to the fixed height. Furthermore, research indicates that rigid models perform better on low-resolution pedestrian patches, whereas part-based models achieve higher accuracy on large pedestrian patches [20]. Thus, a combined approach—where the *best* detector is selected using the spatial location in the image—may increase accuracy.

We performed such experiments to evaluate this hypothesis as follows. The *best* detector is defined as having the highest detection score on the same image patch. Since the range of detection scores for both pedestrian detectors differs, first a normalisation of the detection scores is applied. For this, we warped the 6000 pedestrian image patches mentioned above using the perspective transformation to the optimal height of 160 pixels. These patches are then equally divided in a *training* and a *test set*. Next, both pedestrian detectors were evaluated on the test set, and their detection scores were normalised (subtract the average and divide by standard deviation - determined on the trainingset). This allows for a fair comparison between both detectors. For each patch, the detector with the highest detection score (if both found a detection) is assigned as optimal detector for this patch. These results are visualised in Fig. 8 in function of the position in the image. The colored dots indicate where which pedestrian



Fig. 8. Performance of both pedestrian detectors in function of the position in the image. Colored dots indicate which pedestrian detector performed best. Yellow: DPM. Magenta: FPDW (Color figure online).



Fig. 9. Performance of the two transformation models in function of the position. Colored dots indicate which model performed the detection. Red: perspective model. Blue: similarity model. Green: both models. Yellow indicates missed detections (Color figure online).

detector performed best (yellow: DPM, magenta: FPDW). As visualised, no specific image location is favored by any detector. Therefore, currently we do not perform such a combination, and utilise a single pedestrian detector.

Based on all experimental results mentioned above, we thus opted for the cascaded deformable part-based models as baseline detector in our framework: it achieves excellent accuracy results, lends itself perfect to perform *true* single scale detections and, due to this single scale approach, achieves excellent speed results (as shown further). Figure 9 shows (only for the deformable part-based model) where each transformation model performs best in function of the position in the image. Red dots indicate where the perspective model worked, blue where the similarity model worked and green were both models found the detection. Yellow indicates a missed detection. The perspective model obviously performs much better than the similarity model. The similarity model performs slightly better only at the image border, due to the small calibration error there.

The perspective model performs better close to the truck because of the large amount of viewpoint distortion there. Note that if we analyse positions where both models found the detection, the perspective model achieves the best detection score in 69% of these cases, further indicating its clear advantage over the similarity transformation model.

3.2 Tracking Framework

To further improve the accuracy and detection speed we integrated our warping window approach in a tracking-by-detection framework. This is implemented as follows. Instead of a full frame search, we use initialisation coordinates (which define transformation ROIs) at the border of the image, and initially only perform detection there. See Fig. 10 for an example. If a pedestrian is detected, a linear Kalman filter is instantiated for this detection. As a motion model we use a constant velocity assumption. Our experiments indicate that this assumption holds for our application. The state vector x_k consists of the centre of mass of each detection and the velocity: $x_k = [x \ y \ v_x \ v_y]^T$. Based on the update equation $\hat{x}_k^- = A\hat{x}_{k-1}$ we estimate the next position of the pedestrian. Here, \hat{x}_k^- indicates the *a priori* state estimate and \hat{x}_k indicates the *a posteriori* state estimate, at timestep k .

The process matrix A thus equals:

$$A = \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad (4)$$

Based on this motion model we predict the position (that is, the centre of mass) of the pedestrian in the next frame. Each estimated new pixel coordinate is then used as input for our warping window approach: we calculate the transformation model, warp the ROI and perform pedestrian detection on this ROI. For each



Fig. 10. Example of five initialisation coordinates together with their corresponding transformation ROIs.

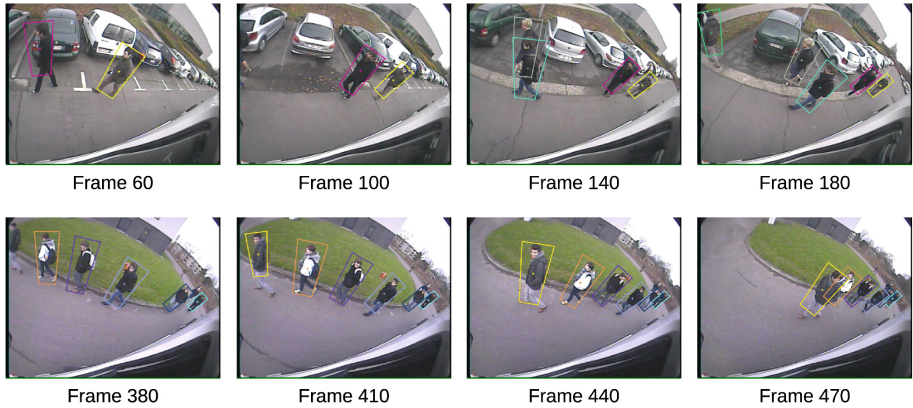


Fig. 11. Qualitative tracking sequences over two of our datasets (top and bottom row) - see <http://youtu.be/gbnysSoSR1Q> for a video.

pedestrian that is being tracked, our algorithm verifies if a new detection is found. This is evaluated by constructing a circular region around the estimated coordinate based on the scale of that tracked instance. If a new detection is found in this region, the Kalman filter is updated and the new position is predicted. If multiple detections are found, we associate the closest based on the Euclidean distance. The bounding box coordinates of tracked instances are averaged to assure smooth transitions between frames. If for tracked pedestrians no new detection is found, the Kalman filter is updated based on the estimated position. In this case we apply a dynamic score strategy, and lower the detection threshold for that instance (within certain boundaries). This ensures that pedestrians which are difficult to detect (e.g. partially occluded or a temporarily low HOG response) can still be tracked. If no detection is found for multiple frames in a row, the tracker is discarded. Evidently, if a new detection is found for which no previous tracker exists, tracking starts from there on. Figure 11 qualitatively illustrates tracking sequences on two of our datasets.

4 Experiments and Results

We performed extensive experiments concerning both speed and accuracy. Our datasets consists of simulated dangerous blind spot scenarios, recorded with a real truck. We used a commercial blind spot camera (Orlaco CCC115°) with a resolution of 640×480 at 15 frames per second. This camera has a 115 degree wide-angle lens. See Fig. 12 for the exact position of the camera. Five different scenarios were recorded, each in which the truck driver makes a right turn and the pedestrians react differently (e.g. the truck driver lets the pedestrians pass, or the truck driver keeps on driving, simulating a near-accident). This resulted in a total of about 11000 frames. For our accuracy and speed experiments we labelled around 3200 pedestrians. Our implementation is CPU-based only, and



Fig. 12. Our test truck with the mounted commercial blind spot camera (circled in red) (Color figure online).

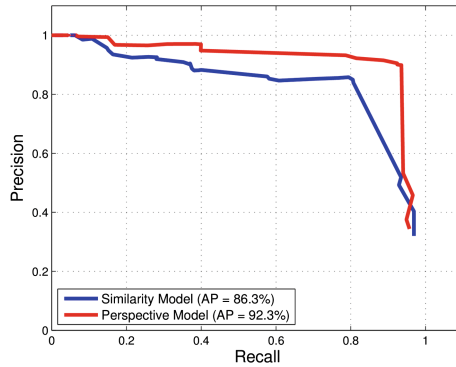


Fig. 13. Precision-recall graph over our dataset.

the hardware consists of an Intel Xeon E5 CPU which runs at 3.1 GHz. Note that all speed experiments are performed on a single core. The algorithm is mainly implemented in Matlab, with time-consuming parts (such as the homography) in OpenCV, using *mexopencv*.

4.1 Accuracy Results

Figure 13 displays the precision-recall graph of our algorithm as calculated over our datasets. The red PR curve indicates our novel perspective transformation approach, while the blue PR curve represents our previous similarity transformation approach. They are calculated as follows. For each detected pedestrian in our algorithm, we look for a labeled instance in a circular region (based on the scale) around the centre of our detection. If such an instance is found, this is counted as being a *true positive*. If this is not the case, this detection is counted as being a *false positive*. Each labeled pedestrian which is not detected accounts for a *false negative*. The PR-graph is then determined as: $precision = \frac{TP}{TP+FP}$ and $recall = \frac{TP}{TP+FN}$. We notice that, although both achieve very good accuracy results, our novel perspective warping window approach has a clear accuracy

advantage over our similarity warping window approach. Indeed, the average precision (AP) for the similarity model equals 86.3%, whereas for the perspective model $AP = 92.3\%$. With the perspective model, at a recall rate of 94%, we still achieve a precision of 90%. Such high accuracy results are due to our warping window approach. Since we know the scale at each position, the number of false positives is minimized. Furthermore this allows us to use a sensitive pedestrian detection threshold.

4.2 Speed Results

As mentioned in Sect. 3.1, if used out-of-the-box the baseline pedestrian detector takes 670 ms (i.e. 1.5 fps). Since in our framework we only need to perform detection at a single scale and ROI, the calculation time drastically decreases. For each default search region and tracked pedestrian in the image we need to perform a warp operation and detection. Thus, the total calculation time evidently depends on the number of tracked pedestrians per image. Figure 14 displays the average calculation time per ROI. Note that if a detection is found, the average calculation time equals 18.3 ms, while if no detection is found the average calculation time drops to 10.8 ms. This calculation time per region is independent of the position in the image. The average detection time per ROI is subdivided into five steps: the calculation of the warp coordinates, the time needed to perform the warp operation, calculation of the HOG features, evaluation of the pedestrian model, and finally the retransformation of the detected coordinates to the input image. The total warp time (*calc. warp coord.* and *perform warping*) only equals about 3 ms. Most time is spent on the actual pedestrian detection. The time needed to perform the retransformation of the coordinates is negligible. Figure 15 displays the frames per second as a function of the number of tracked pedestrians we reached on our datasets. If no pedestrians are tracked we achieve 28.2 fps. On average we achieve 13.0 fps (with an average of 3.4 pedestrians), while our worst-case framerate equals 7.0 fps.

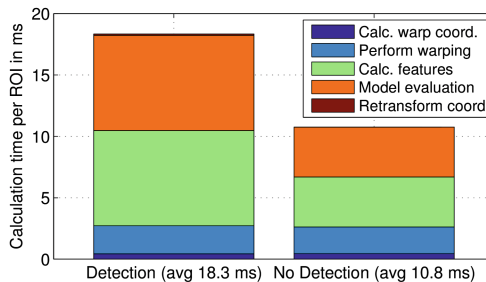


Fig. 14. Calculation time per ROI.

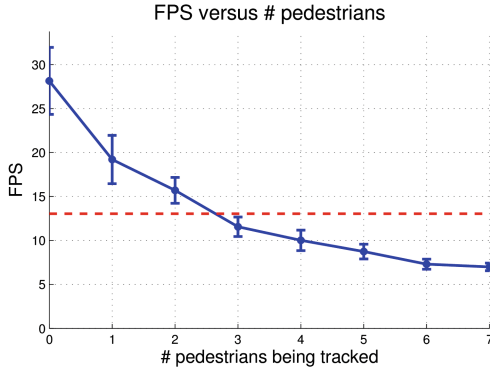


Fig. 15. Speed performance versus the number of tracked pedestrians (dotted red line indicates the average fps) (Color figure online).

5 Conclusions and Future Work

In this work we proposed a multi-pedestrian tracking framework achieving excellent accuracy and speed results on a single-core CPU implementation. The algorithm is based on our novel perspective warping window approach. We proposed this approach to allow for efficient pedestrian detection on the challenging, highly distorted camera images from a blind-spot camera, with minimal CPU resources. However, this approach is easily generalisable to other applications with non-standard camera-viewpoints.

In the future we plan to further extend our framework to multi-class detection: we aim to develop a complete vulnerable road users detection system, starting with bicyclists. Additionally, an efficient combination of multiple pedestrian detectors to further increase the accuracy could be evaluated. Furthermore we aim to investigate if the inclusion of other features (e.g. motion information) could further increase the robustness of our framework.

References

1. Benenson, R., Markus, M., Tuytelaars, T., Van Gool, L.: Seeking the strongest rigid detector. In: Proceedings of CVPR, pp. 3666–3673, Portland, Oregon (2013)
2. Benenson, R., Mathias, M., Timofte, R., Van Gool, L.: Fast stixels computation for fast pedestrian detection, In: ECCV, CVVT workshop, pp. 11–20 (2012)
3. Benenson, R., Mathias, M., Timofte, R., Van Gool, L.: Pedestrian detection at 100 frames per second. In: Proceedings of CVPR, pp. 2903–2910 (2012)
4. Cho, H., Rybski, P., Bar-Hillel, A., Zhang, W.: Real-time pedestrian detection with deformable part models. In: IEEE Intelligent Vehicles Symposium, pp. 1035–1042 (2012)
5. Dalal, N., Triggs, B.: Histograms of oriented gradients for human detection. Proc. CVPR **2**, 886–893 (2005)
6. Dollár, P., Belongie, S., Perona, P.: The fastest pedestrian detector in the west. In: Proceedings of BMVC, pp. 68.1–68.11 (2010)

7. Dollár, P., Tu, Z., Perona, P., Belongie, S.: Integral channel features. In: Proceedings of BMVC, pp. 91.1-91.11 (2009)
8. Dollár, P., Wojek, C., Schiele, B., Perona, P.: Pedestrian detection: a benchmark. In: Proceedings of CVPR, pp. 304–311 (2009)
9. Dollár, P., Wojek, C., Schiele, B., Perona, P.: Pedestrian detection: an evaluation of the state of the art. *IEEE PAMI* **34**, 743–761 (2012)
10. Enzweiler, M., Gavrilu, D.M.: Monocular pedestrian detection: survey and experiments. *IEEE PAMI* **31**, 2179–2195 (2009)
11. Ess, A., Leibe, B., Schindler, K., Van Gool, L.: A mobile vision system for robust multi-person tracking. In: Proceedings of CVPR, pp. 1–8 (2008)
12. EU (22 february 2006). Commission of the european communities, european road safety action programme: mid-term review
13. Felzenszwalb, P., Girschick, R., McAllester, D.: Cascade object detection with deformable part models. In: Proceedings of CVPR, pp. 2241–2248 (2010)
14. Felzenszwalb, P., McAllester, D., Ramanan, D.: A discriminatively trained, multi-scale, deformable part model. In: Proceedings of CVPR (2008)
15. Gavrilu, D., Munder, S.: Multi-cue pedestrian detection and tracking from a moving vehicle. *IJCV* **73**, 41–59 (2007)
16. Huang, C., Ai, H., Li, Y., Lao, S.: Vector boosting for rotation invariant multi-view face detection. In: ICCV, pp. 446–453 (2005)
17. Lampert, C., Blaschko, M., Hoffmann, T.: Efficient subwindow search: a branch and bound framework for object localization. *IEEE PAMI* **31**, 2129–2142 (2009)
18. Martensen, H.: Themarapport vrachtwagenongevallen 2000–2007 (BIVV) (2009)
19. Mathias, M., Timofte, R., Benenson, R., Van Gool, L.: Traffic sign recognition - how far are we from the solution? In: ICJNN (2013)
20. Park, D., Ramanan, D., Fowlkes, C.: Multiresolution models for object detection. In: Daniilidis, K., Maragos, P., Paragios, N. (eds.) ECCV 2010, Part IV. LNCS, vol. 6314, pp. 241–254. Springer, Heidelberg (2010)
21. Pedersoli, M., Gonzalez, J., Hu, X., Roca, X.: Toward real-time pedestrian detection based on a deformable template model. In: IEEE ITS (2013)
22. Prisacariu, V., Reid, I.: fastHOG - a real-time gpu implementation of HOG. Technical report, Department of Engineering Science, Oxford University (2009)
23. Seitner, F., Hanbury, A.: Fast pedestrian tracking based on spatial features and colour. In: Proceedings of CVWW, pp. 105–110 (2006)
24. Van Beeck, K., Goedemé, T., Tuytelaars, T.: Towards an automatic blind spot camera: robust real-time pedestrian tracking from a moving camera. In: Proceedings of MVA, Nara, Japan (2011)
25. Van Beeck, K., Tuytelaars, T., Goedemé, T.: A warping window approach to real-time vision-based pedestrian detection in a truck's blind spot zone. In: Proceedings of ICINCO (2012)
26. Viola, P., Jones, M., Snow, D.: Detecting pedestrians using patterns of motion and appearance. *IJCV* **63**, 153–161 (2005)