

A Video-Based Object Detection System for Improving Safety at Level Crossings

N. Fakhfakh^{1,2}, L. Khoudour¹, E.M. El-Koursi¹, J. Jacot², A. Dufaux²

¹French National Institute for Transport and Safety Research (INRETS): 20, rue Elisée Reclus 59666 Villeneuve d'Ascq, France {nizar.fakhfakh, louahdi.khoudour, el-miloudi.el-koursi}@inrets.fr

²Ecole Polytechnique Fédérale de Lausanne (EPFL): CH-1015 Lausanne, Switzerland {nizar.fakhfakh, jacques.jacot, alain.dufaux}@epfl.ch

Abstract: Improving transport users' safety is one of the main priorities of research into transport system attractiveness. Level crossings are one of the most critical weak point involving road and rail users' infrastructure. They have become increasingly dangerous and unsafe due to road and railway users' behavior. Furthermore, rail and highway safety professionals from several countries must deal with the same subject: providing safer level crossing. Actions are planned in order to exchange and share knowledge on existing level crossings technologies between academic organizations and industrial operators, and provide experiments for improving the management of level crossing safety and performance. This has enabled us to discuss sharing knowledge gained from research in order to improve safety at level crossings. This article provides research results about possible technological solutions to reduce the number of accidents at level crossings. As a main contribution is that we discuss and prove the effectiveness of the use of video sensing for object detection. Furthermore, we have tested and adapted a robust technique for moving object detection, which is followed by a new approach for 3D object localization.

Keywords: Image processing, level crossings, safety, technologies, Vision.

1. INTRODUCTION

Within the past years, railways undertakings became interested in the assessment of level crossings safety. Level crossings have been identified as a particular weak point in road and railway infrastructure, seriously affecting their safety. Statistically, more than four hundred people year die in accidents involving road and vehicles at road-rail level crossings in the European Union (EU) [1]. 90% of these fatalities are linked to errors committed by road vehicle drivers. However, the behavior of road vehicle drivers and pedestrians cannot be previously estimated beforehand. It is important also to note the high cost related to each accident (approximately one hundred million euro per year in the EU for all level crossing accidents). For this purpose, road and highway safety professionals from several countries have dealt with the same subject: providing safer level crossings. Actions are planned in order to exchange information and provide experiments for improving the management of level crossing safety and performance. This has enabled us to discuss sharing knowledge gained from research into improving safety at level crossings. In recent times, the coordinated action for the sixth framework program entitled "Safer European Level Crossing Appraisal and Technology"

(SELCAT) [2] has provided recommendations for further actions, evaluation of possible technological solutions to improve safety at level crossings. The proposed solution must reduce the risk presented by road drivers or railway operators. Another interesting project was led by Japan entitled Intelligent Transport System (ITS) [3]. This project is steadily progressing and aims to make transportation systems safe and convenient. Furthermore, the evolution of technologies allows researchers to test and evaluate the performances of several devices at level crossings. Each device should prove an acceptable level of risk with which a given system could be chosen for implementation.

Nevertheless, the lack of an approved common safety methodology leads to the imposition of the highest safety integrity levels for technical solutions. High safety requirements for level crossing systems mean a high cost which hinders the technological setup of advanced systems. The systems having unacceptable levels of risk must be eliminated and doesn't be implemented. High technology systems are exploited and introduced in order to prevent collisions between trains and automobiles and to help reduce levels of risk from railroad crossings. Several conventional object detection systems have been tested on railroad crossings. These techniques provide more or less significant information accuracy.

An overview of existing technological solutions and a proposed object detection system based on a passive vision system, such as television camera, constitute the

Our thanks go to TL ("Tranports publics de la région Lausannoise"), which is the transportation company of the Lausanne city in Switzerland. By allowing the acquisition of images at several level crossings on their train lines, the TL allowed us progressing in our study, in particular for the evaluation of the proposed solutions."

main part of the research described in this paper. After an introduction covering the problem of level crossings safety and the motivations to improve safety at level crossings, the second part deals with a survey of conventional systems of detection of obstacles at level crossings. We discuss also in this part the advantages and the drawbacks for each system. The third part will be devoted to arguing for and describing the choice of the video sensing and image processing to alleviate the previous drawbacks. This is followed in the fourth section by outlining the proposed new system for object detection based on a passive vision. An overview of the proposed system is presented which will be divided into two parts: a foreground detection based on a single camera and a novel approach for 3D localization of a detected moving foreground. The proposed approach is evaluated in the last part of this paper using conventional datasets of images and tested on real level crossing image dataset. The conclusion part is devoted to a discussion on the obtained results.

2. REQUIREMENTS AND CONVENTIONAL SENSORS AT RAILROAD CROSSING

2.1. REQUIREMENTS

The most reliable solution to decrease the risk and accident rate at level crossings is to eliminate unsafe railroad crossings. This avoids any collisions between trains and road users. Unfortunately, this is impossible, due to location feasibility and cost that would be incurred. To overcome these limits, development of a new obstacle detection system is required. Possible solutions using advanced technologies to improve level crossing safety will be discussed in this part. The proposed system is not intended to replace the present equipment installed on each level crossing. The purpose of such a system is to provide additional information to the human operator. This concerns the detection and tracking of any kind of objects, such as pedestrians, people on two-wheeled transport, wheelchairs and car drivers. Presently, sensors are evaluated relying on their false object detection alert among other. This may increase the risk related to level crossing users. It is important to be noted that risks associated with the use of technology systems are becoming increasingly important in our society. Risk involves notions of failure and consequences of failure. Therefore, it requires an assessment of dependability; this might be expressed, for example, as probability of failure upon demand, rate of occurrence of failures, probability of mission failure, and so on. Each level crossing is equipped with various sensors for tracking and timely detection of potentially hazardous situations. To be reliable, the related information must be shared and transmitted to the train dispatching center, stations, train drivers and road users. A brief summary of the technical characteristics of most devices provides the necessary understanding to design effective obstacle detection system. Generally, most level crossings are fitted with high performance equipment such as lights, automatic full or half barriers, notices. This equipment warns and prevents all users of

the level crossing if a train is approaching the dangerous area.

2.2. CONVENTIONAL SENSORS AT LEVEL CROSSINGS (LC)

Today, there are a number of trigger technologies installed at level crossings, but they all serve the same purpose: they detect moving object when passing at particular points in the LC. However, those conventional obstacle detection systems have been used to prevent collisions between trains and automobiles. Several technologies, such as optical beam provide this (Figure 1). The principle is as follows: optical emitters are placed on one part of the crossing, each one emitting a directed optical beam with a defined field of emission. Then a photon detector having a defined field of view intersects the field of emission of the emitter. If the beams are interrupted, it means that an object is located on the crossing.



Figure 1. Optical beam method.

This technique has the advantage that is easy to replace, but it has many important drawbacks: very expensive cost for installation, need to have several detectors along the crossing, traffic needs to be stopped for installation, and it is unusable in period of heavy snow. Electronic waves passing between transmitters and receivers can detect obstacles in a similar manner to the optical beam method. However, optical beam method cannot detect pedestrians.

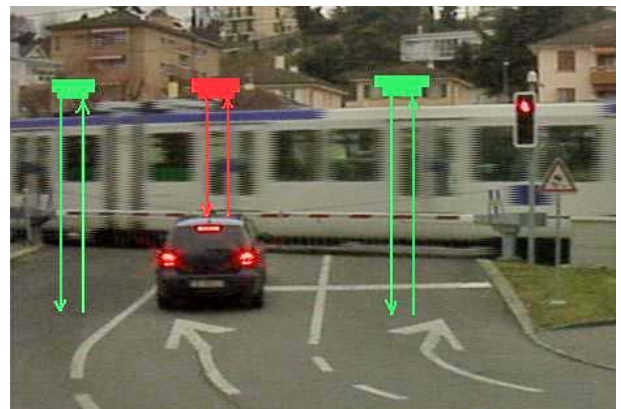


Figure 2. Ultrasonic method.

Another technique consists in using ultrasonic detectors (Figure 2) which rely on differences in ultrasonic reflection times for detection. They transmit pulses of ultrasonic energy towards the roadway. The pulse is then reflected back more quickly when a vehicle passes through. As an advantage, sonic detectors have the advantage that can detect both stationary and moving vehicles. However, the disadvantages resides in their price and installation cost. Additionally, they are extremely sensitive to environmental conditions (inaccurate in congested conditions). Furthermore, the common drawback of these systems is that they do not detect object with low metal content.

Another kind of obstacle detection system has recently been introduced on the railroad crossing such as radar device (Figure 3). The principle of this technique is as follow: Microwaves are sent from a transmitter based at the side of the roadway. The microwaves are reflected back to a receiving antenna with a different frequency. This change is picked up and reflects the presence of an object. Radar presents a lot of advantages such as the fact that traffic does not need to be disrupted for installation and it is immune to electromagnetic interference. The disadvantage is that it is hard to maintain.

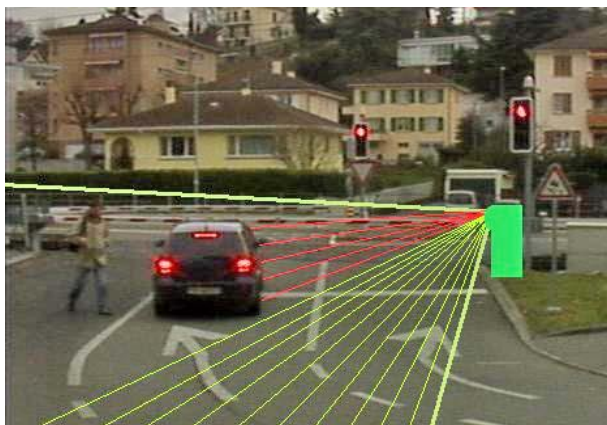


Figure 3. Radar method.

The inductive loop is probably the most common form of vehicle detection. A wire is embedded under the roadway. A magnetic flux is generated around it. When a vehicle is passing over the wire, the flux is cut causing an increase in inductance, and the detection of an obstacle. This is what happens when a car pulls up to the intersection. The huge mass of metal that makes up a car alters the magnetic field around the loop, changing its inductance. This system is easy to install and is not subjected to environmental conditions. The main drawbacks are that it cannot detect pedestrians, and the high cost of installation and maintenance, because it needs a large number of loops to be efficient.

Another object detection system has recently been designed by Hisamitsu [4] based on a 3D laser radar technique in which the principle is the following: 3D laser radar emits a laser pulse to an object, and measures the time that it takes for the reflected laser to return to the radar (time-of-flight method). It acquires the distance

from that object. Figure 4 shows the object detection method used by the 3D laser radar. A laser pulse is emitted in a way that it scans the entire area of a level crossing in the horizontal and vertical directions. The 3-D coordinate's values of each point measured from the reflected laser are returned to the 3-D laser radar. Summarizing, Optical beam and sonic detectors are definitely too expensive and not accurate enough because they need a high number of sensors.

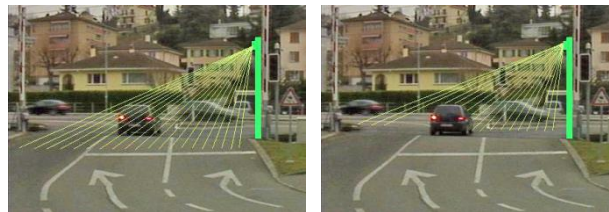


Figure 4. Laser method.

They are also too sensitive to environmental conditions. The most relevant technologies for object detection at level crossings are radar, inductive loops and video imaging. The problem with inductive loops is that they detect only metal objects. Nowadays, many vehicles use fiber glass and aluminum, which are not well detected by the system. Afterwards, the maintenance of such a system is hard because it is located below the road. That is why we are focusing on video and now present some realizations and utilizations of these technologies.

3. VIDEO IMAGING ON LEVEL CROSSINGS

One of the main operational concerns for the introduction of CCTV (Close Circuit Television) at LC for the purpose of automatic detection of specific events. It is envisaged that the new technologies will increase the number of monitored camera images, give the staff earlier warnings to be able to produce more timely responses, and support the recording and collection of evidence of detected events. The level crossing obstacle detection system that uses cameras detects any kind of objects such as pedestrians and wheelchairs that cannot be detected by the conventional systems of the photoelectric (optical beam), ultrasonic, and loop coil types. These particular objects that can be detected using passive vision may be tracked close together in space and time. This allows us to model spatiotemporal behavior from image sequence. Nowadays, most of a conventional object detection systems described above supervises only dangerous crossing areas. According to the technology used, the detected object is tracked when it appears in the critical zone. For this purpose, using artificial vision can provide more relevant information and with such a system, dangerous scenarios can be recognized at the right time. Figure 5 presents a reference scenario of hazardous situations detected from a stereo camera. Interpreting this scenario allows us to deduce the importance of the level crossings environment which gives significant information such as road traffic and the status of red-light. This additional

information can be provided as input to the computational system related to the cameras.

3.1. OVERVIEW OF EXISTING SYSTEMS

Referring to the literature, little research has focused on passive vision to solve the problems at level crossings. Two main systems based on CCD cameras can be distinguished:

- A level crossing obstacle detection system using one single camera: This system using video cameras has been developed by G.L.Foresti [5]. This system uses one single CCD camera placed on a high pole in a corner of the level crossing. First, the object detection is computed based on the difference between current and background images. Then, with the knowledge of intrinsic camera parameters, calibration matrix and ground plane hypothesis, the 3D position of the different objects found are computed. Then objects are tracked with an Extended Kalman Filter (EKF). Finally, object classification based on morphological statistical spectrum is performed in order to classify objects as car, bike, trunk, pedestrian, dog, paper, etc. Tests proved that the system works well in different situations (different camera points of view, bad environmental conditions, noise, object occlusions). But this system proved to be

limited in low illumination and the presence of shadows can lead to false alarm.

- A level crossing obstacle detection system using stereo cameras: M.Ohta [6] has developed a level crossing obstacle detection system using stereo cameras. The problem of using a single TV camera acquiring a wide-area image is car headlights or shadows, which cause errors in detection. The idea of Ohta is to use a stereo camera whose principle will be described later. The object is then projected on each TV screen with a disparity, and with some computation, a 3D shape of the object is obtained. It reduces false detection of car lights or shadows. Then the disparity image (which shows the distance of an object against the camera) is computed. Finally, the 3D shape of the whole scene is extracted and compared with the background shape. The obstacle detection using stereo cameras can detect both vehicles and pedestrians. The tests have proved that the system correctly detects objects during day and night under general weather conditions. But the main problem with this system (and with image processing in general) is that it is extremely sensitive to adverse weather conditions, like heavy rain, fog or snow.



Figure 5. An example of hazardous situation scenario detected from cameras devices. (a) Pedestrians near a crossing, (b) appearance of a car, (c) A car approaching the crossing (first level of red-light sensor), (d) A car approaching the crossing (second level of red-light sensor), (e) A car on the crossing (third level of red-light sensor: a mandatory stop), violation of red-light, (f) A car stopped on the critical area, (g) Start of lowering barriers, (h) Car driver is still on the dangerous zone, (i) A barrier is lowered and jammed on the car, (j) Car driver has forced the barrier: It is full track, (k) The car passing the hazardous area, (l) A free and safe level crossing (barriers are still lowered).

3.2. PROPOSED SYSTEM DESIGN

In order to overcome the limits described above, the passive vision principle as shown in Figure 6 will be adopted in our research which is based on the following two steps: Foreground extraction from image sequences and 3D localization of detected moving foreground. For motion object segmentation, the background subtraction method is implemented. This technique proves effective for outdoor image sequences especially with changes in brightness. The moving objects detected with a single camera will be taken into account for the 3D localization module. This relevant information will be useful for localization and recognition processes. To have safe information for 3D localization, the stereo method is only applied to particular points of the moving object. The remaining depth points will be reassessed on the basis of well matched neighboring points. The newly developed system was tested with a stereoscopic sensor sets in order to cover the whole area of a double-track level crossing. These sensors are placed beside a level crossing for field tests. The rest of this paper will focus on the two previous techniques. That is to say foreground extraction and 3D localization.



Figure 6. A level crossing (in Ecublens-Switzerland): obstacle detection system that utilizes stereo cameras and image processing technology.

4. FOREGROUND EXTRACTION

4.1. MOVING OBJECT DETECTION

In video surveillance, detection of moving objects from an image sequence is very important for target tracking, activity recognition, and behavior understanding. Motion detection aims at segmenting foreground regions corresponding to moving objects from the background. Background subtraction and temporal differencing are two popular approaches to segment moving objects in an image sequence under a stationary camera. Temporal differencing calculates the difference of pixel features between consecutive scene frames in an image sequence. It is very effective to accommodate environmental changes, but generally can only detect partial edge shapes of moving object. In this study, the first techniques family for foreground

subtraction with background model updating strategies was adapted. Background subtraction is a very popular approach for foreground segmentation in a still scene. It detects moving objects in an image by evaluating the difference of pixel features of the current scene image against the reference background image. This approach is very sensitive to illumination changes without adaptively updating the reference background. However, the background in a long image sequence is dynamic even if it is captured by a stationary camera. In this context, the generalized mixture of Gaussians proposed in [7] has been used to model complex, non-static backgrounds. Methods employing a mixture of Gaussians have been widely incorporated into algorithms that utilize Bayesian frameworks [8], dense depth data [9], color and gradient information [10], mean-shift analysis [33], and region-based information [11]. A mixture of Gaussians does have some drawbacks. Backgrounds having fast variations are not easily and accurately modeled with just a few Gaussians, and it may fail to provide sensitive detection. In addition, depending on the learning rate to adapt to background changes, a mixture of Gaussians faces trade-off problems. For a low learning rate, it produces a wide model that has difficulty in detecting a sudden change in the background. If the model adapts too quickly, slowly moving foreground pixels will be absorbed into the background model, resulting in a high false negative rate. This is the foreground aperture problem described in [12]. Recently, a new method has been proposed by Kim [13], who describes a technique for discriminating moving objects from the background. The adopted background subtraction scheme involves two stages, one for training and the other for detection. In the training stage and in order to make the background subtraction adaptive to environmental changes, such as illumination variations, sample background values at each pixel are quantized into codebooks which represent a compressed form of background model for a long image sequence. This allows us to capture structural background variation due to periodic-like motion over a long period of time under limited memory. In the detection stage, a color distance between a given pixel and its correspondent in the trained background is computed to separate the foreground in a scene image with respect to the reference background image.

The codebook algorithm adopts a quantization technique to construct a background model from long observation sequences. For each pixel, it builds a codebook consisting of different states of pixel based on color and intensity. Samples at each pixel are clustered based on a color distortion metric together with brightness bounds. The clusters do not necessarily correspond to single Gaussian or other parametric distributions. The background is encoded on a pixel-by-pixel basis. Detection involves testing the difference of the current image from the background model with respect to color and brightness differences. An incoming pixel must be classified into background or foreground class based on color and brightness distortion referring

to background estimated model. The main advantage of this method is that it works in various environments in which the lighting conditions may change. As a result, the chosen method for foreground extraction is perfectly suited to level crossing application. It is effective for motion detection under a various environmental conditions due to updating background. Moreover, the fundamental challenges is to achieve invariance to illumination changes, and more prominently, to shadows and highlights. The detected moving object includes the real object and its shadow. This over-segmentation will interfere the recognition and classification process which can cause therefore false alarms. It is to be noted that cast shadows (it's created by an object blocking the light source) is the type of shadow that we are interested in.

4.2. SHADOW REMOVAL

Dealing with shadows and highlights is essential in object detection and tracking applications such as automated video surveillance systems. Shadows occur frequently in a wide variety of scenes. It is more problematic for outdoor scenarios subject to variable lighting and weather conditions. However, this is undesirable due to the fact that they often lead to the result of irretrievable processing failures. For instance, the shadow cast by an object results in an improper segmentation result with serious artifacts, or detection of an imaginary object. This might result in shadows misclassified as objects or part of object due to the overestimation in a subsequent matching phase. Hence, the obtained moving object is the sum of the real object and its shadows. We distinguish two types of shadows: Cast-shadows and self-shadows. Only the first type of shadows will be taken into consideration for shadow removal process. Cast-shadows refer to areas in the background projected by objects in the direction of light rays, producing distorted objects silhouettes. Two main techniques for shadow removal exist: model-based and property-based method. The first technique needs certain a priori knowledge (geometry features of the scene, illumination direction,...). The second technique is the most important because it identifies shadows by exploiting their properties in brightness and color. Color information is exploited by means of several color spaces, and by means of photometric invariant features. These methods are mostly based on the fact that shadows change the luminance but color little. Several researches proposed shadow detection methods based on the RGB, HSV and YUV [29] color spaces. Most of research assumes that a shadow belongs to a dark region. By exploiting color information, we overcome the limitation of the dark regions classification process of [30], which assumes that cast shadow pixels are darker than self shadow pixels as an empirical criterion for shadow classification.

Thereby, another property of shadow can be explored since shadow does not alter the value of the invariant color features. On the contrary, a material change modifies their value. This property has been explored by

Salvador [32] using the $c_1c_2c_3$ photometric invariant color features. The identification of shadow pixels is achieved by analyzing the difference in the invariant feature values between the pixel to be classified and it correspondent in the background. Examples of photometric invariant features are Hue and Saturation in the HSV color space and the normalized-RGB color space (rgb) which is tested by Cavallaro [31]. We have adapted the technique proposed by Porikli [33] for shadow removal. Likelihood of being a shadow pixel is evaluated iteratively by observing the color space attributes and local spatial properties. The author carried out the evaluation of the proposed method in the RGB color space, and assumed that shadow decreases the luminance and changes saturation, yet it does not affect the hue. However, the detection of shadows is based on the fact that shadows change significantly the lightness of an area without greatly modifying the color information. The algorithm is designed to be able to work when camera, illumination and scene's characteristics are unknown. Shadows cast on a surface reduce the surface intensities. It is important to be noted that there is not a conventional method for shadows removal which can be applied to any image sequences. Figure 7 represent the results of moving object detection applying the codebook method described above that we have improved. Three image dataset for real environment at different level crossings are used for the experiments.

5. 3-D LOCALIZATION OF THE DETECTED OBJECT: OBJECT DEPTH ESTIMATION

Stereo vision is introduced and has become one of the most extensively researched topics in computer vision. This research has partly solved some problems related to object detection and recognition in various types of scenes. The use of multi-camera systems provides additional information, such as the depth of objects in a given scene. Dense or sparse stereovision techniques can be used to match points. In dense stereovision, all points of an input image are taken into account for matching tasks. Each disparity, determined for each point of the scene, represents the coordinate gap of the same point between the two left- and right-hand images representing the same scene from two different points of view. A depth map is obtained from the two images. Each value of this obtained map is an estimation of the distance between a real point and the stereoscopic sensor and is given using intrinsic and extrinsic parameters related to the used sensors, such as the focal length and the baseline. When a point is imaged from two different viewpoints, its image projection undergoes a displacement from its position in the first image to that in the second image. The amount of displacement, alternatively called disparity, is inversely proportional to distance and may therefore be used to compute 3D geometry. Given a correspondence between imaged points from two known viewpoints, it is possible to compute depth by triangulation. The problem of establishing correspondence is a fundamental difficulty

and is the subject of a large body of literature on stereo vision. One prominent approach is to correlate pixels of similar intensities in two images, using an assumption that each scene point reflects the same intensity of light in the two views. The work carried out by Okutomi and Kanade [14] extends this correlation approach to two, three or more images and demonstrates that using several cameras at different camera separations, or baselines, yields a significant improvement in

reconstruction accuracy.

Most two-frame stereo matching approaches compute disparities and detect occlusions assuming that each pixel in the input image corresponds to a unique depth value. Until recently, computers were much too slow to even dream of having a real-time algorithm implementation of an algorithm based on color stereovision.

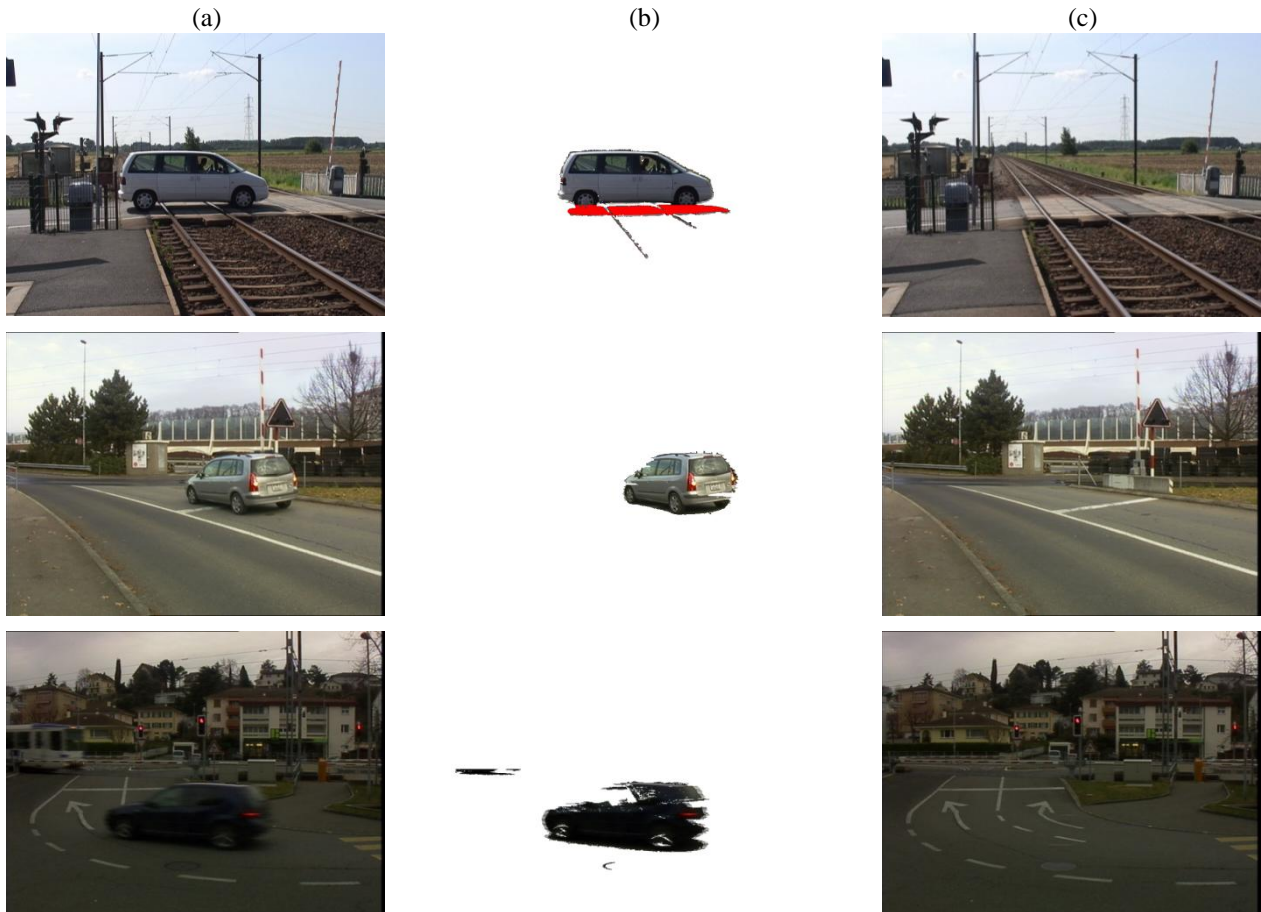


Figure 7. Foreground extraction process. (a) Original image, (b) foreground extraction of moving object applying codebook method. The shadow is detected and represented by red color (c) the updated background.

But costs and processing time are decreasing at a steady pace, and it is becoming realistic to believe that such a thing will be commonplace soon. The stereo algorithm presented later, springs from relaxation and energy minimization fields. Our approach aims to represent a novel framework to improve color dense stereo matching. As a first step, disparity map volume is initialized applying a new local correlation function. After that, an assessment of matching quality is addressed. A confidence measure is attributed to all pairs of matched pixels which are classified into three classes. Then, the disparity value of all unclassified or badly-

matched pixels is updated based only on stable pixels classified as well-matched. The selected pixels are used as input into disparity re-estimation modules to update the remaining points. Our main contribution is that we propose a novel approach to compute a confidence measure for each pair of pixels. The confidence measure is based on a set of original local parameters related to the correlation function used in the first step of our algorithm. This paper will focus on this point. The main goal of our study is to take into account both quality and speed.

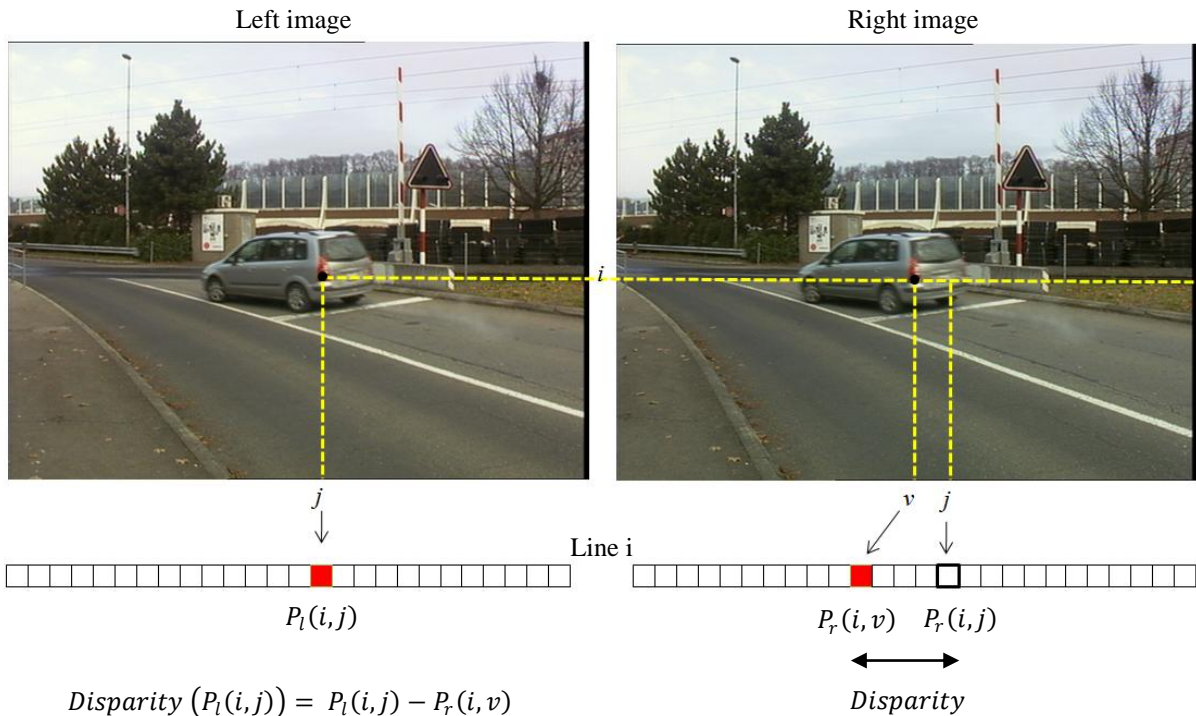


Figure 8. Stereo vision principle for estimation of depth foreground moved object.

5.1 OVERVIEW OF THE PROPOSED STEREO MATCHING ALGORITHM

We present, in this section, an overview of our stereo matching (the principle is shown in Figure 8) algorithm which allows us to improve the accuracy of disparity maps. The algorithm can be divided into three parts: Initialization disparity map, pixel classification and disparity allocation. In the first step, we compute the correlation volume. Thus, we take advantage of local methods by applying a new color window correlation to build the correlation volume. It is called Weighted Average Color Difference (WACD) [15]. Firstly, this dissimilarity function is applied to all the pixels in stereo images using a non-adaptive square correlation window. This local method allows an initial appraisal of the disparity map. The second part of our algorithm aims to classify more accurately the matched pixels. Stereo matching accuracy may be affected by various factors including feature descriptors, similarity measures and matching approach. We assume that depth, generally derived from disparity, varies smoothly within color homogeneity in a given region. Depth discontinuities coincide generally with color boundaries or edges. Under the previous constraints, all matched pixels are classified into three categories according to their location: well-matched, badly-matched pixels and not-classified pixels. A new classification method based on a confidence measure approach is applied in this context. This confidence measure is computed for all matched pixels and is based on a set of local parameters referring to scores obtained from a new dissimilarity function. This method and the associated parameters are detailed later in this paper. This work is based on the principle of

relaxation and the belief propagation method which are based on global criteria. However, the difference is that we consider only candidate pixels to evaluate the matched pair.

5.2. WEIGHTED AVERAGE COLOR DIFFERENCE CORRELATION FUNCTION

Disparity map is initialized applying a local method. The local approach consists in taking into account a set of neighborhoods pixels belonging to a support window centered on the pixels to be matched. Each pixel in the left-hand image has K candidate pixels in its corresponding right-hand image. Each pair of pixels, which is composed of both the pixel to be matched and a candidate pixel, is evaluated based on their neighborhoods pixels. We have exploited additional information such as color components to develop a new local correlation function. It is applied to each pairs of pixels and allows giving a dissimilarity measure. The candidate pixel, with which the correlation function gives the lowest correlation score, is kept as the best candidate. The lower the score, the better the matching. To deal with the image ambiguity problem, area-based local methods use, in most of cases, some kind of statistical correlation between color or intensity patterns in local support windows. This new correlation function is applied with a given support window. Neither the size nor the shape is treated in this study. Among the state-of-art local stereo matching function, we can mention the well-known algorithm named Sum of Absolute Differences (SAD) [16] which employs all neighborhoods pixels belonging to support windows. The particularity of our proposed correlation function is

that only lines belongs to used support window and passed through the central pixel are taken into account to compute dissimilarity measure between two given windows. Lines taken into account are the horizontal, the vertical and the two diagonals. The principle is explained in Figure 9:

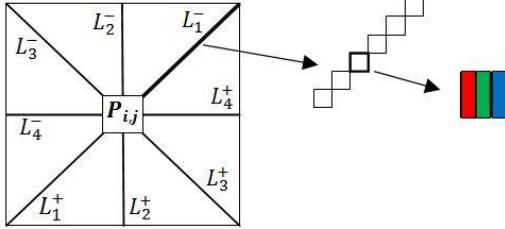


Figure 9. Segments taken into account by the proposed correlation function to compute a dissimilarity measure between two support windows.

The lines are divided into segments and for each among them the average of each color component is computed.

Let denote L_n , $n = 1, 2, 3, 4$ (resp. L_n^+) lines through the central pixel in the left-hand image (resp. right-hand image) having (i, j) coordinates. We take into account as attributes all color components to compute the correlation score.

Let denote also $\psi_{i,j} = \{L_n^+, L_n^-\}$ (resp. $\psi'_{i,j} = \{L_n^+, L_n^-\}$) a set of segments of each lines located on both sides of the central pixel belonging to left (resp. right) support window. In the first time, we are computing the average of each color component for each segment. These average colors are stored as a vector denoted as $v_k = \bar{\psi}_{i,j}$, $k = R, G, B$ (resp. $v'_k = \bar{\psi}'_{i,j}$).

Moreover, we have defined a measure denoted D which represents the distance difference between the k

component of the central pixel and an average of the same component of a segment v_k . This measure is given by the equation 1:

$$D = \text{MAX}(|(P_{i,j,k} - v_k) - (P'_{i,j,k} - v'_k)|, 1) \quad (1)$$

The proposed distance D is weighted by a second measure denoted W which represent the difference between the two quantities δ_l and δ_r . δ_l (resp. δ_r) representing the average of each color component k for a given segment v_k of the left-hand image (resp. right-hand image). It is defined as follows:

$$W = \text{MAX}(|\delta_l - \delta_r|, 1) \quad (2)$$

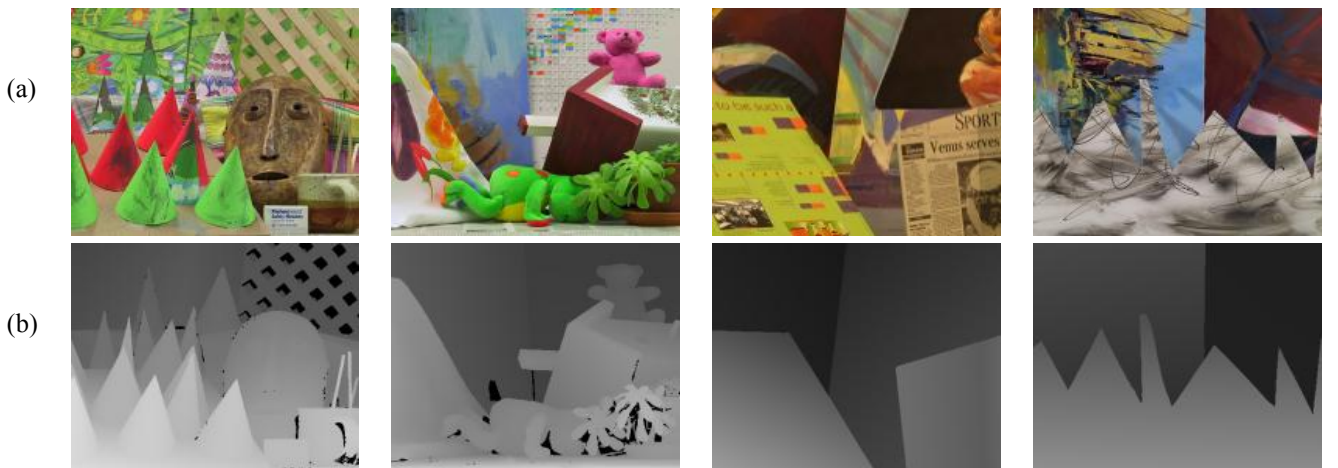
where $\delta_l = \frac{(P_{i,j,k} + v_k)}{2}$ and $\delta_r = \frac{(P'_{i,j,k} + v'_k)}{2}$.

For all primary system such as RGB color space, it has often been proposed notably by Crouzil [17] to use the Euclidian distance to calculate color difference between two pixels. We have adopted in our context this metric for $n = 2$. The color difference between pixel $P_{i,j}$ and $P'_{i,j}$ is expressed as (such as $n = 2$):

$$\Delta^n(P_{i,j}, P'_{i,j}) = \left(\sum_{c \in \{r, g, b\}} (P_{i,j,c} - P'_{i,j,c})^n \right)^{\frac{1}{n}} \quad (3)$$

This correlation function can be seen as weight average color differences (WACD) and it is given by the following equation:

$$\text{DCMP}(P_{i,j,k}, P'_{i,j,k}) = \Delta^2(P_{i,j,k}, P'_{i,j,k}) * \sum (D * W) \quad (4)$$



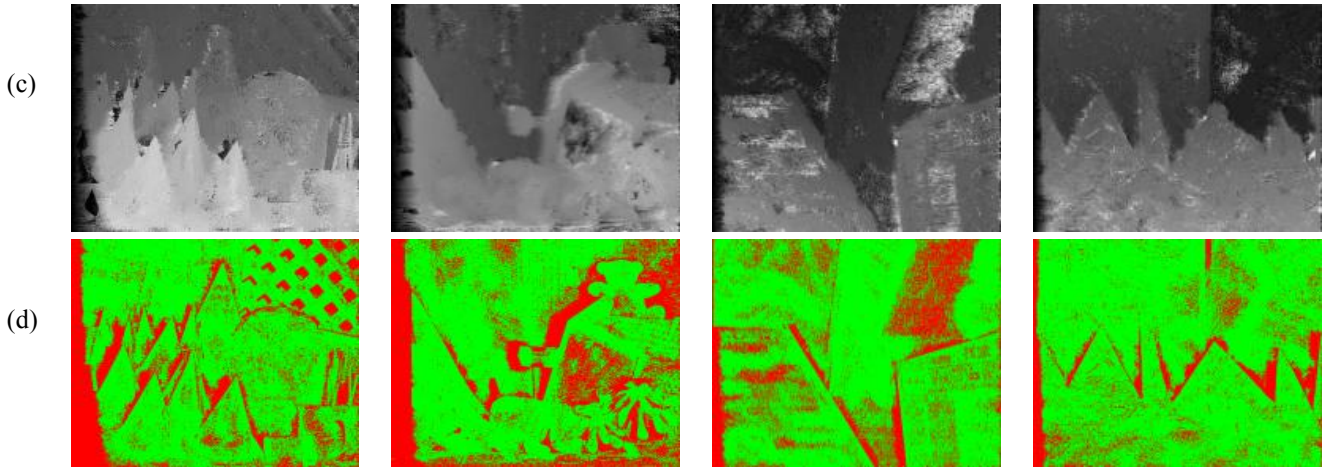


Figure 10. Data sets and output images. (a) Left-hand images used for evaluation (b) ground truth disparity (c) Disparity maps obtained applying only dissimilarity function (d) Well- and badly-matched pixels: Green color for well-matched pixels and red color for badly-matched pixels.

5.3. LOCAL CONFIDENCE MEASURE ESTIMATION THEORY

• RELATED WORKS

In dense stereovision, several well-known stereo algorithms compute an initial disparity map from a pair of images under a known camera configuration. These algorithms are based loosely on local methods, such as window correlation, which take into account only neighborhood points of the pixel to be matched. The disparity map obtained has a lot of noise and erroneous values. This noise concerns mostly the pixels belonging to occluded or textureless image regions. An iterative process is applied then to the initial disparity map in order to improve it. These methods use global primitives. Cost-relaxation approaches, which were invented by Marr and Poggio [18] and which are picked up again by Brockers [19], belong to this family. Some research used a graph-based method [21] and color segmentation based stereo methods [20] which belong to what is called “global approaches”. Other approaches were proposed: they are based on a probabilistic framework optimization, such as expectation-maximization [23] and belief propagation [22, 27]. These methods aim to obtain high-quality and accurate results, but are very expensive in terms of processing time. It is a real challenge to evaluate stereo methods in the case of noise, depth discontinuity, occlusions and non-textured image regions.

Besides, some mixture of local methods is first used to obtain an estimated disparity map which is improved in a second step by global methods. On the basis of local methods, a pixel on the left-hand image is evaluated with candidate pixels on the corresponding right-hand image. Some research carried out by Yoon [24] applies window-based methods in the Lab color space and are coupled with an adaptive window [26] which tries to find an optimal support window for each pixel. These techniques assume that the neighborhood of a pixel to be

matched presents homogeneity in terms of disparity values. In other words, all the pixels in the given correlation window must have very similar disparities.

• DESCRIPTION OF THE PROPOSED APPROACH

In our approach, in order to reduce the processing time and to deal with the problems of ambiguity in the matching process, the correlation function used to evaluate each stereo pair will only be applied on high color variation regions of the images.

The main idea of our approach is to compute a confidence measure for every matched pixel. Furthermore, confidence measure ψ can be seen as a matching probability for pixel P_l (a pixel P in the left-hand image) with pixel P_r (a pixel P in the right-hand image), given some parameters. The way in which confidence measures are calculated is provided by equation 5.

$$\psi(P_l^{i,j}, P_r^{i,v}) = P(P_l^{i,j} / P_r^{i,v}, N, min, \sigma, \omega) \quad (5)$$

The confidence measure with its parameters is given by equation 6:

$$\psi(P_l^{i,j}, P_r^{i,v}) = \left(1 - \frac{min}{\omega}\right)^{N^2 \cdot \log(\sigma)} \quad (6)$$

where:

- *min* : The Best Correlation Score

The output of the dissimilarity function is a measure representing the degree of similarity between two pixels. Then, the candidate pixels are ranked in increasing order according to their corresponding scores. The couple of pixels that has the minimum score is considered as the best-matched pixels. The lower the score, the better the matching. The nearer the minimum score to zero, the

greater the chance of candidate pixel being the right correspondent.

- N : Number of Potential Candidate Pixels

This parameter represents the number of potential candidate pixels having similar scores. N has a big influence because it reflects the behavior of the dissimilarity function. When the value of N is quite large, that means the first potential candidate pixel is located in a uniform color region of the frame.

The lower the value of N , the fewer the potential candidate pixels. In the case where there are a few candidates, the chosen candidate pixel has a greater chance of being the right correspondent. Indeed, the pixel to be matched belongs to a region with high variation of color components.

While establishing the relationship between N and \min values, with a very small value of N and a minimum score \min , near to zero for instance, the pixel to be matched probably belongs to a region of high color variation.

- σ : Disparity variation of N pixels

A disparity value is obtained for each candidate pixel. For the N potential candidate pixels, we compute standard deviation σ on the N disparity values. A small σ means that the N considered pixels are neighbors. In this case the true candidate pixel should belong to a particular region of the frame, such as edge, transition point. Therefore, it increases the confidence measure. A large σ means that the N candidate pixels taken into account are situated in a uniform color region.

- ω : gap value

This parameter represents the difference between the N^{th} and $(N + 1)^{\text{th}}$ scores given with the dissimilarity function used. It is introduced to adjust the impact of the \min score.

To ensure that this function gives a value between 0 and 1, some constraints are introduced. The \min parameter must not be higher than the ω one. If so, parameter ω is forced to $\min + 1$. However, the $\log(\sigma)$ term is used instead of σ alone. It has a big influence on the confidence measure in case of high values of σ , and it is indifferent otherwise. This leads to reducing the impact of high values of σ and to obtaining coherent confidence measures.

The number N of potential candidate pixels is deduced from the k scores obtained with the dissimilarity function previously presented. The main idea is to detect major differences between successive scores. These differences are called main gaps. Let f denote a function which represents all scores given by the dissimilarity function in increasing order. Then, we apply the average rate growth to the f function. This second function can be denoted by η and can be seen as the ratio of the difference between a given score and the first score, and the difference between their ranks. This function is defined in equation 7.

$$\eta(x_m) = \frac{f_{x_m}^{i,j} - f_{x_1}^{i,j}}{x_m - x_1} \quad m = 1 \dots k \quad (7)$$

where $f_{x_m}^{i,j}$ is the m^{th} of k score of the (i, j) coordinate pixel and x_m is the rank of the corresponding score.

$$\xi(x_m) = \frac{\eta_{x_m}^{i,j} - \eta_{x_{m-1}}^{i,j}}{m^2} \quad m = 1 \dots k \quad (8)$$

The previous function (Figure 8) is used to characterize jump scores and is applied only in the case where $(\eta_{x_m}^{i,j} - \eta_{x_{m-1}}^{i,j})$ is a positive value. We have introduced parameter m^2 in order to penalize candidate pixels according to their rank. The number of potential candidate pixels is given by formula 9.

$$N = \underset{m}{\text{Argmax}} \xi(x_m) \quad (9)$$

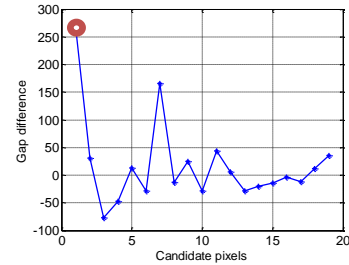


Figure 11. The number N of potential candidate pixels is the rank of global maximum of ξ function.

6. EXPERIMENTAL RESULTS

In this section, we describe the evaluation of the performances of the proposed approach thanks to images with ground truth. Well-known conventional stereo images with available ground truths are employed to test the relevance of the accuracy of our algorithm [25]. In this evaluation, four pairs of stereo images are used: Cones, Teddy, Venus and Sawtooth. As a first step, the disparity map is initialized by applying the dissimilarity function proposed in [15]. This provides a first visual rendering of the disparity map for Cones (Figure 12).

The dissimilarity function used has the particularity of matching well most pixels in regions having high variation of color components. As shown in Figure 12.d, Pixels belonging to uniform color regions or depth discontinuity regions are badly-matched. This can be explained by the presence of several potential candidate pixels for the given pixel to be matched. Therefore, this increases the error of the matching task.

As a second step, and in order to quantify and to automatically identify well-matched pixels from badly-matched pixels, we have exploited the confidence measure method described in section 4. Each couple of matched pixels is evaluated and belongs to one of the

following three categories for a given confidence measure threshold T :

- Well-matched pixels having a confidence measure higher than T ,
- Badly-matched pixels having a confidence measure higher than T and an erroneous disparity value according to the ground truth image,
- Unclassified pixels having a confidence measure lower than T .

The disparity map obtained with an optimal window size (11x11) defined experimentally is taken into consideration. All matched pixels are classified into the three previous categories. Global results are shown in Figure 14. According to Figure 12.d, the badly-matched

pixels represent either occluded pixels or pixels belonging to uniform regions in terms of color. Firstly, in order to reduce the impact of error matching, only pixels belonging to regions of high color variation are considered. In order to reduce the impact of error matching and to have a high gap between the rate of well-matched and badly-matched pixels, only pixels belonging to regions of high color variation are considered. A pixel will be matched only if the sum of all color component variations of neighborhoods pixels belonging to a support window is higher than a threshold. This threshold is computed based on dynamics of color in the image.

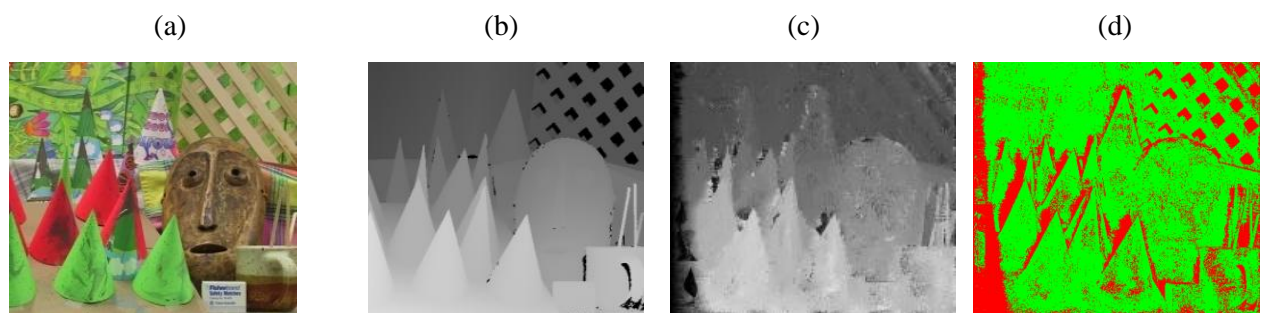


Figure 12. Left-hand “Cones” image and output images. (a) Left-hand images used for evaluation (Cones) (b) ground truth disparity (c) Disparity map obtained with WACD dissimilarity function application only (d) Well- and badly-matched pixels: Green color for well-matched pixels and red color for badly-matched pixels.

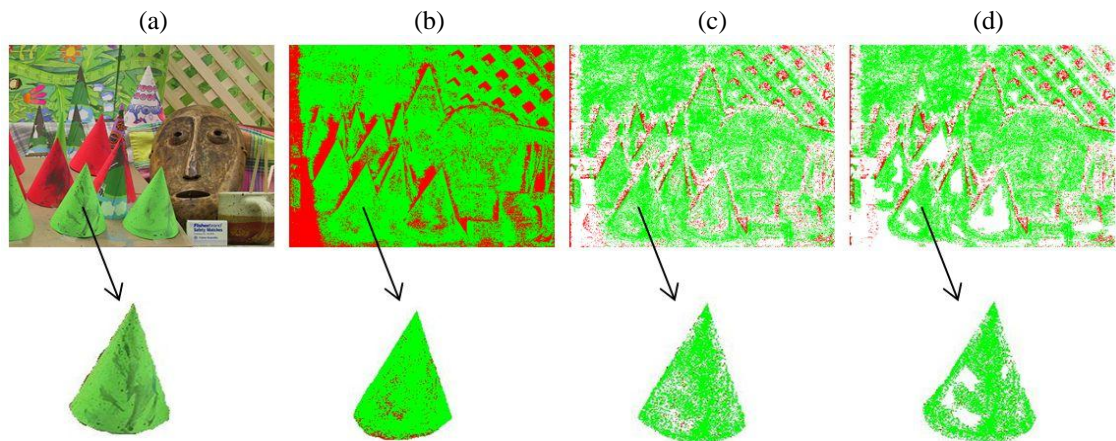


Figure 13. (a) Left-hand Cones image (b) Well-matched and badly-matched pixels with WACD dissimilarity function application only (c) Well-matched in green, badly-matched in red and unclassified pixels in white for a confidence measure of 96% (d) After elimination of pixels belonging to regions having a low color variation.

In figure 14, four diagrams are used to illustrate the matching performances using the WACD correlation function for the lower curves, and using the same function with introduction of confidence measure for the upper curves. We can notice that the use of confidence measure applied to matched pixels allows us to improve the matching rates significantly. Thus, for a correlation window of 11x11 pixels and a 96% confidence measure, the good matching rate passes from:

- 73% to 88% for cones,
- 67% to 86% for Teddy,
- 74% to 92% for Venus,
- 83% to 93% for Sawtooth.

The rate of well-matched pixels is initially evaluated by comparison of the estimated disparity and the corresponding disparity values of ground truth images. This concern all matched pixels of the image. The introduction of the confidence measure principle leads to

decrease the rate of badly-matched pixel that will be considered as unclassified pixels. Moreover, this increases consistently the rate of well-matched pixels for a given region. This rate is calculated based only on well- and badly-matched pixels. A significantly rate of well-matched pixels in a given region (half of pixels minimum) not allows the re-estimation of the remaining pixels which are labeled as unclassified pixels. Otherwise, the modal class disparity of well-matched pixels could be taken for the disparity re-estimation.

It is to be noted that the good matching rates, with the use of the confidence measure, concerns only pixels belonging to high color variation zones (around 80% on average for the four pairs of images tested). The other pixels, those belonging to uniform zones in terms of color, are assessed later. In Figure 13, the different steps of our approach are illustrated. We have extracted a single cone to illustrate visually the improvement for each step

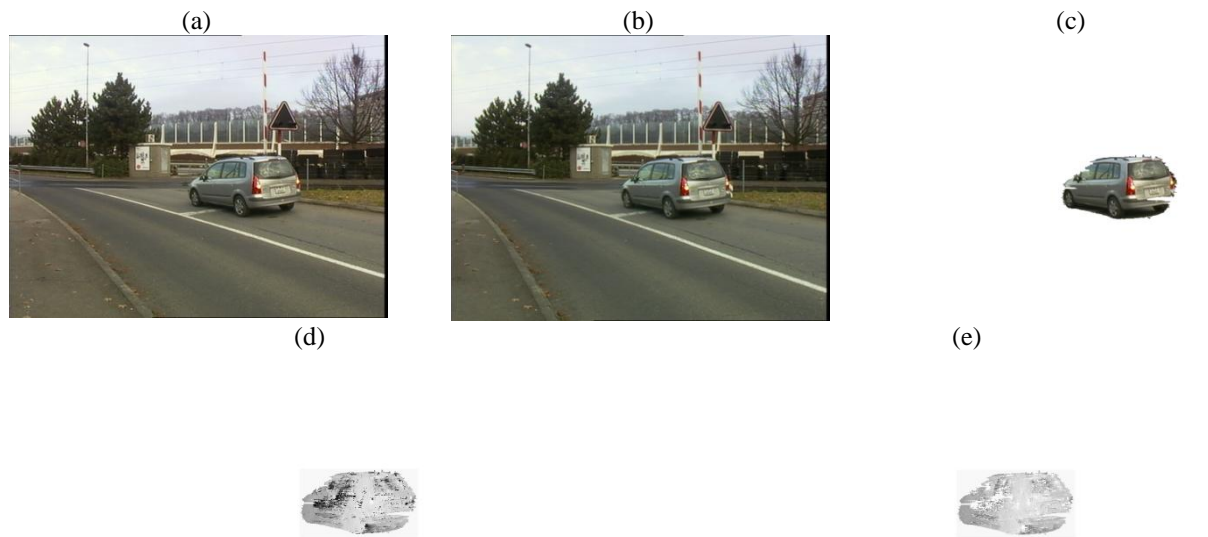


Figure 15. (a) Left-hand image (b) Right-hand image, (c) The extracted moving object, (d) Disparity map obtained applying the proposed (WACD) matching function, (e) Improved disparity map introducing confidence measure theory.

In figure 13.b, for the extracted cone, we can notice that the number of badly-matched pixels (in red) is quite high. These pixels are located in uniform regions of the cone. The application of the confidence measure (passage from 13.b to 13.c) leads to an important decrease in the number of badly-matched pixels. In fact, in figure 13.c we can notice that the number of red pixels has reduced. The white pixels in figure 13.c are the ones with a confidence measure lower than the given threshold (96%). In figure 13.d the unclassified pixels belonging to homogeneous regions in terms of color are identified and marked in white. Their disparities values will be assessed later. Finally, in our method, we have only considered pixels having a high confidence measure. However, in order to update the disparity for

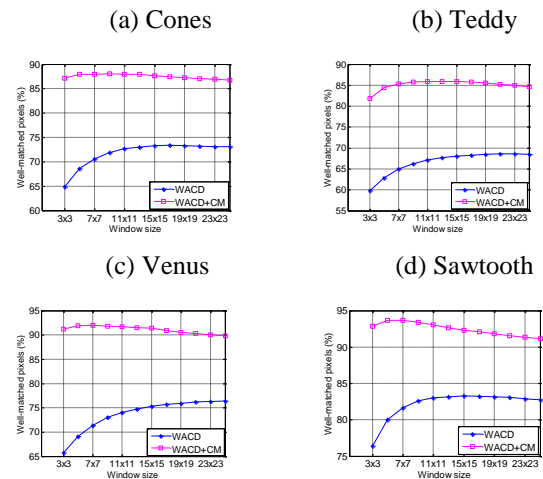


Figure 14. Well-matched pixels rate using only WACD correlation function (blue curves) and introducing both confidence measure approach and homogeneous color regions elimination (pink color).

all the pixels, several methods could be applied. On one hand, most research picked up from the literature uses a color segmentation based method [28]. It considers that pixels belonging to a homogeneous component color also have similar disparity. Statistical methods could also be applied to estimate the disparity values of unclassified pixels: modal disparity class, statistical clustering.

On the other hand, in the foreground extraction process, motion could be used to re-estimate the remaining non-matched pixels. For a given moving region, it is reasonable to consider that the disparity values will stay homogeneous. Then, the 3-D localization is carried out on only moving regions. A

way to re-estimate these remaining disparities is to consider for instance, the modal class disparity of well-matched pixels.

7. CONCLUSION AND PERSPECTIVE

A background subtraction technique for detection of moving objects was applied on natural images sequences. The used non-parametric technique entitled Codebook proves its effectiveness in case of small illumination variation in the scene but fails with a fast illumination background changes. To overcome this problem and to meet the requirements related to level crossings safety, we are developing a novel approach that is based on Independent Component Analyses (ICA) theory. The latter is a statistical and computational technique for revealing hidden factors that underlie sets of random variables, measurements, or signals. It aims to find the source signals from observations data. The observations matrix can be formed by both a random background image and other containing arbitrary foreground objects. The estimated sources signals correspond to the modeled background signal and foreground, which distinctly preserves the foreground object without the detailed contents of the reference background [34].

A novel method was also proposed in section 5 of this paper for 3-D object localization. This aims to obtain an accurate disparity map for 3D localization of objects in level crossing. We have applied a local matching method to initialize the disparity value for each pixel. We have then introduced new local parameter in order to compute a confidence measure for each matched pixel. The main contribution of our approach is twofold. On the one hand, new parameters introduced above can obtain important information for occluded and uniform region detection. On the other hand, unclassified pixels disparities can be updated in a post-processing step in order to obtain a more accurate disparity map. This depends on the disparity measure of each well-matched pixel classified using the confidence measure approach. The choice of the confidence threshold is very important for the success of 3-D estimation process of each moving object. This choice depends on the rate of well-matched pixels among all matched-pixels in a given region. Such a confidence threshold is considered as optimal when it satisfies the following function $T = \operatorname{argmax}_T W_i$ in which W_i is the rate of well-matched pixels in the region i . In order to re-estimate the values of disparities of the pixels whose confidence measure is lower than a given threshold (those belonging to uniform regions), the values of disparities in the surrounding pixels can be used: segmentation, disparity modal class, clustering, etc. Another way is to consider the motion information to re-estimate unclassified pixels.

Our approach is highly dependent on the dissimilarity function used for computing the score of all matched pixels. However, a more extensive study will be carried out in order to enhance the dissimilarity function used. The rate of well-matched pixels will, therefore, be

improved. The results are encouraging in terms of processing time which is compatible for a real-time implementation. The promising result obtained allows us to follow this track.

An extensive evaluation of our algorithms is carried out considering real word situations. Datasets coming from a series of level crossings in Lausanne (Switzerland) are used. Many real scenarios of cars, pedestrians, objects, crossing several different LC are included in the datasets. This will allow us to deeply evaluate the accuracy of our obstacle detection system in terms of objects extraction and then 3D localization. Of course, our detection system will be coupled with a communication one in order to timely warn the car drivers, the approaching train and maybe a control room on the existence of an obstacle in the dangerous zone of the LC. This is under study in the framework of the national work programme entitled PANsafer: towards a safer level crossing.

8. REFERENCES

- [1] ERA (European Railway Agency), "A Summary of 2004-2005 EU Statistics on Railway Safety", Source of data: Eurostat.
- [2] Project SELCAT (Safer European Level Crossing Appraisal and Technology), *A Co-ordination Action of the European Commission's 6th Framework Programme*.
- [3] Government's IT Strategy Headquarters: "New IT Reform Strategy", *The realization of society, in which everyone can benefit from IT, anytime and anywhere*. (In Japanese), January 19, 2006.
- [4] Y. Hisamitsu, K. Sekimoto, K. Nagata, "3-D Laser Radar Level Crossing Obstacle Detection System", *IHI Engineering Review*, vol. 41, no. 2, August 2008.
- [5] G.L. Foresti, "A Real-Time System for Video Surveillance of Unattended Outdoor Environments", *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 8, no. 6, October 1998.
- [6] M. Ohta, "Level Crossings Obstacle Detection System Using Stereo Cameras", *QR of RTRI*, vol. 46, no. 2, June 2005.
- [7] C. Stauffer, W. Grimson, "Adaptive Background Mixture Models for Real-Time Tracking", *IEEE International Conference on Computer Vision and Pattern Recognition*, 2:pp. 246-52, 1999
- [8] D.S. Lee, J.J. Hull, B. Erol. "A Bayesian Framework for Gaussian Mixture Background Modeling", *IEEE International Conference on Image Processing*, 2003.
- [9] M. Harville, "A Framework for High-Level Feedback to Adaptive, Perpixel, Mixture-of-Gaussian Background Models", *European Conference on Computer Vision*, 3:pp. 543-60, 2002.
- [10] O. Javed, K. Shafique, M. Shah, "A Hierarchical Approach to Robust Background Subtraction

- Using Color and Gradient Information”, *IEEE Workshop on Motion and Video Computing (MOTION’02)*; 2002.
- [11] M. Cristani, M. Bicego, V. Murino. ”Integrated Region- and Pixel Based Approach to Background Modeling”, *Proceedings of IEEE Workshop on Motion and Video Computing*, 2002.
- [12] K. Toyama, J. Krumm, B. Brumitt, B. Meyers,” Wallflower: Principles and Practice of Background Maintenance”, *International Conference on Computer Vision*, pp. 255–61. 1999.
- [13] K. Kim, T. H. Chalidabhongse, D. Harwood and L. Davis, ”Real-time Foreground-Background Segmentation Using Codebook Model”, *Real-Time Imaging, Special Issue on Video Object Processing*, vol. 11, Issue 3, pp. 172-185, June 2005.
- [14] M. Okutomi , T. Kanade, “A Multiple-Baseline Stereo”, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 15, no.4, pp.353-363, 1993.
- [15] N. Fakhfakh, L. Khoudour, M. El-Koursi, ”Mise en Correspondance Stéréoscopique d’Images Couleur pour la Détection d’Objets Obstruant la Voie aux Passages à Niveau”, in *TELECOM’09 & 6^{ème} JFMMA*, p. 206 (4 pages), Agadir, Maroc, 2009.
- [16] J. Martin, J-L. Krowley, ”Experimental Comparison of Correlation Techniques”, *International Conference on Intelligent Autonomous Systems (IAS4)*, 2002.
- [17] A. Crouzil, "Perception du Relief et du Mouvement par Analyse d’une Séquence Stéréoscopique d’Images.", *Thèse de doctorat, Université Paul Sabatier, UPS, Toulouse*, pp.58-60, septembre 1997.
- [18] D. Marr, T. Poggio, ”Cooperative Computation of Stereo Disparity”, in *American Association for the Advancement of Science*, vol. 194, Issue 4262, pp. 283—287, 1976.
- [19] R. Brockers, M. Hund, B. Mertsching, ”Stereo Vision Using Cost-Relaxation with 3D Support Regions”, In: *ICVNZ, New Zealand*, 2005.
- [20] Y. Taguchi, B. Wilburn, C. L. Zitnick, ” Stereo Reconstruction with Mixed Pixels Using Adaptive Over-Segmentation”, In: *CVPR*, pp. 1--8, Anchorage, Alaska, 2008.
- [21] P. Foggia, J.M. Jolion, A. Limongiello, M. Vento, ”Stereo Vision for Obstacle Detection: A Graph-Based Approach”, *Lecture Notes in Computer Science, Springer Berlin / Heidelberg*, pp. 37—48, 2007.
- [22] C. Lee, Y. Ho, ”Disparity Estimation Using Belief Propagation for View Interpolation.”, In: *ITC-CSCC*, pp. 21--24, Japan, 2008.
- [23] W. Xiong, H.S. Chung, J. Jia, ” Fractional Stereo Matching Using Expectation-Maximization.”, In: *IEEE TPAMI*, vol. 31, issue 3, pp. 428—443, 2008.
- [24] K.J. Yoon, S. Kweon, ”Adaptative Support-Weight Approach for Correspondence Search.”, In: *IEEE TPAMI*, vol. 28, No. 4, 2006.
- [25] D. Scharstein, R. Szeliski, “ Middlebury stereo vision research page, <http://vision.middlebury.edu/stereo/>
- [26] O. Veksler, “ Fast Variable Window for Stereo Correspondence Using Integral Image.”, In: *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, vol. 1, pp. 556--561, Madison, Wisconsin, 2003.
- [27] J. Sun, N. -N. Zheng, H.Y. Shum,” Stereo Matching Using Belief Propagation.”, In: *IEEE TPAMI*, vol. 25, No. 7 (2003).
- [28] A. Klaus, M. Sormann, K. Karner, “ Segment-Based Stereo Matching Using Belief Propagation and a Self-Adapting Dissimilarity Measure.”, In: *ICPR*, pp. 15—18, 2006.
- [29] O. Schreer, I. Feldmann, U. Goelz, and P. Kauff, “Fast and Robust Shadow Detection in Videoconference Application”, in *Proc. Of VIPromCom 2002, 4th EURASIP IEEE International Symposium on Video Processing and Multimedia Communications*, Zadar, Croatia), June, 2002.
- [30] C. Jian and M. O. Ward, “Shadow Segmentation and Classification in a Constrained Environment”, *CVGIP: Image Understanding*, vol. 59, no. 2, pp.231-225, 1994.
- [31] A. Cavallaro, E. Salvador and T. Ebrahimi, ” Detecting Shadows In Image Sequences”, in *Proc. First European Conference on Visual Media Production*, London (UK), pp.15-16, 2004.
- [32] E. Salvador, A. Cavallaro and T. Ebrahimi, ” Spatio-Temporal Shadow Segmentation and Tracking”, *Proceedings of SPIE's Image and Video Communications and Processing*, vol. 5022, p. 389-400, SPIE, 2003.
- [33] F. Porikli, O. Tuzel, ”Human Body Tracking by Adaptive Background Models and Mean-Shift”, *IEEE International Workshop on Performance Evaluation of Tracking and Surveillance*, PETS-ICVS, 2003.
- [34] D-M. Tsai, S-C. Lai, ” Independent Component Analysis-Based Background Subtraction for Indoor Surveillance”, *IEEE Transactions on Image Processing*, Vol. 18, No. 1, January 2009.