

## Article

# Improving Video Segmentation by Fusing Depth Cues and the ViBe Algorithm

Xiaoqin Zhou <sup>1,2,3,4</sup>, Xiaofeng Liu <sup>2,3,4\*</sup>, Aimin Jiang <sup>2,3,4</sup>, Bin Yan <sup>5</sup> and Chenguang Yang <sup>6</sup>

<sup>1</sup> College of Computer and Information Engineering, Hohai University, Nanjing, China; e-mail: zhouxq@hhu.edu.cn.

<sup>2</sup> College of Internet of Things Engineering, Hohai University, Changzhou, China

<sup>3</sup> Changzhou Key Laboratory of Robotics and Intelligent Technology, Changzhou, China

<sup>4</sup> Jiangsu Key Laboratory of Special Robots (Hohai University), Changzhou 213022, China; e-mail: jiangam@hhuc.edu.cn

<sup>5</sup> College of Electronics, Communication and Physics, Shandong University of Science and Technology, Qingdao, China; e-mail: yanbinhit@hotmail.com.

<sup>6</sup> Zienkiewicz Centre for Computational Engineering, Swansea University, Swansea SA1 8EN, UK; e-mail: cyang@theiet.org

\* Correspondence: xfliu@hhu.edu.cn

**Abstract:** Depth-sensing technology has led to broad applications of inexpensive depth cameras that can capture human motion and scenes in 3D space. Background subtraction algorithms can be improved by fusing color and depth cues, thereby allowing many issues encountered in classical color segmentation to be solved. In this paper, we propose a new fusion method that combines depth and color information for foreground segmentation based on an advanced color-based algorithm. First, a background model and a depth model are developed. Then, based on these models, we propose a new updating strategy that can eliminate ghosting and black shadows almost completely. Extensive experiments have been performed to compare the proposed algorithm with other, conventional RGB-D algorithms. The experimental results suggest that our method extracts foregrounds with higher effectiveness and efficiency.

**Keywords:** object detection; background subtraction; video surveillance; Kinect sensor fusion

## 1. Introduction

In recent years, enormous amounts of data on human behavior have been collected using 2D and 3D cameras, and automated methods for detecting and tracking individuals have begun to play an important role in studies of experimental biology, behavioral science, and related disciplines. The extraction of moving objects from a video sequence is the first step in such video analysis systems. At present, a variety of motion detection methods have been proposed, such as the frame difference method [1], the background subtraction method [2], the optical flow method [3] and the block matching method [4]. The core of a background subtraction algorithm is the modeling of the background. Zones that show notable differences between the current frame and the background model are deemed to correspond to moving objects. Generally, background subtraction algorithms include the average background model (AVG) algorithm, the GMM algorithm [5], the Codebook algorithm [6] and the ViBe algorithm [7–9]. The ViBe algorithm is a fast motion detection algorithm proposed by Olivier Barnich et al. [7]. It is characterized by a high processing efficiency and a good detection effect.

Most of the conventional methods mentioned above were designed for application to color images. However, depth is another interesting cue for segmentation that is less strongly affected by the adverse effects encountered in classical color segmentation, such as shadow and highlight regions. Depth cameras, such as the Microsoft Kinect and the ASUS Xtion Pro, are able to record real-time depth video together with color video. Because of their beneficial depth imaging features and moderate price, such depth cameras are broadly applied in intelligent surveillance, medical

diagnostics, and human-computer interaction applications [10–12]. The Kinect sensor is not sensitive to light conditions; it works well either in a bright room or in a pitch black one. This makes depth images more reliable and easier for a computer program to understand.

Most studies using the Kinect sensor have focused on human body detection and tracking [13–15]. The Histogram of Oriented Depths (HOD) detection algorithm, proposed in [13], can be used to match human body contour information in an image. In [14], a model was presented for detecting humans using a 2-D head contour model and a 3-D head surface model. In these studies, the computational complexity of the feature generation and matching process was relatively high.

Crabb et al. [16] and Schiller et al. [17] focused on combining the depth and color information obtained by low-resolution TOF cameras, but their methods are not well suited for video surveillance. For example, the method of Crabb et al. [16] requires the definition of a distance plane where no foreground object is located behind any part of the background. Fernandez-Sanchez et al. [18] proposed an adaptation of the Codebook background subtraction algorithm that focuses on different sensor channels. In these methods, no emphasis is placed on eliminating ghosting. A ghost is a set of interconnected points that is detected as a moving object but does not correspond to any real object (see Figure 1(c)). Ghosting greatly reduces the effectiveness of motion detection.



**Figure 1.** Ghosting in a foreground segmentation map generated by the ViBe algorithm [7]: (a) color frame, (b) ground truth, (c) foreground extraction result.

In this paper, we propose an adaptive ViBe background subtraction algorithm that fuses the depth and color information obtained by the Kinect sensor to segment foreground regions. First, a background model and a depth model are established. Then, based on these models, we develop a new updating strategy that can efficiently eliminate ghosting and black shadows. The improved algorithm is evaluated using an RGB-D benchmark dataset [19] and achieves good results that provide a perfect basis for subsequent feature extraction and behavior recognition.

The remainder of the paper is organized as follows. In Section 2, we briefly describe the original ViBe algorithm. Then, the improved algorithm is developed in Section 3. In Section 4, experimental results and discussions are presented. Finally, we conclude the paper in Section 5.

## 2. ViBe background subtraction algorithm

In this section, we first review the basic ViBe algorithm. Then, we identify its disadvantages. This technique involves modeling the background based on a set of samples for each pixel. New frames are compared with the background model, pixel by pixel, to determine whether each pixel belongs to the background or the foreground.

### 2.1. Pixel model

Background model construction begins from the first frame. Formally, let  $v(x)$  denote the value in a given Euclidean color space associated with the pixel located at  $x$  in the image, and let  $v_i$  be the background sample value with index  $i$ . Each background pixel  $x$  is modeled based on a collection of  $N$  background sample values  $M(x) = \{v_1, v_2, \dots, v_N\}$ .

## 2.2. Classification process

If the Euclidean distance from a sample  $v_i$  in  $M(x)$  to  $v(x)$  is below a threshold  $R$ , then  $v_i$  is regarded as a neighbor of  $v(x)$ . We define the number of neighbors of the pixel located at  $x$  as  $N_R(x) = \{\|v(x) - v_i\| < R, \forall v_i \in M(x)\}$ . When  $N_R(x)$  is greater than a threshold  $\lambda$ ,  $x$  is a background pixel. Otherwise, it is a foreground pixel.

## 2.3. Updating the background model over time

It is necessary to continuously update the background model with each new frame. This is a crucial step for achieving accurate results over time. When a pixel  $x$  is classified as background, the background model updating strategy is triggered. A sample is chosen randomly. Mathematically, the probability that a sample present in the model at time  $t$  will be preserved is given by  $(N - 1)/N$ . Under the assumption of time continuity, for any later time  $t + dt$ , this probability is equal to

$$P(t, t + dt) = \left(\frac{N - 1}{N}\right)^{(t+dt)-t} \quad (1)$$

which can be rewritten as

$$P(t, t + dt) = e^{-\ln(\frac{N}{N-1})dt} \quad (2)$$

This expression shows that the expected remaining lifespan of any sample value in the model decays exponentially according to a random subsampling strategy. As in [7], we adopt a time subsampling factor of  $\phi$ , meaning that a background pixel value has one chance in  $\phi$  of being selected to update its pixel model.

Reduced pseudo-code for the ViBe construction phase is given in Algorithm 1.

---

### Algorithm 1 Algorithm for ViBe construction

---

```

1: procedure VIBE(image,  $N$ ,  $R$ ,  $\lambda$ ,  $\phi$ )
2:   for each pixel do
3:     while matches <  $\lambda$  and index <  $N$  do
4:       Calculate Euclidean distance between  $v_x$  and  $v_i$ 
5:       if dist <  $R$  then
6:         matches  $\leftarrow$  matches + 1
7:       end if
8:       index  $\leftarrow$  index + 1
9:     end while
10:    if matches  $\geq \lambda$  then
11:      Store that pixel  $\in$  background
12:      Update current pixel background model with probability  $1/\phi$ 
13:      Update neighboring pixel background model with probability  $1/\phi$ 
14:    else
15:      Store that pixel  $\in$  foreground
16:    end if
17:  end for
18: end procedure

```

---

The classical ViBe algorithm has the advantages of simple processing and outstanding performance. Its main drawback is the occurrence of ghosting. A moving object in the first frame often causes ghosting. To resolve this problem, we can take advantage of depth information. The Kinect sensor records the distance to any object that is placed in front of it. This feature can be utilized to determine whether a foreground pixel is a ghost.

### 3. Fusion: depth-extended ViBe (DEVB)

#### 3.1. Depth model

To eliminate ghosting, the ViBe algorithm is improved by enhancing the matching conditions when a pixel is classified as foreground. A depth model  $M_D(x)$  is added. Initially, the pixel values of the first depth frame are saved to  $M_D(x)$ . This depth model also has an updating strategy similar to that for the background model. When the updating strategy is triggered, the depth value  $M_D(x)$  is replaced with that corresponding to the current pixel. If the following condition is satisfied, this pixel will be considered a ghost pixel:

$$v(x, t_0) - v(x, t) > \tau \quad (3)$$

where  $v(x, t_0)$  is the value of the pixel located at  $x$  in the depth model  $M_D(x)$  at time  $t_0$ ,  $v(x, t)$  is the value of the pixel located at  $x$  in the current image at time  $t$ , and  $\tau$  is the threshold on the distance between  $v(x, t_0)$  and  $v(x, t)$ . The parameter  $\tau$  can be tuned to obtain the best result. The value of  $\tau$  is recommended to lie within the range of  $[1, 3]$ .

#### 3.2. Fusion algorithm for color and depth images

A color image conforms to an individual's visual habits and provides detailed information such as color and texture. The most intuitive fusion strategy is to add the depth information as a fourth channel to the ViBe algorithm for color images. The channel combination  $f = (f_r, f_g, f_b, d)$  is formed from the three color channels in RGB space and the depth value  $d$ . In this way, we intend to utilize depth information as a measure of reliability during segmentation. A depth image quantifies the distance from an object to the camera. The higher the value of a pixel is, the more reliable its measurement is in the depth image. Thus, we use the inverse of the depth image to allow it to be used in the same way as a variance image.

The inverted depth values are normalized between zero and one, and the normalized uncertainty is denoted by  $\sigma(x)$ . In a region where the depth uncertainty is high, the depth measurement is considered unreliable. For example, where there are holes in the depth image, the fusion result will depend on the color image. We weight the normalized depth  $\hat{d}(x)$  with the uncertainty  $\sigma(x)$ , resulting in  $w_d$  values that range between zero and one depending on  $\sigma(x)$ . Meanwhile, the color value  $c(x)$  is multiplied by  $\sigma(x)$  and added to itself to obtain the weighted color value, defined as  $w_c$ . The combined normalized image is denoted by  $I(\hat{x})$ , as shown in Equation (4).

$$\begin{cases} w_d = (1 - \sigma(x))\hat{d}(x) \\ w_c = (1 + \sigma(x))c(\hat{x}) \\ I(\hat{x}) = \frac{1}{2}(w_d + w_c) \end{cases} \quad (4)$$

The original algorithm can eliminate ghosts and black shadows in subsequent frames, but the process is relatively slow. We propose a fusion strategy that incorporates an additional depth model to effectively remove ghosting and black shadows. Reduced pseudo-code for the DEVB construction phase is given in Algorithm 2.

Section 4 presents experiments performed without post-processing and the results obtained using three RGB-D algorithms (MOG4D [17], DECB [18] and DEVB) as well as the color-based ViBe algorithm.

### 4. Experiments and results

The program development environment consisted of VC++2010, OpenCV SDK2.4.3 and OpenNI1.5.2.7. The PC was equipped with a Core Duo 2 CPU E7500 with 2.00 GB of RAM. The video

**Algorithm 2** Algorithm for DEVB construction

---

```

1: procedure DEVB( $image_{fusion}, image_{depth}, N, R, \lambda, \phi$ )
2:   for each  $image_{fusion}$  pixel do
3:     while  $matches < \lambda$  and  $index < N$  do
4:       Calculate Euclidean distance between  $v_x$  and  $v_i$ 
5:       if  $dist < R$  then
6:          $matches \leftarrow matches + 1$ 
7:       end if
8:        $index \leftarrow index + 1$ 
9:     end while
10:    if  $matches \geq \lambda$  then
11:      Store current pixel  $\in$  background
12:      Update current pixel background model in  $image_{fusion}$  with probability  $1/\phi$ 
13:      Update neighboring pixel background model in  $image_{fusion}$  with probability  $1/\phi$ 
14:      Update current pixel  $model_{depth}$  using  $image_{depth}$ 
15:    else
16:      Store current pixel  $\in$  foreground
17:      Find ghosts, calculate Euclidean distances between pixels at the same position in
         $image_{depth}$  and  $model_{depth}$ 
18:      if  $dist > \tau$  then
19:        Store that pixel  $\in$  background
20:        Update current pixel background model in  $image_{fusion}$  with probability  $1/\phi$ 
21:        Update neighboring pixel background model in  $image_{fusion}$  with probability  $1/\phi$ 
22:      end if
23:    end if
24:  end for
25: end procedure

```

---

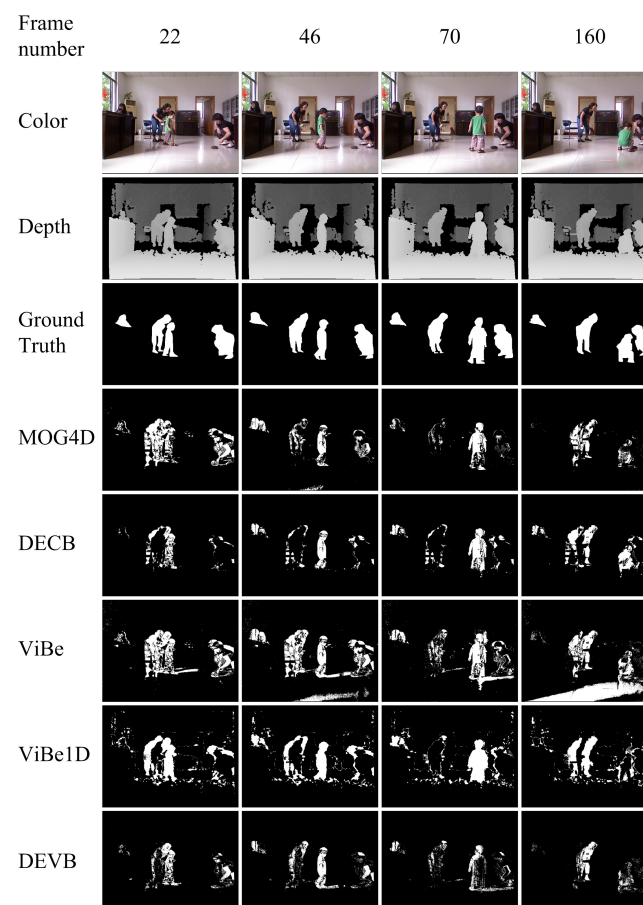
frame rate was 30 fps, and the size of both the color and depth images was  $640 \times 480$ . We compared the results of our method with those obtained using two state-of-the-art RGB-D fusion-based background subtraction algorithms, namely, MOG4D [17] and DECB [18], as well as the ViBe algorithm on the color images [7] (ViBe) and the ViBe algorithm on the depth images (ViBe1D). To evaluate these algorithms objectively through a quantitative analysis, we required a benchmark that would provide information on both color and depth images. The chosen benchmark sequences are publicly available at [19]. However, the depth image sequences provided by this source could not be utilized directly. The depth images are in the 16-bit png format, with the first 3 bits swapped with the last. We needed to swap them back after reading each image to obtain values for each pixel corresponding to the distance from the Kinect sensor to the object in mm. The sequences *child\_no1*, *new\_ex\_occ4*, *walking\_occ1* and *new\_ex\_no\_occ* from [19] were chosen for testing. Each sequence presents important problems that are typically challenging for background subtraction algorithms, such as background clutter and flickering lights.

Many metrics can be used to assess the output of a background subtraction algorithm given a series of ground truths for several frames in each sequence. Various relative metrics can be calculated based on the numbers of true and false positives and negatives (*TP*, *FP*, *TN*, and *FN*). These metrics are most widely used in computer vision to assess the performance of a binary classifier, as in [20]. *PWC* is the percentage of incorrect classifications in the entire image. This measure represents a trade-off between the abilities of an algorithm to detect foreground and background pixels. In general, a lower value of this estimator indicates better performance.

$$PWC = \frac{FN + FP}{TP + TN + FP + FN} \times 100 \quad (5)$$

The proposed approach relies on several parameters originating from the ViBe algorithm:  $N$ ,  $R$ ,  $\lambda$ , and  $\phi$  [7]. Considering our aim of evaluating the overall performance of the algorithms, we chose a unique set of parameters that yielded sufficiently good results on the complete dataset. The parameters were set to  $N = 20$ ,  $R = 20$ ,  $\lambda = 2$ , and  $\phi = 16$ , respectively. To ensure a fair comparison of the performances of the various algorithms, all algorithms were applied without morphological filtering.

The first sequence, *child\_no1*, shows a child and two adults playing in a living room. The main difficulties in this sequence are light reflections and subjects that sometimes remain still or move only slowly. Figure 2 shows the segmentations produced by the four methods as well as the original color and depth frames and the hand-generated segmentations (ground truths). Ghosts appear in some frames for MOG4D, DECB, ViBe and ViBe1D, greatly reducing the effectiveness of foreground detection. The ViBe algorithm yields worse results than the other algorithms in frame 160 because of the reflection in the color image. In general, the DEVB algorithm achieves improvement over ViBe by virtue of the additional depth model, which allows the ghosts and black shadows to be effectively removed.



**Figure 2.** Comparison of background/foreground segmentation images generated by various background subtraction techniques for four frames taken from the *child\_no1* sequence without morphological filtering. The segmented images produced by our method are the closest to the ground-truth references.



Table 1 shows the quantitative *PWC* results obtained by the four approaches on the evaluation frames from the *child\_no1* sequence. All RGB-D approaches achieve improvements with respect to ViBe, obtaining lower *PWC* values. A lower *PWC* value indicates better performance. The proposed DEVB algorithm achieves the lowest average error rate of 7.023% in Table 1, which indicates that our method performs better than the other algorithms.

**Table 1.** Segmentation evaluation for the *child\_no1* sequence. The table shows the *PWC* results for the various approaches on four different evaluation frames and the mean values for this sequence.

<i>child_no1</i>	Evaluation Frame				Global
Approach	22	46	70	160	Mean
DEVB	<b>6.114</b>	<b>6.035</b>	8.826	<b>7.118</b>	<b>7.023</b>
ViBe	7.808	8.932	<b>7.879</b>	16.455	10.269
ViBe1D	11.183	10.686	9.458	10.848	10.544
MOG4D	7.519	6.559	8.003	8.13	7.553
DECB	8.271	8.072	8.672	8.334	8.337

The second sequence, *new\_ex\_occ4*, shows two individuals walking in front of a coffee shop. The main difficulties presented by this sequence are flickering lights and areas where depth information cannot be obtained by the active infrared sensor. Figure 3 shows the segmentations produced by the four approaches. Our DEVB algorithm achieves good results, whereas ghosting greatly reducing the effectiveness of the other algorithms in foreground extraction.



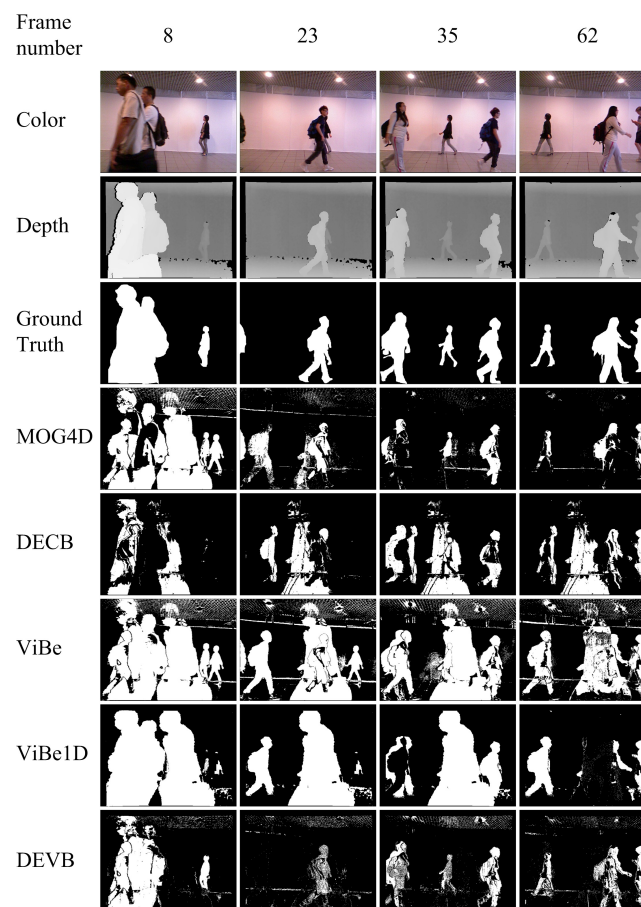
**Figure 3.** Results obtained on the test sequence *new\_ex\_occ4*.

Table 2 shows the quantitative *PWC* results obtained by the four approaches on the evaluation frames from the *new\_ex\_occ4* sequence. The proposed DEVB algorithm achieves the lowest *PWC* of 3.968, which indicates that our method performs better than the other algorithms. ViBe yields the worst result of  $PWC = 12.872$  on this sequence.

**Table 2.** Segmentation evaluation for the *new\_ex\_occ4* sequence. The table shows the *PWC* results for the various approaches on four different evaluation frames as well as the mean values for this sequence.

<i>new_ex_occ4</i>	Evaluation Frame				Global
Approach	15	24	36	42	Mean
DEVB	5.129	2.997	3.261	4.485	3.968
ViBe	12.294	12.394	12.684	14.116	12.872
ViBe1D	12.096	12.554	12.461	11.999	12.278
MOG4D	11.685	11.157	4.845	5.189	8.219
DECB	9.753	8.522	8.706	9.909	9.223

The third sequence, *walking\_occ1*, shows a few people walking in and out of the camera field. In addition, there are flickering lights on the ceiling and sudden illumination changes. Figure 4 shows the segmentations produced by the four approaches. DECB and ViBe1D are less affected by the sudden illumination changes. In addition to ghosting, ViBe1D results in black shadows in frames 35 and 62. The reason for the generation of black shadows is that a new moving object reaches the previous position of an old target. Because the pixel values are similar, the foreground is misclassified as background.



**Figure 4.** Results obtained on the test sequence *walking\_occ1*.

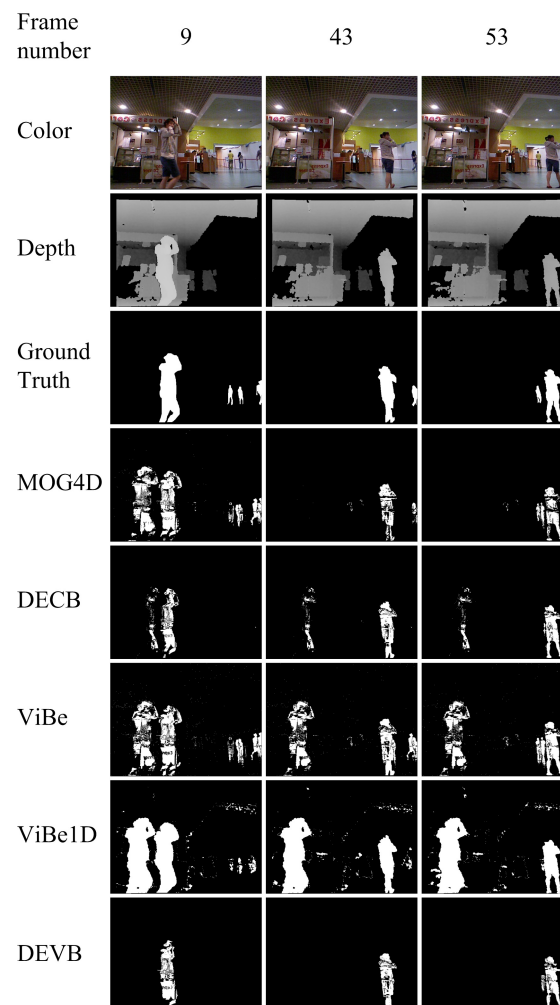


Table 3 shows the quantitative *PWC* results obtained by the four approaches on the evaluation frames from the *walking\_occ1* sequence. The DEVB algorithm achieves  $PWC = 8.847$  in frame 8, whereas the *PWC* values obtained by the other algorithms are higher by more than a factor of three. Moreover, despite being affected by illumination changes in the RGB space, DEVB achieves an average *PWC* of 8.721 (Table 3), indicating that our method is fairly robust to difficult situations.

**Table 3.** Segmentation evaluation for the *walking\_occ1* sequence. The table shows the *PWC* results for the various approaches on four different evaluation frames and the mean values for this sequence.

<i>walking_occ1</i>	Evaluation Frame				Global
Approach	8	23	35	62	Mean
DEVB	8.847	7.728	10.17	8.137	8.721
ViBe	33.287	30.628	30.221	24.078	29.554
ViBe1D	29.944	26.806	31.726	14.332	25.702
MOG4D	38.186	11.195	10.798	9.185	17.341
DECB	30.262	18.616	22.371	21.162	23.103

The fourth sequence, *new\_ex\_no\_occ*, shows a lady walking in front of a coffee shop. The scenario is similar to the *new\_ex\_occ4* sequence discussed above. A large amount of noise is generated by the ViBe1D algorithm because of the holes in the original depth image.



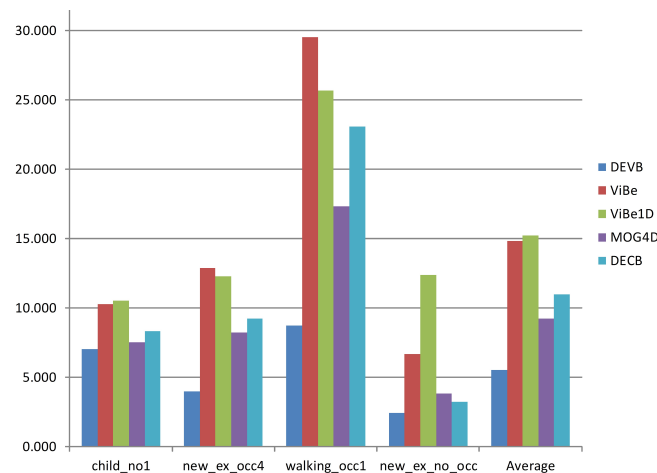
**Figure 5.** Results obtained on the test sequence *new\_ex\_no\_occ*.

Table 4 shows the quantitative *PWC* results obtained by the four approaches on the evaluation frames from the *new\_ex\_no\_occ* sequence. The proposed DEVB algorithm achieves  $PWC = 2.437$ . This value is the lowest in Table 4, which indicates that our method performs better than the other algorithms. ViBe1D obtains the worst result on this sequence, with  $PWC = 12.389$ .

**Table 4.** Segmentation evaluation for the *new\_ex\_no\_occ* sequence. The table shows the *PWC* results for the various approaches on three different evaluation frames as well as the mean values for this sequence.

<i>new_ex_no_occ</i>	Evaluation Frame			Global
Approach	9	43	53	Mean
DEV8	<b>4.102</b>	<b>1.549</b>	<b>1.66</b>	<b>2.437</b>
ViBe	8.188	5.832	6.097	6.706
ViBe1D	14.009	11.249	11.91	12.389
MOG4D	7.977	1.719	1.842	3.846
DECB	5.264	2.222	2.225	3.237

Finally, Figure 6 shows the average *PWC* value obtained by each approach on each benchmark sequence. According to this figure, DEVB yields the best results on every sequence. The *walking\_occ1* sequence is particularly complicated because of the high pedestrian flow; consequently, for each algorithm, the error rate is considerably increased.



**Figure 6.** Average *PWC* values for each of the four sequences and for the entire dataset.

## 5. Conclusion

In this paper, we present an efficient moving object detection algorithm that fuses depth and color information. To incorporate the features of depth images, a depth model is designed in addition to a background model. The inspection mechanism for the classification of foreground pixels is further considered. Finally, we propose a new updating strategy based on the developed background and depth models, which can eliminate ghosting and black shadows almost completely. Experimental results indicate that our method is able to extract foregrounds efficiently, providing an excellent basis for the subsequent motion analysis of a scene. Our proposed method could serve as a convenient research tool for the detection of moving objects captured by the Kinect sensor.

**Acknowledgments:** This work was supported by the National Natural Science Foundation of China (61471157); the Fundamental Research Funds for the Central Universities of China (2011B11114, 2012B07314, 2015B38214); the Natural Science Foundation of Jiangsu Province, China (BK20141157, BK20141159); and the Open Foundation Programs of Changzhou Key Laboratory of Robotics and Intelligent Technology (CZSR2014003).

**Conflicts of Interest:** The authors declare no conflicts of interest.

## References

1. Migliore, D.A.; Matteucci, M.; Naccari, M. A Revaluation of Frame Difference in Fast and Robust Motion Detection. Proceedings of the 4th ACM International Workshop on Video Surveillance and Sensor Networks; ACM: New York, NY, USA, 2006; VSSN '06, pp. 215–218.
2. Elhabian, S.Y.; Elsayed, K.M.; Ahmed, S.H. Moving Object Detection in Spatial Domain using Background Removal Techniques - State-of-Art. *Recent Patents on Computer Science* **2008**, *1*, 32–54.
3. Barron, J.L.; Fleet, D.J.; Beauchemin, S.S. Performance of optical flow techniques. *International journal of computer vision* **1994**, *12*, 43–77.
4. Barjatya, A. Block matching algorithms for motion estimation. *IEEE Transactions Evolution Computation* **2004**, *8*, 225–239.
5. Amamra, A.; Mouats, T.; Aouf, N. GPU based GMM segmentation of kinect data. International Symposium Elmar, 2014, pp. 1 – 4.
6. Murgia, J.; Meurie, C.; Ruichek, Y. An Improved Colorimetric Invariants and RGB-Depth-Based Codebook Model for Background Subtraction Using Kinect. In *Human-Inspired Computing and Its Applications*; Springer, 2014; pp. 380–392.
7. Barnich, O.; Van Droogenbroeck, M. ViBe: A universal background subtraction algorithm for video sequences. *Image Processing, IEEE Transactions on* **2011**, *20*, 1709–1724.
8. Yin, C.; Luo, Y.; Zhang, D. A novel background subtraction for intelligent surveillance in wireless network. Wireless Communications and NETWORKING Conference, 2014, pp. 3017–3021.
9. Kryjak, T.; Gorgon, M. Real-time implementation of the ViBe foreground object segmentation algorithm. Computer Science and Information Systems, 2013, pp. 591–596.
10. Han, J.; Shao, L.; Xu, D.; Shotton, J. Enhanced computer vision with microsoft kinect sensor: A review. *Cybernetics, IEEE Transactions on* **2013**, *43*, 1318–1334.
11. Moeslund, T.B.; Hilton, A.; Krüger, V. A survey of advances in vision-based human motion capture and analysis. *Computer vision and image understanding* **2006**, *104*, 90–126.
12. Ji, X.; Liu, H. Advances in view-invariant human motion analysis: a review. *Systems, Man, and Cybernetics, Part C: Applications and Reviews, IEEE Transactions on* **2010**, *40*, 13–24.
13. Spinello, L.; Kai, O.A. People detection in RGB-D data. IEEE/RSJ International Conference on Intelligent Robots & Systems IEEE/RSJ International Conference on Intelligent Robots & Systems, 2011, pp. 3838–3843.
14. Xia, L.; Chen, C.C.; Aggarwal, J.K. Human detection using depth information by Kinect. *Applied Physics Letters* **2011**, *85*, 5418–5420.
15. Ming, M.; Fangbo, Y.; Qingshan, S.; Yao, S.; Zhizeng, L. Human motion detection based on the depth image of Kinect. *Chinese Journal of Scientific Instrument* **2015**, *2*, 017.
16. Crabb, R.; Tracey, C.; Puranik, A.; Davis, J. Real-time foreground segmentation via range and color imaging. Computer Vision and Pattern Recognition Workshops, 2008. CVPRW '08. IEEE Computer Society Conference on, 2008, pp. 1–5.
17. Schiller, I.; Koch, R., Improved Video Segmentation by Adaptive Combination of Depth Keying and Mixture-of-Gaussians. In *Image Analysis: 17th Scandinavian Conference, SCIA 2011, Ystad, Sweden, May 2011. Proceedings*; Springer Berlin Heidelberg: Berlin, Heidelberg, 2011; pp. 59–68.
18. Fernandez-Sanchez, E.J.; Javier, D.; Eduardo, R. Background Subtraction Based on Color and Depth Using Active Sensors. *Sensors* **2013**, *13*, 8895–915.
19. Song, S.; Xiao, J. Tracking Revisited Using RGBD Camera: Unified Benchmark and Baselines. IEEE International Conference on Computer Vision, 2013, pp. 233–240.
20. Goyette, N.; Jodoin, P.M.; Porikli, F.; Konrad, J.; Ishwar, P. changedetection.net: A New Change Detection Benchmark Dataset. Computer Vision and Pattern Recognition Workshops (CVPRW), 2012 IEEE Computer Society Conference on, 2012, pp. 1–8.



© 2017 by the authors. Licensee Preprints, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC-BY) license (<http://creativecommons.org/licenses/by/4.0/>).