

Event-driven system for fall detection using body-worn accelerometer and depth sensor

ISSN 1751-9632

Received on 21st February 2017

Revised 31st August 2017

Accepted on 10th October 2017

E-First on 27th November 2017

doi: 10.1049/iet-cvi.2017.0119

www.ietdl.org

Michał Kepski¹, Bogdan Kwolek² ✉¹University of Rzeszów, Pigońia 1, 35-310 Rzeszów, Poland²AGH University of Science and Technology, 30 Mickiewicza Av., 30-059 Kraków, Poland

✉ E-mail: bkw@agh.edu.pl

Abstract: The authors present efficient and effective algorithms for fall detection on the basis of sequences of depth maps and data from a wireless inertial sensor worn by a monitored person. A set of descriptors is discussed to permit distinguishing between accidental falls and activities of daily living. Experimental validation is carried out on the freely available dataset consisting of synchronised depth and accelerometric data. Extensive experiments are conducted in the scenario with a static camera facing the scene and an active camera observing the same scene from above. Several experiments consisting of person detection, tracking and fall detection in real-time are carried out to show efficiency and reliability of the proposed solutions. The experimental results show that the developed algorithms for fall detection have high sensitivity and specificity.

1 Introduction

With the rapidly growing aging population on a global scale, the need for improving elderly wellbeing is getting crucial. Smart home technologies can be utilised as means to improve both the quality of care and wellbeing of dependent people. Its form called assistive domotics focuses on making it possible for seniors and people with disabilities to remain at home, safe and comfortable. Smart home technologies are becoming a viable option for older adults who would prefer to stay in the comfort of their homes in place of a move to a retirement home or a healthcare facility [1].

The aim of user-centred ubiquitous computing is to develop solutions for personal assistance, which at the same time sense variations in a human environment and dynamically respond to user needs. It is self-evident that such a technology has strong potential to cope with major societal challenges posed by aging society [2]. Such an increased level of intelligence has a potential to provide improved quality of care in addition to helping elderly people access the knowledge required to offer better decisions when interacting with smart environments [3]. One of the crucial factors that at present pose serious bottlenecks to augment people's lives with ubiquitous computing in a broader scale is a reduced number of affordable energy-saving devices for human activity monitoring and/or energy-efficient units.

Falls are leading cause of morbidity from injury and mortality in the elderly. They are the major reason of injury-related hospitalisation in persons aged 65 years and over and account for a significant fraction of all hospital admissions in this age-group [4]. Even falls that do not lead to physical injuries can result in the so-called post-fall syndrome [5], which typically manifests itself in loss of confidence, loss of muscle and control, problems with balance, and walking disorders leading to loss of mobility and independence. Fear of falling has been identified as one of the key symptoms of this syndrome. The cause for this is that seniors are afraid to lie after the fall on the floor in solitude and without help for a long time [6]. It has also been shown that getting up quickly after the fall can reduce the risk of death even by 80% and the necessity of hospitalisation by 26% [7]. Thus, falls should be detected as early as possible. For this reason, assistive technologies have the strong potential not only to assist in daily activities, but they also have capabilities to reduce risks of fall events.

With the goal to permit prolonged independent living in a secure and homely environment, reliable fall detection is a significant task in the area of ambient-assisted living (AAL) [8].

Medical alert systems were introduced in the 1980s as uncomplicated push-button devices worn around the neck. Afterwards, they were extended about accelerometer-based algorithms to automatically raise an appropriate alert that a fall has occurred. One of the biggest limitations of such automatic fall detection systems is the occurrence of false alarm alerts. When such a device is used, the false alarm could be triggered by an everyday activity such as quickly dropping into a seated position in a chair or even by bending down.

The ambient device-based systems are capable of detecting falls in a non-intrusive way by exploiting audio, vibration, pressure and visual information, to name a few of the most frequently used sources of information in this domain [1]. There are several types of sensors used in this field, including measuring the vibration of the floor to detect falls [9], detecting falls by using pressure mats [10] or impulse-radar sensors [11]. A fall detection system relying on one of the mentioned above sensors typically has a high false alarm rate. None of the sensors mentioned above, if used separately in a fall detection system, is able to meet the requirements of end-users regarding the level of false alarms. Despite the enormous effort of research and the number of IMU-based devices on the market, there is still no system that has sufficient reliability and is accepted by end-users.

To overcome the limitations of these devices, a wide range of vision-based monitoring systems with fall detection functionality have been proposed over recent years [12–14]. Vision-based systems usually use image sequences to analyse motion features of the human body and distinguish the features of fall events from non-fall activities in order to infer about the occurrence of fall. They provide valuable information for assistive monitoring but raise privacy concerns. Besides, such systems fail to work in darkness or when the elderly is outside of the observed area. Their major advantage is that the user does not need to wear any specialised apparatus. However, most of these solutions are energy demanding and expensive. Moreover, their deployment is cumbersome, and only a few of them can meet the demands of the end-users in the detection of the fall in real homes or health care facilities.

Event-driven systems, in which the system activities are triggered in response to events, usually representing a significant change of the state of controlled or monitored physical variables, exhibit certain advantages over other approaches, particularly, in resource-constrained applications. The last few years have witnessed an upsurge in the research interest to harness the

advantages of event-based paradigm applied to a wide spectrum of engineering disciplines including signal processing and control. The application areas of such systems include energy-efficient control, energy-efficient signal processing, rehabilitation [15], event-driven visual attention [16], or frame-free event-driven vision systems [17], to mention a few.

In this work, we present an event-driven system for fall event detection using measurements from a body-worn accelerometer and depth sensor(s). In response to significant motion variation indicated by an accelerometer, the system fetches depth maps from a circular buffer and then processes them to validate the fall event. This way, the most time consuming depth image processing is executed only in the case of a significant change of the person's motion, i.e. high likelihood of fall occurrence. We discuss and compare the efficiency of fall detection on depth maps provided by a ceiling-mounted camera as well as a wall-mounted (facing the user) camera. With the purpose of extending the viewing area, the overhead camera was mounted on a pan-tilt motorised head. The aim of the controller of the pan-tilt unit was to keep the moving person in the centre of the current depth map. We discuss the person delineation as well as feature extraction in both camera settings. We show experimentally that the results achieved in both camera settings are promising. We discuss the advantages and limitations of the considered camera setups. We show experimentally that a two-camera system achieves perfect classification performance on data from the freely available URFD dataset.

2 Background and related work

Recent fall alert devices are usually able to recognise when a person wearing the device has fallen by using accelerometers and optionally gyroscopes, and through detecting changes in the body's orientation and speed. An obvious limitation of such devices is that the senior may be not capable to press the emergency button after the fall due to loss of consciousness or just because of over-excitation. Thus, the applicability of such devices is limited to niche markets as nursing homes. Moreover, today's devices are not widely accepted by primary end-users [7], particularly those who are not impaired. The reason for this is that current systems are not able to guarantee good sensitivity (nearly 100%) with enough specificity to limit the number of false alarms [18, 19]. An alarm is false when it is triggered unnecessarily or for a cause other than fall event. In practice, this means that certain fall-like activities activate alarms, which in turn lead to irritation of the end-users. The reason for this is that the acceleration ranges are overlapping for falls and activities of daily living (ADLs).

Besides the solutions outlined above, more complex systems are now utilised to improve the fall detection accuracy [8, 20]. Such fall detection systems can be divided into two major categories, that is, based on wearable sensors and context-aware systems [4]. Micro-electro-mechanical systems (MEMS) are extensively used in a wide range of applications. MEMS accelerometers are one of the most common types of MEMS sensors, due to their simplicity, ease of fabrication, low price and good usability [21]. In comparison with vision sensors, wearable inertial sensors are lighter, smaller, easier to use, and most importantly, they consume less energy and are far cheaper. They allow collection of data outside of laboratory environments and are perceived as one of the best sensors for AAL [1]. Hence, many different algorithms have been proposed to explore, support or improve fall detection using the only accelerometer(s) [22] or an inertial measurement unit(s) [23, 24].

Usually, approaches relying on a body-worn accelerometer utilise a threshold-based algorithm to examine if a person's movement is higher than some preset threshold [22, 25]. However, as shown in [19], such systems are too sensitive and thus generate a substantial number of false alarms. A similar conclusion has been drawn by an international group of researchers [18], which evaluated the effectiveness of threshold-based algorithms to identify falls on data from real falls. The dataset of 29 real-world falls contains accelerations of persons' movement, each for a period of two days. In the evaluation, 13 different algorithms were examined with respect to their capability of identifying real falls.

Regrettably, none of the examined algorithms gave satisfactory results in terms of both sensitivity (capability of recognising falls that in reality took place) and specificity (ability to properly recognise a movement as a non-fall).

Fall detection methods relying on body-worn accelerometers can be ineffective in the detection of slow falls [26], such as collapsing after a heart attack, which usually does not feature significant accelerations. Moreover, as noticed in [18, 27], some fall phases detected in experimentally simulated falls are usually not detectable in acceleration signals from heterogeneous realworld falls. Thus, having regard to the reduced availability of motion data with real-world falls, the usefulness of machine learning-based methods might be limited in practice. In general, although it is not easy to distinguish between falls and fall-like activities, the inertial sensors are thought to be very useful sensors in fall detection [4]. They are now massively utilised in the mobile smart devices, including smartwatches. Smartwatches are lightweight and waterproof, so they can be kept on when taking a bath. Overall, the inertial sensors can greatly support fall detectors built on other sensors, including vision and depth sensors [20, 28–30].

Among the possible types of context-aware detectors, vision systems offer a promising way of recognising human actions [31] as well as detecting human falls [32]. A variety of vision-based fall detection algorithms have been developed in recent years [13, 14]. One of their advantages is that the monitored person does not need to wear any special apparatus. On the other hand, this form of fall-monitoring is both most intrusive and most expensive. Despite many approaches to preserve privacy, people in the observed spaces still have the feeling of being-watched. In consequence, the usual CMOS/CCD cameras are very often unacceptable, especially in the bedrooms or bathrooms. Moreover, while near-infrared light sources made it possible to record video under low-light conditions and during the night, the quality of the videos might be insufficient to achieve automatic fall detection with high sensitivity and specificity. Nevertheless, thanks to technology progress in the area of the smart camera [33] and smart home, the CMOS/CCD camera-based systems offer many monitoring capabilities.

The cameras providing in real-time the depth maps can considerably enhance the detection and tracking performance making possible to reliably extract the head trajectory, which has been proven to be very useful in fall detection [34]. An entire view of the scene can be very advantageous in fall detection. In [35], it has been demonstrated how the omnidirectional cameras can be utilised to achieve coupled fall detection and tracking. In general, the omnidirectional cameras have been proven to be very useful if big visual field coverage is desired. Thermal video cameras, which detect the amount of thermal radiation emitted/reflected from objects in the scene can also provide very informative information for reliable fall detection [36].

Just a few years ago, the Kinect's sensor was proposed for detection of humans' falls [28, 37, 38]. As shown experimentally [28, 37], the depth maps delivered by the Kinect sensor are enough to extract the person from the background. What is more, owing to the estimation of dense depth maps on the basis of the speckle pattern of infrared laser light, the detection of the person can be done anytime. Despite several approaches to Kinect-based fall detection [38], the existing algorithms do not provide both high sensitivity and specificity. By integrating the acceleration data with video or depth maps [28, 39], the recognition of activities [30], as well as emergency situations, can be noticeably improved.

Our work differs from the research in the area of fall detection (e.g. [38]) in that we do not use only inertial device or depth maps standalone but we use both an inertial unit and depth sensor. The rationale for such an approach is that the current accelerometer-based algorithms being sensitive, generate too many false alarms. Assuming that such algorithms typically produce a few alarms a day [19], our approach is able to reduce the false alarm ratio to almost null. Having regard that accelerometers are frequently available in smartwatches, as well as considering further progress in this area, the obtrusiveness and the discomfort when wearing such a device will be limited. With a view to high computation cost of a vision-based algorithm for fall detection, we apply the event-driven approach to data processing and system design. This allows

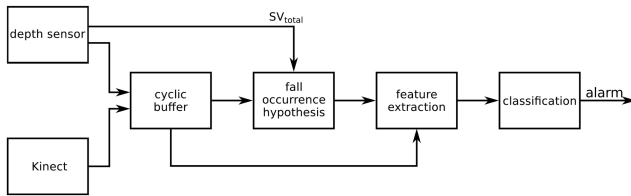


Fig. 1 Event-driven architecture for fall detection

us not only to reduce computational overload but also allows us to determine precisely the time at which the impact took place. Since in the relevant literature [14, 38] there is no detailed comparison of approaches for fall detection on the basis of ceiling-mounted and wall-mounted cameras, we discuss algorithms for both approaches as well as present experimental results on freely available fall detection dataset. Since in the second approach an active camera is recommended in order to extend the observation area, we develop an effective algorithm for person's head detection and tracking. We show experimentally that the results achieved by this algorithm on maps acquired by the active camera are promising.

3 Architecture and main components of the system

While embedded vision is comparatively a novel term, as a technology, it is highly established in a number of domains, and the most successful applications are in the area of factory automation. Smart cameras are another example of successful applications of the embedded vision technology [33]. In addition, one can specify a number of successful applications of this technology in surveillance and transport. Healthcare is one of the main application areas for the embedded vision [3]. The embedded vision technology has significant potential to change the health monitoring in home, for instance through mobile phone applications for monitoring the user's state of health and reporting it to a medical centre. One of the successful examples of embedded vision systems is the Microsoft Kinect game controller [40], which has been designed to perform real-time tracking of the movement of the users. Although Kinect was initially devised only as a motion sensing device for computer games, a strong interest of the computer vision community led to developing several new applications, including applications for activity recognition [31] or rehabilitation [38].

In our approach, a body-worn accelerometer is utilised to indicate a potential fall event and a depth camera is employed to authenticate fall alert. The proposed event-driven sensor data processing method fetches from a circular buffer a sequence of depth maps, which were acquired prior to the fall and then processes it to authenticate the fall alert, instead of processing data frame-by-frame, see Fig. 1. In general, if the person acceleration is higher than a predetermined threshold the algorithm executes a lying pose detector as well as optionally employs a dynamic feature to finally confirm the fall. In consequence, more computationally demanding authentication of the fall is not processed frame-by-frame. Such data stream processing has been designed specifically to operate with the least amount of energy consumed while achieving reliable fall detection in real-time [41].

The presented system can operate in two main modes. In the first mode, the fall authentication is achieved using depth maps acquired by a static depth sensor facing the scene, whereas in the second one the verification of the fall is achieved using depth maps provided by an active ceiling-mounted camera. The main difference between the modes of the system lies in the person detection algorithm. In the first mode with a static depth sensor, the person is extracted by differencing the current depth map from an accommodated depth map of the background. In the second mode with the active camera, the person is delineated using depth region growing followed by a person's head detector. If a person moves the system delineates the person in each frame to extract his/her centroid, which is in turn required by a controller of the active camera to keep the target in the centre of the current depth map.

The fall detection algorithm is executed on a PC or on PandaBoard depending on the configuration. It runs under the Linux operating system. The accelerometric data are acquired by the x-IMU device and then transmitted by Bluetooth to the receiver device of the processing unit. The Xbox Kinect sensor is connected to the processing board via USB. The connection between the microcontroller of the active camera and the board is realised by I2C bus.

4 Person detection in depth maps

In this section, we discuss algorithms for person extraction in depth map sequences. The next subsection is devoted to explaining how a person is delineated in depth maps acquired by a Kinect facing the scene, whereas the subsequent subsection details the method for person detection in depth maps acquired by a Kinect and mounted on the ceiling, i.e. providing the top view of the scene.

4.1 Person detection in frontal depth maps

The key technology behind Kinect is a variant of structured light in which a pseudo-random speckle pattern is projected onto the scene by a laser-based infrared (IR) emitter and then observed by an IR camera. The shift of such a speckle pattern in space is measured and after that mapped to depth through triangulation. However, the depth maps acquired in such a way often contain much noise. Thus, typical detectors when trained from the widely applied image feature descriptors, which demonstrated to be successful in visible images, cannot achieve promising results. As demonstrated in [28, 37], depth information is sufficient to extract human by the use of depth background maps collected by a fixed Kinect. As noticed in [28], person extraction on the basis of depth background maps can be done at low computational cost.

In our event-driven approach, the algorithm extracts the person at low computational cost and then processes the foreground image to prove whether a more costly update of depth background is needed. Moreover, the accommodation of depth background maps is done only in map areas in which the scene changes. The scene changes are detected with low computational cost through extracting coherent depth maps on the foreground map and then examining if the size of the component with person increased considerably, for instance, due to opening a door, or the number of the foreground components is larger than one, i.e. if there is a non-person object of sufficient size in the foreground.

In the person extraction algorithm, we can distinguish a part that is executed every frame and a part, which is evoked when there is a scene change, see lines 1–7, 22 and 8–20 in Algorithm 1 (see Fig. 2), respectively. Let us assume that there is given a background model $B(x, y)$ and a buffer Q consisting of Q_{size} last depth frames. In each frame, the algorithm takes a new frame $D_t(x, y)$ and then extracts the foreground $F_t(x, y)$ through determining the absolute value of the difference between the current depth map and the depth background map, see lines 2 and 3. Then, the algorithm determines the connected components in the binarised foreground images, calculates their number as well as their areas. Afterwards, it determines the blob belonging to a person, see line 5, and stores it on the person image $P_t(x, y)$. Finally, it examines if the number of blobs is larger than one or there is a significant change of the person area on consecutive person images. If both conditions are false the algorithm returns the previous background model, see lines 22 and 23.

A modification of the scene layout requires an update of the depth background map. Fig. 3 illustrates a dynamic scene, where a person closes the door. In such a scenario, the depth model should accommodate scene changes. In our approach, the depth background map is a temporal median over a set of depth maps. The depth background map is updated only in a region of interest (ROI), which is represented by a rectangular sub-image. It contains the foreground objects, see Fig. 3d.

In the case of merging the person blob with a non-person blob, i.e. if there is a change of the area between the current person blob P_t and blob P_{t-1} representing the person in time $t - 1$, see also line 7 in Algorithm 1 (Fig. 2), the algorithm determines the seed for the

- 1: Given a background model $B(x, y)$, and buffer
 $Q = \{D_{t-1}(x, y), D_{t-2}(x, y), \dots, D_{t-Q_{size}}(x, y)\}$
- 2: Acquire new depth map $D_t(x, y)$ ▷ Fig. 2c
- 3: Determine foreground

$$F_t(x, y) = \begin{cases} D_t(x, y), & \text{if } |D_t(x, y) - B(x, y)| \geq B_{th} \\ 0, & \text{otherwise } |D_t(x, y) - B(x, y)| < B_{th} \end{cases}$$
- 4: Determine Blobs on $F_t(x, y)$ using connected comp. ▷ Fig. 2d
- 5: Assign to person image $P_t(x, y)$ the Blob with the greatest similarity to $P_{t-1}(x, y)$
- 6: Determine number of Blobs N_b
- 7: **If** $\frac{area(P_t(x, y))}{area(P_{t-1}(x, y))} > T_a$ or $N_b > 1$
- 8: Determine

$$F_{t-3}(x, y) = \begin{cases} D_{t-3}(x, y), & \text{if } |D_{t-3}(x, y) - B(x, y)| \geq B_{th} \\ 0, & \text{otherwise } |D_{t-3}(x, y) - B(x, y)| < B_{th} \end{cases}$$
- 9: Determine ROI , allocate stack S , determine seed on the basis of F_{t-3} , allocate logical table L and initialize it with false
- 10: $P_t(x, y) = \text{RegionGrowing}(D_t(x, y), \text{seed})$
- 11: $D'_t(ROI) = D_t(ROI) - P_t(ROI)$ ▷ Fig. 2f
- 12: Push $D'_t(ROI)$ on stack S
- 13: $L = L$ or logical($D'_t(ROI)$) ▷ Fig. 2g
- 14: **If** for each (x, y) , $L(x, y) \neq \text{false}$
- 15: $B'(x, y) = \text{Median}(S)$
- 16: return $B'(x, y)$
- 17: **Else**
- 18: acquire new depth map
- 19: determine seed
- 20: go to line #9
- 21: **Else**
- 22: $B'(x, y) = B(x, y)$
- 23: return $B'(x, y)$ ▷ Fig. 2h

Fig. 2 Algorithm 1: Person extraction using depth reference map

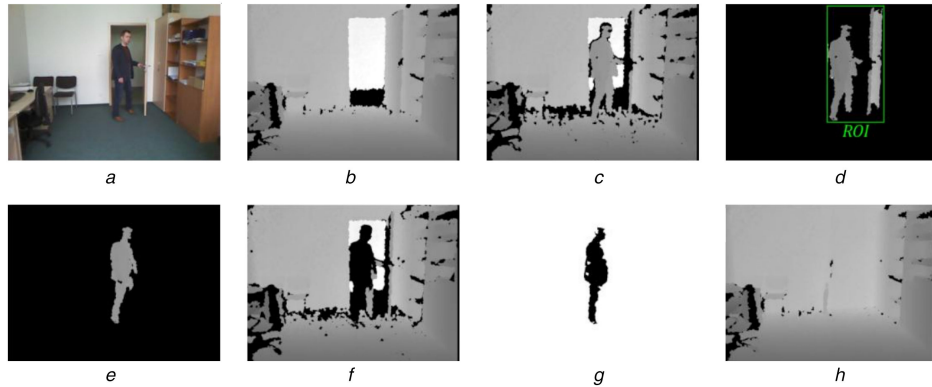


Fig. 3 Person extraction in a dynamic scene

(a) RGB input image, (b) Initial depth background model, (c) Input depth map, (d) Foreground blobs, (e) Segmented person, (f) Input depth map after removing person, (g) Logical table L , (h) Updated depth background model of the scene

region growing. The location of the seed is determined on the basis of an image F , which has been determined before the merging of the blobs. In the current implementation, it is determined in time $t - 3$ since it gave better results in comparison with $t - 2$. After determining the foreground image $F_{t-3}(x, y)$, the algorithm uses it to determine the seed region, which is then used by a region growing procedure, see the tenth line in Algorithm 1 (Fig. 2). Then, in the region constrained by the ROI, the person blob extracted by the region growing is removed from the current depth map, see the 11th line in Algorithm 1 (Fig. 2) and Fig. 3f. The image $D'(x, y)$ extracted in such a way is then pushed on a stack S , which holds a set of depth maps required for determining the temporal median. The temporal median is calculated if all pixels in the logical table L are true. This means that all pixels in the ROI region changed the value from false to true, i.e. that at every (x, y) location, at least one

background pixel D' has been stored in every location (x, y) in S , which in turn is employed in the median filtering. After updating the depth background map, see the 15th line in Algorithm 1 (Fig. 2) as well as Fig. 3h, the algorithm returns the background model. By comparing the depth background map before the scene change, see Fig. 3b, and the depth background map, which has been extracted after the scene change, see Fig. 3h, we can notice that the model accommodated to change the scene. In particular, it contains only objects belonging to the room. As we can notice, the depth background model takes into account the closed door. The thresholds T_a and B_{th} were determined experimentally. The depth maps acquired by the Kinect sensor are stored in a circular buffer of size 15, i.e. Q_{size} is set to 15. In the current implementation, it contains also depth maps for the calculation of dynamical features

```

1: Assign LA to seed pixels, calculate RG_M, insert the neighboring pixels in NHQ and assign
   them IN_NQ
2: While NHQ  $\neq \emptyset$  and QM  $\neq \emptyset$ 
3:   While NHQ  $\neq \emptyset$ 
4:     Delete pixel from NHQ
5:     Calculate its  $\delta$ 
6:     If  $\delta > \delta_{th}$ 
7:       continue
8:     Insert it in QM queue with index  $\delta$ 
9:     Assign him IN_Q label
10:  If QM  $\neq \emptyset$ 
11:    Delete FQ from QM
12:    While FQ  $\neq \emptyset$ 
13:      Delete pixel from FQ
14:      Assign him LA label
15:      Assign its neighbors with NO_L to NHQ
        and assign them IN_NQ label
16:  Actualize RG_M about the pixel values from FQ

```

Fig. 4 Algorithm 2: Depth region growing

Table 1 Notation in the Algorithm 2 (Fig. 4)

RG_M	mean value of the depth region
δ_{th}	threshold for δ value
NHQ	Neighbour holding queue
QM	map of queues holding δ
FQ	queue with smallest δ
Labels	
NO_L	not visited pixel
IN_NQ	pixel is inserted in NHQ
IN_Q	pixel is in QM queue
LA	pixel is assigned to the region

as well as maps for re-initialisation of the background model. The region growing function is discussed below.

4.2 Person detection in depth maps from ceiling-mounted active camera

The Kinect sensor has an angular field of view of 43° vertically and 57° horizontally. The observation area of an overhead Kinect mounted at the altitude of 2.6 m from the floor is about 5.5 m². To increase the field of observation, a home-made pan-tilt head has been utilised to rotate the Kinect sensor. Thanks to the use of such a pan-tilt motorised head the observation area covered by the device is far larger and in effect, the Kinect can cover a typical room, say 15–20 m² [42]. When a person moves, a proportional integral controller rotates the camera in order to keep him/her in the central part of the image. The person is detected in real-time on the basis of a depth region growing. The person's position is represented by the centroid of the delineated blob. To decrease the number of pixels that can be potentially included into the person blob, in advance, the algorithm detects the floor using the RANSAC algorithm [43].

The seeded region growing [44] was originally designed for intensity, i.e. grey-value images. The method takes a set of seeds as starting points along with the image. The regions are then grown from these seed pixels to pixel neighbours depending on a region membership criterion, which determines whether the adjacent point should join a region or not. The magnitude of the difference δ between pixel's intensity value and the region's mean is used as a decision criterion. The order in which the pixels are processed is determined by a global priority queue (PQ), which orders all candidate pixels by their fitness scores. This way the pixel with the smallest measured difference is assigned to the respective region. If the unlabelled pixel meets two or more boundary pixels from adjacent regions, it is joined to a region that has the smallest

similarity distance and then it is marked as a border region. The above process continues until all pixels are assigned to a region.

A disadvantage of the original region growing is that it does not update the previous entries in the sequentially sorted list to reflect new differences from a region whose mean has been updated. In an improved seeded region growing [45] the border pixels, which have the same minimum δ value, are processed in parallel. The improved algorithm employs an ascending PQ and several last in, first out (LIFO) queues, where each LIFO queue contains pixels with the same δ value. When a new pixel is added to the PQ, it is inserted into a LIFO queue that corresponds to the pixel's δ value. This means that instead of removing individual pixels from the PQ, the entire LIFO queue corresponding to the smallest δ value is removed.

In our approach, the person is extracted in the sequence of depth maps using a modified region growing, which starts from a single seed region. The pseudo-code of the algorithm is given in Algorithm 2 (see Fig. 4), whereas the symbols utilised in the pseudo-code are explained below. A neighbour holding queue (NHQ) contains the pointers to pixels neighbouring with the depth region, whereas a QM holds the NHQ indexed according to δ (see Table 1). The queue with the smallest δ is denoted by FQ.

At the beginning, the seed pixels are assigned the LA label, the mean depth value of the pixels belonging to the seed is determined and the neighbouring pixels are inserted in the NHQ. Then, the algorithm iterates until NHQ and QM are non-empty. At the beginning of the iterative process, the algorithm iterates until the NHQ is not empty, see line 3. In the discussed loop, the algorithm deletes the pixel from the NHQ and then calculates its δ . If δ is smaller or equal to a threshold δ_{th} the algorithm inserts the pixel in the QM's queue indexed by the δ value and assigns him the IN_Q label, otherwise, it takes the next pixel from the NHQ. After terminating the loop, the algorithm examines if the QM queue is not empty, see the tenth line in the pseudo-code. If yes, it deletes from QM the queue FQ with the smallest δ and then iterates until the FQ is non-empty. At each step in the loop, the algorithm deletes the pixel from the FQ queue, see the 13th line, assigns him the LA label, and it assigns the neighbouring pixels with NO_L labels to the NHQ as well as changes their label from NO_L to IN_NQ label. After terminating the loop, the algorithm actualises the average value of depth. It continues until NHQ and QM are non-empty.

Fig. 5 demonstrates the extracted person blob by the discussed algorithm together with the images illustrating the region growing stages. Due to the nature of the distribution of depth values on the person's head, i.e. gradually decreasing depth values, the algorithm first extracts the head. This means that if the seed region is located in the head's area, the algorithm will extract the head first, see Fig. 5b, then the arms, see Fig. 5c, and then the remaining body. As we can notice, this creates a possibility for analysis of shapes,

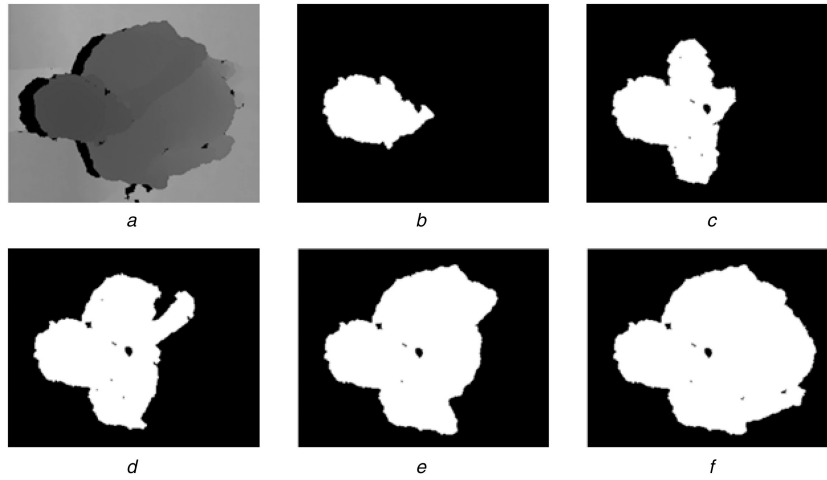


Fig. 5 Person extraction on depth maps acquired by a ceiling-mounted camera using region growing
(a) Depth image, (b)–(f) Binary images illustrating region growing

which arise during region growing to authenticate that the extracted regions belong to the person's area. In particular, the oval shape of the head can be approximated by an ellipse, see Fig. 5b, whereas the head-arm part can be approximated by an ellipse or T-shape like figure, depending on the relative position between a person and the camera, see Fig. 5c.

4.3 Finding person in depth maps

Regular depth region growing suffers from the effect of region chaining (overspill), which takes place when two separate regions are grown into a single region while they are really split. To ameliorate delineation of the subject in such circumstances as well as to improve person following by the active camera, we can configure the system to execute a person's head finder. The detector can also be used for automatic initialisation of person tracking. The head detection is realised by a linear support vector machine (SVM), which operates on the histogram of oriented depth (HOD) features [46]. The HOD features locally to describe the orientation of depth changes. In this work, they are determined in sub-windows of fixed size [47]. The scaling ratio is determined on the basis of the distance between the camera and the head's part, which is closest to the camera. The sub-windows of fixed size undergo a subdivision into cells. The descriptors are extracted from each cell and finally, the oriented depth gradients are assembled into 1D histograms.

5 Feature extraction

The first part of this section is devoted to explaining how we indicate fall events on the basis of motion data. Subsequently, we present recognition of lying pose in depth maps, which are acquired by the frontal as well as the ceiling-mounted depth sensor. Afterwards, we discuss features that describe dynamic transitions of the body. Finally, we explain how person falls are detected.

5.1 Fall indicating using body-worn accelerometer

A lot of various techniques were proposed to achieve reliable fall detection using IMUs [25]. Usually, the accelerometer-based techniques indicate the alarm if the acceleration reaches a certain threshold value. An algorithm proposed in [48] relies on a change in body orientation. It raises alarm if the square root of the sum of the squares of acceleration components exceeds a preset threshold value.

In the discussed system, a fall impact is signalled if the total sum vector SV_{total} is >2.5 g. The SV_{total} value is determined as follows:

$$SV_{total}(t) = \sqrt{A_x^2(t) + A_y^2(t) + A_z^2(t)}, \quad (1)$$

where $A_x(t)$, $A_y(t)$, and $A_z(t)$ stand for the acceleration in the x -, y -, and z -axes at time t , respectively. The SV_{total} value includes both the dynamic and static acceleration components. It equals to 9.81 m/s^2 when the accelerometer has no acceleration, and zero when it is in free fall. Having regard that the potential fall is indicated if the SV_{total} is larger than the experimentally determined threshold, the fall event is signalled only on the basis of body impact, i.e. we do not consider free-fall and post-fall phases. Such phases were considered in approaches aiming at detecting the falls on the basis of acceleration data only [18, 25]. To measure the movements of the whole body, the device was located around the pelvis, which is close to the centre of the body mass. Thus, similar to [18], the device was placed near the spine on the lower back. As pointed out in [49], since the acceleration signal measured from the wrist varies considerably, the signal of ADL samples strongly overlaps with that of the fall events. An analysis [25] of acceleration signals from 240 falls demonstrated that a fall with the smallest trunk magnitude produced a value of 3.5 g. The above-mentioned value provided 100% fall-detection accuracy. In [49], the SV_{total} value from waist was set 2.0 g, i.e. to a smaller value in comparison with the value used in [25]. As noted, this difference might be partly explained by the median filtering, which changes the absolute peak value of the impact signal.

5.2 Fall detection dataset

The classifiers responsible for fall detection were trained on depth map sequences from the URFD dataset [http://fenix.univ.rzeszow.pl/~mkepski/ds/uf.html]. The URFD dataset consists of depth map sequences acquired by two Kinect sensors with the corresponding motion data, which were collected by a body-worn accelerometer. The motion data consisting of the acceleration over time in the x , y , and z axes were acquired by an x-IMU device with a sampling rate of 256 Hz. The frontal depth maps with the corresponding RGB images were acquired by a fixed Kinect that was placed at 1 m altitude from the floor, whereas the top view RGB-D maps were acquired by a second static Kinect, which has been mounted on a ceiling at the height of 3 m from the floor. All depth maps are synchronised with the motion data. The dataset contains 30 image/acceleration sequences with 30 falls, which were simulated by five persons, including one 50+ performer. They simulated falls from standing and from sitting on the chair. A part of the dataset with frontal images contains also 40 image/acceleration sequences that contain typical ADLs like sitting down, picking-up an object from the floor, crouching down, as well as ten data sequences with fall-like activities, consisting of quick lying on the floor and lying on the bed/couch. The sequences with falls consist of 3K images with the corresponding motion exemplars, whereas the total number of images in ADL sequences is equal to 10K.

5.3 Recognition of lying pose

On the basis of features representing the extracted person in the depth maps, we trained classifiers responsible for distinguishing falls from ADLs. The lying pose has been distinguished from ADLs using classifiers trained on features representing the extracted person in the depth maps. For each of the considered camera setting a separate lying pose classifier has been prepared. The training and testing in the facing camera setting have been realised on 1616 and 2425 depth maps from the URFD dataset, respectively. The training and testing in the overhead camera setting have been realised on 350 and 525 depth maps from the URFD dataset as well as additional ADL depth maps, respectively. Such image sets were then employed to build k-nearest neighbors (k-NN) classifiers and to train linear SVM classifiers, whose main task was to check whether the person is lying on the floor. Below we discuss the features that were utilised in both camera settings.

5.3.1 Recognition of lying pose in depth maps from facing camera: The lying pose in frontal depth images has been recognised using both depth features and features expressed as a point cloud. The conversion from the depth data expressed in the 2D array to data expressed as the point cloud has been done using factory calibrated settings. The following features were extracted from the frontal depth maps to recognise the lying pose:

- H/W – the ratio of height to width of the person's bounding box.
- T/T_{\max} – the ratio of the height of the person's surrounding box to the physical height of the person projected onto the depth map.
- D – the distance of the person's centroid to the floor.
- $\max(\sigma_x, \sigma_z)$ – standard deviation from the centroid for the abscissa and the applicate, respectively.
- P_{40} – the ratio of the number of points belonging to the person, contained in a surrounding cuboid of height 40 cm from the floor, with respect to a total number of points belonging to the person.

5.3.2 Lying pose recognition in depth maps from ceiling-mounted camera: The detection of the lying pose in the maps from the ceiling-mounted active camera has been realised on the basis of the following features:

- H/H_{\max} – the ratio expressing the head–floor distance to the person's height.
- $Area$ – the ratio of the person's area in the depth map with respect to the area of the top-view blob of the person in the standing pose.
- l/w – the ratio expressing the major length to major width of person's blob1 in the depth image.

The major length and width (eigenvalues) have been calculated as follows [50]:

$$\begin{aligned} l &= 0.707\sqrt{(a+c) + \sqrt{b^2 + (a-c)^2}}, \\ w &= 0.707\sqrt{(a+c) - \sqrt{b^2 + (a-c)^2}}, \end{aligned} \quad (2)$$

$$a = \frac{M_{20}}{M_{00}} - x_c^2, \quad b = 2\left(\frac{M_{11}}{M_{00}} - x_c y_c\right), \quad c = \frac{M_{02}}{M_{00}} - y_c^2,$$

$$M_{00} = \sum_x \sum_y F(x, y), \quad M_{11} = \sum_x \sum_y xyF(x, y),$$

$$M_{20} = \sum_x \sum_y x^2 F(x, y), \quad M_{02} = \sum_x \sum_y y^2 F(x, y),$$

where $F(x, y)$ indicates a pixel on the binary image representing the extracted person, whereas x, y stand for image coordinates.

5.4 Dynamic transitions

During human fall the head–floor distance changes rapidly due body transition from a vertical orientation to a horizontal one. The distance between the person's centroid and the floor also varies considerably and quickly over the accidental fall. The ratio of the areas occupied by the person in the depth maps from a ceiling-mounted camera also changes meaningfully over the accidental fall. Therefore, by analysing the above-mentioned cues we can settle whether the body transition is intentional or not.

In the setting with the ceiling-mounted sensor, apart from the features discussed in the previous subsection, we employed also a dynamical feature incorporating information about the speed of the falling person's body towards the floor. The speed of the falling body was modelled through the distance between the farthest person points and the floor. The discussed feature was determined in the following manner:

$$h(t) = \frac{H(t)}{H(t - \Delta T)}, \quad (3)$$

where the value of $H(t)$ is determined at the time of the impact and $H(t - \Delta T)$ is calculated ΔT prior the impact. The value of ΔT has been chosen experimentally and it was set to 600 ms. It is worth noting that typical lead times for falls range from 650 to 800 ms. The experimental results showed that in the scenario with a ceiling-mounted depth sensor such a feature quite reliably describes the fall dynamics. In the depth maps from a ceiling-mounted sensor the peak value of $H(t)/H(t - \Delta T)$ for the fall assumes values smaller than one. The accelerometer used as an indicator of the potential simplifies determining this ratio since the impact time t can be determined readily at low computational cost.

5.5 Fall detection

In the setting of the frontal camera, the decision about the fall is taken on the basis of the SV_{total} value (1) and the lying pose detector. If the SV_{total} exceeds the value equal to 2.5 g the lying pose detector is executed to authenticate the fall. This means that the accelerometer filters most of the fall-like activities and as a result, the lying pose detector is executed only in the case of high likelihood of a fall event. In the scenario with the overhead camera, apart from the above mentioned cascade of two classifiers, the dynamical transitions are considered too. If both the accelerometer-based classifier signals the high likelihood of the fall event and the lying pose classifier indicates that the person is lying on the floor, then the threshold-based classifier is executed to examine if the value of dynamic transition (3) of the head assumes a value smaller than 0.6. In consequence, the final decision about the fall in the setting with the ceiling-mounted camera is taken on the basis of a chain consisting of three classifiers. The discussed chain includes a classifier that on the basis of acceleration data signals potential fall, as well as classifiers of lying pose and dynamic transition, which decisions are taken on the basis of depth maps. The final classifier of the fall considers different modalities and has lower false positive rate and high precision.

6 Real-time data acquisition and processing. Tuning and parameters

Having regard that the fall detection system should be inexpensive as well as work anytime and consume low power, we designed an event-driven processing framework in which the body-worn accelerometer is utilised to signal high likelihood of the fall occurrence and depth maps are not analysed frame-by-frame, but instead they are stored in a circular frame buffer. In case of the high likelihood of the fall event, the previously stored depth maps are fetched from the circular buffer and then processed to extract the features. In the setting of the frontal and fixed camera the frame-by-frame calculations comprise person extraction through taking an absolute value of difference between the current depth map and depth background map, determining the connected components as well as calculating the areas and number of the connected

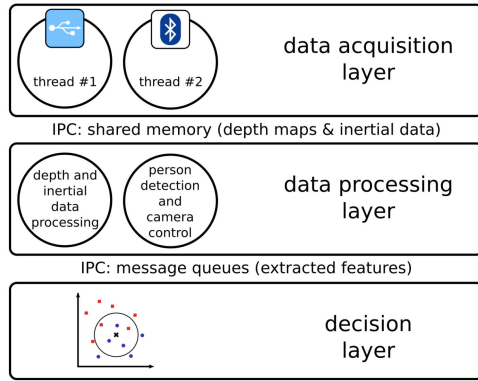


Fig. 6 Data acquisition, processing and communication between main processes

Table 2 Performance of lying pose recognition on frontal depth maps of URFD dataset

Estimated	True		
	Fall	No fall	
SVM, k-NN			
Fall	898	6	F1 score = 99.39%, accuracy = 99.55%
no fall	5	1516	
	FNR = 0.55%	FPR = 0.39%	

components, see lines 1–6 in Algorithm 1 (Fig. 2). These operations do not require a significant computational power. In the setting of the active camera, the person is extracted in every frame using the region growing, which is optionally followed by the person-finder. The extraction of the features is more time consuming due to the extraction of the floor and determining the point clouds. However, owing to event-driven data processing such computations are realised only in the case of the high likelihood of fall. The real-time data processing can be realised on a PC running Linux. The algorithms were also executed on PandaBoard in order to demonstrate their application potential and the possibility of running on low-cost computing boards.

The fall detection application uses an asynchronous message-driven communication model to propagate information throughout four application layers. It runs five main concurrent processes communicating via message queues, which are one of the interprocess communication mechanisms available under Linux and yield asynchronous communication among processes, see Fig. 6. In such a communication model a process usually referred to as the sender writes the generated messages to a queue, while one or more other receiver processes retrieve them from the queue. Once a message has been read, it is deleted by the kernel from the queue. This means that the sender and the recipient of the message do not have to cooperate with the queue at the same time. Even if several receivers are listening to a channel, each message can be retrieved by a single process only. As we can see in Fig. 6, the first process is amenable for acquiring motion data from the wearable device, the second process acquires the maps from the depth sensor, the third process extracts the person, the fourth process is accountable for processing of data and feature extraction, whereas the fifth process is accountable for data classification and triggering the fall alarm. In the case of the use of the PandaBoard, the dual-core processor allows the parallel execution of processing and acquisition processes. The IMU signals were collected with the frequency of 256 Hz and 12-bit resolution.

The algorithms for person detection use parameters that were determined experimentally. The algorithm for person extraction using the depth reference map requires the thresholds T_a and B_{th} , which were determined experimentally. The extraction performance does not drop significantly with the change of parameters mentioned above. The region growing is resistant to changes of experimentally determined δ_{th} values.

7 Experimental results

In this section, we present the experimental results that were obtained from the URFD dataset. In the subsequent subsections, we discuss evaluation results that were obtained from the publicly available URFD dataset using the presented fall detector along with the performance of the person detector and tracker.

7.1 Performance measures

The performance of the fall detector was evaluated with respect to accuracy, F1 score, fall-out (false positive rate – FPR) and miss rate (false negative rate – FNR). They were calculated as follows:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \times 100, \quad (4)$$

$$\text{F1 score} = \frac{2TP}{2TP + FP + FN} \times 100, \quad (5)$$

$$\text{FPR} = \frac{FP}{FP + TN} \times 100, \quad (6)$$

$$\text{FNR} = \frac{FN}{FN + TP} \times 100, \quad (7)$$

where TP stands for true positives (number of detected falls), FN denotes false negatives (number of undetected falls), FP indicates false positives (number of ADL examples giving false fall alarms), and TN stands for true negatives (number of ADL examples not giving fall alarms).

7.2 Evaluation of the fall detector

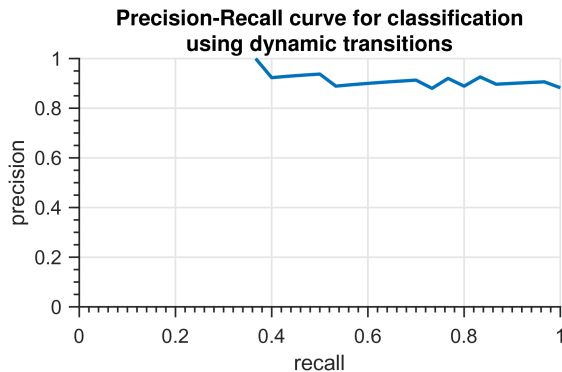
The system can be configured to perform fall detection on the basis of depth data only or both accelerometric and depth data. In the second option, thanks to indicating that person's movement is above some preset threshold, the detection performance is far better. It is superior due to the two-stage decision process, where at the first stage the fall-like activities are filtered out on the basis of acceleration data from the body-worn device, and depth map analysis is conducted on a subset of frames, which likely contain fall events. Thanks to the use of an accelerometer as an indicator of the potential fall the impact time can be determined easily and precisely enough, and thus the dynamic features (3) can be determined without considerable computational overheads.

7.2.1 Threshold selection: Keeping in mind that the accelerometer is used only to signal potential falls, the acceleration threshold was set to 2.5 to indicate all fall and fall-like activities as non-ADLs. The experimental results described below were obtained using the X-IMU accelerometer that was worn near the spine on the lower back, and which was attached to the body using an elastic belt around the waist.

7.2.2 Evaluation of the fall detector for facing camera: The lying pose detector on depth maps from the facing camera was evaluated on 2425 images from the URFD dataset (consisting of 903 lying poses). The detector has been trained on 1616 depth maps of which 602 were lying poses. The samples are selected so that the representative set included images that represent all the possible poses of lying person. We decided to base the system on the SVM and k-NN detectors since they exhibited high fall detection performance. They were built on features discussed in Section 5.3.1. The experimental results that were obtained by a linear SVM and a k-NN with five neighbours are shown in Table 2. Both classifiers gave identical results. The k-NN with three neighbours achieves slightly worse results. As we can notice the presented results are promising both in terms of FPR and FNR rates, which are important performance measures of fall detections systems [4].

Table 3 Performance of lying pose detection on overhead depth maps from URFD dataset [51]

Estimated	True		
	Fall	No Fall	
SVM			
Fall	244	9	F1 score = 97.41%, accuracy = 97.52%
no fall	4	268	
	FNR = 1.61%	FPR = 3.25%	
k-NN			
Fall	244	10	F1 score = 97.21%, accuracy = 97.33%
no fall	4	267	
	FNR = 1.61%	FPR = 3.61%	

**Fig. 7** Precision-recall curve for dynamic feature**Table 4** Performance of lying pose recognition on frontal and overhead depth maps of URFD dataset [%]

Estimated	True		
	Fall	No fall	
SVM, k-NN			
fall	903	0	F1 score = 100.00%, accuracy = 100.0%
no fall	0	1192	
	FNR = 0.00%	FPR = 0.00%	

Table 5 Performance of lying pose detection [%]

	Front. cam.	Top cam.	Front. + Top cam.
accuracy	99.55	97.52	100.0
F1 score	99.39	97.41	100.0
FPR	0.39	3.25	0.0
FNR	0.55	1.61	0.0

7.2.3 Evaluation of the fall detector for overhead camera: The algorithm for lying pose recognition in depth maps from the ceiling-mounted sensor has been evaluated on 875 representative images from the URFD dataset as well as additional images with ADLs. From the above mentioned dataset, a subset of 60% images was chosen for the testing (with 248 images for lying pose and 277 non-lying exemplars), whereas the remaining images were used only in training (with 165 lying and 185 non-lying exemplars). Since the discussed image sequences were acquired by a static depth sensor, the person has been extracted by differencing the depth maps from the depth background map. The discrimination between falls and ADLs has been performed by a linear SVM and a k-NN with five neighbours. The discussed classifiers operated on features discussed in Section 5.3.2. The classification performances, which were obtained by the classifiers mentioned above are shown in Table 3. As we can observe, the results obtained by lying pose detectors operating on features from the ceiling-mounted sensor are promising in terms of both the accuracy

and the miss rate, i.e. FNR. Such a low miss rate of lying pose detection implies that the detector has high sensitivity. We also evaluated the k-NN classifier with three neighbours, which gave slightly worse results. The C parameter of the SVM classifier has been set to a default value, i.e. to one. Experiments consisting of an exhaustive grid search over the parameter space to find the best setting demonstrated that the presented results do not change significantly for the C values differing from the default value.

Subsequently, we conducted evaluations in terms of the usefulness of the dynamic feature for distinguishing between fall and fall-like actions. They were evaluated in the context of improving the distinguishability between the accidental falling and the intentional lying on the floor. Fig. 7 depicts the precision-recall curve of the dynamic feature. It illustrates the classification performance of a binary classifier for different values of the discrimination threshold. The best accuracy was achieved for ΔT equal to 500 ms and threshold equal to 0.525.

Afterwards, we asked two volunteers to act as evaluators of the dynamic feature. It turned out that a simple cascade classifier consisting of the lying pose detector and the dynamic transition detector performs very well in practice as it has an almost null false alarm ratio. In particular, the cascade gives promising results if a moment in which there was the impact is determined precisely. In consequence, the fall detector consisting of an accelerometer-based fall indicator, a depth map-based lying pose detector and a dynamic transition detector achieves the best results for the overhead camera. It is worth noting that such a classifier detected properly all fall in the depth map sequences from the URFD dataset. However, the discussed classifier was not able to completely eliminate the false alarms. The inability to eliminate false alarms by fall detectors working on images from a single camera motivated us to elaborate a fall detector using images both from facing and overhead cameras.

7.2.4 Evaluation of the fall detector using data from facing and overhead cameras: Using the features extracted from both facing and overhead cameras we evaluated a fall detector consisting of an accelerometer-based fall indicator and a k-NN classifier with five neighbours for lying pose detection. The classifier operates on features discussed in Sections 5.3.1 and 5.3.2. It has been trained on features extracted from the maps, which were used in Section 7.2.3, plus the corresponding depth images from the facing sensor. The results are presented in Table 4. As we can observe, the results obtained from the discussed camera setting are superior in comparison with results presented previously.

7.2.5 Comparison of performance of fall detection: Table 5 summarises the classification performances, which were obtained using the discussed camera settings. The camera setup with the frontal camera gives slightly better results in comparison with the setup with the top camera. The best results were obtained using data from both cameras.

7.3 Evaluation of person detector and tracker

If the system is configured to work with a static ceiling-mounted camera, the person can be extracted at a low-computational cost by differencing the current depth map from the continuously updated depth reference map of the scene. On the PandaBoard this operation takes about 10 ms. However, as we already mentioned, the ceiling-mounted and fixed Kinect has a quite limited observation area. By the use of a ceiling-mounted motorised head to rotate the Kinect, the observation area can be expanded noticeably. In such a setup with a pan-tilt sensor, more sophisticated and time consuming algorithms are required to extract the subject and to follow it by the active camera.

We began with an evaluation of the depth region growing in the depth maps from the URFD dataset. It is worth noting that in all sequences the person was extracted correctly. This means that the presented person tracker achieved perfect performance on all fall events registered in URFD. Next, our region growing was examined on five depth map sequences that were acquired by the pan-tilt camera. The subject was moving freely around the room in

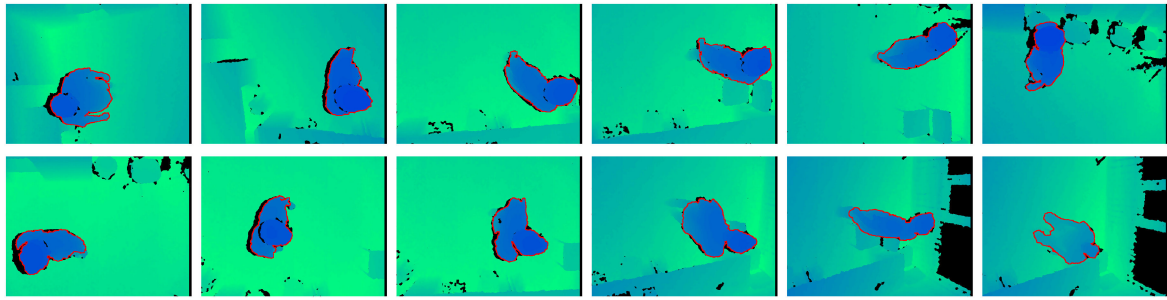


Fig. 8 Region growing-based person detection and tracking on depth maps acquired by a ceiling-mounted active camera

Table 6 Performance of person detection using HOD-SVM on depth maps from ceiling-mounted active camera [%]

	Accuracy	F1 score	FPR	FNR
rotat.	99.45	99.10	0.78	0.00
no rotat.	98.91	98.18	0.78	1.82

an area that could not be observed by the fixed camera. In the discussed experiments, it was required not only to extract the person in real-time but additionally to keep he/she in the central part of the depth map acquired by the active camera. Before starting the delineation of person, the camera was stationary for a while to initially extract the subject through differencing the current depth map from the depth reference map of the scene. It is worth mentioning that in all frames acquired by the active camera, including depth maps containing intentional falls, every main body part has been extracted properly. Fig. 8 presents illustrative results, which were achieved on depth images collected by the active camera. On the PandaBoard, the time needed for region growing-based person delineation depends on the blob size and ranges from 17 to about 25 ms. A video illustrating the person tracked by a pan-tilt depth camera is available under the following link: <http://fenix.univ.rzeszow.pl/~mkepki/demo/act.mp4>.

The person detector, which is discussed in Section 4.3 has been evaluated on 254 positive examples and 638 negative examples of which 60% were used for training and the remaining 40% for testing. The depth maps with the delineated person were rotated to a canonical pose on the basis of the axis of the person's blob. They were also scaled according to the distance of his/her head to the camera. Table 6 contains experimental results that were obtained using the HOD-SVM detector, see also [47, 52]. As we can observe, the detector achieves better performance when the rotation of silhouettes to the canonical pose takes place. On the other hand, the improvement in the performance is not considerable, and this, in turn, indicates that the algorithm is fairly resistant to variations in head poses. This is due to the extraction of the HOD features on gradients forming elliptical like structures on the person's head seen on depth maps from an overhead camera. The discussed results were achieved using the HOD with the cell size equal to 8×8 . The time needed for person detection on the PandaBoard is equal to 41 ms.

8 Discussion

Both camera settings have advantages and disadvantages. As we have already mentioned, ceiling-mounted depth sensors are rarely utilised in fall detection research [8, 38]. One of the reasons for this is the limited monitoring area of depth sensors, including Kinect. Our experimental results demonstrate that owing to mounting the depth sensor on a motorised pan-tilt unit the observation field can be extended considerably. As a result, typical senior rooms of size up to 25 m^2 can be monitored by a single depth sensor. Moreover, we found that person delineation in such a setting with an active camera can be done quite reliably and fast. Our modified depth-region growing for a person extraction demonstrated value in several experiments. The person detection times on the low-cost PandaBoard are close to times needed for real-time processing. The computational power of current PCs is sufficient to execute in real-time our algorithms for person detection and fall recognition. One

of the most significant obstacles to the introduction of fall detection systems and their acceptance by seniors is the barrier of costs of devices and their everyday use [8]. In this context, it is worth noting that our system for fall detection can be built relatively inexpensively, using a low-cost depth camera, a wireless accelerometer and low-cost processing boards like PandaBoard.

One of the advantages of the setting with the active ceiling-mounted camera is that the number of situations in which occlusions impede person extraction is much smaller in comparison with setup with a facing camera. In this context, it is worth noting that there is almost no significant work that deals with fall detection in the case of visual occlusions [8, 13]. Another observation is that for such a camera setting a similar performance of lying pose detection can be obtained with a smaller number yet more discriminative features in comparison with depth map features that are needed for reliable fall detection using a facing camera. Moreover, in a setup with the ceiling-mounted camera the dynamical features, which demonstrated high discrimination power, can be computed quite easy, particularly if a body-worn accelerometer is used. We showed experimentally that a two-camera system achieves perfect detection performance on data from the freely available URFD dataset.

As we have already mentioned, our work differs from the relevant work since we focus on a ceiling-mounted pan-tilt depth sensor, and last but not least, in that we are using a body-worn accelerometer to indicate the context of the event. In this way, an expert knowledge about the specificity of the fall detection problem has been realised in the form of event-driven architecture and a cascade of classifiers. In consequence, a decision about the fall is not taken by a single classifier, trained using a machine learning technique, but it is taken on the basis of carefully designed and evaluated classifiers, which attempt to mimic human experts. Our conclusions are in line with the research findings presented in [18, 27] in that the motion patterns of real-falls might differ from simulated falls, particularly if the falling person is trying to save himself from falling in order to minimise the effects of the fall. The event-based approach to fall detection does not introduce unobtrusiveness due to the possible use of suit-integrated accelerometers, which employ the human body motions to continuously recharge the battery [53].

9 Conclusions

In this work, efficient algorithms for fall detection were developed, implemented and tested using depth map sequences and a wireless inertial sensor worn by a monitored person. A set of descriptors for depth maps has been proposed to permit the classification of person poses as well as his/her actions. The experimental validation was carried out on prepared and then shared data repository consisting of synchronised depth and accelerometric data. Extensive experiments and tests were conducted in the scenario with a static camera facing the scene and an active camera observing the scene from the above. The algorithms were designed with regard to low computational demands and the possibility of their run on ARM platforms. Several experiments consisting of person detection, tracking and fall detection in real-time were carried out to show efficiency and reliability of the proposed solutions. Both camera settings were compared in terms of person detection and fall recognition. The experimental results showed that the developed algorithms for fall detection have both low FPR

and FNR rates. In future work, we will investigate scenarios with two persons as well as scenarios with occlusions. The region growing-based person detection should be extended to deal with activities like sleeping on a coach.

10 Acknowledgment

This work was supported by the Polish National Science Center (NCN) under a research grant 2014/15/B/ST6/02808.

11 References

- [1] Rashidi, P., Mihailidis, A.: 'A survey on ambient-assisted living tools for older adults', *IEEE J. Biomed. Health Inform.*, 2013, **17**, (3), pp. 579–590
- [2] Friedewald, M., Raabe, O.: 'Ubiquitous computing: an overview of technology impacts', *Telemat. Inf.*, 2011, **28**, (2), pp. 55–65
- [3] Acampora, G., Cook, D., Rashidi, P., *et al.*: 'A survey on ambient intelligence in healthcare', *Proc. IEEE*, 2013, **101**, (12), pp. 2470–2494
- [4] Igual, R., Medrano, C., Plaza, I.: 'Challenges, issues and trends in fall detection systems', *Biomed. Eng. Online*, 2013, **12**
- [5] Murphy, J., Isaacs, B.: 'The post-fall syndrome: a study of 36 elderly patients', *Gerontology*, 1982, **28**, pp. 265–270
- [6] Tinetti, M.E.: 'Predictors and prognosis of inability to get up after falls among elderly persons', *JAMA J. Am. Med. Assoc.*, 1993, **269**, (1), p. 65
- [7] Noury, N., Rumeau, P., Bourke, A., *et al.*: 'A proposal for the classification and evaluation of fall detectors', *IRBM*, 2008, **29**, (6), pp. 340–349
- [8] Hamm, J., Money, A.G., Atwal, A., *et al.*: 'Fall prevention intervention technologies: a conceptual framework and survey of the state of the art', *J. Biomed. Inf.*, 2016, **59**, pp. 319–345
- [9] Litvak, D., Zigel, Y., Gannot, I.: 'Fall detection of elderly through floor vibrations and sound', Int. Conf. IEEE Engineering in Medicine and Biology Society, 2008, pp. 4632–4635
- [10] Zhang, Z., Kapoor, U., Narayanan, M., *et al.*: 'Design of an unobtrusive wireless sensor network for nighttime falls detection', Annual Int. Conf. IEEE Engineering in Medicine and Biology Society, August 2011, pp. 5275–5278
- [11] Morawski, R.Z., Yashchysyn, Y., Pirek, M., *et al.*: 'Monitoring of human movements by means of impulse-radar sensors', *Telecommun. Rev. Telecommun. News*, 2015, **6**, pp. 598–602
- [12] Yu, M., Naqvi, S.M., Rhuma, A., *et al.*: 'One class boundary method classifiers for application in a video-based fall detection system', *IET Comput. Vis.*, 2012, **6**, (2), pp. 90–100
- [13] Zhang, Z., Conly, C., Athitsos, V.: 'A survey on vision-based fall detection', Proc. 8th ACM Int. Conf. Pervasive Technologies Related to Assistive Environments, 2015, pp. 1–7
- [14] Sathyanarayana, S., Satzoda, K.R., Sathyanarayana, S., *et al.*: 'Vision-based patient monitoring: a comprehensive review of algorithms and technologies', *J. Ambient Intell. Humanized Comput.*, 2015, pp. 1–27
- [15] Wang, Y., Winters, J.: 'An event-driven dynamic recurrent neuro-fuzzy system for adaptive prognosis in rehabilitation', Int. Conf. IEEE Engineering in Medicine and Biology Society, 2003, pp. II:1256–II:1259
- [16] Rea, F., Metta, G., Bartolozzi, C.: 'Event-driven visual attention for the humanoid robot icub', *Front. Neurosci.*, 2013, **7**, p. 234
- [17] Perez-Carrasco, J., Zhao, B., Serrano, C., *et al.*: 'Mapping from frame-driven to frame-free event-driven vision systems by low-rate rate coding and coincidence processing-application to feedforward convnets', *IEEE Trans. Pattern Anal. Mach. Intell.*, 2013, **35**, (11), pp. 2706–2719
- [18] Bagala, F., Becker, C., Cappello, A., *et al.*: 'Evaluation of accelerometer-based fall detection algorithms on real-world falls', *PLoS ONE*, 2012, **7**, (5), p. e37062
- [19] Kangas, M., Korpelainen, R., Vikman, I., *et al.*: 'Sensitivity and false alarm rate of a fall sensor in long-term fall detection in the elderly', *Gerontology*, 2015, **61**, (1), pp. 61–68
- [20] Chen, C., Jafari, R., Kehtarnavaz, N.: 'A survey of depth and inertial sensor fusion for human action recognition', *Multimedia Tools Appl.*, 2017, **76**, (3), pp. 4405–4425
- [21] Hoflinger, F., Muller, J., Zhang, R., *et al.*: 'A wireless micro inertial measurement unit (IMU)', *IEEE Trans. Instrum. Meas.*, 2013, **62**, (9), pp. 2583–2595
- [22] Kangas, M., Konttila, A., Lindgren, P., *et al.*: 'Comparison of low-complexity fall detection algorithms for body attached accelerometers', *Gait Posture*, 2008, **28**, (2), pp. 285–291
- [23] Li, Q., Stankovic, J., Hanson, M., *et al.*: 'Accurate, fast fall detection using gyroscopes and accelerometer-derived posture information', Sixth Int. Workshop on Wearable and Implantable Body Sensor Networks, June 2009, pp. 138–143
- [24] Jacob, J., Nguyen, T., Lie, D., *et al.*: 'A fall detection study on the sensors placement location and a rule-based multi-thresholds algorithm using both accelerometer and gyroscopes', IEEE Int. Conf. on Fuzzy Systems, 2011, pp. 666–671
- [25] Bourke, A., O'Brien, J., Lyons, G.: 'Evaluation of a threshold-based tri-axial accelerometer fall detection algorithm', *Gait Posture*, 2007, **26**, (2), pp. 194–199
- [26] Lustrek, M., Gjoreski, H., Kozina, S., *et al.*: 'Detecting falls with location sensors and accelerometers', AAAI Conf. on Artificial Intelligence, 2011, pp. 1662–1667
- [27] Kangas, M., Vikman, I., Nyberg, L., *et al.*: 'Comparison of real-life accidental falls in older people with experimental falls in middle-aged test subjects', *Gait Posture*, 2012, **35**, (3), pp. 500–505
- [28] Kepski, M., Kwolek, B.: 'Fall detection on embedded platform using kinect and wireless accelerometer', Int. Conf. on Computers Helping People with Special Needs, 2012, pp. II:407–II:414
- [29] Ma, X., Wang, H., Xue, B., *et al.*: 'Depth-based human fall detection via shape features and improved extreme learning machine', *IEEE J. Biomed. Health Inf.*, 2014, **18**, (6), pp. 1915–1922
- [30] Chen, C., Jafari, R., Kehtarnavaz, N.: 'Improving human action recognition using fusion of depth camera and inertial sensors', *IEEE Trans. Human-Mach. Syst.*, 2015, **45**, (1), pp. 51–61
- [31] Chen, L., Wei, H., Ferryman, J.: 'A survey of human motion analysis using depth imagery', *Pattern Recogn. Lett.*, 2013, **34**, (15), pp. 1995–2006
- [32] Yu, M., Rhuma, A., Naqvi, S., *et al.*: 'A posture recognition-based fall detection system for monitoring an elderly person in a smart home environment', *IEEE Trans. Inf. Technol. Biomed.*, 2012, **16**, (6), pp. 1274–1286
- [33] Wu, C., Aghajan, H.: 'Real-time human pose estimation: a case study in algorithm design for smart camera networks', *Proc. IEEE*, 2008, **96**, (10), pp. 1715–1732
- [34] Rougier, C., Meunier, J., St-Arnaud, A., *et al.*: '3D head tracking for fall detection using a single calibrated camera', *Image Vision Comput.*, 2013, **31**, (3), pp. 246–254
- [35] Demiroz, L.A.B.E., Salah, A.A.: 'Coupling fall detection and tracking in omnidirectional cameras', International Workshop on Human Behavior Understanding, 2014 (LNCS, **8749**), pp. 73–85
- [36] Sokolova, M.V., Serrano-Cuerda, J., Castillo, J.C., *et al.*: 'A fuzzy model for human fall detection in infrared video', *J. Intell. Fuzzy Syst.*, 2013, **24**, (2), pp. 215–228
- [37] Rougier, C., Auvinet, E., Rousseau, J., *et al.*: 'Fall detection from depth map video sequences', Int. Conf. on Smart Homes and Health Telematics, 2011 (LNCS, **6719**), pp. 121–128
- [38] Webster, D., Celik, O.: 'Systematic review of kinect applications in elderly care and stroke rehabilitation', *J. NeuroEng. Rehabil.*, 2014, **11**
- [39] Cuppens, K., Chen, C.-W., Wong, K., *et al.*: 'Integrating video and accelerometer signals for nocturnal epileptic seizure detection', Int. Conf. on Multimodal Interaction, 2012, pp. 161–164
- [40] Shotton, J., Sharp, T., Kipman, A., *et al.*: 'Real-time human pose recognition in parts from single depth images', *Commun. ACM*, 2013, **56**, pp. 116–124
- [41] Kwolek, B., Kepski, M.: 'Improving fall detection by the use of depth sensor and accelerometer', *Neurocomputing*, 2015, **168**, pp. 637–645
- [42] Kepski, M., Kwolek, B.: 'Fall detection using ceiling-mounted 3d depth camera', Int. Conf. on Computer Vision Theory and Appl. (VISAPP), vol. 2, pp. 640–647
- [43] Raguram, R., Frahm, J.-M., Pollefeys, M.: 'A comparative analysis of RANSAC techniques leading to adaptive real-time random sample consensus', Proc. 10th European Conf. on Computer Vision: Part II, 2008, pp. 500–513
- [44] Adams, R., Bischof, L.: 'Seeded region growing', *IEEE Trans. Pattern Anal. Mach. Intell.*, 1994, **16**, (6), pp. 641–647
- [45] Mehnert, A., Jackway, P.: 'An improved seeded region growing algorithm', *Pattern Recogn. Lett.*, 1997, **18**, (10), pp. 1065–1071
- [46] Spinello, L., Arras, K.: 'People detection in RGB-D data', IEEE/RSJ Int. Conf. on Intelligent Robots and Systems (IROS), September 2011, pp. 3838–3843
- [47] Kepski, M., Kwolek, B.: 'Fall detection using body-worn accelerometer and depth maps acquired by active camera', Proc. 11th Int. Conf. on Hybrid Artificial Intelligent Systems, 2016, pp. 414–426
- [48] Chen, J., Kwong, K., Chang, D., *et al.*: 'Wearable sensors for reliable fall detection', IEEE Int. Conf. on Engineering in Medicine and Biology Society, 2005, pp. 3551–3554
- [49] Kangas, M., Konttila, A., Winblad, I., *et al.*: 'Determination of simple thresholds for accelerometry-based parameters for fall detection', 29th Annual Int. Conf. on IEEE Engineering in Medicine and Biology Society, 2007, pp. 1367–1370
- [50] Horn, B.: 'Robot vision' (The MIT Press, Cambridge, MA, 1986)
- [51] Kwolek, B., Kepski, M.: 'Fall detection using kinect sensor and fall energy image', Int. Conf. on Hybrid artificial intelligent systems, 2013 (LNCS, **8073**), pp. 294–303
- [52] Kepski, M., Kwolek, B.: 'Embedded system for fall detection using body-worn accelerometer and depth sensor', IEEE 8th Int. Conf. Intell. Data Acquisition and Advanced Computing Systems, vol. 2, 2015, pp. 755–759
- [53] Jung, S., Hong, S., Kim, J., *et al.*: 'Wearable fall detector using integrated sensors and energy devices', *Sci. Rep.*, 2015, **5**, Article number: 17081. Available at: <http://www.nature.com/articles/srep17081>