

Gesture spotting with body-worn inertial sensors to detect user activities

Holger Junker^a, Oliver Amft^{a,*}, Paul Lukowicz^b, Gerhard Tröster^a

^a*Wearable Computing Lab., ETH Zurich, Gloriastrasse 35, 8092 Zurich, Switzerland*

^b*Embedded Systems Group, University of Passau, Innstrasse 43, 94032 Passau, Germany*

Received 10 January 2007; received in revised form 15 November 2007; accepted 19 November 2007

Abstract

We present a method for spotting sporadically occurring gestures in a continuous data stream from body-worn inertial sensors. Our method is based on a natural partitioning of continuous sensor signals and uses a two-stage approach for the spotting task. In a first stage, signal sections likely to contain specific motion events are preselected using a simple similarity search. Those preselected sections are then further classified in a second stage, exploiting the recognition capabilities of hidden Markov models. Based on two case studies, we discuss implementation details of our approach and show that it is a feasible strategy for the spotting of various types of motion events.

© 2007 Elsevier Ltd. All rights reserved.

Keywords: Natural gesture segmentation; Gesture spotting; Activity recognition; Automatic dietary monitoring; Event detection

1. Introduction

Monitoring and classification of human activity using simple body-worn sensors is emerging as an important research area in machine learning. Activity monitoring itself is motivated by a variety of mobile and ubiquitous computing applications, such as personalisation of the user interface, behavioural monitoring in medicine, medication assessment, assistive systems for the elderly and cognitively disabled or intelligent information delivery and recording systems for industrial assembly and maintenance.

The choice of simple sensors, such as accelerometers instead of computer vision, stems from the limited computational resources of mobile and ubiquitous systems and the very diversified, dynamic environment in which such systems need to operate. The latter often implies varying light conditions, changing backgrounds and a large clutter. This makes extracting relevant information from visual signals difficult and computationally intensive. Body-mounted motion sensors, on the other hand, are influenced by user activity only. The problem with activity recognition using such sensors lies less in the extraction of relevant features than in the fact that the information is

often ambiguous and incomplete. Thus, once a vision system has managed to track, for example, the user's arm, relatively exact trajectories could be obtained for activity recognition. In contrast, arm worn accelerometers react to a combination of earth gravity and arm speed changes. Gyroscopes describe rotational motions of the arm. However, none of the above provides exact trajectory information.

Despite the disadvantages listed above, body-worn motion sensors have been successfully used for a variety of tasks (see related work). One area where little progress has been made so far, is the spotting of sporadically occurring activities in a continuous data stream. This is known to be difficult, even if complete trajectory information is available from a vision system. It is even more difficult in a wearable sensor-based environment.

This paper describes a novel method for tackling this problem based on appropriately adapted machine learning techniques. Focusing on activities associated with distinct arm gestures, the performance of the proposed method is evaluated in two elaborate case studies.

1.1. Paper scope and background

Depending on the specific application, very different types of activity recognition are needed. As an example, consider a

* Corresponding author. Tel.: +41 44 632 5936; fax: +41 44 632 1210.
E-mail address: amft@ife.ee.ethz.ch (O. Amft).

system designed to monitor the overall physical activity level of a person. The idea behind such systems is to provide general information about the effect of certain behavioural recommendations or to estimate energy expenditure without having the patient admitted to stationary care or a laboratory for observation. A wearable system deploying appropriate body-worn sensors can be used to collect this data. Obviously, the type of information that such systems need to deliver is not about single, specific actions, but more about the overall level of activity. Often, the activity level can be assessed by averaging parameters, such as mean acceleration of specific body parts. In a way, this is a very simple form of activity recognition.

On the other side of the spectrum are applications, where reliable recognition on a more fine-grained level is needed. Such applications may include, e.g. the monitoring of specific tasks and/or movements in a rehabilitation scenario, the spotting of specific gestures for novel, more natural human–computer interface or the classification of dietary intake gestures for an automated nutrition monitoring system. Such recognition tasks are particularly difficult, because the relevant activities occur sporadically in between a large variety of other activities. For example, in between the actual activity a user might fetch tools, drink, chat with another person or just scratch the head. As a consequence, the task at hand can be described as *activity spotting*. It is widely recognised as a particularly complex domain of activity recognition and is still an open problem.

The work described in this paper is part of a larger effort of our groups, directed at this problem, e.g. Refs. [1–3]. It focuses on activities that are associated with a characteristic arm gestures. For such activities, the paper presents a novel gesture spotting method based on arm-worn motion sensors. The method uses a natural partitioning of human motions. In order to achieve a balance between precision and recall with reasonable computational effort, the task is partitioned into a fast highly sensitive stage to pick up potentially interesting signal segments and a more complex, highly selective second stage to narrow down the selection and get rid of false positives.

Our method is primarily intended as part of a large activity spotting system that uses additional information such as location, modes of locomotion, e.g. sitting standing, walking [4], supplementary location sensors [5] or information on objects involved in the activity [6]. Nonetheless, we present experiments on activities from two different everyday life domains indicating that even on its own our method achieves reasonable performance.

1.2. Related work

In contrast to isolated motion recognition that has been shown in various areas, the spotting task is much more challenging. The difficulty of spotting specific human motion events stems from a number of sources. These include, among others, co-articulation, where consecutive gestures influence each other [7], as well as intra- and inter-person variability. Another challenge, the system has to deal with, is the fact that the motion events to be spotted may only occur sporadically, in a continuous data stream, while at the same time being

embedded into other, partly arbitrary movements (called *zero-class*). These movements, however, are inherently difficult to model, due to their complexity and unpredictability. As a consequence, conventional recognition schemes for continuous classification, such as hidden Markov models (HMMs), are not directly applicable for our recognition task, since they rely on appropriate zero-class models. Consequently, we cannot take advantage of the implicit data segmentation capabilities that HMMs provide. Moreover, we have to deal with the fact that motion events are typically very short. This means that for any explicit segmentation-based recognition, exact localisation of event boundaries is important.

The recognition of gestures has been studied extensively over years and many approaches have been proposed to tackle the diverse problems. In general, these approaches can be broadly categorised in either of the two following categories: gesture recognition, requiring external infrastructure and gesture recognition, focusing on wearable instrumentation.

The first category is dominated by vision-based motion recognition, using a single or multiple cameras. While an exhaustive review of literature is beyond the scope of this work, we exemplarily indicate related works. Starner [8] proposed an approach for American Sign Language recognition, Campbell and Bobick [9] developed a system for recognising classical ballet steps, Yamato et al. [10] worked on the recognition of different tennis strokes, Brand et al. [11] targeted T'ai Chi movements, Lee and Kim [12] dealt with typical gestures for interacting with a computer and Rao and Shah [13] aimed at manipulative gestures. Further literature on vision-based motion capture and recognition can be found in Refs. [14–16].

More recently the use of wearable instrumentation for gesture recognition has gained much attention mainly due to the success in sensor miniaturisation. Various approaches dealing with the recognition of activities or events have been presented. Chambers et al. [17] targeted Kung Fu moves and Benbasat [18] focused on the recognition of “atomic” gestures. Kern et al. [19] looked at activities, such as keyboard typing, writing on a white-board and shaking hands. Cakmakci et al. [20] tried to identify when a person was looking at the watch. Bao [21] aimed at typical household activities including vacuuming, folding laundry, watching TV or brushing teeth. Lukowicz et al. [22] concentrated on workshop activities including sawing, hammering, drilling and filing. Brashear et al. [23] dealt with gestures for American Sign Language and Lementec and Bajcsy [24] worked on the recognition of gestures used to instruct pilots after landing.

Although many motion recognition approaches exist, few are dealing with the spotting task itself. Deng and Tsui [25] proposed a method for spotting gestures in continuous data. Their approach makes use of an HMM-based accumulation score that supports endpoint detection of a particular gesture in a continuous data stream. Based on a potential endpoint their algorithm searches for a corresponding start point using the viterbi algorithm. While this approach seems promising, it has been evaluated solely for the recognition of two-dimensional trajectories (Arabic numbers). Lee and Kim [12] developed a method deploying HMMs directly, to spot gestures in a

continuous stream of sensor data. They introduced the concept of a threshold model that calculates the likelihood threshold of an input pattern and provides a confirmation mechanism for the provisionally matched gesture patterns. The threshold model is a weak model for all trained gestures and is constructed from all existing gesture models. Lukowicz et al. [22] demonstrated continuous, online motion recognition by partitioning the incoming data using an intensity analysis based on the signals of two microphones exploiting the fact that the movements to be recognised are accompanied with a particular sound.

While the first two approaches made use of the implicit segmentation capabilities of HMMs, the third approach used an explicit segmentation step to facilitate spotting. We believe that explicit gesture segmentation can be very helpful and efficient to facilitate the spotting task. Lee and Yangsheng [26] developed a system for online gesture recognition using HMMs. They were among the first researchers to use segmentation as a pre-processing step to gesture recognition and were able to recognise 14 different gestures online. While they proposed acceleration thresholds for segmentation, they used a simple velocity-based segmentation relying on the fact that there must be short pauses between two consecutive gestures. They successfully demonstrated good recognition performance for the trained gestures; however, they did not deal with the rejection of non-relevant movements. Kahol et al. [27] proposed a gesture segmentation algorithm which employs a hierarchical layered structure to represent the human anatomy. The algorithm used low-level motion parameters to characterise motion in the various layers of this hierarchy and was able to predict segmentation boundaries based on profiles, generated from segmentation results. The segmentation, in turn, was provided by observers, who manually segmented training data. In a recent work, Kahol et al. [28] used the concept to fully document every motion in dance activities using a Vicon camera system. Wang et al. [29] presented an approach for automatically segmenting sequences of natural activities into atomic sections and clustering them. The segmentation was based on finding the local minimum of velocity and local maximum of change in direction. The minimum below and the maximum above the certain threshold were selected as segment points. The limitation of their approach is that it can only segment and label continuous human gestures, but not spot them. Liang and Ouhyoung [30] used a temporal segmentation based on the discontinuity of the movements according to four gesture parameters and HMMs to perform real-time continuous gesture recognition of sign language. Their approach allows the recognition of gestures that were defined in vocabularies only; thus rejection of non-gesture patterns is not considered. Morguet [31] proposed a two-step approach to the continuous recognition of gestures in video sequences. In a first step, a simple segmentation algorithm was used to identify start and endpoints of potentially meaningful segments. This segmentation process used a threshold on a specific motion parameter in conjunction with simple rules to obtain valid segments. These segments were then classified in isolation. However, this approach cannot reject non-gesture patterns that are falsely retrieved in the first stage.

1.3. Paper contributions and organisation

As stated in the introduction, the work presented in this paper is part of a large effort towards reliable spotting of complex activities using simple on-body sensors. It focuses on the recognition of gestures that build the basis for the inference of more abstract activities. The primary aim is to support complex activity spotting systems rather than to develop an activity spotting system based solely on arm gestures. Nonetheless, we show how, for suitable domains, good performance can be achieved without any additional information.

Within this scope the paper makes the following contributions:

- (1) It presents a novel, two-stage gesture spotting method based on body-worn motion sensors. The method is specifically designed towards the needs and constraints of activity recognition in wearable and pervasive systems. This includes a large null class, lack of appropriate models for the null class, large variability in the way gestures are performed and a variable gesture length. It also refrains from excessively computationally intensive operations such as correlations over large data sets or complex searches. Instead, it uses a natural partitioning of human motions, combined with a simple parametrisation scheme as a computationally cheap preselection stage (PS) that identifies potentially interesting data sections. These sections are then reevaluated using HMMs to reduce the number of classification errors. This combination of a cheap, highly sensitive initial stage with a highly selective second stage is what makes our approach unique and well suited to the intended domain.
- (2) The paper describes the verification of the proposed method on two scenarios that together comprise nearly a thousand relevant gestures. The first one, interaction with different everyday objects, is part of a wide range of wearable systems applications. The second one, nutrition intake, is a highly specialised application motivated by the needs of a large industry dominated health monitoring project. In both cases studies we arrive at recall values between 80% and 90% and a precision of over 70%. The significance of these case studies and results are twofold. First, they confirm the soundness of our approach. Second they are a strong indication for the feasibility of reliable activity spotting using wearable sensors, in particular, since the approach presented in this paper is meant to be used as part of a large system that uses other information to further improve the results.

As indicated in the related work section, two-stage activity spotting approaches have been tried before. However, to our knowledge, the specific approach described in this paper, with its focus on the peculiarities of activity spotting using simple sensors and wearable systems is novel. Taking into account the results achieved in our case studies it represents a significant contribution to the field.

1.3.1. Paper organisation

In Section 2 of the paper, we introduce our two-stage spotting approach. In Section 3, we describe the case studies used to validate our approach. In Section 4, we focus on implementation details of the spotting algorithm and in Section 5, we detail the experimental setup to acquire sensor data for the case studies. In Section 6, we finally present our evaluation results. In Sections 7 and 8, we discuss the results and highlight future work, respectively.

2. Spotting approach

Our two-stage spotting approach consists of a PS stage (first stage) and a classification stage (CS) (second stage) as shown in Fig. 1. The task of the PS stage is to localise and preselect sections in the continuous signal stream, likely to contain relevant motion events. These candidate sections are then passed on to the CS and are classified in isolation using appropriate classifiers.

The preselection of sections in a continuous signal stream can be considered as a search problem. In a naive approach, the search may be performed on all possible sections in the data stream. The major problem of this exhaustive approach is its computational effort. To reliably capture human motions with inertial sensors, the sensors are usually sampled with up to 100 Hz. Considering that a relevant motion event may take several seconds, the above-mentioned search strategy would require to check a large number of sections.

Obviously, one solution to reduce the complexity is to apply a coarse search, where not all but only certain sections in the continuous signal stream are considered for the search. One way to implement such a coarse search is to partition the signal stream into segments which are significantly longer than a single sampling interval and to consider the segment boundaries as possible start/endpoints of the sections to be searched. However, an artificial partitioning is likely to miss the exact boundaries of the relevant motion events contained in the data stream. This makes the recognition more complicated, since sections may contain only parts of the relevant motion event as well as other motions.

We propose to use a natural partitioning of the data into “motion segments”. Inspired by the taxonomy of Bobick [32] these motion segments are described as non-overlapping, atomic units of human movement, characterised by their spatio-temporal trajectories. Assuming that a motion event can be subdivided into a sequence of motion segments, we can obtain a natural, non-ambiguous partitioning of the overall motion with the start and end of the motion events corresponding to the start/end of a specific motion segment. Thus, the search can be constrained to those sections, whose boundaries coincide with the boundaries

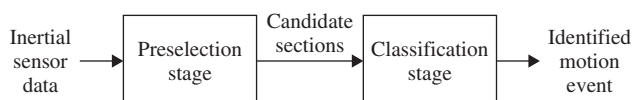


Fig. 1. Sensor data flow through the two-stage recognition framework.

Table 1
Applied terminology of human motion in this work

Term	Description
Motion segment	Represents atomic, non-overlapping unit of human motion that can be characterised by their spatio-temporal trajectory
Motion event	Span a sequence of motion segments. A gesture can be considered as a particular class of motion events, mainly involving movements of the arms and trunk
Activity	Describes a situation that may consist of various motion events. Thus, it refers to higher-level context
Signal segment	A slice of sensor data that corresponds to a motion segment
Candidate section	A slice of sensor data that may contain a gesture

of the motion segments. Table 1 summarises the terminology used in this work.

For the partitioning task, motion parameter(s) were used that represent the motion event closely. The number and types of motion parameters to be used is specific to the motion event to be recognised. For arm-related motion events, they include, e.g. relative orientation information, such as joint angles between the lower and the upper arm, absolute orientation information of the arm segments to an earth-fixed reference frame or simply the raw signals from the sensors attached to the arm segments.

While the PS stage identifies potential candidate sections, the classification stage is used to eliminate those sections that have been falsely retrieved in the PS stage. This is achieved by individually classifying the candidate sections using HMMs and comparing the classification result to the result of the PS stage.

The main motivation behind this two-stage approach is to reduce the complexity of the spotting task, by constraining the search space within the continuous data stream and by applying a simple similarity analysis to preselect potential sections. The subsequent CS is used to make the recognition more robust and retain only relevant sections.

3. Case studies

In order to discuss the implementation of our approach, we considered the spotting of typical, everyday life gestures in a continuous data stream from body-worn inertial sensors. Specifically, we investigated two different case studies.

Case study 1 deals with the spotting of diverse object interaction gestures, reflecting common activities of daily living. The detection of such gestures is considered as key component in a context recognition system, to monitor complex human activities. Furthermore, such gestures may facilitate more natural human–computer interfaces.

Case study 2 focuses on dietary intake gestures. The spotting of body motions related to food intake is expected to become one sensing domain of an automated dietary monitoring system [33]. Although the automatic determination of exact type

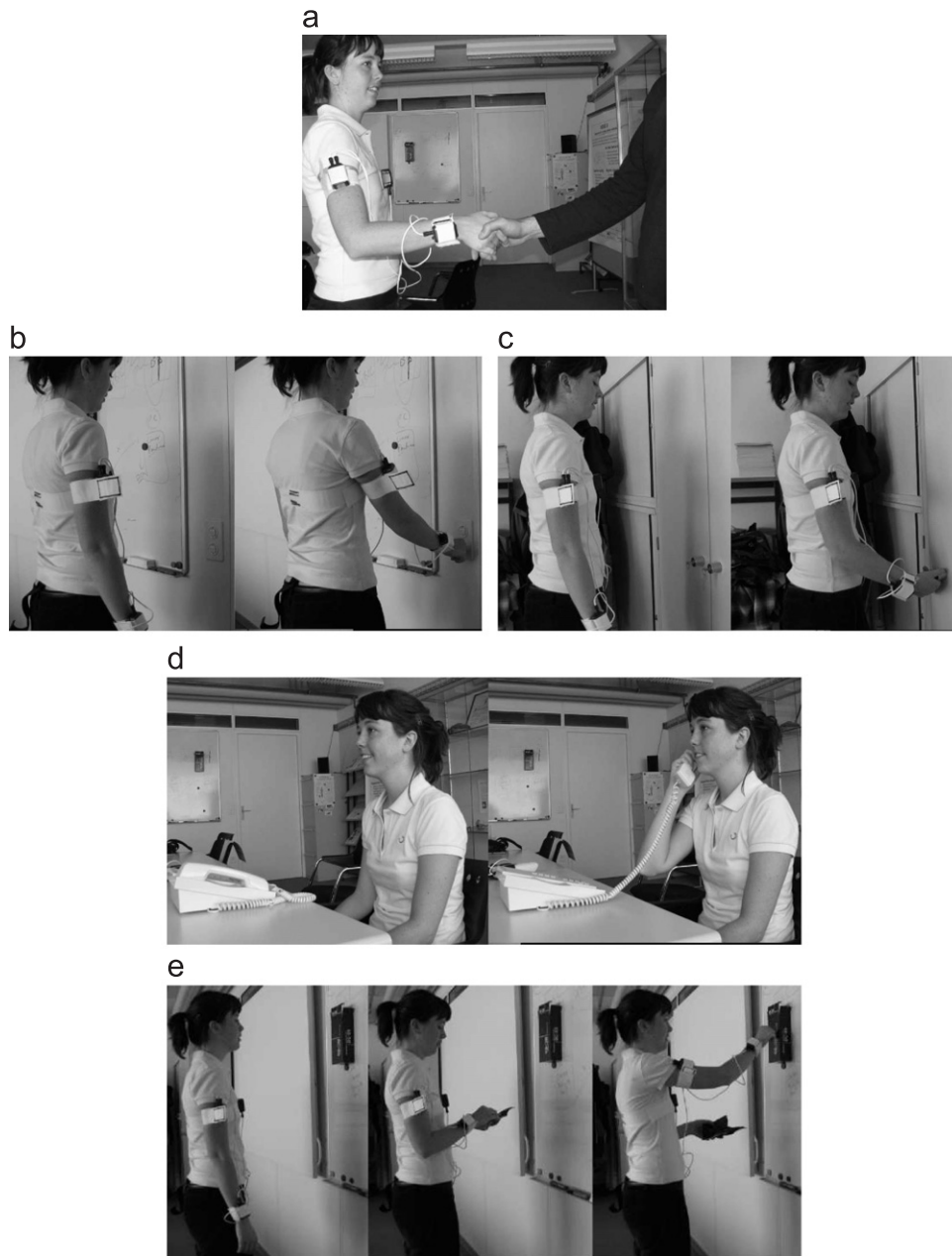


Fig. 2. Visualisation of the relevant gestures (acted) as performed in case study 1: (a) Handshake, (b) light button, (c) door, (d) phone and (e) coin.

and amount of all foods is rather visionary, we believe that an assistive system based on different sensors is conceivable. Hence, the gestures included in this study refer to frequently used human feeding motions. Detecting such gestures reveals information about the timing of nutrition events, e.g. breakfast or lunch and on the category of the food item, e.g. a soup is fed with a spoon, a glass, cup or bottle is usually used for drinking.

Figs. 2 and 3 illustrate the gestures that we aimed to recognise (relevant gestures) in each case study (see Table 2 for a brief description). All relevant gestures are characterised by distinctive movements of the left or right arm. While in case study 1 only movements from the right arm and trunk were used to detect the gestures; case study 2 uses information from both arms as well as from the trunk.

4. Spotting implementation

The implementation of our two-stage spotting approach is detailed in this section. The first stage preselects candidate sections and the second stage refines the preselection (see Fig. 4).

4.1. Preselection stage

This section details the segmentation scheme used for the initial partitioning of the continuous signal stream into motion segments, the search strategy and similarity measure applied to identify potential sections and, finally, the selection of candidate sections.

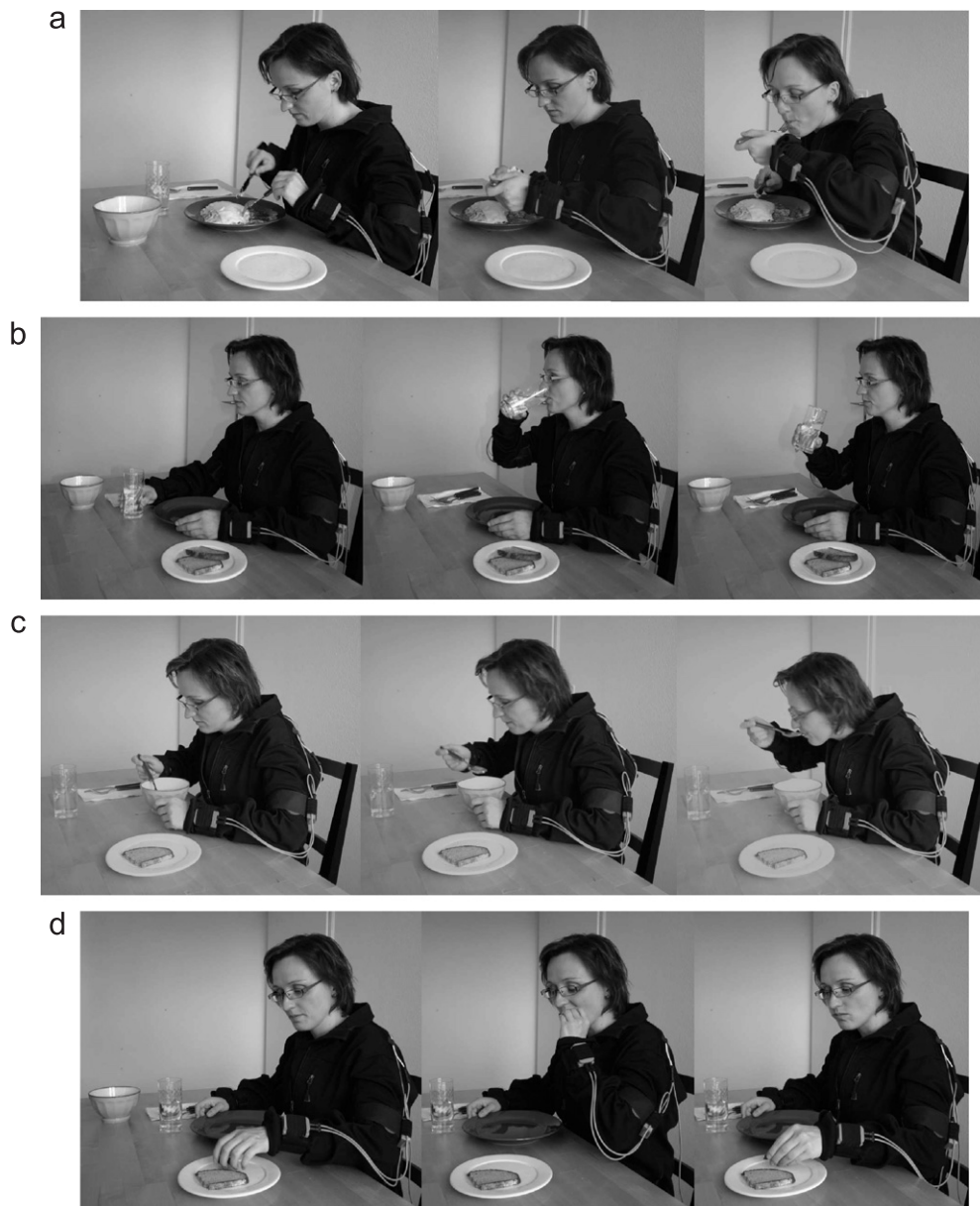


Fig. 3. Visualisation of the relevant gestures (acted) as performed in case study 2: (a) cutlery, (b) drink, (c) spoon and (d) handheld.

4.1.1. Motion segment partitioning

The task of the segmentation algorithm is to partition a motion parameter into non-overlapping, meaningful segments. This task can be considered as a time-series segmentation problem, which has been extensively studied in many application domains. An excellent review of time-series segmentation approaches was provided by Keogh et al. [34].

As motion parameter, we used the pitch and the roll of the lower arm, which are the angle of the lower arm segment to the horizontal plane and the rotation angle with the rotation axis along the limb of the lower arm (see Fig. 5).

These angles have been chosen mainly for the following reasons: Many movements of the entire arm typically involve movements of the lower arm as well. Furthermore, the signals of the lower arm orientation (and in particular pitch and roll)

correlated well with our visual perception of the gestures. Despite good initial results by using the pitch in case study 1, the roll was additionally investigated in case study 2. For certain gestures the segmentation based on the roll matched the gesture boundaries better. This can be explained by the typical feeding motion (moving the hand with a tool to the mouth) involved in the gestures of case study 2.

Although relative orientation information between the lower arm and the upper arm segment, such as joint angles, would generally be well suited for the partitioning of the signal streams, we found that the estimation of those angles using inertial sensors attached to the arm segments can be prone to large errors. The two major sources of errors were inaccurate orientation estimation of the involved sensors (mainly due to magnetic disturbances) and the loose attachment of the sensors to the arm

segments. Attachment issues make the sensors susceptible to displacement while moving the arm. Conversely the pitch and roll of the lower arm could be derived very robustly. The estimation of these angles from raw sensor data was less prone to magnetic disturbances than other orientation angles, specifically the orientation in the horizontal plane.

For the segmentation task, we used the sliding-window and bottom-up algorithm (SWAB) introduced by Keogh et al. [34]. Based on the evaluation of typical test data, we found the algorithm to be well suited for our application. SWAB combines the advantages of a precise bottom-up segmentation scheme with those of a sliding-window algorithm. This allows the algorithm to be used online while keeping a global view on the data.

The algorithm kept a small buffer of the signal data. A bottom-up segmentation was applied to the data in the buffer. From the resulting signal segmentation, the segment with the oldest data was extracted from the buffer and new data were added using the sliding-window approach. This procedure

was repeated as long as new data were available, potentially forever.

The bottom-up partitioning of each buffer of length n started from the arbitrary segmentation of the signal into $n/2$ segments. Next, the cost of merging each pair of adjacent segments was calculated and the lowest cost pair in the buffer was iteratively merged. As the algorithm iterates, more signal segments were merged until all adjacent segments in the buffer exceeded a cost threshold when merged. Figs. 6(a) and (b) depicts the segmentation process of the buffered signal for different segmentation steps (iterations). At iteration 0 the fine-grained initial partitioning can be seen. The final state is depicted in Fig. 6(b). The sliding-window algorithm of SWAB reported the left-most segment from the bottom-up buffer and added new data accordingly. The procedure was restarted with this new data in the bottom-up buffer.

The cost metric for merging two segments was based on the error of approximating the signal with its linear regression (residuals) in the bounds defined by the merged segment. This method can be explained as follows: When the pair of segments differ strongly in its signal shape, the approximation of the merged segments incurs large residuals. Hence it is less likely that the segments belong to the same motion segment. We used the squared sum of the residuals in the bounds of the merged segment as cost function.

To ensure that the algorithm provided a good approximation of the signal, a small cost threshold was required, typically leading to a large number of segments for any of the relevant gestures. These segments did not correspond well to the small number of visually perceived sub-movements of the gesture. As a solution to this problem, we merged adjacent segments, as created by the SWAB algorithm, if their linear regressions had similar slopes. As a result of this extension we obtained motion segment boundaries. Fig. 7 depicts an example of the

Table 2
Description of the relevant gestures in case studies 1 and 2

Gesture	Description
<i>Case study 1</i>	
Light button (LB)	Press light button to turn lights on
Handshake (HS)	Greet person by shaking hands
Phone up (PU)	Pick up receiver. Start position: arm resting on leg, end position: hold receiver to ear
Phone down (PD)	Put down telephone receiver. End position: arm resting on leg
Door (DR)	Turn door knob and open door of cabinet
Coin (CN)	Take out purse from right back pocket of trousers—open purse with right hand—take coin and insert it into slot of vending machine—close purse with right hand—put purse back into pocket
<i>Case study 2</i>	
Cutlery (CL)	Meal intake of Lasagne using fork and knife. Fork tap, loading and manoeuvring to mouth and back with left hand
Drink (DK)	Pick up cup from table—take a sip—put cup back on the table
Spoon (SP)	Meal intake of cereals or soup using a spoon. Spoon loading and manoeuvring to mouth and back
Handheld (HH)	Meal intake of bread slice or chocolate bar using the hand only: moving the left hand to mouth and back.

Unless otherwise noted, all gestures were conducted with the right arm pointing downwards at start/end in case study 1 and with both arms at rest on table at start/end in case study 2.

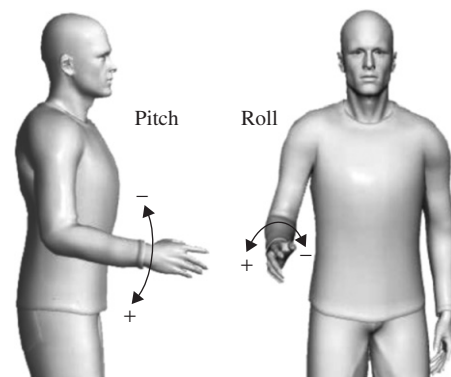


Fig. 5. Orientation angle “pitch” and “roll” of the lower arm segment.

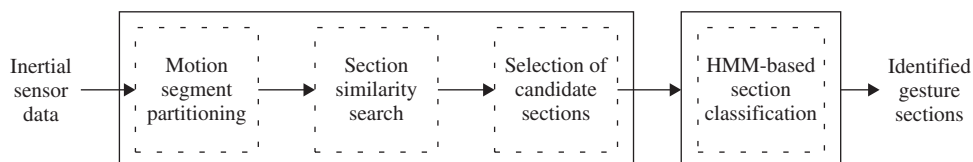


Fig. 4. Detailed structure of the two-stage recognition framework.

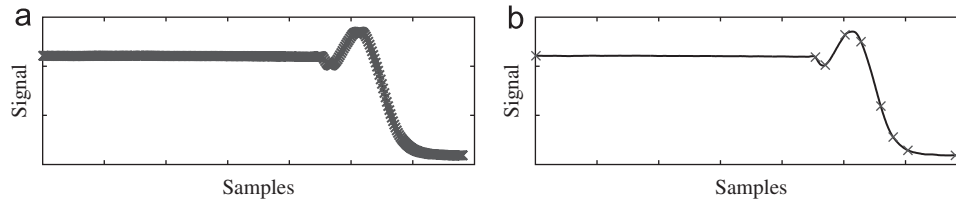


Fig. 6. Segmentation of an example signal stored in the bottom-up buffer at different algorithm iterations. The “cross”-symbols indicate segment boundaries. (a) Initialisation (iteration 0) and (b) termination (iteration 332).

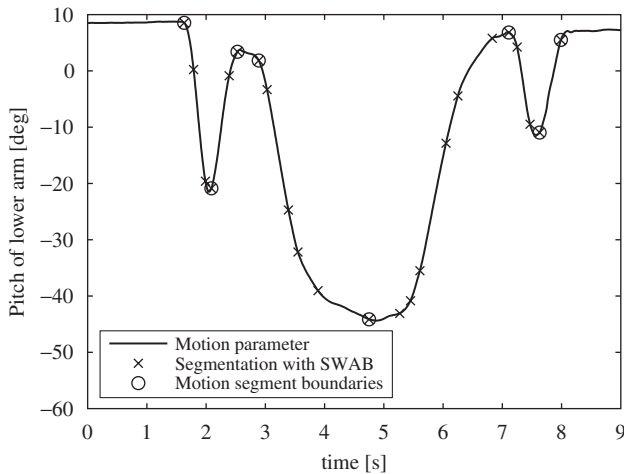


Fig. 7. Segmentation of the “DK” gesture using the pitch of the lower arm as segmentation signal. The cross symbols (“x”) correspond to segmentation boundaries obtained from the SWAB algorithm. The circles (“o”) highlight the remaining segmentation points (motion segment boundaries) based on the proposed extension of the segmentation algorithm.

Table 3
Motion parameter selection for the SWAB algorithm

Gesture	SWAB motion parameter	Body side used in studies
<i>Study 1</i>		
Light button (LB), handshake (HS), phone up (PU), phone down (PD), door (DR), coin (CN)	$pitch_{LA}(t)^a$	Right
<i>Study 2</i>		
Cutlery (CL)	$roll_{LA}(t)^a$	Left
Drink (DK), spoon (SP)	$pitch_{LA}(t)$	Right
Handheld (HH)	$pitch_{LA}(t)$	Left

^aPitch, roll and yaw are Euler angles representing rotations of an object in three-dimensional Euclidean space. The orientation of pitch and roll angles are described in Section 4.1 and Fig. 5. The yaw angle corresponds to absolute orientation in the horizontal plane.

segmentation steps, based on the “DK” gesture that uses the pitch angle as segmentation signal.

For each gesture an individual segmentation parameter could be chosen. Person-specific training was used to accommodate for the dominant body side. In the investigated case studies, the body side was fixed. Table 3 summarises the final choices made in our implementation.

The mean number of segmentation points per gesture for the data sets of both case studies are shown in Table 4. The ratio

Table 4
SWAB segmentation results

Segmentation category	Case study 1	Case study 2
Mean number of SWAB segmentation points per gesture	15 506	13 020
Ratio of segmentation points per gesture by total recorded samples (%)	2.2	0.77

The SWAB segmentation points correspond to the total number of segmentation points for the entire data sets. The ratio of segmentation points by total recorded samples indicates the reduction in search effort achieved by the preselection stage.

of segmentation points to the total recorded samples indicates the achieved reduction in search effort.

4.1.2. Section similarity search

A coarse search based on the motion segment boundaries was used to find sections that contain relevant gestures. The search was performed by considering each motion segment endpoint as potential end of a gesture. For each endpoint, potential start points were derived from preceding motion segment boundaries. The search was performed for each gesture separately. To confine the search space, we introduced two constraints on the sections to be searched. These constraints were adapted to the gesture by training data:

1. For the actual length T of the section we constrained $T_{min} \leq T \leq T_{max}$, where T_{min} and T_{max} denote the minimum and maximum length of the section to be considered.
2. For the number of motion segments n_{MS} in the actual section, we selected $N_{MS,min} \leq n_{MS} \leq N_{MS,max}$, where $N_{MS,min}$ and $N_{MS,max}$, correspond to the minimum and maximum number of motion segments to be contained in the section, respectively.

As search criterion, we used the normalised Euclidean distance¹ given in Eq. (1), where f_{PS} denotes the N_F -dimensional feature vector of the PS, derived from the section under consideration.

We used simple single-value features, such as minimum and maximum signal values of the lower and upper arm pitch and roll, sum of signal samples, the duration of the gesture and the

¹ The normalised Euclidean distance corresponds to the Mahalanobis distance where the covariance matrix is a diagonal matrix.

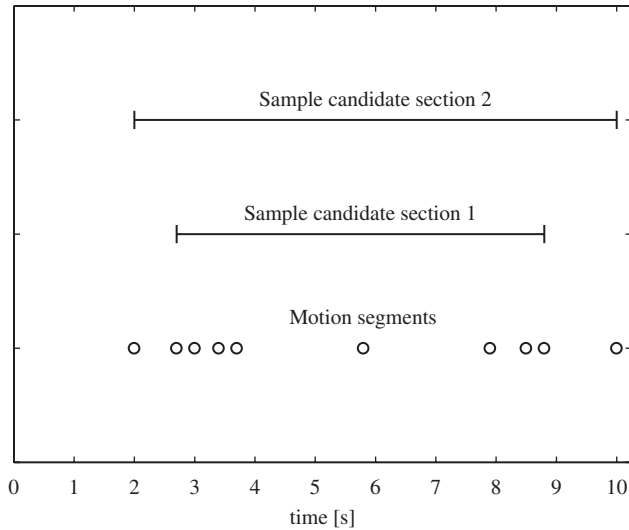


Fig. 8. Overlapping candidate sections.

number of motion segments in the section under consideration. In case study 2, the minimal distance between the hand and estimated head position was additionally used.

The parameters μ_{ik} and σ_{ik} represent the mean and the standard deviation of the i th element of the feature vector of gesture G_k . These were computed from training data:

$$d(\mathbf{f}_{PS}; G_k) = \sqrt{\sum_{i=1}^{N_F} \left(\frac{f_{PS_i} - \mu_{ik}}{\sigma_{ik}} \right)^2},$$

$$\mathbf{f}_{PS} = [f_{PS_1}, \dots, f_{PS_{N_F}}]. \quad (1)$$

The normalised Euclidean distance provided a measure of how similar the motion pattern given in the section were to a specific gesture. During the evaluation of all possible start points for one endpoint, only the section with the minimal distance as retained.

If the distance $d(\mathbf{f}_{PS}, G_k)$ was smaller than a gesture-specific threshold value $d_{min}(G_k)$, the section under investigation was considered to contain gesture G_k . If the condition was satisfied for more than one gesture, the section was considered to contain either one of the corresponding gestures. Depending on the application such collisions need to be checked and handled.

4.1.3. Selection of candidate sections

Fig. 8 schematically shows the collision of two sections obtained by the section search procedure. These overlapping candidate sections were resolved by selecting sections with the smallest similarity values for every occurring collision. In this way non-overlapping candidate sections were obtained for a particular gesture.

4.2. Classification stage

In the CS, we used HMMs, which have long been used in speech recognition, due to their ability to cope with temporal and spatial variations of input patterns [35].

For our evaluation, we considered left–right models with eight continuous features. The features used for the classification differ from the features used in the PS. While in the PS, data sections were characterised by single-valued features, such as the minimum and maximum signal value and the duration of the section, the HMMs were fed with time-series features derived from the candidate sections. Moreover, a separate definition of the feature set was useful to address the classification goal.

The following features were used for the HMM-based CS:

- Pitch and roll angles from the lower arm sensors.
- Pitch and roll angles from the upper arm sensors.
- Derivative of the acceleration signal from the lower arm sensor, with the measurement orientation along the pitch angle measurement.
- The cumulative sum of the acceleration from the lower arm (orientation as before).
- Derivative of the rate of turn signal from the lower arm sensor, with the measurement orientation along the roll angle measurement.
- The cumulative sum of the rate of turn from the lower arm (orientation as before).

We found that all gestures could be modelled using single Gaussian models. Our gesture models consisted of 4–10 states. The choice of the states for each gesture model reflects a trade-off between the complexity of the gesture on the one hand, and available training data, which is necessary to estimate the model parameters properly, on the other. Although some gestures may require more states, we achieved good recognition results with our models, as shown below.

5. Experiments

For the experimental evaluation of our approach, we recorded a variety of different data sets using a commercially available measurement system² with five inertial sensors placed on the body (see Fig. 9). Sensors were attached to the wrists, upper arms and on the upper torso.

Using this setup, we independently recorded continuous data sets from one female and three male right-handed subjects, aged 25–35 years in both case studies. In case study 2 food intake was recorded in two sessions on different days. The subject data sets (S1.1–S1.4 for case study 1 and S2.1–S2.4 for case study 2) were used for testing of our spotting approach. Additional person-specific data were used for training purposes. The purpose of the studies was explained to the subjects. However, the subjects were asked to perform the movements as natural as possible while wearing the sensors.

In order to obtain data sets with a realistic zero-class, we did not set constraints to the movements of the subjects, except that we asked the subject to perform the relevant gestures according to the descriptions given in Table 2. Moreover, to enrich the diversity of movements and to avoid wide intervals constituting

² <http://www.xsens.com>.

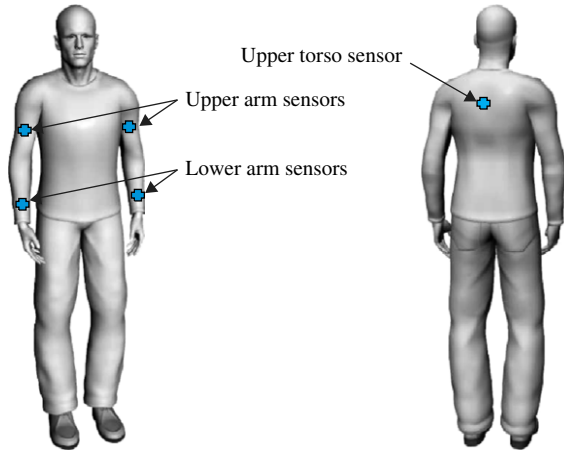


Fig. 9. Sensor placement for gesture recording.

Table 5
Statistics of the recorded data sets

Feature	Case study 1	Case study 2
Total duration of all data sets	7185 s (2.00h)	16848 s (4.68h)
Share of relevant gestures in data sets	25.4% (1826 s)	34.7% (5846 s)

no motion, we defined eight additional gestures to be carried out during the recording which were similar to those gestures we intended to spot. In total 2 h of motion data were recorded for case study 1, and 4.7 h for case study 2, with only 25.4% and 34.7% of the data sets containing relevant gestures for case studies 1 and 2, respectively (see Table 5).

6. Results

For the evaluation of our approach, the evaluation metrics *Precision* and *Recall* were used. These metrics were derived as follows:

$$Recall = \frac{\text{Recognised gestures}}{\text{Relevant gestures}},$$

$$Precision = \frac{\text{Recognised gestures}}{\text{Retrieved gestures}}.$$

Relevant gestures are those gestures that have been conducted by the subject, while retrieved gestures represent the sections that have been reported in either PS or CS. A recognised gesture is a relevant gesture that has been retrieved. Furthermore, we derived the number of insertions (sections that have been retrieved but do not contain a relevant gesture), and the number of deletions (relevant gestures that have not been reported). Fig. 10 illustrates the different evaluation metrics schematically. Set A corresponds to the relevant gestures, set B to the retrieved gestures after the PS and set C to gestures retained after the

CS. The depicted subsets (1–5) reflect the metrics used in this paper.

6.1. Preselection stage

For the spotting of sections likely to contain motion events, appropriate threshold values $d_{min}(G_k)$ were identified for each gesture G_k , by evaluating the performance of the PS on the training data. In general, we observed that the larger the threshold value, the more relevant gestures were retrieved. However, at the same time, the total number of falsely retrieved gestures increased. The precision–recall curves given in Fig. 11 for the gesture “HS” from case study 1 and Fig. 12 for the gesture “SP” from case study 2 illustrate this trade-off for the test data sets (S1.1–S2.4), respectively. Moreover, the individual curves in figures indicate the variation of the detection performance among the subjects.

The vertical lines towards the maximum recall in Fig. 12 can be seen as limitation of the similarity search. For these gestures, some instances were not successfully detected due to variation between training and testing gestures.

Based on such precision–recall curves derived from training data, appropriate threshold values can be chosen considering application-specific requirements. For further evaluation of our approach, we set the thresholds for the individual gestures such that at least 90% of the relevant gestures (gestures that have been conducted) were retrieved in case study 1 and 70% of the relevant in case study 2. This corresponds to a recall value of 0.90 and 0.70, respectively.

Table 6 finally summarises the results of the PS stage for both case studies. For an overall recall value larger than 0.90, we obtained an overall precision value of 0.47 in case study 1. In case study 2, with a recall larger than 0.70, precision dropped to 0.57. The low precision indicated many falsely retrieved sections that did not contain a relevant gesture (insertions). As can be seen, the spotting of simple gestures such as “HS” and “HH” tend to cause more insertions (smaller precision values) than the others.

6.2. Classification stage

We used HMMs to refine the spotting results from the PS.

6.2.1. Model training and initial testing

To accommodate for varying quality in the training process that is due to random initialisation of certain HMM parameters, we trained 10 instantiations of each model and kept the one with the highest score.

For initial model validation, isolated recognition was performed on the test data based on manually added labelling information. From 258 gestures in case study 1, 254 were classified correctly, leading to a recognition rate of 98.4%. For case study 2, a recognition rate of 97.4% was reached from 784 gestures. The results indicate that the models represented the gestures well and were able to recognise the different gestures in the test set accurately.

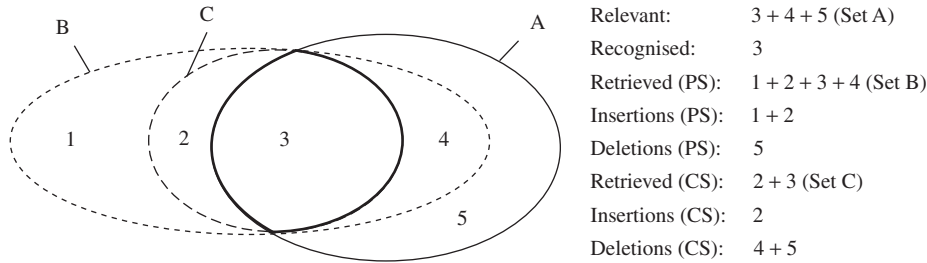


Fig. 10. Visualisation of the applied evaluation metrics for the preselection (PS) and classification (CS) stages.

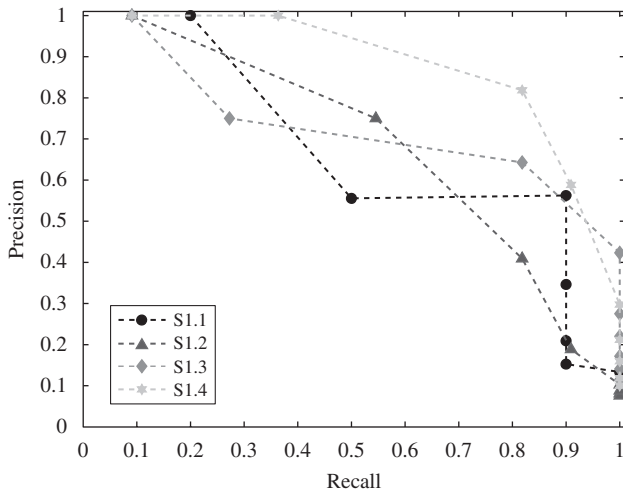


Fig. 11. Precision–recall curves for the “HS” gesture from case study 1, based on the evaluation of data sets from four different test subjects (S1.1–S1.4).

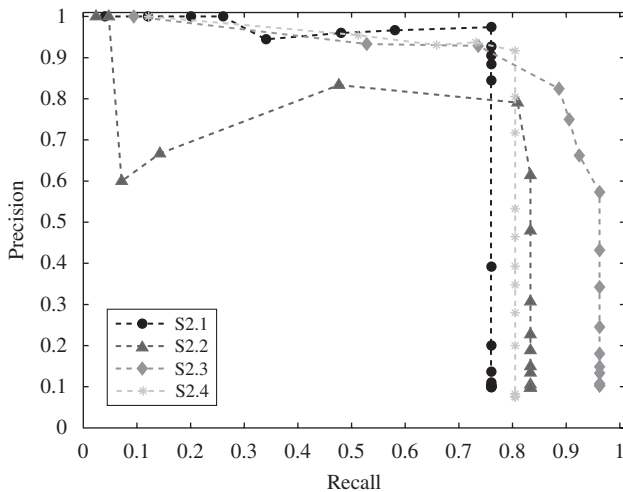


Fig. 12. Precision–recall curves for the “SP” gesture from case study 2, based on the evaluation of data sets from four different test subjects (S2.1–S2.4).

6.2.2. Classification of candidate sections

The trained models were used to classify the candidate sections that have been retrieved in the PS. Only those sections

were retained, for which the recognition of PS and CS agreed. Table 7 presents the final results of this stage for both case studies and all subjects.

The CS correctly recognised most of the relevant gestures that have been retrieved in the PS (the average recall value was slightly reduced from 0.96 to 0.93 for case study 1 and from 0.80 to 0.79 for case study 2). The CS discarded many sections that have been falsely retrieved, leading to much higher precision values, especially in case of the “HH”, “HS” and “LB” gestures. Finally, Fig. 13 depicts the summarised spotting results for all gestures of the case studies 1 and 2.

6.3. Extensions of the CS

Several options exist in which our spotting approach can be extended. One possibility is to include a zero-class model in the CS. The modelling of the zero-class is a challenging and yet unsolved problem. We evaluated the use of two different zero-class models as extension of the CS. These extensions propose no viable elements of our spotting approach, but rather indicate directions of further research. The preliminary results of this investigation are shown in this section.

In case study 1 we evaluated the performance of a zero-class model that is extracted from all relevant gesture models, following the approach presented by Lee and Kim [12]. This modified CS yields a total recall performance of 0.81 (without threshold model: 0.93) and a total precision of 0.82 (without threshold model: 0.74). In direct comparison to the classification without the threshold model, a further increase of the precision was achieved, however, at the cost of decreased recall. Fig. 13 shows the results graphically.

In case study 2, we evaluated the spotting performance using a zero-class model that was constructed on the basis of additional gestures that were carried out by the subjects. An equal number of the gestures was used to build one additional HMM. This garbage model was included in the CS. The modified CS yielded a total recall performance of 0.78 (without garbage model: 0.79) and a total precision of 0.77 (without garbage model: 0.73). Compared to the results of the classification without the garbage model this indicates an improvement of precision at almost constant recall.

Both concepts indicate that classification improvements with zero-class models can be achieved; however, further work in these area is needed.

Table 6
Evaluation results of preselection stage

	Case study 1							Case study 2				
	HS	CN	DR	LB	PU	PD	Total	CL	DK	SP	HH	Total
Relevant ^a	43	43	43	43	43	43	258	196	165	186	153	700
Retrieved ^a	159	64	90	97	63	65	473	278	199	196	310	983
Recognised ^a	41	41	41	40	42	43	248	146	138	154	125	563
Insertions ^a	118	23	49	57	21	22	290	132	61	42	185	420
Deletions ^a	2	2	2	3	1	0	10	50	27	32	28	137
Recall	0.95	0.95	0.95	0.93	0.98	1.0	0.96	0.74	0.84	0.83	0.82	0.80
Precision	0.26	0.64	0.46	0.41	0.67	0.66	0.47	0.53	0.69	0.79	0.40	0.57

^aSee Fig. 10 for corresponding description of evaluation metrics.

Table 7
Spotting results after classification (second stage)

	Case study 1							Case study 2				
	HS	CN	DR	LB	PU	PD	Total	CL	DK	SP	HH	Total
Relevant ^a	43	43	43	43	43	43	258	196	165	186	153	700
Retrieved ^a	57	61	58	41	47	65	329	225	155	163	209	752
Recognised ^a	41	41	41	31	42	43	239	146	137	145	124	552
Insertions ^a	16	20	17	10	5	22	90	79	18	18	85	200
Deletions ^a	2	2	2	12	1	0	19	50	28	41	29	148
Recall	0.95	0.95	0.95	0.72	0.98	1.0	0.93	0.74	0.83	0.78	0.81	0.79
Precision	0.72	0.67	0.71	0.76	0.89	0.66	0.74	0.65	0.88	0.89	0.59	0.73

^aSee Fig. 10 for corresponding description of evaluation metrics.

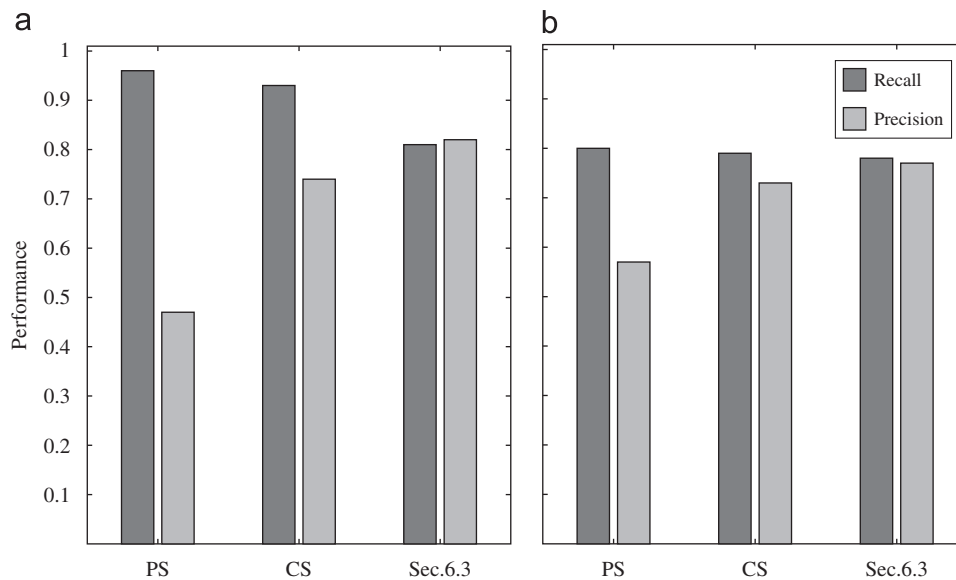


Fig. 13. Summary of the total spotting results for the preselection (PS) and classification (CS) stages in case studies (a) 1 and (b) 2. Additionally, the results of two extensions, discussed in Section 6.3, are shown.

7. Discussion

HMMs have proven to be applicable for recognition tasks in a variety of application domains, including gesture classification from inertial body-worn sensors. However, the spotting of

gestures in a continuous data stream with HMMs is problematic due to their complexity and requirement for a zero-class. The similarity-based search in the PS of our approach presents an elegant way to avoid the explicit modelling of a zero-class. In the HMM-based classification we exploit the competition of

all trained models to select the most probable one. This requires that more than one gesture needs to be included in the CS, which can be seen as a limitation of our approach. However, for most applications, the spotting of several different motion events is aimed. Moreover, an explicit zero-class model can be added, when available, to improve the recognition. Initial results for two different zero-class extensions have been presented in this work.

The similarity-based search procedure used in the PS permits different feature sets for individual gestures. Thus, the search can be tailored to the individual characteristic of a gesture. For example, consider game control gestures, as in Ref. [36], that are conducted in the horizontal or vertical plane only. Such gestures could be described more precisely by specific feature sets. This is an advantage over many established classification procedures, such as k -nearest neighbour classifiers or HMMs, which use the same features for all gestures to be recognised.

The section similarity search can be regarded as a natural extension of the frequently used sliding-window approach for motion and activity detection, as, e.g. in Ref. [2]. We introduced a size-variable search window to accommodate for the variability in the length of gestures and used a dynamic step size given by the segmentation points. While the trivial sliding window was mainly used for the detection of repetitive motions, such as hammering, the approach presented in this work was successfully evaluated for non-cyclic motion events in the two case studies.

The problem of human gesture recognition depends largely on the application domain: In contrast to artificial gestures used for human–computer control or repetitive motions in very specific activities, natural motions in activities of daily living are more challenging to spot. This is due to the fact that control gestures can be constructed to provide strong discrimination, that is typically not the case for gestures being part of activities of daily living. Hence, such gestures contain more intra- and inter-person variability, making the spotting task more challenging. However, the presented results indicate that our spotting procedure performs well for these types of gestures.

In a related work of the authors, one-hand gestures, specifically constructed for game control, were investigated [36]. The approach in that work differed from the current work: firstly, raw inertial sensor signals were used from a sensor attached to a glove at the hand and secondly, the gestures were designed to aim at discrimination and detection in a gaming scenario. In contrast, the current work aimed at recognising natural everyday life gestures with large fluctuations in length and execution from using sensor data from the lower and upper arm. Consequently, with the use of HMMs, a more complex approach was deployed in the current work to achieve the recognition.

The focus of the current work was to analyse the recognition performance using person-specific training. The case studies were designed to incorporate additional motions and gestures and maintain a low share of relevant gestures: 25.4% in case study 1 and 34.7% in case study 2. Both case studies evaluated four subjects each. An initial insight into the subject-specific

variability was obtained from reviewing the precision–recall curves. However, a larger number of users should be evaluated in future works to study the fluctuation in recognition performance and investigate non-personalised detection models in more depth.

The temporal phases of a gestures are onset, core and conclusion. Typically, onset and conclusion are variable transfer states between consecutive gestures. However, the core part is specific for a gesture. In the evaluated case studies most of the gestures were acquired with a defined start and ending position, but all contained a core part. For example, in the phone pick up gesture, the user’s hand moved towards the receiver, picked it up and moved the receiver to the ear. While the movement may commence with the hand at an arbitrary position, the core is preserved in order to successfully complete the activity. The motion segments in the core phase and during the transitions involve direction changes in the segmentation signal. In our approach, segmentation points were created at these positions. Based on the preselection feature set, the section similarity search was used to test for gestures cores at every segmentation point. Hence, we expect that by using the segmentation and search procedure, gestures embedded in arbitrary transitions can be detected.

Looking at the individual results of case studies 1 and 2, we observe lower spotting performance for those gestures included in case study 2. We assume that this is due to higher intra-person variability of those gestures. More specifically, we observed the following additional challenges for the spotting of gestures: (1) differences in the size and consistency of food pieces, (2) additional degrees of freedom produced by the tools used for the food intake and (3) temporal aspects, such as the temperature change of the food and the natural satiety of the subject developed during the intake session. To overcome potential weaknesses in the spotting of gestures related to food intake, we argue that the recognition of such gestures can be enhanced by combining different sensing modalities to develop a dietary monitoring system [37].

We expect that the presented spotting approach can be applied to other types of motion events. At the implementation level an appropriate motion parameter must be selected. This motion parameter shall describe the major properties of the motion event and lead to a reproducible and distinctive motion segmentation. We believe that this can be achieved for many applications.

8. Conclusion and outlook

We conclude that our spotting scheme based on the concept of motion segments is a feasible strategy for the identification of motion events in a continuous signal stream. We demonstrated that our approach works well for arm-based motions that are particularly difficult to recognise due to the inherent complexity of arm motions. Moreover, we have shown that our approach simplifies the rejection of non-relevant gestures. We argue that our method is likely to facilitate a wide range of real-life applications of context and activity recognition.

Acknowledgement

This work was partly supported by the Swiss State Secretariat for Education and Research (SER).

References

- [1] J.A. Ward, P. Lukowicz, G. Tröster, Gesture spotting using wrist worn microphone and 3-axis accelerometer, in: Proceedings of the 2005 Joint Conference on Smart Objects and Ambient Intelligence: Innovative Context-Aware Services: Usages and Technologies, 2005, pp. 99–104.
- [2] J. Ward, P. Lukowicz, G. Tröster, T. Starner, Activity recognition of assembly tasks using body-worn microphones and accelerometers, *IEEE Trans. Pattern Anal.* 28 (10) (2006) 1553–1567.
- [3] M. Stäger, P. Lukowicz, N. Perera, T. von Büren, G. Tröster, T. Starner, Soundbutton: design of a low power wearable audio classification system, in: ISWC 2003: Proceedings of the Seventh IEEE International Symposium on Wearable Computers, 2003, pp. 12–17.
- [4] H. Junker, P. Lukowicz, G. Tröster, Locomotion analysis using a simple feature derived from force sensitive resistors, in: Proceedings of the Second International Conference on Biomedical Engineering, 2004.
- [5] G. Ogris, T. Stiefmeier, H. Junker, P. Lukowicz, G. Troster, Using ultrasonic hand tracking to augment motion analysis based recognition of manipulative gestures, in: B. Rhodes, K. Mase (Eds.), ISWC 2005: Proceedings of the Ninth IEEE International Symposium on Wearable Computers, IEEE Press, New York, 2005, pp. 152–159.
- [6] D. Patterson, D. Fox, H. Kautz, M. Philipose, Fine-grained activity recognition by aggregating abstract object usage, in: B. Rhodes, K. Mase (Eds.), ISWC 2005: Proceedings of the Ninth IEEE International Symposium on Wearable Computers, IEEE Press, New York, 2005, pp. 44–51.
- [7] W. Gao, J. Ma, J. Wu, C. Wang, Sign language recognition based on HMM/ANN/DP, *Int J Pattern Recognition Artif. Intell.* 14 (5) (2000) 587–602.
- [8] T. Starner, Visual recognition of American sign language using hidden markov models, Master's Thesis, Massachusetts Institute of Technology, Boston, USA, 1995.
- [9] L. Campbell, A. Bobick, Recognition of human body motion using phase space constraints, in: Proceedings of the Fifth International Conference on Computer Vision, 1995, pp. 624–630.
- [10] J. Yamato, J. Ohya, K. Ishii, Recognizing human action in time-sequential images using hidden Markov model, in: CVPR 1992: Proceedings of the Conference on Computer Vision and Pattern Recognition, 1992, pp. 379–385.
- [11] M. Brand, N. Oliver, A. Pentland, Coupled hidden Markov models for complex action recognition, in: CVPR 1997: Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 1997, pp. 994–999.
- [12] H.-K. Lee, J.H. Kim, An HMM-based threshold model approach for gesture recognition, *IEEE Trans. Pattern Anal.* 21 (10) (1999) 961–973.
- [13] C. Rao, M. Shah, View-invariance in action recognition, in: CVPR 2001: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, vol. 2, 2001, pp. 316–322.
- [14] M. Shah, R. Jain, Motion-Based Recognition, Kluwer Academic Publishers, Dordrecht, 1997.
- [15] T. Moeslund, E. Granum, A survey of computer vision-based human motion capture, *Comput. Vision Image Understanding* 81 (3) (2001) 231–268.
- [16] Y. Wu, T. Huang, Vision-based gesture recognition: a review, in: Proceedings of the International Gesture Workshop, France, 1999.
- [17] S. Chambers, S. Venkatesh, G. West, H. Bui, Hierarchical recognition of intentional human gestures for sports video annotation, in: R. Kasturi, D. Laurendeau, C. Suen (Eds.), Proceedings of the 16th International Conference on Pattern Recognition, vol. 2, IEEE Press, New York, 2002, pp. 1082–1085.
- [18] A.Y. Benbasat, An inertial measurement unit for user interfaces, Master's Thesis, Massachusetts Institute of Technology, Boston, USA, 2000.
- [19] N. Kern, B. Schiele, A. Schmidt, Multi-sensor activity context detection for wearable computing, in: Proceedings of the European Symposium on Ambient Intelligence, Eindhoven, The Netherlands, 2003, pp. 220–232.
- [20] O. Cakmakci, J. Coutaz, K. Van Laerhoven, H. Gellersen, Context awareness in systems with limited resources, in: ECAI 2002: Proceedings of the Third Workshop on Artificial Intelligence in Mobile Systems (AIMS), 2002.
- [21] L. Bao, Physical activity recognition from acceleration data under semi-naturalistic conditions, Master's Thesis, Massachusetts Institute of Technology, Boston, USA, 2003.
- [22] P. Lukowicz, J.A. Ward, H. Junker, M. Stäger, G. Tröster, A. Atrash, T. Starner, Recognizing workshop activity using body worn microphones and accelerometers, in: Pervasive 2004: Proceedings of the International Conference on Pervasive Computing, Lecture Notes in Computer Science, vol. 3001, Springer, Berlin, 2004, pp. 18–32.
- [23] H. Brashear, T. Starner, P. Lukowicz, H. Junker, Using multiple sensors for mobile sign language recognition, in: ISWC2003: Proceedings of the Seventh IEEE International Symposium on Wearable Computers, 2003, pp. 45–52.
- [24] J.C. Lementec, P. Bajcsy, Recognition of arm gestures using multiple orientation sensors: gesture classification, in: Proceedings of the Seventh International IEEE Conference on Intelligent Transportation Systems, 2004, pp. 965–970.
- [25] J. Deng, H. Tsui, An HMM-based approach for gesture segmentation and recognition, in: Proceedings of the 15th International Conference on Pattern Recognition, vol. 2, 2000, pp. 679–682.
- [26] C. Lee, X. Yangsheng, Online, interactive learning of gestures for human/robot interfaces, in: N. Caplan, C.G. Lee (Eds.), ICRA 1996: Proceedings of the IEEE International Conference on Robotics and Automation, of IEEE Robotics and Automation Society, vol. 4, IEEE Press, New York, 1996, pp. 2982–2987.
- [27] K. Kahol, P. Tripathi, S. Panchanathan, T. Rikakis, Gesture segmentation in complex motion sequences, in: ICIP2003: Proceedings of the International Conference on Image Processing, vol. 2, 2003, pp. 105–108.
- [28] K. Kahol, K. Tripathi, S. Panchanathan, Documenting motion sequences with a personalized annotation system, *IEEE Multimedia* 13 (1) (2006) 37–45.
- [29] T.S. Wang, Y. Shum, Y. Xu, N. Zheng, Unsupervised analysis of human gestures, in: IEEE Pacific Rim Conference on Multimedia, 2001, pp. 174–181.
- [30] R.-H. Liang, M. Ouhyoung, A real-time continuous gesture recognition system for sign language, in: Third IEEE International Conference on Automatic Face and Gesture Recognition, 1998, pp. 558–567.
- [31] P. Morguet, Stochastic modeling of image sequences for the segmentation and recognition of dynamic gestures, Ph.D. Thesis, Technische Universität München, 2000.
- [32] A. Bobick, Movement, activity, and action: the role of knowledge in the perception of motion, *Philos. Trans. R. Soc. London Ser. B* 352 (1358) (1997) 1257–1265.
- [33] O. Amft, H. Junker, G. Tröster, Detection of eating and drinking arm gestures using inertial body-worn sensors, in: B. Rhodes, K. Mase (Eds.), ISWC 2005: IEEE Proceedings of the Ninth International Symposium on Wearable Computers, IEEE Press, New York, 2005, pp. 160–163.
- [34] E. Keogh, S. Chu, D. Hart, M. Pazzani, An online algorithm for segmenting time series, in: Proceedings of the IEEE International Conference on Data Mining, 2001, pp. 289–296.
- [35] L. Rabiner, A tutorial on hidden Markov models and selected applications in speech recognition, *Proc. IEEE* 77 (2) (1989) 257–286.
- [36] D. Bannach, O. Amft, K.S. Kunze, E.A. Heinz, G. Tröster, P. Lukowicz, Waving real hand gestures recorded by wearable motion sensors to a virtual car and driver in a mixed-reality parking game, in: CIG 2007: Proceedings of the 2nd IEEE Symposium on Computational Intelligence and Games, IEEE Press, New York, 2007, pp. 32–39.
- [37] O. Amft, M. Stäger, P. Lukowicz, G. Tröster, Analysis of chewing sounds for dietary monitoring, in: M. Beigl, S. Intille, J. Rekimoto, H. Tokuda (Eds.), UbiComp 2005: Proceedings of the Seventh International Conference on Ubiquitous Computing, Lecture Notes in Computer Science, vol. 3660, Springer, Berlin, Heidelberg, 2005, pp. 56–72.

About the Author—HOLGER JUNKER received the Dipl.-Ing. (M.Sc.) degree in electrical engineering from the Technical University at Brunswick, Germany, in 2000 and the Dr. sc. ETH Zurich (Ph.D.) degree in information technology and electrical engineering from ETH Zurich, Switzerland, in 2005. He joined the Electronics Laboratory at ETH Zurich in 2000 as a research and teaching assistant in the Wearable Computing Group. His research interests included wearable computing, context modelling and recognition, and hardware design of context-aware system architectures.

About the Author—OLIVER AMFT is a Ph.D. candidate at the Wearable Computing Lab., ETH Zurich, planning to complete his thesis before March 2008. He received his Dipl.-Ing. (M.S.) in electrical engineering from Technical University Chemnitz, Germany, in 1999. Before joining ETH in 2004, he was five years with ABB Inc., Switzerland. Oliver led the development of embedded communication systems in different ranks, including the Senior Development Engineer and the R&D Project Manager. He continues to support early development stages at ABB as external consultant. Oliver's research focuses on pervasive healthcare and assistive systems. This includes embedded systems as well as pervasive sensing and pattern recognition for physiology, activity and behavioural analyses. His Ph.D. thesis investigates solutions for on-body automatic dietary monitoring.

About the Author—PAUL LUKOWICZ received the M.Sc. degree in computer science, the M.Sc. degree in physics and the Ph.D. degree in computer science from the University of Karlsruhe, Germany, in 1992, 1993 and 1999, respectively. From 1999 to 2004 he led the Wearable Computing and Computer Architecture Groups at the Electronics Laboratory, ETH Zurich. In 2003, he was appointed Professor of Computer Science at the Institute for Computer Systems and Networks, University of Health Informatics and Technology Tirol, Innsbruck, Austria. Since 2006 he is a professor for Computer Science at the University of Passau, Germany. His research interests include wearable and mobile computer architecture, context and activity recognition.

About the Author—GERHARD TRÖSTER received the M.Sc. degree from the Technical University of Karlsruhe, Germany, in 1978, and the Ph.D. degree from the Technical University of Darmstadt, Germany, in 1984, both in electrical engineering. He is a Professor and head of the Electronics Laboratory, ETH Zurich, Switzerland. During the eight years he spent at Telefunken Corporation, Germany; he was responsible for various national and international research projects focused on key components for ISDN and digital mobile phones. His field of research includes wearable computing, reconfigurable systems, signal processing, mechatronics and electronic packaging. He authored and coauthored more than 100 articles and holds five patents.