ORIGINAL PAPER

# Invariant gait continuum based on the duty-factor

**Preben Fihl · Thomas B. Moeslund**

**Abstract** In this paper, we present a method to describe the continuum of human gait in an invariant manner. The gait description is based on the duty-factor which is adopted from the biomechanics literature. We generate a database of artificial silhouettes representing the three main types of gait, i.e. walking, jogging, and running. By generating silhouettes from different camera angles we make the method invariant to camera viewpoint and to changing directions of movement. Silhouettes are extracted using the Codebook method and represented in a scale- and translation-invariant manner by using shape contexts and tangent orientations. Input silhouettes are matched to the database using the Hungarian method. We define a classifier based on the dissimilarity between the input silhouettes and the gait actions of the database. This classification achieves an overall recognition rate of 87.1% on a diverse test set, which is better than that achieved by other approaches applied to similar data. We extend this classification and results show that our representation of the gait continuum preserves the main features of the duty-factor.

**Keywords** Computer vision · Human motion · Gait analysis · Action recognition · Gait continuum

P. Fihl (✉) · T. B. Moeslund
Laboratory of Computer Vision and Media Technology,
Aalborg University, Niels Jernes Vej 14,
9220 Aalborg, Denmark
e-mail: pfa@cvmt.aau.dk

T. B. Moeslund
e-mail: tbm@cvmt.aau.dk

## 1 Introduction

The human gait[1] contains large amounts of information and has been actively investigated in different research areas like psychology [24], biomechanics [2], robotics [31] and clinical diagnostics [28]. Computer vision research has also focused on analyzing human gait. Since the human gait is a very distinctive type of motion it can be used in many contexts to detect the presence of people, e.g. from surveillance cameras [6,18,26]. Gait as a biometric measure has also received much attention because it is non-intrusive [5,13,25,27,30]. Finally, there has been considerable interest in the computer vision community in the classification of gait types or, more generally, of different types of human action [4,7,20]. This paper is concerned with gait, but its focus is on action recognition rather than on the use of gait in personal identification and we will use the term gait analysis to describe the process of recognizing gait actions.

Action classification has many applications in advanced user interfaces, annotation of video data, intelligent vehicles, and surveillance. More specifically, one of the main human activities that surveillance cameras observe is that of gait, i.e. walking, jogging, or running. Gait analysis is therefore a natural part of automatic surveillance systems.

Used in vast and increasing numbers, surveillance cameras are often applied in unconstrained environments. Many different cameras are used and placements are defined by the surroundings rather than ideal positions for e.g. action recognition. The movements of people in the surveyed areas are also rarely constrained resulting in for example changing moving directions and significant changes in scale. Successful gait analysis in surveillance video will have to take these factors into account, i.e. a general method is needed which

---

[1] By gait is meant bipedal locomotion: walking, jogging and running.

is *invariant* to camera frame rate and calibration, view point, moving speeds, scale change, and non-linear paths of motion. This paper presents such a method.

Other papers have presented systems invariant to one or more of these factors within the area of classification of gait types, but so far none has considered all these factors simultaneously. [14] presents good results on classification of different types of human motion but the system is limited to motion parallel to the image plane. [19] describes a method for behavior understanding by combining actions into human behavior. The method handles rather unconstrained scenes but uses the moving speed of people to classify the action being performed. The moving speed cannot be used for gait-type classification. A person jogging along could easily be moving slower than another person walking fast and human observers distinguishing jogging from running do typically not use the speed as a feature. Furthermore, estimation of speed would require scene knowledge that is not always accessible in surveillance video. [4] uses space-time shapes to recognize actions independently of speed. The method is robust to different viewpoints but cannot cope with non-linear paths created by changes in direction of movement. Other state-of-the-art approaches are mentioned in Sect. 9 along with a comparison of results.

Current approaches to action classification and gait-type analysis consider two or three distinct gait classes, e.g. [14,19] who consider walking and running, or [4,7,20] who consider walking, jogging, and running. However, this distinct classification is not always possible, not even to human observers, and we therefore extend the gait analysis with a more appropriate gait continuum description. Considering gait as a continuum seems intuitive correct for jogging and running, and including walking in such a continuum makes it possible to apply a single descriptor for the whole range of gait types. In this paper, we introduce a formal description of a gait continuum based on a visual recognizable physical feature instead of e.g. a mixture of probabilities of walking, jogging, and running. The next section will describe the basis of this continuum.

## 2 The duty-factor

When a human wants to move fast he/she will run. Running is not simply walking done fast and the different types of gaits are in fact different actions. This is true for vertebrates in general. For example, birds and bats have two distinct flying actions and horses have three different types of gaits. Which action to apply to obtain a certain speed is determined by minimizing some physiological property and physiological research has shown that the optimum action changes discontinuously with changing speed. For example, turtles seem to optimize with respect to muscle power, horses and humans
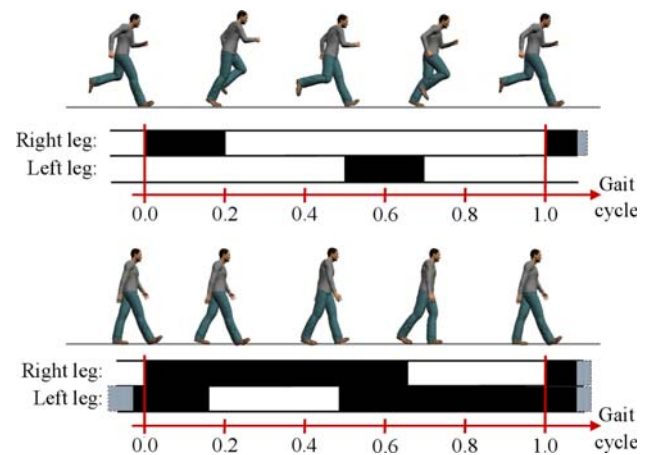
**Fig. 1** Illustration of the duty-factor. The duration of a gait cycle where each foot is on the ground is marked with the *black* areas. The duty-factor for the depicted run cycle (*top*) is 0.2 and 0.65 for the depicted walk cycle (*bottom*)
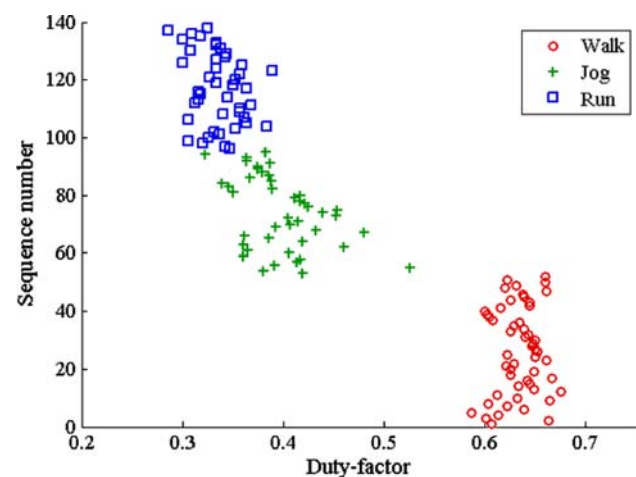


**Fig. 2** The manually annotated duty-factor and gait type for 138 different sequences. Note that the sole purpose of the *y* axis is to spread out the data

with respect to oxygen consumption and other animals by minimizing metabolic power [1].

From a computer vision point of view the question is now if *one* (recognizable) descriptor exist, which can represent the continuum of gait. For bipedal locomotion, in general, the *duty-factor* can do exactly this. The duty-factor is defined as *the fraction of the duration of a stride for which each foot remains on the ground* [2]. Figure 1 illustrates the duty-factor in a walk cycle and a run cycle.

To illustrate the power of this descriptor, we have manually estimated the duty-factor in 138 video sequences containing humans walking, jogging, or running, see Fig. 2. These sequences come from four different sources and contain many different individuals entering and exiting at different angles. Some not even following a straight line (see example frames in Fig. 10).

Figure 2 shows a very clear separation between walking and jogging/running which is in accordance with the fact that those types of gait are in fact different ways of moving. Jogging and running, however, cannot be separated as clearly and there is a gradual transition from one gait type to the other. In fact, the classification of jogging and running is dependent on the observer when considering movements in the transition phase and there exists no clear definition of what separates jogging from running. This problem is apparent in the classification of the sequences used in Fig. 2. Each sequence is either classified by us or comes from a data set where it has been labeled by others. By having more people classify the same sequences it turns out that the classification of some sequences is ambiguous which illustrates the subjectivity in evaluation of jogging and running.[2] Patron and Reid [17] reports classification results from 300 video sequences of people walking, jogging, and running. The sequences are classified by several people resulting in classification rates of 100% for walking, 98% for jogging, and only 81% for running, which illustrates the inherent difficulty in distinguishing the two gait types.

With these results in mind, we will not attempt to do a traditional classification of walking, jogging, and running which in reality has doubtful ground truth data. Rather, we will use the duty-factor to describe jogging and running as a continuum. This explicitly handles the ambiguity of jogging and running since a video sequence that some people will classify as jogging and other people will classify as running simply map to a point on the continuum described by the duty-factor. This point will not have a precise interpretation in terms of jogging and running but the duty-factor will be precise.

As stated earlier, walking and jogging/running are two different ways of moving. However, to get a unified description for all types of gait usually observed in surveillance videos, we also apply the duty-factor to walking and get a single descriptor for the whole gait continuum.

Even though jogging and running are considered as one gait type in the context of the duty-factor they still have a visual distinction to some extend. This visual distinction is used with some success in the current approaches which classify gait into walking, jogging, and running. We acknowledge the results obtained by this type of approaches and we also propose a new method to classify gait into walking, jogging, and running. In our approach however, this is only an intermediate step to optimize the estimation of the duty-factor which we believe to be the best way of describing gait.

---

[2] The problem of ambiguous classification will be clear when watching for example video sequences from the KTH data set [20], e.g. person 4 jogging in scenario 2 versus person 2 running in scenario 2

## 3 Our approach

The method presented in this paper use the duty-factor to describe the major gait types in a unified gait continuum. To enhance the precision in estimation of the duty-factor, we use an effective gait-type classifier to reduce the solution space and then calculate the duty-factor within this subspace. The following will elaborate on our approach.

A current trend in computer vision approaches that deal with analysis of human movement is to use massive amounts of training data, which means spending a lot of time on extracting and annotating the data and temporally aligning the training sequences. To circumvent these problems an alternative approach can be applied in which computer graphics models are used to generate training data. The advantages of this are very fast training plus the ability to easily generate training data from new viewpoints by changing the camera angle.

In classifying gait types, it is not necessary to record a person's exact pose, and silhouettes are therefore sufficient as inputs. Silhouette-based methods have been used with success in the area of human identification by gait [5,13,27]. The goal in human identification is to extract features that describe the personal variation in gait patterns. The features used are often chosen so that they are invariant to the walking speed and in [30] the same set of features even describe the personal variation in gait patterns of people no matter whether they are walking or running. Inspired by the ability of the silhouette-based approaches to describe details in gait, we propose a similar method. Our goal is however quite different from human identification since we want to allow personal variation and describe the different gait types through the duty-factor.

A silhouette-based approach does not need a completely realistic looking computer graphics model as long as the shape is correct and the 3D rendering software Poser [23], which has a built-in Walk Designer, can be used to animate human gaits.

### 3.1 Contributions

To sum up, our approach offers the following three main contributions:

1. The methods applied are chosen and developed to allow for classification in an unconstrained environment. This results in a system that is invariant to more factors than other approaches, i.e. invariant in regard to camera frame rate and calibration, viewpoint, moving speeds, scale change, and non-linear paths of motion.
2. The use of the computer graphics model decouples the training set completely from the test set. Usually methods are tested on data similar to the training set, whereas we
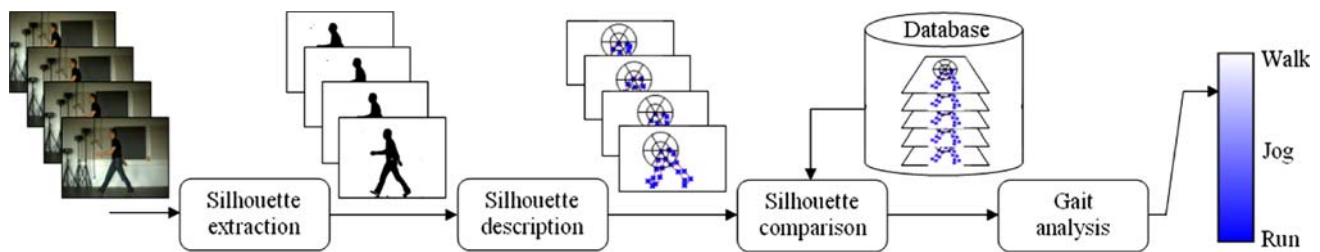
**Fig. 3** An overview of the approach. The main contributions of this paper are the computer generated silhouette database, the gait analysis resulting in a gait continuum, and the ability to handle unconstrained environments achieved by the methods applied throughout the system. The gait-type analysis is further detailed in Fig. 8

train on computer-generated images and test on video data from several different data sets. This is a more challenging task and it makes the system more independent of the type of input data and therefore increases the applicability of the system.

3. The gait continuum is based on a well-established physical property of gait. The duty-factor allows us to describe the whole range of gait types with a single parameter and to extract information that is not dependant on the partially subjective notion of jogging and running.

The framework described in this paper is shown in Fig. 3. The human silhouette is first extracted (Sect. 4) and represented efficiently (Sect. 5). We then compare the silhouette with computer graphics silhouettes (Sect. 7) from a database (Sect. 6). The results of the comparison are calculated for an entire sequence and the gait type and duty-factor of that sequence is extracted (Sect. 8). Results are presented in Sect. 9 and Sect. 10 contain the discussion. Section 11 presents a multi-view extension of the system and Sect. 12 concludes the paper.

## 4 Silhouette extraction

To extract silhouettes from video sequences we do foreground segmentation using the Codebook background subtraction method as described in [8,10]. This method has been shown to be robust in handling both foreground camouflage and shadows. This is achieved by separating intensity and chromaticity in the background model. Moreover, the background model is multi modal and multi layered which allows it to model moving backgrounds such as tree branches and objects that become part of the background after staying stationary for a period of time. To maintain good background subtraction quality over time, it is essential to update the background model and [8] describes two different update mechanisms to handle rapid and gradual changes respectively. By using this robust background subtraction method we can use a diverse set of input sequences from both indoor and outdoor scenes.

## 5 Silhouette description

When a person is moving around in a typical surveillance setup his or her arms will not necessarily swing in a typical "gait" manner; the person may be making other gestures, such as waving, or he/she might be carrying an object. To circumvent the variability and complexity of such scenarios we choose to classify the gait solely on the silhouette of the legs. Furthermore, [12] shows that identification of people on the basis of gait, using the silhouette of legs alone, works just as well as identification based on the silhouette of the entire body.

To extract the silhouette of the legs, we find the height of the silhouette of the entire person and use the bottom 50% as the leg silhouette. Without loss of generality this approach avoids errors from the swinging hands below the hips, although it may not be strictly correct from an anatomic point of view. To reduce noise along the contour we apply morphological operations to the silhouette. Some leg configurations cause holes in the silhouette, e.g. the middle image of Fig. 5. Such holes are descriptive for the silhouette and we include the contour of these holes in the silhouette description.

To allow recognition of gait types across different scales we use shape contexts [3] to describe the leg silhouettes. $n$ points are sampled from the contour of the leg silhouette
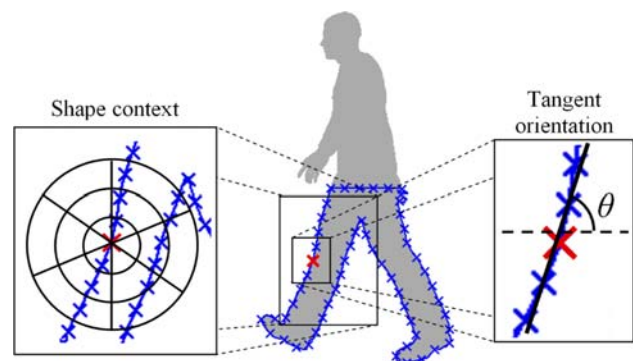


**Fig. 4** Illustration of the silhouette description. The *crosses* illustrate the points sampled from the silhouette. Shape contexts and tangent orientations are used to describe the silhouette

and for each point we determine the shape context and the
tangent orientation at that point, see Fig. 4. With $K$ bins
in the log-polar histogram of the shape context, we get an
$n \times (K + 1)$ matrix describing each silhouette. Scale inva-
riance is achieved with shape contexts by normalizing the
size of the histograms according to the mean distance bet-
ween all point pairs on the contour. Specifically, the normali-
zing constant $q$ used for the radial distances of the histograms
are defined as follows:

$$q = \frac{1}{n^2} \sum_{i=1}^{n} \sum_{j=1}^{n} |p_i - p_j| \qquad (1)$$

where $n$ is the number of points $p$ sampled from the contour.

## 6 Silhouette database

We create a database of human silhouettes performing one
cycle of each of the main gait types: walking, jogging, and
running. To make our method invariant to changes in view-
point we generate database silhouettes from three different
camera angles. With 3D-rendering software this is an easy
and very rapid process that does not require us to capture new
real life data for statistical analysis. The database contains
silhouettes of the human model seen from a side view and
from cameras rotated 30° to both sides. The combination
of the robust silhouette description and three camera angles
enable the method to handle diverse moving directions and
oblique viewing angles. Specifically, database silhouettes can
be matched with silhouettes of people moving at angles of
at least ±45° with respect to the viewing direction. People
moving around in open spaces will often change direction
while in the camera's field of view (creating non-linear paths
of motion), thus we cannot make assumptions about the direc-
tion of movement. To handle this variability each new input
silhouette is matched to database silhouettes taken from all
camera angles. Figure 10, row 1, shows a sequence with a
non-linear motion path where the first frames will match data-
base silhouettes from a viewpoint of −30° and the last frames
will match database silhouettes from a viewpoint of 30°. The
silhouettes generated are represented as described in Sect. 5.
We generate $T$ silhouettes of a gait cycle for each of the three
gait types. This is repeated for the three viewpoints, i.e. $T \cdot 3 \cdot 3$



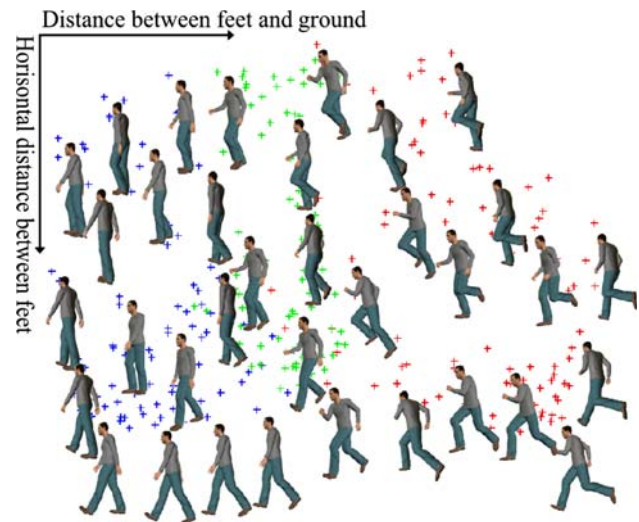**Fig. 5** The database contains silhouettes generated by 3D rendering
software



**Fig. 6** Illustration of the ISOMAP embedding of a representative
subset of the database silhouettes

silhouettes in total. Figure 5 shows five of the poses used for
the description of a side view of a running man.

Each silhouette in the database is annotated with the num-
ber of feet in contact with the ground which is the basis of
the duty-factor calculation.

To analyze the content of the database with respect to the
ability to describe gait we created an Isomap embedding [21]
of the shape context description of the silhouettes. Based on
the cyclic nature of gait and the great resemblance between
gait types we expect that gait information can be described
by some low dimensional manifold. Figure 6 shows the two-
dimensional embedding of our database with silhouettes des-
cribed by shape contexts and tangent orientations and using
the costs resulting from the Hungarian method (described in
Sect. 7) as distances between silhouettes.

According to Fig. 6, we can conclude that the first two
intrinsic parameters of the database represent (1) the total dis-
tance between both feet and the ground and (2) the horizontal
distance between the feet. This reasonable two-dimensional
representation of the database silhouettes shows that our
description of the silhouettes and our silhouette compari-
son metric does capture the underlying manifold of gait sil-
houettes in a precise manner. Hence, gait-type analysis based
on our silhouette description and comparison seems promi-
sing.

## 7 Silhouette comparison

To find the best match between an input silhouette and data-
base silhouettes we follow the method of [3]. We calculate the
cost of matching a sampled point on the input silhouette with
a sampled point on a database silhouette using the $\chi^2$ test

statistics. The cost of matching the shape contexts of point $p_i$ on one silhouette and point $p_j$ on the other silhouette is denoted $c_{i,j}$. The normalized shape contexts at points $p_i$ and $p_j$ are denoted $h_i(k)$ and $h_j(k)$ respectively with $k$ as the bin number, $k = 1, 2, \ldots, K$. The $\chi^2$ test statistics is given as

$$c_{i,j} = \frac{1}{2} \sum_{k=1}^{K} \frac{\left[ h_i(k) - h_j(k) \right]^2}{h_i(k) + h_j(k)} \tag{2}$$

The normalized shape contexts gives $c_{i,j} \in [0; 1]$.

The difference in tangent orientation $\phi_{i,j}$ between points $p_i$ and $p_j$ is normalized and added to $c_{i,j}$ ($\phi_{i,j} \in [0; 1]$). This gives the final cost $C_{i,j}$ of matching the two points:

$$C_{i,j} = a \cdot c_{i,j} + b \cdot \phi_{i,j} \tag{3}$$

where $a$ and $b$ are weights. Experiments have shown that $\phi_{i,j}$ effectively discriminates points that are quite dissimilar whereas $c_{i,j}$ expresses more detailed differences which should have a high impact on the final cost only when tangent orientations are alike. According to this observation we weight the difference in tangent orientation $\phi_{i,j}$ higher than shape context distances $c_{i,j}$. Preliminary experiments show that the method is not too sensitive to the choice of these weights but a ratio of 1 to 3 yields good results, i.e. $a = 1$ and $b = 3$

The costs of matching all point pairs between the two silhouettes are calculated. The Hungarian method [16] is used to solve the square assignment problem of identifying which one-to-one mapping between the two point sets that minimizes the total cost, see Fig. 7. All point pairs are included in the cost minimization, i.e. the ordering of the points is not considered. This is because points sampled from a silhouette with holes will have a very different ordering compared to points sampled from a silhouette without holes but with similar leg configuration, see row three of Fig. 10 (second and third image) for an example.

By finding the best one-to-one mapping between the input silhouette and each of the database silhouettes we can now identify the best match in the whole database as the database silhouette involving the lowest total cost.
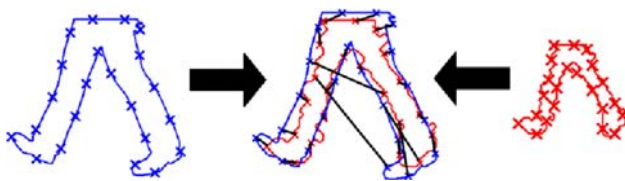
**Fig. 7** An example of one-to-one matching between sampled points from a database silhouette (*left*) and sampled points from an input silhouette (*right*)
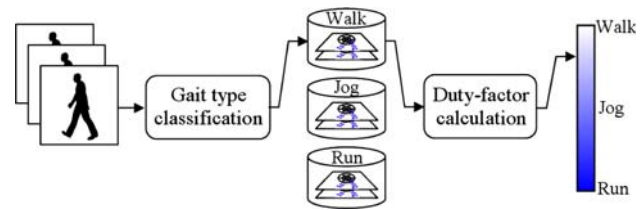
**Fig. 8** An overview of the gait analysis. The figure shows the details of the block "Gait analysis" in Fig. 3. The output of the silhouette comparison is a set of database silhouettes matched to the input sequence. In the gait-type classification these database silhouettes are classified as a gait type which defines a part of the database to be used for the duty-factor calculation

## 8 Gait analysis

The gait analysis consists of two steps. First, we do classification into one of the three gait types, i.e. walking, jogging, or running. Next we calculate the duty-factor $D$ based on the silhouettes from the classified gait type. This is done to maximize the likelihood of a correct duty-factor estimation. Figure 8 illustrates the steps involved in the gait-type analysis. Note that the silhouette extraction, silhouette description, and silhouette comparison all process a single input frame at a time whereas the gait analysis is based on a sequence of input frames.

To get a robust classification of the gait type in the first step we combine three different types of information. We calculate an *action error E* for each action and two associated weights: *action likelihood* $\alpha$ and *temporal consistency* $\beta$. The following subsections describe the gait analysis in detail starting with the action error and the two associated weights followed by the duty-factor calculation.

### 8.1 Action error

The output of the silhouette comparison is a set of distances between the input silhouette and each of the database silhouettes. These distances express the difference or error between two silhouettes. Figure 9 illustrates the output of the silhouette comparison. The database silhouettes are divided into three groups corresponding to walking, jogging, and running. We accumulate the errors in the best matches within each group of database silhouettes to find the difference between the action being performed in the input video and each of the three actions in the database. These accumulated errors constitute the *action error E*.

### 8.2 Action likelihood

When silhouettes of people are extracted in difficult scenarios and at low resolutions the silhouettes can be noisy. This may result in large errors between the input silhouette
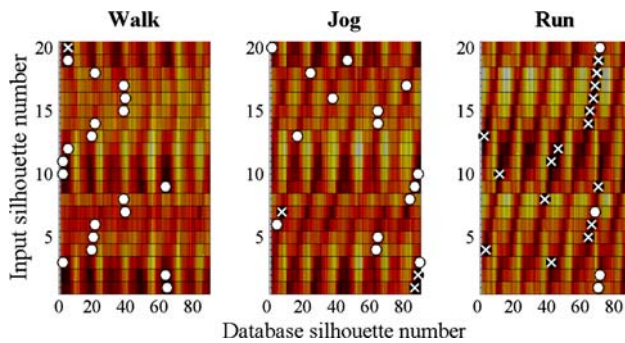
**Fig. 9** Illustration of the silhouette comparison output. The distances between each input silhouette (*y* axes) and the database silhouettes (*x* axes) of each gait type are found. The 90 database silhouettes of the *x* axes are 30 silhouettes ordered in accordance with a gait cycle, i.e. $T = 30$, repeated for the three view points. *Dark colors* illustrate small errors and *bright colors* illustrate large errors. For each input silhouette the best match among silhouettes of the same action is marked with a *white dot* and the best overall match is marked with a *white cross*. The input sequence in the example should be interpreted as follows: the silhouette in the first input frame is closest to walking silhouette number 64, to jogging silhouette number 86, and to running silhouette number 70. The distances to these silhouettes are used when calculating the action error. When all database silhouettes are considered together, the silhouette in the first input frame is closest to jogging silhouette number 86. This is used in the calculation of the two weights

and a database silhouette, even though the actual pose of the person is very similar to that of the database silhouette. At the same time, small errors may be found between noisy input silhouettes and database silhouettes with quite different body configurations (somewhat random matches). To minimize the effect of the latter inaccuracies we weight the action error by the likelihood of that action. The action likelihood of action *a* is given as the percentage of input silhouettes that match action *a* better than the other actions. Since we use the minimum action error the actual weight applied is one minus the action likelihood:

$$\alpha_a = 1 - \frac{n_a}{N} \tag{4}$$

where $n_a$ is the number of input silhouettes in a sequence with the best overall match to a silhouette from action *a*, and *N* is the total number of input silhouettes in that video sequence. This weight will penalize actions that have only a few overall best matches, but with small errors, and will benefit actions that have many overall best matches, e.g. the running action in Fig. 9.

### 8.3 Temporal consistency

When considering only the overall best matches, we can find sub-sequences of the input video where all the best matches are of the same action *and* in the right order with respect to a gait cycle. This is illustrated in Fig. 9, where the running action has great temporal consistency (silhouette numbers

14–19). The database silhouettes are ordered in accordance with a gait cycle. Hence, the straight line between the overall best matches for input silhouettes 14–19 shows that each new input silhouette matches the database silhouette that corresponds to the next body configuration of the running gait cycle.

Sub-sequences with correct temporal ordering of the overall best matches increase our confidence that the action identified is the true action. The temporal consistency describes the length of these sub-sequences. Again, since we use the minimum action error, we apply one minus the temporal consistency as the weight $\beta_a$:

$$\beta_a = 1 - \frac{m_a}{N} \tag{5}$$

where $m_a$ is the number of input silhouettes in a sequence in which the best overall match has correct temporal ordering within action *a*, and *N* is the total number of input silhouettes in that video sequence.

Our definition of temporal consistency is rather strict when you consider the great variation in input silhouettes caused by the unconstrained nature of the input. A strict definition of temporal consistency allows us to weight it more highly than action likelihood, i.e. we apply a scaling factor *w* to $\beta$ to increase the importance of temporal consistency in relation to action likelihood:

$$\beta_a = 1 - w \cdot \frac{m_a}{N} \tag{6}$$

### 8.4 Gait-type classification

The final classifier for the gait type utilizes both the action likelihood and the temporal consistency as weights on the action error. This yields:

$$\text{Action} = \arg\min_a (E_a \cdot \alpha_a \cdot \beta_a) \tag{7}$$

where $E_a$ is the action error, $\alpha_a$ is the action likelihood, $\beta_a$ is the weighted temporal consistency.

### 8.5 Duty-factor calculation

As stated earlier, the duty-factor is defined as the fraction of the duration of a stride for which each foot remains on the ground so we need to identify the duration of a stride and for how long each foot is in contact with the ground.

A stride is defined as one complete gait cycle and consists of two steps. A stride can be identified as the motion from a left foot takeoff (the foot leaves the ground) and until the next left foot takeoff. Accordingly a step can be identified as the motion from a left foot takeoff to the next right foot takeoff. Given this definition of a step it is natural to identify steps in the video sequence by use of the silhouette width. From a side view the silhouette width of a walking person will

oscillate in a periodic manner with peaks corresponding to silhouettes with the feet furthest apart. The interval between two peaks will (to a close approximation) define one step [5]. This also holds for jogging and running and can furthermore be applied to situations with people moving diagonally with respect to the viewing direction. By extracting the silhouette width from each frame of a video sequence we can identify each step (peaks in silhouette width) and hence determine the mean duration of a stride $t_s$ in that sequence.

For how long each foot remains on the ground can be estimated by looking at the database silhouettes that have been matched to a sequence. We do not attempt to estimate ground contact directly in the input videos which would require assumptions about the ground plane and camera calibrations. For a system intended to work in unconstrained surveillance setups such requirements will be a limitation to the system. Instead of estimating the feet's ground contact in the input sequence we infer the ground contact from the database silhouettes that are matched to that sequence. Since each database silhouette is annotated with the number of feet supported on the ground this is a simple lookup in the database. The ground support estimation is based solely on silhouettes from the gait type found in the gait-type classification which maximize the likelihood of a correct estimate of the ground support.

The total ground support $G$ of both feet for a video sequence is the sum of ground support of all the matched database silhouettes within the specific gait type. To get the ground support for each foot we assume a normal moving pattern (not limping, dragging one leg, etc.) so the left and right foot have equal ground support and the mean ground support $g$ for each foot during one stride is $\frac{G}{2 \cdot n_s}$, where $n_s$ is the number of strides in the sequence. The duty-factor $D$ is now given as $D = \frac{g}{t_s}$. In summary we have

$$\text{Duty-factor} \quad D = \frac{G}{2 \cdot n_s \cdot t_s} \tag{8}$$

where $G$ is the total ground support, $n_s$ is the number of strides, and $t_s$ is the mean duration of a stride in the sequence.

The manual labeled data of Fig. 2 allows us to further enhance the precision of the duty-factor description. It can be seen from Fig. 2 that the duty-factor for running is in the interval [0.28;0.39] and jogging is in the interval [0.34;0.53]. This can not be guarantied to be true for all possible executions of running and jogging but the great diversity in the manually labeled data allows us to use these intervals in the duty-factor estimation. Since walking clearly separates from jogging and running and since no lower limit is needed for running we infer the following constraints on the duty factor of running and jogging:

$$D_{\text{running}} \in [0; 0.39]$$
$$D_{\text{jogging}} \in [0.34; 0.53]$$

We apply these bounds as a post-processing step. If the duty-factor of a sequence lies outside one of the appropriate bounds then the duty-factor will be assigned the value of the exceeded bound.

## 9 Results

To emphasize the contributions of our two-step gait analysis we present results on both steps individually and on the gait continuum achieved by combining the two steps.

A number of recent papers have reported good results on the classification of gait types (often in the context of human action classification). To compare our method to these results and to show that the gait-type classification is a solid base for the duty-factor calculation we have tested this first step of the gait analysis on its own. After this comparison we test the duty-factor description with respect to the ground truth data shown in Fig. 2, both on its own and in combination with the gait-type classification.

The tests are conducted on a large and diverse data set. We have compiled 138 video sequences from 4 different data sets. The data sets cover indoor and outdoor video, different moving directions with respect to the camera (up to $\pm45°$ from the viewing direction), non-linear paths, different camera elevations and tilt angles, different video resolutions, and varying silhouette heights (from 41 pixels to 454 pixels). Figure 10 shows example frames from the input
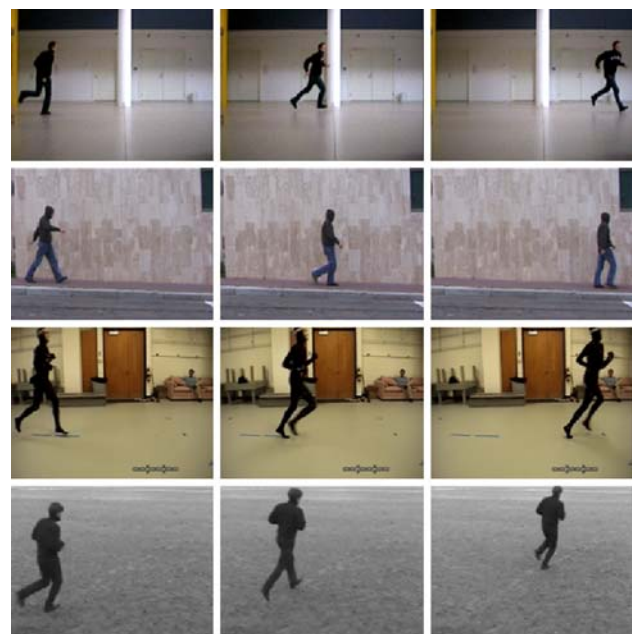


**Fig. 10** Samples from the four different data sets used in the test. *First row* data from our own data set. *Second row* data from the Weizmann data set [4]. *Third row* data from the CMU data set [22]. *Last row* data from the KTH data set [20]

**Table 1** Confusion matrix for the gait-type classification results

|        | Walk | Jog  | Run  |
| ------ | ---- | ---- | ---- |
| Walk   | 96.2 | 3.8  | 0.0  |
| Jog    | 0.0  | 65.9 | 34.1 |
| Run    | 0.0  | 2.6  | 97.4 |

**Table 2** Best reported classification results on the KTH data set

| Methods   |      | Classification results in % | | | |
| --------- | ---- | ----- | ----- | ---- | ---- |
|           |      | Total | Walk  | Jog  | Run  |
| Our method |      | 92.0  | 100.0 | 80.6 | 96.3 |
| [11][a]   | 2008 | 89.3  | 99    | 89   | 80   |
| [17]      | 2007 | 84.3  | 98    | 79   | 76   |
| [29][a]   | 2007 | 80.0  | 88    | 75   | 77   |
| [9][a]    | 2006 | 77.8  | 93.3  | 80.0 | 73.3 |
| [7][a]    | 2005 | 77.3  | 90    | 57   | 85   |
| [20][a]   | 2004 | 75.0  | 83.8  | 60.4 | 54.9 |
| [15][a]   | 2006 | 73.0  | 79    | 88   | 52   |

The matching results of our method are based on the 87 KTH sequences included in our test set

[a] Method work on all actions of the KTH data set

videos. Ground truth gait types were adopted from the data sets when available and manually assigned by us otherwise.

For the silhouette description the number of sampled points $n$ was 100 and the number of bins in the shape contexts $K$ was 60. Thirty silhouettes were used for each gait cycle, i.e. $T = 30$. The temporal consistency was weighted by a factor of four determined through quantitative experiments, i.e. $w = 4$.

### 9.1 Gait-type classification

When testing only the first step of the gait analysis we achieve an overall recognition rate of 87.1%. Table 1 shows the classification results in a confusion matrix.

The matching percentages in Table 1 cannot directly be compared to the results of others since we have included samples from different data sets to obtain more diversity. However, 87 of the sequences originate from the KTH data set and a comparison is possible on this subset of our test sequences. In Table 2, we therefore list the matching results of different methods working on the KTH data set. We acknowledge that the KTH data set contains three additional actions (boxing, hand waving, and hand clapping) and that some of the listed results include these. However, for the results reported in the literature the gait actions are in general not confused with the three hand actions. This means that our results on the KTH data set are comparable to the results of the other approaches.

Another part of our test set is taken from the Weizmann data set [4]. They classify nine different human actions
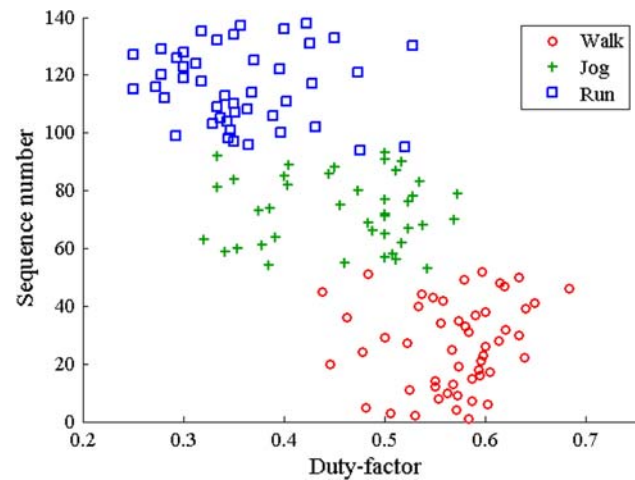


**Fig. 11** The automatically estimated duty-factor from the 138 test sequences without the use of the gait-type classification. The *y* axis solely spreads out the data

including walking and running but not jogging. They achieve a near perfect recognition rate for running and walking. To compare our results to this we remove the jogging silhouettes from the database and leave out the jogging sequences from the test set. In this walking/running classification we achieve an overall recognition rate of 98.9% which is slightly lower. Note, however, that the data sets we are testing on include sequences with varying moving directions where the results in [4] are based on side view sequences.

In summary, the recognition results of our gait-type classification provides a very good basis for the estimation of the duty-factor.

### 9.2 Duty-factor

To test our duty-factor description, we estimate it automatically in the test sequences. To show the effect of our combined gait analysis we first present results for the duty-factor estimated without the preceding gait-type classification to allow for a direct comparison.

Figure 11 shows the resulting duty-factors when the gait-type classification is not used to limit the database silhouettes to just one gait type.

Figure 12 shows the estimated duty-factors with our two-step gait analysis scheme. The estimate of the duty-factor is significantly improved by utilizing the classification results of the gait-type classification. The mean error for the estimate is 0.050 with a standard deviation of 0.045.

## 10 Discussion

When comparing the results of the estimated duty-factor (Fig. 12) with the ground truth data (Fig. 2) it is clear that the overall tendency of the duty-factor is reproduced with the
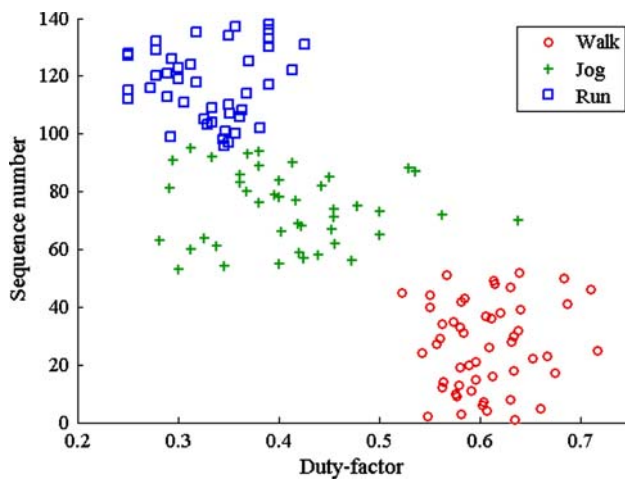
**Fig. 12** The automatically estimated duty-factor from the 138 test sequences when the gait-type classification has been used to limit the database to just one gait type. The $y$ axis solely spreads out the data

**Table 3** The effect of different video characteristics on the classification errors, e.g. 43% of the sequences have a non-side view and these sequences account for 41% of the errors

| Video characteristic | Percentage of sequences | Percentage of errors |
|---|---|---|
| Non-side view | 43 | 41 |
| Small silhouettes[a] | 58 | 59 |
| Low-resolution images[b] | 63 | 65 |
| Non-linear path | 3 | 0 |
| Significant scale change[c] | 41 | 41 |
| Less than two strides | 43 | 41 |

The results are based on 138 test sequences out of which 17 sequences were erroneously classified

[a] Mean silhouette height of less than 90 pixels

[b] Image resolution of 160 × 120 or smaller

[c] Scale change larger than 20% of the mean silhouette height during the sequence

automatic estimation. The estimated duty-factor has greater variability mainly due to small inaccuracies in the silhouette matching. A precise estimate of the duty-factor requires a precise detection of when the foot actually touches the ground. However, this detection is difficult because silhouettes of the human model are quite similar just before and after the foot touches the ground. Inaccuracies in the segmentation of the silhouettes in the input video can make for additional ambiguity in the matching.

The difficulty in estimating the precise moment of ground contact leads to considerations on alternative measures of a gait continuum, e.g. the Froude number [1] that is based on walking speed and the length of the legs. However, such measures requires information about camera calibration and the ground plane which is not always accessible with surveillance video in unconstrained environments. The processing steps involved in our system and the silhouette database all contributes to the overall goal of creating a system that is invariant to usual challenges in surveillance video and a system that can be applied in surveillance setups without requiring additional calibrations.

The misclassifications of the three-class classifier also affect the accuracy of the estimated duty-factor. The duty-factor of the four jogging sequences misclassified as walking disrupt the perfect separation of walking and jogging/running expected from the manually annotated data. All correctly classified sequences however maintain this perfect separation.

A further analysis of the classification errors in Table 1 shows no significant correlation between the classification errors and the camera viewpoint (pan and tilt), the size and quality of the silhouette extracted, the image resolution, the linearity of the path, and the amount of scale change. Furthermore, we also evaluated the effect of the number of frames

(number of gait cycles) in the sequences and found that our method classifies gait types correctly even when there are only a few cycles in the sequence. This analysis is detailed in Table 3 which shows the result of looking at a subset of the test sequences containing a specific video characteristic.

A number of the sequences in Table 3 have more than one of the listed characteristics (e.g. small silhouettes in low resolution images) so the error percentages are somewhat correlated. It should also be noted that the gait-type classification results in only 17 errors which gives a relatively small number of sequences for this analysis. However, the number of errors in each subset corresponds directly to the number of sequences in that subset which is a strong indication that our method is indeed invariant to the main factors relevant for gait classification.

The majority of the errors in Table 1 occur simply because the gait type of jogging resembles that of running which supports the need for a gait continuum.

## 11 Multi-camera setup

The system has been designed to be invariant towards the major challenges in a realistic surveillance setup. Regarding invariance to view point, we have achieved this for gait classification of people moving at an angle of up to ±45° with respect to the view direction. The single-view system can however easily be extended to a multi-view system with synchronized cameras which can allow for gait classification of people moving at completely arbitrary directions. A multi-view system must analyze the gait based on each stride rather than a complete video sequence since people may change both moving direction and type of gait during a sequence.

The direction of movement can be determined in each view by tracking the people and analyzing the tracking data. Tracking is done as described in [8]. If the direction of movement is outside the $\pm 45°$ interval then that view can be excluded. The duration of a stride can be determined as described in Sect. 8.5 from the view where the moving direction is closest to a direct side-view. The gait classification results of the remaining views can be combined into a multi-view classification system by extending Eqs. 7 and 8 into the following and doing the calculations based on the last stride in stead of the whole sequence.

$$\text{Action} = \arg\min_a \left( \sum_V E_a \cdot \alpha_a \cdot w\beta_a \right) \quad (9)$$

$$D = \frac{1}{n_V} \cdot \sum_V D_v \quad (10)$$

where $V$ is the collection of views with acceptable moving directions, $E_a$ is the action error, $\alpha_a$ is the action likelihood, $\beta_a$ is the temporal consistency, $w$ is the scaling factor, $D$ is the duty-factor, $n_V$ is the number of views, and $D_v$ is the duty-factor from view $v$.

Figure 13 illustrates a two-camera setup where the gait classification is based on either one of the cameras or a combination of both cameras.
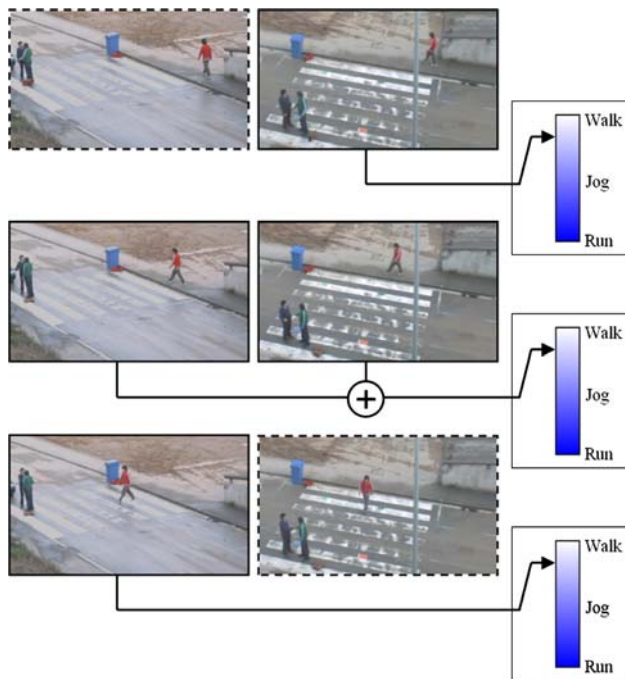


**Fig. 13** A two-camera setup. The figure shows three sets of synchronized frames from two cameras. The multi-camera gait classification enables the system to do classification based on either one view (*top* and *bottom* frames) or a combination of both views (*middle* frame)

## 12 Conclusion

In this paper, we have presented a method for describing gait types with a gait continuum. The method extends the notion of having three gait types, namely running, jogging, and walking. The method is *not* based on statistical analysis of training data but rather on a general gait motion model synthesized using a computer graphics human model. This makes training (from different views) very easy and separates the training and test data completely. The method has been evaluated on different data sets containing all the important factors which such a method should be able to handle. The method performs well (both in its own right and in comparison to related methods) and we therefore conclude that it can be characterized as an *invariant* method for gait description. Furthermore, a multi-view extension of the method is presented.

**References**

1. Alexander, R.: Optimization and gaits in the locomotion of vertebrates. Physiol. Rev. **69**(4), 1199–1227 (1989)
2. Alexander, R.: Energetics and optimization of human walking and running: the 2000 Raymond Pearl Memorial Lecture. Am. J. Hum. Biol. **14**(5), 641–648 (2002)
3. Belongie, S., Malik, J., Puzicha, J.: Shape matching and object recognition using shape contexts. PAMI **24**(4), 509–522 (2002)
4. Blank, M., Gorelick, L., Shechtman, E., Irani, M., Basri, R.: Actions as space-time shapes. In: ICCV (2005)
5. Collins, R., Gross, R., Shi, J.: Silhouette-based human identification from body shape and gait. In: FGR (2002)
6. Cutler, R., Davis, L.S.: Robust real-time periodic motion detection, analysis, and applications. PAMI **22**(8), 781–796 (2000)
7. Dollár, P., Rabaud, V., Cottrell, G., Belongie, S.: Behavior Recognition via Sparse Spatio-Temporal Features. In: VS-PETS (2005)
8. Fihl, P., Corlin, R., Park, S., Moeslund, T., Trivedi, M.: Tracking of individuals in very long video sequences. In: Int. Symposium on Visual Computing. Lake Tahoe, Nevada, USA (2006)
9. Jiang, H., Drew, M.S., Li, Z.N.: Successive convex matching for action detection. In: CVPR (2006)
10. Kim, K., Chalidabhongse, T., Harwood, D., Davis, L.: Real-time foreground–background segmentation using codebook model. Real-Time Imaging **11**(3), 172–185 (2005)
11. Laptev, I., Marszalek, M., Schmid, C., Rozenfeld, B.: Learning realistic human actions from movies. In: CVPR. Alaska, USA (2008)
12. Liu, Z., Malave, L., Osuntugun, A., Sudhakar, P., Sarkar, S.: Towards Understanding the limits of gait recognition. In: Int. Symposium on Defense and Security. Orlando, Florida, USA (2004)
13. Liu, Z., Sarkar, S.: Improved gait recognition by gait dynamics normalization. PAMI **28**(6), 863–876 (2006)
14. Masoud, O., Papanikolopoulos, N.: A method for human action recognition. Image Vis. Comput. **21**(8), 729–743 (2003)
15. Niebles, J.C., Wang, H., Fei-Fei, L.: Unsupervised learning of human action categories using spatial-temporal words. In: BMVC (2006)
16. Papadimitriou, C., Steiglitz, K.: Combinatorial Optimization: Algorithms and Complexity. Courier Dover Publications, Mineola, NY, USA (1998)

17. Patron, A., Reid, I.: A probabilistic framework for recognizing similar actions using spatio-temporal features. In: BMVC (2007)
18. Ran, Y., Weiss, I., Zheng, Q., Davis, L.S.: Pedestrian detection via periodic motion analysis. IJCV **71**(2), 143–160 (2007)
19. Robertson, N., Reid, I.: Behaviour understanding in video: a combined method. In: ICCV (2005)
20. Schüldt, C., Laptev, I., Caputo, B.: Recognizing human actions: a local SVM approach. In: ICPR (2004)
21. Tenenbaum, J., de Silva, V., Langford, J.: A global geometric framework for nonlinear dimensionality reduction. Science **290**(5500), 2319–2323 (2000)
22. CMU Graphics Lab Motion Capture Database (2007). http://mocap.cs.cmu.edu/
23. Poser (ver. 6.0.3.140) (2007). http://www.e-frontier.com/go/poser/
24. Troje, N.F.: Decomposing biological motion: a framework for analysis and synthesis of human gait patterns. J. Vis. 2(5):371–387 (2002)
25. Veeraraghavan, A., Roy-Chowdhury, A., Chellappa, R.: Matching shape sequences in video with applications in human movement analysis. PAMI **27**(12), 1896–1909 (2005)
26. Viola, P., Jones, M.J., Snow, D.: Detecting pedestrians using patterns of motion and appearance. IJCV **63**(2), 153–161 (2005)
27. Wang, L., Tan, T.N., Ning, H.Z., Hu, W.M.: Fusion of static and dynamic body biometrics for gait recognition. IEEE Trans. Circuits Syst. Video Technol. **14**(2), 149–158 (2004)
28. Whittle, M.W.: Gait Analysis: An Introduction. Butterworth-Heinemann Ltd., London (2001)
29. Wong, S.F., Cipolla, R.: Extracting spatiotemporal interest points using global information. In: ICCV. Rio de Janeiro, Brazil (2007)
30. Yam, C., Nixon, M., Carter, J.: On the relationship of human walking and running: automatic person identification by gait. In: ICPR (2002)
31. Yang, H.D., Park, A.Y., Lee, S.W.: Human–robot interaction by whole body gesture spotting and recognition. In: ICPR (2006)