

Real-Time Human Motion Analysis by Image Skeletonization

Hironobu FUJIYOSHI[†], *Member*, Alan J. LIPTON^{††}, *Nonmember*,
and Takeo KANADE^{†††}, *Member*

SUMMARY In this paper, a process is described for analysing the motion of a human target in a video stream. Moving targets are detected and their boundaries extracted. From these, a “star” skeleton is produced. Two motion cues are determined from this skeletonization: body posture, and cyclic motion of skeleton segments. These cues are used to determine human activities such as walking or running, and even potentially, the target’s gait. Unlike other methods, this does not require an *a priori* human model, or a large number of “pixels on target”. Furthermore, it is computationally inexpensive, and thus ideal for real-world video applications such as outdoor video surveillance.

key words: *image skeletonization, human motion analysis, activity recognition*

1. Introduction

Using video in machine understanding has recently become a significant research topic. One of the more active areas is activity understanding from video imagery [7]. Understanding activities involves being able to detect and classify targets of interest and analyze what they are doing. Human motion analysis is one such research area. There have been several good human detection schemes, such as [8] which use static imagery. But detecting and analyzing human motion in real time from video imagery has only recently become viable with algorithms like *Pfinder* [10] and *W⁴* [5]. These algorithms represent a good first step to the problem of recognizing and analyzing humans, but they still have some drawbacks. In general, they work by detecting features (such as hands, feet and head), tracking them, and fitting them to some *a priori* human model such as the *cardboard model* of Ju *et al.* [6].

There are two main drawbacks of these systems in their present forms: they are completely human specific, and they require a great deal of image-based information in order to work effectively. For general video applications, it may be necessary to derive motion analysis tools which are not constrained to human models,

but are applicable to other types of targets, or even to classifying targets into different types. In some real video applications, such as outdoor surveillance, it is unlikely that there will be enough “pixels on target” to adequately apply these methods. What is required is a fast, robust system which can make broad assumptions about target motion from small amounts of image data.

This paper proposes the use of the “star” skeletonization procedure for analyzing the motion of targets — particularly, human targets. The notion is that a simple form of skeletonization which only extracts the broad internal motion features of a target can be employed to analyze its motion.

Once a skeleton is extracted, motion cues can be determined from it. The two cues dealt with in this paper are: cyclic motion of “leg” segments, and the posture of the “torso” segment. These cues, when taken together can be used to classify the motion of an erect human as “walking” or “running”.

This paper is organized as follows: section 2 describes how moving targets are extracted in real-time from a video stream, Sect. 3 describes the processing of these target images and Sect. 4 describes human motion analysis. System analysis and conclusions are presented in Sects. 5 and 6.

2. Real-Time Target Extraction

The initial stage of the human motion analysis problem is the extraction of moving targets from a video stream. There are three conventional approaches to moving target detection: temporal differencing (two-frame or three-frame) [1], background subtraction [5], [10] and optical flow (see [2] for an excellent discussion). Temporal differencing is very adaptive to dynamic environments, but generally does a poor job of extracting all relevant feature pixels. Background subtraction provides the most complete feature data, but is extremely sensitive to dynamic scene changes due to lighting and extraneous events. Optical flow can be used to detect independently moving targets in the presence of camera motion, however most optical flow computation methods are very complex and are inapplicable to real-time algorithms without specialized hardware.

The approach presented here is similar to that taken in [5] and is an attempt to make background

Manuscript received March 28, 2003.

Manuscript revised July 31, 2003.

[†]The author is with the Department of Computer Science, Chubu University, Kasugai-shi, 487-8501 Japan.

^{††}The author is with the Diamondback Vision, Inc., Reston, VA, USA.

^{†††}The author is with the Robotics Institute, Carnegie Mellon University, Pittsburgh, PA, USA.

subtraction more robust to environmental dynamism. The notion is to use an adaptive background model to accommodate changes to the background while maintaining the ability to detect independently moving targets.

Consider a stabilized video stream or a stationary video camera viewing a scene. The returned image stream is denoted I_n where n is the frame number. There are four types of image motion which are significant for the purposes of moving target detection: slow dynamic changes to the environment such as slowly changing lighting conditions; “once-off” independently moving false alarms such as tree branches breaking and falling to the ground; moving environmental clutter such as leaves blowing in the wind; and legitimate moving targets.

The first of these issues is dealt with by using a statistical model of the background to provide a mechanism to adapt to slow changes in the environment. For each pixel value p_n in the n^{th} frame, a running average \bar{p}_n and a form of standard deviation σ_{p_n} are maintained by temporal filtering. Due to the filtering process, these statistics change over time reflecting dynamism in the environment.

The filter is of the form

$$F(t) = e^a, \quad a = -\frac{t}{\tau} \quad (1)$$

where τ is a time constant which can be configured to refine the behavior of the system. When the a is minus value, $F(t)$ will work as an attenuator.

The filter is implemented at the discrete domain:

$$\begin{aligned} \bar{p}_{n+1} &= (1 - \alpha)p_{n+1} + \alpha\bar{p}_n \\ \bar{\sigma}_{n+1} &= (1 - \alpha)|p_{n+1} - \bar{p}_{n+1}| + \alpha\bar{\sigma}_n \\ (0 < \alpha < 1) \end{aligned} \quad (2)$$

where $\alpha = \tau \times f$, and f is the frame rate. Unlike the models of both [5] and [10], this statistical model incorporates noise measurements to determine foreground pixels, rather than a simple threshold. This idea is inspired by [4].

If a pixel has a value which is more than 2σ from \bar{p}_n , then it is considered a foreground pixel. At this point a multiple hypothesis approach is used for determining its behavior. A new set of statistics (\bar{p}', σ') is initialized for this pixel and the original set is remembered. If, after time $t = 3\tau$, the pixel value has not returned to its original statistical value, the new statistics are chosen as replacements for the old.

“Moving” pixels are aggregated using a connected component approach so that individual target regions can be extracted. Transient moving objects will cause short term changes to the image stream that will not be included in the background model, but will be continually tracked, whereas more permanent changes will (after 3τ) be absorbed into the background.

3. Target Pre-Processing

No motion detection algorithm is perfect. There will be spurious pixels detected, holes in moving features, “interlacing” effects from video digitization processes, and other anomalies. Foreground regions are initially filtered for size to remove spurious features, and then the remaining targets are pre-processed before motion analysis is performed.

3.1 Pre-Processing

The first pre-processing step is to clean up anomalies in the targets. This is done by a morphological dilation followed by an erosion. This removes any small holes in the target and smoothes out any interlacing anomalies. In this implementation, the target is dilated twice followed by a single erosion. This effectively robustifies small features such as thin arm or leg segments.

After the target has been cleaned, its outline is extracted using a border following algorithm. The process is shown in Fig. 1.

3.2 “Star” Skeletonization

An important cue in determining the internal motion of a moving target is the change in its boundary shape over time and a good way to quantify this is to use skeletonization. There are many standard techniques for skeletonization such as thinning and distance transformation. However, these techniques are computationally expensive and moreover, are highly susceptible to noise in the target boundary. The method proposed here provides a simple, real-time, robust way of detecting extremal points on the boundary of the target to produce a “star” skeleton. The “star” skeleton consists of only the gross extremities of the target joined to its centroid in a “star” fashion.

1. The centroid of the target image boundary (x_c, y_c)

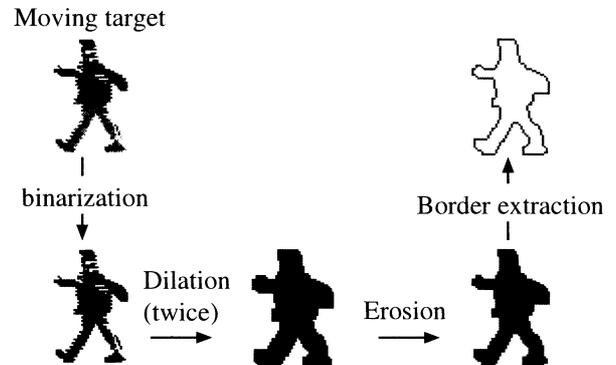


Fig. 1 Target pre-processing. A moving target region is morphologically dilated (twice) then eroded. Then its border is extracted.

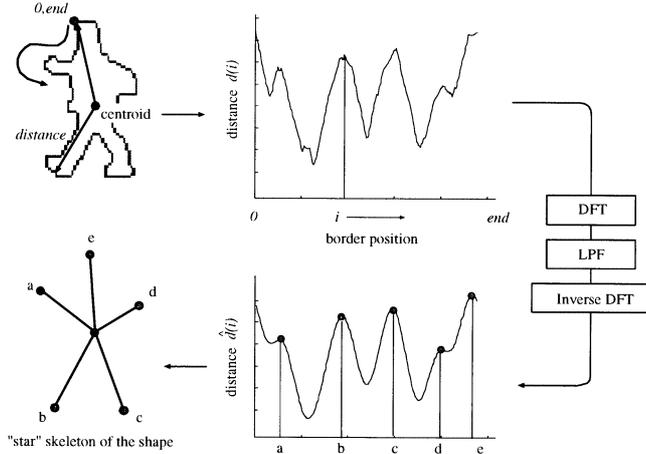


Fig. 2 The boundary is “unwrapped” as a distance function from the centroid. This function is then smoothed and extremal points are extracted.

is determined.

$$\begin{aligned} x_c &= \frac{1}{N_b} \sum_{i=1}^{N_b} x_i, \\ y_c &= \frac{1}{N_b} \sum_{i=1}^{N_b} y_i \end{aligned} \quad (3)$$

where (x_c, y_c) is the *average* boundary pixel position, N_b is the number of boundary pixels, and (x_i, y_i) is a pixel on the boundary of the target.

- The distances d_i from the centroid (x_c, y_c) to each border point (x_i, y_i) are calculated

$$d_i = \sqrt{(x_i - x_c)^2 + (y_i - y_c)^2} \quad (4)$$

These are expressed as a one dimensional discrete function $d(i) = d_i$. Note that this function is periodic with period N_b .

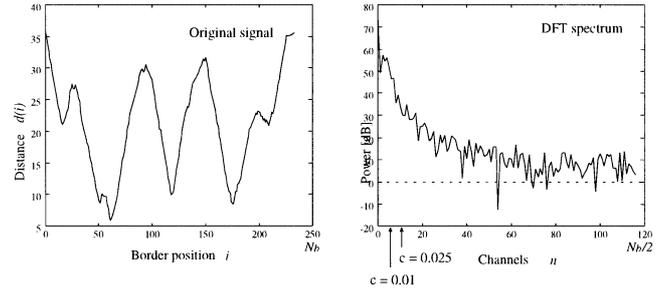
- The signal $d(i)$ is then smoothed for noise reduction, becoming $\hat{d}(i)$. This can be done using a linear smoothing filter or low pass filtering in the Fourier domain.
- Local maxima of $\hat{d}(i)$ are taken as extremal points, and the “star” skeleton is constructed by connecting them to the target centroid (x_c, y_c) . Local maxima are detected by finding zero-crossings of the difference function

$$\delta(i) = \hat{d}(i) - \hat{d}(i - 1) \quad (5)$$

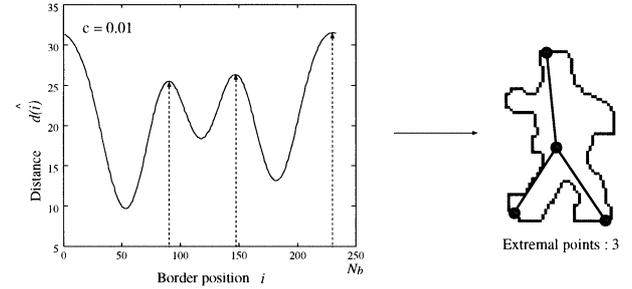
This procedure for producing “star” skeletons is illustrated in Fig. 2.

3.3 Advantages of “Star” Skeletonization

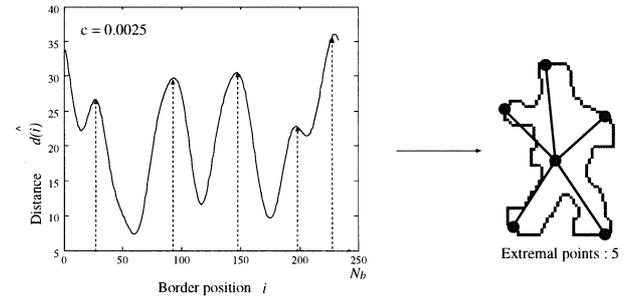
There are three main advantages of this type of skeletonization process. It is not iterative and is, therefore, computationally cheap. It also explicitly provides



(a) Original signal and its DFT spectrum



(b) smoothed signal 1 ($c=0.01$)



(c) smoothed signal 2 ($c=0.025$)

Fig. 3 Effect of cut-off value c . When c is small only gross features are extracted, but larger values of c detect more extremal points.

a mechanism for controlling scale sensitivity. Finally, it relies on no *a priori* human model.

The scale of features which can be detected is directly configurable by changing the cutoff frequency c of the low-pass filter. Figure 3 shows two smoothed versions of $d(i)$ for different values of c : $c = 0.01 \times N_b$ and $c = 0.025 \times N_b$. For the higher value of c , more detail is included in the “star” skeleton because more of the smaller boundary features are retained in $\hat{d}(i)$. So the method can be scaled for different levels of target complexity.

An interesting application of this scalability is the ability to measure the complexity of a target by examining the number of extremal points extracted as a function of smoothing.

Other analysis techniques [5], [6], [10], require *a priori* models of humans — such as the *cardboard model* in order to analyze human activities. Using the skeletonization approach, no such models are required, so the method can be applied to other objects like ani-

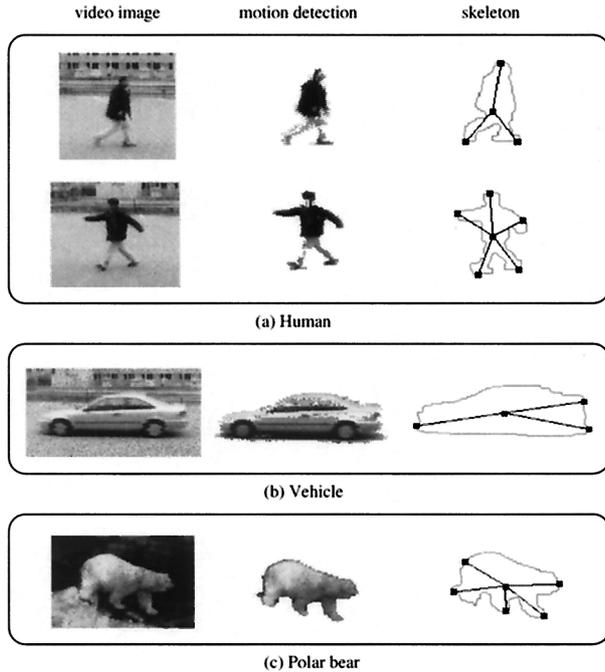


Fig. 4 Skeletonization of different moving targets. It is clear the structure and rigidity of the skeleton is significant in analyzing target motion.

mals and vehicles (see Fig. 4). It is clear that the structure and rigidity of the skeleton are important cues in analysing different types of targets. However, in this implementation, only human motion is considered. Also, unlike other methods which require the tracking of specific features, this method uses only the object's boundary so there is no requirement for a large number of "pixels on target".

4. Human Motion Analysis

One technique often used to analyze the motion or gait of an individual target is the cyclic motion of skeletal components [9]. However, in this implementation, the knowledge of individual joint positions cannot be determined in real-time. So a more fundamental cyclic analysis must be performed.

Another cue to the gait of the target is its posture. Using only a metric based on the "star" skeleton, it is possible to determine the posture of a moving human.

4.1 Significant Features of the "Star" Skeleton

For the cases in which a human is moving in an upright position, it can be assumed that the lower extremal points are legs, so choosing these as points to analyze cyclic motion seems a reasonable approach. In particular, the left-most lower extremal point (l_x, l_y) is used as the cyclic point. Note that this choice does not guarantee that the analysis is being performed on the same physical leg such as a right/left leg at all times.

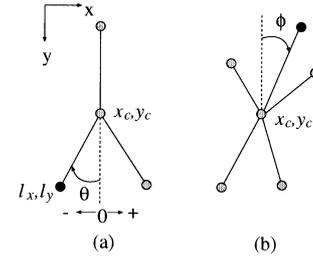


Fig. 5 Determination of skeleton features. (a) θ is the angle the left cyclic point (leg) makes with the vertical, and (b) ϕ is the angle the torso makes with the vertical.

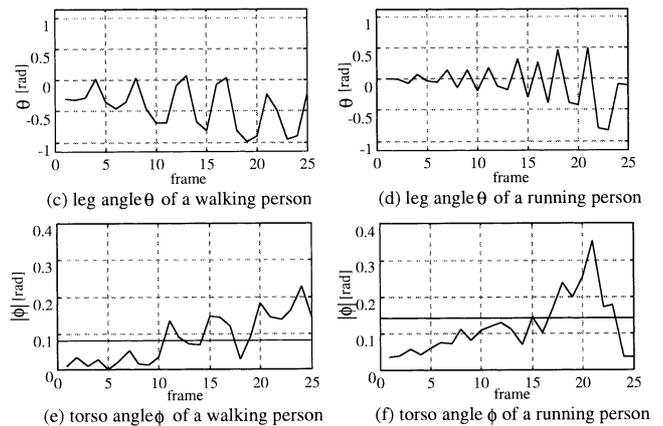
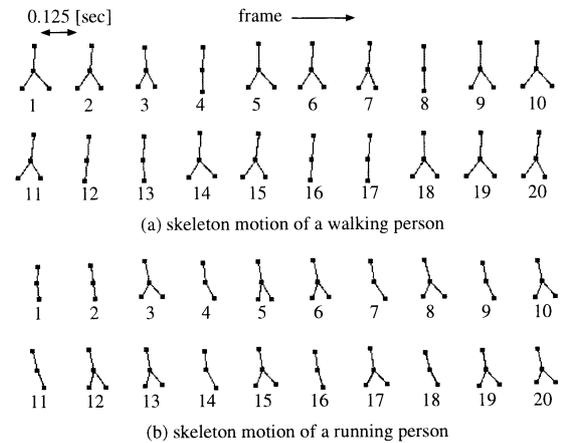


Fig. 6 Skeleton motion sequences. Clearly, the periodic motion of θ_n provides cues to the target's motion as does the mean value of ϕ_n .

However, it is not necessary that the same leg are detected at all times, because the cyclic structure of the motion will still be evident from this point's motion. If $\{(x_i^s, y_i^s)\}$ is the set of extremal points, (l_x, l_y) is chosen according to the following condition:

$$(l_x, l_y) = (x_i^s, y_i^s) : x_i^s = \min_{y_i^s < y_c} x_i^s \quad (6)$$

Then, the angle (l_x, l_y) makes with the vertical θ is calculated as

$$\theta = \tan^{-1} \frac{l_x - x_c}{l_y - y_c} \quad (7)$$

Figure 5 (a) shows the definition of (l_x, l_y) and θ .

One cue to determining the posture of a moving human is the inclination of the torso. This can be approximated by the angle of the upper-most extremal point of the target. This angle ϕ can be determined in exactly the same manner as θ . See Fig. 5 (b).

Figures 6 (a)–(b) show human target skeleton motion sequences for walking and running when the cutoff frequency c was set as 0.001. Figures 6 (c)–(d) show the values of θ_n for the cyclic point. These data were acquired in real-time from a video stream with frame rate 8 Hz. This value is not a constant in this technique but depends on the amount of processing which is required to perform motion analysis and target pre-processing.

Note that in Fig. 6 (c), there is an offset in the value of θ_n in the negative direction. This is because only the leftmost leg (from a visual point of view) is used in the calculation and the calculation of θ is therefore biased towards the negative. There is also a bias introduced by the gait of the person. If s/he is running, the body tends to lean forward, and the values of θ_n tend to reflect this overall posture. Another feature which can clearly be observed is that the frequency of the cyclic motion point is clearly higher in the case of the running person, so this can be used as a good metric for classifying the speed of human motion.

Comparing the average values $\bar{\phi}_n$ in Figs. 6 (e)–(f) show that the posture of a running target can easily be distinguished from that of a walking one using the angle of the torso segment as a guide.

4.1.1 Cycle Detection

Figures 6 (c)–(d) display a clear cyclical nature in θ_n . To quantify these signals, it is useful to move into the Fourier domain. However, there is a great deal of signal noise, so a naive Fourier transform will not yield useful results — see Fig. 7 (b). Here, the power spectrum of θ_n shows a great deal of background noise.

To emphasize the major cyclic component, an autocorrelation is performed on θ_n providing a new signal R_i .

$$R_i = \frac{1}{N+1-i} \sum_{n=1}^N \theta_n \theta_{n-i} \quad (8)$$

where N is number of frames. This is shown in Fig. 7 (c).

This autocorrelation process introduces a new source of noise due to the bias (or DC component) of the θ_n signal. When low frequency components are autocorrelated, they remain in the signal and show up in

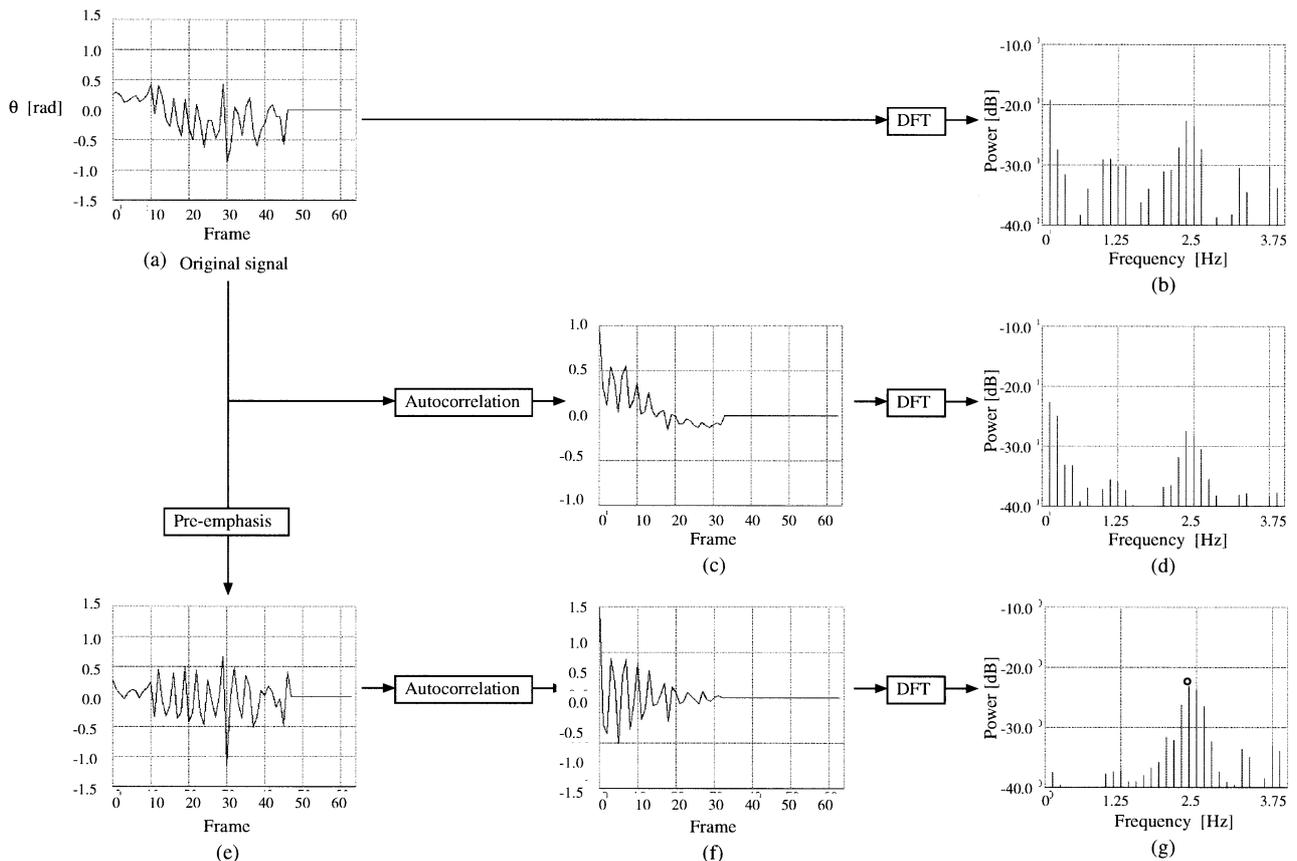


Fig. 7 Process for detecting cyclic motion.

the power spectrum as a large peak in the low frequencies with a degeneration of 6 [dB/oct] in the case of Fig. 7 (d). To alleviate this problem, a high frequency pre-emphasis filter $H(z)$ is applied to the signal before autocorrelation. The filter used is:

$$H(z) = 1 - az^{-1} \tag{9}$$

with a chosen empirically to be ≈ 1.0 . This yields the figure shown in Fig. 7 (e).

Finally, Fig. 7 (g) shows that the major cyclic component of the cyclic point can be easily extracted from the power spectrum of this processed signal.

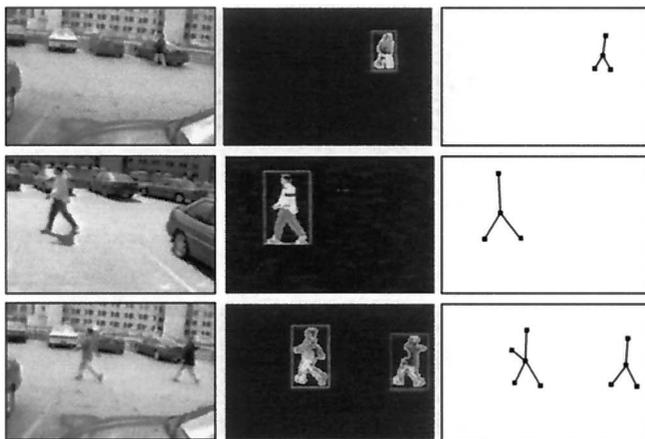
5. Analysis

This motion analysis scheme has been tried on a database of video sequences of people walking and running at outdoor. There are approximately 20 video sequences in each category, with pixels on target ranging from ≈ 50 to ≈ 400 caused by the distance from camera to the target. The targets are a mixture of adults and children. The end-to-end process of MTD, target

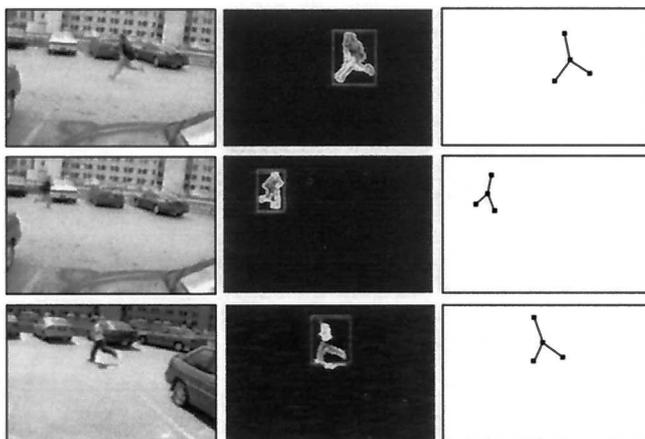
pre-processing, and motion analysis was performed on an SGI O2 machine containing an R10000 175 MHz processor. Figure 8 shows examples of detected region by the MTD and its skeleton for walking and running.

Figure 9 shows histograms of the peaks of the power spectrum for each of the video streams. It is clear from Fig. 9 (a) that the low frequency noise would cause a serious bias if motion classification were attempted. However, Fig. 9 (b) shows how effective the pre-emphasis filter is in removing this noise. It also shows how it is possible to classify motion in terms of walking or running based on the frequency of the cyclic motion. The average walking frequency is 1.75 [Hz] and for running it is 2.875 [Hz]. A threshold frequency of 2.0 [Hz] correctly classifies 97.5% of the target motions. Note that these frequencies are twice the actual footstep frequency because only the visually leftmost leg is considered. Another point of interest is that the variance of running frequencies is greater than that of walking frequencies, so it could be possible to classify different “types” of running such as jogging or sprinting.

For each video sequence, the average inclination $\bar{\phi}$

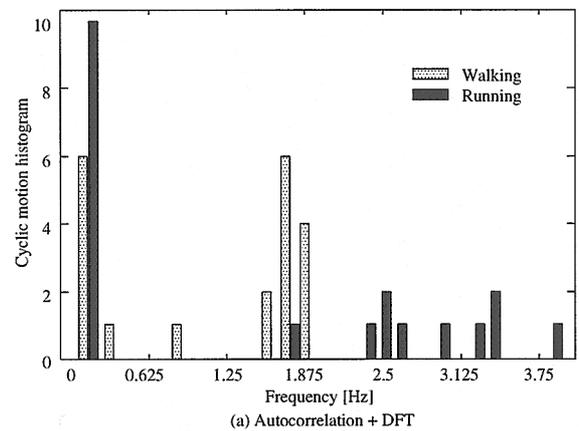


(a) walking

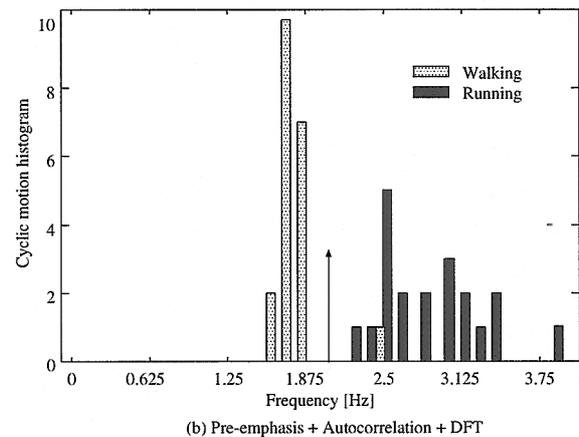


(b) running

Fig. 8 Target skelton for walking and running.



(a) Autocorrelation + DFT



(b) Pre-emphasis + Autocorrelation + DFT

Fig. 9 Histogram of cyclic motion frequency peaks. (a) The bias in θ_n often produces a frequency peak which is significantly higher than the peak produced by cyclic motion. (b) The pre-emphasis filter effectively removes this noise.

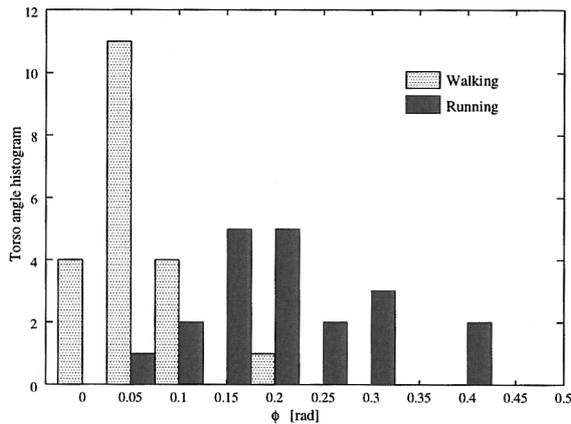


Fig. 10 Average inclination histogram of torso for classification.

of the upper extremal point (or torso) was determined. These values are shown in Fig. 10. It can be seen that the forward leaning of a running figure can be clearly distinguished from the more vertical posture of a walking one. A threshold value of 0.15 [rads] correctly classifies 90% of the target motions.

6. Conclusion

Analyzing human motion for video applications is a complex problem. Real-world implementations will have to be computationally inexpensive and be applicable to real scenes in which targets are small and data is noisy. The notion of using a target's boundary to analyze its motion is a useful one under these conditions. Algorithms need only be applied to a small number of pixels and internal target detail, which may be sketchy, becomes less important.

This paper presents the approach of "star" skeletonization by which the component parts of a target with internal motion may easily, if grossly, be extracted. Further, two analysis techniques have been investigated which can broadly classify human motion. Body inclination can be measured from the "star" skeleton to determine the posture of the human, which derives clues as to the type of motion being executed. In addition, cyclic analysis of extremal points provides a very clean way of broadly distinguishing human motion in terms of walking and running and potentially even different types of gait.

In the future, this analysis technique will be applied to more complex human motions such as crawling, jumping, and so on. It may even be applied to the gaits of animals.

Acknowledgement

The authors would like to thank the CMU VSAM team members: Rober T. Collins, David Duggins, Yanghai Tsin, Raju Patil, David Tolliver, Osamu Hasegawa,

Nobuyoshi Enomoto, Yong-Tae Do and Alan Lee, for their tireless efforts and good humor.

References

- [1] C. Anderson, P. Burt, and G. van der Wal, "Change detection and tracking using pyramid transformation techniques," *SPIE — Intelligent Robots and Computer Vision*, vol.579, pp.72–78, 1985.
- [2] J. Barron, D. Fleet, and S. Beauchemin, "Performance of optical flow techniques," *Int. J. Comput. Vis.*, vol.12, no.1, pp.42–77, Jan. 1994.
- [3] J. Davis and A. Bobick, "The representation and recognition of human movement using temporal templates," *Proc. IEEE CVPR 97*, pp.928–934, 1997.
- [4] E. Grimson and P. Viola, "A forest of sensors," *DARPA - VSAM Workshop*, Nov. 1997.
- [5] I. Haritaoglu, D. Harwood, and L.S. Davis, " W^4 who? when? where? what? A real time system for detecting and tracking people," *FGR98*, pp.222–227, 1998.
- [6] S. Ju, M. Black, and Y. Yacoob, "Cardboard people: A parameterized model of articulated image motion," *Int. Conference on Face and Gesture Analysis*, 1996.
- [7] T. Kanade, R. Collins, A. Lipton, P. Anandan, and P. Burt, "Cooperative multisensor video surveillance," *Proc. DARPA Image Understanding Workshop 1997*, vol.I, pp.3–10, 1997.
- [8] M. Oren, C. Papageorgiou, P. Sinha, E. Osuna, and T. Poggio, "Pedestrian detection using wavelet templates," *Proc. IEEE CVPR 97*, pp.193–199, 1997.
- [9] P. Tsai, M. Shah, K. Ketter, and T. Kasparis, "Cyclic motion detection for motion based recognition," *Pattern Recognition*, vol.27, no.12, pp.1591–1603, 1994.
- [10] C. Wren, A. Azarbayejani, T. Darrell, and A. Pentland, "Pfinder: Real-time tracking of the human body," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol.19, no.7, pp.780–785, July 1997.



Hironobu Fujiyoshi received the Ph.D. degree in electrical engineering from Chubu University, Japan, in 1997. For his thesis, he developed a fingerprint verification method using spectrum analysis, which has been incorporated into a manufactured device sold by a Japanese security company. He is a Member of Faculty at the Department of Computer Science, Chubu University. From 1997 to 2000 he was a post-doctoral fellow at the

Robotics Institute of Carnegie Mellon University, Pittsburgh, PA, USA, working on the DARPA Video Surveillance and Monitoring (VSAM) effort and the humanoid vision project for the Honda Humanoid Robot. He performs research in the areas of real-time object detection, tracking, and recognition from video.



Alan J. Lipton received the Ph.D. degree in electrical and computer systems engineering from Monash University, Melbourne, Australia, in 1996. For his thesis, he studied the problem of mobile robot navigation by natural landmark recognition using on-board vision sensing. He is currently with Diamondback Vision, Inc., Reston, VA. From 1997 to 2000, he was on the faculty at the Robotics Institute of Carnegie Mellon University, Pittsburgh,

PA, USA, where technical co-director of the DARPA Video Surveillance and Monitoring (VSAM) effort. Under this program, he performed research in the areas of real-time object detection, tracking, and recognition from video.



Takeo Kanade received the B.E. degree in electrical engineering in 1968, the M.E. degree in 1970, and the Ph.D. degree in 1973 from Kyoto University, Japan. Currently, he is Helen Whitaker Professor of Computer Science at Carnegie Mellon University, Pittsburgh, PA, USA. Dr. Kanade was a recipient of the Robotics Industry Association, Joseph F. Engelberger Award in 1995, the Japan Robotics Association, JARA Award in 1997, the

Yokogawa Prize at the International Conference on Multi Sensor Fusion and Integration for Intelligent Systems in 1997, the Hip Society, Otto AuFranc Award in 1998, the Hosokawa Kikin Foundation Award in 1994, and the Marr Prize at The Third International Conference on Computer Vision in December 1990. He was also selected as the author of one of the most influential papers that appeared in the Artificial Intelligence journal in the last ten years in 1992. He is Founding Chief Editor of the *International Journal of Computer Vision*. He is a Member of the National Academy of Engineering and a Fellow of the American Association for Artificial Intelligence (AAAI).