

Velocity adaptation of spatio-temporal receptive fields for direct recognition of activities: an experimental study

Ivan Laptev*, Tony Lindeberg

Computational Vision and Active Perception Laboratory (CVAP), Department of Numerical Analysis and Computer Science, KTH, SE-100 44 Stockholm, Sweden

Received 26 September 2002; received in revised form 27 June 2003; accepted 2 July 2003

Abstract

This article presents an experimental study of the influence of velocity adaptation when recognizing spatio-temporal patterns using a histogram-based statistical framework. The basic idea consists of adapting the shapes of the filter kernels to the local direction of motion, so as to allow the computation of image descriptors that are invariant to the relative motion in the image plane between the camera and the objects or events that are studied. Based on a framework of recursive spatio-temporal scale-space, we first outline how a straightforward mechanism for local velocity adaptation can be expressed. Then, for a test problem of recognizing activities, we present an experimental evaluation, which shows the advantages of using velocity-adapted spatio-temporal receptive fields, compared to directional derivatives or regular partial derivatives for which the filter kernels have not been adapted to the local image motion.

© 2003 Elsevier B.V. All rights reserved.

Keywords: Motion; Spatio-temporal filtering; Scale-space; Recognition

1. Introduction

A recent approach for recognition consists of computing statistical descriptors of receptive field responses. In particular, histogram-based schemes of derivative operators have emerged as an interesting alternative for formulating recognition schemes for static as well as time-dependent image data [1–7]. Computing responses of local spatio-temporal receptive fields involves filtering in both space and time. This naturally rises the question of how to express filtering operations in space–time.

When analysing spatio-temporal image data, one observation that can be made is that temporal events can often be characterized by their extents over time in a similar manner as spatial structures have their characteristic scales in space. This motivates and emphasizes the need for analysing spatio-temporal data at different scales, both with respect to time and space [8–14].

The temporal domain, however, also has a number of specific properties, which differ from spatial data, and which

must be taken into account explicitly. A basic constraint on real-time processing is that the time direction is causal, and real-time algorithms may only access information from the past [10,13]. Another difference concerns the classes of characteristic transformations that influence the data. Whereas perspective transformations have a high influence on the image data in the spatial image domain, one of the most important sources of changes in the temporal dimension is due to motion between the observer and the patterns that are studied. This is shown in Fig. 1, where the spatio-temporal pattern of a walking person is influenced by the relative motion of the camera (Fig. 1b and c). If separable spatial filtering is extended to the temporal domain, we observe that the filter responses are highly dependent on the relative motion between the person and the camera (Fig. 1d and f).

When interpreting image data, it is important to base the analysis on image representations that are invariant to the external imaging conditions. Hence, it is important to construct representations of spatio-temporal patterns that are independent of the relative motion between the patterns and the observer. Previous work has addressed this problem by first stabilizing patterns of interest in the field of view, and then computing spatio-temporal descriptors using

* Corresponding author.

E-mail addresses: laptev@nada.kth.se (I. Laptev); tony@nada.kth.se (T. Lindeberg).

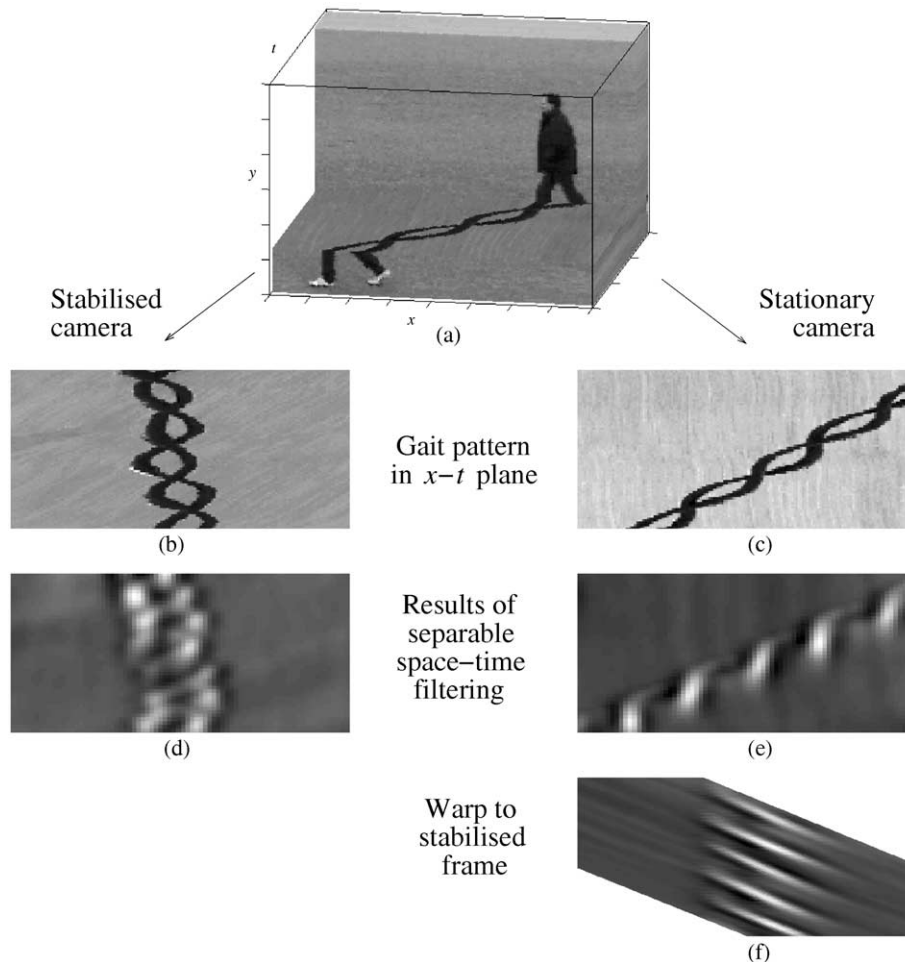


Fig. 1. Spatio-temporal image of a walking person (a) depends on the relative motion between the person and the camera (b)–(c). If this motion is not taken into account, spatio-temporal filtering (here, the second order spatial derivative) results in highly different responses as illustrated in (d) and (e). Manual stabilization of the pattern in (e) shown in (f) makes the difference more explicit for comparisons with (d).

a fixed set of filters [7,15] for related stabilization approaches. Camera stabilization, however, may not always be available, for example, in situations with multiple moving objects, moving backgrounds or in cases where initial segmentation of the patterns of interest cannot be done without (preliminary) recognition.

The main aim of this work is to define and compute spatio-temporal descriptors that compensate for the relative motion between the pattern and the observer and do not rely on external camera stabilization. This is achieved by local velocity adaptation of receptive fields. In Section 2 we first introduce velocity-adapted filtering using the framework of spatio-temporal scale-space. Then in Section 3, a mechanism for performing local velocity adaptation is described. By integration with a histogram-based statistical framework in Section 4, we then consider a test problem of recognizing activities and show how velocity adaptation results in a considerable increase in recognition performance compared to two other receptive field representations not involving velocity adaptation. Section 5 concludes the paper with a summary and discussion.

1.1. Related work

Velocity adaptation of spatio-temporal receptive fields follows the idea of shape adaptation in the spatial domain, which has previously been considered in Refs. [16–22]. In the spatio-temporal domain, adaptive spatio-temporal filters have been studied in Refs. [12,23–26]. Nagel and Gehrke [24] proposed an adaptation scheme close to ours and used it for robust estimation of optic flow.

With regard to recognition, this work relates to histogram-based methods first proposed in the spatial domain by Swain and Ballard [1] using color histograms computed from single pixel responses. Extensions to receptive field histograms were later presented in Refs. [2,3,5,6]. Specifically, combinations of automatic scale selection in the spatial domain [27] with Gaussian derivative-based recognition schemes have been presented in Refs. [3,5]. In the spatio-temporal domain, histogram-based approaches have been used for the recognition of activities in Refs. [4,7]. Here, we build upon this work and

show how the performance of spatio-temporal recognition schemes can be increased by velocity adaptation.

2. Spatio-temporal scale-space representation

The image data we analyse is a spatio-temporal image sequence, in the continuous case modeled as a function $f : \mathbb{R}^2 \times \mathbb{R} \rightarrow \mathbb{R}$ or in the discrete case as $f : \mathbb{Z}^2 \times \mathbb{Z} \rightarrow \mathbb{Z}$. From this signal, a separable spatio-temporal scale-space representation $L(x, y, t; \sigma^2, \tau^2)$ is defined¹ by separable convolution of f with a set of spatial smoothing kernels $g(x, y; \sigma^2)$ with variances σ^2 and a set of temporal smoothing kernels $h(t; \tau^2)$ with variances τ^2 . Hence, L is a function that represents the image sequence at different scales of observations over both space and time.

For continuous data, the natural choice of a spatial smoothing kernel is the Gaussian kernel [8,9,11,14]. Regarding continuous time, we may model the temporal smoothing operation either by a non-causal Gaussian kernel, or as a causal Gaussian kernel on a logarithmically transformed temporal domain [10,13]. For discrete data, a canonical spatial scale-space concept originates from the discrete analogue of the Gaussian kernel [11]

$$T(x, y; \sigma^2) = e^{-2\sigma^2} I_x(\sigma^2) I_y(\sigma^2) \quad (x, y) \in \mathbb{Z}^2 \quad (1)$$

where I_x and I_y denote the modified Bessel functions of integer order [28]. Regarding discrete time, a natural and computationally efficient scale-space representation can be computed by coupling first-order recursive filters in cascade [11,13]

$$L^{(k+1)}(x, y, t) = \frac{1}{1 + \mu} (L^{(k)}(x, y, t) + \mu L^{(k+1)}(x, y, t - 1)), \quad (2)$$

where k denotes the number of temporal smoothing stages. The corresponding temporal smoothing kernel with coefficients $c_n \geq 0$ obeys temporal causality by only accessing data from the past. Moreover, this kernel is normalized to $\sum_{n=-\infty}^{\infty} c_n = 1$ and has mean value $m = \sum_{n=-\infty}^{\infty} n c_n = \mu$ and variance $\tau^2 = \sum_{n=-\infty}^{\infty} (n - m)^2 c_n = \mu^2 + \mu$. By coupling k such recursive filters in Eq. (2) in cascade, we obtain a filter with mean $m_k = \sum_{i=1}^k \mu_i$ and variance $\tau_k^2 = \sum_{i=1}^k \mu_i^2 + \mu_i$.

It can be shown that if for a given variance τ^2 we let $\mu_i = \tau^2/K$ become successively smaller by increasing the number of filtering steps K , then the filter kernel approaches the Poisson kernel [12], which corresponds to the canonical temporal scale-space concept having a continuous scale parameter on a discrete temporal domain. Another practical advantage of the recursive filtering scheme in Eq. (2) is that

¹ Here, $(x, y) \in \mathbb{R}^2$ (or \mathbb{Z}^2) denote the spatial coordinates, $t \in \mathbb{R}$ (or \mathbb{Z}) denotes time, $\sigma^2 \in \mathbb{R}^+$ is the spatial scale parameter and $\tau^2 \in \mathbb{R}^+$ is the temporal scale parameter.

it enables the computation of temporal scale-space representations without need of buffering previous time frames.

2.1. Transformation properties under motion

To describe the spatio-temporal smoothing step, we will henceforth use covariance matrices of filter kernels. For a separable smoothing kernel, with a spatial variance σ^2 and a temporal variance τ^2 , the covariance matrix is diagonal:

$$\Sigma = \begin{pmatrix} C_{xx} & C_{xt} & C_{xt} \\ C_{xy} & C_{yy} & C_{yt} \\ C_{xt} & C_{yt} & C_{tt} \end{pmatrix} = \begin{pmatrix} \sigma^2 & & \\ & \sigma^2 & \\ & & \tau_k^2 \end{pmatrix}. \quad (3)$$

A limitation of using a separable scale-space for analysing motion patterns, however, originates from the fact that this scale-space concept is not closed under 2D motions in the image plane. For a 2D Galilean motion

$$\begin{pmatrix} x' \\ y' \\ t' \end{pmatrix} = \begin{pmatrix} 1 & 0 & v_x \\ 0 & 1 & v_y \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} x \\ y \\ t \end{pmatrix} \quad (4)$$

the covariance matrix of the smoothing kernel transforms as [12,23]

$$\begin{pmatrix} C'_{xx} & C'_{xt} & C'_{xt} \\ C'_{xy} & C'_{yy} & C'_{yt} \\ C'_{xt} & C'_{yt} & C'_{tt} \end{pmatrix} = \begin{pmatrix} 1 & 0 & v_x \\ 0 & 1 & v_y \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} C_{xx} & C_{xt} & C_{xt} \\ C_{xy} & C_{yy} & C_{yt} \\ C_{xt} & C_{yt} & C_{tt} \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ v_x & v_y & 1 \end{pmatrix} \quad (5)$$

and spatio-temporal derivatives transform according to

$$\partial_{x'} = \partial_x \quad \partial_{y'} = \partial_y \quad \partial_{t'} = -v_x \partial_x - v_y \partial_y + \partial_t. \quad (6)$$

Hence, if we consider separable smoothing kernels only and if we do not take the transformation property of spatio-temporal derivatives into explicit account, it will not be possible to perfectly match the spatio-temporal scale-space representations for different amounts of motion.

2.2. Scale-space with velocity adaptation

A natural way of defining a scale-space that is closed under Galilean motion in the image plane, is by considering a scale-space representation that is parameterized by the full family of (positive definite) covariance matrices [12,14,23]. In terms of implementation, there are two basic ways of computing such a scale-space—either by transforming the smoothing kernels themselves, or by transforming the input image prior to smoothing (see Fig. 2). In this work, the latter approach is taken, and for reasons of simplicity and computational efficiency, we restrict the set of image

$$\begin{array}{ccc}
\partial_{x^\alpha y^\beta t^\gamma} L(x, y, t; \Sigma) & \xrightarrow{\left\{ \begin{array}{l} \text{Galilean transformation} \\ \text{of receptive fields} \end{array} \right\}} & \partial_{x'^\alpha y'^\beta t'^\gamma} L'(x', y', t'; \Sigma) \\
\uparrow & & \uparrow \\
* g(x, y, t; \Sigma) & & * g'(x', y', t'; \Sigma') \\
| & & | \\
f(x, y, t) & \xrightarrow{\left\{ \begin{array}{l} \text{Galilean transformation} \\ \text{of space-time} \end{array} \right\}} & f'(x', y', t')
\end{array}$$

Fig. 2. A pre-requisite for perfect matching of spatio-temporal receptive field responses for different amounts of motion is that the image representation is closed under motions in the image domain. The aim of the velocity adaptation mechanism is to allow for such closedness, and to permit the construction of a velocity invariant recognition scheme.

velocities to integer multiples of the pixel size. Thus, in combination with a spatial smoothing step

$$L^{(0)}(x, y, t; \sigma^2) = T(x, y; \sigma^2) f(x, y, t), \quad (7)$$

a set of velocity-adapted time-recursive smoothing steps is computed according to

$$\begin{aligned}
L^{(k+1)}(x, y, t; \sigma^2) &= \frac{1}{1 + \mu_k} (L^{(k)}(x, y, t; \sigma^2) \\
&\quad + \mu_k L^{(k+1)}(x - v_x, y - v_y, t - 1; \sigma^2)),
\end{aligned} \quad (8)$$

where k represents the level of temporal smoothing corresponding to the convolution with a set of temporal kernels with variances τ_k^2 . The scale-space concept we make use of, will hence be parameterized by a spatial scale parameter σ^2 , a temporal scale parameter τ^2 and a set of discrete image velocities $(v_x, v_y)^T$.

The result of applying such velocity-adapted filters to spatio-temporal image data is shown in Fig. 3. Here, a synthetic pattern with one spatial and one temporal dimension has been filtered using different values of velocity parameter v . As can be seen, depending on the value of v , the filtering is able to emphasize either the moving pattern (Fig. 3b) or the stationary background (Fig. 3c).

3. A mechanism for local velocity adaptation

If we want to interpret events independently of their relative motion to the camera, one approach is to adapt the receptive fields *globally* with respect to the velocity of the events in the field of view. This approach also corresponds to camera stabilization followed by non-adapted filtering. As shown in Fig. 3b, the result of filtering with globally adapted receptive fields with $v = -1$ indeed enhances the structure of the moving pattern. However, the stationary pattern is suppressed and it follows that global velocity adaptation is not able to handle multiple motions. Moreover, global velocity adaptation is likely to fail if the external velocity information is incorrect (Fig. 3d).

To address these problems, we propose to make use of *local* velocity adaptation of receptive fields. The main idea is to obtain information about motion in the local

neighborhood and to use this information for velocity adaptation of receptive fields in the same neighborhood.

Before proceeding to specific schemes for local velocity adaptation in space-time, however, let us observe that there are two main approaches for handling multiple image velocities. One approach is to consider the entire ensemble of receptive fields over image motions as the representation, while the other is to select receptive field outputs corresponding to a single motion estimate. From basic arguments, the first approach can be expected to be more robust in critical situations (compare with biological vision systems), while the second approach followed in this work could be expected to be more accurate and also computationally more efficient on a serial architecture.

The mechanism we will use for accomplishing local velocity adaptation is inspired by related work on automatic scale selection [27] extended to a multi-parameter scale-space [12] as well as by motion energy approaches for computing optic flow [29,30]. Given a set of image velocities, the normalized Laplacian response is computed for each image velocity in a motion compensated frame (8)

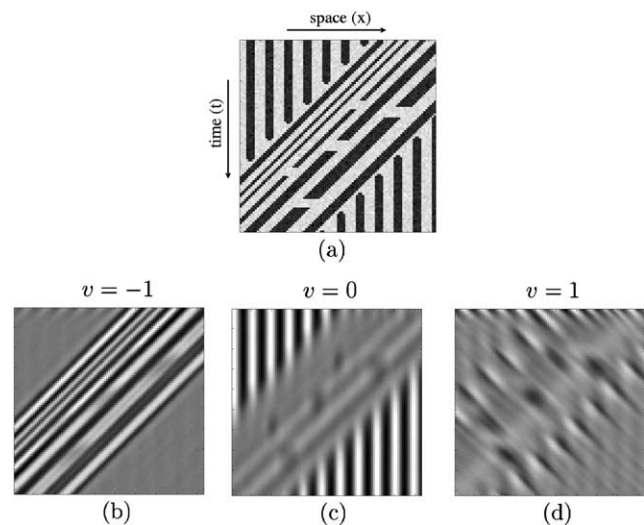


Fig. 3. The effect of global velocity adaptation for a synthetic spatio-temporal pattern in (a). (b)–(d) Convolution of (a) with spatio-temporal second-order derivative operators with $\sigma^2 = 32$, $\tau^2 = 32$ and velocity parameters $v = -1, 0, 1$, respectively. Note, that depending on the velocity parameter, global velocity adaptation emphasizes either the moving pattern (b) or the stationary pattern (c).

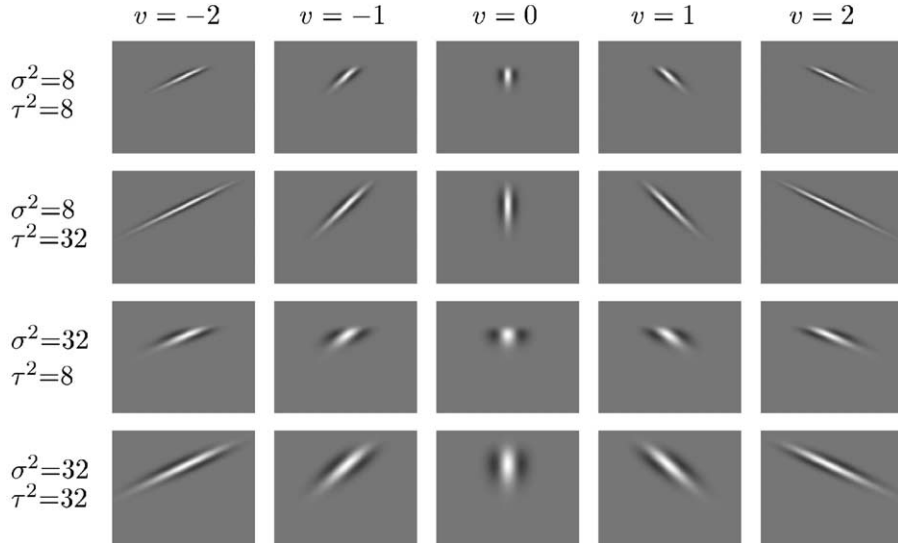


Fig. 4. Spatio-temporal filters L_{xx} computed from a velocity-adapted spatio-temporal scale-space for a 1 + 1D image pattern, for different values of the velocity parameter v , the spatial scale σ^2 and the temporal scale τ^2 .

in the spatio-temporal scale-space. Then, for each scale, a motion estimate is computed from the velocity $(v_x, v_y)^T$ that maximizes the normalized derivative response

$$(\hat{v}_x, \hat{v}_y)^T(x, y, t)^{(k)} = \operatorname{argmax}_{v_x, v_y} (\nabla_{\text{norm}}^2 L^{(k)}(x, y, t; \sigma^2, v_x, v_y))^2, \quad (9)$$

where $\nabla_{\text{norm}}^2 = \sigma^2(\partial_{xx} + \partial_{yy})$ is a scale-normalized Laplacian operator in space. This approach is equivalent to

the application of a set of velocity-adapted Laplacian operators (Fig. 4) at each spatio-temporal scale, and selecting the motion estimate from the spatio-temporal filter parameters that gives the maximum response. While one could also consider the use of optic flow estimation schemes for computing the velocity estimates [24], a main reason why we here consider maximization of normalized receptive field responses over image velocities is that

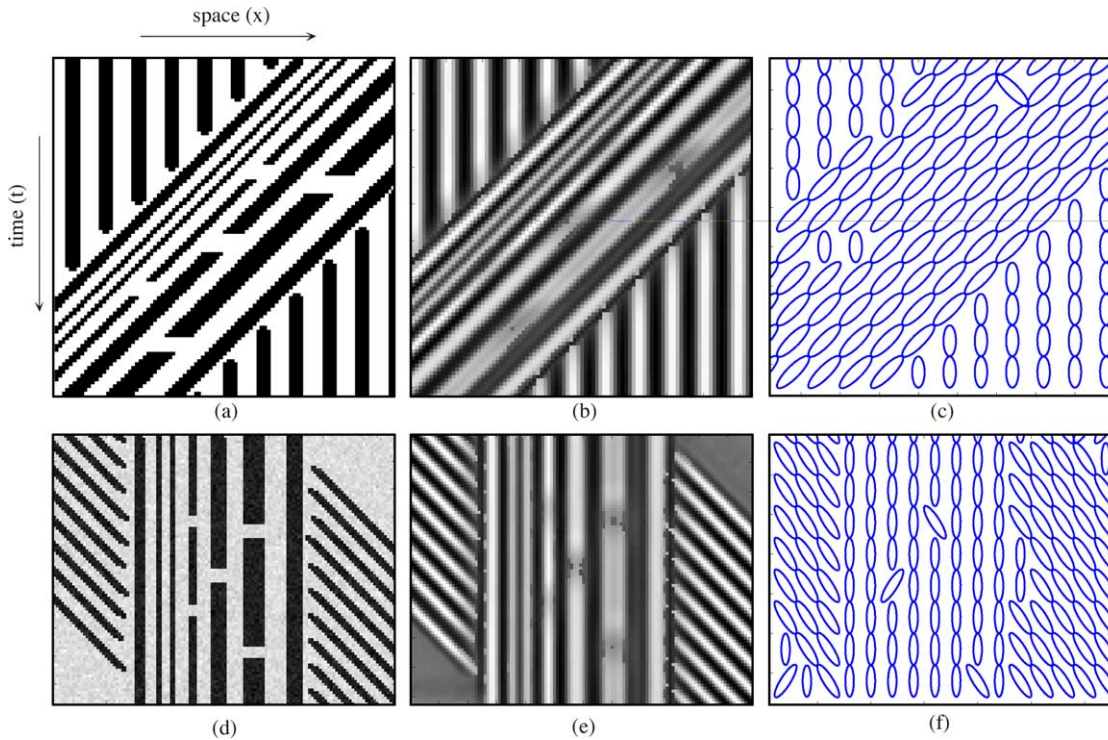


Fig. 5. Results of filtering original patterns in (a) and (d) using the proposed *local velocity adaptation* are illustrated in (b) and (e), respectively. The orientation of the ellipses in (c) and (f) show the chosen velocity at each point of the pattern. Note that filtering with local velocity adaptation preserves the details of the moving and stationary pattern. The similarity of the filter responses in (b) and (e) also illustrates the independence of the filtering results with respect to the amount of camera motion.

a similar mechanism, when extended to maximization over spatial scales and temporal scales, can also be used for performing simultaneous automatic selection of spatial scales and temporal scales [27,31].

Fig. 5 shows the results of local velocity adaptation for a synthetic spatio-temporal pattern (Fig. 5a) and its Galilean transformation (Fig. 5d). From the responses of velocity-adapted receptive fields and from the ellipses displaying the selected orientation of filters in space-time, it is apparent that the proposed filtering scheme adapts to the local motion and enhances structures both in the moving pattern and in the static background. Moreover, by comparing the results

in Fig. 5e and f, we can visually confirm the invariance of locally adapted receptive field responses with respect to the Galilean transformation of the pattern or, equivalently, to the relative motion between the pattern and the camera.

Application of the local velocity adaptation to a sequence with a walking person is shown in Fig. 6. Note, that filtering here has been done in three dimensions while for the purpose of demonstration, the results are shown only for one $x-t$ -slice of a spatio-temporal cube (see Fig. 1). As for the synthetic pattern above, we observe successful adaptation of filter kernels to the motion structure of a gait pattern (Fig. 6c and d). The results in Fig. 6e–g also demonstrate approximative

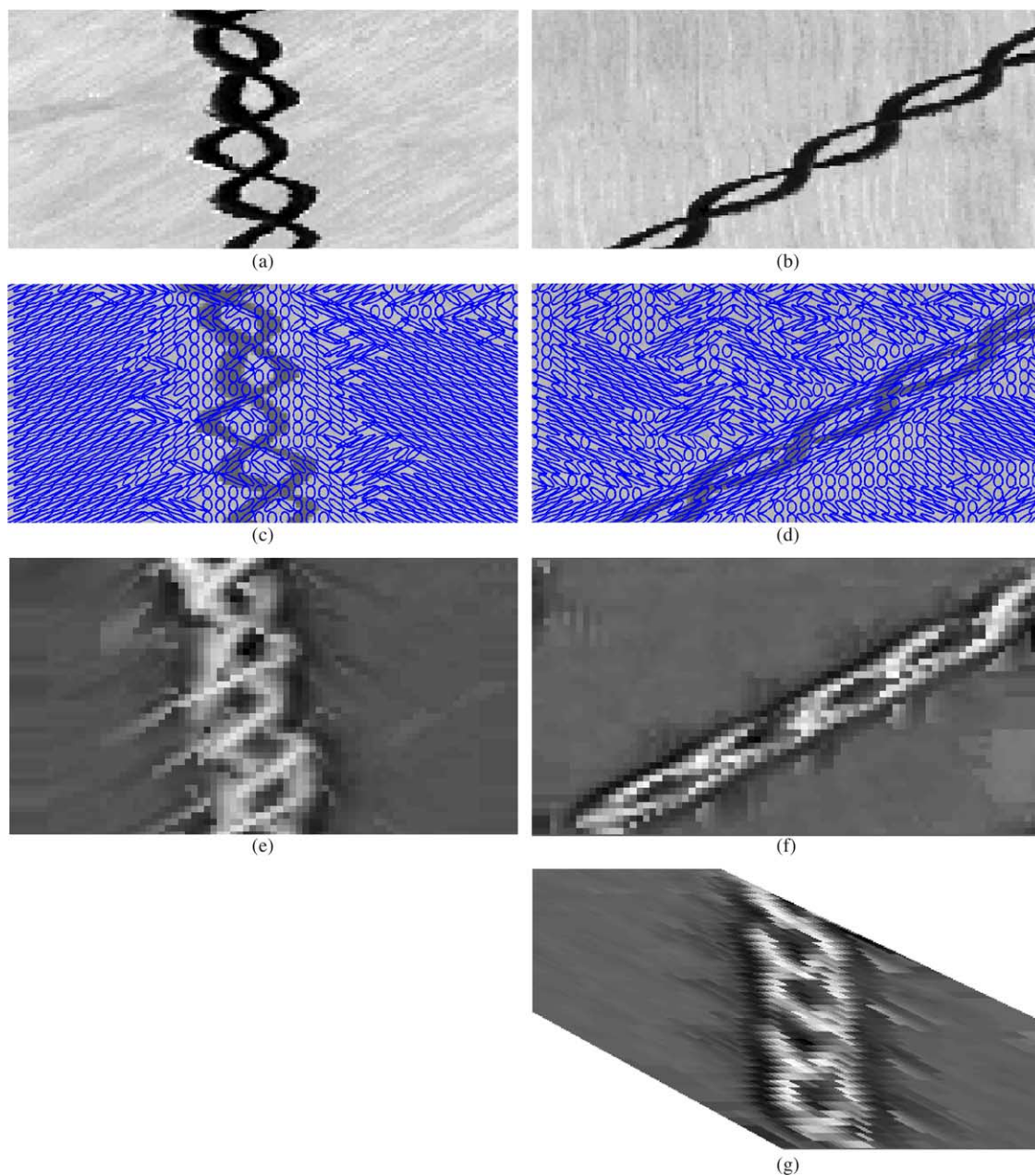


Fig. 6. Spatio-temporal filtering with local velocity adaptation applied to a gait pattern recorded with a stabilized camera (a) and a stationary camera (b) (see Fig. 1 for comparison); (c) and (d) velocity adapted shape of filter kernels; (e) and (f) results of filtering with a second-order derivative operator; (g) warped version of (f) showing high similarity with (e).

invariance of filter responses with respect to camera motion. The desired effect of the proposed local velocity adaptation is especially evident when these results are compared to the results of separable filtering as shown in Fig. 1d–f.

3.1. Comparison with steerable filters

When computing spatio-temporal derivatives, we perform velocity adaptation of both the shapes of smoothing kernels and the derivatives according to Eqs. (5) and (6). An alternative approach that is more efficient but less accurate consists of separable smoothing step followed by adaptation of the derivatives only. Such a scheme is closely related to steerable filters [32] for computing higher-order spatial derivatives in a rotationally invariant way. To differentiate these two approaches, we will refer to them as *velocity-adapted filtering* and *velocity-steered filtering*.

To compare these two alternatives and to illustrate the importance of shape adaptation of filter kernels, we will first compare the results of filtering a synthetic prototype of a moving spatio-temporal impulse. The original signal is shown in Fig. 7a in two spatial and one temporal dimensions. Fig. 7b shows the result of computing a partial spatio-temporal derivative ∂_{xxt} using velocity-adapted filtering. With positive and negative filter values represented by different colors, we can visually confirm the correctness of the resulting shape. On the contrary, computation of the same derivative using velocity-steered filtering (Fig. 7c)

results in a different and incorrect shape. A similar result is obtained when filtering is performed without adaptation of neither the smoothing kernels nor the derivatives (Fig. 7d).

In Section 4, we apply these filtering schemes to a recognition task and give their quantitative comparison as well as emphasize the importance of velocity-adapted filtering in practice.

4. Histogram-based recognition

The responses of spatio-temporal derivatives describe the structure of local spatio-temporal neighborhoods and therefore can be used for discriminating between motion patterns with different spatio-temporal structure. Higher order derivatives provide for a more rich and discriminative representation while lower order derivatives are less sensitive to noise and other sources of variations in the pattern. Moreover, the velocity adaptation of derivatives makes them independent of the first order motion but still enables to capture and represent the motion of higher order. Since the relative motion between the camera and the pattern can be approximated by the constant velocity (at least for a short period of time), the velocity adaptation enables to compute descriptors independently of the relative camera motion.

Computing the statistics of derivative responses over all points of the image sequence is attractive due to

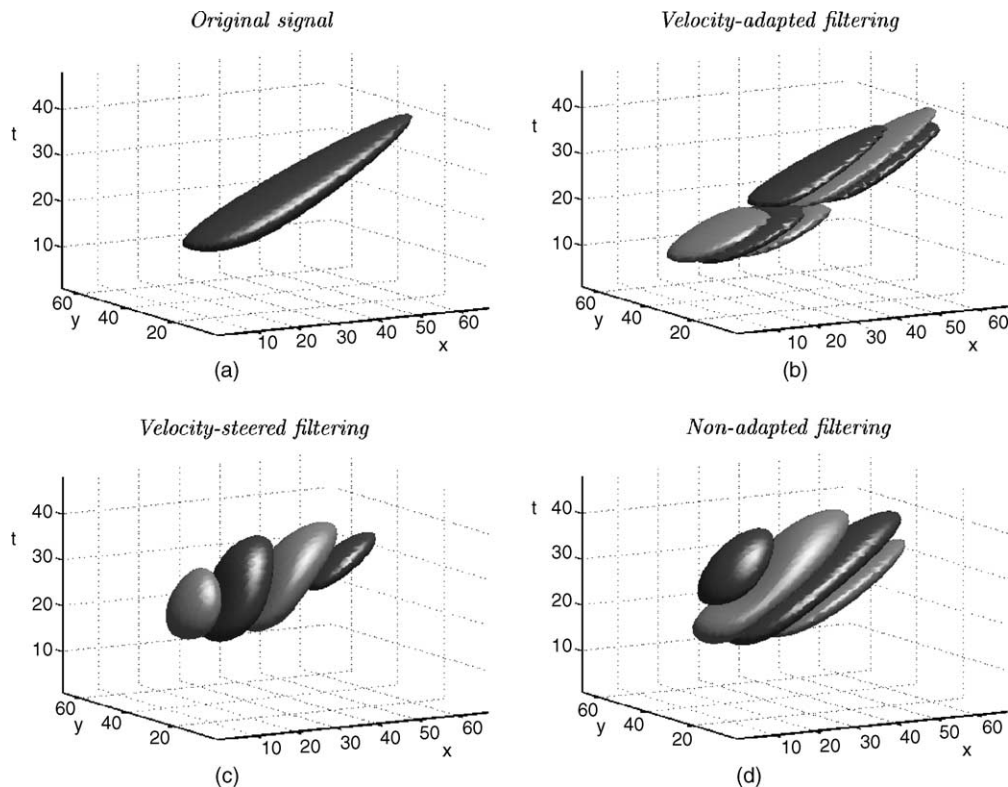


Fig. 7. (a) Prototype spatio-temporal blob signal with velocity $v_x = 2$. (b)–(d) Responses to the ∂_{xxt} -derivative operator when using (b): velocity-adapted filters; (c): velocity-steered filters; (d): non-adapted filters. A correct shape of the filter response is obtained only for the case of velocity-adapted filtering.

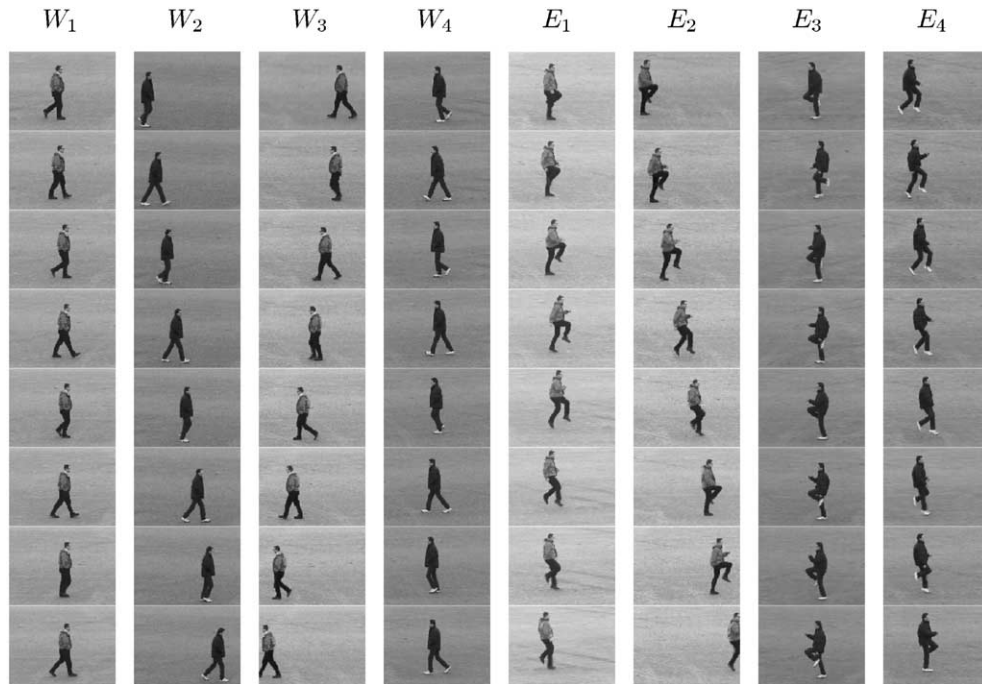


Fig. 8. Test sequences of people walking W_1 – W_4 and people performing an exercise E_1 – E_4 . Whereas the sequences W_1 , W_4 , E_1 , E_3 were taken with a manually stabilized camera, the other four sequences were recorded using a stationary camera.

the invariance of such a descriptor with respect to the translations of the pattern in space and time. Hence, following Refs. [2–4,7], we represent image patterns by histograms of receptive field responses. For this purpose, we

use velocity-adapted spatio-temporal derivative operators up to order four and collect histograms of these at different spatial and temporal scales. For simplicity, we restrict ourselves to 1D histograms for each type of filter response.

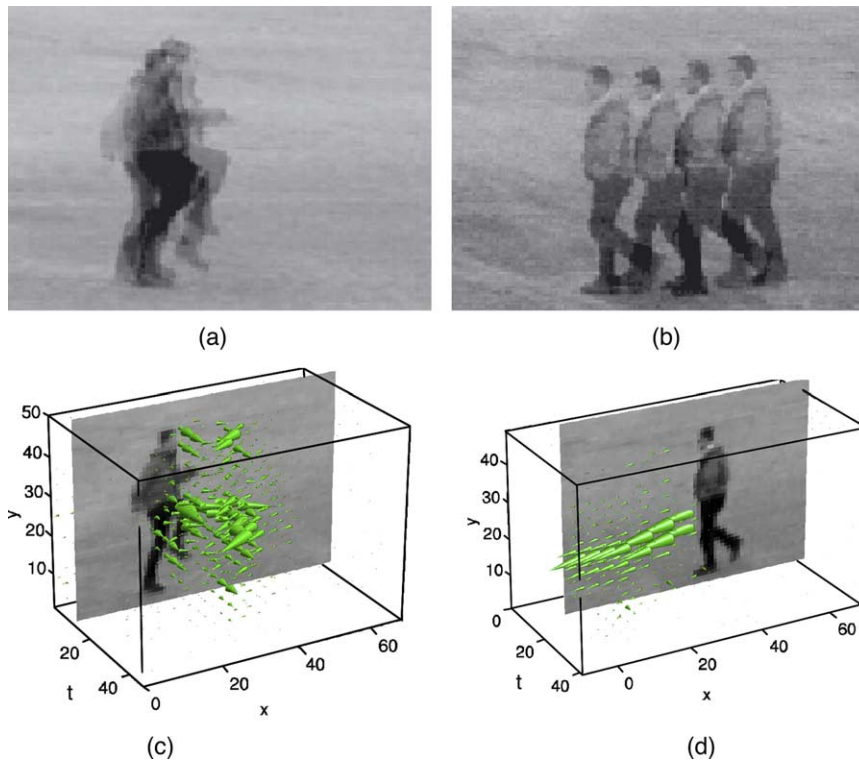


Fig. 9. Results of local velocity adaptation for image sequences recorded with a manually stabilized camera (a), and with a stationary camera (b). Directions of cones in (c) and (d) correspond to the velocity chosen by the proposed adaptation algorithm. The size of the cones corresponds the value of the squared Laplacian $((\partial_{xx} + \partial_{yy})L(x, y, t; \sigma, \tau))^2$ at the selected velocities.

To achieve independence with respect to the direction of motion (left/right or up/down) and the sign of the spatial grey-level variations, we simplify the problem by only considering the absolute values of the filter responses. Moreover, to emphasize the parts of the histograms that correspond to stronger spatio-temporal responses, we use heuristics and weight the accumulated histograms $H(i)$ by a function $f(i) = i^2$ resulting in $h(i) = i^2 H(i)$.

4.1. Experimental setup

As a test problem we have chosen a data set with image sequences containing people performing actions of type *walking* $W_1 \dots W_4$ and *exercise* $E_1 \dots E_4$ as shown in Fig. 8. Some of the sequences were taken with a stationary camera, while the others were recorded with a manually stabilized

camera. Each of these 4 s long sequences were subsampled to a spatio-temporal resolution of $80 \times 60 \times 50$ pixels and convolved with a set of spatio-temporal smoothing kernels for all combinations of seven velocities $v_x = -3 \dots 3$, five spatial scales $\sigma^2 = \{2, 4, 8, 16, 32\}$ and five temporal scales $\tau^2 = \{2, 4, 8, 12, 16\}$.

For each spatial scale σ_i , velocity adaptation was performed according to Eq. (9) at scale level σ_{i+1} . Since in our examples the relative camera motion was mostly horizontal, we maximized Eq. (9) over v_x only. The result of this adaptation for the sequences W_2 and E_1 is shown in Fig. 9.

To represent the patterns, we accumulated histograms of derivative responses for each combination of scales and each type of derivative. For the purpose of evaluation, separate histograms were accumulated over (i) velocity-adapted derivative responses; (ii) velocity-steered

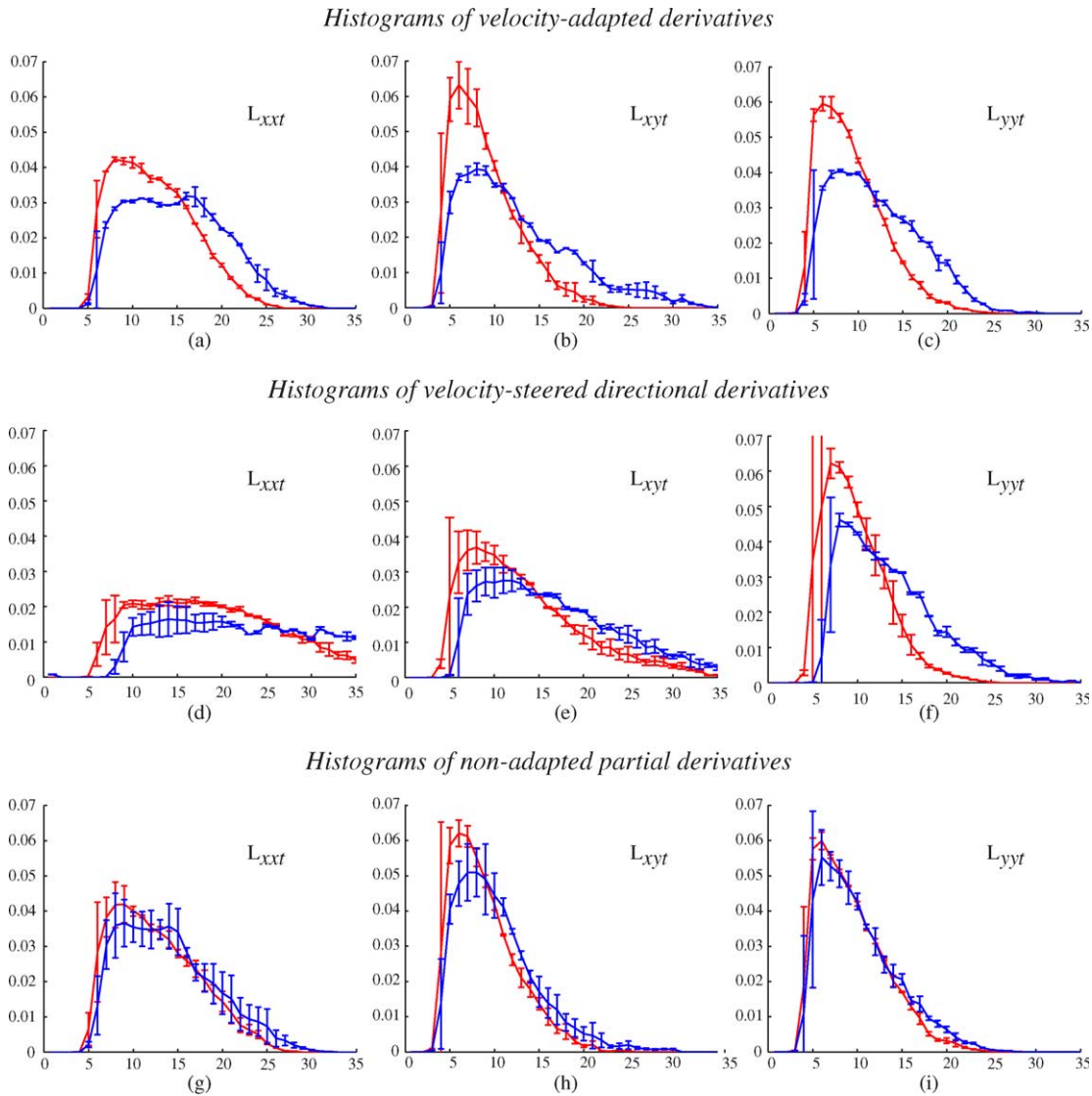


Fig. 10. Means and variances of histograms for the activities ‘walking’ (red) and ‘exercise’ (blue). (a)–(c) Histograms of velocity-adapted derivatives L_{xxt} , L_{xyt} , L_{yyt} ; (d)–(f) histograms of velocity-steered directional derivatives L_{xxt} , L_{xyt} , L_{yyt} ; (g)–(i) histograms of non-adapted partial derivatives L_{xxt} , L_{xyt} , L_{yyt} . As can be seen, the velocity-adapted filter responses give considerably better possibility to discriminate the motion patterns compared to velocity-steered or non-adapted filters.

directional derivative responses and (iii) non-adapted partial derivative responses computed at velocity $v = 0$.

4.2. Discriminability of histograms

Fig. 10 shows the means and the variances of the histograms computed separately for both of the classes. As can be seen from Fig. 10a–c, velocity adaptation of receptive fields results in discriminative class histograms and low variation of histograms computed for the same class of activities. On the contrary, the high variations in the histograms in Fig. 10d–i clearly indicate that activities are much harder to recognize when using velocity-steered or non-adapted receptive fields.

Whereas Fig. 10 presents histograms for three types of derivatives L_{xxt} , L_{xyt} and L_{yyt} at scales $\sigma^2 = 4$, $\tau^2 = 4$ only, we have observed a similar behavior for other derivatives at most of the other scales considered.

4.3. Discriminability measure

To quantify these results, let us measure the distance between pairs of histograms (h_1, h_2) defined according to the χ^2 -divergence measure

$$D(h_1, h_2) = \sum_i \frac{(h_1(i) - h_2(i))^2}{h_1(i) + h_2(i)}, \quad (10)$$

where i is the index to the histogram bin. To evaluate the distance between a pair of sequences, we accumulate differences of histograms over different spatial and temporal scales as well as over different types of receptive fields according to $d(h_1, h_2) = \sum_{l, \sigma, \tau} D(h_1, h_2)$, where l denotes the type of the spatio-temporal filters, σ^2 the spatial scale and τ^2 the temporal scale.

To measure the degree of discrimination between different actions, we compare the distances between pairs of sequences that belong to the same class d_{same} with distances between sequences of different classes d_{diff} . Then, to quantify the average performance of the velocity adaptation algorithm, we compute the mean distances \bar{d}_{same} , \bar{d}_{diff} for all valid pairs of examples and define a *distance ratio* according to $r = \bar{d}_{\text{same}} / \bar{d}_{\text{diff}}$. Hence, low values of r indicate good discriminability, while r close to one corresponds to a performance no better than chance.

Fig. 11 shows distance ratios computed separately for different types of receptive fields. The lower values of the curve corresponding to velocity adaptation clearly indicate the better recognition performance obtained by using velocity-adapted filters compared to velocity-steered or non-adapted filters. Computing distance ratios over all types of derivatives and scales used, results in the following distance ratios: $r_{\text{adapt}} = 0.64$ when using velocity-adapted filters, $r_{\text{steered}} = 0.81$ using velocity-steered filters, and $r_{\text{non-adapt}} = 0.92$ using non-adapted filters.

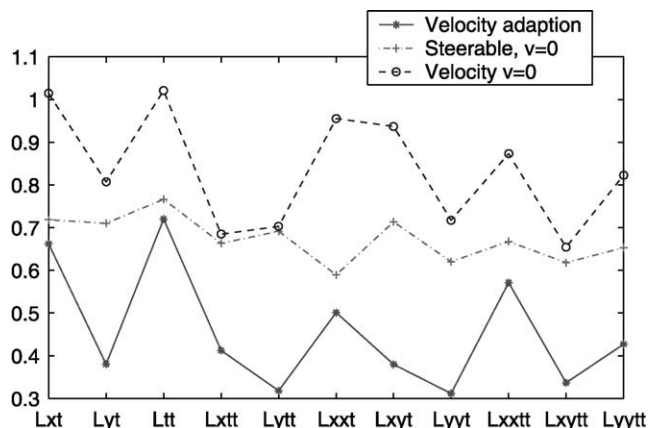


Fig. 11. Distance ratios computed for different types of derivatives and for velocity-adapted (solid lines), velocity-steered (point-dashed lines) and non-adapted (dashed lines) filter responses. As can be seen, local velocity adaptation results in lower values of the distance ratio and therefore better recognition performance compared to steered or non-adapted filter responses.

4.4. Dependency on scales

When analysing discrimination performance for different types of derivatives and different scales, we have observed an interesting dependency of the distance ratio on the spatial and the temporal scales. Fig. 12a and b shows how the distance ratio has a clear minimum over scales at $\sigma^2 = 2$, $\tau^2 = 8$ indicating that these scales give rise to the best discrimination for patterns considered here. In particular, it can be noted that $\tau^2 = 8$ approximately corresponds to the temporal extent of one gait cycle in our examples.

Computation of distance ratios for the selected scale values results in $r_{\text{adapt}} = 0.41$ when using velocity-adapted filters, $r_{\text{steered}} = 0.71$ using velocity-steered filters and $r_{\text{non-adapt}} = 0.79$ using non-adapted filters (see Fig. 13). The existence of such preferred scales motivates approaches for automatic selection of both spatial [27] and temporal [31] scales.

5. Summary and discussion

We have addressed the problem of representing and recognizing events in video in situations where the relative motion between the camera and the observed events is unknown. Experiments on a test problem of recognizing activities show that the use of a velocity adaptation scheme results in a clear improvement in the recognition performance compared to using either (steerable) directional derivatives or regular partial derivatives computed from a non-adapted spatio-temporal filtering step. Whereas for the treated set of examples, recognition could also have been accomplished by using a camera stabilization approach, a major aim here has been to consider a filtering scheme that can be extended to

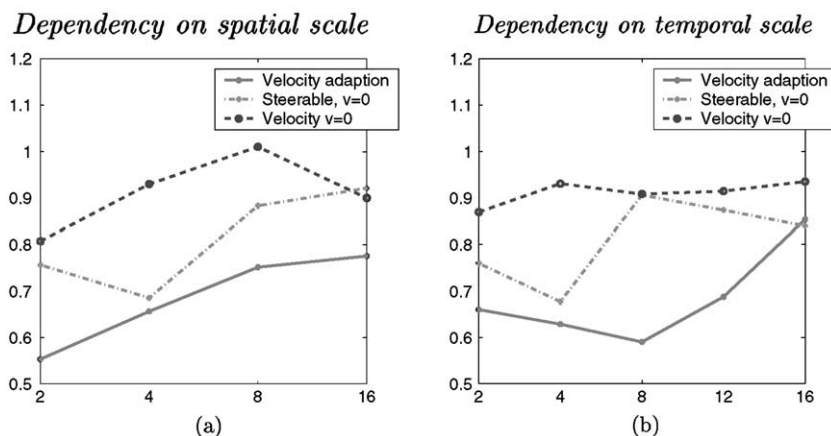


Fig. 12. Evolution of the distance ratio r over spatial scales (a) and temporal scales (b). Minima over scales indicate scale values with the highest discrimination ability.

recognition in complex scenes, where reliable camera stabilization may not be possible, i.e. scenes with complex non-static backgrounds or multiple events of interest. Full-fledged recognition in such situations, however, requires more sophisticated statistical methods for recognition than the present histogram-based scheme. We plan to investigate such extensions in future work.

Less restricted to this specific visual task, the results of our investigation also indicate how, when dealing with filter-based representations of spatio-temporal image data, velocity adaptation appears as an essential complement to more traditional approaches of using separable filtering in space-time. For the purpose of performing a clean experimental investigation, we have in this work made use of an explicit velocity-adapted spatio-temporal filtering for each image velocity. While such an implementation has interesting qualitative similarities to biological vision systems (where there are two main classes of receptive fields in space-time—separable filters and non-separable ones [33]), there is a need for developing more sophisticated multi-velocity filtering schemes for efficient implementations in practice.

Finally, future work should also address the problem of selecting appropriate scales in both the spatial and the temporal domains. The preliminary results in Section 4.4 indicate the potential of performing joint scale selection in space-time for increasing the recognition performance.

	velocity-adapted filtering	steerable filtering	non-adapted filtering
Average over all considered scales	0.64	0.81	0.92
At (manually) selected scales $\sigma^2 = 2$, $\tau^2 = 8$	0.41	0.71	0.79

Fig. 13. Values of distance ratios when averaged over all scales and at the manually selected scales that give best discrimination performance.

Acknowledgements

The support from the Swedish Research Council for Engineering Sciences (TFR), the Swedish Research Council (VR), as well as the Royal Swedish Academy of Sciences (KVA) and the Knut and Alice Wallenberg Foundation is gratefully acknowledged.

References

- [1] M. Swain, D. Ballard, Color indexing, *International Journal of Computer Vision* 7 (1) (1991) 11–32.
- [2] B. Schiele, J. Crowley, Recognition without correspondence using multidimensional receptive field histograms, *International Journal of Computer Vision* 36 (1) (2000) 31–50.
- [3] O. Chomat, V. de Verdiere, D. Hall, J. Crowley, Local scale selection for Gaussian based description techniques, *Proceedings of the Sixth European Conference on Computer Vision, Lecture Notes in Computer Science*, vol. 1842, Springer, Berlin, 2000, pp. 117–133.
- [4] O. Chomat, J. Martin, J. Crowley, A probabilistic sensor for the perception and recognition of activities, *Proceedings of the Sixth European Conference on Computer Vision, Dublin, Ireland (2000)* I:487–I:503.
- [5] D. Hall, V. de Verdiere, J. Crowley, Object recognition using coloured receptive fields, *Proceedings of the Sixth European Conference on Computer Vision, Lecture Notes in Computer Science*, vol. 1842, Springer, Berlin, 2000, pp. 164–177.
- [6] H. Schneiderman, T. Kanade, A statistical method for 3D object detection applied to faces and cars, *Proceedings of the Computer Vision and Pattern Recognition, Hilton Head, SC*, vol. I, 2000, pp. 746–751.
- [7] L. Zelnik-Manor, M. Irani, Event-based analysis of video, *Proceedings of the Computer Vision and Pattern Recognition, Kauai Marriott, Hawaii (2001)* II:123–II:130.
- [8] A.P. Witkin, Scale-space filtering, *Proceedings of the Eighth International Joint Conference on Artificial Intelligence, Karlsruhe, Germany (1983)* 1019–1022.
- [9] J.J. Koenderink, The structure of images, *Biological Cybernetics* 50 (1984) 363–370.
- [10] J.J. Koenderink, Scale-time, *Biological Cybernetics* 58 (1988) 159–162.
- [11] T. Lindeberg, *Scale-space Theory in Computer Vision*, The Kluwer

- International Series in Engineering and Computer Science, Kluwer, Dordrecht, 1994.
- [12] T. Lindeberg, Linear spatio-temporal scale-space, in: B.M. ter Haar Romeny, L.M.J. Florack, J.J. Koenderink, M.A. Viergever (Eds.), *Scale-space Theory in Computer Vision: Proceedings of the First International Conference on Scale-space'97*, Lecture Notes in Computer Science, vol. 1252, Springer, New York, 1997, pp. 113–127, Extended version available as Technical Report ISRN KTH/NA/P-01/22-SE from KTH (<http://www.nada.kth.se/cvap/abstracts/cvap257.html>).
- [13] T. Lindeberg, D. Fagerström, Scale-space with causal time direction, *Proceedings of the Fourth European Conference on Computer Vision*, vol. 1064, Springer, Berlin, 1996, pp. 229–240.
- [14] L.M.J. Florack, *Image Structure*, Series in Mathematical Imaging and Vision, Kluwer, Dordrecht, 1997.
- [15] M. Irani, P. Anandan, S. Hsu, Mosaic based representations of video sequences and their applications, *Proceedings of the Fifth International Conference on Computer Vision*, Cambridge, MA (1995) 605–611.
- [16] T. Lindeberg, J. Garding, Shape-adapted smoothing in estimation of 3D depth cues from affine distortions of local 2D structure, *Proceedings of the Third European Conference on Conference Vision*, Stockholm, Sweden (1994) A:389–A:400.
- [17] C. Ballester, M. Gonzalez, Affine invariant texture segmentation and shape from texture by variational methods, *Journal of Mathematical Imaging and Vision* 9 (1998) 141–171.
- [18] L. Florack, W. Niessen, M. Nielsen, The intrinsic structure of optic flow incorporating measurement duality, *International Journal of Computer Vision* 27 (3) (1998) 263–286.
- [19] J. Weickert, *Anisotropic Diffusion in Image Processing*, Teubner, Stuttgart, 1998.
- [20] A. Almansa, T. Lindeberg, Fingerprint enhancement by shape adaptation of scale-space operators with automatic scale-selection, *IEEE Transactions on Image Processing* 9 (12) (2000) 2027–2042.
- [21] F. Schaffalitzky, A. Zisserman, Viewpoint invariant texture matching and wide baseline stereo, *Proceedings of the Eighth International Conference on Computer Vision*, Vancouver, Canada (2001) II:636–II:643.
- [22] K. Mikolajczyk, C. Schmid, An affine invariant interest point detector, *Proceedings of the Seventh European Conference on Computer Vision*, Lecture Notes in Computer Science, vol. 2350, Springer, Berlin, 2002, pp. I:128–I:142.
- [23] T. Lindeberg, Time-recursive velocity-adapted spatio-temporal scale-space filters, *Proceedings of the Seventh European Conference on Computer Vision*, Lecture Notes in Computer Science, vol. 2350, Springer, Berlin, 2002, pp. I:52–I:67.
- [24] H. Nagel, A. Gehrke, Spatiotemporal adaptive filtering for estimation and segmentation of optical flow fields, *Proceedings of the Fifth European Conference on Computer Vision*, Freiburg, Germany (1998) II:86–II:102.
- [25] M. Black, Recursive non-linear estimation of discontinuous flow fields, *Proceedings of the Third European Conference on Computer Vision*, Lecture Notes in Computer Science, vol. 801, Springer, Berlin, 1994, pp. 138–145.
- [26] F. Guichard, A morphological, affine, and Galilean invariant scale-space for movies, *IEEE Transactions on Image Processing* 7 (3) (1998) 444–456.
- [27] T. Lindeberg, Feature detection with automatic scale selection, *International Journal of Computer Vision* 30 (2) (1998) 77–116.
- [28] M. Abramowitz, A. Stegun (Eds.), *Handbook of Mathematical Functions*, Applied Mathematics Series, 55th ed., National Bureau of Standards, 1964.
- [29] E. Adelson, J. Bergen, Spatiotemporal energy models for the perception of motion, *Journal of the Optical Society of America A* 2 (1985) 284–299.
- [30] D. Heeger, Optical flow using spatiotemporal filters, *International Journal of Computer Vision* 1 (1988) 279–302.
- [31] T. Lindeberg, On automatic selection of temporal scales in time-casual scale-space, in: G. Sommer, J.J. Koenderink (Eds.), *Proceedings of the AFPAC'97 Algebraic Frames for the Perception–Action Cycle*, Lecture Notes in Computer Science, vol. 1315, Springer, Berlin, 1997, pp. 94–113.
- [32] W.T. Freeman, E.H. Adelson, The design and use of steerable filters, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 13 (9) (1991) 891–906.
- [33] G.C. DeAngelis, I. Ohzawa, R.D. Freeman, Receptive field dynamics in the central visual pathways, *Trends in Neuroscience* 18 (10) (1995) 451–457.