

Feature Tracking with Automatic Selection of Spatial Scales

Lars Bretzner and Tony Lindeberg

Computational Vision and Active Perception Laboratory (CVAP),
Department of Numerical Analysis and Computing Science,
KTH, S-100 44 Stockholm, Sweden.
Email: bretzner@@bion.kth.se, tony@@bion.kth.se

Technical report ISRN KTH NA/P-96/21-SE.

Abstract

When observing a dynamic world, the size of image structures may vary over time. This article emphasizes the need for including explicit mechanisms for automatic scale selection in feature tracking algorithms in order to: (i) adapt the local scale of processing to the local image structure, and (ii) adapt to the size variations that may occur over time.

The problems of corner detection and blob detection are treated in detail, and a combined framework for feature tracking is presented in which the image features at every time moment are detected at locally determined and automatically selected scales. A useful property of the scale selection method is that the scale levels selected in the feature detection step reflect the spatial extent of the image structures. Thereby, the integrated tracking algorithm has the ability to adapt to spatial as well as temporal size variations, and can in this way overcome some of the inherent limitations of exposing fixed-scale tracking methods to image sequences in which the size variations are large.

In the composed tracking procedure, the scale information is used for two additional major purposes: (i) for defining local regions of interest for searching for matching candidates as well as setting the window size for correlation when evaluating matching candidates, and (ii) stability over time of the scale and significance descriptors produced by the scale selection procedure are used for formulating a multi-cue similarity measure for matching.

Experiments on real-world sequences are presented showing the performance of the algorithm when applied to (individual) tracking of corners and blobs. Specifically, comparisons with fixed-scale tracking methods are included as well as illustrations of the increase in performance obtained by using multiple cues in the feature matching step.

Keywords: feature, tracking, motion, blob, corner, scale, scale-space, scale selection, similarity, computer vision

Contents

1	Introduction	1
2	The need for automatic scale selection in feature tracking	3
3	Feature detection with automatic scale selection	4
3.1	Normalized derivatives	4
3.2	Corner detection with automatic scale selection	5
3.3	Blob detection with automatic scale selection	5
4	Tracking and prediction in a multi-scale context	6
5	Matching on multi-cue similarity	7
6	Combined tracking algorithm	8
7	Experimental results	10
7.1	Corner tracking	10
7.2	Blob tracking	12
8	Summary and Discussion	14
8.1	Spatial consistency and statistical evaluation.	15
8.2	Multi-cue tracking	15
8.3	Temporal consistency	16
A	Algorithmic details	20
A.1	Prediction	20
A.2	Feature detection	20
A.3	Matching	20

1 Introduction

Being able to track image structures over time is a useful and sometimes necessary capability for vision systems intended to interact with a dynamic world. There are several computer vision algorithms in which tracking arises as an important subproblem. Some situations are:

- Fixation means maintaining a relationship between a physical point or region in the world and some (usually central) region in a camera system. To maintain such a relationship over time, we have to relate some characteristic properties of the physical point to entities that are measurable from the available image data.
- Object recognition in a dynamically varying environment gives rise to the same type of problem, including the case when the visual agent is active and moves relative to the scene. Examples of the latter are navigation as well as active scene exploration. When objects move relative to the observer, feature tracking is a useful processing step for preserving the identity of image features over time.
- The identity problem is also essential in algorithms for motion segmentation and structure from motion. To compute structural properties or invariant descriptors which depend on the temporal variation of a geometric configuration, some mechanism is needed for matching corresponding image features over time.

There is an extensive literature on tracking methods operating without specific *a priori* knowledge about the world, such as object models or highly restricted domains. Without any aim of giving an extensive survey, the work in this direction can be classified into three main categories:

Correlation based tracking The presumably earliest approach to image matching is the correlation technique based on the similarity between corresponding grey-level patches over time. Given a window of some size, which covers an image detail at a certain time moment, the corresponding detail at the next time moment is defined as the position of the window (of the same size) that gives the highest correlation score when compared to the previous patch.

Optical flow based tracking The definition of an optic flow field gives rise to a motion field in the image domain, which can be interpreted as the result of tracking all image points simultaneously. With respect to the tracking problem, the motion of coherently moving (and possibly segmented) regions computed from optic flow algorithms can be used for guiding tracking procedures, as shown by [Thompson *et al.*, 1993] and [Meyer and Bouthemy, 1994].

Feature tracking Over the years a large number of approaches have been developed for tracking image features such as edges and corners over time. Essentially, what characterizes a feature tracking method is that image features are first extracted in a bottom-up processing step and then these features are used as the

main primitives for the tracking and matching procedures. Concerning corner tracking, [Shapiro *et al.*, 1992b] detect and track corners individually in an algorithm originally aimed at applications such as videoconferencing. [Smith and Brady, 1995] track a large set of corners and use the results in a flow-based segmentation algorithm. [Zheng and Chellappa, 1995] have studied feature tracking when compensating for camera motion, and [Gee and Cipolla, 1995] track locally darkest points with applications to pose estimation. In contour tracking, [Blake *et al.*, 1993, Curwen *et al.*, 1991] use snakes to track moving, deforming image features. [Cipolla and Blake, 1992] apply such an approach to estimate time-to-contact, and [Koller *et al.*, 1994] track combined motion and grey-level boundaries in traffic surveillance. An overview of different approaches to edge tracking can be found in the recent book by [Faugeras, 1993].

The subject of this article is to consider the domain of feature tracking and to complement previous works on this subject by addressing the problem of scale and scale selection in the spatial domain and by introducing new similarity measures in the matching step. In most previous works, the analysis is performed at a single predetermined scale. Here, we will emphasize and show by examples why it is useful to include an explicit mechanism for automatic scale selection to be able to handle situations in which the size variations are large. Besides avoiding explicit setting of scale levels for feature detection, and thus overcoming some of the fundamental limitations of processing image sequences at a single scale, it will be demonstrated how scale levels selected by a scale selection procedure can constitute a useful source of information when defining a similarity measure over time, as well as for adapting the window size for correlation to the local image structure.

Moreover, since the resulting matching algorithm we will arrive at is based on a similarity measure defined as the combination of different discriminative properties, and with small modifications can be applied to tracking of both corners and blobs, we will emphasize this multi-cue aspect as an important component for increasing the robustness of feature tracking algorithms.

The presentation is organized as follows: Section 2 illustrates the need for adaptive scale selection in feature tracking. It gives a hands-on demonstration of the improvement in performance that can be obtained by including a scale selection mechanism when tracking features in image sequences in which the size variations over time are large. Section 3 describes the feature detection step and reviews the basic components in a general principle for scale selection. Sections 4 and 5 explain how the scale information obtained from these processing modules can be used in the prediction step and in the evaluation of matching candidates. Section 6 summarizes how these components can be combined with a classical feature tracking scheme with prediction followed by detection and matching. Section 7 shows the performance of the algorithm when applied to real-world data. Feature tracking using adaptive scales is compared to tracking at one, fixed scale. Comparisons are also made between single-cue and multi-cue similarity measures. Finally, we conclude in section 8 by summarizing the main properties of the method and by outlining natural extensions.

2 The need for automatic scale selection in feature tracking

To extract features from an image, we have to apply some operators to the data. The type of features that can be extracted are largely determined by the spatial extent of these operators. When dealing with real-world data about which no or very little information is available, we can hardly expect to know in advance what scales are relevant for processing a given image. Therefore, a reasonable approach is to consider a large number of scales simultaneously, and this is one of the major motivations for using a multi-scale representation when automatically processing measurement data such as images.

Despite this now rather well-spread insight, most work on feature tracking still performs the analysis at one scale only. For correlation based tracking methods, this corresponds to using a fixed-size window over time, and concerning feature tracking to detecting image features at the same scale at all time moments. Such an approach will, however, suffer from inherent limitations when applied to real-life image sequences in which the size variations are large. This basic property constitutes one illustration of why a mechanism for automatic scale selection is an essential complement to traditional multi-scale processing in general, and to feature detection and feature tracking in particular.

In an image sequence, the size of image structures may change over time due to expansions or contractions. A typical example of the former is when the observer approaches an object as shown in figure 1. The left column in this figure shows a few snapshots from a tracker which follows a corner on the object over time using a standard feature tracking technique with a fixed scale for corner detection and a fixed window size for hypothesis evaluation by correlation. After a number of frames, the algorithm fails to detect the right feature and the corner is lost. The reason why this occurs, is simply the fact that the corner no longer exists at the predetermined scale. As a comparison, the right column shows the result of incorporating a mechanism for adaptation of the scale levels to the local image structure (details will be given in later sections). As can be seen, the corner is correctly tracked over the whole sequence. (The same initial scale was used in both experiments.)

Another motivation to this work originates from the fact that all feature detectors suffer from localization errors due to e.g noise and motion blur. When detecting rigid body motion or recovering 3D structure from feature point correspondences in an image sequence, it is important that the motion in the scene is large compared to the localization errors of the feature detector. If the inter-frame motion is small, we therefore have to track features over a large number of frames to obtain accurate results. This requirement constitutes a key motivation for including a scale selection mechanism in the feature tracker, to obtain longer trajectories of corresponding features as input to algorithms for motion estimation and recovery of 3D structure.

Concerning the common use of fixed scale levels in tracking methods, it is worth pointing out that in situations where the image features are distinct (e.g. sharp corners on a smooth background), traditional methods using fixed scales might be sufficient. The main advantages of having a mechanism for automatic

scale selection in such situations are that: (i) the actual tuning of the scale parameter can be avoided, (ii) as will be illustrated later, stability over time of the selected scale levels turns out to be a useful discriminative constraint to include in a matching criterion.

3 Feature detection with automatic scale selection

A natural framework to use when extracting features from image data is to define the image features from multi-scale differential invariants expressed in terms of Gaussian derivative operators [Koenderink and van Doorn, 1992, Florack *et al.*, 1992], or more specifically, as maxima or zero-crossings of such entities [Lindeberg, 1994c]. In this way, image features such as corners, blobs, edges and ridges can be computed at any level of scale.

A basic problem that arises for any such feature detector concerns how to determine at what scales the image features should be extracted, or if the feature detection is performed at several scales simultaneously, what image features should be regarded as significant. A framework addressing this problem has been developed in [Lindeberg, 1993, Lindeberg, 1994c]. In summary, one of the main results from this work is a general principle for scale selection, which states that scale levels for feature detection can be selected from the scales at which normalized differential invariants assume maxima over scales. In this section, we shall give a brief review of how this methodology applies to the detection of features such as blobs and corners. The image features so obtained, with their associated attributes resulting from the scale selection method, will then be used as basic primitives for the tracking procedure.

3.1 Normalized derivatives

The scale-space representation [Witkin, 1983, Koenderink, 1984] of a signal f is defined as the result of convolving f

$$L(\cdot; t) = g(\cdot; t) * f \quad (1)$$

with Gaussian kernels having different values of the scale parameter t

$$g(x; t) = \frac{1}{2\pi t} e^{-(x^2+y^2)/(2t)} \quad (2)$$

In this representation, γ -normalized derivatives [Lindeberg, 1996a] are defined by

$$\partial_{\xi} = t^{\gamma/2} \partial_x \quad (3)$$

where t is the variance of the Gaussian kernel. From this construction, a normalized differential invariant is then obtained by replacing all spatial derivatives by corresponding normalized derivatives according to (3).

3.2 Corner detection with automatic scale selection

A common way to define a corner in a grey-level image in differential geometric terms is as a point at which both the curvature of a level curve

$$\kappa = \frac{-(L_{yy}L_x^2 + L_{xx}L_y^2 - 2L_xL_yL_{xy})}{(L_x^2 + L_y^2)^{3/2}} \quad (4)$$

and the gradient magnitude

$$|\nabla L| = \sqrt{L_x^2 + L_y^2} \quad (5)$$

are high [Kitchen and Rosenfeld, 1982, Koenderink and Richards, 1988, Deriche and Giraudon, 1990, Blom, 1992]. If we consider the product of κ and the gradient magnitude raised to some power, and choose the power equal to three, we obtain the essentially affine invariant expression

$$\tilde{\kappa} = L_{yy}L_x^2 + L_{xx}L_y^2 - 2L_xL_yL_{xy} \quad (6)$$

with its corresponding γ -normalized differential invariant

$$\tilde{\kappa}_{\gamma\text{-norm}} = t^{2\gamma}\tilde{\kappa} \quad (7)$$

In [Lindeberg, 1994a] it is shown how a junction detector with automatic scale selection can be formulated in terms of the detection of *scale-space maxima* of $\tilde{\kappa}_{\gamma\text{-norm}}^2$, *i.e.*, by detecting points in scale-space where $\tilde{\kappa}_{\gamma\text{-norm}}^2$ assumes maxima with respect to both scale and space. When detecting image features at coarse scales it turns out that the localization can be poor. Therefore, this detection step is complemented by a second localization stage, in which a modified Förstner operator [Förstner and Gülch, 1987], is used for iteratively computing new localization estimates using scale information from the initial detection step (see the references for details).

A useful property of this corner detection method is that it leads to selection of coarser scales for corners having large spatial extent. Figure 2 illustrates this property by showing the result of applying the corner detection method to two different images, and graphically illustrating each detected and localized corner by a circle with the radius proportional to the detection scale. Notably, the support regions of these blobs serve as natural regions of interest around the detected corners. As we shall demonstrate later, such regions of interest and context information turn out to be highly useful for a feature tracking procedure.

3.3 Blob detection with automatic scale selection

As shown in the abovementioned references, a straightforward method for blob detection can be formulated in an analogous manner by detecting scale-space maxima of the square of the normalized Laplacian

$$\nabla_{\text{norm}}^2 L = t(L_{xx} + L_{yy}) \quad (8)$$

This operator gives a strong response for blobs that are brighter or darker than their background, and in analogy with the corner detection method, the selected scale levels provide information about the characteristic size of the blob.

Figure 3 shows the result of applying this blob detection method to the same images as used in figure 2. As can be seen, a representative set of blob features at different scales is extracted. Moreover, it can be noted how well the blob circles reflect the size variations, in particular, considering how simple operations the blob detection algorithm is based on (Gaussian smoothing, derivative computation, and detection of scale-space maxima).

4 Tracking and prediction in a multi-scale context

When tracking features over time, both the position of the feature and the appearance of its surrounding grey-level pattern can be expected to change. To relate features over time, we shall throughout this work make use of the common assumption about small motions between successive frames.

There are several ways to predict the position of a feature in the next frame based on its positions in previous frames. Whereas the Kalman filtering methodology has been commonly used in the computer vision literature, this approach suffers from a fundamental limitation if the motion direction suddenly changes. If a feature moving in a certain direction has been tracked over a long period of time, then the built-in temporal smoothing of the feature trajectory in the Kalman filter, implies that the predictions will continue to be in essentially the same direction, although the actual direction of the motion changes. If the covariance matrices in the Kalman filter have been adapted to small oscillations around the previously smooth trajectory, it will hence be likely that the feature is lost at the discontinuity.¹

For this reason, we shall make use of simpler first-order prediction, which uses the motion between the previous two successive frames as a prediction to the next frame.²

Within a neighbourhood of each predicted feature position, we detect new features using the corner (or blob) detection procedure with automatic scale selection. The support regions associated with the features serve as natural regions of interest when searching for new corresponding features in the next frame. In this way, we can avoid the problem of setting a global threshold on the distance between matching candidates. There is, of course, a certain scaling factor between the detection scale and the size of the support region. The important property of this method, however, is that it will automatically select smaller regions of interest for small-size image structures, and larger search regions for larger size structures. Here, we shall make use of this scale information for three main purposes:

¹As will be shown in the experiments in section 7, the resulting feature trajectories may be quite irregular. Enforced temporal smoothing of the image positions of the features, leading to smoother trajectories, would not be appropriate for such data.

²Both constant acceleration and constant velocity models have been used, but the latter has given better performance in most cases.

- Setting the search region for possible matching candidates.
- Setting the window size for correlation matching.
- Using the stability of the detection scale as a matching condition.

We set the size of the search region to the spatial extent of the previous image feature, multiplied by a safety factor. Within this window, a certain number of candidate matches are selected. Then, an evaluation of these matching candidates is made based on a combined similarity measure to be defined in the next section.

5 Matching on multi-cue similarity

Based on the assumption of small inter-frame image motions, we use a multiple cue approach to the feature matching problem. Instead of evaluating the matching candidates using a correlation measure on a local grey-level patch only, as done in most feature tracking algorithms, we combine the correlation measure with significance stability, scale stability and proximity measures as defined below.

Patch similarity. This measure is a normalized Gaussian-weighted intensity cross-correlation between two image patches. Here, we compute this measure over a square centered at the feature and with its size set from the detection scale. The measure is derived from the cross-correlation of the image patches, see [Shapiro *et al.*, 1992a], computed using a Gaussian weight function centered at the feature. The motivation for using a Gaussian weight function is that image structures near the feature center should be regarded as more significant than peripheral structures. Given two brightness functions I_A and I_B , and two image regions $D_A \subset \mathbb{R}$ and $D_B \subset \mathbb{R}$ of the same size $|D| = |D_A| = |D_B|$ centered at p_A and p_B respectively, the weighted cross-correlation between the patches is defined as:

$$C(A, B) = \frac{1}{|D|} \sum_{x \in D_A} e^{-(x-p_A)^2} I_A(x) I_B(x - p_A + p_B) - \frac{1}{|D|^2} \sum_{x_A \in D_A} e^{-(x-p_A)^2} I_A(x_A) \sum_{x_B \in D_B} e^{-(x-p_B)^2} I_B(x_B) \quad (9)$$

and the normalized weighted cross-correlation is

$$S_{patch}(A, B) = \frac{C(A, B)}{\sqrt{C(A, A) C(B, B)}} \quad (10)$$

where

$$C(A, A) = \frac{1}{|D|} \sum_{x \in D_A} (e^{-(x-p_A)^2} I_A(x))^2 - \frac{1}{|D|^2} \left(\sum_{x \in D_A} e^{-(x-p_A)^2} I_A(x) \right)^2 \quad (11)$$

and $C(B, B)$ is defined analogously. As is well-known, this similarity measure is invariant to superimposed linear illumination gradients. Hence, first-order effects of scene lightning do not affect this measure, and the measure only accounts for changes in the structure of the patches.

Significance stability. A straightforward significance measure of a feature detected according to the method described in section 3 is the normalized response at the local scale-space maximum. For corners, this measure is the normalized level curve curvature according to (7) and for blobs it is the normalized Laplacian according to (8). To compare significance values over time, we measure similarity by relative differences instead of absolute, and define this measure as

$$S_{sign} = \left| \log \frac{R_B}{R_A} \right| \quad (12)$$

where R_A and R_B are the significance measures of the corresponding features A and B .

Scale stability. Since the features are detected at different scales, the ratio between the detection scales of two features constitutes a measure of stability over scales. To measure relative scale variations, we use the absolute value of the logarithm of this ratio, defined as

$$S_{scale} = \left| \log \frac{t_B}{t_A} \right| \quad (13)$$

where t_A and t_B are the detection scales of A and B .

Proximity We measure how well the position x_A of feature A corresponds to the position x_{pred} predicted from feature B

$$S_{pos} = \frac{\|x_A - x_{pred}\|}{\sqrt{t_B}} \quad (14)$$

where t_B is the detection scale feature B .

Combined similarity measure. In summary, the similarity measure we make use of a weighted sum of (10), (12) and (13),

$$S_{comb} = c_{patch} S_{patch} + c_{sign} S_{sign} + c_{scale} S_{scale} + c_{pos} S_{pos} \quad (15)$$

where c_{patch} , c_{sing} , c_{scale} and c_{pos} are tuning parameters to be determined.

6 Combined tracking algorithm

By combining the components described in the previous sections, we obtain a feature tracking scheme based on a traditional predict-detect-update loop. In addition, the following processing steps are added:

- *Quality measure.* Each feature is assigned a quality measure indicating how stable it is over time.
- *Bidirectional matching.* To provide additional information to later processing stages about the reliability of the matches, the matching can be done bidirectionally. Given a feature F_1 from the feature set, we first compute its winning matching candidate F_2 in the current image. If then F_1

is the winning candidate of F_2 in the backward matching direction, the match between F_1 and F_2 is registered as safe. This processing step is useful for signalling possible matching errors.

During the tracking procedure each feature is associated with the following attributes:

- its detection scale t_{det} ,
- its estimated size $D = k_{size} * \sqrt{t_{det}}$ bounded from below to D_{min} ,
- its position,
- its quality value.

An overview of the tracking algorithm is given in figure 4. At a more detailed level, each individual module operates as follows:

Prediction The prediction is performed as described in section 4. For each feature in the feature set, a linear prediction of the position in the current frame is computed based on the positions of the corresponding feature in the two previous frames. The size of the search window is computed as $k_{w1} * D$ (with the size D bounded from below). When a trajectory is initiated, there is no feature history to base the prediction on, so we use a larger search window of size $k_{w2} * D$ ($k_{w2} > k_{w1}$) and use the original feature position as the predicted position.

Detection In each frame, image features are detected as described in section 3. The window obtained from the prediction step is searched for the same kind of features over a locally adapted range of scales $[t_{min}, t_{max}]$, where $t_{max} = k_{range} * t_{det}$ and $t_{min} = t_{det}/k_{range}$. The number n of detected candidates depends on which feature extraction method we use in the detection step.

Matching The matching is based on the similarity measures described in section 5. The original feature is matched to the candidates obtained from the detection step and the winner is the feature having the highest combined similarity value above a fixed threshold T_{comb} and a patch correlation value above a threshold T_{patch} . These thresholds are necessary to suppress false matches when features disappear due to e.g occlusion.

If a feature is matched, the quality value is increased by dq_i and its position, its scale descriptor, its significance value and its grey-level patch are updated.

If no match is found, the feature is considered unmatched, its quality value is decreased by dq_d and its position is set to the predicted position.

Finally for each frame, the feature set is parsed to detect feature merges and to remove features having quality values below a threshold T_q . When two features merge, their trajectories are terminated and a new trajectory is initiated. In this way, we obtain more reliable feature trajectories for further processing.

7 Experimental results

7.1 Corner tracking

Let us first demonstrate the performance of the algorithm when applied to an image sequence consisting of 60 frames. In this sequence, the camera moves in a fairly complex way relative to a static scene. The objects of interest on which the features (here corners) are detected are a telephone and a package on a table. From the junctions detected in the initial frame, a subset of 14 features were selected manually as shown in figure 5.

Figure 6 shows the situation after 30, 50 and 60 frames. In the illustrations, black segments on the trajectories indicate matched positions, while white segments show unmatched (predicted) positions. The matching is based on the combined similarity measure incorporating patch correlation, scale stability, significance stability and proximity. The detection scales of the features are illustrated by the size of the circles in the images, and we see how all corners are detected at fine scales in the initial frame. As time evolves, the detection scales adapt to the size changes of the image structures; tracked sharp corners are still detected at fine scales while blunt corners are detected at coarser scales when the camera approaches the scene.

Figure 7 shows the result of an attempt to track the same corners at fixed scales, using the automatically determined detection scales from the initial image. As can be seen, the sharpest corners are correctly tracked but the blunt corners are inevitably lost. This effect is similar to the initial illustration in section 2.

Figure 8 shows another example for a camera tracking a toy train on a table. In the initial frame, 29 corners were selected manually; 25 on the train and 4 on an object in the background. Some of these corners are enumerated and will be referred to when discussing the performance below.

<i>Corner no</i>	<i>Patch similarity only</i>	<i>Combined similarity measure</i>
1	lost in frame 29	lost in frame 29
2	mismatched in 18	mismatched in 18
3	mismatched in 16	mismatched in 16
4	lost in 83	—
5	mismatched in 63	—
6	lost in 81	lost in 75
7	lost in 33	—
8	lost in 46	lost in 46

Table 1: Table showing when eight of the enumerated corners in the train sequence are lost. Note that out of the corners which are lost when matching on patch similarity only, three corners are tracked during the whole sequence when using the combined similarity measure.

Figure 9 shows the situation after 60, 100 and 140 frames, using the combined similarity measure in the matching step. The white parts of the tracks

show when the algorithm failed to match the corners (stressing the importance of keeping unmatched features over a certain number of frames). Noisy image data and motion blur will increase the number of matching failures. Corners no 2, 3, 6 and 8 are lost due to moving structures in the background causing accidental views. In the last frames of the sequence, corner no 9 has poor localization, since the corner edges are aligned causing the corner to disappear. The importance of using the combined similarity measure in the matching step is illustrated in the train sequence in figure 10, showing the result of matching on patch correlation only. We see that corners no 4, 5, and 7, which were all tracked using the combined similarity measure, now are lost. Table 1 shows, for both experiments, when the enumerated corners in the train sequence are lost.

7.2 Blob tracking

Let us now apply the same framework for blob tracking. In the train sequence, we manually selected 11 blobs on the train and 2 blobs in the background in the initial frame shown in figure 11. Figure 12 shows the situation after 30, 90 and 150 frames. The size of the circles in the figures correspond to the detection scales of the blobs. Note how the detection scale adapts to the local image structure when the blobs undergo expansion followed by contraction. All visible blobs except one are tracked during the whole sequence.

Referring to the need for automatic scale selection in feature tracking, as advocated in section 2, it is illustrative to show the results of attempting blob tracking with feature detection at a fixed scale. The scale level for detecting each blob was automatically selected in the first frame and was then kept fixed throughout the sequence. Figure 13 shows the result after 30 and 150 frames. Clearly, the tracker has severe problems due to the expansion and contraction in the sequence.

As a further illustration of the capability of the algorithm to track blobs under large size changes we applied it to a sequence of 87 images where a person, dressed in a spotted shirt, approaches the camera. In a rectangular area in the initial frame, the 20 most significant blobs were automatically detected, as shown in figure 14. Figure 15 shows the results after 25, 50 and 87 frames when matching on the combined similarity measure. All blobs except one are correctly tracked over the entire sequence.

Figure 16 shows the situation after 25 frames when matching on patch similarity only. Compared to figure 15, three more blobs are now lost, and one blob is mismatched. In scenes like this one, with repetitive, similar structures, the rate of mismatches is considerably higher if we match on patch correlation only instead of using the combined similarity measure.

When trying to track the blobs at a fixed scale, as can be seen in figure 17, most of the blobs are lost already after 25 frames. The last correctly tracked blob is lost after about 50 frames.

In summary, these experiments show that similar qualitative properties hold for blob tracking and for junction tracking: (i) By including the significance values and the selected scale levels in the matching criterion, we obtain a better performance than when matching on grey-level correlation only. (ii) The performance of tracking at adaptively determined scale levels is superior compared to similar tracking at a fixed scale.

Let us finally illustrate how feature tracking with automatic scale selection over a large number of frames is likely to give us trajectories which correspond to reliable and stable physical scene points or regions of interest on objects. By explicitly registering the features that are stable over time, we are able to suppress spurious feature responses due to noise, temporary occlusions etc. Figure 18 shows the initial frame of a sequence in which the 10 most significant blobs have been tracked in a region around the face of the subject. The subject first approaches the camera and then moves back to the initial position. Figure 19 shows the situation after 20, 45 and 90 frames. We can see that after a while only four features remain in the feature set and these are the stable features corresponding to the nostrils and the eyes. This ability to register stable image structures over time is clearly a desirable quality in many computer vision applications. Notably, for general scenes with large expansions or contractions, a scale selection mechanism is essential to allow for such registrations.

8 Summary and Discussion

We have presented a framework for feature tracking in which a mechanism for automatic scale selection has been built into the feature detection stage and the additional attributes of the image features obtained from the scale selection module are used for guiding the other processing steps in the tracking procedure.

We have argued that such a mechanism is essential for any feature tracking procedure intended to operate in a complex environment, in order to adapt the scale of processing to the size variations that may occur in the image data as well as over time. If we attempt to track features by processing the image data at one single scale only, we can hardly expect to be able to follow the features over large size variations. This property is a basic consequence of the inherent multi-scale nature of image structures, which means that a given object may appear in different ways depending on the scale of observation.

Specifically, based on a previously developed feature detection framework with automatic scale selection, we have presented a scheme for tracking corners and blobs over time in which:

- the image features at any time moment are detected using a feature detection method with automatic scale selection, and
- this information is used for
 - guiding the detection and selection of new feature candidates,
 - providing context information for the matching procedure,
 - formulating a similarity measure for matching features over time.

Besides avoiding explicit selection of scale levels for feature detection, the feature detection procedure with automatic scale selection allows us to track image features over large size variations. As demonstrated in the introductory example in section 2, we can in this way obtain a substantial improvement in the performance relative to a fixed-scale feature tracker.

Since the scale levels obtained from the scale selection procedure reflect the spatial extent of the image structures, we can also use this context information for avoiding explicit settings of distance thresholds and predefined window sizes for matching. Moreover, by including the scale and significance information associated with the image features from the scale selection procedure into a multi-cue similarity measure, we showed how we in this way can improve the reliability of the low-level matching procedure.

Of course, there are inherent limitations in tracking each feature individually as done in this work, and as can be seen from the examples, there are a number of situations where the tracking algorithm fails. Typically, this occurs because of rapid changes in the local grey-level pattern around the corner, corresponding to violations of the assumption about small inter-frame motions.

A notable conclusion that can be made in this context, is that despite these limitations, we have shown by examples that the resulting tracking procedure is able to track most of the visible features that can be followed over time in the sequences presented in this article. By this we argue that the type of

framework presented here provides an important step towards overcoming some of the limitations in previous feature tracking algorithms.

8.1 *Spatial consistency and statistical evaluation.*

In the scheme presented so far, each feature is tracked *individually*, without any explicit notion of coherently moving clusters. It is obvious that the performance of a tracking method can be improved if the latter notion can be introduced, and the overall motion of the clusters can be used for generating better predictions, as well as more refined evaluation criteria of matching candidates. To investigate if the motions of the tracked features possibly correspond to the same rigid body motion, we might compute descriptors such as affine 3-D coordinates. Interesting work in this direction have been presented by [Reid and Murray, 1993, Wiles and Brady, 1995, Shapiro, 1995].

It is also natural to include a statistical evaluation of the reliability of matches as well as their possible agreement with different clusters, as done in [Shapiro, 1995]. Whereas such an approach has not been explored in this work, this should not be interpreted as implying that the scale selection method excludes the usefulness of a statistical evaluation. The main intention behind this work has been to explore how far it is possible to reach by using a bottom-up construction of feature trajectories and by including a mechanism for automatic scale selection in the feature detection step. Then, the intention is that these two approaches should be applied in a complementary manner, where the scale selection method serves as a pre-conditioner for generating more reliable hypotheses with more reliable input data. The scale selection method can also provide context information over what domains statistical evaluations should be made.

8.2 *Multi-cue tracking*

A tracking method based on a single visual cue, like those reviewed in section 1 may have a rather good performance under certain conditions but may fail in more complex scenes. In this context, a multi-cue approach to the tracking problem is natural, i.e a system in which several types of algorithms operate simultaneously and the algorithm most suitable to a given situation dominates. This means that the vision system must have the ability to evaluate the reliability of the various tracking methods and to switch between them in an appropriate way.

Initial work in this direction, combining disparity cues with optical flow based object segmentation, has been performed by [Uhlin *et al.*, 1995]. The approach developed here lends itself naturally to integration with such techniques, in which such cues can be used for evaluating candidate feature clusters, and the feature tracking module in turn can be used as a more refined processing mechanism for maintaining object hypotheses over time. Of course, this leads to basic problems of feature selection. One possible approach for addressing such problems has been presented by [Shi and Tomasi, 1994].

8.3 Temporal consistency

As a final remark it is worth pointing out that in this work, the image features in each frame have been extracted *independently* from each other and without any other explicit use of temporal consistency than the heuristic condition that a feature hypothesis is allowed to survive over a few frames. To make more explicit use of temporal consistency, it is natural to incorporate the notion of a temporal scale-space representation [Lindeberg and Fagerström, 1996] and to include scale selection over the temporal scale domain as well [Lindeberg, 1996b].

In this context, it is also natural to combine the feature tracking approach with a simultaneous calculation of optical flow estimates and to integrate these two approaches so as to make use of their relative advantages. These subjects, including the integration of multiple tracking techniques into a multi-cue framework, constitute major goals of our continued research.

References

- [Blake *et al.*, 1993] Blake et al. “Affine-invariant contour tracking with automatic control of spatiotemporal scale”. In *Proc. 4th International Conference on Computer Vision*, Berlin, Germany, 1993. IEEE Computer Society Press.
- [Blom, 1992] J. Blom. *Topological and Geometrical Aspects of Image Structure*. PhD thesis, Dept. Med. Phys. Physics, Univ. Utrecht, NL-3508 Utrecht, Netherlands, 1992.
- [Cipolla and Blake, 1992] R. Cipolla and A. Blake. “Surface orientation and time to contact from image divergence and deformation”. In G. Sandini, editor, *Proc. 2nd European Conference on Computer Vision*, pages 187–202, Santa Margherita Ligure, Italy, 1992. Springer Verlag, Berlin.
- [Curwen *et al.*, 1991] Curwen et al. “Parallel implementation of Lagrangian dynamics for real-time snakes”. In *Proc. British Machine Vision Conference*. Springer Verlag, Berlin, 1991.
- [Deriche and Giraudon, 1990] R. Deriche and G. Giraudon. “Accurate Corner Detection: An Analytical Study”. In *Proc. 3rd Int. Conf. on Computer Vision*, pages 66–70, Osaka, Japan, 1990.
- [Faugeras, 1993] O. Faugeras. *Three-dimensional computer vision*. MIT Press, Cambridge, Massachusetts, 1993.
- [Florack *et al.*, 1992] L. M. J. Florack; B. M. ter Haar Romeny; J. J. Koenderink, and M. A. Viergever. “Scale and the Differential Structure of Images”. *Image and Vision Computing*, 10(6):376–388, Jul. 1992.
- [Förstner and Gülch, 1987] W. A. Förstner and E. Gülch. “A Fast Operator for Detection and Precise Location of Distinct Points, Corners and Centers of Circular Features”. In *Proc. Intercommission Workshop of the Int. Soc. for Photogrammetry and Remote Sensing*, Interlaken, Switzerland, 1987.
- [Gee and Cipolla, 1995] A.H. Gee and R. Cipolla. “Fast visual tracking by temporal consensus”. Technical Report CUED/F-INFENG/TR207, Dept of Engineering, University of Cambridge, England, 1995.
- [Kitchen and Rosenfeld, 1982] L. Kitchen and A. Rosenfeld. “Gray-Level Corner Detection”. *Pattern Recognition Letters*, 1(2):95–102, 1982.
- [Koenderink and Richards, 1988] J. J. Koenderink and W. Richards. “Two-Dimensional Curvature Operators”. *J. of the Optical Society of America*, 5:7:1136–1141, 1988.
- [Koenderink and van Doorn, 1992] J. J. Koenderink and A. J. van Doorn. “Generic neighborhood operators”. *IEEE Trans. Pattern Analysis and Machine Intell.*, 14(6):597–605, Jun. 1992.
- [Koenderink, 1984] J. J. Koenderink. “The structure of images”. *Biological Cybernetics*, 50:363–370, 1984.
- [Koller *et al.*, 1994] D. Koller; J. Weber, and J. Malik. “Robust multiple car tracking with occlusion reasoning”. In J.-O. Eklundh, editor, *Proc. 3rd European Conference on Computer Vision*, pages 189–196, Stockholm, Sweden, 1994. Springer Verlag, Berlin.
- [Lindeberg and Fagerström, 1996] T. Lindeberg and D. Fagerström. “Scale-Space with causal time direction”. In *Proc. 4th European Conference on Computer Vision*, volume 1064, pages 229–240, Cambridge, UK, April 1996. Springer Verlag, Berlin.

- [Lindeberg, 1993] T. Lindeberg. “On Scale Selection for Differential Operators”. In K. Heia K. A. Høgdra, B. Braathen, editor, *Proc. 8th Scandinavian Conf. on Image Analysis*, pages 857–866, Tromsø, Norway, May. 1993. Norwegian Society for Image Processing and Pattern Recognition.
- [Lindeberg, 1994a] T. Lindeberg. “Junction detection with automatic selection of detection scales and localization scales”. In *Proc. 1st International Conference on Image Processing*, volume I, pages 924–928, Austin, Texas, Nov. 1994. IEEE Computer Society Press.
- [Lindeberg, 1994b] T. Lindeberg. “Scale Selection for Differential Operators”. Technical Report ISRN KTH/NA/P--94/03--SE, Dept. of Numerical Analysis and Computing Science, KTH, Stockholm, Sweden, Jan. 1994. (Submitted).
- [Lindeberg, 1994c] T. Lindeberg. *Scale-Space Theory in Computer Vision*. The Kluwer International Series in Engineering and Computer Science. Kluwer Academic Publishers, Dordrecht, Netherlands, 1994.
- [Lindeberg, 1996a] T. Lindeberg. “Edge detection and ridge detection with automatic scale selection”. In *Proc. IEEE Comp. Soc. Conf. on Computer Vision and Pattern Recognition, 1996*, pages 465–470, San Francisco, California, June 1996. IEEE Computer Society Press.
- [Lindeberg, 1996b] T. Lindeberg. “On automatic selection of temporal scales”. 1996.
- [Meyer and Bouthemy, 1994] F. G. Meyer and P. Bouthemy. “Region-based tracking using affine motion models in long image sequences”. *Computer Vision, Graphics, and Image Processing :Image Understanding*, 60(2):119–140, 1994.
- [Reid and Murray, 1993] I. D. Reid and D. W. Murray. “Tracking foveated corner clusters using affine structure”. In *Proc. 4th International Conference on Computer Vision*, pages 76–83, Berlin, Germany, 1993. IEEE Computer Society Press.
- [Shapiro, 1995] L. S. Shapiro. *Affine analysis of image sequences*. Cambridge University Press, Cambridge, England, 1995.
- [Shapiro *et al.*, 1992a] L. S. Shapiro; H. Wang, and J. M. Brady. “A corner matching and tracking strategy applied to videophony”. Technical Report OUEL 1933/92, Robotics Research Group, University of Oxford, 1992.
- [Shapiro *et al.*, 1992b] L. S. Shapiro; H. Wang, and J. M. Brady. “A matching and tracking strategy for independently moving objects”. In *Proc. British Machine Vision Conference*, pages 306–315. Springer Verlag, Berlin, 1992.
- [Shi and Tomasi, 1994] J. Shi and C. Tomasi. “Good features to track”. In *Proc. IEEE Comp. Soc. Conf. on Computer Vision and Pattern Recognition*, pages 593–600. IEEE Computer Society Press, 1994.
- [Smith and Brady, 1995] S. M. Smith and J. M. Brady. “ASSET-2: Real-time motion segmentation and shape tracking”. *IEEE Trans. Pattern Analysis and Machine Intell.*, 17(8):814–820, 1995.
- [Thompson *et al.*, 1993] W. B. Thompson; P. Lechleider, and E. R. Stuck. “Detecting moving objects using the rigidity constraint”. *IEEE Trans. Pattern Analysis and Machine Intell.*, 15(2):162–165, 1993.
- [Uhlin *et al.*, 1995] T. Uhlin; P. Nordlund; A. Maki, and J.-O. Eklundh. “Towards an Active Visual Observer”. In *Proc. 5th International Conference on Computer Vision*, pages 679–686, Cambridge, MA, June 1995.

- [Wiles and Brady, 1995] C. S. Wiles and M. Brady. “Closing the loop on multiple motions”. In *Proc. 5th International Conference on Computer Vision*, pages 308–313. IEEE Computer Society Press, 1995.
- [Witkin, 1983] A. P. Witkin. “Scale-space filtering”. In *Proc. 8th Int. Joint Conf. Art. Intell.*, pages 1019–1022, Karlsruhe, West Germany, Aug. 1983.
- [Zheng and Chellappa, 1995] Q. Zheng and R. Chellappa. “Automatic feature point extraction and tracking in image sequences for arbitrary camera motion”. *International Journal of Computer Vision*, 15(1):31–76, 1995.

A Algorithmic details

This appendix gives a detailed listing of the parameters that influence the algorithm as well as the parameter settings that have been used for generating the experiments.

A.1 Prediction

The parameters determining the size of the search window (see section 6) were

$$\begin{aligned} k_{size} &= 5 \\ k_{w1} &= 1.5 \\ k_{w2} &= 2 * k_{w1} \\ D_{min} &= 16 \end{aligned}$$

A.2 Feature detection

When detecting features with automatic scale selection, the following scale ranges were used in the initial frame:

<i>Junction detection</i>	<i>Blob detection</i>
$t_{min} = 4.0$	$t_{min} = 4.0$
$t_{max} = 256.0$	$t_{max} = 512.0$

and the parameter γ in the normalized derivative concept (see section 3) was set to:

<i>Junction detection</i>	<i>Blob detection</i>
$\gamma = 0.875$	$\gamma = 1$

When searching for new image features, the search for matching candidates to a feature detected at scale t_{det} was performed in the interval $[t_{det}/k_1, t_{det}k_1]$, where $k_{range} = 3$.

In all experiments, the sampling density in the scale direction was set to correspond to a minimum of 5 scale levels per octave. In all other aspects, the feature detection algorithms followed the default implementation of junction and blob detection with automatic scale selection described in [Lindeberg, 1994b]. The maximum number of matching candidates evaluated for each feature was:

<i>Junction detection</i>	<i>Blob detection</i>
$n = 8$	$n = 20$

A.3 Matching

The following thresholds were used in the matching step

<i>Junction detection</i>	<i>Blob detection</i>
$T_{patch} = 0.75$	$T_{patch} = 0.6$
$T_{comb} = 0.65$	$T_{comb} = 0.5$

and the parameters for controlling the quality measure over time (see section 6)

$$\begin{aligned} dq_i &= 0.2 \\ dq_d &= 0.1 \\ T_q &= 0 \end{aligned}$$

Similarity measures: Relative weights In the experiments presented here, the following relative weights (see section 5) were used in the combined significance measure (15):

<i>Junction detection</i>	<i>Blob detection</i>
$c_{patch} = 1.0$	$c_{patch} = 1.0$
$c_{sign} = -0.08$	$c_{sign} = -0.25$
$c_{scale} = -0.08$	$c_{scale} = -0.08$
$c_{pos} = -0.1$	$c_{pos} = -0.1$

To give a qualitative motivation for using these orders of magnitude for the relative weights, let us first estimate the ranges in which these descriptors will vary:

- For the cross-correlation measure, it trivially holds that $|S_{patch}| < 1$. By the thresholding operation on this value, $|T_{patch}| = 0.7$, the variation of this entity is confined to the interval $|S_{patch}| \in [0.7, 1.0]$. In practice, the relative variations are usually in the interval $|S_{patch}| \in [0.8, 1.0]$.
- Concerning the significance measure, the significance values of corners computed from an image with grey-level values in the range $[0, 255]$ typically vary in the interval $\log R < 25$. Empirically, the relative variations are usually of the order of $\Delta \log R < 3$. For blob features, the corresponding values are $\log R < 8$ and $\Delta \log R < 1$.
- Concerning the stability of the scale values, the restricted search range given by k_{range} , implies that the relative variation of this descriptor will always be less than $\Delta \log t \approx 1$.
- For the proximity measure the maximum value is $\sqrt{2} * 0.5 * k_{range} * k_{w1} \approx 5$. With smooth scene motions the value is normally considerably smaller.

Motivated by the fact that the relative variation in S_{patch} is about a factor of ten smaller than the other entities, the relative weights of the components in S_{comb} were set according to the table above.

Note that the correlation measure is the dominant component, and the relative influence of the other components corresponds to about half that variation.

The reason why c_{sign} is increased in blob detection, is that the dimension of the significance measures are different:

$$\begin{aligned} [\tilde{k}_{\gamma-norm}^2] &= [\text{brightness}]^6 \\ [(\nabla_{norm}^2 L)^2] &= [\text{brightness}]^2 \end{aligned}$$

Hence, it is natural to increase the coefficient of $S_{sign} = |\log \frac{R_B}{R_A}|$ by a factor of three in blob detection compared to junction detection. As a general rule, we have not performed any fine-tuning of the parameters, and all parameter values have been the same in all experiments.

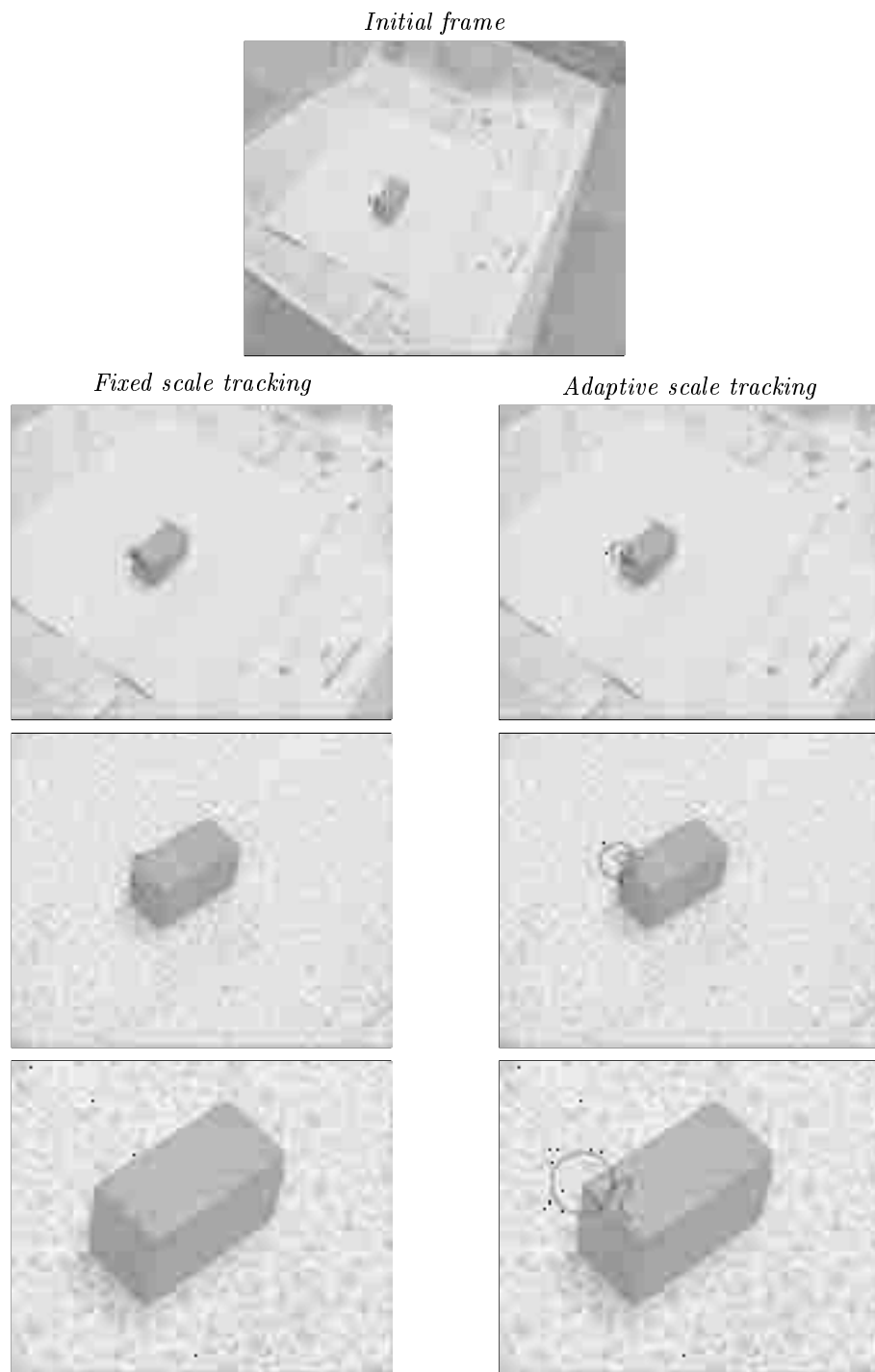


Figure 1: Illustration of the importance of automatic scale selection when tracking image structures over time. The corner is lost using detection at a fixed scale (left column), whereas it is correctly tracked using adaptive scale selection (right column). The size of the circles correspond to the detection scales of the corner features.

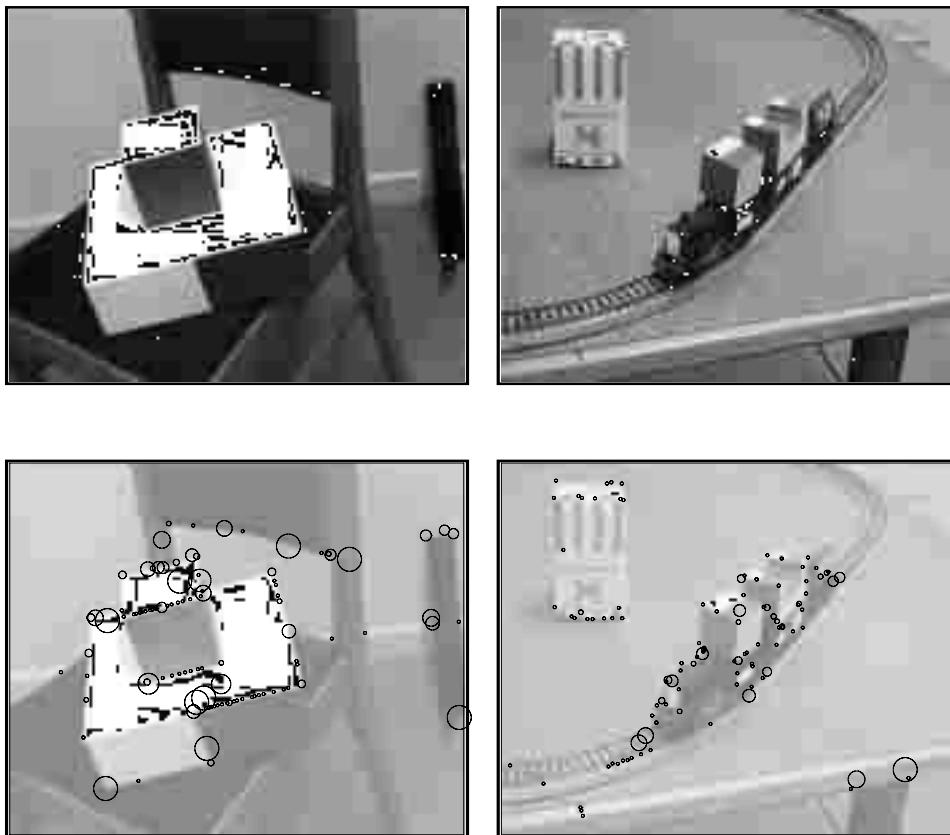


Figure 2: The result of applying the corner detection algorithm with automatic scale selection to two different grey-level images. (top row) Original grey-level images. (bottom row) The 100 most significant corners superimposed onto a bright copy of the original image. Graphically, each corner is illustrated by circle with the radius reflecting the detection scale. Observe that a reasonable set of junction candidates is obtained, and that the circles serve as natural regions of interest around the corners to be used in further processing.

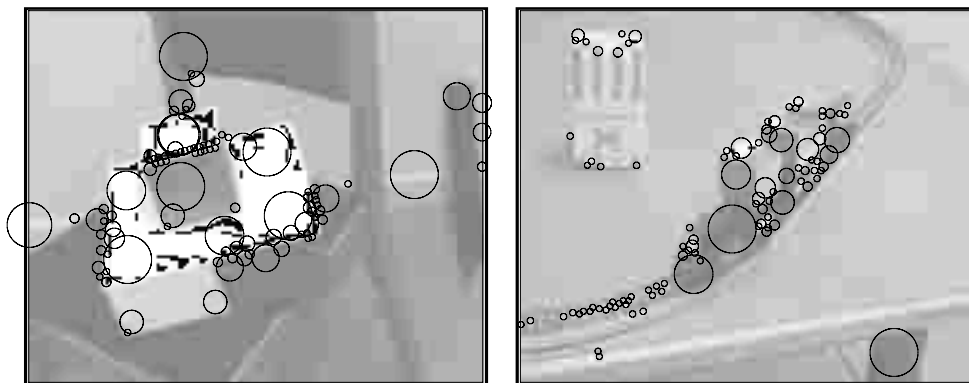


Figure 3: The result of applying the blob detection algorithm with automatic scale selection to the same images as used for corner detection in figure 2. The 100 most significant blobs have been graphically illustrated by circles with their radius proportional to the detection scale.

Algorithm:

For each frame:

For each feature F in the feature set:

1. Prediction

- 1.1 Predict the position of the feature F in the current frame based on information from the previous frames.
- 1.2 Compute the search region in the current frame based on information from the previous frames and the scale of the feature.

2. Detection

Detect n candidates C_k over a reduced set of scales in the region of interest in the current frame.

3. Matching

- 3.1 Match every candidate C_k to the feature F and find the best match using the combined similarity measure.
- 3.2 Optionally, perform bidirectional matching to register safe matches.
- 3.3 Compare the similarity value to a predetermined threshold:
 - If above: consider the feature as matched; update its position, its scale descriptor, its significance value, its grey-level patch and increase its quality value.
 - If below: consider the feature as unmatched; update its position to the predicted position and decrease its quality value.

Parse the feature set to detect feature merges and remove features having quality values below a certain threshold.

Figure 4: Overview of the feature tracking algorithm.

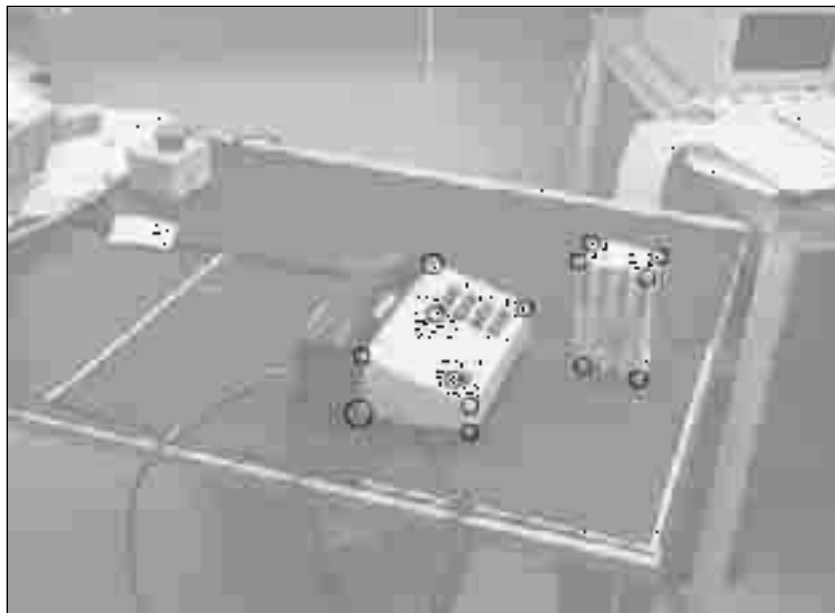


Figure 5: The phone sequence: The initial frame with 14 detected corners.

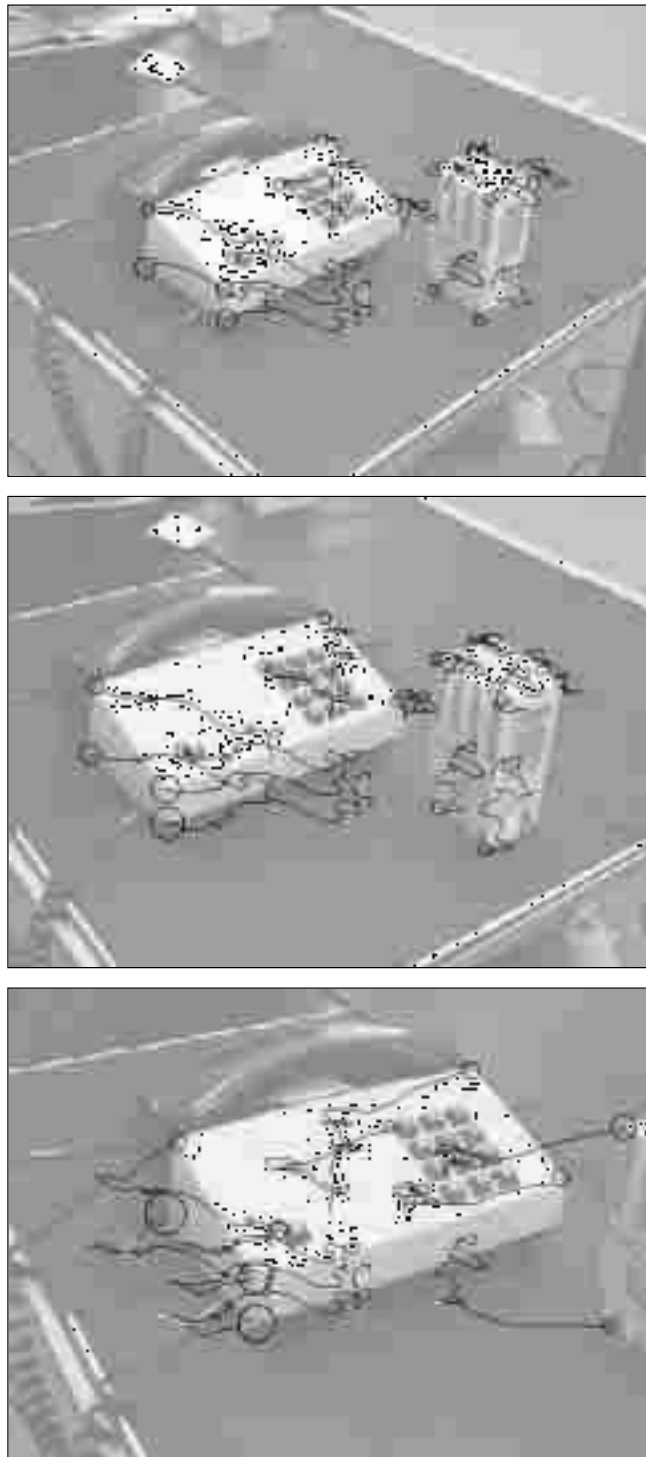


Figure 6: Corner tracking with adaptive scale selection and matching on combined similarity: the tracked corners in the phone sequence after 30 frames (top), 50 (middle) and 60 frames (bottom). As can be seen, all corners are correctly tracked.

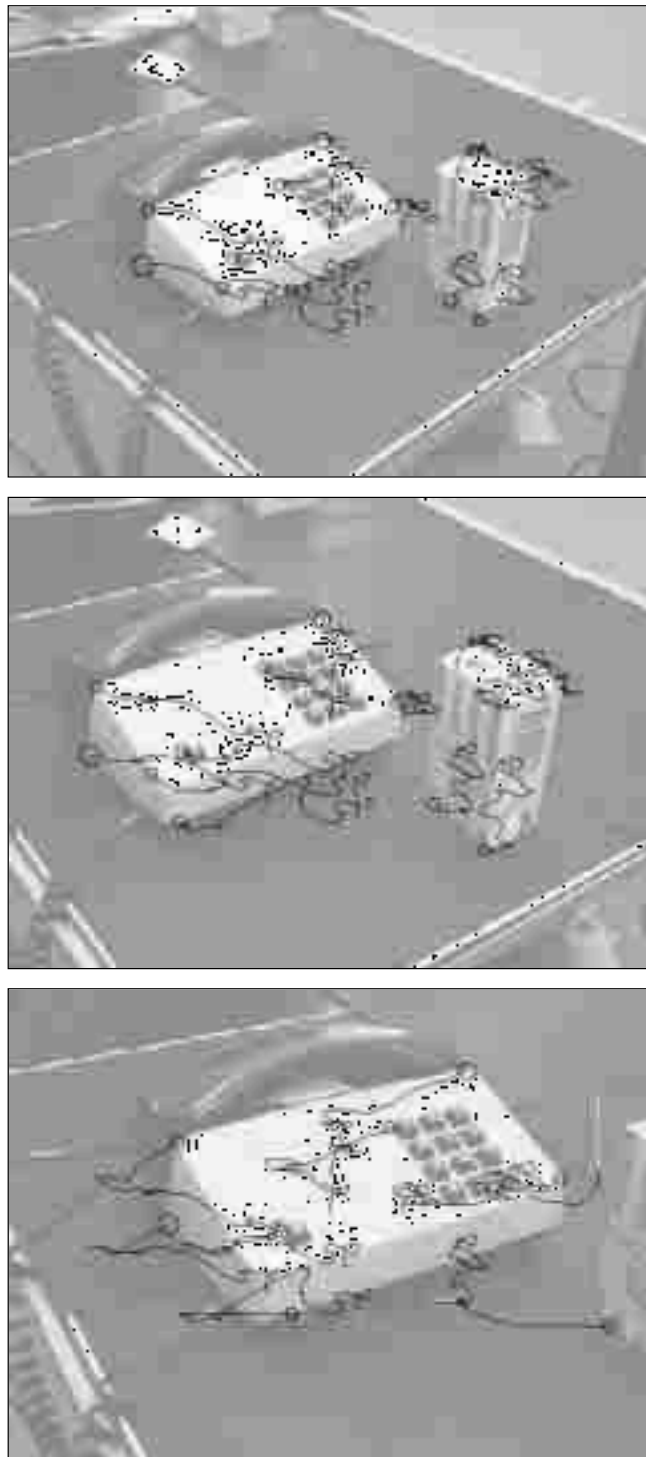


Figure 7: Corner tracking with fixed scales over time: the tracked corners in phone sequence after 30 frames (top), 50 (middle) and 60 frames (bottom). Note that the blunt corners are lost compared to the adaptive scale tracking in figure 6.

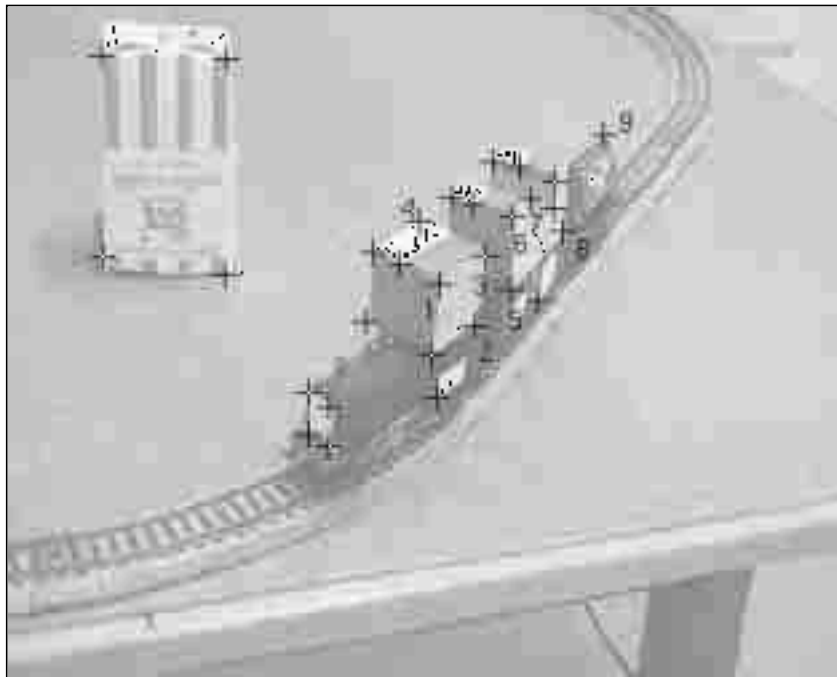


Figure 8: The train sequence: The initial frame with 29 detected corners.

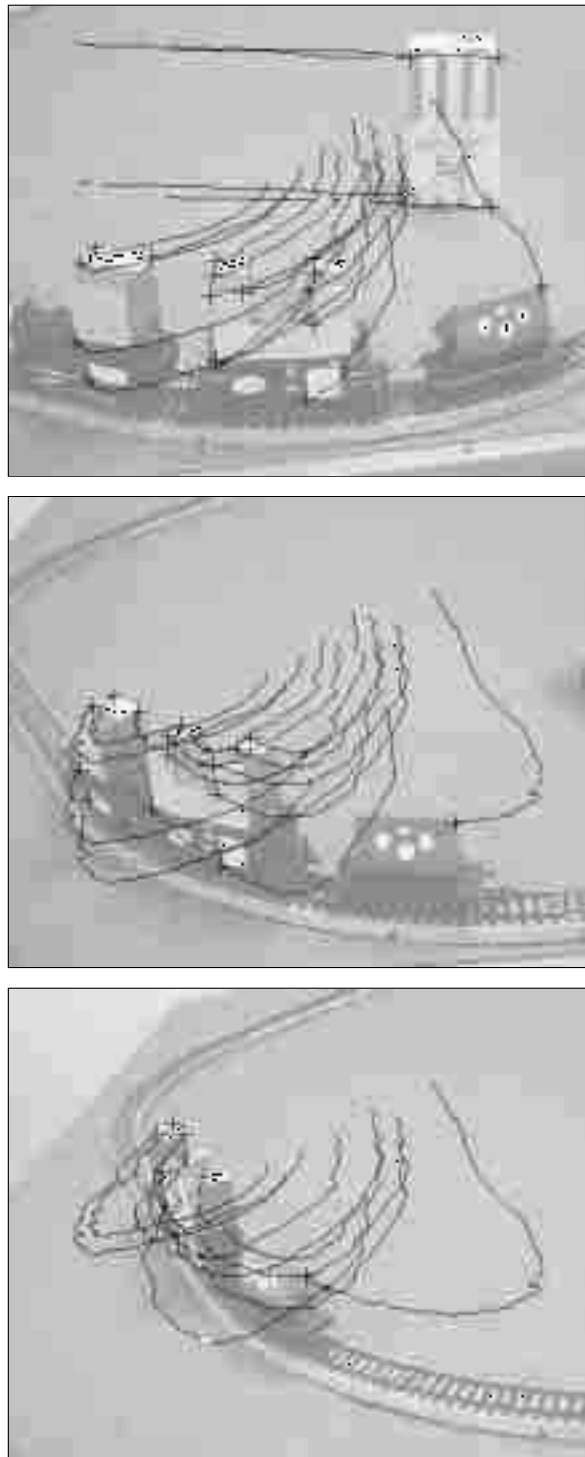


Figure 9: Corner tracking with adaptive scale selection and matching on combined similarity: the tracked corners in the train sequence after 60 frames (top), 100 (middle) and 140 frames (bottom).

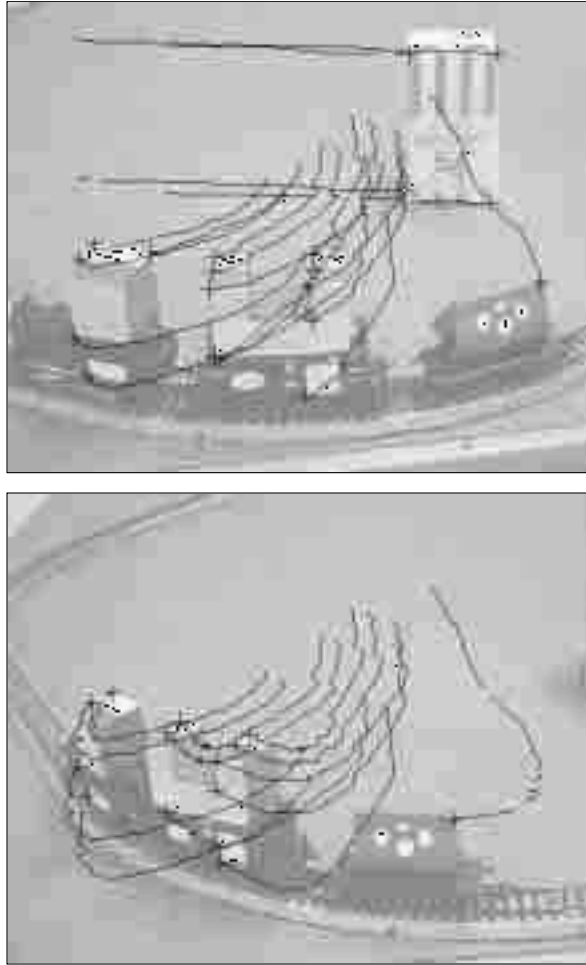


Figure 10: Matching candidates on patch correlation only: the tracked corners in the train sequence after 60 frames (top) and 100 frames (bottom). Three more corners are lost as compared to figure 9.



Figure 11: The train sequence: The initial frame with 13 detected blobs. (The size of the circles correspond to the detection scales of the blob features.)

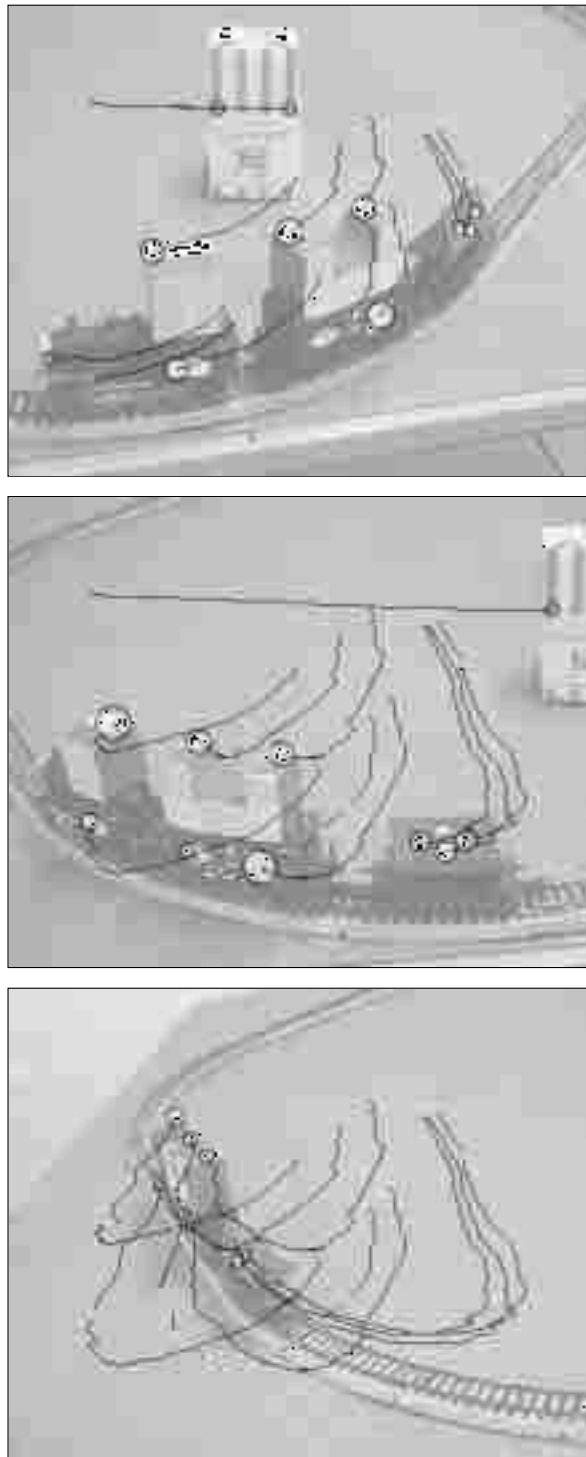


Figure 12: Blob tracking with adaptive scale selection and matching on combined similarity: the tracked blobs in the train sequence after 30 frames (top), 90 (middle) and 150 frames (bottom). All blobs are correctly tracked.

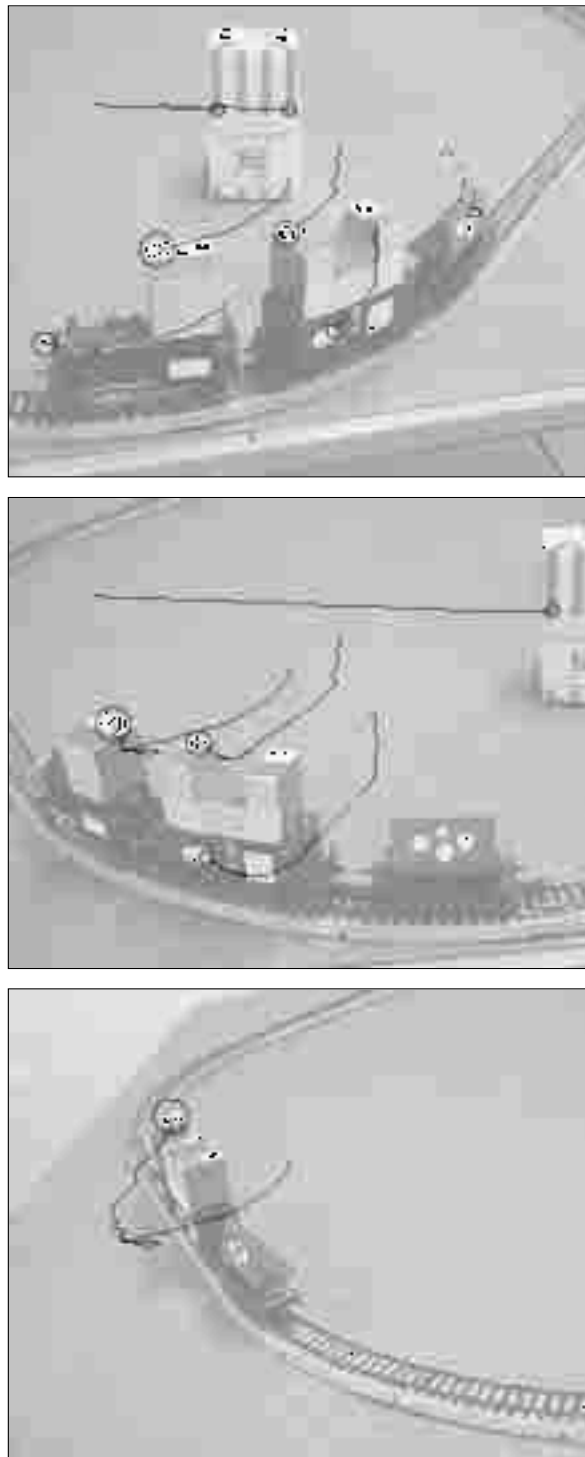


Figure 13: Blob tracking using fixed scales in the detection procedure: the tracked blobs in train sequence after 30 frames (top), 90 (middle) and 150 frames (bottom). Only one blob is correctly tracked over the whole sequence.



Figure 14: The initial frame of the shirt sequence with the 20 strongest blobs detected in a rectangular window. The size of the circles correspond to the detection scales of the blob features.)

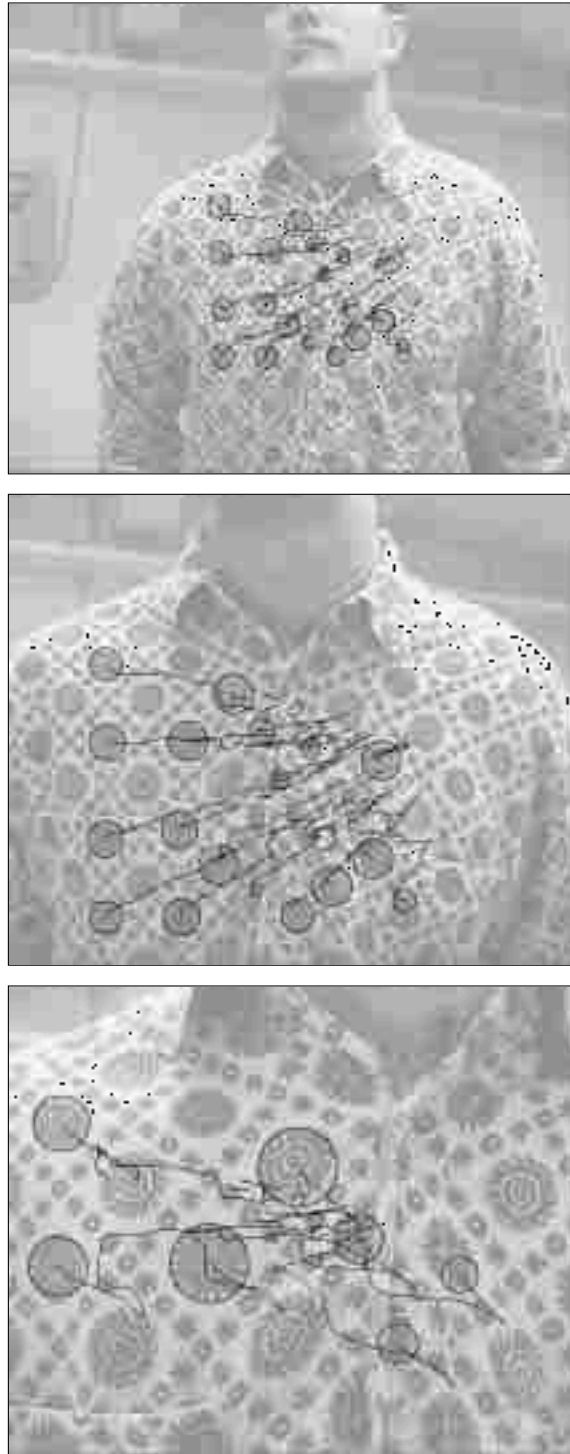


Figure 15: Blob matching using the combined similarity measure: the tracked blobs in the shirt sequence after 25 frames (top), 50 frames (middle) and 87 frames (bottom). Note how the scales, illustrated by the size of the circles, adapt to the size changes of the image structures.



Figure 16: Matching the candidates on patch similarity only: the tracked blobs in the shirt sequence after 25 frames. Compared to the top image in figure 15, three more blobs are lost and one is mismatched.



Figure 17: Blob tracking using fixed scales in the detection procedure: the tracked blobs in the shirt sequence after 25 frames. Most blobs are already lost because they no longer exist at the initially chosen scale.



Figure 18: The initial frame of the face sequence with the 10 most significant blobs detected in a region around the face of the subject.

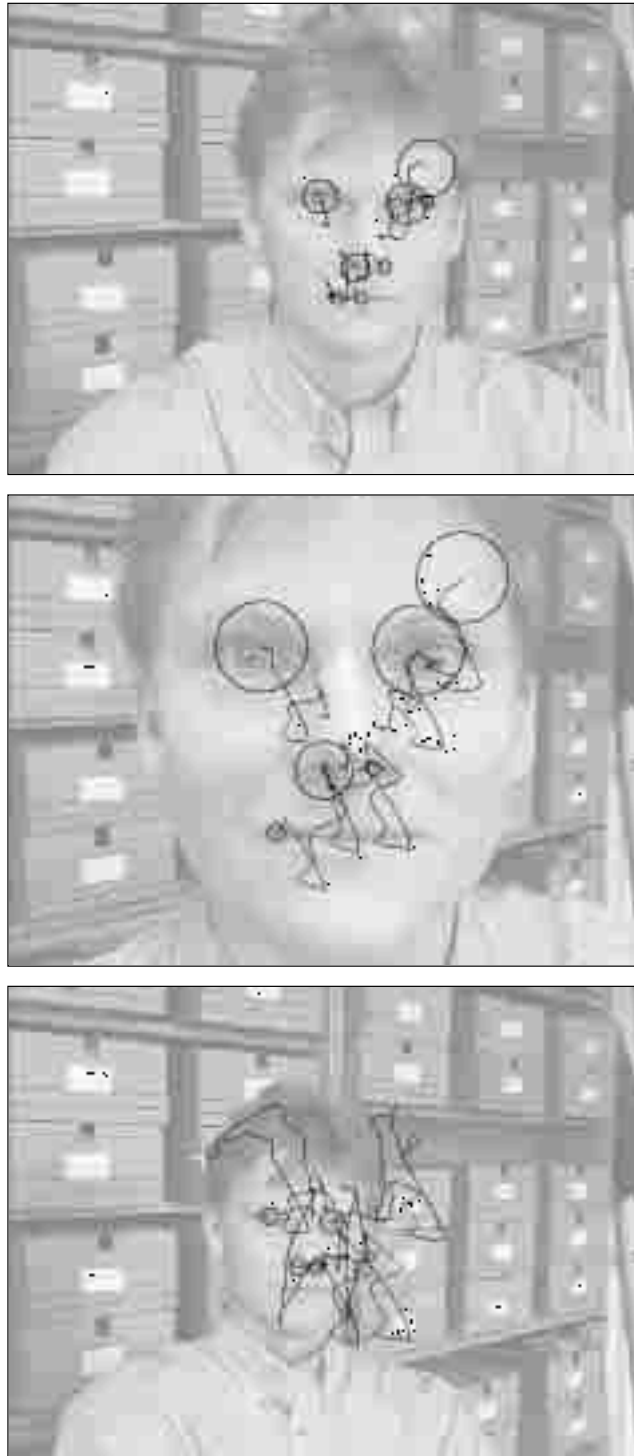


Figure 19: Tracking the blobs in the face sequence with automatic scale selection; the situation after 20, 45 and 90 frames. After about 60 frames only the 4 most stable blobs remain in the feature set.