

# Detecting and Solving Template Ambiguities in Motion Segmentation

Pedro M. Q. Aguiar \*

José M. F. Moura

Instituto de Sistemas e Robótica  
Instituto Superior Técnico  
Lisboa, Portugal  
E-mail : aguiar@isr.ist.utl.pt

Electrical and Computer Engineering  
Carnegie Mellon University  
Pittsburgh, U.S.A.  
E-mail : moura@ece.cmu.edu

## Abstract

*When the color or gray level of a moving object is very similar to that of the background, motion-based segmentation methods fail. This leads to ambiguous templates for the moving objects. In this paper, we propose a method that segments unambiguously from motion the templates of moving objects. Our method, which we call **Incremental Motion Segmentation** integrates over time the small differences between the gray level of the moving object and that of the background. Our experiments with segmenting a robot soccer video clip show the quality of our results.*

## 1 Introduction

We introduced in [1] **Incremental Motion Segmentation (IMS)** and detailed its use in **Generative Video (GV)** [2] analysis. **IMS** achieves motion-based segmentation for general scenes by integrating over time the information content of an image sequence. **GV** reduces video sequences to world images and ancillary data. The world images are augmented views of the world - background world image - and complete views of moving objects - figure world images. The ancillary data registers the world images, stratifies them at each time instant, and positions the camera with respect to the layering of world images. A major task in **GV** analysis is the segmentation from motion of the (background and figures) world images.

When a moving object has color or gray-level that is very similar to the background, current motion-based segmentation methods that use only a few consecutive frames fail. To resolve these problems, approaches like in [3, 4] use statistical regularization techniques or combine motion with other attributes, like color or texture. In general, these methods lead to complex and time consuming algorithms. In contrast, the

method we propose here integrates over time any existing small differences. Our results show that **IMS** has the ability to detect and solve unambiguously the templates for these low contrast objects. We present experiments with real data that construct the **GV** representation for a robot soccer video sequence.

The paper is organized as follows. Section 2 formulates the problem, leading to an ML-based cost-function. The minimization procedure is described in section 3. The **IMS** algorithm is described in section 4. Experimental results and conclusions are in sections 5 and 6.

## 2 Problem formulation

**IMS** is derived as a computationally simple approximation to *Maximum Likelihood* (ML) estimation.

**Observation model.** We motivate **IMS** by considering a single object moving in front of a background, the scene being captured by a moving camera. The image  $I_i$  is modeled as

$$I_i = \left\{ \mathbf{B}(\mathbf{p}_i^\#) \left[ 1 - \mathbf{T}(\mathbf{q}_i^\#) \right] + \mathbf{O}(\mathbf{q}_i^\#) \mathbf{T}(\mathbf{q}_i^\#) + \mathbf{W}_i \right\} \mathbf{H} \quad (1)$$

where  $\mathbf{B}$  and  $\mathbf{O}$  are the background and object world images [2],  $\mathbf{T}$  is the object template,  $\mathbf{p}_i$  and  $\mathbf{q}_i$  are the camera pose and the object position, and  $\mathbf{W}_i$  is zero mean white Gauss noise. We assume that  $I_i(x, y) = 0$  for  $(x, y)$  outside the region observed by the camera.  $\mathbf{H}$  is such that  $\mathbf{H}(x, y) = 1$  for  $(x, y)$  in the region observed in the image and  $\mathbf{H}(x, y) = 0$  otherwise. We denote by  $\mathbf{I}(\mathbf{p})$  the registration of the image  $\mathbf{I}$  according to the position vector  $\mathbf{p}$ . The pixel  $(x, y)$  of the registered image is denoted by  $\mathbf{I}(\mathbf{p}; x, y)$ . The registration of  $\mathbf{I}(\mathbf{p})$  according to the position vector  $\mathbf{q}$  is denoted by  $\mathbf{I}(\mathbf{q}, \mathbf{p})$ . We denote by  $\mathbf{I}(\mathbf{p}^\#)$  the image that registered according to  $\mathbf{p}$  equals  $\mathbf{I}$ , so we have  $\mathbf{I}(\mathbf{p}, \mathbf{p}^\#) = \mathbf{I}$ .

**Maximum likelihood estimate.** The problem is to estimate from the  $N$  images  $\mathbf{I}_1, \mathbf{I}_2, \dots, \mathbf{I}_N$  the back-

---

\*The work of the first author was partially supported by NATO.

ground world image  $\mathbf{B}$ , the object world image  $\mathbf{O}$ , the object template  $\mathbf{T}$ , the camera pose  $\{\mathbf{p}_i\}$ , and the object position  $\{\mathbf{q}_i\}$ . These quantities define the  $\mathbf{GV}$  representation of the sequence. The ML estimate minimizes the functional  $J$  over all  $\mathbf{GV}$  parameters

$$J = \sum_{x,y} \sum_{i=1}^N \left\{ \mathbf{I}_i - \mathbf{B}(\mathbf{p}_i^\#) \left[ \mathbf{1} - \mathbf{T}(\mathbf{q}_i^\#) \right] - \mathbf{O}(\mathbf{q}_i^\#) \mathbf{T}(\mathbf{q}_i^\#) \right\}^2 \mathbf{H} \quad (2)$$

where the inner sum is over the full set of  $N$  images and the outer sum is over all pixels observed in each image. We omit the dependence on  $(x, y)$  to simplify the notation.

Expression (1) models explicitly the occlusion of the background by the moving object. The estimation of the parameters of (1) using  $N$  images rather than a single pair of images is a feature that distinguishes our work from other techniques which use only two or three consecutive frames.

The simultaneous minimization of the functional  $J$  over the estimation parameters  $\mathbf{B}, \mathbf{O}, \mathbf{T}, \{\mathbf{p}_i\}$ , and  $\{\mathbf{q}_i\}$  is very complex. To simplify this task and motivated by our experience with real video sequences, we decouple the estimation of the motion of the camera and the motion of the moving object from the estimation of the remaining parameters. Any method that estimates motions of multiple objects can be used with our **IMS** framework. See reference [5] for a survey. We estimate motion vectors by an affine model. To cope with large displacements, we use a spatial multi-resolution pyramid. To deal with multiple moving objects, we use a *quad-tree* decomposition.

### 3 Minimization Procedure

The position vectors  $\{\mathbf{p}_i\}$  and  $\{\mathbf{q}_i\}$  are estimated as explained in section 2, so that in this section we assume that they are known. The problem becomes the minimization of  $J$ , given by expression (2), with respect to the background world image  $\mathbf{B}$ , the world image  $\mathbf{O}$  of the moving object, and its template  $\mathbf{T}$ .

We express the estimate  $\hat{\mathbf{O}}$  of the object world image that minimizes  $J$  in terms of the template  $\mathbf{T}$ . Then we replace  $\hat{\mathbf{O}}$  in expression (2), getting  $J$  in terms of  $\mathbf{B}$  and  $\mathbf{T}$ . To solve for  $\hat{\mathbf{B}}$  and  $\hat{\mathbf{T}}$ , we use a two-step iterative method. The steps involved are: i) solve for  $\hat{\mathbf{B}}$  with fixed  $\hat{\mathbf{T}}$ , and ii) solve for  $\hat{\mathbf{T}}$  with fixed  $\hat{\mathbf{B}}$ . We obtain closed-form solutions for both steps i) and ii) in contrast to the solution for the simultaneous minimization of  $J$  with respect to  $\mathbf{B}$  and  $\mathbf{T}$ .

As it is usual with this kind of methods, although the value of  $J$  decreases along the iterative process, the convergence to the global minimum is not guaranteed,

being necessary to provide a good initialization. We compute a first estimate of the moving object template by using the method described in [1]. It estimates  $\hat{\mathbf{T}}$  as the union of the regions detected as being in movement over a set of frames.

**Object / background world images.** We consider now the estimation of the world image of the moving object and of the world image of the background, given the template  $\mathbf{T}$ .

The estimate  $\hat{\mathbf{O}}$  of the object world image is found by setting  $\nabla_{\mathbf{O}} J = \mathbf{0}$ , with  $\nabla$  the gradient operator

$$\hat{\mathbf{O}} = \mathbf{T} \frac{1}{N} \sum_{i=1}^N \mathbf{I}_i(\mathbf{q}_i) \quad (3)$$

This expression averages the observations  $\mathbf{I}_i$  registered according to the motion  $\mathbf{q}_i$  of the object in the region corresponding to the template  $\mathbf{T}$  of the moving object.

The background world image estimate  $\hat{\mathbf{B}}$  is found by minimizing  $J$  in (2), given the template  $\mathbf{T}$ . After algebraic manipulations, we obtain

$$\hat{\mathbf{B}} = \frac{\sum_{i=1}^N \left[ \mathbf{1} - \mathbf{T}(\mathbf{p}_i, \mathbf{q}_i^\#) \right] \mathbf{I}_i(\mathbf{p}_i)}{\sum_{i=1}^N \left[ \mathbf{1} - \mathbf{T}(\mathbf{p}_i, \mathbf{q}_i^\#) \right] \mathbf{H}(\mathbf{p}_i)} \quad (4)$$

The estimate averages the observations  $\mathbf{I}_i$  registered with the background motion  $\mathbf{p}_i$ , in the regions  $\{(x, y)\}$  not occluded by the moving object, i.e., when  $\mathbf{T}(\mathbf{p}_i, \mathbf{q}_i^\#; x, y) = 0$ . In the denominator,  $\mathbf{H}(\mathbf{p}_i)$  weighs differently regions that, due to camera motion, are not seen in all images.

**Template estimation.** Replacing the object world image estimate  $\hat{\mathbf{O}}$  given by (3) in (2), we express the functional  $J$  in terms of the background world image  $\mathbf{B}$  and the object template  $\mathbf{T}$ . The complete derivation is given in [6]. We obtain

$$J = \sum_{x,y} \mathbf{T}(x, y) \mathbf{Q}(x, y) + \text{Constant} \quad (5)$$

$$\mathbf{Q}(x, y) = \mathbf{Q}^1(x, y) - \mathbf{Q}^2(x, y) \quad (6)$$

$$\mathbf{Q}^1(x, y) = \frac{1}{N} \sum_{i=2}^N \sum_{k=1}^{i-1} \left[ \mathbf{I}_i(\mathbf{q}_i; x, y) - \mathbf{I}_k(\mathbf{q}_k; x, y) \right]^2 \quad (7)$$

$$\mathbf{Q}^2(x, y) = \sum_{i=1}^N \left[ \mathbf{I}_i(\mathbf{q}_i; x, y) - \mathbf{B}(\mathbf{q}_i, \mathbf{p}_i^\#; x, y) \right]^2 \quad (8)$$

Consider the minimization of  $J$  given by (5) over the template  $\mathbf{T}$ , given the background world image  $\mathbf{B}$ . It is clear from (5), that the minimization of  $J$  with respect to each spatial location of  $\mathbf{T}$  is independent

from the minimization over the other locations. The template  $\hat{\mathbf{T}}$  that minimizes  $J$  is given by the following test evaluated at each pixel:

$$\mathbf{Q}^1(x, y) \begin{array}{c} \hat{\mathbf{T}}(x, y) = 0 \\ > \\ < \\ \hat{\mathbf{T}}(x, y) = 1 \end{array} \mathbf{Q}^2(x, y) \quad (9)$$

In the spatial locations where the differences between each frame  $\mathbf{I}_i(\mathbf{q}_i)$  and the background  $\mathbf{B}(\mathbf{q}_i, \mathbf{p}_i^\#)$  are greater than the differences between each pair of co-registered frames  $\mathbf{I}_n(\mathbf{q}_n)$  and  $\mathbf{I}_k(\mathbf{q}_k)$ , we estimate  $\hat{\mathbf{T}}(x, y) = 1$ , these regions belong to the moving object. If not, the regions belong to the background.

For regions with low contrast between the moving object and the background, the test may initially, for a small number of frames, be inconclusive ( $\mathbf{Q}^1(x, y) \simeq \mathbf{Q}^2(x, y)$ ). As more frames are processed and small differences accumulate, the initial test ambiguity is resolved and the region is assigned to the moving object or to the background.

#### 4 IMS Algorithm

After processing a number of images, the estimate  $\hat{\mathbf{T}}$  does not change significantly as new images become available. To see this, we replace the image model given by (1) in (7) and (8), and compute the expected value  $\mathbb{E}\{\mathbf{Q}\} = \mathbb{E}\{\mathbf{Q}_1 - \mathbf{Q}_2\}$  with respect to the observation noise. See [6] for the complete derivation. We obtain

$$\mathbb{E}\{\mathbf{Q}\} = [\mathbf{1} - \mathbf{T}] \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^{i-1} \left[ \mathbf{B}(\mathbf{q}_i, \mathbf{p}_i^\#) \mathbf{H}(\mathbf{q}_i) - \mathbf{B}(\mathbf{q}_k, \mathbf{p}_k^\#) \mathbf{H}(\mathbf{q}_k) \right]^2 - \mathbf{T} \sum_{i=1}^N \left[ \mathbf{O} - \mathbf{B}(\mathbf{q}_i, \mathbf{p}_i^\#) \right]^2 + \text{residual independent of } \mathbf{T} \quad (10)$$

As we process more images,  $\mathbb{E}\{\mathbf{Q}(x, y)\}$  becomes more negative in the spatial locations where  $\mathbf{T}(x, y) = 1$ , and more positive where  $\mathbf{T}(x, y) = 0$  (note that the coefficients that multiply  $\mathbf{T}$  and  $[\mathbf{1} - \mathbf{T}]$  in (10) are non-negative). After processing a sufficiently large number of frames the sign of  $\mathbf{Q}$  at each pixel  $(x, y)$  stabilizes, which in turn, through test (9), stabilizes the estimate  $\hat{\mathbf{T}}$  of the template. These considerations motivate splitting the algorithm in two phases - template acquisition and recursive world image updating. **Template acquisition.** This phase detects and solves the template ambiguities mentioned before. It implements the sequential method described in section 3 with an increasing larger set of frames. The background world image estimate  $\hat{\mathbf{B}}$  is given by expression (4), replacing the previous estimate  $\hat{\mathbf{T}}$ . The template estimate  $\hat{\mathbf{T}}$  is given by

the test (9), replacing the previous estimate  $\hat{\mathbf{B}}$ . The template acquisition phase lasts until the test  $\mathbf{Q}_n^1(x, y) \gtrless \mathbf{Q}_n^2(x, y)$  is conclusive at all spatial locations, leading to the template estimate  $\hat{\mathbf{T}}(x, y)$ .

**Recursive world image generation.** From expressions (3) and (4), we obtain the following recursive updates, for fixed  $\mathbf{T} = \hat{\mathbf{T}}$ :

$$\hat{\mathbf{O}}_n = \frac{n-1}{n} \hat{\mathbf{O}}_{n-1} + \hat{\mathbf{T}} \frac{1}{n} \mathbf{I}_n(\mathbf{q}_n) \quad (11)$$

$$\mathbf{S}_n^b = \mathbf{S}_{n-1}^b + \left[ \mathbf{1} - \hat{\mathbf{T}}(\mathbf{p}_n, \mathbf{q}_n^\#) \right] \mathbf{H}(\mathbf{p}_n) \quad (12)$$

$$\hat{\mathbf{B}}_n = \frac{\mathbf{S}_{n-1}^b}{\mathbf{S}_n^b} \hat{\mathbf{B}}_{n-1} + \frac{[\mathbf{1} - \hat{\mathbf{T}}(\mathbf{p}_n, \mathbf{q}_n^\#)]}{\mathbf{S}_n^b} \mathbf{I}_n(\mathbf{p}_n) \quad (13)$$

where we define the weights  $\mathbf{S}_n^b$  as the denominator of (4), and the element  $\mathbf{S}_n^b(x, y)$  is the number of times the pixel  $\mathbf{B}(x, y)$  is observed in the first  $n$  frames.

#### 5 Experimental Results

We test the algorithm with a sequence of 20 images from a robot soccer game, see [7], showing a robot pursuing a ball, see figure 1. Due to the white stripes on the field, the robot template is ambiguous during the first frames of the sequence. Only after the robot rotates, is it possible to determine its template.

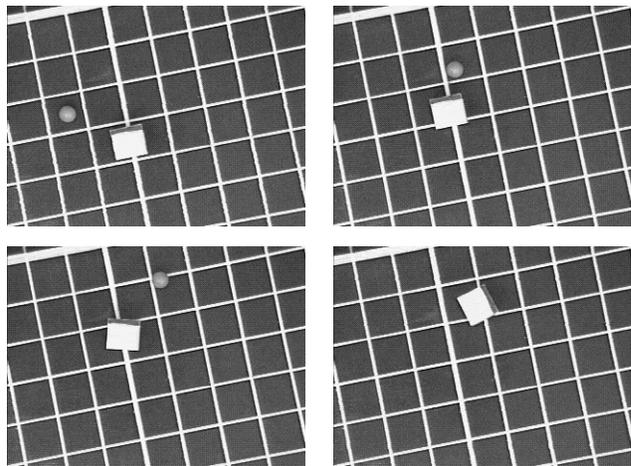


Figure 1: Image sequence. Frames 1, 4, 8, 16.

The initialization phase takes 5 frames. It detects three motions, the static background, the moving ball, and the moving robot. According to test (9), the ball template is unambiguous after the 5 frames used in the initialization step. Figure 2 shows the evolution of the robot template. Regions where the test is inconclusive are grey, regions classified as being part of the robot template are white. The black regions are

classified as either background or as belonging to the ball template. The acquisition phase is ended after 10 frames. The final robot template estimate is shown on the right side of figure 2.

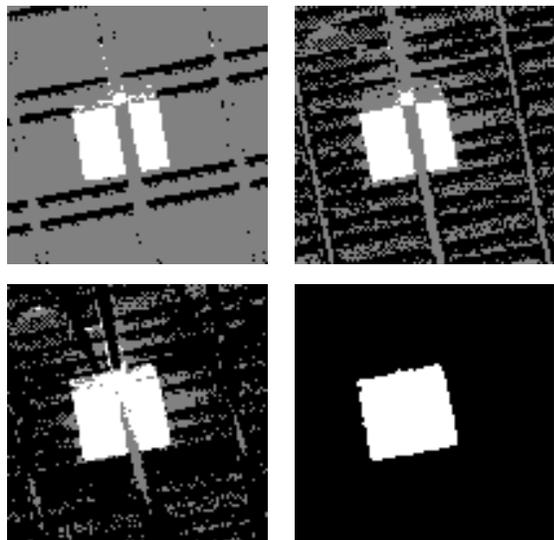


Figure 2: Robot template after frames 2, 4, 6, 10.

Figure 3 illustrates the evolution of the matrix  $Q$ . The curves on the left side represent the value of  $Q(x, y)$  for several pixels  $(x, y)$  that belong to the template of the robot. These curves start close to zero and decrease with the number of frames processed, as predicted by expression (10). The curves on the right side of figure 3 were obtained for several pixels that do not belong to the template of the robot. For these pixels,  $Q(x, y)$  increases with the number of frames.

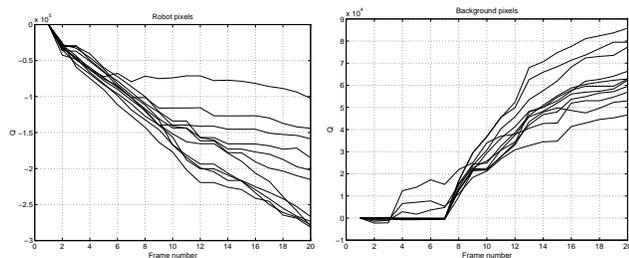


Figure 3: Evolution of  $Q$ .

Figure 4 shows the world images for the two moving objects and background, after processing the entire sequence of 20 frames.

## 6 Conclusion

The **Incremental Motion Segmentation** algorithm described here achieves segmentation from mo-

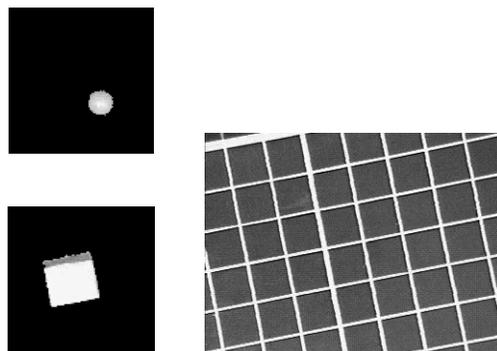


Figure 4: Ball, robot and background world images.

tion of objects with low contrast. It is computationally simple, resolving the ambiguities that arise by integrating over time the information contents of the image sequence. The experimental results show that IMS is suitable for 2D object-based video analysis.

## References

- [1] P. M. Q. Aguiar and J. M. F. Moura. Incremental motion segmentation in low texture. In *IEEE Int. Conf. Image Processing*, volume I, pages 233–236, Switzerland, Sep. 1996.
- [2] R. S. Jasinschi and J. M. F. Moura. Content-based video sequence representation. In *IEEE Int. Conf. Image Processing*, volume II, pages 229–232, U.S.A., Oct. 1995.
- [3] P. Bouthemy and E. François. Motion segmentation and qualitative dynamic scene analysis from an image sequence. *Int. J. Computer Vision*, 10(2):157–182, 1993.
- [4] M-P. Dubuisson and A. K. Jain. Contour extraction of moving objects in complex outdoor scenes. *Int. J. Computer Vision*, 14(1):83–105, 1995.
- [5] S. Ayer. *Sequential and Competitive Methods for Estimation of Multiple Motions*. PhD thesis, École Polytechnique Fédérale de Lausanne, 1995.
- [6] P. M. Q. Aguiar and J. M. F. Moura. Moving objects segmentation in low texture / low contrast video. To be submitted.
- [7] M. Veloso, P. Stone, S. Achim, and M. Bowling. A layered approach for an autonomous robotic soccer system. In *1st Int. Conf. Autonomous Agents*, U.S.A., Feb. 1997.