

JOINT SEGMENTATION OF MOVING OBJECT AND ESTIMATION OF BACKGROUND IN LOW-LIGHT VIDEO USING RELAXATION

Pedro M. Q. Aguiar

Institute for Systems and Robotics / IST
Lisboa, Portugal
E-mail: aguiar@isr.ist.utl.pt

José M. F. Moura

ECE Dep., Carnegie Mellon University
Pittsburgh PA, USA
E-mail: moura@ece.cmu.edu

ABSTRACT

When the scene background is known and the intensity of moving objects contrasts with the intensity of the background, the objects are easily captured by exploiting occlusion, *e.g.*, background-subtraction. However, when processing general scenes, the background is not known and researchers have mostly attempted to segment moving objects by using motion cues rather than occlusion. Since motion can only be accurately computed at highly textured regions, current motion segmentation methods either fail to segment low textured objects, or require expensive regularization techniques. We present a computationally simple algorithm and test it with segmentation of moving objects in low texture / low contrast videos that are obtained in low-light scenes. The images in the sequence are modeled taking into account the *rigidity* of the moving object and the *occlusion* of the background. We formulate the problem as the minimization of a *penalized likelihood* cost. Relaxation of the weight of the penalty term leads to a simple solution to the nonlinear minimization. We describe experiments that illustrate the good performance of our method.

1. INTRODUCTION

Modern content-based video representations demand efficient methods to infer the contents of video sequences, like the shape and texture of objects and their motions [1]. In this paper we address this problem of segmenting out moving objects from video. Although methods that require human interaction have lead to good results, fully automatic methods are still being investigated. In particular, we seek methods capable of dealing with challenging video sequences, as those obtained in low-light scenarios. In fact, when the scene is poorly illuminated, which happens in many situations, *e.g.*, evening shots, only low contrast images can be obtained¹.

Related work Among the approaches to segmentation of moving objects, so-called background-subtraction methods are very appealing due to their simplicity. These methods capture the moving objects by subtracting the input image from a previously stored background [2, 3]. Although background-subtraction succeeds in relevant situations, *e.g.*, surveillance applications, the background knowledge requirement limits its application to general scenes.

Some approaches in the computer vision literature cope with low texture by using prior knowledge about the scenes, *e.g.*, statistical

This work was partially supported by ONR grant N000 14-00-1-0593 and FCT ISR/IST plurianual funding, POSC, FEDER.

¹Naturally, image contrast can be artificially increased through post-processing but the signal-to-noise ratio remains the same, *i.e.*, low-light images are inherently very noisy.

regularization, or by combining motion with other attributes. In general, these methods lead to complex and time consuming algorithms. Another popular trend uses active contours—the contour of the object is computed by minimizing a global cost function, thus leading to robust estimates [4]. The drawback of these approaches is that the global cost minimization requires calculus of variations, making the algorithms computationally expensive.

Layered models brought new approaches to the segmentation of moving objects. For example, the work in [5] uses an offline approach to infer flexible templates. To cope with the very large dimensionality of the search space, the authors restrict the motions to single pixel translations.

Reference [6] uses temporal integration by averaging the images registered according to the motion of the objects in the scene. After processing a number of frames, each of these integrated images is expected to show only one sharp region corresponding to the tracked object. The object is found by detecting the stationary regions between the integrated image and the current frame. Unless the background is textured enough to blur completely the averaged images, some regions of the background can be misclassified as stationary. In this situation, the method of [6] overestimates the template of the moving object. This is particularly likely to happen when the background has large regions with almost constant color.

In [7], we segmented moving objects using *Maximum Likelihood* (ML) estimation. ML estimation succeeds even when there is little contrast between the moving object and the background because it integrates small differences over time. However, the minimum of the ML cost function is not always sharply defined. In fact, the likelihood that some region belongs to the background may be very similar to the likelihood that the same region belongs to the moving object, when that region has low texture.

Approach Like the appealing and simple background-subtraction algorithms, the approach we present in this paper exploits the fact that the moving object occludes the background. However, unlike these, we do not assume that the background is known *a priori*. Like some of the elegant and robust computer vision approaches outlined above, we formulate segmentation in a global way, as a parameter estimation problem. However, unlike these, we do not use complex and computationally expensive algorithms to compute the object shape. Since in several situations the shape of the moving object does not change across a number of frames, *e.g.*, moving cars, we also exploit the object rigidity. In the paper we show how *occlusion+rigidity* enable a computationally simple algorithm to jointly estimate the unknown background and shape of the moving object, directly from the image intensity values.

The segmentation algorithm is derived as an approximation to

a *penalized likelihood* (PL) [8] estimate of the parameters involved in the video model: the motions, the template of the moving object, and the intensity levels of the object pixels (object texture) and of the background pixels (background texture). Our PL cost function balances two terms. The first term is the ML cost introduced in [7]. It measures the error between the observed data and the model. The second term measures the size of the moving object, *i.e.*, the area of its template. By incorporating the penalization term, we make the segmentation problem well-posed: we look for the *smallest* template that describes well the observed data.

To minimize the PL cost function, we describe a computationally simple algorithm that performs alternatively two estimation steps for which we derive closed form solutions. The penalization term has a second very relevant feature—it improves the convergence of the two-step iterative algorithm. As for any iterative minimization algorithm, the initial guess is a very relevant issue to the good convergence of the process. By using a relaxation strategy for the weight of the penalization term, we avoid using computationally expensive ad-hoc methods to compute initial estimates. Our experience shows that this strategy makes the behavior of the algorithm quite insensitive to the initial guess, so much so that it suffices to initialize the process by the trivial guess of having no moving object, *i.e.*, every pixel belonging to the background.

2. PROBLEM FORMULATION

We consider 2-D parallel motions. We represent the motions by specifying position vectors that code rotation-translation pairs. The image obtained by applying the rigid motion coded by the vector \mathbf{p} to the image \mathbf{I} is denoted by $\mathcal{M}(\mathbf{p})\mathbf{I}$ (registration). The registration of $\mathcal{M}(\mathbf{p})\mathbf{I}$ according to the vector \mathbf{q} is denoted by $\mathcal{M}(\mathbf{q}\mathbf{p})\mathbf{I}$. We denote the inverse of \mathbf{p} by $\mathbf{p}^\#$, thus the registration of $\mathcal{M}(\mathbf{p})\mathbf{I}$ according to $\mathbf{p}^\#$ obtains the original image \mathbf{I} , *i.e.*, $\mathcal{M}(\mathbf{p}^\#\mathbf{p})\mathbf{I} = \mathbf{I}$.

Observation model Consider a scene with a moving object in front of a moving camera. Each pixel of each image belongs either to the background or to the moving object. Thus, the frame \mathbf{I}_f is

$$\mathbf{I}_f = \left\{ \mathcal{M}(\mathbf{p}_f^\#)\mathbf{B} \left[\mathbf{1} - \mathcal{M}(\mathbf{q}_f^\#)\mathbf{T} \right] + \mathcal{M}(\mathbf{q}_f^\#)\mathbf{O} \mathcal{M}(\mathbf{q}_f^\#)\mathbf{T} + \mathbf{W}_f \right\} \mathbf{H}, \quad (1)$$

where \mathbf{T} is the object template (a binary image defining the region occupied by the moving object), \mathbf{B} and \mathbf{O} represent the patterns of intensity levels, *i.e.*, the textures, of the background and of the moving object, \mathbf{p}_f and \mathbf{q}_f are the camera pose and the object position, \mathbf{W}_f stands for the observation noise, assumed Gaussian, zero mean, and white, \mathbf{H} is a binary window that defines the rectangular field of view of each image, and $\mathbf{1}$ is constant with value 1.

Penalized likelihood estimation Given F frames $\{\mathbf{I}_f\}$, we want to estimate the background \mathbf{B} , the object \mathbf{O} , the template \mathbf{T} , the camera poses $\{\mathbf{p}_f\}$, and the object positions $\{\mathbf{q}_f\}$. The problem as just stated may be ill-posed. As an example, consider that the object moves in front of a constant intensity background, *i.e.*, the background has no texture. This image sequence is indistinguishable from an image sequence where the object template is arbitrarily enlarged with pixels whose intensity equals the intensity of the background. Without additional knowledge, it is not possible to decide whether a pixel with intensity equal to the background intensity belongs to the background or to the moving object, *i.e.*, no algorithm can segment unambiguously the moving object. Although extreme, this example illustrates the difficulties of segmenting objects from backgrounds that have large patches with low texture.

To address this issue, we assume that the object is small. This is in agreement with what the human visual system usually implicitly

assumes. We incorporate this constraint into the segmentation problem by minimizing a cost function given by

$$C_{\text{PL}} = C_{\text{ML}} + \alpha \text{Area}(\mathbf{T}), \quad (2)$$

where C_{ML} is the ML cost function studied in [7], α is a non-negative weight, and $\text{Area}(\mathbf{T})$ is the area of the template. Minimizing the cost C_{PL} balances the agreement between the observations and the model (term C_{ML}) with minimizing the area of the template. The term $\alpha \text{Area}(\mathbf{T})$ can be interpreted as a Bayesian prior and the cost function (2) as the negative log posterior probability whose minimization leads to the Maximum a Posteriori estimate, as usual in Bayesian inference approaches. It can also be motivated through information-theoretic criteria like Akaike's AIC or the Minimum Description Length principle. Different basic principles lead to different choices for the parameter α but the structure of the cost function is still as in (2). Statisticians usually call the generic form (2) a *penalized likelihood* (PL) cost function [8]. Our choice for the weight α is discussed below.

The minimization of the functional C_{PL} in (2) with respect to (wrt) $\{\mathbf{B}, \mathbf{O}, \mathbf{T}, \{\mathbf{p}_f, \mathbf{q}_f\}\}$ is a highly complex task. To obtain a computationally feasible algorithm, we decouple the estimation of the motions $\{\mathbf{p}_f, \mathbf{q}_f\}$ from the determination of $\mathbf{B}, \mathbf{O}, \mathbf{T}$. This is reasonable from a practical point of view and is well supported by our experimental results with real videos. The rationale behind the simplification is that the motions can be usually inferred without knowing precisely the object template. To make this point clearer, consider an image sequence with no prior knowledge available, except that an object moves wrt an unknown background. Even with no spatial cues, *e.g.*, if the background and object textures are white noise random fields, the human visual system can easily infer the motion of the background and the motion of the object from only two consecutive frames. However, this is not the case wrt the template of the moving object: to infer an accurate template we need a much higher number of frames that enables us to easily capture the *rigidity* of the object across time. This observation motivated our approach of decoupling the estimation of the motions from the estimation of the remaining parameters. We compute the motions, frame by frame, using a simple sequential method. We first compute the dominant motion, which corresponds to the motion of the background. Then, after compensating for the background motion, we compute the object motion. We estimate the parameters describing both motions by using standard LS techniques. After estimating the motions, we introduce the motion estimates into the PL cost (2) and minimize wrt the remaining parameters. Clearly, this solution is sub-optimal, in the sense that it is an approximation to the PL estimate of the entire set of parameters, and it can be thought of as an initial guess for the minimizer of C_{PL} . This initial estimate can then be refined by using a greedy approach. We emphasize that the key problem here is finding the initial guess in an expedite way, not the final refinement.

3. MINIMIZATION PROCEDURE

We now address the minimization of C_{PL} in (2). Re-write C_{PL} as

$$C_{\text{PL}} = C_{\text{ML}} + \alpha \iint \mathbf{T}(x, y) dx dy. \quad (3)$$

To carry out the minimization, first note that the second term in (3) does not depend on \mathbf{O} , neither on \mathbf{B} , so, for fixed \mathbf{T} , we get $\hat{\mathbf{O}}_{\text{PL}} = \hat{\mathbf{O}}_{\text{ML}}$ and $\hat{\mathbf{B}}_{\text{PL}} = \hat{\mathbf{B}}_{\text{ML}}$. In [7], we address ML estimation, *i.e.*, the minimization of C_{ML} . There we concluded that $\hat{\mathbf{O}}_{\text{ML}}$ averages the

observations \mathbf{I}_f registered according to the motion \mathbf{q}_f of the object in the region corresponding to the template \mathbf{T} ,

$$\widehat{\mathbf{O}}_{\text{PL}} = \widehat{\mathbf{O}}_{\text{ML}} = \mathbf{T} \frac{1}{F} \sum_{f=1}^F \mathcal{M}(\mathbf{q}_f) \mathbf{I}_f, \quad (4)$$

and $\widehat{\mathbf{B}}_{\text{ML}}$ is the average of the observations \mathbf{I}_f , registered according to the background motion \mathbf{p}_i , in the regions $\{(x, y)\}$ not occluded by the moving object, *i.e.*, when $\mathcal{M}(\mathbf{p}_f \mathbf{q}_f^\#) \mathbf{T}(x, y) = 0$,

$$\widehat{\mathbf{B}}_{\text{PL}} = \widehat{\mathbf{B}}_{\text{ML}} = \frac{\sum_{f=1}^F \left[\mathbf{1} - \mathcal{M}(\mathbf{p}_f \mathbf{q}_f^\#) \mathbf{T} \right] \mathcal{M}(\mathbf{p}_f) \mathbf{I}_f}{\sum_{i=f}^F \left[\mathbf{1} - \mathcal{M}(\mathbf{p}_f \mathbf{q}_f^\#) \mathbf{T} \right] \mathcal{M}(\mathbf{p}_f) \mathbf{H}}. \quad (5)$$

Two-step iterative algorithm If we replace the estimates $\widehat{\mathbf{O}}_{\text{PL}}$, $\widehat{\mathbf{B}}_{\text{PL}}$ given by (4,5) in (3), we get an expression for $C_{\text{PL}}(\mathbf{T})$ in which the minimization wrt each different spatial location $\mathbf{T}(x, y)$ is not independent from the other locations. Solving this binary minimization problem by a conventional method is extremely time consuming. In contrast, if we replace only $\widehat{\mathbf{O}}_{\text{PL}}$, the minimization of $C_{\text{PL}}(\mathbf{B}, \mathbf{T})$ over \mathbf{T} for fixed \mathbf{B} , results in a local binary test. We exploit this fact to derive a computationally simple two-step iterative minimization algorithm: (i) solve for the background \mathbf{B} while the template \mathbf{T} is kept fixed; and (ii) solve for \mathbf{T} while \mathbf{B} is kept fixed. The solution for step (i) is given by (5).

To find the solution for step (ii), we replace $\widehat{\mathbf{O}}_{\text{PL}}$ (4) in (3). By manipulating C_{PL} as we did in [7] for C_{ML} , we obtain

$$C_{\text{PL}} = \iint \mathbf{T}(x, y) [\mathbf{Q}(x, y) + \alpha] dx dy + \text{Constant}, \quad (6)$$

where \mathbf{Q} , which we call the *segmentation matrix* [7], is given by

$$\mathbf{Q} = \frac{1}{F} \sum_{f,g} \left[\mathcal{M}(\mathbf{q}_f) \mathbf{I}_f - \mathcal{M}(\mathbf{q}_g) \mathbf{I}_g \right]^2 - \sum_f \left[\mathcal{M}(\mathbf{q}_f) \mathbf{I}_f - \mathcal{M}(\mathbf{q}_f \mathbf{p}_f^\#) \mathbf{B} \right]^2.$$

We estimate the template by minimizing C_{PL} in (6) over \mathbf{T} , given the background \mathbf{B} . It is clear from (6), that the minimization of C_{PL} wrt each spatial location of \mathbf{T} is independent from the minimization over the other locations. The template $\widehat{\mathbf{T}}_{\text{PL}}$ that minimizes C_{PL} is given by the following test evaluated at each pixel:

$$\begin{aligned} \widehat{\mathbf{T}}_{\text{PL}}(x, y) &= 0 \\ \mathbf{Q}(x, y) &\begin{matrix} > \\ < \end{matrix} -\alpha \\ \widehat{\mathbf{T}}_{\text{PL}}(x, y) &= 1 \end{aligned} \quad (7)$$

Note that by describing the shape of the moving object by the binary template \mathbf{T} , we were able to express the cost function (6) in terms of an integral whose region of integration is independent of the unknown shape. This is what enabled developing a computationally simple algorithm to estimate the shape of the object. The same type of idea has been used in the context of the single-image intensity-based segmentation problem, for example, Ambrosio and Tortorelli [9] adapted Mumford and Shah's theory [10] by using a binary field instead of an edge process. Reference [11] presents a detailed description of the two-step iterative algorithm.

Choosing $\alpha = 0$ leads to ML estimation ($C_{\text{PL}} = C_{\text{ML}}$). In this case, as anticipated above, it may happen that, after processing the F available frames, the test (7) with $\alpha = 0$ remains inconclusive at a given pixel (x, y) , *i.e.*, $\mathbf{Q}(x, y) \simeq 0$. In other words, it is not possible to decide if this pixel belongs to the moving object or to the background. This ambiguity comes from the fact that the available observations are in agreement with both hypothesis. We make the decision unambiguous by looking for the *smallest* template that describes well the observations, through PL estimation.

Relaxation The initial guess in iterative algorithms is very relevant to their convergence—a bad initial guess may lead to a local optimum. Instead of using computationally complex methods to compute an initial guess, we use a continuation method—we relax the cost function. We start from a cost for which we know we can find the global minimum, and then we gradually change the cost, keeping track of the minimum, to end at the desired cost function. Due to the structure of the PL cost (2), the continuation method is easily implemented by relaxing the weight α , as in annealing schedules. We start with a high value for α such that the minimum of (2) occurs at $\mathbf{T}(x, y) = 0, \forall x, y$. Then, we gradually decrease α and minimize the corresponding intermediate costs, till we reach the desired cost and the correct segmentation.

To emphasize the advantage of relaxation, consider using ML as in [7] initialized by estimating the background as the average of the co-registered input images, *i.e.*, the initial estimate of the background is contaminated by the moving object intensity values. It may happen that the next estimate of the template, obtained from (7) with $\alpha = 0$, is, erroneously, so large that, in the next step, the estimate of the background can not be computed at all pixels and the algorithm freezes and can not proceed. Consider now using the same initialization but with a relaxation scheme for α . Using (7) with a large value α , the next estimate of the template will be very small (α can even be set to a value such that the template estimate will contain a single pixel). Using this template estimate, the next estimate of the background will be less contaminated by the moving object intensity values and thus closer to the true background. The next estimate of the template, obtained from (7) with a slightly smaller α , will then be slightly larger and closer to the true template of the moving object. This relaxation proceeds until α reaches either zero, leading to the ML estimate, or a value previously chosen, leading to the PL estimate minimizing (3).

Stopping criteria To stop the relaxation process we could adopt as strategy to stop as soon as the estimate of the template stabilizes, *i.e.*, as soon as no more pixels are added to it. However, to resolve the problems with low contrast background that motivated the use of penalized likelihood estimation, we stop the relaxation when α reaches a pre-specified minimum value α_{MIN} . This α_{MIN} can be chosen by experimentation, but we can actually predict from the model (1) what are good choices for it. If the minimum value α_{MIN} is chosen very high, we risk that some pixel (x, y) of the moving object, *i.e.*, with $\mathbf{T}(x, y) = 1$, is erroneously classified as belonging to the background, since from (7), $\mathbf{Q}(x, y) > -\alpha_{\text{MIN}} \Rightarrow \widehat{\mathbf{T}}_{\text{PL}}(x, y) = 0$. We show elsewhere that the expected value of the entry $\mathbf{Q}(x, y)$ for a pixel (x, y) of the moving object, *i.e.*, with $\mathbf{T}(x, y) = 1$, can be approximated as

$$E \{ \mathbf{Q}(x, y) \} \simeq - \sum_{f=1}^F \left[\mathbf{O}(x, y) - \mathcal{M}(\mathbf{q}_f \mathbf{p}_f^\#) \mathbf{B}(x, y) \right]^2. \quad (8)$$

Then, as we process more frames, $E \{ \mathbf{Q}(x, y) \}$ becomes more negative, reducing the probability of $\mathbf{Q}(x, y) > -\alpha_{\text{MIN}}$, and so of misclassifying the pixel (x, y) as belonging to the background. Good choices for α_{MIN} are then in the interval $]0, -E \{ \mathbf{Q} \} [$. Since in practice we can not compute $E \{ \mathbf{Q} \}$ because we do not know before hand what are the intensity levels of the object and the background, we assume a value S^2 for their average square difference and chose α_{MIN} in the middle of the interval, $]0, FS^2 [$, where F is a constant. With gray-level intensities in $[0, 255]$, we used $\alpha_{\text{MIN}} = 20$, obtained by setting $S = 2$ and $F = 10$. Our experience has shown that any other value α_{MIN} not too close to the extremes of the above interval would lead to equivalent estimates.

4. EXPERIMENTS

Challenging synthetic sequence By rotating and translating the object shown in the left image of Fig. 1, we synthesized 20 frames, two of which are shown in the middle and right images of Fig. 1. As these images clearly show, the noise and the similarity between the textures of the background and the object makes it very challenging to obtain an accurate segmentation. Fig. 2 describes the evolution of the template estimate. The final estimate, shown in the bottom-right image, shows that our algorithm was able to recover the true shape of the moving object (left image of Fig. 1).

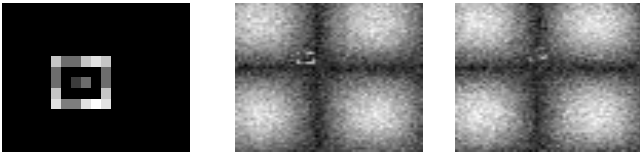


Fig. 1. Left: moving object. Middle and right: noisy video frames.

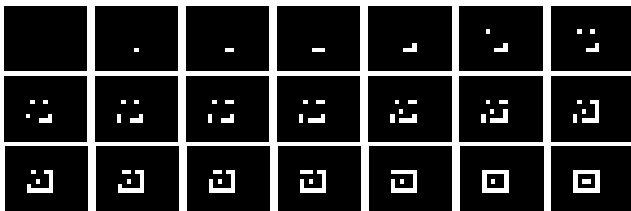


Fig. 2. Relaxation. The final estimate of the template (bottom-right) coincides with the true shape of the moving object in Fig. 1.

Low-light video Fig. 3 shows an illustrative frame of a real-life traffic video sequence taped in the evening, *i.e.*, in a low-light situation. Note that the contrast between the moving car and the road is so small that it is hardly perceived from a single image, even by the human visual system. In Fig. 4, we represent, from left to right, three stages of the evolution of the relaxation algorithm when segmenting this video sequence. The final estimate of the template of the moving car (on the rightmost image of Fig. 4) is visually correct.



Fig. 3. Low-light scenes lead to video sequences with very low contrast between moving objects and background.



Fig. 4. Three stages of the relaxation process for the video sequence illustrated in Fig. 3. The final estimate of the moving object template is on the rightmost image.

5. CONCLUSION

We describe an algorithm for segmenting moving objects from video. Our method models the *rigidity* of the moving object and the *occlusion* of the background. We use relaxation to approximate a penalized likelihood estimate. Experiments show that our algorithm succeeds in recovering complex templates in low-light scenes, *i.e.*, from low contrast videos.

6. REFERENCES

- [1] P. Aguiar, R. Jasinschi, J. Moura, and C. Pluempitwiriyawej, "Content-based image sequence representation," in *Digital Image Sequence Processing, Compression, and Analysis*, Todd Reed, Ed. CRC Press, 2004.
- [2] B. Li and M. Sezan, "Adaptive video background replacement," in *IEEE Int. Conf. on Multimedia and Expo*, Tokio, Japan, 2001.
- [3] J. Pan, C.-W. Lin, C. Gu, and M.-T. Sun, "A robust spatio-temporal video object segmentation scheme with prestored background information," in *IEEE Int. Symp. on Circuits and Systems*, Arizona, USA, 2002.
- [4] T. Chan and L. Vese, "Active contours without edges," *IEEE Trans. on Image Processing*, vol. 10, no. 2, 2001.
- [5] N. Jojic and B. Frey, "Learning flexible sprites in video layers," in *IEEE Int. Conf. on Computer Vision and Pattern Recognition*, Hawaii, 2001.
- [6] M. Irani, B. Rousso, and S. Peleg, "Computing occluding and transparent motions," *Int. Journal of Computer Vision*, vol. 12, no. 1, 1994.
- [7] P. Aguiar and J. Moura, "Detecting and solving template ambiguities in motion segmentation," in *IEEE Int. Conf. on Image Processing*, Santa Barbara CA, USA, 1997.
- [8] P. Green, "Penalized likelihood," in *Encyclopedia of Statistical Sciences*. John Wiley & Sons, New York, 1998.
- [9] L. Ambrosio and V. Tortorelli, "Approximation of functionals depending on jumps by elliptic functionals," *Comm. Pure and Applied Math.*, vol. 43, no. 8, 1990.
- [10] D. Mumford and J. Shah, "Boundary detection by minimizing functionals," in *IEEE Int. Conf. on Computer Vision and Pattern Recognition*, San Francisco, CA, USA, 1985.
- [11] P. Aguiar and J. Moura, "Figure-ground segmentation from occlusion," *IEEE Trans. on Image Processing*, vol. 14, no. 8, 2005.