

OCCLUSION TRACKING USING LOGIC MODELS

James H. von Brecht*
Department of Mathematics
University of California, Los Angeles
520 Portola Plaza
Los Angeles, CA 90095
jub@ucla.edu

Sheshadri R. Thiruvankadam
Department of Mathematics
University of California, Los Angeles
520 Portola Plaza
Los Angeles, CA 90095
sheshad@math.ucla.edu

Tony F. Chan
Department of Mathematics
University of California, Los Angeles
520 Portola Plaza
Los Angeles, CA 90095
chan@math.ucla.edu

ABSTRACT

We present a variational PDE based model for tracking objects under occlusion. Here, prior shape information is used to segment object boundaries that are occluded. The novelty in this work is that the shape prior is combined with the image term using logical operations pertaining to a unique occlusion scenario, thus avoiding locally optimal solutions. The model was tested on real and synthetic image sequences with promising results.

KEY WORDS

Image segmentation, tracking, variational methods

1. Introduction

Occlusion tracking [6, 7, 8, 13] presents a difficult problem since most, if not all, of the information which characterizes a particular object becomes unreliable under occlusions. A successful algorithm for occlusion tracking, then, must have some means of identifying the object of interest when such information is lacking, or inaccurate. Typically, this is accomplished by using information from other frames of a video sequence to aid in the identification of the object in an occluded frame.

In this work, we assume that the boundary of the object of interest is available, and use this *prior shape* information to track the object in occluded frames. There have been previous works [10, 11, 12] that use prior knowledge of the shape of objects to facilitate segmentation specially under low contrasts, occlusions and other undesirable noisy conditions. Most of these works incorporate the shape term additively within the segmentation energy which results in locally optimal solutions. The novelty in this work is that the shape prior is combined with the image term using *logical operations* pertaining to a unique occlusion scenario,

thus leading to “meaningful” solutions. Our work is based on Sandberg et. al. [1] algorithm for logical segmentation of multi channel images, and related to the joint segmentation and registration framework of Moelich et. al. [2].

2. Description of the Model

Our model is a forward tracking algorithm; we rely solely upon data from previous frames in order to identify and segment the object in the current frame. Throughout the remainder of the discussion, we make three assumptions: the first frame of the video sequence, which we refer to as the “template” frame, contains the entire, un-occluded boundary of the object (*shape prior*) we wish to track; second, the object boundary does not deform; and third, the object undergoes only affine movement between frames. The first assumption will allow us to not have to use any *a priori* about the particular occlusion scenario for the algorithm to succeed. The second assumption results from the desire to incorporate prior shape information into the algorithm. We make the last assumption for simplicity, in that more complicated motions amount to a more sophisticated registration model than we use here.

In the following discussion, we present the focus of this paper; how to logically interpret the moving object in the current frame based on the occlusion scenario. This translates into how one could use shape information within a tracking algorithm to avoid the commonly encountered local minima issues.

2.1 Logic Models

Before introducing our model, we briefly describe the region-based logic models [1, 2] based upon the Chan-Vese (C-V) segmentation energy [4]. While dealing with multi-channel images, often there are disparities in the appearance of an object in each of the channels. In such cases, there is more than one valid interpretation of the actual ob-

ject. These alternate interpretations correspond to different logical interpretations of the images. The *logic models*, developed by Sandberg et. al. [1], are designed to segment multi-channel images according to such logical interpretations. For example, given two images f_1 and f_2 in Figure 1(a) which contain two different instances of a particular object of interest, logic models allow us to interpret the actual object (white curve in Figure 1 (b) and (c)) by combining the segmentation in each frame according to a pre-selected logical operation.

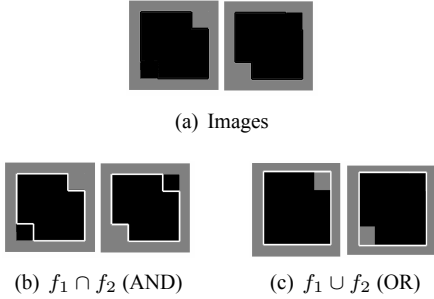


Figure 1. Logical Segmentation on Two Frames

In particular, we consider two logic segmentation models, denoted as $f_1 \cap f_2$ (AND) and $f_1 \cup f_2$ (OR). The AND model interprets the actual object as the intersection of the object regions that appear in the two frames. Similarly, the OR model, then, is the union of the object regions which appear in the frames. As discussed in [1], a segmentation energy for a single channel can be easily recast in to a logical framework when dealing with multi channels. We briefly review the discussion from [1] here.

Given two frames f_1 and f_2 , each which contain an object of interest (that might appear different in each frame), we first define the functions

$$\begin{aligned} z_1^{in} &= \frac{(f_1 - c_1^{in})^2}{M_1} & z_1^{out} &= \frac{(f_1 - c_1^{out})^2}{N_1} \\ z_2^{in} &= \frac{(f_2 - c_2^{in})^2}{M_2} & z_2^{out} &= \frac{(f_2 - c_2^{out})^2}{N_2}. \end{aligned} \quad (1)$$

As in the standard C-V model, c_i^{in} ($i = 1, 2$) represents the average intensity inside the object in frame f_i , respectively. Similarly, c_i^{out} represents the average intensities of the background. The constants M_i and N_i ensure that each z function takes values only between 0 and 1. The segmentation energy takes the familiar form

$$E = \int_{\Omega} f_{in} H(\phi) + \int_{\Omega} f_{out} (1 - H(\phi)) + \int_{\Omega} (|\nabla \phi| - 1)^2, \quad (2)$$

where the last term is a regularization term which also prevents the level-sets of ϕ from becoming too flat. The functions f_{in} and f_{out} vary depending upon which logical combination is desired. $H(t)$ is the heaviside function.

Based on the definitions of the z -functions (1), we

have $z_1^{in} \approx 0/1$ inside/outside the object in f_1 and similarly $z_1^{out} \approx 0/1$ outside/inside the object in f_1 . The case is similar for the functions z_2 . When taking the OR model, we desire $f_{in} = 0$ for all points inside the object in *at least* one of the frames (see Figure 1(c)), and $f_{out} = 0$ for all points that lie outside the object in *both* frames, which we achieve by defining f_{in} and f_{out} as

$$\begin{aligned} f_{in}^{\cup} &= \sqrt{z_1^{in} z_2^{in}} \\ f_{out}^{\cup} &= 1 - \sqrt{(1 - z_1^{out})(1 - z_2^{out})}. \end{aligned} \quad (3)$$

For the AND model (Figure 1(b)), the case is reversed. We desire $f_{in} = 0$ for all points inside the object in both frames, and $f_{out} = 0$ for points outside the object in either frame. The resulting definitions are

$$\begin{aligned} f_{in}^{\cap} &= 1 - \sqrt{(1 - z_1^{in})(1 - z_2^{in})} \\ f_{out}^{\cap} &= \sqrt{z_1^{out} z_2^{out}}. \end{aligned} \quad (4)$$

2.2 Logic Models within Tracking

Now we motivate the use of logic models within a tracking algorithm. In Figure 2, a simple demonstration is shown on how the logic models as discussed above can be used to recover the boundary of an artificially occluded object (the white turning car). The first image (a), or template f_{τ} , is used to segment the other two occluded frames (b and c). In the second image (b), a region of different intensity occludes the object, and hence we take the segmentation model $f_{\tau} \cup f_1$ to recover the boundary of the object (curve shown in (d)). In the image (c), the occlusion is of similar intensity to the car, so $f_{\tau} \cap f_2$ yields the desired result (curve shown in (e)). Note how the appropriate logic model depends upon the intensity of the occlusion *relative* to the object. When no occlusion is present, either logic model will give the desired result, since the object appears identical in each frame. Thus we see that the application of the shape prior (through the template image f_{τ}) depends on the occlusion type, which allows our algorithm to *avoid local minima* problems of models that just additively introduce the shape term.

To summarize, in this paper, we deal with two types of occlusion scenario, depending on the intensity of the occlusion relative (similar/different) to the object being tracked. We use an appropriate logic segmentation model (AND/OR) for the occlusion scenario, to correctly segment the object from the current frame, using the template image. Finally, we discuss a technique which we use to automatically switch between the logic models to deal with changing occlusion scenario across frames.

2.3 Joint Registration and Segmentation

The logic models as presented in [1] assume pre-registered images. Consequently, for use in a tracking algorithm, we

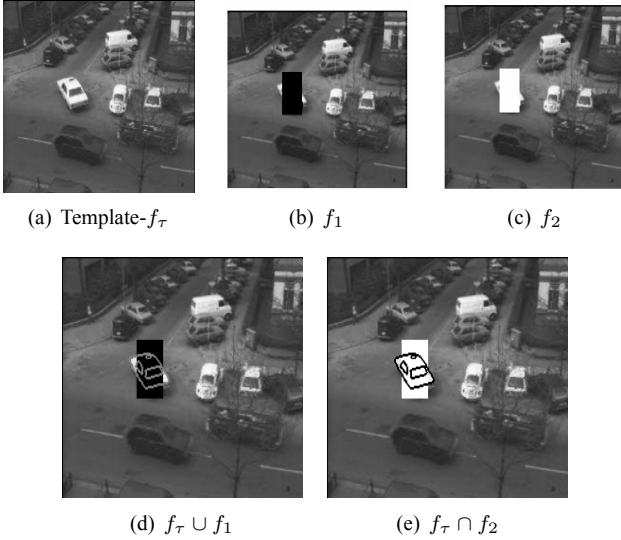


Figure 2. Occlusion Segmentation

must incorporate a registration model into our algorithm to reflect object motion. For our purposes, we employ a Joint Registration and Segmentation algorithm similar to [3] and extended to the logic models by Moelich and Chan in [2]. Given a template frame $f_\tau : \Omega_1 \rightarrow \mathbb{R}^+$ and an occluded frame $f_i : \Omega \rightarrow \mathbb{R}^+$, we register the two images by introducing a spatial correspondence between the domains Ω_1 and Ω , denoted by g , with the parameters $\{\Delta x \ \Delta y \ \theta\}$. Here, Δx and Δy represent translation and θ represents rotation. The selection of the transformation g is arbitrary; many other valid choices exist which allow for more general object motions. For further details, see [3] or [2].

2.4 Automation of Logic Models

As discussed earlier, while segmenting the boundary of an occluded object (e.g. Figure 2), the correct logical model to be used depends upon the similarity of the intensities of the object and the occlusion. The application of the incorrect logic model will lead to not only an error when segmenting the images (see Fig. 3 (d)), but can also conceivably cause an error in the registration of the images as well. Consequently, in order to employ the logic models in a tracking algorithm, we introduce a method by which we can *automatically* determine the appropriate choice of logical segmentation.

To determine the appropriate logic model, we make use of the prior shape of the object, given by the contour C_τ in the template frame f_τ . Also, we denote the contour given by logical segmentation of the current frame by C . Regardless of the intensity of the occluding object, the correct logic model is that which gives the least shape dissimilarity between C_τ and C . In this work, we use area difference as the shape dissimilarity measure. Therefore, the correct logic model is that which minimizes the quan-

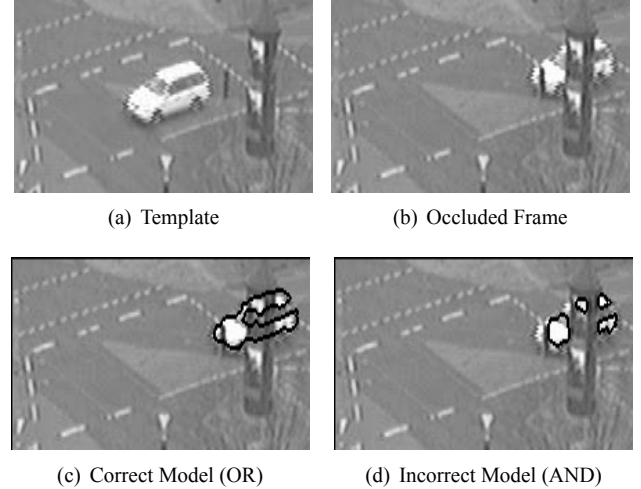


Figure 3. Need for Automation

tity

$$(AREA(inside(C_\tau)) - AREA(inside(C)))^2. \quad (5)$$

In the level-set framework, we denote by ψ the function used to implicitly represent the contour C_τ , and ϕ represents C . The variational form of (5) then becomes

$$\int_{\Omega} (H(\psi_g) - H(\phi))^2, \quad (6)$$

where $\psi_g = \psi(g^{-1})$, and H is the Heaviside function. To enforce this constraint in practice, we compute two functions, ϕ^\cap using the AND model, and ϕ^\cup using the OR model, check the quantity (6) in each case, and select as ϕ that which produces a minimum.

2.5 Variational Framework

We now describe the level-set formulation of our model. Given a template frame f_τ and the i^{th} frame from a video sequence f_i , define $F_\tau = f_\tau(g^{-1})$. Extending the logic models to include a registration component is then straightforward. The z -functions (1) become

$$z_\tau^{in} = \frac{(F_\tau - c_\tau^{in})^2}{M_\tau} \quad z_\tau^{out} = \frac{(F_\tau - c_\tau^{out})^2}{N_\tau} \quad (7)$$

$$z_i^{in} = \frac{(f_i - c_i^{in})^2}{M_i} \quad z_i^{out} = \frac{(f_i - c_i^{out})^2}{N_i}.$$

Note that in practice, to target the correct object in the i^{th} frame, we fix $c_i^{in} = c_\tau^{in}$ and only update c_i^{out} along with the contour. The functions $f_{in,out}$ then become

$$f_{in}^\cap = 1 - \sqrt{(1 - z_\tau^{in})(1 - z_i^{in})}$$

$$f_{out}^\cap = \sqrt{z_\tau^{out} z_i^{out}}. \quad (8)$$

$$f_{in}^{\cup} = \sqrt{z_{\tau}^{in} z_i^{in}}$$

$$f_{out}^{\cup} = 1 - \sqrt{(1 - z_{\tau}^{out})(1 - z_i^{out})}. \quad (9)$$

And finally in our formulation, for each of the logic models (AND/OR), we add the shape term (5) to the segmentation energy. The addition of the shape term has several benefits. Foremost, it helps to ensure a correct registration between frames. Also, it helps prevent unwanted portions (usually similar to the object's intensity) of the image from being included in the final segmentation. Thus, we may introduce the variational form of our model,

$$E(\phi, \Delta x, \Delta y, \theta) = \int_{\Omega} f_{in} H(\phi) + f_{out} (1 - H(\phi)) dx$$

$$+ \beta \int_{\Omega} (H(\psi_g) - H(\phi))^2 dx + \lambda \int_{\Omega} (|\nabla \phi| - 1)^2 dx. \quad (10)$$

Again, H is the Heaviside function, and ψ represents the boundary of the object in the (unregistered) template frame. The function ψ_g then, is defined as $\psi_g = \psi(g^{-1})$. λ and β are parameters to balance the terms.

In this context, our algorithm reduces to finding sequentially, for each frame f_i of the sequence, the function ϕ_i and the parameters $p_i = \{\Delta x_i \Delta y_i \theta_i\}$ which minimize the energy (10). Since the functions $f_{in,out}$ differ depending upon the logic model, we minimize the energy separately for each case, to produce two sets of functions with coupled parameters, $\{\phi_i^{\cap}, p_i^{\cap}\}$ and $\{\phi_i^{\cup}, p_i^{\cup}\}$, then select as $\{\phi_i, p_i\}$ that which minimizes (6). This selection gives the desired segmentation; the segmentation closest to the shape prior.

In summary, we combine both shape information and the correct logical interpretation of images to achieve the desired result. We can thus avoid many local minima that other models, which just additively introduce shape, may encounter. In Fig. 4, the first frame demonstrates the result of our algorithm, which combines prior shape and the automated choice of logic model to achieve the desired segmentation. The final frame demonstrates that prior shape alone is not sufficient, and might result in a local minimum.

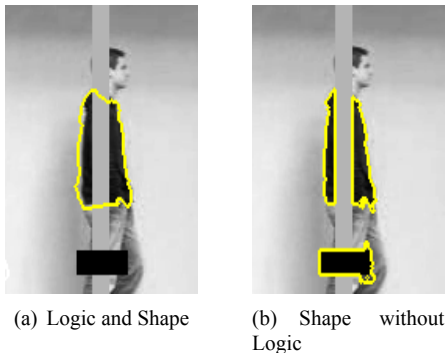


Figure 4. Various Combinations of Logic Models and Prior Shape

3. Numerical Implementation

When implementing the algorithm, we begin with an initial ϕ_0 and an initial set of parameters $\{\Delta x_0 \Delta y_0 \theta_0\}$, and evolve them according to the Euler-Lagrange equations of (10) until a minimum is reached. The equations for gradient descent are given by

$$\frac{\partial \phi}{\partial t} = [f_{out} - f_{in} + 2\beta (H(\psi_g) - H(\phi))] \delta(\phi)$$

$$+ 2\lambda \left(\nabla^2 \phi - \text{div} \left(\frac{\nabla \phi}{|\nabla \phi|} \right) \right) \quad (11)$$

and

$$\frac{\partial \Delta x}{\partial t} = - \frac{\partial E_{in}(\phi, \Delta x, \Delta y, \theta)}{\partial \Delta x}$$

$$\frac{\partial \Delta y}{\partial t} = - \frac{\partial E_{in}(\phi, \Delta x, \Delta y, \theta)}{\partial \Delta y}$$

$$\frac{\partial \theta}{\partial t} = - \frac{\partial E_{in}(\phi, \Delta x, \Delta y, \theta)}{\partial \theta} \quad (12)$$

To minimize (10), we follow the procedure outlined in [2, 3]. Beginning from the initial set $\{\Delta x_0 \Delta y_0 \theta_0 \phi_0\}$, we first hold ϕ fixed and perform one iteration of (12), and update the functions F_{τ} and ψ_g with the new parameters. The registration parameters are then fixed and ϕ is updated via one iteration of (11). This process is continued until convergence.

In practice, selecting the initial set of registration parameters for the i^{th} frame presents the largest barrier for our algorithm to track successfully. That is, how to select $p_0^i = \{\Delta x_0^i \Delta y_0^i \theta_0^i\}$. In simpler cases, it suffices simply to take $p_0^i = p_{final}^{i-1}$. However, when the sequence has moderate to severe occlusions or low contrast (see Fig. 7), the algorithm becomes more sensitive to local minima, and hence we must select the parameters more carefully. From our experience, we have developed several techniques to combat this sensitivity. Based upon the definitions of f_{in}^{\cup} and f_{in}^{\cap} , we see that the OR model allows for more possible registration local minima than the AND model, and hence the AND model is less sensitive to the prediction of p_0^i . Consequently we can first compute ϕ^{\cap} and use those final parameters as the initial parameters for the computation of ϕ^{\cup} . Alternatively, we might hypothesize a motion trajectory for the tracked object (in the case of Fig. 7-a linear trajectory), and use this to generate each p_0^i . Finally, to reduce the computational burden we run the algorithm only locally around the object of interest in each frame. That is, we simply crop from each frame a small region around the object of interest (this process is automated via use of the shape prior and the prediction of p_0 for each frame) and run the algorithm on the reduced frame. Once we have the desired contour, we bring the result back to the original image. Due to such difficulties, and to try and develop a more robust and computationally efficient algorithm, we plan to incorporate the described logical framework into a particle filtering algorithm, such as that described in [13].

4. Experimental Results

We now give results of our algorithm on both synthetic and real examples. In all cases, the first frame f_1 was used as the template, thus the sequences begin with frame f_2 . The first example (Figure 5) demonstrates the need to automate the choice of logic model via (6). Without some means of determining the appropriate logical interpretation of the images, an undesirable segmentation can result. In this example, the arbitrary choice was made to use the OR model across all frames. While the algorithm tracks successfully through the first occlusion, when the person reaches the second occlusion, the OR model is no longer correct, and so the algorithm fails. In the later frames of the sequence, when the person has passed completely through the second occlusion, the correct segmentation is once again realized since in such frames, when no occlusion is present, the object appears identical in each channel and hence either logic model gives the desired result. However, in more severe cases, since the registration prediction also depends upon the accuracy of the final segmentation in the previous frame, if the incorrect model is used, the algorithm can completely lose track of the object.

The second example (Figure 6) shows the full algorithm on the same sequence as in (Figure 5). It demonstrates the capability of the algorithm to handle occlusions of both types when the automation method is utilized. The intersection model was taken automatically as the person passes through the black-line occlusion, and union automatically through the second, gray-line occlusion. The current model is unable to cope with the intermediate case, in which the object of interest is simultaneously occluded by regions of both similar intensity and different intensity to the object itself. Such a scenario requires a combination AND/OR model, and we are currently experimenting with a multi-phase level-set method to handle this final case.

The final example, (Figure 7) demonstrates the algorithm on a real video sequence, and was the most challenging. We employed the full algorithm as described, which selected only the OR model across each frame. As the sequence progresses, poor image contrast and more severe occlusions make the tracking more difficult, but with a careful choice of the target intensity c_τ^{in} and a careful prediction of the initial parameters at each step, our algorithm succeeded.

Acknowledgement

This research supported by ONR grant N00014-06-1-0345, NSF grants DMS-0601395, DMS-0610079, and ARO MURI grant 50363-MA-MUR.

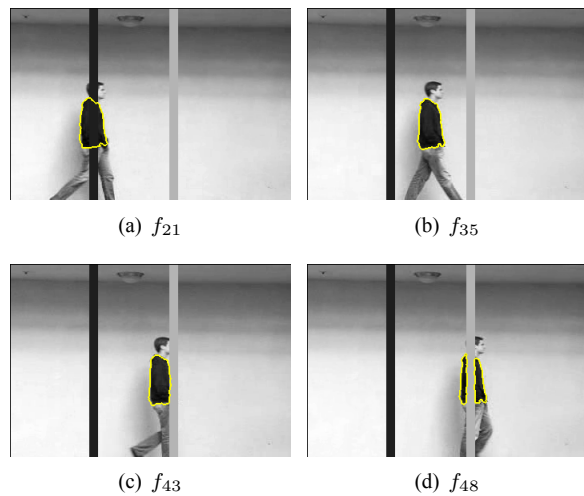


Figure 5. Result without using (6)

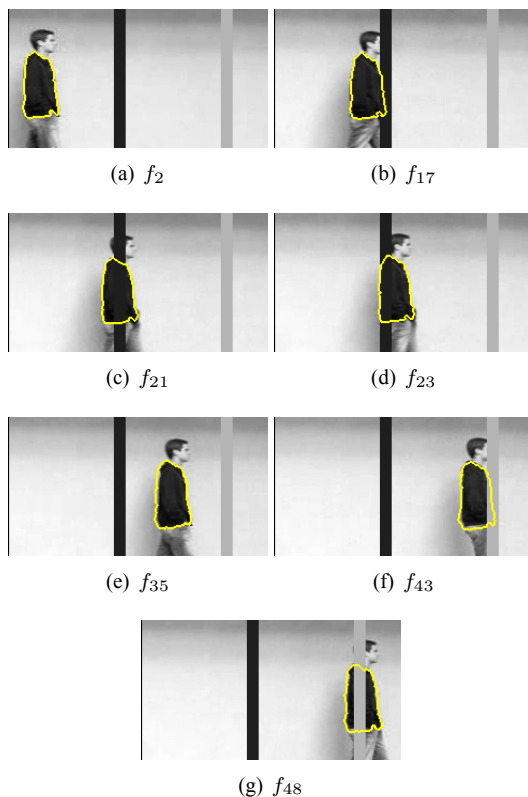
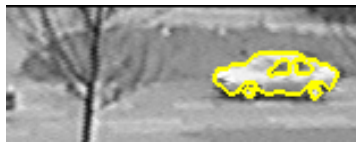


Figure 6. Tracking through both types of occlusions

References

- [1] B. Sandberg and T. Chan, Logic Operators for Active Contours on Multi-Channel Images, *UCLA CAM Report 02-12*, 2002.
- [2] M. Moelich and T. Chan, Joint Segmentation and Registration Using Logic Models, *UCLA CAM Report 03-06*, 2003.
- [3] A. Yezzi, L. Zollei, and T. Kapur, A variational approach to joint segmentation and registration, *Proc.*



(a) f_2



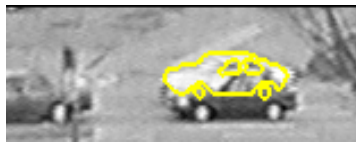
(b) f_{30}



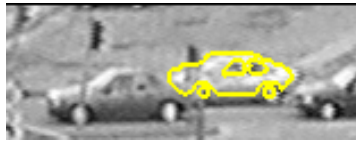
(c) f_{40}



(d) f_{49}



(e) f_{55}



(f) f_{76}



(g) f_{83}



(h) f_{90}

Figure 7. Real sequence with poor image contrast

IEEE Conf. on Comp. Vision and Pattern Recogn., 2001.

- [4] T. Chan and L. Vese, Active Contours Without Edges, *IEEE Transactions on Image Processing*, vol. 10, no. 8, pp. 266-277, 2001.
- [5] T. Chan, B. Sandberg, and L. Vese, Active contours without edges for vector-valued images, *Journal of Visual Communication and Image Representation*, 11:130-141, 1999.
- [6] Y. Shi, J. Konrad, and W. Karl, Multiple Motion and Occlusion Segmentation with a Multiphase Level Set Method, 2004.
- [7] A. Bartesaghi and G. Sapiro, Tracking of Moving Objects under Severe and Total Occlusions, *IMA Preprint Series 2015*, 2005.
- [8] J. Jackson, A. Yezzi, and S. Soatto, Tracking Deformable Moving Objects Under Severe Occlusions, *IEEE Conf. on Decision and Control*, 2004.
- [9] M. Moelich and T. Chan, Tracking Objects with the Chan-Vese Algorithm, *UCLA CAM Report 03-14*, 2003.
- [10] Y. Chen, H. Tagare, S. Thiruvenkadam, F. Huang, D. C. Wilson, K.S. Gopinath, R W. Briggs, and E. A. Geiser, Using prior shapes in geometric active contours in a variational framework, *IJCV*, 50(3):315-328, 2002.
- [11] Timothy F. Cootes, Christopher J. Taylor, David H. Cooper, and Jim Graham, Active shape models-their training and application, *CVPR*, 61(1):38-59, 1995.
- [12] M. Leventon, W. L. Grimson, and O. Faugeras, Statistical shape influence in geodesic active contours, *CVPR*, volume 1, pages 316-323, 2000.
- [13] Y. Rathi, N. Vaswani, A. Tannenbaum, and A. Yezzi, Particle Filtering for Geometric Active Contours with Application to Tracking Moving and Deforming Objects.