

COMBINING AN ACTIVE SHAPE AND MOTION MODELS FOR OBJECT SEGMENTATION IN IMAGE SEQUENCES

Carlos Santiago, Jacinto C. Nascimento, Jorge S. Marques

Institute for Systems and Robotics (ISR/IST), LARSyS, Instituto Superior Técnico,
Universidade Lisboa, Portugal

ABSTRACT

Obtaining the segmentation of an object in a sequence of images is usually achieved using a tracking methodology. However, in some applications, the whole sequence is available beforehand. This means that the segmentations can be determined simultaneously for all the frames in the sequence and taking into account the motion of the object. This paper proposes a new framework to incorporate motion information in the segmentation of image sequences using an active shape model (ASM). The motion of the object is modeled using a vector field, which is learned and refined online as the segmentation algorithm proceeds. The vector field is determined from the trajectories described by ASM points throughout the sequence. The vector field, in turn, influences the estimation of the ASM parameters by acting as a regularizer, ensuring that the segmentations are in agreement with the expected motion. The results show that coupling these models during the segmentation leads to an increase in performance, in particular by guaranteeing more consistent segmentations and by avoiding gross errors in more challenging frames.

Index Terms— Segmentation, Active Shape Model, Motion Model, Vector Field

1. INTRODUCTION

Taking into account the motion of an object is often beneficial to achieve better results in its segmentation in sequences of images. The traditional approach is to adopt a tracking methodology, *e.g.*, the Kalman filter, in which the segmentation in a specific frame is influenced by the segmentation obtained in the previous frame, through the dynamical model. In some applications, however, the sequence is fully available beforehand. This means that the segmentations do not have to be determined sequentially, but can be estimated simultaneously. The advantage is that this allows each frame to influence (and be influenced by) all the remaining frames, and not just the following frame. This has the potential to im-

prove the results by making the segmentations more robust in frames where the image conditions are more challenging.

An example where this approach is useful is the segmentation of cardiac magnetic resonance (CMR) data. A CMR sequence typically covers one cardiac cycle and the segmentation of the frames in which the heart is contracted is harder because boundaries are not so clearly defined [1]. Therefore, motion information allows the segmentation of these frames to be less dependent on image information and more reliant on the segmentations of the remaining frames.

In this work, we propose a framework to combine an active shape model (ASM) [2] with a motion model to determine the object position in the whole sequence simultaneously. This motion model acts as a regularizer in the estimation of the ASM parameters, ensuring that the segmentations are in agreement with the expected dynamics. Here, the motion model is described by a vector field that is learned and refined online, as the segmentation algorithm iterates, by analyzing the trajectories of the model points throughout the sequence. This information is then used to update the segmentations in an alternating scheme.

2. RELATED WORK

Adding motion information to an ASM-based methodology has been accomplished using two approaches: i) *embedding* this information within the deformation modes of a high-dimensional shape model [3, 4, 5, 6]; and ii) *combining* the ASM with a specific motion model [7, 8, 9, 10, 11, 12].

In the first approach (*embedding*), the traditional ASM is extended to jointly model the position of an object in an image sequence (instead of a single image). Similarly to the principal component analysis (PCA) performed to extract the deformation modes, in the extended version, the same analysis is used to extract not only local shape variations but also the variation along time. The downside of this approach is that it requires a large amount of training data, a problem that is typically called the *curse of dimensionality* [6, 13]. To overcome this limitation, most works resort to hierarchical ASMs [3, 5, 6, 9]. This variation of the traditional ASM divides the model into several patches and learns the shape statistics independently for each patch. This leads to a significant reduction

This work was supported by the FCT PhD grant [SFRH/BD/87347/2012], and FCT project and plurianual funding: [PTDC/EEIPRO/0426/2014], [UID/EEA/50009/2013].

in the data dimensionality. However, splitting the model into patches induces a loss of notion about the connectivity between patches, which may cause unexpected segmentations to be obtained. Furthermore, the model is also unable to capture the patterns of variation along time that the patches might share.

Regarding the second type of approach (*combining*), several motion models have been proposed to represent the object dynamics. A popular approach is to model the motion between consecutive frames using optical flow or registrations techniques [10, 11, 12]. The above methods provide a vector field between consecutive frames based on the comparison of intensity images (or other features derived from them). This type of approach is similar to the one proposed in this work. However, instead of using image features to determine the vector fields, we rely on the trajectories described by the estimated position of the ASM points along time. The following section describes the proposed methodology in more detail.

3. PROPOSED APPROACH

The proposed framework consists of alternating between two main steps: 1) **Estimating the motion model** based on the estimated ASM position in all frames; 2) **Updating the shape model parameters** based on observation points extracted from the image and according to the motion model estimated in 1). These steps are described next.

3.1. Estimating the motion model

We represent the motion model as a vector field (VF), defined on a regular grid within the image domain, where each vector in the grid determines the motion in a specific position of the image [14]. Let $\mathbf{V} = [\mathbf{v}_1^\top, \dots, \mathbf{v}_N^\top]^\top \in \mathbb{R}^{2N}$ denote the collection of the all the motion vectors in the grid with N nodes, where $\mathbf{v}_n \in \mathbb{R}^2$ is the vector associated with the n -th node. We define $V : \mathbb{R}^2 \mapsto \mathbb{R}^2$ as a function that maps a position in the image to the corresponding motion vector. For a generic position $\mathbf{x} \in \mathbb{R}^2$ within the image domain, the corresponding motion vector, $V(\mathbf{x}) \in \mathbb{R}^2$, is given by

$$\begin{aligned} V(\mathbf{x}) &= \sum_{n=1}^N \phi_n(\mathbf{x}) \mathbf{v}_n \\ &= \Phi(\mathbf{x}) \mathbf{V}, \end{aligned} \quad (1)$$

where $\Phi(\mathbf{x}) \in \mathbb{R}^{2 \times 2N}$ is a sparse matrix that determines the contribution of each node in the grid to the computation of $V(\mathbf{x})$, such that

$$\Phi(\mathbf{x}) = \begin{bmatrix} \phi_1(\mathbf{x}) & 0 & \dots & \phi_N(\mathbf{x}) & 0 \\ 0 & \phi_1(\mathbf{x}) & \dots & 0 & \phi_N(\mathbf{x}) \end{bmatrix}. \quad (2)$$

We adopt a bilinear interpolation scheme, where only the four closest grid nodes contribute to the computation of

the motion vector, *i.e.*, at most, only four elements of $\{\phi_1(\mathbf{x}), \dots, \phi_N(\mathbf{x})\}$ are non-zero, and they satisfy the constraints $0 \leq \phi_n(\mathbf{x}) \leq 1$ and $\sum_{n=1}^N \phi_n(\mathbf{x}) = 1$.

Following [14], the VF, \mathbf{V} , is estimated from the trajectories described by each model point in the sequence. Let us assume we are given a set of K independent trajectories, such that the k -th trajectory is given by $\mathcal{X}^k = \{\mathbf{x}^k(1), \dots, \mathbf{x}^k(T)\}$, where $\mathbf{x}(f) \in \mathbb{R}^2$ denotes the position of a point at a specific frame $f \in \{1, \dots, T\}$.¹ The VF describes the trajectory of a point through the following dynamical model

$$\mathbf{x}(f) = \mathbf{x}(f-1) + V(\mathbf{x}(f-1)) + \mathbf{w} \quad (3)$$

where $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \sigma \mathbf{I})$ is white noise with Gaussian distribution. Replacing (1) into (3) we obtain

$$\mathbf{x}(f) = \mathbf{x}(f-1) + \Phi(\mathbf{x}(f-1)) \mathbf{V} + \mathbf{w}, \quad (4)$$

The maximum posterior estimate of the VF is given by (see details in [14])

$$\begin{aligned} \mathbf{V}^* &= \arg \min_{\mathbf{V}} \frac{1}{\sigma^2} \sum_{k=1}^K \sum_{f=2}^T \|\Delta \mathbf{x}^k(f) - \Phi(\mathbf{x}^k(f-1)) \mathbf{V}\|^2 + \\ &+ \alpha \mathbf{V}^\top \Lambda^{-1} \mathbf{V}, \end{aligned} \quad (5)$$

where $\Delta \mathbf{x}^k(f) = \mathbf{x}^k(f) - \mathbf{x}^k(f-1)$, and the second term, $\alpha \mathbf{V}^\top \Lambda^{-1} \mathbf{V}$, is related to a Gaussian prior. This prior acts as a smoothness regularizer by penalizing large differences between the vectors of neighboring nodes, encoded in Λ . The parameter α determines the strength of the prior.

The solution of (5), \mathbf{V}^* , is obtained by computing the derivative of the objective function with respect to \mathbf{V} and equating to zero, which leads to the following linear equation

$$\begin{aligned} \left(\alpha \Lambda^{-1} + \frac{1}{\sigma^2} \sum_{k=1}^K \sum_{f=2}^T \Phi^\top(\mathbf{x}^k(f-1)) \Phi(\mathbf{x}^k(f-1)) \right) \mathbf{V}^* &= \\ = \frac{1}{\sigma^2} \sum_{k=1}^K \sum_{f=2}^T \Phi^\top(\mathbf{x}^k(f-1)) \Delta \mathbf{x}^k(f). \end{aligned} \quad (6)$$

3.2. Updating the shape model parameters

The ASM models the position of the contour point by

$$\mathbf{x}^i(f) = \mathbf{A}(f)(\bar{\mathbf{x}}^i + \mathbf{D}^i \mathbf{b}(f)) + \mathbf{t}(f), \quad (7)$$

such that $\bar{\mathbf{x}}^i$ and \mathbf{D}^i are the mean position and the deformation modes, respectively, learned from a training set. The segmentation of the object in a frame f is determined by the set of ASM parameters $\Theta(f) = \{\mathbf{A}(f), \mathbf{t}(f), \mathbf{b}(f)\}$. The estimation of these parameters is accomplished using an extension of

¹Without loss of generality, in the following equations the trajectories are all assumed to be sampled at the same instants and the same number of times.

the EM-RASM [15] (expectation-maximization robust ASM) formulation for sequences of images. Given a set of observation points, $\{\mathbf{y}^i(f)\}_{i=1,\dots,M}, f = 1, \dots, T$, extracted from the image sequence, the shape model parameters are updated iteratively with the EM algorithm. In the E-step, the weight of each observation point is computed as follows

$$w^i(f) \propto \mathcal{N}(\mathbf{y}^i(f); \hat{\mathbf{x}}^i(f), \Sigma^i), \quad (8)$$

where $\hat{\mathbf{x}}^i(f)$ depends on the ASM parameters estimate in the current iteration t , $\hat{\Theta}_{(t)}(f)$. These weights correspond to the likelihood of that observation belonging to the target object. In the M-step, the ASM parameters in all the frame are updated by minimizing a weighted least squares fit between the observations and the corresponding model points,

$$\{\hat{\Theta}(1), \dots, \hat{\Theta}(T)\}_{(t+1)} = \arg \min_{\Theta(1), \dots, \Theta(T)} \sum_{f=1}^T \sum_{i=1}^M w^i(f) \|\mathbf{x}^i(f) - \mathbf{y}^i(f)\|_{\Sigma^i}^2, \quad (9)$$

where $\|\mathbf{v}\|_{\Sigma}^2 = \mathbf{v}^\top \Sigma^{-1} \mathbf{v}$, Σ^i is a diagonal covariance matrix associated with $\mathbf{x}^i(f)$, and the position of $\mathbf{x}^i(f)$

In order to combine this estimation with the motion model, an additional VF term is included. This new term penalizes large deviations between each point, $\mathbf{x}^i(f)$, and its expected position according to the VF, given by

$$\mathbf{y}_V^i(f) = \mathbf{x}^i(f-1) + V(\mathbf{x}^i(f-1)), \quad (10)$$

where $V(\mathbf{x}^i(f-1))$ is computed by (1). Formally, we add the following new term,

$$\sum_{f=2}^T \sum_{i=1}^M \|\mathbf{x}^i(f) - \mathbf{y}_V^i(f)\|_{\Sigma^i}^2, \quad (11)$$

to the objective function in (9), leading to

$$\{\hat{\Theta}(1), \dots, \hat{\Theta}(T)\}_{(t+1)} = \arg \min_{\Theta(1), \dots, \Theta(T)} \sum_{f=1}^T \sum_{i=1}^M w^i(f) \|\mathbf{x}^i(f) - \mathbf{y}^i(f)\|_{\Sigma^i}^2 + \lambda_V \|\mathbf{x}^i(f) - \mathbf{y}_V^i(f)\|_{\Sigma^i}^2, \quad (12)$$

where λ_V is a constant that determines the importance of the VF term, and, for the first frame, $\mathbf{y}_V^i(1) = \mathbf{x}^i(1)$ (*i.e.*, the motion model does not influence the first frame).

We impose that $\mathbf{y}_V^i(f-1)$, computed using (10), is fixed given the current ASM parameters. Therefore, it does not depend on the $\hat{\Theta}(f-1)$. In practice, this can be seen as considering new observation points given by the expected motion, which are all weighted by λ_V . This will introduce a bias in the estimation of the ASM parameters that will make the segmentations combine the information from the images with the expected motion.

The solution of (12) is approximated by first minimizing with respect to the transformation parameters, $\mathbf{a}(f), \mathbf{t}(f)$, $f = 1, \dots, T$, and then minimizing with respect to the deformation coefficients, $\mathbf{b}(f)$, $f = 1, \dots, T$ (see details in [15]).

4. EXPERIMENTAL SETUP

The proposed approach is evaluated on two problems: 1) the segmentation of the left ventricle in CMR; and 2) the segmentation of lips in face images. For the first problem, we use the CMR sequence dataset [6], which comprises 33 sequences of volumes of healthy and disease patients. Each sequence contains 20 volumes, covering the systole and diastole phases of the cardiac cycle. Each volume contains 5-10 slices, with a spacing of 6-13 mm. Each slice is a 256×256 image, with a resolution of 0.93-1.64 mm, with a total of 5011 images. For the second problem, we use 61 sequences of 490×640 face images from the ‘‘happy’’ expression of the Cohn-Kanade expression database [16], each with 10-45 frames, for a total of 1241 images. The ground truth (GT) segmentations of both datasets is also provided.

The ASM (mean shape and modes of deformation) is learned using a leave-one-sequence-out strategy, *i.e.*, to test on a specific sequence, the model is learned using the remaining sequences. The quantitative evaluation of the segmentations is performed using two metrics: (*i*) the Dice coefficient, d_{Dice} , and (*ii*) the average perpendicular distance, d_{AV} .

In order to evaluate the advantage of the proposed approach, we compare our results with the segmentations obtained without the information from the motion model, *i.e.*, each frame is analyzed independently and without any notion of temporal dependency.

5. RESULTS

5.1. LV segmentation in CMR sequences

In the case of the LV segmentation, two VFs have to be learned: one for the contraction phase and one for the dilation phase. To separate these two phases, the trajectories of the model points along the sequence are divided in two parts. The first set of trajectories, associated to the frames $f = 1, \dots, f_s$, are used to compute the VF of the contraction phase, whereas the trajectories from frames $f = f_s, \dots, T$ are used to compute the VF of the dilation phase. The switching frame f_s is determined as the frame in which the LV area was the smallest (*i.e.*, corresponding to the end of the contraction phase). An example of the VFs obtained is shown in Fig. 1.

Table 1 shows the average accuracy across the entire dataset with both metrics. It is possible to see that there is an increase in the performance when using information from the motion model. In particular, this information makes the segmentation algorithm give more consistent results across time and avoid many of the gross errors obtained when analyzing

Table 1. Mean (std) accuracy on LV segmentation in CMR sequences from [6].

Motion Model	None	VF
d_{Dice}	84.2 (9.2)	85.6 (7.4)
d_{AV}	2.5 (1.5)	2.3 (1.2)

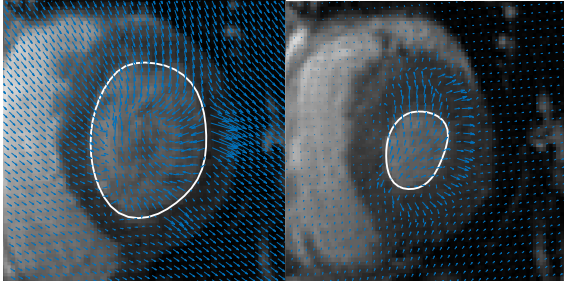


Fig. 1. Example of the VFs obtained to represent the motion of the LV during contraction (left) and dilation (right).

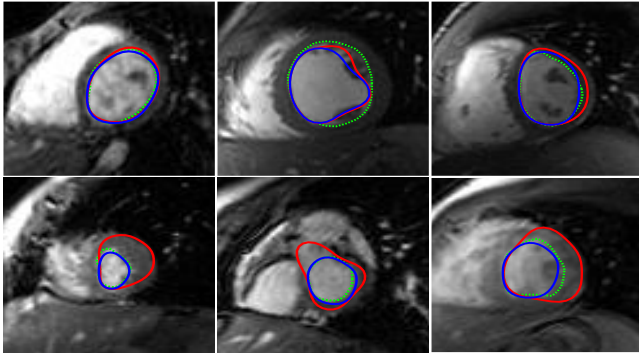


Fig. 2. Examples of LV segmentations with (blue) and without (red) the motion model and comparison with the GT (green).

each frame independently. Examples of the segmentations obtained with and without the motion model are shown in Fig. 2. In the examples shown in the top row, both methods have a similar performance. However, in the bottom row we show examples where there is a clear advantage in using the motion model.

5.2. Lip segmentation in face image sequences

In the lip segmentation problem, we observed the same advantages discussed above. The accuracy of the segmentation is increased when the motion model is included, leading to the overall improvement shown in Table 2. An example of the VF obtained in a particular sequence is shown in Fig. 3, where it is possible to see that it represents the motion of the lips when a person starts smiling. Fig. 4 shows a comparison between the segmentation obtained with and without the motion model. As previously, the top row shows examples in

Table 2. Mean (std) accuracy on lip segmentation in face image sequences from [16].

Motion Model	None	VF
d_{Dice}	85.1 (6.9)	86.2 (5.6)
d_{AV}	3.1 (1.4)	2.9 (1.1)

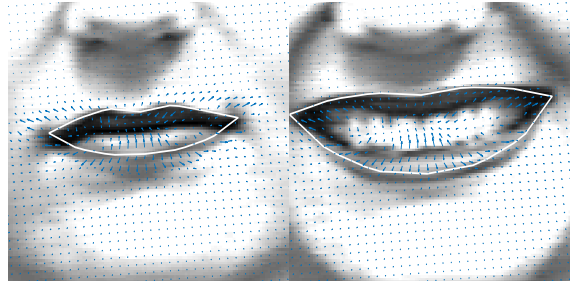


Fig. 3. Example of the VF obtained to represent the motion of the lips.



Fig. 4. Examples of lip segmentations with (blue) and without (red) the motion model and comparison with the GT (green).

which the performance is similar in both approaches, but the bottom row shows examples where it was beneficial to use the VF constraint.

6. CONCLUSIONS

This paper proposes a framework to incorporate motion information for the segmentation of sequences using an ASM-based model. In this new framework, the segmentation of each frame in the sequence occurs simultaneously. The motion of the object is modeled using a VF, which is learned and refined online, as the segmentation algorithm iterates. The VF acts as a regularizer in the estimation of the ASM parameters, ensuring that the segmentations obtained are in agreement with the expected motion of the object. The results show that there is a clear advantage in using this information during the segmentation, in particular to guarantee consistency between segmentations in different frames and by avoid gross errors in more challenging frames.

7. REFERENCES

- [1] A. Suinesiaputra, B. R. Cowan, A. O. Al-Agamy, M. A. Elattar, N. Ayache, A. S. Fahmy, A. M. Khalifa, P. Medrano-Gracia, M. Jolly, A. H. Kadish, et al., “A collaborative resource to build consensus for automated left ventricular segmentation of cardiac MR images,” *Medical image analysis*, vol. 18, no. 1, pp. 50–62, 2014.
- [2] T. F. Cootes, C. J. Taylor, D. H. Cooper, and J. Graham, “Active shape models-their training and application,” *Computer vision and image understanding*, vol. 61, no. 1, pp. 38–59, 1995.
- [3] C. Davatzikos, X. Tao, and D. Shen, “Hierarchical active shape models, using the wavelet transform,” *IEEE transactions on medical imaging*, vol. 22, no. 3, pp. 414–423, 2003.
- [4] G. Hamarneh and T. Gustavsson, “Deformable spatio-temporal shape models: extending active shape models to 2D+ time,” *Image and Vision Computing*, vol. 22, no. 6, pp. 461–470, 2004.
- [5] J. Lötjönen, K. Antila, E. Lamminmäki, J. Koikkalainen, M. Lilja, and T. Cootes, “Artificial enlargement of a training set for statistical shape models: Application to cardiac images,” in *International Workshop on Functional Imaging and Modeling of the Heart*. Springer, 2005, pp. 92–101.
- [6] A. Andreopoulos and J. K. Tsotsos, “Efficient and generalizable statistical models of shape and appearance for analysis of cardiac MRI,” *Medical Image Analysis*, vol. 12, no. 3, pp. 335–357, 2008.
- [7] F. Billet, M. Sermesant, H. Delingette, and N. Ayache, “Cardiac motion recovery and boundary conditions estimation by coupling an electromechanical model and cine-MRI data,” in *Functional Imaging and Modeling of the Heart*, pp. 376–385. Springer, 2009.
- [8] C. Casta, P. Clarysse, J. Schaerer, and J. Pousin, “Evaluation of the dynamic deformable elastic template model for the segmentation of the heart in MRI sequences,” *MIDAS J-Card MR Left Ventricle Segmentation Challenge*, 2009.
- [9] S-W Lee, J. Kang, J. Shin, and J. Paik, “Hierarchical active shape model with motion prediction for real-time tracking of non-rigid objects,” *IET Computer Vision*, vol. 1, no. 1, pp. 17–24, 2007.
- [10] L. A. Schwarz, A. Mkhitarian, D. Mateus, and N. Navab, “Human skeleton tracking from depth data using geodesic distances and optical flow,” *Image and Vision Computing*, vol. 30, no. 3, pp. 217–226, 2012.
- [11] O. Arif, G. Sundaramoorthi, B. Hong, and A. Yezzi, “Tracking using motion estimation with physically motivated inter-region constraints,” *IEEE transactions on medical imaging*, vol. 33, no. 9, pp. 1875–1889, 2014.
- [12] L. Wang, A. Basarab, P. R. Girard, P. Croisille, P. Clarysse, and P. Delachartre, “Analytic signal phase-based myocardial motion estimation in tagged mri sequences by a bilinear model and motion compensation,” *Medical image analysis*, vol. 24, no. 1, pp. 149–162, 2015.
- [13] S. P. O’Brien, O. Ghita, and P. F. Whelan, “A Novel Model-Based 3D Time Left Ventricular Segmentation Technique,” *Medical Imaging, IEEE Transactions on*, vol. 30, no. 2, pp. 461–474, 2011.
- [14] J.C. Nascimento, M.A.T. Figueiredo, and J.S. Marques, “Activity Recognition Using a Mixture of Vector Fields,” *Image Processing, IEEE Transactions on*, vol. 22, no. 5, pp. 1712–1725, May 2013.
- [15] Carlos Santiago, Jacinto C Nascimento, and Jorge S Marques, “2D Segmentation Using a Robust Active Shape Model With the EM Algorithm,” *Image Processing, IEEE Transactions on*, vol. 24, no. 8, pp. 2592–2601, Aug 2015.
- [16] P. Lucey, J.F. Cohn, T. Kanade, J. Saragih, Z. Ambadar, and I. Matthews, “The Extended Cohn-Kanade Dataset (CK+): A complete dataset for action unit and emotion-specified expression,” in *Computer Vision and Pattern Recognition Workshops (CVPRW), 2010 IEEE Computer Society Conference on*, 2010, pp. 94–101.