

A Closed-Form Solution for Image Sequence Segmentation with Dynamical Shape Priors

Frank R. Schmidt and Daniel Cremers

Computer Science Department
University of Bonn, Germany

Abstract. In this paper, we address the problem of image sequence segmentation with dynamical shape priors. While existing formulations are typically based on hard decisions, we propose a formalism which allows to reconsider all segmentations of past images. Firstly, we prove that the marginalization over all (exponentially many) reinterpretations of past measurements can be carried out in closed form. Secondly, we prove that computing the optimal segmentation at time t given all images up to t and a dynamical shape prior amounts to the optimization of a convex energy and can therefore be optimized globally. Experimental results confirm that for large amounts of noise, the proposed reconsideration of past measurements improves the performance of the tracking method.

1 Introduction

A classical challenge in Computer Vision is the segmentation and tracking of a deformable object. Numerous researchers have addressed this problem by introducing statistical shape priors into segmentation and tracking [1–7].

While in earlier approaches every image of a sequence was handled independently, Cremers [8] suggested to consider the correlations which characterize many deforming objects. The introduction of such *dynamical* shape priors allows to substantially improve the performance of tracking algorithms: The dynamics are learned via an auto-regressive model and segmentations of the preceding images guide the segmentation of the current image. Upon a closer look, this approach suffers from two drawbacks:

- The optimization in [8] was done in a level set framework which only allows for *locally* optimal solutions. As a consequence, depending on the initialization the resulting solutions may be suboptimal.
- At any given time the algorithm in [8] computed the currently optimal segmentation and only retained the segmentations of the two preceding frames. Past measurements were never reinterpreted in the light of new measurements. As a consequence, any incorrect decision would not be corrected at later stages of processing. While dynamical shape priors were called *priors with memory* in [8], what is memorized are only the decisions the algorithm took on previous frames – the measurements are instantly lost from memory, a reinterpretation is not considered in [8].

The reinterpretation of past measurements in the light of new measurements is a difficult computational challenge due to the exponential growth of the solution space: Even if a tracking system only had k *discrete* states representing the system at any time t , then after T time steps, there are k^T possible system configurations explaining all measurements. In this work silhouettes are represented by k *continuous* real-valued parameters: While determining the silhouette for time t amounts to an optimization in \mathbb{R}^k , the optimization over *all* silhouettes up to time T amounts to an optimization over $\mathbb{R}^{k \cdot T}$.

Recent works tried to address the above shortcomings. Papadakis and Memin suggested in [9] a control framework for segmentation which aimed at a consistent sequence segmentation by forward- and backward propagation of the current solution according to a dynamical system. Yet this approach is entirely based on level set methods and local optimization as well. Moreover, extrapolations into the past and the future rely on a sophisticated partial differential equation. In [10] the sequence segmentation was addressed in a convex framework. While this allowed to compute globally optimal solutions independent of initialization, it does not allow a reinterpretation of past measurements. Hence incorrect segmentations will negatively affect future segmentations.

The contribution of this paper is it to introduce a novel framework for image sequence segmentation which overcomes both of the above drawbacks. While [8, 10] compute the best segmentation given the current image and past *segmentations* here we propose to compute the best segmentation given the current image and all previous *images*. In particular we propose a statistical inference framework which gives rise to a marginalization over all possible segmentations of all previous images. The theoretical contribution of this work is therefore two-fold. Firstly, we prove that the marginalization over all segmentations of the preceding images can be solved in closed form which allows to handle the combinatorial explosion analytically. Secondly, we prove that the resulting functional is *convex*, such that the maximum a posteriori inference of the currently best segmentation can be solved globally. Experimental results confirm that this marginalization over preceding segmentations improves the accuracy of the tracking scheme in the presence of large amounts of noise.

2 An Implicit Dynamic Shape Model

In the following, we will briefly review the dynamical shape model introduced in [10]. It is based on the notion of a *probabilistic shape* u defined as a mapping

$$u : \Omega \rightarrow [0, 1] \tag{1}$$

that assigns to every pixel x of the shape domain $\Omega \subset \mathbb{R}^d$ the probability that this pixel is inside the given shape. While our algorithm will compute such a relaxed shape, for visualization of a silhouette we will simply threshold u at $\frac{1}{2}$. We present a general model for shapes in arbitrary dimension. However, the approach is tested for planar shapes ($d = 2$).

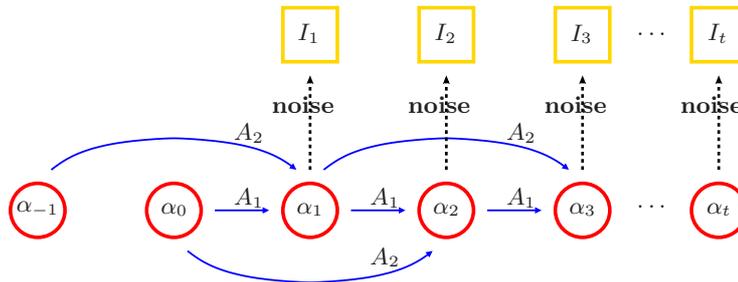


Fig. 1. Model for image sequence segmentation. We assume that all information about the observed images I_τ (top row) is encoded in the segmentation variables α_τ (bottom row) and that the dynamics of α_τ follow the autoregressive model (3) learned beforehand. If the state space was discrete with N possible states per time instance, then one would need to consider N^t different states to find the optimal segmentation of the t -th image. In Theorem 1, we provide a closed-form solution for the integration over all preceding segmentations. In Theorem 2, we prove that the final expression is convex in α_t and can therefore be optimized globally.

The space of all *probabilistic shapes* forms a convex set, and the space spanned by a few training shapes $\{u_1, \dots, u_N\}$ forms a convex subset. Any shape u can be approximated by a linear combination of the first n principal components Ψ_i of the training set:

$$u(x) \approx u_0(x) + \sum_{i=1}^n \alpha^i \cdot \Psi_i(x) \quad (2)$$

with an average shape u_0 . Also, the set

$$\mathcal{Q} := \{\alpha \in \mathbb{R}^n \mid \forall x \in \Omega : 0 \leq u_0 + \sum_{i=1}^n \alpha^i \cdot \Psi_i(x) \leq 1\}$$

of feasible α -parameters is convex [10].

Any given sequence of shapes u_1, \dots, u_N can be reduced to a sequence of low dimensional coefficient vectors $\alpha_1, \dots, \alpha_N \in \mathcal{Q} \subset \mathbb{R}^n$. The evolution of these coefficient vectors can be modeled as an autoregressive system

$$\alpha_i = \sum_{j=1}^k A_j \alpha_{i-j} + \eta_{\Sigma^{-1}} \quad (3)$$

of order $k \in \mathbb{N}$, where the transition matrices $A_j \in \mathbb{R}^{n \times n}$ describe the linear dependency of the current observation on the previous k observations. Here $\eta_{\Sigma^{-1}}$ denotes Gaussian noise with covariance matrix Σ^{-1} .

3 A Statistical Formulation of Sequence Segmentation

In the following, we will develop a statistical framework for image sequence segmentation which for any time t determines the most likely segmentation u_t given all images $I_{1:t}$ up to time t and given the dynamical model in (3). The goal is to maximize the conditional probability $\mathcal{P}(\alpha_t|I_{1:t})$, where $\alpha_t \in \mathbb{R}^n$ represents the segmentation $u_t := u_0 + \Psi \cdot \alpha_t$.

For the derivation we will make use of four concepts from probabilistic reasoning:

- Firstly, the conditional probability is defined as

$$\mathcal{P}(A|B) := \frac{\mathcal{P}(A, B)}{\mathcal{P}(B)}. \quad (4)$$

- Secondly, the application of this definition leads to the Bayesian formula

$$\mathcal{P}(A|B) = \frac{\mathcal{P}(B|A) \cdot \mathcal{P}(A)}{\mathcal{P}(B)} \quad (5)$$

- Thirdly, we have the concept of marginalization:

$$\mathcal{P}(A) = \int \mathcal{P}(A|B) \cdot \mathcal{P}(B) \, dB \quad (6)$$

which represents the probability $\mathcal{P}(A)$ as a weighted integration of $\mathcal{P}(A|B)$ over all conceivable states B . In the context of time-series analysis this marginalization is often referred to as the Chapman-Kolmogorov equation [11]. In particle physics it is popular in the formalism of path integral computations.

- Fourthly, besides these stochastic properties we make the assumption that for any time τ the probability for measuring image I_τ is completely characterized by its segmentation α_τ as shown in Figure 1:

The segmentation α_τ contains all information about the system in state τ . The rest of the state τ is independent noise. Hence, I_τ contains no further hidden information, its probability is uniquely determined by α_τ .

(7)

With these four properties, we can now derive an expression for the probability $\mathcal{P}(\alpha_t|I_{1:t})$ that we like to maximize. Using Bayes rule with all expressions in (5) conditioned on $I_{1:t-1}$, we receive

$$\mathcal{P}(\alpha_t|I_{1:t}) \propto \mathcal{P}(I_t|\alpha_t, I_{1:t-1}) \cdot \mathcal{P}(\alpha_t|I_{1:t-1}) \quad (8)$$

Due to property (7), we can drop the dependency on the previous images in the first factor. Moreover, we can expand the second factor using Bayes rule again:

$$\mathcal{P}(\alpha_t|I_{1:t}) \propto \mathcal{P}(I_t|\alpha_t) \cdot \mathcal{P}(I_{1:t-1}|\alpha_t) \cdot \mathcal{P}(\alpha_t) \quad (9)$$

Applying the Chapman-Kolmogorov equation (6) to (9), we obtain

$$\mathcal{P}(\alpha_t|I_{1:t}) \propto \mathcal{P}(I_t|\alpha_t) \int \mathcal{P}(I_{1:t-1}|\alpha_{1:t}) \cdot \underbrace{\mathcal{P}(\alpha|\alpha_t) \cdot \mathcal{P}(\alpha_t)}_{\mathcal{P}(\alpha_{1:t})} d\alpha_{1:t-1} \quad (10)$$

This expression shows that the optimal solution for α_t requires an integration over all conceivable segmentations $\alpha_{1:t-1}$ of the preceding images.

To evaluate the right hand side of (10), we will model the probabilities $\mathcal{P}(I_t|\alpha_t)$, $\mathcal{P}(I_{1:t-1}|\alpha_{1:t})$ and $\mathcal{P}(\alpha_{1:t})$. Assuming a spatially independent pre-learned color distribution \mathcal{P}_{ob} of the object and \mathcal{P}_{bg} of the background, we can define $p(x) := -\log(\mathcal{P}_{\text{ob}}(x)/\mathcal{P}_{\text{bg}}(x))$ which is negative for every pixel that is more likely to be an object pixel than a background pixel. By introducing an exponential weighting parameter γ for the color distributions, $\mathcal{P}(I_t|\alpha_t)$ becomes

$$\begin{aligned} \mathcal{P}(I_t|\alpha_t) &= \prod_{x \in \Omega} \mathcal{P}_{\text{ob}}(x)^{\gamma u_t(x)} \mathcal{P}_{\text{bg}}(x)^{\gamma(1-u_t(x))} \propto \exp\left(\sum_{x \in \Omega} \gamma u_t(x) \log\left(\frac{\mathcal{P}_{\text{ob}}(x)}{\mathcal{P}_{\text{bg}}(x)}\right)\right) \\ &\propto \exp\left(-\sum_{i=1}^n \gamma \cdot (\alpha_t)_i \cdot \underbrace{\left(\sum_{x \in \Omega} \Psi_i(x) \cdot p(x)\right)}_{f_{t,i}}\right) = \exp(-\gamma \langle a_t, f_t \rangle). \end{aligned}$$

To compute $\mathcal{P}(I_{1:t-1}|\alpha_{1:t})$, we use the assumption (7). Besides the information encoded in $\alpha_{1:t}$, the images I_τ contain no further informations and are therefore pairwise independent:

$$\mathcal{P}(I_{1:t-1}|\alpha_{1:t}) = \prod_{\tau=1}^{t-1} \mathcal{P}(I_\tau|\alpha_{1:t}) = \prod_{\tau=1}^{t-1} \mathcal{P}(I_\tau|\alpha_\tau) = \prod_{\tau=1}^{t-1} \exp(-\gamma \langle a_\tau, f_\tau \rangle)$$

The second equation holds again due to (7): Since the probability for I_τ is uniquely determined by α_τ , the dependency on the other states can be dropped.

Now, we have to address the probability $\mathcal{P}(\alpha_{1:t})$ which can recursively be simplified via (4):

$$\mathcal{P}(\alpha_{1:t}) = \mathcal{P}(\alpha_t|\alpha_{1:t-1}) \cdot \mathcal{P}(\alpha_{1:t-1}) = \dots = \prod_{\tau=1}^{t-1} \mathcal{P}(\alpha_\tau|\alpha_{1:\tau-1}) \quad (11)$$

Using the dynamic shape prior (3), this expression becomes

$$\mathcal{P}(\alpha_{1:t}) \propto \prod_{\tau=1}^{t-1} \exp\left(-\left\|\alpha_\tau - \sum_{i=1}^k A_i \alpha_{\tau-i}\right\|_{\Sigma^{-1}}^2\right)$$

To make this formula more accessible, we introduced k additional segmentation parameters $\alpha_{1-k}, \dots, \alpha_0$. These parameters represent the segmentation of the past prior to the first observation I_1 (cf. Figure 1). To simplify the notation, we will introduce $\boldsymbol{\alpha} := \alpha_{1-k:t-1}$. These are the parameters that represent all segmentations prior to the current segmentation α_t .

Combining all derived probabilities, we can formulate the image segmentation as the following minimization task

$$\arg \min_{\boldsymbol{\alpha}} \int \exp \left(- \sum_{\tau=1}^t \gamma \cdot \langle f_{\tau}, \alpha_{\tau} \rangle - \sum_{\tau=1}^t \left\| \alpha_{\tau} - \sum_{j=1}^k A_j \alpha_{\tau-j} \right\|_{\Sigma^{-1}}^2 \right) d\boldsymbol{\alpha} \quad (12)$$

Numerically computing this $n \cdot (t + k - 1)$ -dimensional integral of (12) leads to a combinatorial explosion. Even for a simple example of $t = 25$ frames, $n = 5$ eigenmodes and an autoregressive model size of $k = 1$, a 100-dimensional integral has to be computed. In [8], this computational challenge was circumvented by the crude assumption of a Dirac distribution centered at precomputed segmentation results – i.e. rather than considering all possible trajectories the algorithm only retained for each previous time the one segmentation which was then most likely.

In this paper, we will compute this integral explicitly and receive a closed-form expression for (12) described in Theorem 1. This closed-form formulation has the important advantage that for any given time it allows an optimal reconsideration of all conceivable previous segmentations.

To simplify (12), we write the integral as $\int \exp(Q(\boldsymbol{\alpha}, \alpha_t)) d\boldsymbol{\alpha}$. Note that Q is a quadratic expression that can be written as

$$Q(\boldsymbol{\alpha}, \alpha_t) = \underbrace{\gamma \cdot \langle f_t, \alpha_t \rangle}_I + \underbrace{\| \alpha_t \|_{\Sigma^{-1}} + \langle \boldsymbol{\alpha}, M \boldsymbol{\alpha} \rangle}_{II} - \underbrace{\langle b, \boldsymbol{\alpha} \rangle}_{III} \quad (13)$$

with the block vector b and the block matrix M :

$$b_i = \underbrace{-\gamma \cdot f_i}_{i \geq 1} + \underbrace{2A_{t-i}^T \Sigma^{-1} \alpha_t}_{i \geq t-k}$$

$$M_{i,j} = \underbrace{\Sigma A_{t-i}^T \Sigma^{-1} A_{t-j}}_{i,j \geq t-k} + \underbrace{\mathbb{1}}_{i=j \geq 1} - \underbrace{2A_{i-j}}_{\substack{i \geq 1 \\ k \geq i-j \geq 1}} + \sum_{\substack{1 \leq l \leq k \\ 1 \leq i+l \leq t-1 \\ 1 \leq i-j+l \leq k}} \Sigma A_l^T \Sigma^{-1} A_{i-j+l}$$

Despite their complicated nature, the three terms in (13) have the following intuitive interpretations:

- I assures that the current segmentation encoded by α_t optimally segments the current image.
- II assures that the *segmentation path* $(\alpha_{-1}, \dots, \alpha_t)$ is consistent with the learned autoregressive model encoded by (A_i, Σ^{-1}) .

- *III* assures that the current segmentation α_t also consistently segments all previous images when propagated back in time according to the dynamical model. In dynamical systems such backpropagation is modeled by the adjoints A^T of the transition matrices.

In the next theorem we will provide a closed form expression for (12) that is freed of any integration process and can therefore be computed more efficiently. Additionally, we will come up with a convex energy functional. Therefore, to compute the global optimum of the image sequence problem is an easy task.

Theorem 1. *The integration over all conceivable interpretations of past measurements can be solved in the following closed form:*

$$\mathcal{P}(\alpha_t|I_{1:t}) = \exp \left[-\gamma \langle \alpha_t, f_t \rangle - \|\alpha_t\|_{\Sigma^{-1}}^2 + \frac{1}{4} \langle M_s^{-1} b, b \rangle + \text{const} \right] \quad (14)$$

Proof.

$$\begin{aligned} \mathcal{P}(\alpha_t|I_{1:t}) &\propto \int e^{-\gamma \langle \alpha_t, f_t \rangle - \|\alpha_t\|_{\Sigma^{-1}}^2 - \langle \alpha, M_s \alpha \rangle + \langle b, \alpha \rangle} d\alpha \\ &= \int e^{-\langle \alpha_t, f_t \rangle - \|\alpha_t\|_{\Sigma^{-1}}^2 - \|\alpha - \frac{1}{2} M_s^{-1} b\|_{M_s}^2 + \frac{1}{4} \|M_s^{-1} b\|_{M_s}^2} d\alpha \\ &\propto \exp \left[-\gamma \langle \alpha_t, f_t \rangle - \|\alpha_t\|_{\Sigma^{-1}}^2 + \frac{1}{4} \langle M_s^{-1} b, b \rangle \right] \end{aligned}$$

□

Theorem 2. *The resulting energy $E(\alpha_t) = -\log(\mathcal{P}(\alpha_t|I_{1:t}))$ is convex and can therefore be minimized globally.*

Proof. The density function $\mathcal{P}(\alpha_t|I_{1:t})$ is the integral of a log-concave function, i.e., their logarithm is a concave function. It was shown in [12] that integrals of log-concave functions are log-concave. Hence, E is convex. Therefore, the global optimum can be computed using, for example, a gradient descent approach. □

In [10], discarding all preceding images and merely retaining the segmentations of the last frames gave rise to the simple objective function:

$$E_1(\alpha_t) = \gamma \cdot \langle \alpha_t, f_t \rangle + \|\alpha_t - v\|_{\Sigma^{-1}}^2 \quad (15)$$

where v is the prediction obtained using the AR model (3) on the basis of the last segmentations.

The proposed optimal path integration gives rise to the new objective function

$$E_2(\alpha_t) = \gamma \cdot \langle \alpha_t, f_t \rangle + \|\alpha_t\|_{\Sigma^{-1}}^2 - \frac{1}{4} \langle M_s^{-1} b, b \rangle \quad (16)$$

In the next section, we will experimentally quantify the difference in performance brought about by the proposed marginalization over preceding segmentations.

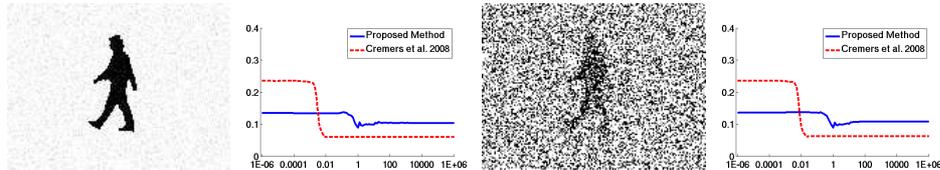


Fig. 2. Optimal Parameter Estimation. The tracking error averaged over all frames (plotted as a function of γ) shows that $\gamma = 1$ produces the best results for both methods at various noise levels (shown here are $\sigma = 16$ and $\sigma = 256$).

4 Experimental Results

In the following experiments, the goal is to track a walking person in spite of noise and missing data. To measure the tracking accuracy, we handsegmented the sequence (before adding noise) and measured the relative error with respect to this ground truth. Let $T : \Omega \rightarrow \{0, 1\}$ be the true segmentation and $S : \Omega \rightarrow \{0, 1\}$ be the estimated one. Then we define the scaled relative error ϵ as

$$\epsilon := \frac{\int_{\Omega} |S(x) - T(x)| dx}{2 \cdot \int_{\Omega} T(x) dx}.$$

It measures the area difference relative to twice the area of the ground truth. Thus we have $\epsilon=0$ for a perfect segmentation and $\epsilon=1$ for a completely wrong segmentation (of the same size).

Optimal parameter estimation.

In order to estimate the optimal parameter γ for both approaches, we added Gaussian noise of standard deviation σ to the training images. As we can see in Figure 2, the lowest tracking error ϵ (averaged over all frames) is obtained at $\gamma = 1$ for both approaches. Therefore, we will fix $\gamma = 1$ for the test series in the next section.

Robust tracking through prominent noise.

The proposed framework allows to track a deformable silhouette despite large amounts of noise. Figure 3 shows segmentation results obtained with the proposed method for various levels of Gaussian noise. The segmentations are quite accurate even for high levels of noise.

Quantitative comparison to the method in [10].

For a quantitative comparison of the proposed approach with the method of [10], we compute the average error ϵ of the learned input sequence $I_{1:151}$ for different levels of Gaussian noise. Figure 4 shows two different aspects. While the method in [10] exhibits slightly lower errors for small noise levels, the proposed method shows less dependency on noise and exhibits substantially better performance at larger noise levels. While the difference in the segmentation results for low noise level are barely recognizable (middle row), for high noise level, the method in [10] clearly estimates incorrect poses (bottom row).

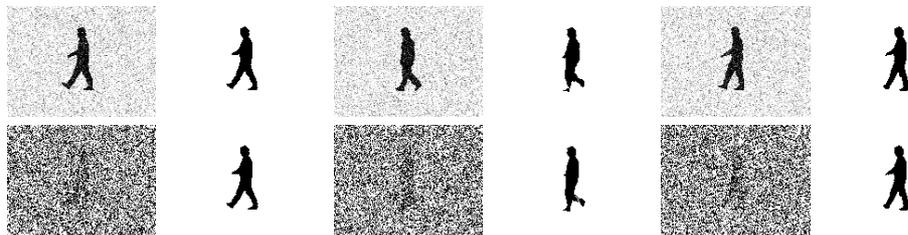


Fig. 3. Close-ups of segmentation results. The proposed method gets correct segmentation results. Even at the presence of high Gaussian noise ($\sigma \in \{64, 512\}$).

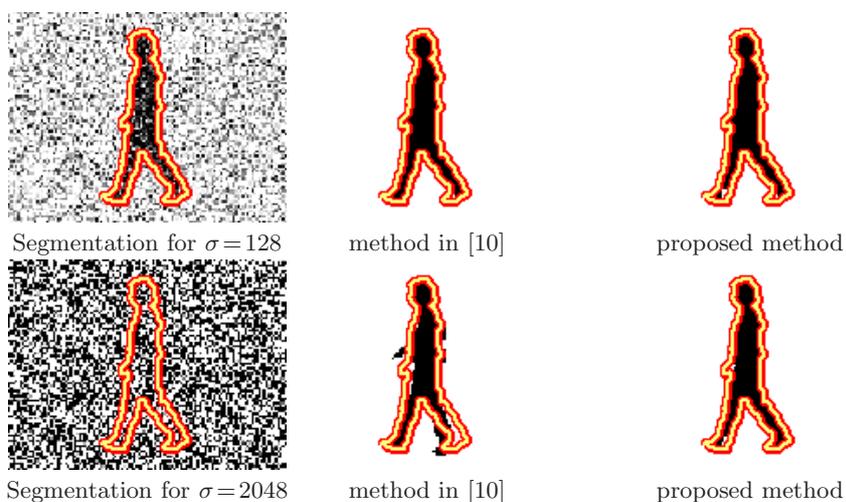
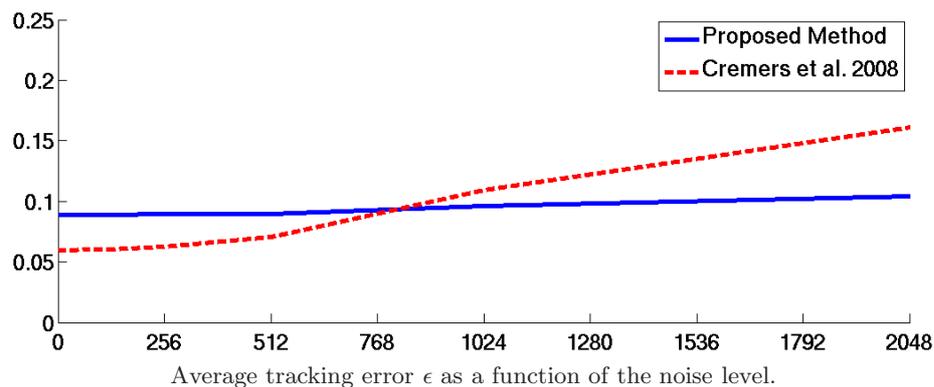


Fig. 4. Robustness with respect to noise. Tracking experiments demonstrate that in contrast to the approach in [10], the performance of the proposed algorithm is less sensitive to noise and outperforms the former in the regime of large noise. While for low noise, the resulting segmentations are qualitatively similar (middle row), for high noise level, the method in [10] provides an obviously wrong pose estimate (bottom row).

5 Conclusion

In this paper we presented the first approach for variational object tracking with dynamical shape priors which allows to marginalize over *all* previous segmentations. Firstly, we proved that this marginalization over an exponentially growing space of solutions can be solved analytically. Secondly, we proved that the resulting functional is convex. As a consequence, one can efficiently compute the globally optimal segmentation at time t given all images up to time t .

In experiments, we confirmed that the resulting algorithm allows to reliably track walking people despite prominent noise. In particular for very large amounts of noise, it outperforms an alternative algorithm [10] that does not include a marginalization over the preceding segmentations.

References

1. Leventon, M., Grimson, W., Faugeras, O.: Statistical shape influence in geodesic active contours. In: IEEE Int. Conf. on Computer Vision and Pattern Recognition. Volume 1., (2000) 316–323
2. Tsai, A., Yezzi, A., Wells, W., Tempany, C., Tucker, D., Fan, A., Grimson, E., Willsky, A.: Model-based curve evolution technique for image segmentation. In: Comp. Vision Patt. Recog., (2001) 463–468
3. Cremers, D., Kohlberger, T., Schnörr, C.: Nonlinear shape statistics in Mumford–Shah based segmentation. In Heyden, A., et al., eds.: Europ. Conf. on Comp. Vis. Volume 2351 of LNCS., Springer (2002) 93–108
4. Riklin-Raviv, T., Kiryati, N., Sochen, N.: Unlevel sets: Geometry and prior-based segmentation. In Pajdla, T., Hlavac, V., eds.: European Conf. on Computer Vision. Volume 3024 of LNCS., Springer (2004) 50–61
5. Rousson, M., Paragios, N., Deriche, R.: Implicit active shape models for 3d segmentation in MRI imaging. In: MICCAI. Volume 2217 of LNCS., Springer (2004) 209–216
6. Kohlberger, T., Cremers, D., Rousson, M., Ramaraj, R.: 4d shape priors for level set segmentation of the left myocardium in SPECT sequences. In: Medical Image Computing and Computer Assisted Intervention. Volume 4190 of LNCS. (2006) 92–100
7. Charpiat, G., Faugeras, O., Keriven, R.: Shape statistics for image segmentation with prior. In: IEEE Int. Conf. on Computer Vision and Pattern Recognition. (2007)
8. Cremers, D.: Dynamical statistical shape priors for level set based tracking. IEEE PAMI **28**(8) (August 2006) 1262–1273
9. Papadakis, N., Mémin, E.: Variational optimal control technique for the tracking of deformable objects. In: IEEE Int. Conf. on Comp. Vis. (2007)
10. Cremers, D., Schmidt, F.R., Barthel, F.: Shape priors in variational image segmentation: Convexity, Lipschitz continuity and globally optimal solutions. In: IEEE Int. Conf. on Computer Vision and Pattern Recognition (2008)
11. Papoulis, A.: Probability, Random Variables, and Stochastic Processes. McGraw-Hill, New York (1984)
12. Prékopa, A.: Logarithmic concave measures with application to stochastic programming. Acta Scientiarum Mathematicarum **34** (1971) 301–316