

# Spatiotemporal Smoothing as a Basis for Facial Tissue Tracking in Thermal Imaging

Yan Zhou, Panagiotis Tsiamyrtzis, Peggy Lindner, Ilya Timofeyev, and Ioannis Pavlidis\*, *Senior Member, IEEE*

**Abstract**—Accurate tracking of facial tissue in thermal infrared imaging is challenging because it is affected not only by positional but also physiological (functional) changes. This paper presents a particle filter tracker driven by a probabilistic template function with both spatial and temporal smoothing components, which is capable of adapting to abrupt positional and physiological changes. The method was tested on tracking facial regions of subjects under varying physiological and environmental conditions in 25 thermal clips. It demonstrated robustness and accuracy, outperforming other strategies. This new method promises improved performance in a number of biomedical applications that involve physiological measurements on the face, such as unobtrusive sleep and stress studies.

**Index Terms**—Facial tracking, matte, sleep studies, stress studies, thermal imaging.

## I. INTRODUCTION

IN the last few years, facial tracking in the thermal infrared spectrum received increasing attention. Initially, applications in surveillance and face recognition were the driving force, where thermal imaging has the distinct advantage of being insensitive to lighting conditions [1] [2]. Later, physiological variables, such as vital signs, proved measurable in this modality [3]–[6], which gave rise to applications in human–computer interaction (HCI) [7], medicine [8], and psychology [9]. The degree of success of such measurements depends on a tracking method that can reliably follow the tissue of interest over time. For example, in sleep studies, if the tracker momentarily loses the nasal region of interest (ROI), the generated breathing signal is far from accurate (see Fig. 1), which affects the ensuing analysis. Thus, the specification of a facial tracker in thermal infrared needs to be quite stringent.

Manuscript received November 5, 2012; accepted November 14, 2012. Date of publication December 11, 2012; date of current version April 15, 2013. This work was supported in part by the Defense Academy for Credibility Assessment and the National Science Foundation (NSF) under Grant # IIS-0812526 entitled “Do Nintendo Surgeons Defy Stress.” *Asterisk indicates corresponding author.*

Y. Zhou was with the Department of Computer Science, University of Houston, Houston, TX 77004 USA. She is now with Elekta Inc., Maryland Heights, MO 63043 USA (e-mail: yan.zhou@elekta.com).

P. Tsiamyrtzis is with the Department of Statistics, Athens University of Economics and Business, Athens 10434, Greece (e-mail: pt@aueb.gr).

P. Lindner is with the Department of Computer Science, University of Houston, Houston, TX 77004 USA (e-mail: plindner@uh.edu).

I. Timofeyev is with the Department of Mathematics, University of Houston, Houston, TX 77004 USA (e-mail: ilya@math.uh.edu).

\*I. Pavlidis is with the Department of Computer Science, University of Houston, Houston, TX 77004 USA (e-mail: ipavlidis@uh.edu).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TBME.2012.2232927

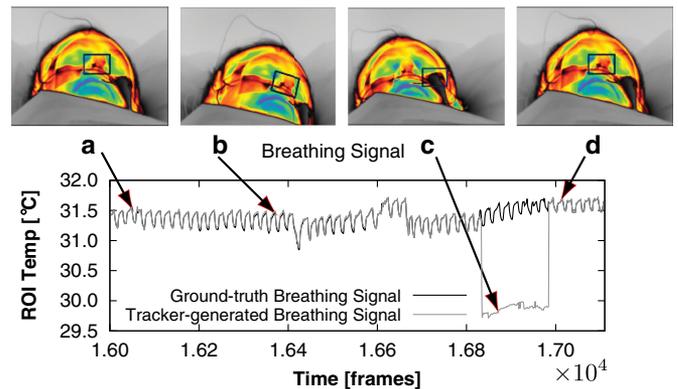


Fig. 1. Tracker generated compared with ground-truth breathing signal. (a) Initial frame with the rectangular ROI centered on the nostrils. (b) and (d) When the tracker works well, the generated breathing signal is good. (c) When the tracker loses the ROI, the generated breathing signal is inaccurate.

The proposed method uses a particle-filter tracker, which is driven by a template-based objective function. The choice has to do with the peculiarities of thermal imaging and the needs of the targeted applications. Tracking based on shape models [10] is not very appealing in facial thermal imaging, because the modality images function rather than structure.

To give an example, consider the case of tracking nasal tissue in thermal infrared imagery for the purpose of computing breathing function. Under normal conditions, the nose is colder than the surrounding tissue due to convection from nasal air flow. This translates to a characteristic thermal shape similar to the one appearing in visual (structural) images. At some point, an irritant reaches the subject’s nasal cavity, there is an allergic reaction that blocks air flow in the nostrils and breathing continues mainly through the mouth. Because air flow is severely curtailed, the temperature over the nasal tissue rises and the nose blends with the surrounding tissue in the imagery. The nose’s functionality has dramatically changed and so its characteristic shape (see Fig. 2). In such cases where stochastic physiological changes affect thermal emission, a shape model tracker may encounter significant difficulties.

Note that spatial resolution in thermal imagery is typically lower than that in visual imagery ( $640 \times 512$  pixels in our case). Edges in thermal imagery are fuzzy due to diffusion, complicating matters further. None of these factors is conducive to shape-based tracking. Finally, many of the targeted applications are in medicine and HCI. This necessitates a computationally “light” tracker for real-time performance. For example, if the nasal tracking and signal extraction method were not real time, then it would be impossible for a medical technician

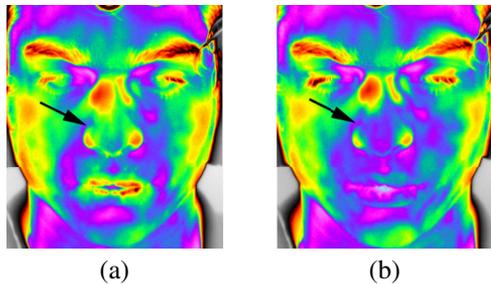


Fig. 2. (a) Thermal image of a subject with normal breathing function. The nasal area is at a contrast with the surrounding tissue. (b) Thermal image of the same subject a few minutes later, when his breathing function is impaired by an irritant. Notice the blending of the nasal area with the neighboring tissue.

to optimize the continuous positive airway pressure in a sleep intervention [11].

Although trackers based on shape models are not appropriate for the problem at hand, tracking methods based on statistical filtering, such as Kalman or particle filters, are quite appealing. In particular, we chose to proceed with particle filtering because in the context of sleep and stress studies of interest, the subjects exhibit infrequent and abrupt turns of the head, which are highly nonlinear. Indeed, particle filtering is not only a general mechanism free of explicit modeling, but it can also handle nonlinear motion in the predict-update loop. We opted to implement the update operation in the predict-update loop through a probabilistic template algorithm. We will demonstrate that this combination of particle filtering with a probabilistic template produces a fast, flexible, and accurate tracker, fulfilling the specifications of the application domain.

#### A. Previous Work

The literature on templates and particle filter tracking is vast and well known. This section focuses on a few representative methods, some of which have been used as comparative yardsticks in the experimental part. It is by no means an exhaustive literature account.

In the visual imaging domain, Matthews and Ishikawa [12] developed a drift correction algorithm by computing principal components of previous image ROIs. This method prevented tracker drifts but at a high computational cost, prohibiting real-time applications.

Also in the visual imaging domain, Jepson *et al.* [13] proposed a statistical appearance model that weighted heavily pixels with stable behavior. This model served as the template and worked well when the subject’s facial visual appearance exhibited small variations. It failed, however, when the subject’s face exhibited large appearance changes providing very few stable pixels. A different version of appearance modeling paired with particle filter tracking is described in [14].

In general, appearance modeling is a powerful template mechanism but has two weak points: it does not cope well with sudden change and saturates after a long tracking period (needs restarting). Both are problematic in the context of thermal fa-

cial imaging because 1) physiology can produce abrupt changes (e.g., perspiration) and so does head motion. 2) Sleep and stress studies (which are the applications of interest) sometimes last hours.

In the thermal imaging domain, Eveland *et al.* [1] proposed a particle filter facial tracking algorithm driven by a Bayesian formula. The method required training for forming the likelihood distributions. Dowdall *et al.* [15] proposed a network of particle filter trackers driven by static template functions. This was an extremely rigid assumption, which was ameliorated by the multitrapper composition. In other words, most often than not, one or two tracking regions in the tracking network remained relatively stable and provided the necessary support to keep the network tracker on target. In a later incarnation [16], this “network tracker” adopted a dynamic template scheme by computing successive differences of intensity, pixelwise. Each template pixel was updated or not, based on whether the respective difference exceeded or not a predetermined threshold. The threshold was set to the maximum plausible physiological difference. Such a zero-one approach (i.e., either do not or do update a pixel) could not handle well certain changes because it was either overcommitting or undercommitting, depending on the magnitude of its guiding threshold. It was also plagued by the drifting problem, due to its deterministic nature.

This paper describes a particle filter tracking method driven by a novel probabilistic template mechanism. This mechanism is based in part on a fuzzy mask (Matte) [17], which was originally developed for segmentation purposes. To the best of our knowledge, it is the first time that it is adopted for tracking purposes. The strong point of Matte for the problem at hand is that it is based on pixel dependence (spatial smoothness). Indeed, there is spatial smoothness in thermophysiological imagery of the face. Muscular areas are relatively homogeneous and so are vascular areas. This is in contrast to the pixel independence assumption of appearance modeling, which is not realistic here. We have also introduced a temporal smoothness assumption, by modifying the Matte formula accordingly. This assumption holds true for appropriately small time windows and reduces oscillation.

The rest of this paper is organized as follows: Section II describes the methodology. In Section III, the experimental results demonstrate the relative advantage of the method with respect to other plausible approaches. Finally, Section IV concludes this paper. An early, short version of this paper appeared in the 2009 proceedings of the International Conference on Medical Image Computing and Computer Assisted Intervention (MICCAI 2009) [18]. This paper provides a thorough description of the methodology and reports results on a much larger and more diverse dataset with respect to [18].

## II. METHODOLOGY

We use particle filtering to track the ROI’s position in the current frame, based on template matching. We denote the motion state of the tracker at time  $t$  by  $\mathbf{X}_t$  and its observations by  $\mathbf{Z}_t$ . The state of the tracker consists of three variables  $\mathbf{X}_t = (x_t, y_t, \theta_t)$ ,

where  $x_t$  and  $y_t$  are the spatial coordinates of the ROI's centroid and  $\theta_t$  is the ROI's rotation angle. The observation  $\mathbf{Z}_t$  refers to the pixels in the current frame. The particle filter tracker uses  $N = 100$  particles (candidate ROIs) in a single iteration per frame. We resample particles for each frame according to

$$\mathbf{X}_t = \mathbf{X}_{t-1} + \mathbf{V}_t \quad (1)$$

where  $\mathbf{V}_t$  is a 2-D independent and identically distributed Gaussian noise process. The mean of  $\mathbf{V}_t$  is zero and the variance is one of the parameters that can be adjusted to ensure an efficient filter performance.

The particle filter algorithm approximates the posterior distribution  $p(\mathbf{X}_t|\mathbf{Z}_{1:t})$  via a set of weighted particles  $\Omega_t = \{\mathbf{Y}_t^r, \omega_t^r\}_{r=1}^N$ , where  $\sum_{r=1}^N \omega_t^r = 1$ ; the particles  $\mathbf{Y}_t^r$  are weighted with respect to  $p(\mathbf{X}_t|\mathbf{Z}_{1:t})$ . The posterior distribution  $p(\mathbf{X}_t|\mathbf{Z}_{1:t})$  involves the entire time history of observations  $\mathbf{Z}_{1:t}$ , which is computationally expensive. For this reason, we use instead the posterior distribution  $p(\mathbf{X}_t|\mathbf{Z}_t)$  that depends on the most recent observation only. Then, the maximum *a posteriori* (MAP) estimate to determine the state of the tracker becomes

$$\hat{\mathbf{X}}_t = \arg \max_{\mathbf{Y}_t^r} p(\mathbf{Y}_t^r|\mathbf{Z}_t), \quad \text{where } p(\mathbf{Y}_t^r|\mathbf{Z}_t) = \omega_t^r. \quad (2)$$

More specifically, the weight values are proportional to the correlation coefficients between the template  $\mathbf{T}_{t-1}$  and the corresponding candidate ROIs  $\mathcal{F}_t^r$

$$\omega_t^r \propto \frac{\sum_i (\mathbf{T}_{t-1}[i] - \mu(\mathbf{T}_{t-1})) (\mathcal{F}_t^r[i] - \mu(\mathcal{F}_t^r))}{\sigma(\mathbf{T}_{t-1})\sigma(\mathcal{F}_t^r)} \quad (3)$$

where  $\mu(\mathbf{T}_{t-1})$  and  $\mu(\mathcal{F}_t^r)$  are the means of  $\mathbf{T}_{t-1}$  and  $\mathcal{F}_t^r$ , respectively;  $\sigma(\mathbf{T}_{t-1})$  and  $\sigma(\mathcal{F}_t^r)$  denote the standard deviations of  $\mathbf{T}_{t-1}$  and  $\mathcal{F}_t^r$ , respectively; and, the index  $i$  denotes the  $i$ th pixel in the  $r$  candidate ROI (particle) or in the existing template.

The MAP estimate selects the current ROI  $\mathcal{F}_t^m$  from the candidate ROIs  $\mathcal{F}_t^r$  ( $r = 1, \dots, 100$ ). For notational simplicity, we will denote the current ROI as  $\mathcal{F}_t$  from this point onward. Following the selection of the current ROI [see Fig. 3(a)], the computations for the time step  $t$  complete with the updating of the template [see Fig. 3(b)].

The template updating strategy aims 1) to update pixels that exhibit significant intensity variations, in order to adapt to physiological and orientation changes; and 2) to preserve pixels with insignificant intensity alterations, in order to prevent drifting. The former pixels constitute the unstable, while the latter pixels the stable category. At each time step  $t$ , the key issue is to determine the degree of updating for each pixel in the template. To model this process, we assume that the pixel intensity is a convex combination of a stable and an unstable map

$$\mathcal{F}_t[i] = \alpha_t[i]\mathbf{S}_t[i] + (1 - \alpha_t[i])\mathbf{U}_t[i] \quad (4)$$

where  $\mathcal{F}_t[i]$  is the intensity of the  $i$ th pixel in the ROI under consideration and  $\alpha_t^i$  is the mask value of the  $i$ th pixel that determines the updating ratio;  $\mathbf{S}_t$  and  $\mathbf{U}_t$  denote the stable and unstable maps, respectively. The parameters on the right-hand side of (4) are unknown and the goal is to solve for  $\alpha_t[i]$ .

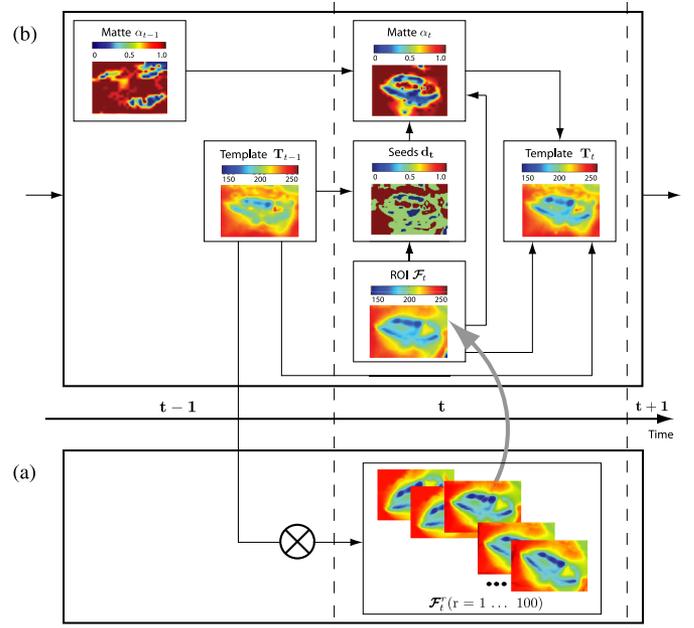


Fig. 3. Illustration of algorithmic flow. (a) First, the likely current ROI is selected based on a MAP estimate (particle filtering). (b) Then, the template is updated based on the formation of a Matte from stable/unstable seeds, the previous template, and the MAP ROI estimate. The updated template is used in the next time step.

To this end, a mask matrix  $\alpha_t$  is constructed with the following properties.

- 1) It has the same size as the template. This one-to-one correspondence between the template and the mask is used to identify which template pixels need to be updated.
- 2) Each entry in the mask matrix is a number in the range  $[0, 1]$  indicating the probability that the corresponding pixel is stable. The larger the probability is, the less the pixel needs to be updated. The most stable and the most unstable pixels in the Matte provide the seed values, which initialize the current Matte computation.
- 3) The mask takes into account spatial information (i.e., changes occur in regions and not in isolated pixels), providing a smooth outcome on the image plane.
- 4) The mask also takes into account temporal information (i.e., changes occur in finite time windows), providing a smooth outcome along the time line.

The user manually inputs the initial template in the first frame, by selecting a rectangular ROI with the mouse. Next, for each incoming frame, the method automatically performs the template updating according to the following steps (see Fig. 3).

*Step 1:* Extraction of stable and unstable seeds.

*Step 2:* Computation of the spatiotemporal Matte (STM).

*Step 3:* Update of the template.

The thus formed template correlates with candidate ROIs in the next time step to form the weights in the particle filtering process [see (3)].

### A. Extraction of Stable and Unstable Seeds

Initially, the method selects the most stable and unstable pixels in the ROI under consideration. These pixels constitute the seeds for the Matte computation step. The criteria for extracting maximally stable and unstable pixels are met when pixelwise intensity differences of the current frame from the template exceed predetermined thresholds. Thus, if we denote by  $\mathcal{F}_t[i]$  the ROI's  $i$ th pixel intensity at the current frame  $t$  and by  $\mathbf{T}_{t-1}[i]$  the template's  $i$ th pixel intensity at the previous frame  $t-1$ , then pixel  $i$  is

$$\begin{aligned} \text{stable if } |\mathcal{F}_t[i] - \mathbf{T}_{t-1}[i]| &< \lambda_1 \\ \text{unstable if } |\mathcal{F}_t[i] - \mathbf{T}_{t-1}[i]| &> \lambda_2 \end{aligned} \quad (5)$$

where,  $\lambda_1 < \lambda_2$  are predetermined thresholds that delineate the fuzzy range of physiologically plausible temperature differences. Values to the left of this range are certainly noise, while to the right are manifestations of strong physiological or other change. There is relative flexibility in the choice of  $\lambda_1$  and  $\lambda_2$  values, an issue that is thoroughly investigated in Section III-D.

The Matte values are set to 0 or 1 at the locations of extremely unstable or stable pixels (seeds), respectively. Initially, only the seed entries are known in the Matte; the next section describes how the method estimates nonseed entries.

### B. Computation of the STM

Equation (4) is similar to the one appearing in [17] and various methods to solve for  $\alpha_t^i$  have been proposed in [17], and [19]–[21]. There are three unknowns on the right-hand side of (4) and additional minimization constraints need to be introduced to solve for the parameters  $\alpha_t[i]$ ,  $\mathbf{S}_t[i]$ , and  $\mathbf{U}_t[i]$ . In [17], the authors introduce the parameters  $\mathbf{a}[i] = 1/(\mathbf{S}[i] - \mathbf{U}[i])$ ,  $\mathbf{b}[i] = \mathbf{U}[i]/(\mathbf{S}[i] - \mathbf{U}[i])$ ; then, solve for  $\alpha[i]$ ,  $\mathbf{a}[i]$ , and  $\mathbf{b}[i]$  by minimizing the cost function

$$J_S(\alpha, \mathbf{a}, \mathbf{b}) = \sum_{i \in \mathcal{F}} \left( \sum_{j \in \mathbf{w}^i} (\alpha[j] - \mathbf{a}[i]\mathcal{F}[j] - \mathbf{b}[i])^2 + \epsilon \mathbf{a}[i]^2 \right). \quad (6)$$

It is demonstrated in [17] that the minimization problem with respect to  $\mathbf{a}$  and  $\mathbf{b}$  can be solved explicitly, i.e., optimal  $\mathbf{a}$  and  $\mathbf{b}$  can be obtained explicitly in terms of  $\alpha$ . Therefore, the minimization problem can be recast as a quadratic programming problem in terms of  $\alpha$  alone.

In this paper, we formulate the cost function as

$$J_T(\alpha_t, \mathbf{a}_t, \mathbf{b}_t) = \sum_{i \in \mathcal{F}_t} \left[ \sum_{j \in \mathbf{w}_t^i} \left( (\alpha_t[j] - \mathbf{a}_t[i]\mathcal{F}_t[j] - \mathbf{b}_t[i])^2 + \epsilon \mathbf{a}_t[i]^2 \right) + (\alpha_t[i] - \alpha_{t-1}[i])^2 \right] \quad (7)$$

where,  $\mathbf{a}_t[i] = 1/(\mathbf{S}_t[i] - \mathbf{U}_t[i])$  and  $\mathbf{b}_t[i] = \mathbf{U}_t[i]/(\mathbf{S}_t[i] - \mathbf{U}_t[i])$ ;  $\mathbf{w}_t^i$  is a small image window (usually  $3 \times 3$ ) centered at pixel  $i$ .  $\mathbf{a}_t[i]$  and  $\mathbf{b}_t[i]$  are related to  $\alpha_t[j]$  by (4); and,  $\epsilon$  is a small constant used for numerical stability. The first part of

the cost function is identical to the cost function  $J_S$  introduced in [17]. The term  $(\alpha_t[i] - \alpha_{t-1}[i])^2$  is the temporal smoothing term. By adding this term, the original segmentation method of a single image is able to accommodate a sequence of images with temporal consistency.

The goal is to solve for  $\alpha_t$ ,  $\mathbf{a}_t$ , and  $\mathbf{b}_t$  minimizing the cost function  $J_T(\alpha_t, \mathbf{a}_t, \mathbf{b}_t)$

$$\hat{\alpha}_t = \arg \min_{\alpha_t, \mathbf{a}_t, \mathbf{b}_t} J_T(\alpha_t, \mathbf{a}_t, \mathbf{b}_t). \quad (8)$$

The minimization problem with respect to  $\mathbf{a}_t$  and  $\mathbf{b}_t$  is similar to the discussion in [17] and, thus,  $\mathbf{a}_t$  and  $\mathbf{b}_t$  can be obtained explicitly in terms of  $\alpha_t$ . Therefore, (8) can be recast as a quadratic optimization problem with respect to  $\alpha_t$  alone

$$\hat{\alpha}_t = \arg \min_{\alpha_t} (\alpha_t^T \mathbf{L}_t \alpha_t + (\alpha_t - \alpha_{t-1})^T (\alpha_t - \alpha_{t-1})) \quad (9)$$

where  $\alpha_t$  is a  $M \times 1$  vector ( $M$  is the number of pixels in the current ROI) and  $\mathbf{L}_t$  is a  $M \times M$  matting Laplacian matrix with its  $(i, j)$ th element given by

$$\sum_{k: (i, j) \in \mathbf{w}_t^k} \left( \delta[i, j] - \frac{1}{|\mathbf{w}_t^k|} \left( 1 + \frac{(\mathcal{F}_t[i] - \mu_k)(\mathcal{F}_t[j] - \mu_k)}{\frac{\epsilon}{|\mathbf{w}_t^k|} + \sigma_k^2} \right) \right).$$

Here,  $\delta[i, j]$  is the Kronecker delta,  $\mu_k$  and  $\sigma_k^2$  are the mean and variance of the intensities in the window  $\mathbf{w}_t^k$  centered at pixel  $k$ , respectively, and  $|\mathbf{w}_t^k|$  is the number of pixels in this window. The  $(i, j)$ th element of matrix  $\mathbf{L}_t$  measures the similarity between pixels  $i$  and  $j$ .

The  $\alpha_t$  value for the stable and unstable seeds is 1 and 0, respectively. Given the seed values at time  $t$  and the Matte values  $\alpha_{t-1}$  at the previous time step, the cost function (9) becomes

$$\begin{aligned} \hat{\alpha}_t = \arg \min_{\alpha_t} (\alpha_t^T \mathbf{L}_t \alpha_t + \mathcal{M}(\alpha_t - \mathbf{b}_t)^T \mathbf{D}_t (\alpha_t - \mathbf{d}_t) \\ + (\alpha_t - \alpha_{t-1})^T (\alpha_t - \alpha_{t-1})) \end{aligned} \quad (10)$$

where  $\mathcal{M}$  is some large number,  $\mathbf{D}_t$  is a diagonal matrix whose diagonal elements are 1 for the seeds, and 0 for all other pixels.  $\mathbf{d}_t$  is the vector containing the prespecified  $\alpha_t$  values for the seeds, and 0 for all other pixels. Since the aforementioned function is quadratic in  $\alpha_t$ , the global minimum  $\hat{\alpha}_t$  can be found by differentiating (10) and setting the derivatives to zero. This amounts to solving the following sparse linear system:

$$\begin{bmatrix} (\mathbf{L}_t + \mathcal{M}\mathbf{D}_t) \\ \mathbf{I}\alpha_{t-1} \end{bmatrix} \alpha_t = \begin{bmatrix} \mathcal{M}\mathbf{d}_t \\ \alpha_{t-1} \end{bmatrix} \quad (11)$$

where  $\mathbf{I}$  is the identity matrix and  $\alpha_{t-1}$  is the  $\alpha$  values at time  $t-1$ . We used bandwidth sparse matrix storage format and iterative generalized minimal residual method linear equation solver to solve the equation. The solution of (11) provides values in  $[0, 1]$ , where each value indicates the stability probability of the corresponding pixel in the template.

### C. Update of the Template

The more unstable the pixels are, the more aggressive updating they need. The estimated Matte values indicate the necessary degree of updating for each pixel. More precisely, the pixel of

the updated template, at time  $t$ , will arise as a weighted sum of the previous template  $\mathbf{T}_{t-1}[i]$ , and the ROI pixel  $\mathcal{F}_t[i]$  from the current frame; the weight  $\alpha_t[i]$  is the Matte value (note that we have dropped the  $\hat{\cdot}$  that denotes optimality, for notational simplicity)

$$\mathbf{T}_t[i] = \alpha_t[i]\mathbf{T}_{t-1}[i] + (1 - \alpha_t[i])\mathcal{F}_t[i]. \quad (12)$$

We can deduce from (12) that for a stable seed, the template value will not change (since  $\alpha_t[i] = 1$ ), while for an unstable seed, the template value will update to the corresponding pixel in the current ROI (since  $\alpha_t[i] = 0$ ). Given the computed Matte, the new template will not only update the unstable seeds and reserve the stable seeds, but will also proportionally update their surrounding pixels based on the Matte values.

The new template resembles the newest version of the subject's regional appearance, so that it stays relevant and representative. At the same time, it reserves the stable pixels of the previous template preventing the tracker from drifting.

### III. EXPERIMENTAL DESIGN, RESULTS, AND ANALYSIS

For the purpose of testing the STM template update method in the context of particle filter tracking, we used 25 thermal clips from 24 subjects. The clips were generated as part of sleep studies (subjects identified as Sxx and Lxx) [8], a stress study related to mock-crime interrogation (subjects identified as Dxxx<sub>IS</sub>) [9], and a stress study related to inanimate laparoscopic surgical training (subjects identified as Dxxx<sub>SS</sub>), as per the approval of the appropriate institutional review boards. The set included clips that had at minimum  $\sim 6500$  and at maximum  $\sim 49500$  frames. At the recoding speed of 30 frames/s, these clips ranged from  $\sim 3.5$  to  $\sim 27.5$  min in duration. From each thermal clip, challenging segments that featured significant positional and/or physiological change were selected for facial tracking; in some cases, these segments were as long as  $\sim 8000$  frames ( $\sim 5$  min). The targeted facial areas included the nostrils, where vital physiological function is resident, or the periorbital, supraorbital, or maxillary regions where sympathetic activation is manifested.

The STM particle filter tracker was compared with the on-line appearance model (OAM) tracker reported in [13] and the zero-one particle filter tracker reported in [16]. The trackers optimized three state variables, which served as ROI descriptors. These were  $(x, y)$  for translation and  $\phi$  for rotation on the image plane. The templates in all three trackers were formed out of normalized thermal values. All three tracking methods achieved real-time ( $>25$  frames/s) performance on a PentiumIV 4-core computer, with 4 GB memory. A short commentary about each method and the rationale for its inclusion in the benchmarking set is given in the following.

- 1) STM method: This is the method proposed in this paper that combines the agility of the particle filter framework with the sophistication of a spatiotemporal smoothing template.
- 2) OAM method: This is an advanced template method applied in visual facial tracking [13]. Thus, it can serve as a representative of noteworthy approaches from the relevant visual imaging literature. The method's statistical tem-

plate is formulated as a mixture of three components [13], namely a stable component ( $S$ ), a wandering component ( $W$ ), and an outlier ( $L$ ) component. The stable component captures the portion that is stable over time. It follows a normal distribution, the mean and variance of which are updated at every time step. The wandering component represents sudden appearance change. The outlier component is for short time occlusions. The method is probabilistic in nature but has no assumptions about spatial and temporal dependence.

- 3) Zero-One method: This is a method that combines the agility of the particle filter framework with the simplicity of a deterministic template [16]. It was recently used for facial tracking in thermal infrared. Therefore, it demarcates progress in the particular domain and can demonstrate the clear benefit of using more sophisticated templating to drive the particle filter loop.

The particle filter mechanisms of the Zero-One and the STM methods featured identical parameterizations. For every subject, all three trackers were tasked to track a selected facial tissue (ROI) from the exact same initial frame.

#### A. Qualitative Results

In thermal facial imaging, there are two major factors that affect the tracker's performance: head motion and physiological changes. The first factor alters the ROI location, while the second factor affects the pixel values within the ROI. For this reason, our dataset features subjects that exhibited large/small changes in the position or/and the physiology of the ROI. Accordingly, we split the dataset into three groups reflective of three distinct operational scenarios.

- 1) *Scenario 1*: Subjects that exhibited large changes in position and small changes in physiology.
- 2) *Scenario 2*: Subjects that exhibited small changes in position and large changes in physiology.
- 3) *Scenario 3*: Subjects that exhibited large changes both in position and physiology.

Fig. 4 shows a panorama of all subjects in the dataset categorized per operational scenario. Two characteristic thermal shots are shown for each subject with the STM tracker reliably tracking a facial tissue of interest. This figure gives a visual insight to the diversity of experimental circumstances that STM can negotiate.

Fig. 5 shows comparative tracker performance for a case representative of Scenario 3, where large positional and physiological changes occur. The signals represent the evolution of the translational and rotational errors of STM, OAM, and zero-one regarding the tracking of subject's D009<sub>IS</sub> nasal ROI. STM performs flawlessly in terms of translational accuracy and exhibits only small rotational errors in a short interval. OAM exhibits moderate translational errors and large rotational errors for extended periods of time. Zero-one maintains translational accuracy but exhibits significant rotational inaccuracy. This figure gives a dynamic sense of tracker performance, associating error numbers with visual impressions.

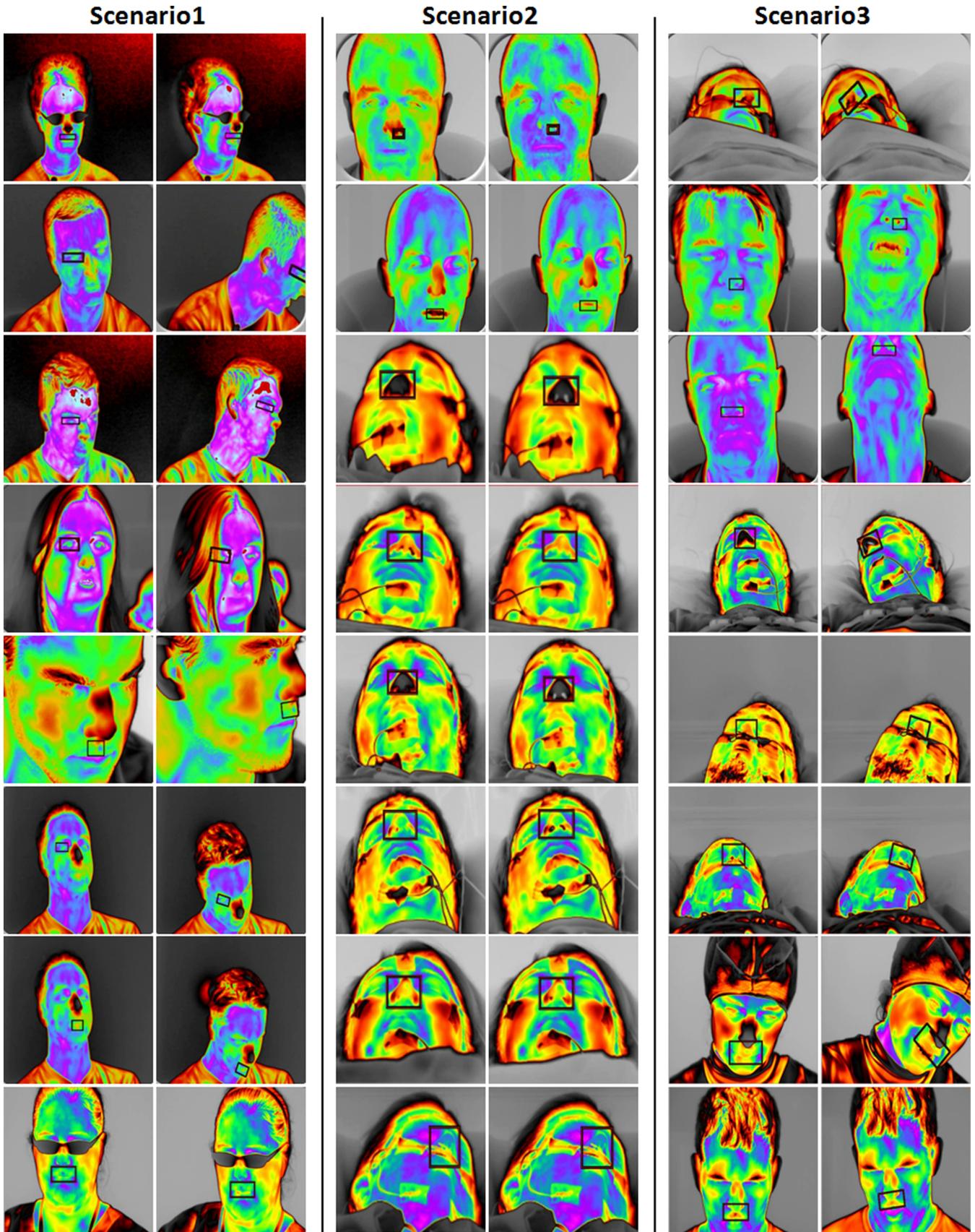


Fig. 4. Panorama of characteristic thermal shots (two per subject) from 24 out of the 25 thermal clips in the dataset, categorized per operational scenario.

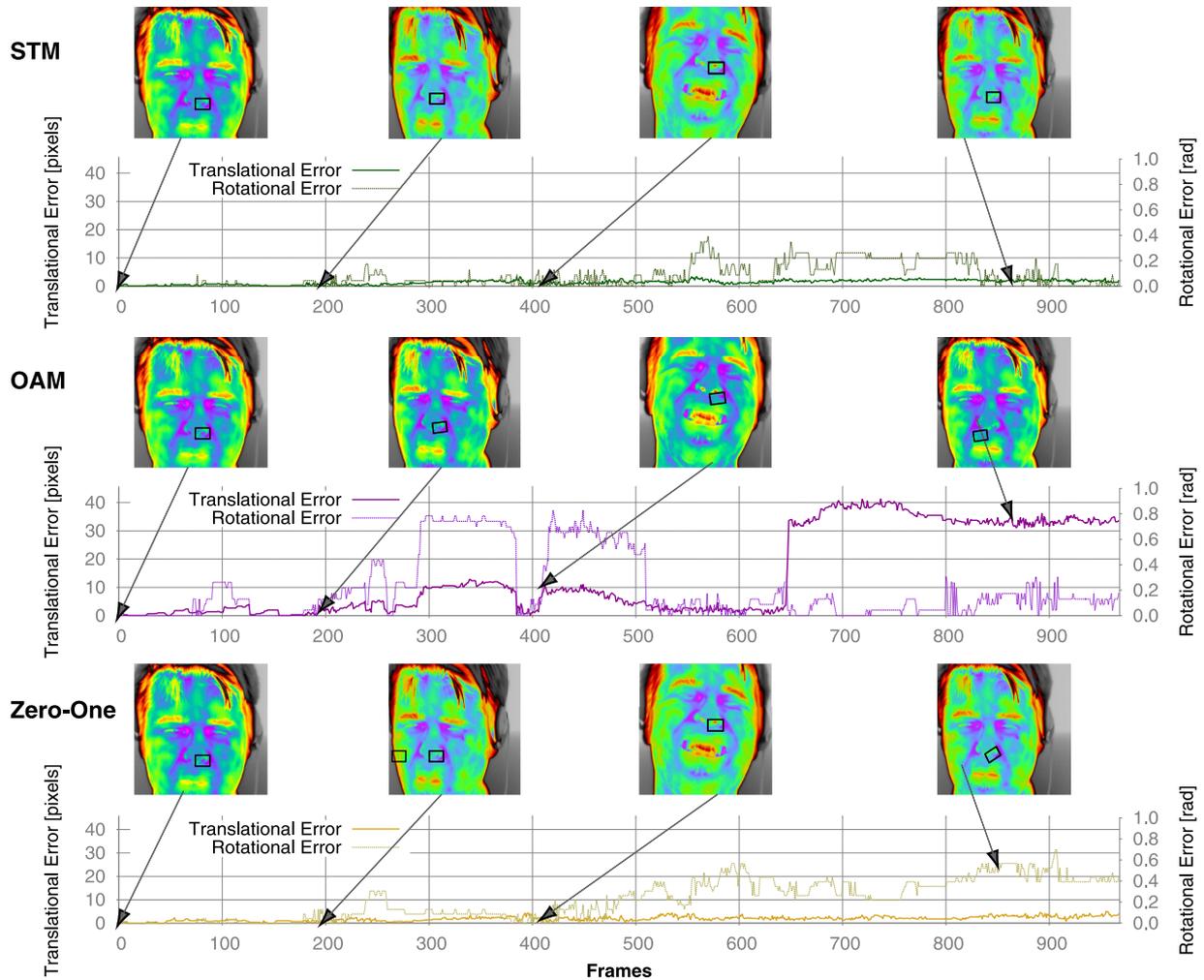


Fig. 5. Evolution of translational and rotational tracking errors for STM, OAM, and zero-one for subject D009<sub>1S</sub>. Thermal snapshots are indexed to representative points in the error signals to enable association with the nasal tracker's (black box) position in the actual runs.

## B. Quantitative Results

To quantify tracker performance, we need to compare the tracker's ROI with the ground-truth ROI throughout the time line. In medical imaging, the ground-truth data are usually obtained by manually segmenting the ROI in each frame. With thousands of frames in the dataset, however, manual ground truthing was not practical. For this reason, we adopted a different strategy. We used each of the three trackers to generate tracking results. We examined the results and where each tracker appeared to have failed, we manually repositioned the tracker and reinitiated tracking from that point onward, to correct the error. At the end, we formed ground-truth trackers as the means of the individual corrected trackers.

Tracking performance correlates to the Euclidean distance and angular difference between the ground-truth ROI and the ROI that each of the three competing methods produces. The smaller the Euclidean distance and angular difference are, the better.

The 25 clips of the dataset when partitioned according to Scenario 1, Scenario 2, and Scenario 3 provide 8, 9, and 8 clips, respectively. Fig. 6(a) shows a graphical representation of the distribution of translational (Euclidean) errors per tracking method and operational scenario. As the plot indicates, the STM approach outperforms the other two template update strategies in all scenarios. Both the OAM and Zero-One methods appear to have particular difficulties with large combined positional and physiological changes (Scenario 3). Fig. 6(b) shows a graphical representation of the distribution of rotational errors per tracking method and operational scenario. STM still outperforms the other two methods across the spectrum, but its relative error magnitude increased with respect to the translational error.

To statistically verify these indications, a series of hypotheses tests (at 0.05 level of significance) were performed to check how the means of the translational and rotational error distributions differ between methods for each of the scenarios in the database. Let us denote by  $\mu_{TS}$ ,  $\mu_{TO}$ , and  $\mu_{TZ}$  the means

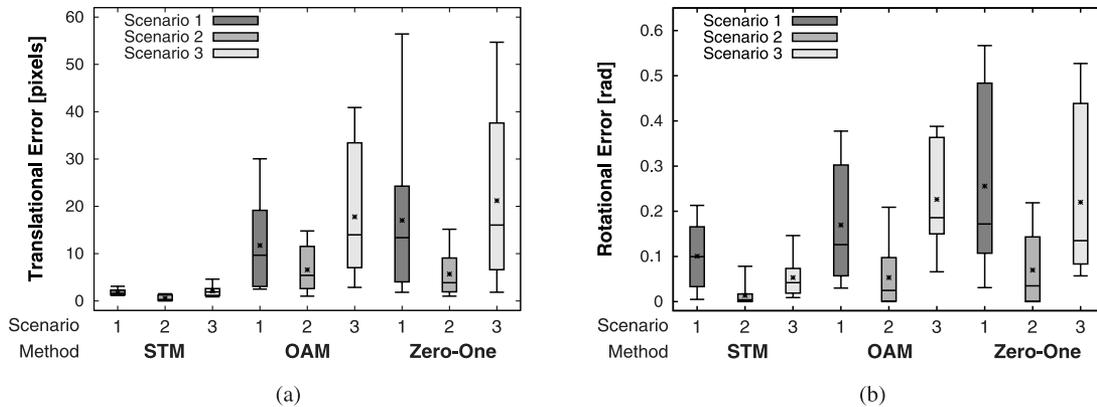


Fig. 6. (a) Boxplots of the translational (Euclidean) error distributions per tracking method and operational scenario for all subjects. (b) Boxplots of the rotational error distributions per tracking method and operational scenario for all subjects. (Scenario 1: large position and small physiology changes, Scenario 2: small position and large physiology changes, Scenario 3: large position and large physiology changes).

of the translational error distributions for the STM, OAM, and zero-one strategies, respectively. Accordingly, let us denote by  $\mu_{RS}$ ,  $\mu_{RO}$ , and  $\mu_{RZ}$  the means of the rotational error distributions for the STM, OAM, and zero-one strategies, respectively. The only instances where the tests fail to reject the null hypotheses are for subjects L02, L06, and L08, where the three competing trackers appear to perform on par. In all other cases, mean performance differs significantly among trackers, with STM clearly outperforming the other two. Note that L02, L06, and L08 are all sleep study cases, where tracking is not as challenging as in stress study cases. Fig. 7 gives a graphical representation of the mean positional and rotational tracking errors per subject for the three competing tracking methods; it visually correlates with the outcome of the statistical tests.

### C. Benefit of Temporal Smoothing

To specify the beneficial effect of temporal smoothing, a simulation was run where a thermal nasal region was translated only in the  $x$ -direction, while the  $y$ -direction and angle of rotation  $\phi$  were kept constant. The region featured semiperiodic fluctuation in temperature akin to the effect of breathing. This region was tracked first with a particle filter tracker driven by the classical Matte formula with spatial smoothing only. Then, it was tracked with the same particle filter tracker but driven by STM, i.e., the modified Matte formula with both spatial and temporal smoothing. The trajectory results in Fig. 8 demonstrate the fault oscillation introduced in the  $y$  and rotational dimensions by the classical Matte method.

### D. Sensitivity Analysis

In the Matte template strategy, the values of the nuisance parameters  $\lambda_1$  and  $\lambda_2$  ( $\lambda_1 < \lambda_2$ ) determine the most stable and unstable pixels in the current frame, which are used as seeds [see (5)]. In this study, the following parameterization was used:  $\lambda_1 = 5$  and  $\lambda_2 = 20$ . Note that the temperature values were normalized in the 0–255 range. An experiment was per-

TABLE I  
MEAN TRANSLATIONAL ERROR FOR VARIOUS CHOICES OF  $(\lambda_1, \lambda_2)$

$(\lambda_1, \lambda_2)$	D011 <sub>IS</sub> -nasal	D225 <sub>IS</sub> -periorbital	D016 <sub>IS</sub> -nasal
(2, 15)	1.614	2.815	1.162
(2, 20)	1.605	2.745	1.177
(2, 25)	1.605	2.718	1.171
(5, 15)	1.606	2.712	1.174
(5, 20)	1.515	2.947	1.088
(5, 25)	1.611	2.643	1.178
(8, 15)	1.613	3.245	1.144
(8, 20)	1.612	2.773	1.170
(8, 25)	1.599	2.707	1.163

formed to test the sensitivity of the STM method with respect to the  $\lambda$  parameters. Specifically, a representative clip from each of the three operational scenarios (Scenario 1: D011<sub>IS</sub>-nasal, Scenario 2: D225<sub>IS</sub>-periorbital, and Scenario 3: D016<sub>IS</sub>-nasal) was selected. The tracker of the STM template update method was run on these clips multiple times, varying at each iteration the values of  $\lambda_1$  and/or  $\lambda_2$ . Values for the tuple  $(\lambda_1, \lambda_2)$  were drawn from the sets  $\lambda_1 \in \{2, 5, 8\}$  and  $\lambda_2 \in \{15, 20, 25\}$ , providing nine different pairs. Thus, nine runs for every selected subject were produced, each one providing an error distribution of the Euclidean distance and rotational difference between the tracked and ground-truth ROIs. Tables I and II give the mean translational and rotational errors for each of the three clips and every case of the  $(\lambda_1, \lambda_2)$  parameter choices, respectively.

The mean distance value of the tracker’s ROI from the ground-truth position appears to be rather stable for a wide range of parameter choices, indicating that the method is not very sensitive. The sensitivity increases (and the performance deteriorates) when the values for the  $\lambda$  parameters get close. In this case, the method identifies many of the ROI pixels as seeds and starts losing its probabilistic (smoothing) advantage. An example is the case of the pair (8,15) in Table I that features the highest mean errors in all selected subjects.

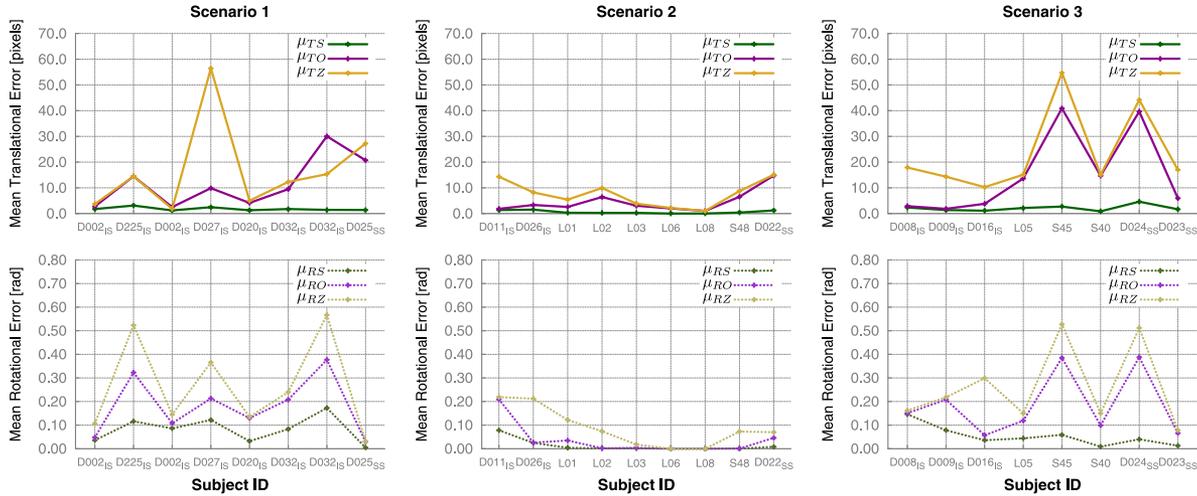


Fig. 7. Mean translational and rotational errors per subject for each competing tracking method grouped by operational scenario. The STM method consistently yields the smallest mean errors.

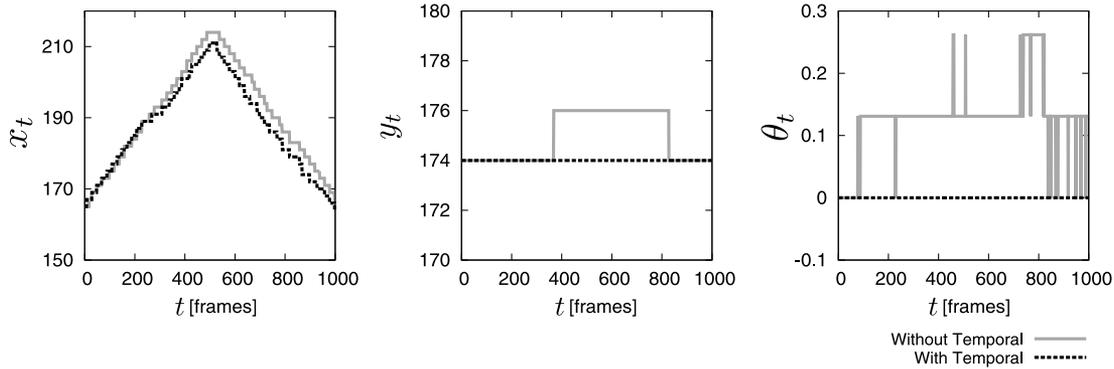


Fig. 8. Comparative nasal ROI trajectories (decomposed into the  $x$ ,  $y$ , and  $\phi$  dimensions) in a controlled simulation experiment. Light solid trajectories were produced by a particle filter tracker operating on a legacy Matte with spatial smoothing only. Dark dotted trajectories were produced by a particle filter tracker operating on a Matte with both spatial and temporal smoothing.

TABLE II  
MEAN ROTATIONAL ERROR FOR VARIOUS CHOICES OF  $(\lambda_1, \lambda_2)$

$(\lambda_1, \lambda_2)$	D011 <sub>IS</sub> -nasal	D225 <sub>IS</sub> -periorbital	D016 <sub>IS</sub> -nasal
(2, 15)	0.0028	0.0812	0.0370
(2, 20)	0.0028	0.0737	0.0383
(2, 25)	0.0029	0.0780	0.0373
(5, 15)	0.0028	0.0703	0.0380
(5, 20)	0.0028	0.1087	0.0361
(5, 25)	0.0028	0.0722	0.0379
(8, 15)	0.0028	0.1084	0.0376
(8, 20)	0.0028	0.0720	0.0382
(8, 25)	0.0028	0.0722	0.0370

IV. CONCLUSION

This paper presents a new probabilistic template update method that when drives a particle filter tracker is capable of producing sophisticated tracking behavior in thermal facial imaging. Specifically, the method can cope with both large positional and physiological changes, something that other methods from the thermal or visual domain fail to do. The power of the method stems from the spatial and temporal smoothness components of the template that capture well natural thermophysiological characteristics. The new approach was tested on a dataset con-

sisting of 25 thermal clips, thousands of frames each, featuring a variety of conditions that naturally occur in practice. The method promises improved performance in a number of biomedical applications, where unobtrusive physiological measurements on the face are preferred (e.g., sleep and stress studies).

Note that in the case the subject’s head exhibits frequent and significant out of plane rotation and movement, the resulting physiological signal will not be accurate as this is outside the tracker’s operational scenario. Fortunately, for the targeted applications, this is rarely the case. This is obvious for sleep studies. It is also true for stress studies, where the subject stays put, maintaining (by design) directional attention to the stimulus or the interviewer.

V. ACKNOWLEDGMENT

Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the funding agencies.

REFERENCES

[1] C. Eveland, D. Socolinsky, and L. Wolff, “Tracking human faces in infrared video,” *Imag. Vis. Comput.*, vol. 21, pp. 579–590, Aug. 2003.

- [2] S. Kong, J. Heo, B. Abidi, J. Paik, and M. Abidi, "Recent advances in visual and infrared face recognition—a review," *Comput. Vis. Imag. Understand.*, vol. 97, pp. 103–135, 2005.
- [3] N. Sun, M. Garbey, A. Merla, and I. Pavlidis, "Imaging the cardiovascular pulse," in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Patt. Recognit.*, San Diego, CA, 2005, vol. 2, pp. 416–421.
- [4] N. Sun, I. Pavlidis, M. Garbey, and J. Fei, "Harvesting the thermal cardiac pulse signal," in *Medical Image Computing and Computer-Assisted Intervention*, (Lecture Notes in Computer Science Series). vol. 4191, New York: Springer-Verlag, 2006, pp. 569–576.
- [5] R. Murthy and I. Pavlidis, "Noncontact measurement of breathing function," *IEEE Eng. Med. Biol. Mag.*, vol. 25, no. 3, pp. 57–67, May/Jun. 2006.
- [6] S. Chekmenev, A. Farag, and E. Essock, "Thermal imaging of the superficial temporal artery: An arterial pulse recovery model," in *Proc. IEEE Conf. Comput. Vis. Patt. Recognit.*, Minneapolis, MN, Jun. 17–22, 2007, pp. 1–6.
- [7] I. Pavlidis, J. Dowdall, N. Sun, C. Puri, J. Fei, and M. Garbey, "Interacting with human physiology," *Comput. Vis. Imag. Understand.*, vol. 108, no. 1–2, pp. 150–170, Oct./Nov. 2007.
- [8] J. Murthy, S. Faiz, J. Fei, I. Pavlidis, A. Abeulhagia, and R. Castriota, "Remote infrared imaging: A novel non-contact method to monitor airflow during polysomnography," in *Proc. Chest Meet. Abstr.*, Chicago, IL, Oct. 20–25, 2007, vol. 132, no. 4, p. 464.
- [9] P. Tsiamyrtzis, J. Dowdall, D. Shastri, I. Pavlidis, M. Frank, and P. Ekman, "Imaging facial physiology for the detection of deceit," *Int. J. Comput. Vis.*, vol. 71, no. 2, pp. 197–214, Oct. 2006.
- [10] B. Faser and J. Luetin, "Automatic facial expression analysis: A survey," *Pattern Recognit.*, vol. 36, no. 1, pp. 259–275, 2003.
- [11] M. Bureau and F. Sérès, "Comparison of two in-laboratory titration methods to determine effective pressure levels in patients with obstructive sleep apnea," *Thorax*, vol. 55, pp. 741–745, 2000.
- [12] L. Matthews and T. Ishikawa, "The template update problem," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 26, no. 6, pp. 810–815, Jun. 2004.
- [13] A. Jepson, D. Fleet, and T. E.-Maraghi, "Robust online appearance models for visual tracking," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 25, no. 10, pp. 415–422, Oct. 2003.
- [14] S. Zhou, R. Chellapa, and B. Moghaddam, "Visual tracking and recognition using appearance-adaptive models in particle filters," *IEEE Trans. Imag. Process.*, vol. 13, no. 11, pp. 1491–1506, Nov. 2004.
- [15] J. Dowdall, I. Pavlidis, and P. Tsiamyrtzis, "Coalitional tracking," *Comput. Vis. Imag. Understand.*, vol. 106, no. 2–3, pp. 205–219, May 2007.
- [16] J. Dowdall, "Tracking tissue in thermal infrared video," Ph.D. dissertation, Dept. Comput. Sci., Univ. Houston, Houston, TX, 2006.
- [17] A. Levin, D. Lischinski, and Y. Weiss, "A closed form solution to natural image matting," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 30, no. 2, pp. 228–242, Feb. 2008.
- [18] Y. Zhou, P. Tsiamyrtzis, and I. Pavlidis, "Tissue tracking in thermo-physiological imagery through spatio-temporal smoothing," in *Medical Image Computing and Computer-Assisted Intervention*, (Lecture Notes in Computer Science Series). vol. 5762, New York: Springer-Verlag, 2009, pp. 1092–1099.
- [19] Y. Li, J. Sun, C. Tang, and H. Shum, "Lazy snapping," *ACM Spec. Interest Group Comput. Graph. Interactive Tech.*, vol. 23, no. 3, pp. 303–308, Aug. 2004.
- [20] C. Rother, V. Kolmogorov, and A. Blake, "GrabCut—interactive foreground extraction using iterated graph cuts," *ACM Spec. Interest Group Comput. Graph. Interactive Tech.*, vol. 23, no. 3, pp. 309–314, Aug. 2004.
- [21] J. Shi and J. Malik, "Normalized cuts and image segmentation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 22, no. 8, pp. 888–905, Aug. 2000.



**Yan Zhou** received the B.S. degree with a major in electrical engineering from Fudan University, Shanghai, China, in 2005, and the M.S. and Ph.D. degrees from the Department of Computer Science, University of Houston, Houston, TX, in 2009 and 2011, respectively.

She was a Research Intern at the Motorola R&D Center, Shanghai, China, in 2006, and at Siemens Healthcare, Malvern, PA, in 2010. In 2012, she joined Elekta Inc., Maryland Heights, MO, as a Research Scientist. She is currently involved in research and

development of autosegmentation tools in treatment planning systems for radiation therapy. She has published in the areas of medical imaging and the physiological basis of human behavior.



**Panagiotis Tsiamyrtzis** received the B.S. degree in mathematics from the Aristotle University of Thessaloniki, Thessaloniki, Greece, and the M.S. and Ph.D. degrees in statistics from the University of Minnesota, Minneapolis.

In Fall 2000, he was a visiting faculty member in the School of Statistics, University of Minnesota. In 2004, he joined the Department of Statistics, Athens University of Economics and Business, Athens, Greece, where he is currently an Assistant Professor. His research interests include statistical aspects of computer vision problems, statistical process control, and applications of Bayesian statistics.

Dr. Tsiamyrtzis was the recipient of the Best Student Paper Award and Best Contributed Paper in the 2000 Joint Statistical Meetings of the Risk Analysis Section of the American Statistical Association. He also received the Best Talk Award in the 2007 Annual Conference of the European Network for Business and Industrial Statistics.

Dr. Tsiamyrtzis was the recipient of the Best Student Paper Award and Best Contributed Paper in the 2000 Joint Statistical Meetings of the Risk Analysis Section of the American Statistical Association. He also received the Best Talk Award in the 2007 Annual Conference of the European Network for Business and Industrial Statistics.



**Peggy Lindner** received the B.S. and M.S. degrees in geotechnology/mining from the University of Freiberg, Freiberg, Germany, in 1996 and 2000, respectively, and the Ph.D. degree in mechanical engineering from the University of Stuttgart, Stuttgart, Germany, in 2007.

Her interdisciplinary background has shaped her primary research interests in high-performance computing and grid computing. She is a Research Assistant Professor in the Department of Computer Science, University of Houston, Houston, TX.

Her current research interests include stress research and data management/representation of large physiological datasets.



**Ilya Timofeyev** received the B.S. and Ph.D. degrees in mathematics from the Rensselaer Polytechnic Institute, Troy, NY, in 1994 and 1998, respectively.

He spent four years as a Postdoctoral Fellow at the New York University. He is currently an Associate Professor in the Department of Mathematics, University of Houston, Houston, TX. His research interest includes applied mathematics.



**Ioannis Pavlidis** (S'85–M'87–SM'00) received the B.E. degree in electrical engineering from Democritus University of Thrace, Xanthi, Greece, in 1987, the first M.S. degree in robotics from the University of London, London, U.K., in 1989, and the second M.S. and Ph.D. degrees in computer science from the University of Minnesota, Minneapolis, in 1995 and 1996, respectively.

He is currently the Eckhard-Pfeiffer Professor of computer science and the Director of the Computational Physiology Laboratory, University of Houston,

Houston, TX. He is the author of many scientific articles on computational physiology and affective computing. He is well known for his research on facial signs of stress, which first appeared in *Nature* and *Lancet*.