



HAL
open science

Human Motion Tracking using a Color-Based Particle Filter Driven by Optical Flow

Tony Tung, Takashi Matsuyama

► **To cite this version:**

Tony Tung, Takashi Matsuyama. Human Motion Tracking using a Color-Based Particle Filter Driven by Optical Flow. The 1st International Workshop on Machine Learning for Vision-based Motion Analysis - MLVMA'08, Oct 2008, Marseille, France. inria-00326720

HAL Id: inria-00326720

<https://inria.hal.science/inria-00326720>

Submitted on 5 Oct 2008

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Human Motion Tracking using a Color-Based Particle Filter Driven by Optical Flow

Tony Tung and Takashi Matsuyama

Graduate School of Informatics, Kyoto University, Japan
{tung,tm}@vision.kuee.kyoto-u.ac.jp

Abstract. In this paper, we present a new formulation for the problem of human motion tracking in video. Tracking is still a challenging problem when strong appearance changes occur as in videos of human in motion. Most trackers rely on a predefined template or on a training dataset to achieve detection and tracking. Therefore they are not efficient to track objects which appearance is not known in advance. A solution is to use an online method that updates iteratively a subspace of reference target models. In addition, we propose to integrate color and motion cues in a particle filter framework to track human body parts. The algorithm process consists in two modes, switching between detection and tracking. Detection steps involve trained classifiers to update estimated positions of the tracking windows, whereas tracking steps rely on an adaptative color-based particle filter coupled with optical flow estimations. The Earth Mover distance is used to compare color models in a global fashion, and constraints on flow features avoid drifting effects. The proposed method has revealed its efficiency to track body parts in motion and can cope with full appearance changes. Experiments were lead on challenging real world videos with poorly textured models and non-linear motions.

1 Introduction

Human motion tracking is a common requirement for many real world applications, such as video surveillance, games, cultural and medical applications (e.g. for motion and behavior study). The literature has provided successful algorithms to detect and track objects of a predefined class in image streams or videos. Simple object can be detected and tracked using various image features such as color regions, edges, contours, or texture. In the other hand, complex objects such as human faces require more sophisticated features to handle the multiple possible instances of the object class. For this purpose, statistical methods are a good alternative. First, a statistical model (or classifier) learns different patterns related to the object of interest (e.g. different views of human faces), including good and bad samples. And then the system is able to estimate whether a region contains an object of interest or not. This kind of approach has become very popular (e.g. the face detector of [1] is well-known for its efficiency). The main drawbacks are that it requires prior knowledge on the object class which

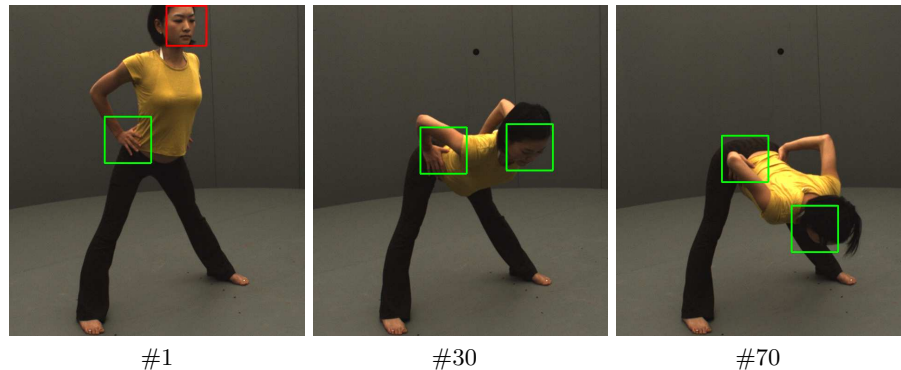


Fig. 1. Body part tracking with color-based particle filter driven by optical flow. The proposed approach is robust to strong occlusion and full appearance change. Detected regions are denoted by red squares, and tracked regions by green squares.

is usually a huge training dataset, and then it is somehow constrained to it. As a matter of fact, most of the tracking methods were not designed to keep the track of an object which appearance could strongly change. If there is no a priori knowledge on its multiple possible appearances, then the detection fails and the track is lost. Hence, tracking a head which turns completely, or tracking a hand in action remain challenging problems, as appearance changes occur quite frequently for human body parts in motion. Therefore we propose a new formulation dedicated to the problem of appearance changes for object tracking in video. Our approach integrates color cues and motion cues to establish a robust tracking. As well, an online iterative process updates a subspace of reference templates so that the tracking system remains robust to occlusions. The method workflow contains two modes, switching between detection and tracking. Detection steps involve trained classifiers to update estimated positions of the tracking windows. In particular, we use the cascade of boosted classifiers of Haar-like features by [1] to perform head detection. Other body parts can be either detected using this technique with ad-hoc training samples, or chosen by users at the initialization step, or as well can be deduced based on prior knowledge on human shape features and constraints. The tracking steps rely on an adaptive color-based particle filter [2] coupled with optical flow estimations [3, 4]. The Earth Mover distance [5] has been chosen to compare color models due to its robustness to small color variations. Drift effects inherent to adaptive tracking methods are handled using optical flow estimations (motion features). Our experiments show the accuracy and robustness of the proposed method on challenging video sequences of human in motion. For example, videos of yoga performances (stretching exercises at various speed) with poorly textured models and non-linear motions were used for testing (cf. Fig. 1).

The rest of the paper is organized as follows. The next section discusses work related to the techniques presented in this paper. Section 3 presents an

overview of the algorithm (initialization step and workflow). Section 4 describes the tracking process based on our color-based particle filter driven by optical flow. Section 5 presents experimental results. Section 6 concludes with a discussion on our contributions.

2 Related work

In the last decade, acquisition devices have become even more accurate and accessible for non-expert users. This has led to a rapid growth of various imaging applications. In particular, the scientific community has shown a real interest to human body part detection and tracking. For example, face detection in images is nowadays a popular and well explored topic [1, 6, 7]. In [1], Viola and Jones proposed a cascade of boosted tree classifiers of Haar-like features. The classifier is first trained on positive and negative samples, and then the detection is performed by sliding a search window through candidate images and checking whether a region contains an object of interest or not. The technique is known to be fast and efficient, and can be tuned to detect any kind of object class if the classifier is trained on adapted samples.

Tracking in video is a popular field of research as well. Recognition from video is still challenging because frames are often of low quality, and details can be small (e.g. in video surveillance). Various approaches were proposed to track image features [3, 4, 8–10]. Lucas, Tomasi and Kanade [3, 4] first select the good features which are optimal for tracking, and then keep the tracks of these features in consecutive frames. The KLT feature tracker is often used for optical flow estimation to estimate the deformations between two frames. As a differential method, it assumes that the pixel intensity of objects is not significantly different between two frames.

Techniques based on prediction and correction as Kalman filter, and more recently particle filters have become widely used [2, 11–19]. Particle filters (or sequential Monte Carlo or Condensation) are Bayesian model estimation techniques based on simulation. The basic idea is to approximate a sequence of probability distributions using a large set of random samples (called particles). Then the particles are propagated through the frames based on importance sampling and resampling mechanisms. Usually, the particles converge rapidly to the distributions of interest. The algorithm allows robust tracking of objects in cluttered scene, and can handle non-linear motion models more complex than those commonly used in Kalman filters. The major differences between the different particle filter based approaches rely on the design of the sampling strategies, which make particles having higher probability mass in regions of interest.

In [20, 21, 16, 18, 19], linear dimension reduction methods (PCA, LDA) are used to extract feature vectors from the regions of interest. These approaches suit well for adaptative face tracking and can be formulated in the particle filtering framework as well [16, 18, 19]. Nevertheless they require a big training data set to be efficient [22], and still cannot cope with unpredicted change of appearance. In the other hand, color-based models of regions can capture larger

appearance variations [23, 24]. In [12], the authors integrate a color-based model tracker (as in the Meanshift technique [24]) within a particle filter framework. The model uses color histograms in the HSV space and the Bhattacharyya distance for color distribution comparisons. Nevertheless these methods usually fail to track objects in motion or have an increasing drift on long video sequences due to strong appearance changes or important lighting variations [25]. Indeed most algorithms assume that the model of the target object does not change significantly over time. To adapt the model to appearance changes and lighting variations, subspace of the target object features are extracted [21, 16, 18, 19]. In [18], a subspace of eigenvectors representing the target object is incrementally updated through the tracking process. Thus, offline learning step is not required and tracking of unknown objects is possible. Recently, [19] proposed to extend this approach with additional terms in the data likelihood definition. In particular, the drift error is handled using an additional dataset of images. However, these approaches are particularly tuned for face tracking, and still require training datasets for every different view of faces.

The core of our approach divides into two steps which are detection and tracking, as [13, 17]. Switching between the two modes allows to dynamically update the search window to an accurate position whenever the detection is positive. In this paper, we propose to run a color-based particle filter to achieve the tracking process. Our tracker uses a subspace of color models of regions of interest extracted from the previous frames, and relies on them to estimate the position of the object in the current frame. The subspace is iteratively updated through the video sequence, and dynamically updated by the detection process. The detection is performed by a cascade of boosted classifiers [1] and thus can be trained to detect any object class. We also introduce the Earth Mover distance to improve the robustness of tracking with lighting variations, and constraints based on optical flow estimations to cope with drift effects.

3 Algorithm overview

This section describes the algorithm workflow. The proposed approach combines two modes, switching between detection mode and tracking mode. The tracking process can be run independently if no detector is available for the class of the object of interest. Besides the tracking process, the detection improves the response accuracy online, and is used as well as initialization step. A subspace of color-based models is used to infer the object of interest location.

3.1 Initialization

The initialization step consists in defining the objects to track. In our framework, we focused on human body parts because of the wide range of possible applications. Basically, there are three straightforward ways to define the regions of interest:

1. Automatic: this can be achieved by a detection process using statistical machine learning method (e.g. the face detector of [1]).
2. Manual: regions of interest are defined by the user (e.g. by picking regions in the first frame). This allows to track any body part without having any prior knowledge.
3. Deduction: as the human body has self-constrained motions, its structure can be deduced using a priori knowledge and fuzzy rules. For example, a detected face gives some hints to deduce the torso position, etc.

In some of our experiments (cf. Sect. 5), we have combined the three approaches (e.g. head is automatically detected, torso is deduced, and hands are picked). Afterwards, the regions of interest are used as reference templates on which the tracker relies on to process the next frames.

3.2 Workflow

Assuming the initialization occurs at time t_0 , then for every frame at $t, t > t_0$, the tracker estimates the positions of M objects of interest $\{A_i\}_{i=1\dots M}$ based on the color-model subspace $\mathbf{S}_t^i = \{h_{t-k}^i, \dots, h_{t-1}^i\}$, where h_j^i denotes the color-model of A_i at time j , and k is the size of the subspaces (which in fact can be different for every object). Assuming a Bayesian framework (cf. Sect. 4), the hidden state \mathbf{x}_t^i corresponding to the estimated position of A_i at time t by the tracker, is inferred by \mathbf{S}_t^i and \mathbf{x}_{t-1}^i . We denote by \mathbf{y}_t^i the data corresponding to the detection of A_i at time t , and \mathbf{z}_t^i the response of the algorithm. Thus, if the detection of A_i at t is positive, then $\mathbf{z}_t^i = \mathbf{y}_t^i$, else $\mathbf{z}_t^i = \mathbf{x}_t^i$. Indeed if the detection of A_i at t is positive, then \mathbf{S}_{t+1}^i will be updated with the color model corresponding to \mathbf{y}_t^i . And if not, then \mathbf{S}_{t+1}^i will be updated with the color model corresponding to \mathbf{x}_t^i . The workflow is illustrated on Figure 2 with $M = 1$ and $k = 1$.

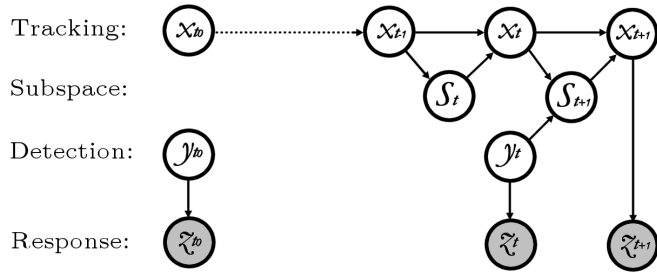


Fig. 2. Algorithm workflow. In this example, the detection state is positive at time t , and negative at $t + 1$. Hence, the algorithm response \mathbf{z}_t at time t is updated with \mathbf{y}_t , and \mathbf{z}_{t+1} is updated with the tracker estimation \mathbf{x}_{t+1} at time $t + 1$. \mathbf{S}_t denotes the subspace of color models used to infer \mathbf{x}_t . As the detection is positive at t , \mathbf{S}_{t+1} takes into account the color model corresponding to \mathbf{y}_t .

4 Particle filtering driven by optical flow

In this section we present our algorithm formulation based on color-based particle filtering [2, 12] and optical flow estimations [4]. We propose to use the Earth Mover Distance [5] to compare color models, and extracted motion features to improve tracking accuracy. Moreover our method updates iteratively a subspace of template models to handle appearance changes and partial occlusions.

4.1 Particle filtering

We denote by \mathbf{x}_t a target state at time t , \mathbf{z}_t the observation data at time t , and $\mathbf{Z}_t = \{\mathbf{z}_1, \dots, \mathbf{z}_t\}$ all the observations up to time t . Assuming a non-Gaussian state space model, the prior probability $p(\mathbf{x}_t|\mathbf{Z}_{t-1})$ at time t in a Markov process is defined as:

$$p(\mathbf{x}_t|\mathbf{Z}_{t-1}) = \int p(\mathbf{x}_t|\mathbf{x}_{t-1})p(\mathbf{x}_{t-1}|\mathbf{Z}_{t-1})d\mathbf{x}_{t-1} , \quad (1)$$

where $p(\mathbf{x}_t|\mathbf{x}_{t-1})$ is a state transition distribution, and $p(\mathbf{x}_{t-1}|\mathbf{Z}_{t-1})$ stands for a posterior probability at time $t-1$. The posterior probability which the tracking system aims to estimate at each time is defined as:

$$p(\mathbf{x}_t|\mathbf{Z}_t) \propto p(\mathbf{z}_t|\mathbf{x}_t)p(\mathbf{x}_t|\mathbf{Z}_{t-1}) , \quad (2)$$

where $p(\mathbf{z}_t|\mathbf{x}_t)$ is the data likelihood at time t . According to the particle filtering framework, the posterior $p(\mathbf{x}_t|\mathbf{Z}_t)$ is approximated by a Dirac measure on a finite set of P particles $\{\mathbf{x}_t^i\}_{i=1\dots P}$ following a sequential Monte Carlo framework [11]. Candidate particles are sampled by a proposal transition kernel $q(\tilde{\mathbf{x}}_t^i|\mathbf{x}_{t-1}^i, \mathbf{z}_{t-1})$. The new filtering distribution is then approximated by a new sample set of particles $\{\tilde{\mathbf{x}}_t^i\}_{i=1\dots P}$ having the importance weights $\{w_t^i\}_{i=1\dots P}$, where

$$w_t^i \propto \frac{p(\mathbf{z}_t|\tilde{\mathbf{x}}_t^i)p(\tilde{\mathbf{x}}_t^i|\mathbf{x}_{t-1}^i)}{q(\tilde{\mathbf{x}}_t^i|\mathbf{x}_{t-1}^i, \mathbf{z}_{t-1})} \quad \text{and} \quad \sum_{i=1}^P w_t^i = 1 . \quad (3)$$

The sample set $\{\mathbf{x}_t^i\}_{i=1\dots P}$ can then be obtained by resampling $\{\tilde{\mathbf{x}}_t^i\}_{i=1\dots P}$ with respect to $\{w_t^i\}_{i=1\dots P}$. By default, the Bootstrap filter is chosen as proposal distribution: $q(\tilde{\mathbf{x}}_t^i|\mathbf{x}_{t-1}^i, \mathbf{z}_{t-1}) = p(\tilde{\mathbf{x}}_t^i|\mathbf{x}_{t-1}^i)$. Hence the weights can be computed by evaluating the corresponding data likelihood. Finally, \mathbf{x}_t is estimated upon the Monte Carlo approximation of the expectation $\hat{\mathbf{x}}_t = \frac{1}{P} \sum_{i=1}^P \mathbf{x}_t^i$.

We denote by E , the overall energy function: $E = E_s + E_m + E_d$, where E_s is an energy related to color cues (cf. Sect. 4.2), E_m and E_d are energies related to motion features (cf. Sect. 4.4). E has lower values as the search window is close to the target object. Thus, to favor candidate regions which color distribution is similar to the reference model at time t , the data likelihood $p(\mathbf{z}_t|\mathbf{x}_t)$ is modeled as a Gaussian function:

$$p(\mathbf{z}_t|\tilde{\mathbf{x}}_t^i) \propto \exp\left(-\frac{E}{\sigma^2}\right) , \quad (4)$$

where σ is a scale factor, and therefore a small E returns a large weight.

4.2 Color-based model

The efficiency of color distributions to track color content of regions that match a reference color model has been demonstrated in [23, 24, 12]. They are represented by histograms to characterize the chromatic information of regions. Hence they are robust against non-rigidity and rotation. In addition, the Hue-Saturation-Value (HSV) color space has been chosen due to its low sensitivity to lighting condition. In our approach, color distributions are discretized into three histograms of N_h , N_s , and N_v bins for the hue, saturation, and value respectively.

Let α be h , s , or v , $\mathbf{q}_t(\mathbf{x}_t) = \frac{1}{3} \sum_{\alpha} \mathbf{q}_t^{\alpha}(\mathbf{x}_t)$, and $\mathbf{q}_t^{\alpha}(\mathbf{x}_t) = \{q_t^{\alpha}(i, \mathbf{x}_t)\}_{i=1 \dots N_{\alpha}}$. $\mathbf{q}_t^{\alpha}(\mathbf{x}_t)$ denotes the kernel density estimate of the color distribution in the candidate region $R(\mathbf{x}_t)$ of the state \mathbf{x}_t at time t , and is composed by:

$$q_t^{\alpha}(i, \mathbf{x}_t) = K_{\alpha} \sum_{\mathbf{u} \in R(\mathbf{x}_t)} \delta[h_{\alpha}(\mathbf{u}) - i] , \quad (5)$$

where K_{α} is a normalization constant so that $\sum_{i=1}^{N_{\alpha}} q_t^{\alpha}(i, \mathbf{x}_t) = 1$, h_{α} is a function assigning the pixel color at location \mathbf{u} to the corresponding histogram bin, and δ is the Kronecker delta function.

At time t , $\mathbf{q}_t(\mathbf{x}_t)$ is compared to a set of reference color model templates $\mathbf{S}_t = \{h_{t-k}, \dots, h_{t-1}\}$, where k is the number of templates. The templates are extracted iteratively from the detected regions at each frame. We recall that color model subspaces help to handle appearance changes and partial occlusions, and we define the energy function:

$$E_s[\mathbf{S}_t, \mathbf{q}_t(\mathbf{x}_t)] = \min_{h \in \mathbf{S}_t} (D^2[h, \mathbf{q}_t(\mathbf{x}_t)]) , \quad (6)$$

where D is a distance between color distributions (cf. Sect. 4.3).

4.3 Earth Mover distance

We propose to use the Earth Mover distance (EMD) [26, 5] to strengthen the property of invariance to lighting of the HSV color space. EMD allows to make global comparison of color distributions relying on a global optimization process. This method is more robust than approaches relying on histogram bin-to-bin distances that are more sensitive to quantization and small color changes. The distributions are represented by sets of weighted features called *signatures*. The EMD is then defined as the minimal amount of *work* needed to match a signature to another one. The notion of work relies on a metric (e.g. a distance) between two features. In our framework we use the L_1 norm as distance, and histogram bins as features.

Assuming two signatures to compare $P = \{(p_1, w_1), \dots, (p_m, w_m)\}$ and $Q = \{(q_1, u_1), \dots, (q_n, u_n)\}$, P having m components p_i with weight w_i , and Q having n components q_j with weight u_j . The global optimization process consists in finding the amount of data f_{ij} of a signature to be transported from the component i

to the component j that minimizes the work W :

$$W = \min_{f_{ij}} \left(\sum_{i=1}^m \sum_{j=1}^n d_{ij} f_{ij} \right) , \quad (7)$$

where d_{ij} is the distance between the components p_i and q_j assuming the following constraints:

$$\begin{aligned} f_{ij} &\geq 0 & 1 \leq i \leq m, 1 \leq j \leq n , \\ \sum_{j=1}^n f_{ij} &\leq w_i & 1 \leq i \leq m , \\ \sum_{i=1}^m f_{ij} &\leq u_j & 1 \leq j \leq n , \\ \sum_{i=1}^m \sum_{j=1}^n f_{ij} &= \min(\sum_{i=1}^m w_i, \sum_{j=1}^n u_j) . \end{aligned} \quad (8)$$

The first constraint allows only the displacements from P to Q . The two following constraints bound the amount of data transported by P , and the amount of data received by Q to their respective weights. The last constraint sets the maximal amount of data that can be displaced.

The EMD distance D between two signatures P and Q is then defined as:

$$D(P, Q) = \frac{W}{\mathcal{N}} = \frac{\sum_{i=1}^m \sum_{j=1}^n d_{ij} f_{ij}}{\sum_{i=1}^m \sum_{j=1}^n f_{ij}} , \quad (9)$$

where the normalization factor \mathcal{N} ensures a good balance when comparing signatures of different size (\mathcal{N} is the smallest sum of the signature weights). Note that EMD computation can be approximated in linear time with guaranteed error bounds [27].

4.4 Motion cues

We propose to use motion features to guide the search window through the tracking process. Motion features are extracted using the KLT feature tracker [3, 4]. The method detects feature windows and matches the similar ones between consecutive frames.

Assuming the set $\mathbf{Y}_{t-1} = \{y_{t-1}^j\}_{j=1\dots m}$ of m motion features detected in the neighborhood region of the state \mathbf{x}_{t-1} (cf. Sect. 4) at time $t-1$, and the set $\mathbf{Y}_t = \{y_t^j\}_{j=1\dots m}$ of matching features extracted at time t , then $(\mathbf{Y}_{t-1}, \mathbf{Y}_t) = \{(y_{t-1}^j, y_t^j)\}_{j=1\dots m}$ forms a set of m motion vectors (optical flow field) between the frames at time $t-1$ and t . As well, we denote by $\tilde{\mathbf{Y}}_t^i$ the set of features detected in the neighborhood region of the particle $\tilde{\mathbf{x}}_t^i$, and $\tilde{\mathbf{y}}_t$ the position of the search window estimated by optical flow as: $\tilde{\mathbf{y}}_t = \mathbf{x}_{t-1} + \text{median}(\{y_{t-1}^j - y_t^j\}_{j=1\dots m})$. Thus we define the following energy functions:

$$E_m(\tilde{\mathbf{x}}_t^i) = \alpha \cdot \|\tilde{\mathbf{x}}_t^i - \tilde{\mathbf{y}}_t\|_2 \quad \text{and} \quad E_d = \beta \cdot C(\tilde{\mathbf{Y}}_t^i, \mathbf{Y}_t) , \quad (10)$$

where α and β are two constant values, and C is the following function:

$$C(\tilde{\mathbf{Y}}_t^i, \mathbf{Y}_t) = 1 - \frac{\text{card}(\tilde{\mathbf{Y}}_t^i \cap \mathbf{Y}_t)}{\text{card}(\tilde{\mathbf{Y}}_t^i)} . \quad (11)$$

The data energy E_m aims to favor the particles located around the object target position estimated by optical flow, whereas E_d aims to prevent the drift effect. E_d works as a constraint which attracts the particles near the estimated search window (cf. Fig. 3). E_m and E_d are introduced in the overall energy formulation as described in Sect. 4.2.

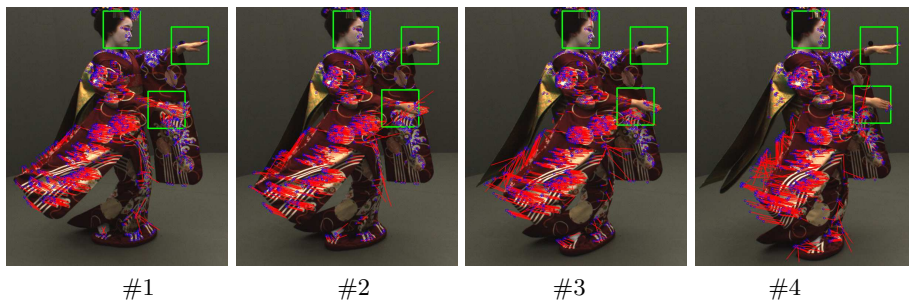


Fig. 3. Using optical flow to improve tracking. The combination of color cues and motion cues allows to perform robust tracking and prevent drift effects. The tracking of hands is efficient even with a changing background.

5 Experimental results

Our algorithm has been tested on various real video sequences. For example, we have tracked the body parts of a lady practicing yoga (head, hands, torso, and feet) in different video sequences and from different viewpoints. The model wears simple clothes with no additional features (cf. Fig. 1 and Fig. 5). As well, we have tested the tracker on a model wearing traditional Japanese clothes which are more much complex and contain a lot of features (cf. Fig. 3). In this paper, the video frame sizes are 640x480 and 720x576 pixels and were acquired at 25 fps. The algorithm was run on a Core2Duo 3.0GHz with 4GB RAM.

The following parameters were identical for all the experiments: we have used $N_h = 10$, $N_s = 10$ and $N_v = 10$ for the quantization of color models, $P = 200$ particles, $k = 5$ for the color model subspace size, and $\sigma^2 = 0.1$ as scale factor of the likelihood model. The constant values α and β weight the contribution of the motion cues, and are tuned regarding to the frame size. He have defined a square window size of 40 pixels to determine the regions of interest. The proposed formulation has shown promising results. The Figures 1 and 4 illustrate the robustness to appearance change, lighting variation and partial occlusion, thanks to the online update of the color-based model subspace combined with the Earth Mover distance and motion cues. For example, the system can track a head even if the face in no more visible (e.g. hidden by hair or due to changing viewpoint). Figure 3 illustrates an accurate tracking

with free-drift effect of a hand with a varying background under the guidance of optical flow as motion cues. Figure 5 illustrates the robustness of our approach in comparison to a color-based particle filter (Condensation) [12] that does not include our features. We show that the Condensation mixes regions having the same color shape and distribution whereas our tracker is not confused by the similar regions. This is due in particular to the addition of motion cues.

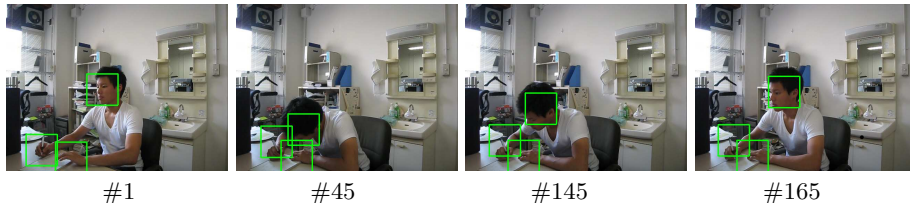


Fig. 4. Tracking with appearance change. The proposed approach integrates motion cues and a subspace of color models which is updated online through the video sequence. The system can track objects in motion with appearance change.

6 Conclusion

Human motion tracking in video is an attractive research field due to the numerous possible applications. The literature has provided powerful algorithms based on statistical methods especially dedicated to face detection and tracking. Nevertheless, it is still challenging to handle complex object classes such as human body parts which appearance changes occur quite frequently while in motion. In this paper, we propose to integrate color cues and motion cues in a tracking process relying on a particle filter framework. We have used the Earth Mover distance to compare color-based model distribution in the HSV color space in order to strengthen the invariance to lighting condition. Combined with an online iterative update of color-based model subspace, we have obtained robustness to partial occlusion. We have also proposed to integrate extracted motion features (optical flow) to handle strong appearance changes and prevent drift effect. In addition, our tracking process is run jointly with a detection process that dynamically updates the system response. Our new formulation has been tested on real videos, and results on different sequences were shown. For future work, we believe our approach can be easily extended to handle an online manifold learning process. This would improve both detection and tracking modes.

Acknowledgments

The authors would like to thank Dr. Hiroaki Kawashima and Dr. Shohei Nobuhara for helpful discussions.

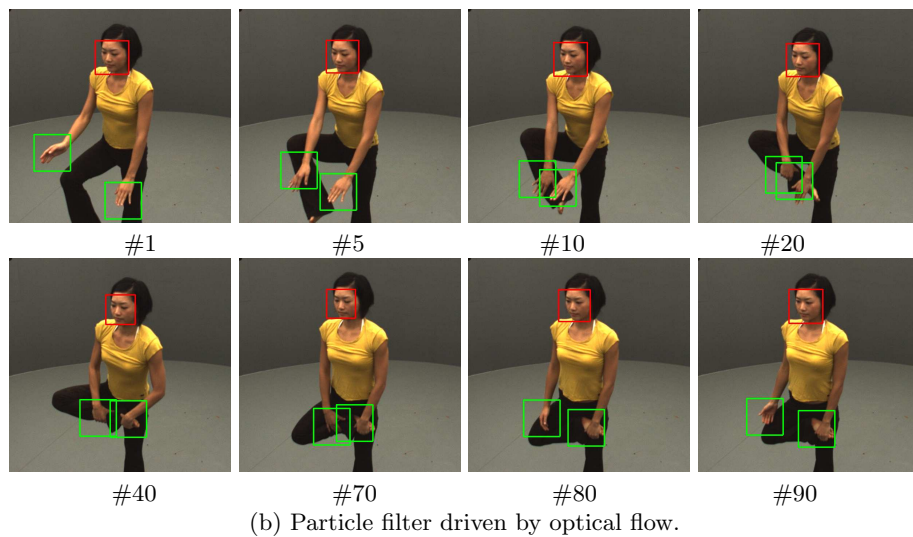
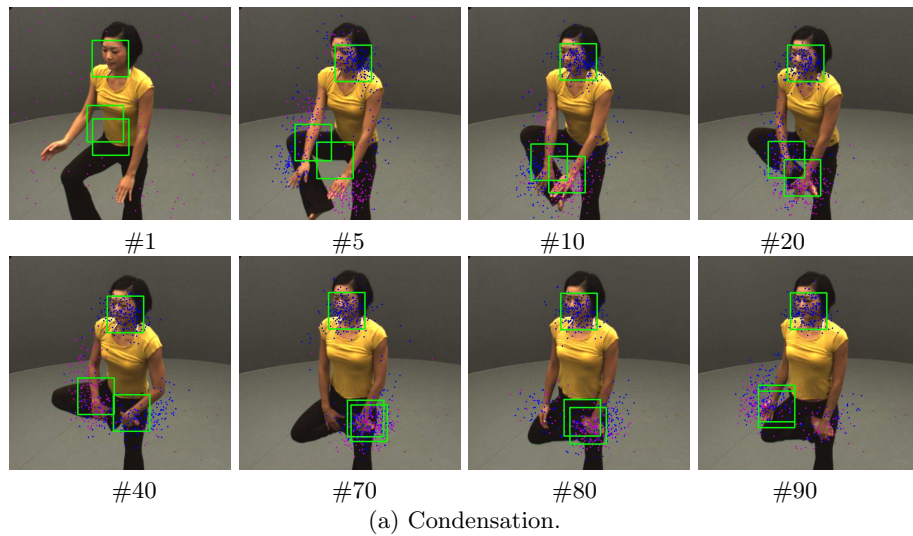


Fig. 5. Robust body part tracking. (a) Classical Condensation methods [2, 12] are confused by regions with similar color and shape content. (b) In frame #20, both hands are almost included in a same tracking window, but afterwards motion cues have helped to discriminate the different tracks.

References

1. Viola, P., Jones, M.: Rapid object detection using a boosted cascade of simple features. CVPR (2001) 511–518
2. Isard, M., Blake, A.: Condensation - conditional density propagation for visual tracking. IJCV **29**(1) (1998) 5–28

3. Lucas, B., Kanade, T.: An iterative image registration technique with an application to stereo vision. *IJCAI* (1981) 674–679
4. Tomasi, C., Kanade, T.: Detection and tracking of point features. Technical Report CMU-CS-91-132, Carnegie Mellon University (1991)
5. Rubner, Y., Tomasi, C., Guibas, L.J.: A metric for distributions with applications to image databases. *ICCV* (1998) 59–66
6. Hjelmas, E., Low, B.K.: Face detection: a survey. *CVIU* **83** (2002) 236–274
7. Choudhury, R., Schmid, C., Mikolajczyk, K.: Face detection and tracking in a video by propagating detection probabilities. *PAMI* **25** (2003)
8. Lowe, D.G.: Distinctive image features from scale-invariant keypoints. *IJCV* **60**(2) (2004) 91–110
9. Lucena, M., Fuertes, J.M., de la Blanca, N.P.: Evaluation of three optical flow-based observation models for tracking. *ICPR* (2004) 236–239
10. Tola, E., Lepetit, V., Fua, P.: A fast local descriptor for dense matching. *CVPR* (2008)
11. Doucet, A., Godsill, S., Andrieu, C.: On sequential monte carlo sampling methods for bayesian filtering. *Statistics and Computing* **10**(3) (2000) 197–208
12. Perez, P., Hue, C., Vermaak, J., Gangnet, M.: Color-based probabilistic tracking. *ECCV* (2002) 661–675
13. A. Sugimoto, K.Y., Matsuyama, T.: Tracking human heads based on interaction between hypotheses with certainty. The 13th Scandinavian Conference on Image Analysis (2003)
14. Okuma, K., Taleghani, A., de Freitas, N., Kakade, S., Little, J., Lowe, D.: A boosted particle filter: multitarget detection and tracking. *ECCV* (2004) 28–39
15. Dorkania, F., Davoine, F.: Simultaneous facial action tracking and expression recognition using a particle filter. *ICCV* (2005)
16. Wang, J., Chen, X., Gao, W.: Online selecting discriminative tracking features using particle filter. *CVPR* (2005)
17. Li, Y., Ai, H., Yamashita, T., Lao, S., Kawade, M.: Tracking in low frame rate video: A cascade particle filter with discriminative observers of different lifespans. *CVPR* (2007)
18. Ross, D., Lim, J., Lin, R., Yang, M.: Incremental learning for robust visual tracking. *IJCV* (2007)
19. Kim, M., Kumar, S., Pavlovic, V., Rowley, H.: Face tracking and recognition with visual constraints in real-world videos. *CVPR* (2008)
20. Black, M., Jepson, A.: Eigentracking: Robust matching and tracking of articulated objects using a view-based representation. *IJCV* **26** (1998) 63–84
21. Collins, R., Liu, Y., Leordeanu, M.: On-line selection of discriminative tracking features. *PAMI* (2005)
22. Martinez, A.M., Kak, A.C.: Pca versus lda. *PAMI* **23**(2) (2001) 228–233
23. Bradski, G.: Computer vision face tracking as a component of a perceptual user interface. In *Workshop on Applications of Computer Vision* (1998) 214–219
24. Comaniciu, D., Ramesh, V., Meeh, P.: Real-time tracking of non-rigid objects using mean shift. *CVPR* **2** (2000) 142–149
25. Matthews, I., Ishikawa, T., Baker, S.: The template update problem. *PAMI* **26**(6) (2004) 810–815
26. Hillier, F.S., Lieberman, G.J.: Introduction to mathematical programming. McGraw-Hill (1990)
27. Shirdhonkar, S., Jacobs, D.W.: Approximate earth mover’s distance in linear time. *CVPR* (2008)