

Appearance-Guided Particle Filtering for Articulated Hand Tracking

Wen-Yan Chang^{1,2}, Chu-Song Chen^{1*}, and Yi-Ping Hung^{1,2}

¹ *Institute of Information Science, Academia Sinica, Taipei, Taiwan*

² *Dept. of Computer Science and Information Engineering, National Taiwan University, Taiwan*

Email: {wychang | song}@iis.sinica.edu.tw, hung@csie.ntu.edu.tw

Abstract

We propose a model-based tracking method, called appearance-guided particle filtering (AGPF), which integrates both sequential motion transition information and appearance information. A probability propagation model is derived from a Bayesian formulation for this framework, and a sequential Monte Carlo method is introduced for its realization. We apply the proposed method to articulated hand tracking, and show that it performs better than methods that only use either sequential motion transition information or only use appearance information.

1. Introduction

High degree-of-freedom (DOF) tracking, such as articulated hand tracking in arbitrary situations, is a challenging task. In general, there are two approaches suggested for high DOF tracking. One is the appearance-based approach [2][9][13][14] that estimates articulated motion states directly from images by learning the mapping from an image feature space to the object state space [9]. The other is the model-based approach [3][7][12][15][16] that estimates articulated motion states by projecting a 3D model on to the image space and then compares the projections with the observations. One advantage of the former is that observations from arbitrary viewpoints can be processed. However, large and dense reference images should be collected in advance to get an accurate estimation. Also, effective learning or retrieval in a large image set is very demanding. The latter approach can provide an accurate estimation when a 3D model is well initialized, but searching in a high-dimensional space is very complex.

In the model-based approach, the motion state is recovered from the 3D configuration with the maximal similarity. This problem has been formulated as an optimization problem [8][15], and can also be treated in a probabilistic framework as the state estimation of a dynamic system. Since closed-form solutions of a highly non-linear dynamic system are intractable, sequential Monte Carlo methods such as particle filtering were introduced to solve this problem. In the past, Isard

and Blake [5] introduced the concept of particle filtering to visual tracking, named CONDENSATION. Rui and Chen [10] integrated unscented Kalman filter and particle filtering to generate a better proposal distribution. Particle filtering has been widely used in articulated hand tracking [3][7][16] recently. Wu *et al.* [16] suggested a method to represent the motion state in low-dimensional space by a set of linear manifolds constructed from base configurations and used particle filtering to track. Bray *et al.* [3] integrated the stochastic meta-descent optimization into particle filtering to find good particles for tracking, while Lin *et al.* [7] proposed a stochastic simplex search algorithm by combining the Nelder-Mead algorithm with particle filtering in a feasible space.

However, most of the particle filtering-based high DOF hand-tracking methods only use visual information from previous time steps. Although applying state estimators to a dynamic system has been shown to be effective for visual tracking, it has certain limitations. First, only initial states are employed; thus, the tracking process may get trapped in local minimums. Second, existing state estimation methods find it difficult to apply known object appearance information to boost the tracking performance, even when such information is easy to acquire. To overcome these difficulties, we study the state estimation of a dynamic system under the assumption that there are some known *attractors*, in addition to the initial state, in the state space. In this paper, an attractor is referred to as a state space vector whose observation is known. For a visual tracking problem, attractors are some reference images of the objects with known motion states, and serve as prior knowledge to guide the tracking in a high-dimensional space.

2. Statistical Model and Its Derivations

2.1. Observation Model for Hand Tracking

Let the state parameter vector of a target at time t be denoted as x_t , and its observation as z_t . The history of observations from time 1 to t is denoted as $Z_t = \{z_1, \dots, z_t\}$. A generic 3D hand model that has 22 DOFs is used for

* The corresponding author.

hand tracking in our work, where each finger has 4 DOFs and the palm has 2 DOFs of rotation. In our current implementation, the likelihood $p(z_t | x_t)$ is measured by using hand silhouettes. A hand area H_R is rendered from the 3D hand model under the state x_t . For each input hand area H_I , we define the likelihood by calculating the difference area between H_I and H_R . In practice, two forces are formulated to minimize the difference area. One is the *shrink force*, F_{shrink} , which has a tendency to shrink fingers of the 3D model to minimize the area, $E_{shrink} = H_R - (H_I \cap H_R)$. The other force is the *stretch force*, $F_{stretch}$, which has a tendency to expand fingers of the 3D hand model to minimize the difference area, $E_{stretch} = H_I - (H_I \cap H_R)$. The likelihood is defined by combining these two forces. Then, the likelihood is defined as

$$p(z_t | x_t) \propto \exp(-[w_1 \cdot Area(E_{shrink}) + w_2 \cdot Area(E_{stretch})]^2 / 2\sigma^2), \quad (1)$$

where $Area(E)$ denotes the area of E , σ is a variance constant, and w_1 and w_2 are two weights.

Note that when $w_1 = w_2 = 0.5$, $p(z_t | x_t) \propto$

$$\exp(-Area[(H_I \cup H_R) - (H_I \cap H_R)]^2 / 8\sigma^2),$$

which is the total area of the non-overlapping regions.

2.2. Bayesian Formulation of AGPF

Assume that there are n attractors, A_1, \dots, A_n , in the state space. These attractors are pre-collected training samples whose observations are known and affect the states in $X_t = \{x_1, \dots, x_t\}$ in such a way that the state x_t is not only influenced by its previous state, x_{t-1} , but also by A_1, \dots, A_n . With such prior information, the problem we focus on is as follows:

Given a set of observations from time 1 to t , $Z_t = \{z_1, \dots, z_t\}$, and a set of attractors $A = \{A_1, \dots, A_n\}$, find the maximal posterior (MAP) estimation of the state x_t according to the observations Z_t and the attractors A .

To achieve this, we investigate the probability:

$$p(x_t | Z_t, A). \quad (2)$$

Suppose that

- (i) The history of observations, Z_t , is conditionally independent of the attractors, A , given X_t .
- (ii) The state at time t , x_t , is conditionally independent of the previous states X_{t-2} , given x_{t-1} and A .

The dynamic Bayesian network (BN) structure considered in our work is shown in Fig. 1(a). Note that if the nodes

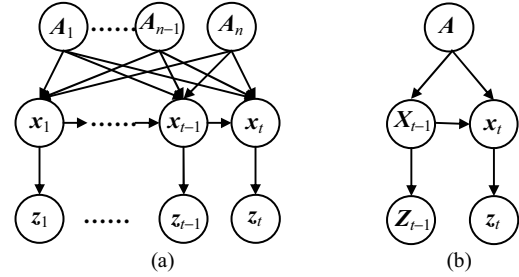


Figure 1. (a) The dynamic BN structure of AGPF. (b) A concise representation.

and links in associations with A are all removed, it degenerates to a BN structure used for the classical particle filtering, where simply a first-order Markov chain is concerned; thus, the states are only influenced by previous time steps.

By condensing $\{A_1, \dots, A_n\}$, $\{x_1, \dots, x_{t-1}\}$, and $\{z_1, \dots, z_{t-1}\}$ to the super nodes A , X_{t-1} , and Z_{t-1} in the BN, respectively, (as shown in Fig. 1(b)), it is easy to observe from the D-separation property [11] that there are some other conditional independencies inherent in the BN as shown below:

- (iii) The observation at time t , z_t , is conditionally independent of Z_{t-1} and X_{t-1} , given the state x_t .
- (iv) The state at time t , x_t , is conditionally independent of Z_{t-1} , given the previous states X_{t-1} .

From (i) to (iv), equation (2) can be resolved as

$$\begin{aligned} p(x_t | Z_t, A) &= \int_{x_1 \dots x_{t-1}} p(x_t | Z_t, A) \propto \int_{x_1 \dots x_{t-1}} p(x_t, Z_t, A) \\ &= \int_{x_1 \dots x_{t-1}} p(A | x_t, Z_t) p(x_t, Z_t) \\ &= \int_{x_1 \dots x_{t-1}} p(A | X_t) p(x_t, Z_t) \\ &= \int_{x_1 \dots x_{t-1}} p(A | X_t) p(z_t | x_t) p(x_t | X_{t-1}, Z_{t-1}) p(X_{t-1}, Z_{t-1}) \\ &= \int_{x_1 \dots x_{t-1}} p(A | X_t) p(z_t | x_t) p(x_t | X_{t-1}) p(X_{t-1}, Z_{t-1}). \end{aligned} \quad (3)$$

In (4), the term $p(A | X_t)$ can be rewritten as

$$\begin{aligned} p(A | X_t) &= p(A, X_{t-1}, x_t) / p(X_{t-1}, x_t) \\ &= p(x_t | X_{t-1}, A) p(A | X_{t-1}) / p(x_t | X_{t-1}) \\ &= p(x_t | x_{t-1}, A) p(A | X_{t-1}) / p(x_t | X_{t-1}). \end{aligned} \quad (5)$$

Note that from (3),

$$p(x_{t-1} | Z_{t-1}, A) \propto \int_{x_1 \dots x_{t-2}} p(A | X_{t-1}) p(X_{t-1}, Z_{t-1}).$$

Then, by substituting (5) into (4), we have

$$p(\mathbf{x}_t | \mathbf{Z}_t, \mathbf{A}) \propto p(\mathbf{z}_t | \mathbf{x}_t) \int_{\mathbf{x}_{t-1}} p(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{A}) \cdot p(\mathbf{x}_{t-1} | \mathbf{Z}_{t-1}, \mathbf{A}). \quad (6)$$

Equation (6) relates the posterior probabilities $p(\mathbf{x}_t | \mathbf{Z}_t, \mathbf{A})$ to $p(\mathbf{x}_{t-1} | \mathbf{Z}_{t-1}, \mathbf{A})$ recursively, which shows how the posterior probabilities propagate given the prior probabilities \mathbf{z}_t and \mathbf{A} . Like the original particle filtering, the MAP estimation of \mathbf{x}_t can be iteratively obtained from previous time steps. A major distinction is that \mathbf{x}_t is further affected by the attractors in \mathbf{A} which contain prior appearance information.

3. Realization of AGPF

The Bayesian formulation of classical particle filtering is expressed as

$$p(\mathbf{x}_t | \mathbf{Z}_t) \propto p(\mathbf{z}_t | \mathbf{x}_t) \cdot p(\mathbf{x}_t | \mathbf{Z}_{t-1}), \quad (7)$$

To precisely compute the posterior probability $p(\mathbf{x}_t | \mathbf{Z}_t)$ for obtaining a Bayesian optimal solution is computationally infeasible for nonlinear/non-Gaussian systems. In particle filtering, sequential Monte Carlo methods using importance sampling or re-sampling have been adopted to realize the computations [1][6][10]. The use of importance sampling has been shown to be a powerful methodology for sequential signal processing, since it can cope with difficult nonlinear and/or non-Gaussian problems.

However, tracking with particle filtering that depends only on the previous time step and the current frame suffers from undesirable drifting effects, particularly for long image sequences and/or high DOF moving objects. Unlike classical particle filtering in which tracking employs only sequential state probability propagation information, we have shown (in Section 2.2) how to obtain Bayesian optimal solutions when further prior knowledge of object appearances is imposed. A suitable importance sampling strategy for finding the solutions, referred to as AGPF, is derived below.

To realize (6), like classical particle filtering, we represent the posterior distribution $p(\mathbf{x}_{t-1} | \mathbf{Z}_{t-1}, \mathbf{A})$ by a set of weighted samples $\{(s_{t-1}^i, \pi_{t-1}^i), i = 1, \dots, N\}$ at time-step $t-1$, where s_{t-1}^i is the sample drawn from $p(\mathbf{x}_{t-1} | \mathbf{Z}_{t-1}, \mathbf{A})$ and π_{t-1}^i is the associated weight. The weights are chosen using the principle of importance sampling [4] with $\pi_t \propto p(\mathbf{X}_t | \mathbf{Z}_t, \mathbf{A}) / q(\mathbf{X}_t | \mathbf{Z}_t, \mathbf{A})$, where $q(\mathbf{X}_t | \mathbf{Z}_t, \mathbf{A})$ is an importance distribution (or proposal function) from which it is easier to draw samples than from the probability density $p(\mathbf{X}_t | \mathbf{Z}_t, \mathbf{A})$.

To derive a useful iteration scheme in which prior appearance information is considered, similar to the derivation of classical particle filtering, we use an

importance function that can be factorized as $q(\mathbf{X}_t | \mathbf{Z}_t, \mathbf{A}) = q(\mathbf{x}_t | \mathbf{X}_{t-1}, \mathbf{Z}_t, \mathbf{A}) \cdot q(\mathbf{X}_{t-1} | \mathbf{Z}_{t-1}, \mathbf{A})$, where s_t^i can be generated from $q(\mathbf{x}_t | \mathbf{X}_{t-1}, \mathbf{Z}_t, \mathbf{A})$ and $\mathbf{S}_{t-1}^i = \{s_{t-1}^1, \dots, s_{t-1}^N; i = 1, \dots, N\}$ can be generated from $q(\mathbf{X}_{t-1} | \mathbf{Z}_{t-1}, \mathbf{A})$. Consider that

$$\pi_t \propto p(\mathbf{X}_t | \mathbf{Z}_t, \mathbf{A}) / q(\mathbf{X}_t | \mathbf{Z}_t, \mathbf{A}) \propto p(\mathbf{X}_t, \mathbf{Z}_t, \mathbf{A}) / q(\mathbf{X}_t | \mathbf{Z}_t, \mathbf{A}). \quad (8)$$

In addition, according to (3) and (4), the term $p(\mathbf{X}_t, \mathbf{Z}_t, \mathbf{A})$ can be rewritten as

$$\begin{aligned} p(\mathbf{X}_t, \mathbf{Z}_t, \mathbf{A}) &= p(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{A}) \cdot p(\mathbf{A} | \mathbf{X}_{t-1}) \cdot p(\mathbf{z}_t | \mathbf{x}_t) \cdot p(\mathbf{X}_{t-1}, \mathbf{Z}_{t-1}) \\ &\propto p(\mathbf{z}_t | \mathbf{x}_t) \cdot p(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{A}) \cdot p(\mathbf{X}_{t-1} | \mathbf{Z}_{t-1}, \mathbf{A}). \end{aligned} \quad (9)$$

By substituting (9) into (8), we have

$$\begin{aligned} \pi_t &\propto \pi_{t-1} \cdot [p(\mathbf{z}_t | \mathbf{x}_t) \cdot p(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{A})] / [q(\mathbf{x}_t | \mathbf{X}_{t-1}, \mathbf{Z}_t, \mathbf{A})] \\ &= \pi_{t-1} \cdot [p(\mathbf{z}_t, \mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{A})] / [q(\mathbf{x}_t | \mathbf{X}_{t-1}, \mathbf{Z}_t, \mathbf{A})]. \end{aligned} \quad (10)$$

The expectation value of π_t , conditional upon $\mathbf{X}_{t-1}, \mathbf{Z}_t$ and \mathbf{A} , is

$$\begin{aligned} E_{q(\mathbf{x}_t | \mathbf{X}_{t-1}, \mathbf{Z}_t, \mathbf{A})}[\pi_t] &= \int_{\mathbf{x}_t} \pi_t \cdot q(\mathbf{x}_t | \mathbf{X}_{t-1}, \mathbf{Z}_t, \mathbf{A}) \\ &\propto \int_{\mathbf{x}_t} \pi_{t-1} \cdot p(\mathbf{z}_t, \mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{A}) = \pi_{t-1} \cdot p(\mathbf{z}_t | \mathbf{x}_{t-1}, \mathbf{A}). \end{aligned}$$

Theoretically, an optimal importance distribution can be chosen by minimizing the variance of the importance weights conditional upon $\mathbf{X}_{t-1}, \mathbf{Z}_t$ and \mathbf{A} [4]. That is,

$$\begin{aligned} \text{Var}_{q(\mathbf{x}_t | \mathbf{X}_{t-1}, \mathbf{Z}_t, \mathbf{A})}[\pi_t] &= E_{q(\mathbf{x}_t | \mathbf{X}_{t-1}, \mathbf{Z}_t, \mathbf{A})}[\pi_t^2] - (E_{q(\mathbf{x}_t | \mathbf{X}_{t-1}, \mathbf{Z}_t, \mathbf{A})}[\pi_t])^2 \\ &\propto (\pi_{t-1})^2 \left[\int_{\mathbf{x}_t} \frac{p^2(\mathbf{z}_t, \mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{A})}{q(\mathbf{x}_t | \mathbf{X}_{t-1}, \mathbf{Z}_t, \mathbf{A})} - p^2(\mathbf{z}_t | \mathbf{x}_{t-1}, \mathbf{A}) \right]. \end{aligned}$$

Since $p(\mathbf{z}_t, \mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{A}) = p(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{z}_t, \mathbf{A}) \cdot p(\mathbf{z}_t | \mathbf{x}_{t-1}, \mathbf{A})$, the variance is minimized to zero when the importance distribution is chosen as $q(\mathbf{x}_t | \mathbf{X}_{t-1}, \mathbf{Z}_t, \mathbf{A}) = p(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{z}_t, \mathbf{A})$. In this case,

$$\begin{aligned} q(\mathbf{x}_t | \mathbf{X}_{t-1}, \mathbf{Z}_t, \mathbf{A}) &= p(\mathbf{z}_t | \mathbf{x}_t, \mathbf{x}_{t-1}, \mathbf{A}) \cdot p(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{A}) / p(\mathbf{z}_t | \mathbf{x}_{t-1}, \mathbf{A}) \\ &= p(\mathbf{z}_t | \mathbf{x}_t) \cdot p(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{A}) / p(\mathbf{z}_t | \mathbf{x}_{t-1}, \mathbf{A}). \end{aligned} \quad (11)$$

Substituting (11) into (10), the optimal weight is

$$\pi_t \propto \pi_{t-1} \cdot p(\mathbf{z}_t | \mathbf{x}_{t-1}, \mathbf{A}) = \pi_{t-1} \int_{\mathbf{x}'_t} p(\mathbf{z}_t | \mathbf{x}'_t) \cdot p(\mathbf{x}'_t | \mathbf{x}_{t-1}, \mathbf{A}). \quad (12)$$

However, the optimal importance weight (12) is difficult to evaluate since an integral is needed. In practice, we choose $q(\mathbf{x}_t | \mathbf{X}_{t-1}, \mathbf{Z}_t, \mathbf{A}) = p(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{A})$ instead. This is similar to the common choice for classical particle filtering, $q(\mathbf{x}_t | \mathbf{X}_{t-1}, \mathbf{Z}_t) = p(\mathbf{x}_t | \mathbf{x}_{t-1})$, but the influence of \mathbf{A} is imposed. In this case, $\pi_t \propto \pi_{t-1} \cdot p(\mathbf{z}_t | \mathbf{x}_t)$, which is the same as that in classical particle filtering.

4. AGPF Filtering Distributions

We set the distribution of $p(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{A})$ as a mixture of $p(\mathbf{x}_t | \mathbf{x}_{t-1})$ and $p(\mathbf{x}_t | \mathbf{A}_i)$, where $i = 1, \dots, n$.

$$p(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{A}) = \alpha_0 p(\mathbf{x}_t | \mathbf{x}_{t-1}) + \sum_{i=1, \dots, n} \alpha_i p(\mathbf{x}_t | \mathbf{A}_i), \quad (13)$$

where $\sum_{i=0, \dots, n} \alpha_i = 1$.

In (13), the probability $p(\mathbf{x}_t | \mathbf{x}_{t-1})$ is the probability caused by state transition, and $p(\mathbf{x}_t | \mathbf{A})$ is the probability caused by appearance similarity. By substituting (13) into (6), we have

$$p(\mathbf{x}_t | \mathbf{Z}_t, \mathbf{A}) \propto p(\mathbf{z}_t | \mathbf{x}_t) \cdot [\sum_{i=1, \dots, n} \alpha_i p(\mathbf{x}_t | \mathbf{A}_i) + \int_{\mathbf{x}_{t-1}} \alpha_0 p(\mathbf{x}_t | \mathbf{x}_{t-1}) \cdot p(\mathbf{x}_{t-1} | \mathbf{Z}_{t-1}, \mathbf{A})]. \quad (14)$$

When α_0 is set to one, we obtain the classical particle filtering. In contrast, if α_0 is set to zero, our method degenerates to a pure appearance-based approach.

In AGPF, a sample set $\{\mathbf{s}_t^i, i = 1, \dots, N\}$ is randomly selected and generated from $\{\mathbf{s}_{t-1}^i\}$ via the transition model $p(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{A})$. Since $p(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{A})$ is defined as a mixture distribution of $p(\mathbf{x}_t | \mathbf{x}_{t-1})$ and $p(\mathbf{x}_t | \mathbf{A})$ in our approach, we generate particles for both of them. Let $\{\hat{\mathbf{s}}_t^k, k = 1, \dots, M_1\}$ and $\{\mathbf{a}^j, j = 1, \dots, M_2\}$ be sets of samples generated from $p(\mathbf{x}_t | \mathbf{x}_{t-1})$ and $p(\mathbf{x}_t | \mathbf{A})$, respectively, where $M_1 + M_2 = N$. To draw the mixture distribution in (6), M_1 is set as $\alpha_0 \cdot N$ and M_2 is set as $(1 - \alpha_0) \cdot N$. Hence, the sample set at time-step t $\{\mathbf{s}_t^i, i = 1, \dots, N\}$ is $\{\hat{\mathbf{s}}_t^k\} \cup \{\mathbf{a}^j\}$. The weight π_t^i of each sample in $\{\mathbf{s}_t^i\}$ is calculated from the observation distribution $p(\mathbf{z}_t | \mathbf{x}_t = \mathbf{s}_t^i)$ by comparing the similarity between observation \mathbf{z}_t and the state \mathbf{x}_t . Finally, the desired posterior distribution $p(\mathbf{x}_t | \mathbf{Z}_t, \mathbf{A})$ can be represented by the set of weighted samples $\{(\mathbf{s}_t^i, \pi_t^i)\}$. The state \mathbf{x}_t at time-step t is estimated by

$$\mathbf{x}_t^* = \mathbf{s}_t^*, \text{ where } \pi_t^* = \max(\pi_t^i). \quad (15)$$

5. AGPF-based Hand Tracking

A primary difficulty of articulated hand tracking is that the motion DOF is too high, resulting in too many possible appearances. Although there are some approaches that use a large quantity of pre-collected appearances [2][9], collecting all the appearances for comparison is still too hard to be feasible. On the other hand, by starting from the initial state in association with an articulated configuration, all the configurations (in association with their appearances) can be reached by gradually changing the articulated motion parameters. This is a significant reason why state-space methods (such as particle filtering) have been widely used in recent

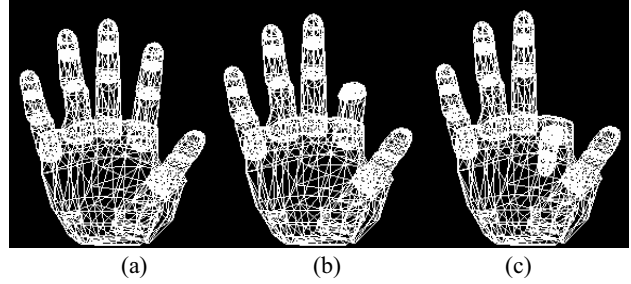


Figure 2. Three appearances of the index finger.

studies. However, since articulated configurations are estimated only from sequential updating information, state-space methods easily mis-track, especially when motions between two consecutive images are large. The AGPF method integrates both the motion transition model and appearance information, and thereby avoids both difficulties.

To construct the appearances for AGPF, only a limited number of the appearances are needed, since sequential motion-transition information is also available. In our work, three appearances for each finger and nine appearances for global hand motion are pre-collected by bending the finger and rotating the palm to different levels. An example of the appearances of the index finger is shown in Fig. 2.

In our framework, all the attractors affect the state to be estimated with certain probabilities. Nevertheless, to make our implementation more efficient, we use the most significant K attractors which have the largest probability values of $p(\mathbf{z}_t | \mathbf{A}_i)$ instead. As the attractors far away from the current state usually have small probabilities that can be neglected, our method becomes more efficient by using this strategy.

6. Experimental Results

In our experiments, we compare the tracking results of the AGPF method with those obtained by either using classical particle filtering, or by using appearance information only. Image sequences with a large range of motions are captured. Twenty particles are used for each experiment, i.e., $N = 20$. The probabilities $p(\mathbf{x}_t | \mathbf{x}_{t-1})$ and $p(\mathbf{x}_t | \mathbf{A}_i)$ are modeled by Gaussian distributions,

$$p(\mathbf{x}_t | \mathbf{x}_{t-1}) \sim \mathcal{N}(\mathbf{x}_{t-1}, \mathbf{\Sigma}_1) \text{ and } p(\mathbf{x}_t | \mathbf{A}_i) \sim \mathcal{N}(\mathbf{A}_i, \mathbf{\Sigma}_2),$$

respectively, where $\mathbf{\Sigma}_1$ and $\mathbf{\Sigma}_2$ are diagonal covariance matrices.

In the first experiment (Fig. 3), K is set as one and an 8-DOF model obtained by bending two fingers is used. There are $3^2 = 9$ attractors, where 3 attractors are

generated for each finger. Fig. 3(a) shows part of the input images. The tracking results of classical particle filtering, the pure appearance-based approach (by setting α_0 as zero), and the AGPF method are shown in Figs. 3(b), 3(c), and 3(d), respectively. From Fig. 3(b), one can see that classical particle filtering fails to track the third, fifth, and sixth images. In addition, tracking also is not very accurate when using pure appearance information, as shown in Fig. 3(c), particular when no pre-collected images match the current motion state (e.g., the second image in Fig. 3(c)). By contrast, the proposed AGPF successfully recovers these motions, as shown in Fig. 3(d).

In the second experiment (Fig. 4), a 14-DOF model and $9 \times 3^3 = 243$ attractors (containing 3 attractors for each finger and 9 attractors for global hand rotation) are used, and K is set to one. As with the first experiment, the tracking results from using classical particle filtering or pure appearance information are not very accurate in this experiment. Fig. 4 (b) shows that classical particle filtering fails to track all of the images, except the first one. Furthermore, the fourth image in Fig. 4(c) shows that tracking also is not very accurate when using pure appearance information. However, our proposed method provides better results, as shown in Fig. 4(d).

Some other experiments were performed with different DOFs and different numbers of attractors. Fig. 5 shows the results of tracking four fingers with 16 DOFs. In this experiment, 81 attractors are used. Another four-finger tracking with 81 attractors is shown in Fig. 6. Finally, Fig. 7 shows the tracking results of global hand motion with five-finger articulation, where the DOF is 22 with 2,187 attractors. The K value sets for the above three experiments are one, three, and one, respectively.

To demonstrate how the tracking trajectory is influenced by the attractors, we use the first experiment as an example. Fig. 8 shows the tracking trajectories of classical particle filtering and AGPF, where principal component analysis (PCA) is used to reduce the high dimensional representation of samples to a 2D space for visualization. In Fig. 8, the blue points and line respectively represent the samples and trajectory generated by classical particle filtering, and the red points and line respectively represent those generated by the AGPF method. The point sign, \bullet , represents the samples generated from the motion transition model and the star sign, $*$, represents the samples generated by attractors. In addition, the attractor with the maximum probability $p(\mathbf{z}_i | \mathcal{A}_i)$ is shown as a green circle. From this figure, it is obvious that attractors can guide the motion trajectory toward a more accurate tracking result.

One can see from the experiments that, with a limited number of pre-collected appearances, our approach

significantly refines the performance of the approach of classical particle filtering in which only a motion transition model is used, and also performs better than the approach when only appearances are used.

7. Conclusions

In this paper, we give a statistically optimized framework for tracking that considers both pre-collected appearance and on-line motion transition information. A probability propagation model that integrates both types of information is derived. With the pre-collected appearances, the proposed method can still be performed efficiently since a sequential Monte Carlo method is provided. By applying prior appearance information, our method can recover rapid motions that are difficult to solve by using simple particle filtering. In addition, unlike pure appearance approaches, the proposed method can handle unseen images by taking advantage of the motion transition model, since it is possible to recover motion states even if they are not pre-selected in the appearance database. Promising tracking results are obtained by using the proposed method in our experiments.

Acknowledgment: This work was supported in part under Grants NSC 93-2752-E-002-007-PAE, NSC 93-2213-E-001-022, and 93-EC-17-A-02-S1-032.

References

- [1] M. S. Arulampalam, S. Maskell, N. Gordon, and T. Clapp, "A tutorial on particle filters for online nonlinear/non-Gaussian Bayesian tracking," *IEEE Trans. on Signal Processing*, pp. 174-188, Vol. 50, No. 2, 2002.
- [2] V. Athitsos and S. Sclaroff, "Estimating 3D Hand Pose from a Cluttered Image," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, pp. 432-439, Vol. 2, 2003.
- [3] M. Bray, E. Koller-Meier, and L. Van Gool, "Smart particle filtering for 3D hand tracking," *Proc. IEEE Int. Conf. Automatic Face & Gesture Recognition*, pp. 675-680, 2004.
- [4] A. Doucet, S. Godsill, and C. Andrieu, "On sequential Monte Carlo sampling methods for Bayesian filtering," *Statist. Comput.*, Vol. 10, No. 3, pp. 197-208, 2000.
- [5] M. Isard and A. Blake, "CONDENSATION-conditional density propagation for visual tracking," *Int. J. of Computer Vision*, pp. 5-28, Vol. 29, No. 1, 1998.
- [6] M. Isard and A. Blake, "ICONDENSATION: unifying low-level and high-level tracking in a stochastic framework," *Proc. European Conf. Computer Vision*, pp. 893-908, 1998.
- [7] J. Y. Lin, Y. Wu, and T. S. Huang, "3D model-based hand tracking using stochastic direct search method," *Proc. IEEE Int. Conf. Automatic Face & Gesture Recognition*, pp. 693-698, 2004.
- [8] S. Lu, D. Metaxas, D. Samaras, and J. Oliensis, "Using multiple cues for hand tracking and model refinement," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, pp. 443-450, Vol. 2, 2003.
- [9] R. Rosales, V. Athitsos, L. Sigal, and S. Sclaroff, "3D hand pose reconstruction using specialized mappings," *Proc. IEEE Int.*

Conf. Computer Vision, pp. 378-385, Vol. 1, 2001.

[10] Y. Rui and Y. Chen, "Better proposal distributions: object tracking using unscented particle filter," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, pp. 786-793, Vol. 2, 2001.

[11] S. Russell and P. Norvig, *Artificial intelligence- A modern approach*, Prentice-Hall Inc., 1995.

[12] B. Stenger, P. R. S. Mendonca, and R. Cipolla, "Model-based 3D tracking of an articulated hand," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, pp. 310-315, Vol. 2, 2001.

[13] B. Stenger, A. Thayananthan, P. H. S. Torr, and R. Cipolla,

"Filtering using a tree-based estimator," *Proc. IEEE Int. Conf. Computer Vision*, pp. 1063-1070, Vol. 2, 2003.

[14] K. Toyama and A. Blake, "Probabilistic tracking in a metric space," *Proc. IEEE Int. Conf. Computer Vision*, pp. 50-57, Vol. 2, 2001.

[15] Y. Wu and T. S. Huang, "Capturing articulated human motion: a divide-and-conquer approach," *Proc. IEEE Int. Conf. Computer Vision*, pp. 606-611, 1999.

[16] Y. Wu, J. Y. Lin, and T. S. Huang, "Capturing natural hand articulation," *Proc. IEEE Int. Conf. Computer Vision*, pp. 426-432, Vol. 2, 2001.

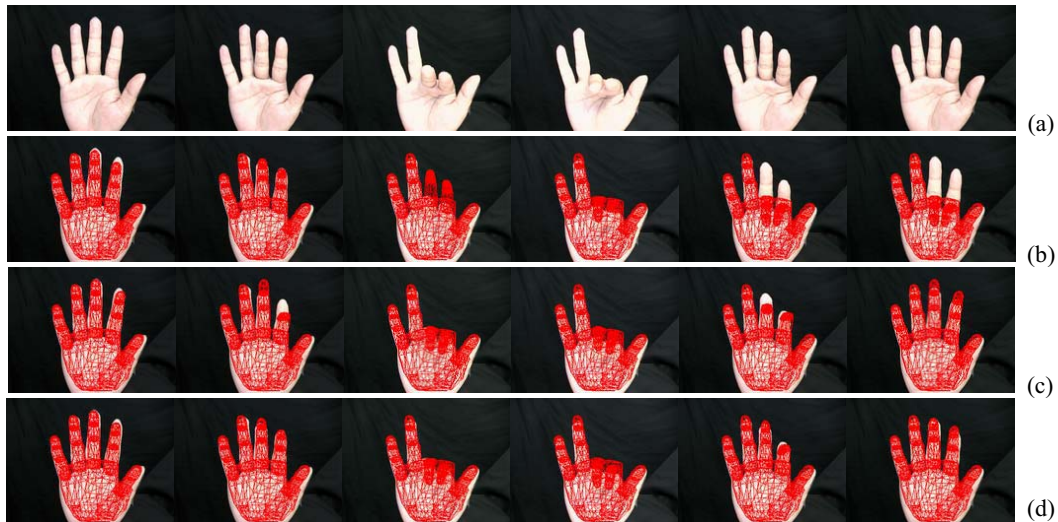


Figure 3. Two-finger tracking with 8 DOFs. (a) Input sequence. (b) Tracking results by using classical particle filtering. (c) Tracking results by using appearance information only. (d) Tracking results by using the AGPF method.

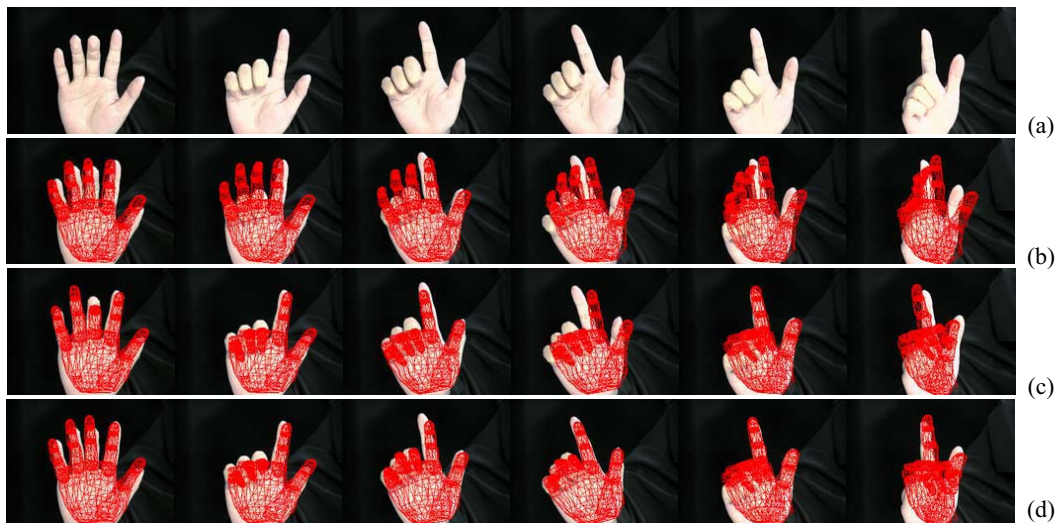


Figure 4. Global hand motion with three-finger articulation in 14 DOFs. (a) Input sequence. (b) Tracking results by using classical particle filtering. (c) Tracking results by using appearance information only. (d) Tracking results by using the AGPF method.

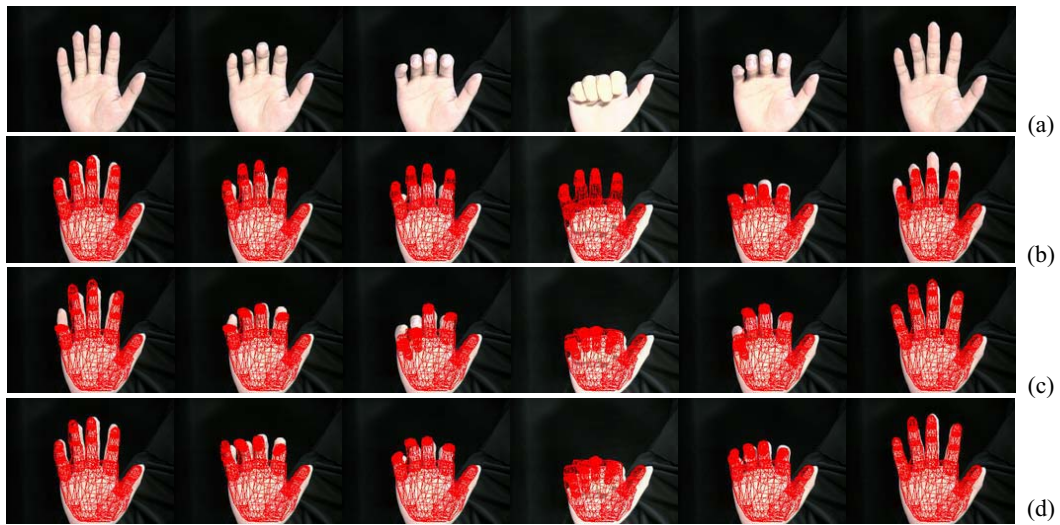


Figure 5. Four-finger tracking with 16 DOFs. (a) Input sequence. (b) Tracking results by using classical particle filtering. (c) Tracking results by using appearance information only. (d) Tracking results by using the AGPF method.

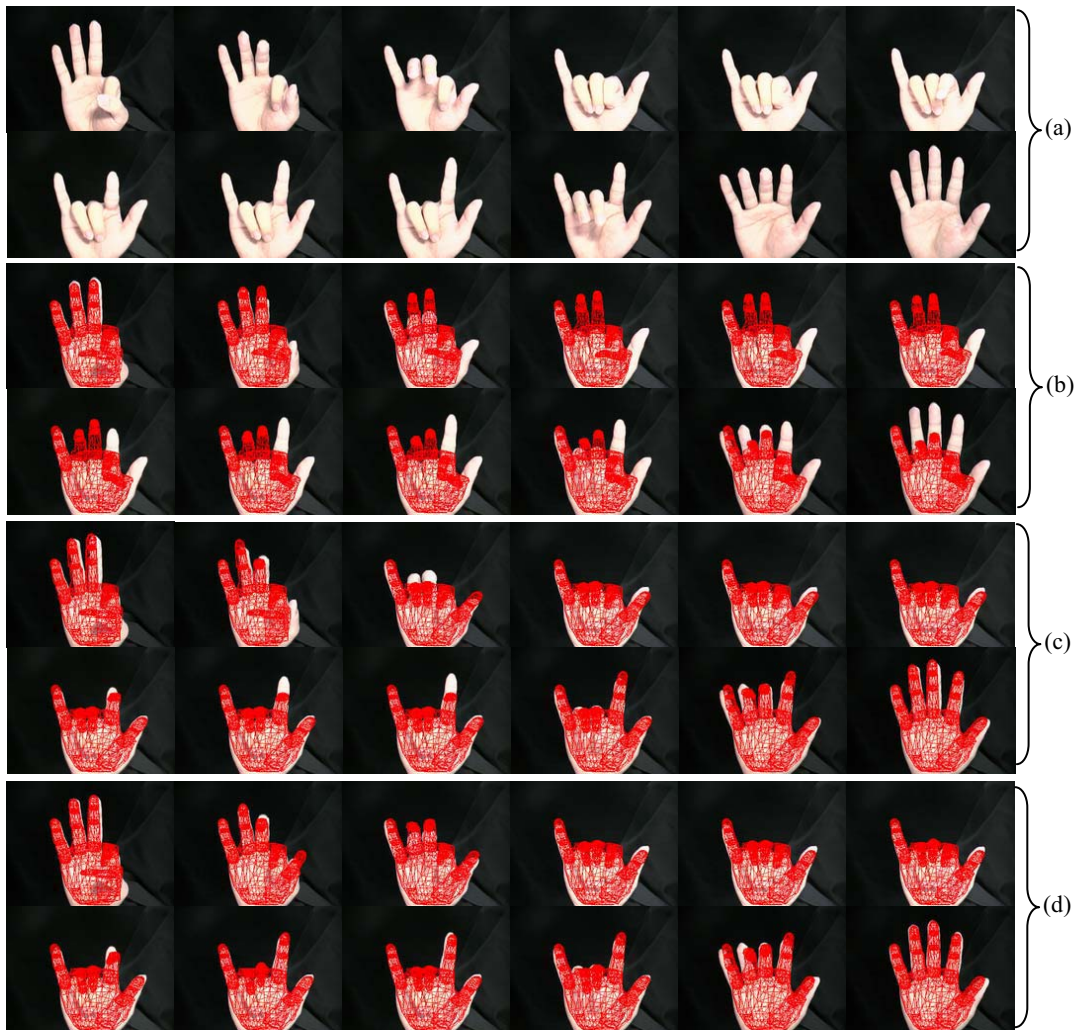


Figure 6. Four-finger tracking with 16 DOFs. (a) Input sequence. (b) Tracking results by using classical particle filtering. (c) Tracking results by using appearance information only. (d) Tracking results by using the AGPF method.

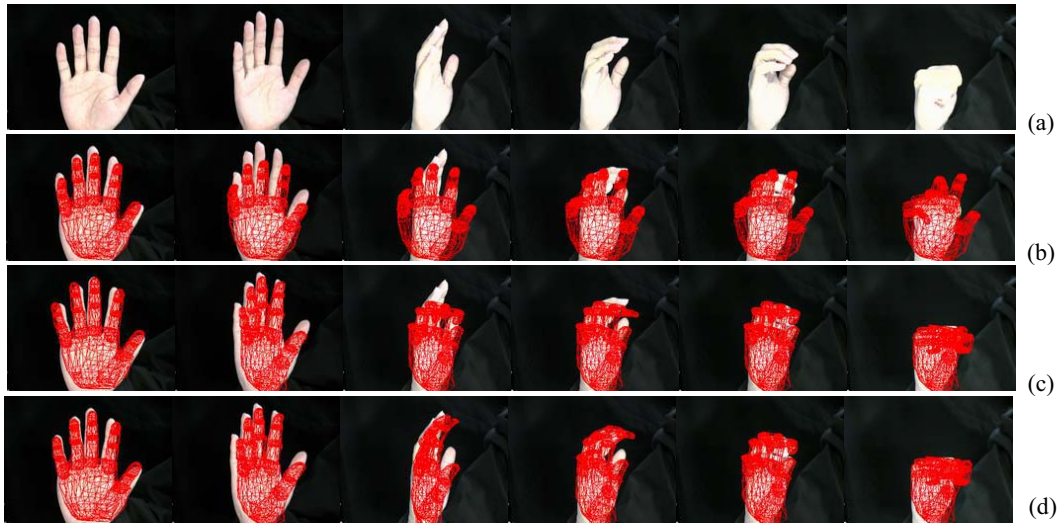


Figure 7. Global hand motion with five-finger articulation in 22 DOFs. (a) Input sequence. (b) Tracking results by using classical particle filtering. (c) Tracking results by using appearance information only. (d) Tracking results by using the AGPF method.

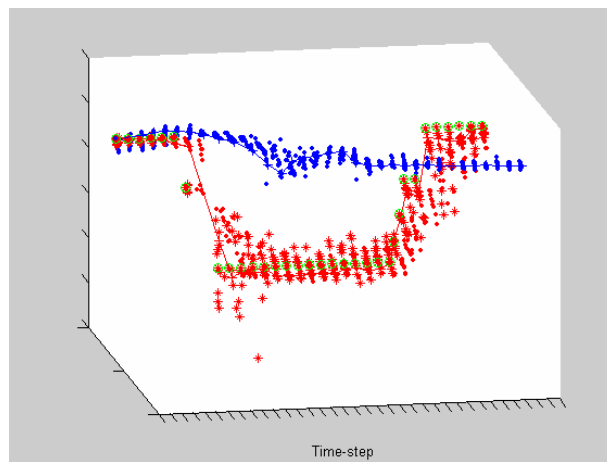


Figure 8. Tracking trajectories reduced to a 2D PCA space of classical particle filtering and the AGPF of the first experiment (Figure 3). The blue points and line respectively represent the samples and trajectory generated by classical particle filtering, and the red points and line respectively represent those generated by the AGPF method. The point sign, \bullet , represents the samples generated from the motion transition model and the star sign, $*$, represents the samples generated by attractors. In addition, the attractor with the maximum probability $p(z_i | A_i)$ is shown as a green circle. From this figure, it is obvious that attractors can guide the motion trajectory toward a more accurate tracking result.