# Style Translation Filter to Change Attribute of Motion

Akihiko Yamaguchi, Shiori Sato, Kentaro Takemura, Jun Takamatsu, and Tsukasa Ogasawara

Graduate School of Information Science

Nara Institute of Science and Technology

8916-5, Takayama, Ikoma, Nara 630-0192, JAPAN

{akihiko-y, j-taka, ogasawar}@is.naist.jp

*Abstract*—In this paper, we propose a style translation filter that changes the attribute (style) of the motion coming from the actors' ages, genders, and so on. Using this filter, we can diversify the motions. Specifically, this filter is modeled by the Gaussian process regression that estimates the difference of pose (joint angles) between a neutral motion and the motion of a target attribute. In learning this filter, a key technique is to find pairs of corresponding posed from the sample motions. We solve this problem by employing the Multifactor Gaussian Process Model (MGPM) proposed by Wang *et al.* [1]. In the experiments, we constructed multiple style translation filters from several attributes of walking motions, such as genders, ages, and emotions. The obtained filters were applied to some kinds of testing motions, such as walking, jumping, kicking, and dancing. The acquired motions were verified by a questionnaire study; the most of their attributes were changed to the filters' target attributes.

## I. INTRODUCTION

Motion capture system is widely used as a tool for generating motions of robots and CG (Computer Graphics) characters because of its ease of utilization. However, in order to generate variety of motions, we need to capture many motions. Even in capturing a same kind of motions, there is variation; for instance, walking motions of a child and an elder person are different. Rather than capturing a number of motions, creating new motions from a small number of captured motions is useful.

For such a purpose, methods to *synthesize* a new motion from *sample motions* are studied [2]∼[6]. Broadly speaking, there are two approaches: (1) dividing motions into several segments and producing a new motion as a sequence of the segments (e.g. the Motion Graphs by Kovar *et al.* [2]), and (2) constructing a generative model from sample motions and generating a new motion with manipulated parameters of the model (e.g. the Style Machine by Brand *et al.* [3]). On the other hand, though there are few researches, it is useful to directly model an *attribute* of the motions coming from the actors' ages, genders, and so on, which enables us to diversify an input motion easily.

Therefore, in this paper, we propose a method to directly model such an attribute of motions. Specifically, we propose the *style translation filter* that changes the attribute (style) of an input motion. This filter is constructed from the sample motions consisting of neutral motions and motions of the attribute; the filter models their difference. Rather than modeling the entire motions and their difference, we simply model the difference between corresponding poses in the motions. Here, a pose means the joint angles of the whole body at a frame. In addition, the difference of speed caused by the attribute is also estimated. For an input motion, frame by frame, we calculate the difference of pose and output the sum of the input pose and the difference, then by adjusting the cycle of the motion, a new motion is obtained.

This filter assumes that among the different kinds of motions, the differences of pose caused by an attribute have a similar tendency. For instance, an elder person bends at waist during a motion, which can be seen in many kinds of motions. Under this assumption, the style translation filter trained with a few motions is applicable to the other kinds of motions widely. Thus, we consider that this filter is applicable to a variety of motions that are not used to construct the filter; we verify that through the experiments.

A key technique in making this filter is to calculate the differences of pose; namely, we need to find pairs of corresponding poses between the neutral motion and the motion of the attribute. We refer to this problem as a phase matching problem. In order to solve this problem, we employ the Multifactor Gaussian Process Model (MGPM) [1]. The calculated differences of pose are modeled using the Gaussian process regression [7].

In the experiments, we construct multiple style translation filters from motions performed by professional actors and actresses. The modeled attributes include genders, ages, and emotions; these attributes are virtually performed by the actors and the actresses. The obtained filters are applied to several kinds of testing motions, such as walking, running, jumping, kicking, picking, and dancing, which are verified by a questionnaire study. The experimental results demonstrate that the proposed style translation filter can convert the attribute of a motion as the filter have been designed.

The rest of this paper is organized as follows. Section II describes the related works. Section III proposes the style translation filter. Section IV describes about the experiments. Section V concludes this paper.

## II. RELATED WORK

Iwai *et al.* proposed a method to extract differences of feature between the dancing motions of Japanese and Latin American [4]; the extracted features can be used for synthesizing a new motion. However, the features are extracted manually, which is difficult to use in our scenario.

Hsu *et al.* proposed a style translation model that is constructed from two motions of the same kind [5]. This method models a translation using a dynamical system.

However, the translation uses a linear time-invariant model; our method can represent much more complex differences.

Inamura *et al.* proposed a method to model motions with the Hidden Markov Models (HMMs), and synthesize a new motion by interpolating or extrapolating the parameters of HMMs [6]. In the CG field, Brand *et al.* proposed the Style Machine which is a kind of HMM [3]. In the style machine, a variable to represent a style is introduced, which is a generative model, so can produce a new motion by setting a new value for the style variable. Wang *et al.* proposed the Multifactor Gaussian Process Models (MGPM), where the kernel includes style parameters that encode the style of motions [1]. MGPM is also a generative model, which can synthesize a new motion.

These approaches are considered to be difficult in changing the attribute of a motion that is not contained in the sample motions. For instance, when we have some walking motions of several genders as the sample motions, changing the attribute of gender of a kicking motion may produce a strange motion, since these methods do not model an attribute directly. Our proposed filter directly models an attribute in a frame-by-frame manner, which makes the filter applicable to an input motion other than the sample motions used to train the filter. In addition, the proposed filter estimates the difference to be near zero for an input pose that is far from those in the sample motions, which hardly produce a strange motion.

## III. STYLE TRANSLATION FILTER

The style translation filter models an attribute of a motion. Specifically, the filter is a regression model that inputs a pose, the joint angles of the whole body at a frame, and outputs a difference of pose between a pose in a neutral motion and a pose in a motion of the attribute. We introduce MGPM [1] for calculating the differences. The calculated differences of pose are modeled using the Gaussian process regression [7]. In this section, we first describe how to learn filter from the sample motions, and then, we describe how to apply the filter to an input motion.

### A. Learning Filter

Style translation filter is learned from a set of sample motions. Each style translation filter is constructed for a *target attribute*. We assume that each motion is labeled by a *kind* of motion, such as walking and jumping, and by an attribute of motion including "neutral". For each pair of kind and attribute, several motions are prepared in order to remove the individuality and the experimental noise.

The learning stage consists of the following steps (Fig. 1):

(1) Calculating the mean of the neutral motions and the mean of the motions of the target attribute by using MGPM respectively.

(2) Calculating the differences of pose between the two mean motions frame by frame.

Step (1) and (2) are performed for each kind of motion. Thus, we obtain data for training the regression model. Then,
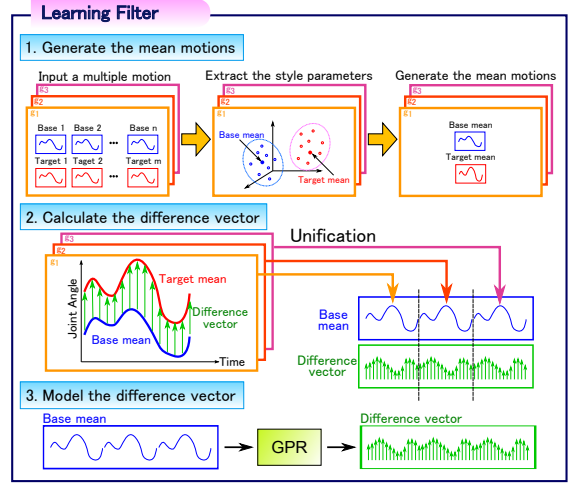


Fig. 1. Learning the style translation filter.

(3) Training the regression model to approximate the difference from a pose in a neutral motion.

In the rest of this section, we describe the detail of learning.

*1) Sample Motions:* Here, we define the set of sample motions used for learning. We denote the set as $\mathcal{M}$. Each motion $m \in \mathcal{M}$ has a label $g_m \in \mathcal{G}$ denoting the kind of motion, and a label $s_m \in \{B, A\}$ denoting the attribute of motion. $\mathcal{G}$ denotes the set of kinds, $B$ denotes the neutral attribute ($B$ comes from Base), and $A$ denotes the target attribute. Each motion $m$ consists of a sequence of states $m = \{\mathbf{y}_{m1}, \dots, \mathbf{y}_{mT_m}\}$ where $T_m$ denotes the number of frames. Each state $\mathbf{y}_m$ consists of the translational velocity of the body link $\mathbf{v}_m$, the joint angles (pose) $\mathbf{q}_m$ including the rotation of the body link, and the joint angular velocity $\boldsymbol{\omega}_m$; that is $\mathbf{y}_{mt} = (\mathbf{v}_{mt}, \mathbf{q}_{mt}, \boldsymbol{\omega}_{mt})$. The rotation of the body link and the rotation of each joint are encoded by Exponential Map [8].

*2) Calculating Mean Motion:* In calculating the differences, we need to find pairs of corresponding posed from the sample motions, i.e. the phase matching problem. To solve this, we first model a set of motions of kind $g$ by using the Multifactor Gaussian Process Model (MGPM). These motions are categorized into two sets according to the attribute $s_m$. In each set of motions, its mean motion is generated by using MGPM.

MGPM is a kind of Gaussian process regression model that inputs multiple latent variables and outputs an observable target variable. In our case, the latent variables are an internal state $\mathbf{x}$ which is time-variant and a style parameter $\boldsymbol{\xi}$ which is time-invariant; the target variable is a pose $\mathbf{y}$. The style parameter $\boldsymbol{\xi}$ is a low-dimensional vector whose dimensionality is chosen for each data set. The internal state $\mathbf{x}$ is constrained by a Circle Dynamics Model (CDM). Namely, the internal state $\mathbf{x}_t$ at time $t$ is given by

$$\mathbf{x}_t = (\cos\theta_t, \sin\theta_t), \quad \theta_t = \theta_0 + t\Delta\theta, \quad (1)$$

where $\theta_0$ denotes an initial phase, and $\Delta\theta$ denotes a step-size of phase. For each motion $m$, the style parameter $\boldsymbol{\xi}_m$ and the parameters of CDM $\theta_{0m}, \Delta\theta_m$ are assigned; they are trained from the sample motions so that the likelihood is maximized.

Next, we produce mean motions by using MGPM with averaged style parameters. Specifically, we first average the style parameters of the neutral motions and the style parameters of the motions of the target attribute as follows:

$$\overline{\boldsymbol{\xi}}_{Bg} = \mathrm{avr}\{\boldsymbol{\xi}_m \mid m \in \mathcal{M}, g_m = g, s_m = B\}, \qquad (2)$$

$$\overline{\boldsymbol{\xi}}_{Ag} = \mathrm{avr}\{\boldsymbol{\xi}_m \mid m \in \mathcal{M}, g_m = g, s_m = A\}, \qquad (3)$$

where $\mathrm{avr}$ denotes a function that averages the elements of a set.

For each $\overline{\boldsymbol{\xi}}_{Bg}$ and $\overline{\boldsymbol{\xi}}_{Ag}$, we produce a motion by MGPM where we use the same sequence of internal states generated by CDM of the same $\theta_0$ and $\Delta\theta$. Let $\overline{m}_{Bg}$ and $\overline{m}_{Ag}$ denote the obtained mean motions respectively, and let $T_{\Delta\overline{m}_g}$ denote the number of frames in each motion.

*3) Calculating Differences of Pose:* The differences of pose are calculated from the mean motions frame by frame as follows:

$$\{\Delta\mathbf{q}_{Agt} \mid \Delta\mathbf{q}_{Agt} = \mathbf{q}_{\overline{m}_{Ag}t} - \mathbf{q}_{\overline{m}_{Bg}t}, t = 1, \ldots, T_{\Delta\overline{m}_g}\}. \quad (4)$$

Corresponding poses in the mean of neutral motions are used in training the filter: $\{\mathbf{q}_{\overline{m}_{Bg}t} \mid t = 1, \ldots, T_{\Delta\overline{m}_g}\}$. If multiple kinds of motions are available as the sample motions, we repeat the same calculation for each motion kind $g$, and unify the results. Let $\{\Delta\mathbf{q}_A\}$ denote the differences of pose, and $\{\mathbf{q}_B\}$ denote the corresponding poses in the mean of neutral motions.

*4) Training Filter:* We construct the style translation filter by training a Gaussian process regression model with $\{\Delta\mathbf{q}_A\}$ and $\{\mathbf{q}_B\}$ where the input is a pose $\mathbf{q}_B$ and the output is a corresponding difference $\Delta\mathbf{q}_A$. Let $\mathbf{f}_A$ denote the regression model which predicts the difference for a given pose as follows: $\Delta\mathbf{q}_A = \mathbf{f}_A(\mathbf{q}_B)$. As the kernel of model, we employ a Gaussian kernel:

$$k_d(\mathbf{q}, \mathbf{q}') = \exp\left(\lambda - \frac{\mu\|\mathbf{q} - \mathbf{q}'\|^2}{2}\right) + b, \qquad (5)$$

where $d$ denotes the dimensionality of pose. The parameters of each kernel $\lambda$, $\mu$, and $b$ are trained by a scaled conjugate gradient method.

In addition, the ratio of the average step-sizes is defined in order to model the translation of motion speed. Specifically, the step-size ratio is given as follows:

$$\overline{\Delta\theta}_B = \mathrm{avr}\{\Delta\theta_m \mid m \in \mathcal{M}, s_m = B\}, \qquad (6)$$

$$\overline{\Delta\theta}_A = \mathrm{avr}\{\Delta\theta_m \mid m \in \mathcal{M}, s_m = A\}, \qquad (7)$$

$$\delta\theta_A = \overline{\Delta\theta}_A / \overline{\Delta\theta}_B. \qquad (8)$$

*B. Applying Filter*

We describe how to apply the style translation filter to an input motion. Here, we assume that the attribute of the input motion is neutral. The filter is applied frame by frame; for
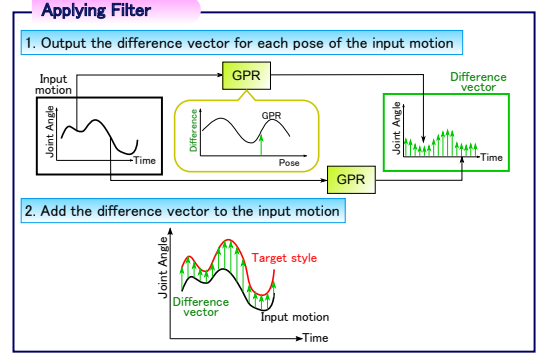


Fig. 2. Applying the style translation filter.

each frame, the difference is calculated for the pose, then is added to the pose with a strength factor $\alpha$ (Fig. 2).

Let $\{\mathbf{y}_{m_{in}t} \mid \mathbf{y}_{m_{in}t} = (\mathbf{v}_{m_{in}t}, \mathbf{q}_{m_{in}t}, \boldsymbol{\omega}_{m_{in}t}), t = 1, \ldots, T_{m_{in}}\}$ denote the input motion. For $t$-th frame, the filter is applied for the pose:

$$\mathbf{q}_{m_{out}t} = \mathbf{q}_{m_{in}t} + \alpha\mathbf{f}_A(\mathbf{q}_{m_{in}t}), \qquad (9)$$

where $\alpha > 0$ is the strength factor to adjust the strength of the filter (its typical value is 1). The other elements of the output motion are the same as those of the input motion; the output motion can be denoted as $\{\mathbf{y}_{m_{out}t} \mid \mathbf{y}_{m_{out}t} = (\mathbf{v}_{m_{in}t}, \mathbf{q}_{m_{out}t}, \boldsymbol{\omega}_{m_{in}t}), t = 1, \ldots, T_{m_{in}}\}$. The step-size ratio $\delta\theta_A$ is applied by changing the FPS (frame per second) of the motion.

## IV. EXPERIMENTS:
## MODELING THE ATTRIBUTE OF ACTING

We construct multiple style translation filters from motions acted by professional actors and actresses. As the sample motions for learning, we use walking motions whose attributes include genders, ages, and emotions; these attributes are virtually performed by the actors and the actresses. Each attribute is modeled by a separate filter. The obtained filters are applied to some testing motions, which are verified by a questionnaire study.

In order to obtain the motions, we use an inertial motion capture system MVN made by Xsens co. A raw motion is captured in 120 Hz which is down-sampled to 30 Hz for training the filters. Each state $\mathbf{y}_{mt}$ consists of the 3-dimensional translational velocity of the body link $\mathbf{v}_{mt}$, the 63-dimensional joint angles (pose) $\mathbf{q}_{mt}$ including the rotation of the body link, and the 63-dimensional joint angular velocity $\boldsymbol{\omega}_{mt}$; that is, $\mathbf{y}_{mt}$ is a 129-dimensional vector.

*A. Training Filters*

Every kind of the sample motions is walking, and there are seven attributes: neutral, masculine, feminine, elderish, childish, sad-looking, and happy-looking. The neutral motion is performed as a natural walking by the actor or the actress. Each motion of these attributes is performed by two actors and two actresses; they are professional, have acting

Table I
STEP-SIZES OF PHASE OF THE MEAN MOTIONS.

| Attribute | Step-size [rad] | Attribute | Step-size [rad] |
|-----------|-----------------|-----------|-----------------|
| Neutral | 0.17 | Childish | 0.19 |
| Masculine | 0.17 | Sad-looking | 0.11 |
| Feminine | 0.17 | Happy-looking | 0.18 |
| Elderish | 0.12 | | |

experience more than eight years, and are aged between 23 and 36.

These motions are modeled by MGPM[1]. Table I shows the step-sizes $\Delta\theta$ of phase of the mean motions. We can find that the step-sizes of the elderish motion and the sad-looking motion are two-thirds of that of the neutral motion.

Fig. 3 shows the acquired mean motions of the attributes where the same sequence of internal state is used in order to compare the corresponding poses. Let us compare the motion of each attribute with the neutral motion. In the masculine motion (Fig. 3(b)), the stride is longer, each arm opens along the coronal plane, and the motion looks duckfooted. In the feminine motion (Fig. 3(c)), the stride is shorter, each arm is along the body, and the motion looks pigeon-toed. In the elderish motion (Fig. 3(d)), the character bends at waist, swings the arms shortly, and bends the knees throughout walking. In the childish motion (Fig. 3(e)), the character opens the arms along the coronal plane, swings the arms widely, and walks bending a knee widely. In the sad-looking motion (Fig. 3(f)), the character bends the neck, swings the arms shortly, and the stride is shorter. In the happy-looking motion (Fig. 3(g)), the character swings the arms widely, and the stride is longer.

From these motions, we calculate the differences of pose, and train six filters. Each filter models an attribute; we refer to each of them as the masculine, the feminine, the elderish, the childish, the sad-looking, and the happy-looking filter respectively. As the implementation of the Gaussian process regression, we use the NETLAB toolbox for Matlab[2].

### B. Applying Filters

In this experiment, we apply the trained six filters to six kinds of testing motions: walking, running, jumping, kicking, picking, and dancing. Note that the testing walking motion is not included in the sample motions. Each motion is input to the filters independently.

We acquired 36 motions; some of them are shown in Fig. 4 to 7. Fig. 4 shows an input walking motion, the output of the masculine filter, and the output of the feminine filter. This input motion was performed by a female subject. In each filter, the strength factor $\alpha$ was set to be $1.5$. Comparing the input and each output at the 30th frame, we find that by the masculine filter, the pose changed to duckfooted one

---

[1]In Section III, the algorithm is defined for a single attribute. In order to handle multiple attributes, we use MGPM in two stages; first, the motions of each attribute are modeled by MGPM and their mean motion is generated, then, the mean motions of all attributes are modeled by MGPM.
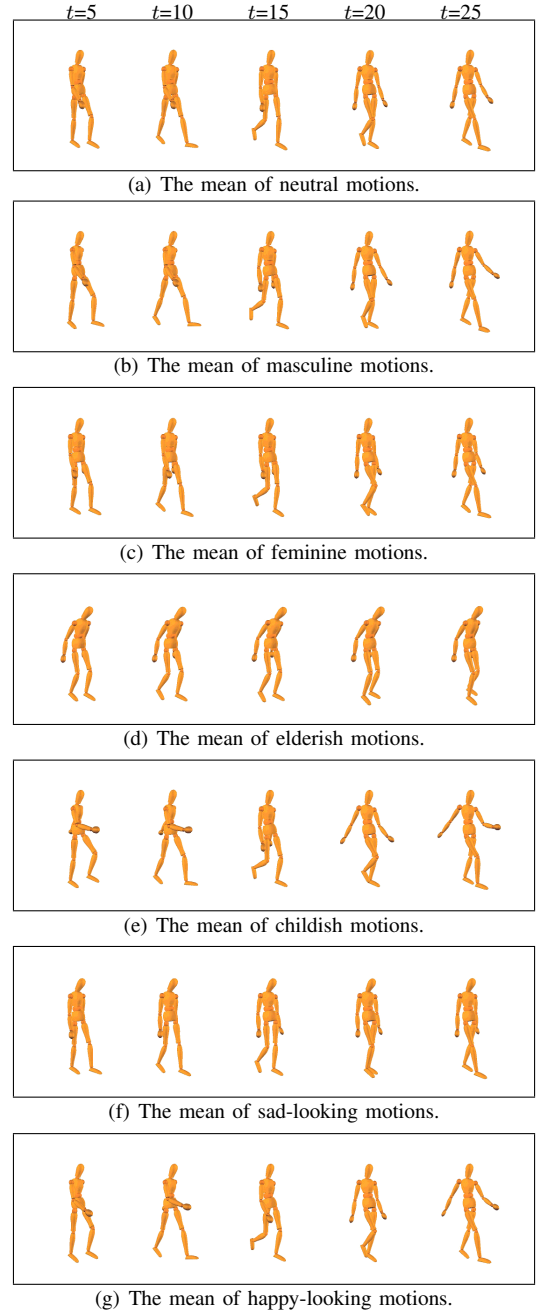
[2]http://www1.aston.ac.uk/eas/research/groups/ncrg/resources/netlab/

---

$t$=5  $t$=10  $t$=15  $t$=20  $t$=25



(a) The mean of neutral motions.



(b) The mean of masculine motions.



(c) The mean of feminine motions.



(d) The mean of elderish motions.



(e) The mean of childish motions.



(f) The mean of sad-looking motions.



(g) The mean of happy-looking motions.

Fig. 3. Snapshots of the mean walking motions of the attributes; $t$ denotes a frame index.

and each arm opened along the coronal plane. On the other hand, by the feminine filter, the pose changed to pigeon-toed. Thus, these filters have a generalization ability to walking motions other than the sample motions. Note that the same conclusions are obtained in using the other filters.

Fig. 5 shows an input picking motion, the output of the elderish filter, and the output of the childish filter. This input motion was performed by a female subject. In each filter, the strength factor $\alpha$ was set to be $1.5$. Comparing the input and each output at the 40th frame, we find that by the elderish filter, the character bent at waist. By the childish
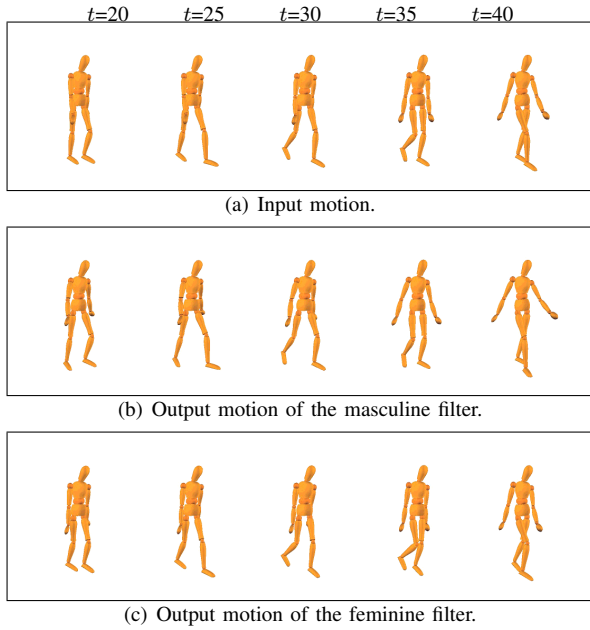
(a) Input motion.



(b) Output motion of the masculine filter.



(c) Output motion of the feminine filter.

Fig. 4. Snapshots of the input and filtered walking motions; $t$ denotes a frame index.



(a) Input motion.



(b) Output motion of the elderish filter.



(c) Output motion of the childish filter.

Fig. 5. Snapshots of the input and filtered picking motions; $t$ denotes a frame index.



(a) Input motion.



(b) Output motion of the sad-looking filter.

Fig. 6. Snapshots of the input and filtered kicking motions; $t$ denotes a frame index.



(a) Input motion.



(b) Filtered motion: translated to the happy style.

Fig. 7. Snapshots of the input and filtered dancing motions; $t$ denotes a frame index.

Fig. 7 shows an input dancing motion, and the output of the happy-looking filter. The input motion was performed by a female subject. In this filter, the strength factor $\alpha$ was set to be 1. Comparing the input and the output, we find that by the happy-looking filter, the character swung the arms widely and the steps was slightly big. From these results, the applicability to various kinds of motions is verified.

### C. Questionnaire Study

In this section, we conduct a questionnaire study of the 6 input and the 36 output motions of the previous section. The subjects consist of 10 males and 10 females who are aged between 23 to 25. In this experiment, we ask each subject to (1) watch the movie of each motion, then (2) rate the movie about gender (masculine/feminine), age (elderish/childish), and emotion (sad-looking/happy-looking) in five levels respectively.

Fig. 8 and 9 show the results of the questionnaire of the 42 motions where each score is rated between $-1$ to $+1$ in five levels. In each graph, each point denotes the mean of score over the subjects; in Fig. 8, the points are plotted on the gender-age graph, and in Fig. 9, the points

filter, the character swung the arms widely. The reason of the difference of height at picking is that in this experiment, we do not constrain the endeffector's position.

Fig. 6 shows an input kicking motion, and the output of the sad-looking filter. This input motion was performed by a male subject. In this filter, the strength factor $\alpha$ was set to be 1.5. Comparing the input and the output at the 60th frame, we find that by the sad-looking filter, the character bent the neck and swung the leg shortly.
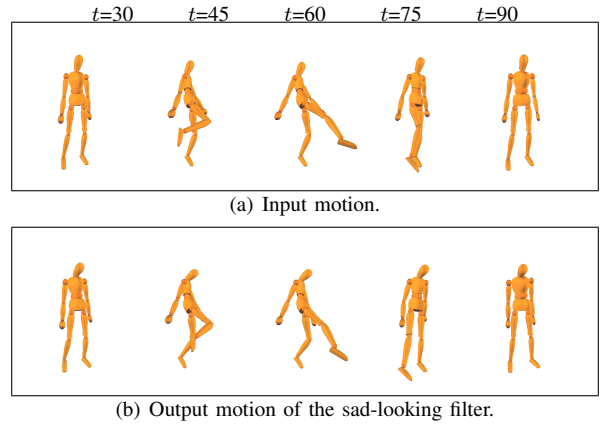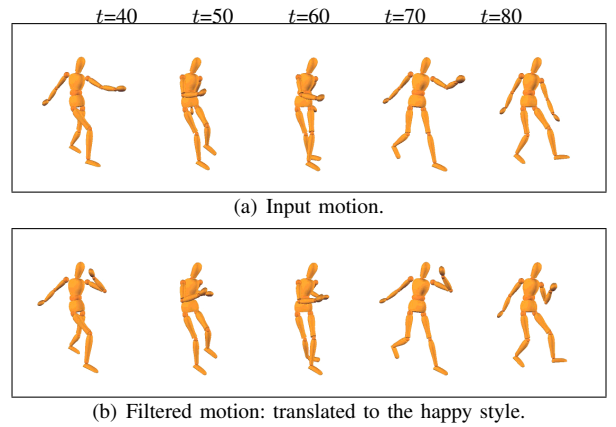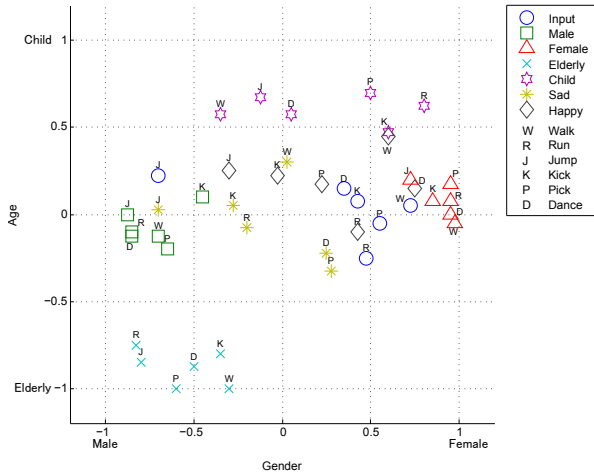
Fig. 8. Evaluation of the degree of gender and age for input and filtered motions.
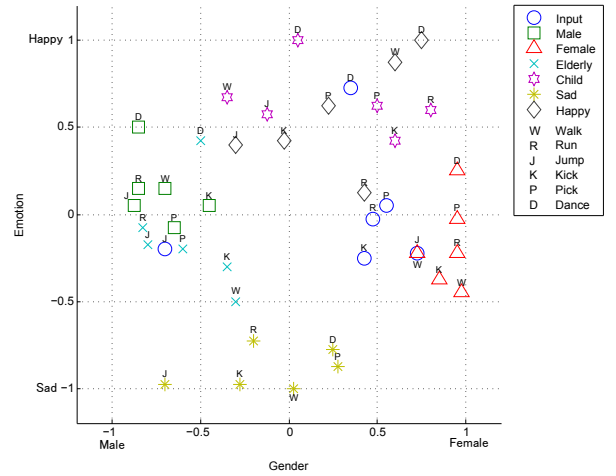


Fig. 9. Evaluation of the degree of gender and emotion for input and filtered motions.

are plotted on the gender-emotion graph. By comparing the points of the output motions with those of the input motions, we find that the most of their attributes were changed to the filters' target attributes. In addition, by using an one-tailed t-test, 34 out of 36 output motions differs significantly from the corresponding input motions ($p < 0.05$). The two motions that do not have significant differences are the jumping motion translated by the masculine filter and the running motion translated by the happy-looking filter. The reason is considered as follows: the input jumping motion originally had a rate of gender close to masculine. About the running motion, it is considered that in such a dynamic motion, recognizing the effect of the happy-looking filter was difficult. Nevertheless, our style translation filters could generally change the attribute of the motion as the filter had been designed.

Finally, we discuss about the side effects of the filters; that is to say, the effects of the filter to the attributes other than the target attribute. From Fig. 8 and 9, we discover the following: (i) the elderish filter changed the attribute of gender toward masculine (Fig. 8), (ii) the childish filter changed the attribute of emotion toward happy (Fig. 9), and (iii) the filters of gender and emotion did not change the attribute of age (Fig. 8, 9). About (i) and (ii), there may be a certain relation between being elder and being masculine, or between being childish and being happy-looking. For instance, we may have an image of children as being happy.

## V. CONCLUSION

In this paper, we proposed the *style translation filter* that directly models the attribute of motions. This filter changes the attribute of an input motion, resulting in diversifying the motions. Specifically, the filter estimates a difference caused by the attribute for each pose of the input motion, where a pose means the joint angles of the whole body at a frame. In learning the filter, we employed MGPM [1] in order to solve

the phase matching problem. The calculated differences of pose were modeled by the Gaussian process regression [7].

In order to verify the filter, we conducted experiments where we constructed multiple style translation filters from motions acted by professional actors and actresses. Specifically, we modeled six attributes: genders (masculine, feminine), ages (childish, elderish), and emotions (sad-looking, happy-looking). These filters are applied to the testing motions that consist of walking, running, jumping, kicking, picking, and dancing motions. From an observation analysis and a questionnaire study, we found that each filter changed the attribute of the input motion as the filter had been designed.

## REFERENCES

[1] J. M. Wang, D. J. Fleet, and A. Hertzmann, "Multifactor gaussian process models for style-content separation," in *the Twenty-Fourth International Conference on Machine Learning (ICML 2007)*, 2007, pp. 975–982.

[2] L. Kovar, M. Gleicher, and F. Pighin, "Motion graphs," in *the 29th annual conference on Computer graphics and interactive techniques SIGGRAPH*, vol. 21, 2002, pp. 473–482.

[3] M. Brand and A. Hertzmann, "Style machines," in *Proceedings of SIGGRAPH 2000*, 2000, pp. 183–192.

[4] D. Iwai, T. Felipe, N. Nagata, and S. Inokouchi, "Identification of motion features affecting perceived rhythmic sense of virtual characters through comparison of latin american and japanese dances." *Information and Media Technologies*, vol. 65, no. 2, pp. 203–210, 2011.

[5] E. Hsu, K. Pulli, and J. Popović, "Style translation for human motion," in *Proceedings of the 2005 ACM SIGGRAPH*, 2005, pp. 1082–1089.

[6] T. Inamura and T. Shibata, "Interpolation and extrapolation of motion patterns in the proto-symbol space," in *International Conference on Neural Information Processing*, 2007.

[7] C. M. Bishop, *Pattern Recognition and Machine Learning*. Springer, 2006.

[8] F. S. Grassia, "Practical parameterization of rotations using the exponential map," *Journal of Graphics Tools*, vol. 3, pp. 29–48, March 1998.