# Segmenting Continuous Motions with Hidden Semi-markov Models and Gaussian Processes

*Tomoaki Nakamura[1]\*, Takayuki Nagai[1], Daichi Mochihashi[2], Ichiro Kobayashi[3], Hideki Asoh[4] and Masahide Kaneko[1]*

[1] *Department of Mechanical Engineering and Intelligent Systems, The University of Electro-Communications, Chofu-shi, Japan,* [2] *Department of Mathematical Analysis and Statistical Inference, Institute of Statistical Mathematics, Tachikawa, Japan,* [3] *Department of Information Sciences, Faculty of Sciences, Ochanomizu University, Bunkyo-ku, Japan,* [4] *Artificial Intelligence Research Center, National Institute of Advanced Industrial Science and Technology, Tsukuba, Japan*

Humans divide perceived continuous information into segments to facilitate recognition. For example, humans can segment speech waves into recognizable morphemes. Analogously, continuous motions are segmented into recognizable unit actions. People can divide continuous information into segments without using explicit segment points. This capacity for unsupervised segmentation is also useful for robots, because it enables them to flexibly learn languages, gestures, and actions. In this paper, we propose a Gaussian process-hidden semi-Markov model (GP-HSMM) that can divide continuous time series data into segments in an unsupervised manner. Our proposed method consists of a generative model based on the hidden semi-Markov model (HSMM), the emission distributions of which are Gaussian processes (GPs). Continuous time series data is generated by connecting segments generated by the GP. Segmentation can be achieved by using forward filtering-backward sampling to estimate the model's parameters, including the lengths and classes of the segments. In an experiment using the CMU motion capture dataset, we tested GP-HSMM with motion capture data containing simple exercise motions; the results of this experiment showed that the proposed GP-HSMM was comparable with other methods. We also conducted an experiment using karate motion capture data, which is more complex than exercise motion capture data; in this experiment, the segmentation accuracy of GP-HSMM was 0.92, which outperformed other methods.

Keywords: motion segmentation, Gaussian process, hidden semi-Markov model, motion capture data

## 1. INTRODUCTION

Human beings typically divide perceived continuous information into segments to enable recognition. For example, humans can segment speech waves into recognizable morphemes. Similarly, continuous motions are segmented into recognizable unit actions. In particular, motions are divided into smaller components called motion primitives, which are used for imitation learning and motion generation (Argall et al., 2009; Lin et al., 2016). It is possible for us to divide continuous information into segments without using explicit segment points. This capacity for unsupervised segmentation is also useful for robots, because it enables them to flexibly learn languages, gestures, and actions.

However, segmentation of time series data is a difficult task. When time series data is segmented, the data points in the sequence must be classified, and each segment's start and end points must be determined. Moreover, each segment affects other segments because of the nature of time series data. Hence, segmentation of time series data requires the exploration of all possible segment lengths and classes. However, this exploration process is difficult; in many studies, the lengths are not estimated explicitly or heuristics are used to reduce computational complexity. Furthermore, in the case of motions, the sequences vary because of dynamic characteristics, even though the same movements are performed. For segmentation of actual human motions, we must address such variations.

In this paper, we propose GP-HSMM (Gaussian process–hidden semi-Markov model), a novel method to divide time series motion data into unit actions by using a stochastic model to estimate their lengths and classes. The proposed method involves a hidden semi-Markov model (HSMM) with a Gaussian process (GP) emission distribution, where each state represents a unit action. **Figure 1** shows an overview of the proposed GP-HSMM. The observed time series data is generated by connecting segments generated by each class. The segment points and segment classes are estimated by learning the parameters of the model in an unsupervised manner. Forward filtering-backward sampling (Uchiumi et al., 2015) is used for the learning process; the segment lengths and segment classes are determined by sampling them simultaneously.
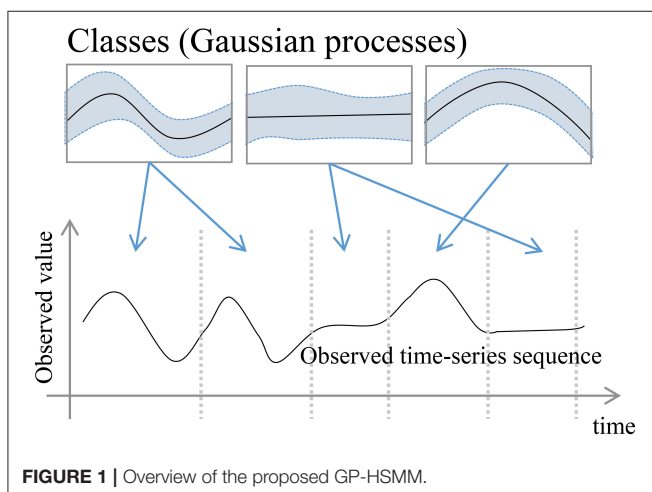
## 2. RELATED WORK

Various studies have focused on learning motion primitives from manually segmented motions (Gräve and Behnke, 2012; Manschitz et al., 2015). Manschitz et al. proposed a method to generate sequential skills by using motion primitives that are learned in a supervised manner. Gräve et al. proposed segmenting motions using motion primitives that are learned by a supervised hidden Markov model. In these studies, the motions

are segmented and labeled in advance. However, we consider that it is difficult to segment and label all possible motion primitives.

Additionally, some studies have proposed unsupervised motion segmentation. However, these studies rely on heuristics. For instance, Wächter et al. have proposed a method to segment human manipulation motions based on contact relations between the end-effectors and objects in a scene (Wachter and Asfour, 2015); in their method, the points at which the end-effectors make contact with an object are determined as boundaries of motions. We believe this method works well in limited scenes; however, there are many motions, such as gestures and dances, in which objects are not manipulated. Lioutikov et al. proposed unsupervised segmentation; however, to reduce computational costs, this technique requires the possible boundary candidates between motion primitives to be specified in advance (Lioutikov et al., 2015). Therefore, the segmentation depends on those candidates, and motions cannot be segmented correctly if the correct candidates are not selected. In contrast, our proposed method does not require such candidates; all possible cutting points are considered by use of forward filtering-backward sampling, which uses the principles of dynamic programming. In some methods (Fod et al., 2002; Shiratori et al., 2004; Lin and Kulić, 2012), motion features (such as the zero velocity of joint angles) are used for motion segmentation. However, these features cannot be applied to all motions. Takano et al. use the error between actual movements and predicted movements as the criteria for specifying boundaries (Takano and Nakamura, 2016). However, the threshold must be manually tuned according to the motions to be segmented. Moreover, they used an HMM that is a stochastic model. We consider such an assumption to be unnatural from the viewpoint of stochastic models, and boundaries should be determined based on a stochastic model. In our proposed method, we do not use such heuristics and assumptions, and instead formulate the segmentation based on a stochastic model.

Fox et al. have proposed unsupervised segmentation for the discovery of a set of latent, shared dynamical behaviors in multiple time series data (Fox et al., 2011). They introduce a beta process, which represents a share of motion primitives in multiple motions, into autoregressive HMM. They formulate the segmentation using a stochastic model, and no heuristics are used in their proposed model. However, in their proposed method, continuous data points that are classified into the same states are extracted as segments, and the lengths of the segments are not estimated. The states can be changed in the short term, and therefore shorter segments are estimated. They reported that some true segments were split into two or more categories, and that those shorter segments were bridged in their experiment. On the other hand, our proposed method classifies data points into states, and uses HSMM to estimate segment lengths. Hence, our proposed method can prevent states from being changed in the short term.

Matsubara et al. proposed an unsupervised segmentation method called AutoPlait (Matsubara et al., 2014). This method uses multiple HMMs, each of which represents a fixed pattern; moreover, transitions between the HMMs are allowed. Therefore,



**FIGURE 1 |** Overview of the proposed GP-HSMM.

time series data is segmented at points at which the state is changed to another HMM's state. However, we believe that HMMs are too simple to represent complicated sequences such as motions. **Figure 2** illustrates an example of representation of time series data by HMM. The graph on the right in **Figure 2** represents the mean and standard deviation learned by HMM from data points shown in the graph on the left. HMM represents time series data using only the mean and standard deviation; therefore, details of time series data can be lost. Therefore, we use Gaussian processes, which are non-parametric methods that can represent complex time series data.

The field of natural language processing has also produced literature related to sequence data segmentation. For example, unsupervised morphological analysis has been proposed for segmenting sequence data (Goldwater, 2006; Mochihashi et al., 2009; Uchiumi et al., 2015). Goldwater et al. proposed a method to divide sentences into words by estimating the parameters of a 2-gram language model based on a hierarchical Dirichlet process. The parameters are estimated in an unsupervised manner by Gibbs sampling (Goldwater, 2006). Mochihashi et al. proposed a nested Pitman-Yor language model (NPYLM) (Mochihashi et al., 2009). In this method, parameters of an $n$-gram language model based on the hierarchical Pitman-Yor process are estimated via the forward filtering-backward sampling algorithm. NPYLM can thus divide sentences into words more quickly and accurately than the method proposed in (Goldwater, 2006). Moreover, Uchiumi et al. extended the NPYLM to a Pitman-Yor hidden semi-Markov model (PY-HSMM) (Uchiumi et al., 2015) that can divide sentences into words and estimate the parts of speech (POS) of the words by sampling not only words, but also POS in the sampling phase of the forward filtering-backward sampling algorithm. However, these relevant studies aimed to divide symbolized sequences (such as sentences) into segments, and did not consider analogous divisions in continuous sequence data, such as that obtained by analyzing human motion.

Taniguchi et al. proposed a method to divide continuous sequences into segments by utilizing NPYLM (Taniguchi and Nagasaka, 2011). In their method, continuous sequences are discretized and converted into discrete-valued sequences using the infinite hidden Markov model (Fox et al., 2007). The discrete-valued sequences are then divided into segments by

using NPYLM. In this method, motions can be recognized by the learned model, but cannot be generated naively because they are discretized. Moreover, segmentation based on NPYLM does not work well if errors occur in the discretization step.

Therefore, we propose a method to divide a continuous sequence into segments without using discretization. This method divides continuous motions into unit actions. Our proposed method is based on HSMM, the emission distribution of which is GP, which represents continuous unit actions. To learn the model parameters, we use forward filtering-backward sampling, and segment points and classes are sampled simultaneously. However, our proposed method also has limitations. One limitation is that the method requires the number of motion classes to be specified in advance. It is estimated automatically in methods such as (Fox et al., 2011) and (Matsubara et al., 2014). Another limitation is that computational costs are very high, owing to the numerous recursive calculations. We discuss these limitations in the experiments.

# 3. GAUSSIAN PROCESS-HIDDEN SEMI-MARKOV MODEL

**Figure 3** shows a graphical representation of the proposed GP-HSMM. In this figure, $c_j (j = 1, 2, \cdots, J)$ denotes classes
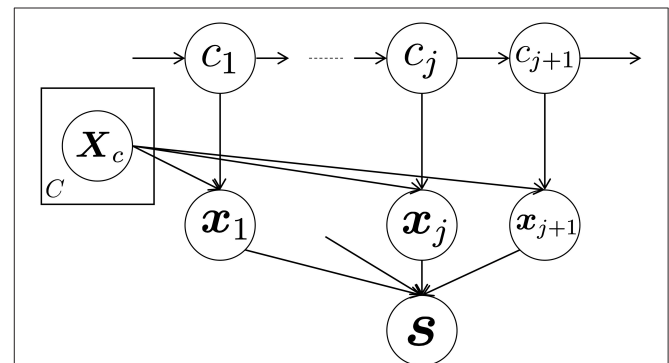


**FIGURE 3 |** Graphical representation of the proposed GP-HSMM.
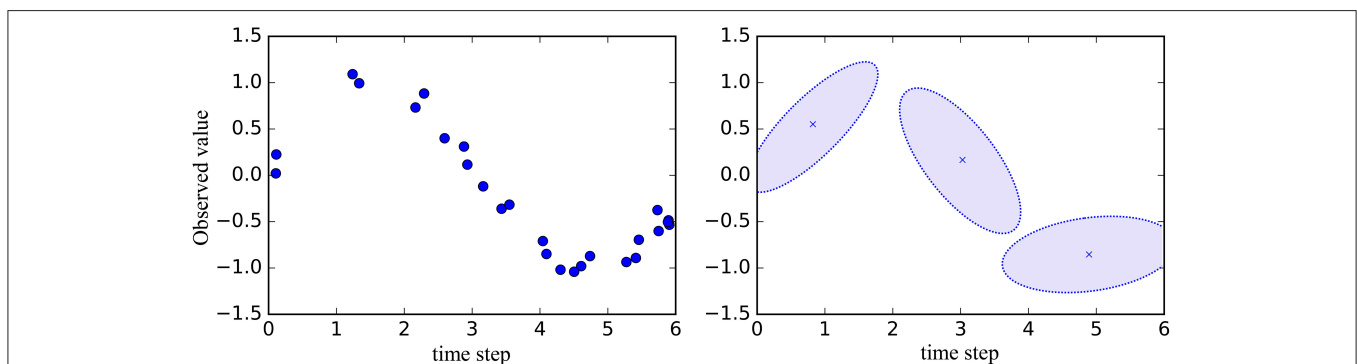


**FIGURE 2 |** Example of representation of time series data by HMM. **Left:** Data points for learning HMM. **Right:** Mean and standard deviation learned by HMM.

of segments, and each segment $x_j$ is generated by a Gaussian process, with parameters denoted by $X_c$ and given by the following generative process:

$$c_j \sim P(c|c_{j-1}), \tag{1}$$

$$x_j \sim \mathcal{GP}(x|X_{c_j}), \tag{2}$$

where $X_c$ represents a set of segments classified into class $c$. Segments are generated by this generative process, and the observed time-series data $s$ is generated by connecting the segments.

## 3.1. Gaussian Process

In this study, we utilize Gaussian process regression, which learns emission $x_i$ of time step $i$ in a segment. This makes it possible to represent each unit action as part of a continuous trajectory. If we obtain pairs $(i, X_c)$ of emissions $x_i$ of time step $i$ of segments belonging to the same class $c$, a predictive distribution whereby the emission of time step $i$ becomes $x$ follows a Gaussian distribution.

$$p(x|i, X_c, i) \propto \mathcal{N}(k^T C^{-1} i, c - k^T C^{-1} k), \tag{3}$$

where $k(\cdot, \cdot)$ represents the kernel function and $C$ is a matrix whose elements are

$$C(i_p, i_q) = k(i_p, i_q) + \beta^{-1} \delta_{pq}. \tag{4}$$

$\beta$ is a hyperparameter that represents noise in the observation. In Equation (3), $k$ is a vector containing the elements $k(i_p, i)$, and $c$ is a scalar value $k(i, i)$. Using the kernel function, GP can learn a time-series sequence that contains complex changes. We use the following Gaussian kernel, which is generally used for Gaussian process regression:

$$k(i_p, i_q) = \theta_0 \exp(-\frac{1}{2}||i_p - i_q||^2 + \theta_2 + \theta_3 i_p i_q), \tag{5}$$

where $\theta_*$ represents parameters of the kernel. **Figure 4** shows examples of Gaussian processes. The left graph in each pair of graphs represents learning data points $(i, X_c)$, and the right graph shows the learned probabilistic distribution $p(x|i, X_c, i)$. One can see that the standard deviation decreases with an increase in the number of learning data points. If the emission of time step $i$ is multidimensional vector $x = (x_0, x_1, \cdots)$, we assume that each dimension is generated independently, and a predictive distribution $\mathcal{GP}(x|X_c)$ is computed as follows:

$$
\begin{aligned}
\mathcal{GP}(x|X_c) = \ & p(x_0|i, X_{c,0}, i_c) \\
& \times p(x_1|i, X_{c,1}, i_c) \\
& \times p(x_2|i, X_{c,2}, i_c) \cdots .
\end{aligned}
\tag{6}
$$

Based on this probability, similar segments can be classified into the same class.

## 3.2. Learning of GP-HSMM
### 3.2.1. Blocked Gibbs Sampler
Segments and classes of segments in the observed sequences are estimated based on dynamic programming and sampling. For efficient sampling, we use the blocked Gibbs sampler, which samples segments and their classes in an observed sequence. In the initialization phase, all observed sequences are first randomly divided into segments. Segments $x_{nj}(j = 1, 2, \cdots, J_n)$ in observed sequence $s_n$ are then removed from the learning data, and parameter $X_c$ of the Gaussian process and transition probability $P(c|c')$ of HSMM are updated. Segments $x_{nj}(j = 1, 2, \cdots, J_n)$ and their classes $c_{nj}(j = 1, 2, \cdots, J_n)$ are then estimated as follows:

$$(x_{n1}, \cdots, x_{nJ_n}), (c_{n1}, \cdots, c_{nJ_n}) \sim P(X, c|s_n), \tag{7}$$

where $X$ is a set of segments into which $s_n$ is divided, and $c$ denotes classes of the segments. To carry out this sampling efficiently, the probability of all possible segments $X$ and
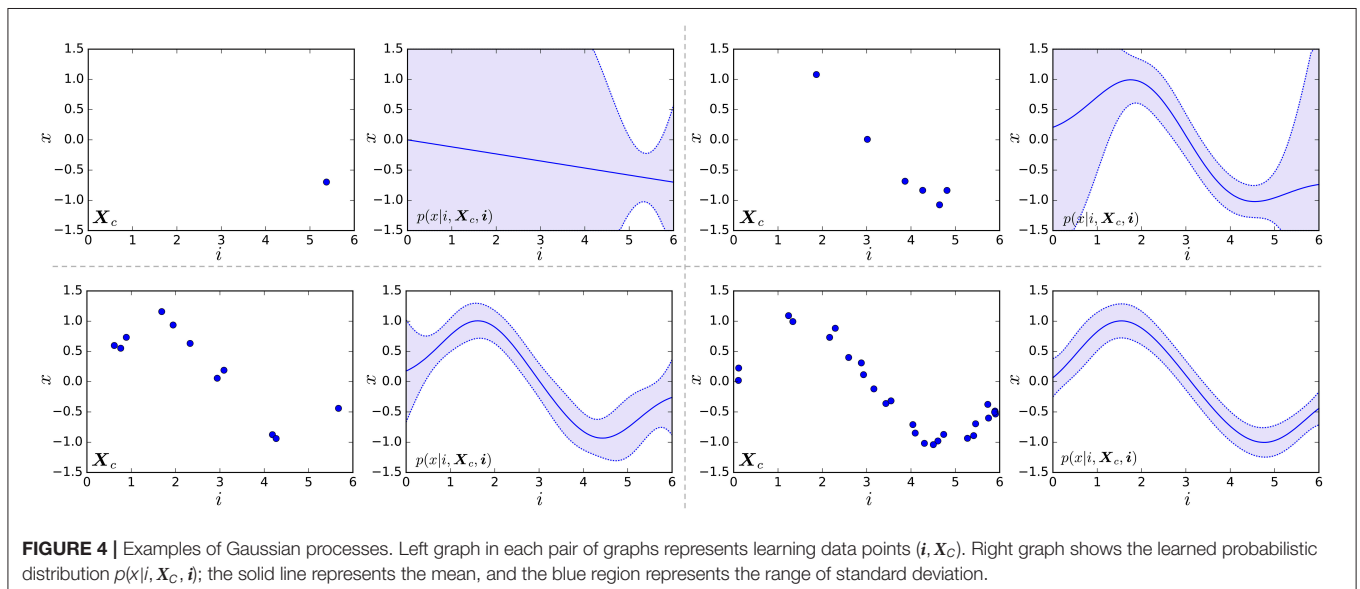


**FIGURE 4** | Examples of Gaussian processes. Left graph in each pair of graphs represents learning data points $(i, X_c)$. Right graph shows the learned probabilistic distribution $p(x|i, X_c, i)$; the solid line represents the mean, and the blue region represents the range of standard deviation.

**Algorithm 1** Blocked Gibbs Sampler

```
1:  // Iterate the following procedure until convergence
2:  for  n = 1 to N do
3:      for  j = 1 to J_n do
4:          N_{c_{nj}} − = 1
5:          N_{c_{nj}, c_{n,j+1}} − = 1
6:          if  j ≠ 0 then
7:              Delete segments x_{nj} from X_{c_{nj}}
8:          end if
9:      end for
10:
11:     // Sample segments and their classes
12:     (x_{n1}, · · · , x_{nJ_n}), (c_{n1}, · · · , c_{nJ_n}) ∼ P(X, c|s_n)
13:
14:     for  j = 1 to J_n do
15:         N_{c_{nj}} + +
16:         N_{c_{nj}, c_{n,j+1}} + +
17:         if  j ≠ then
18:             Add segments x_{nj} into X_{c_{nj}}
19:         end if
20:     end for
21: end for
```

**Algorithm 2** Forward filtering-backward sampling

```
1:  // Forward filtering
2:  for  t = 1 to T do
3:      for  k = 1 to K do
4:          for  c = 1 to C do
5:              Compute α[t][k][c]
6:          end for
7:      end for
8:  end for
9:
10: // Backward sampling
11: t = T, j = 1
12: while  t > 0 do
13:     k, c ∼ α[t][k][c]
14:     x_j = s_{t−k:t}
15:     c_j = c
16:     t = t − k
17:     j = j + 1
18: end while
19: return (x_{J_n}, x_{J_n−1}, · · · , x_1), (c_{J_n}, c_{J_n−1}, · · · , c_1)
```

classes $c$ must be computed; however, these probabilities are difficult to compute simply because the number of potential combinations is very large. Thus, we utilize forward filtering-backward sampling, which we presently explain. After sampling $x_{nj}$ and $c_{nj}$, parameter $X_c$ of the Gaussian process and transition probability $P(c|c')$ of HSMM are updated by adding them to the learning data. The segments and parameters of Gaussian processes are optimized alternately by iteratively performing the above procedure. Algorithm 1 shows the pseudocode of the blocked Gibbs sampler. $N_{c_{nj}}$ and $N_{c_{nj}, c_{n, j+1}}$ represent parameters for computing the transition probability in Equation (10).

### 3.2.2. Forward Filtering-Backward Sampling

In this study, we regard segments and their classes as latent variables that are sampled by forward filtering-backward sampling (Algorithm 2). In forward filtering, as shown in
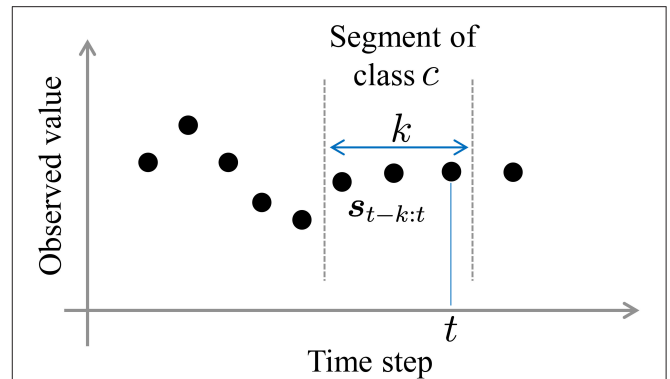


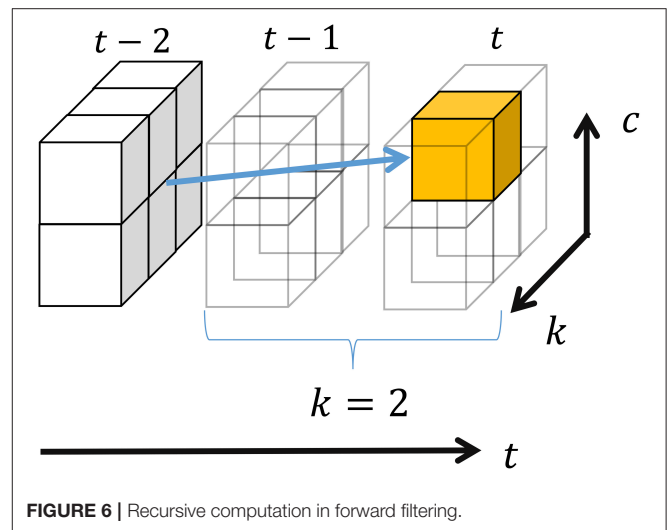**FIGURE 5 |** A segment whose probability is computed during forward filtering.



**FIGURE 6 |** Recursive computation in forward filtering.

Figure 5, the probability that $k$ samples $s_{t−k:t}$ prior to time step $t$ in observed sequence $s$ form a segment, and that the resulting segment belongs to class $c$, is computed as follows:

$$\alpha[t][k][c] = P(s_{t−k:t}|X_c)$$
$$\times \sum_{k'=1}^{K} \sum_{c'=1}^{C} p(c|c')\alpha[t−k][k'][c'], \quad (8)$$

where $C$ and $K$ denote the number of classes and the maximum length of segments, respectively. $P(s_{t−k:t}|X_c)$ represents the probability that $s_{t−k:t}$ is generated from a class $c$; this is computed as follows:

$$P(s_{t−k:t}|X_c) = \mathcal{GP}(s_{t−k:t}|X_c)P_{len}(k|\lambda). \quad (9)$$

where $P_{len}(k|\lambda)$ represents a Poisson distribution with a mean parameter $\lambda$; this corresponds to the distribution of the segment lengths. $p(c|c')$ in Equation (8) represents a transition probability computed as follows:

$$p(c|c') = \frac{N_{c'c} + \alpha}{N_{c'} + C\alpha}, \quad (10)$$

where $N_{c'}$ and $N_{c'c}$ denote the number of segments whose classes are $c'$ and the number of transitions from $c'$ to $c$, respectively, and $k'$ and $c'$ respectively denote the length and class of the segment preceding $s_{t-k:t}$; these are marginalized out in Equation (8). Moreover, $\alpha[t][k][*] = 0$ if $t - k < 0$, and $\alpha[0][0][*] = 1.0$. All elements of $\alpha[*][*][*]$ in Equation (8) can be recursively computed from $\alpha[1][1][*]$ by dynamic programming. **Figure 6** depicts the computation of a three-dimensional array $\alpha[t][k][c]$. In this example, the probability that two samples before time step $t$ become a segment is computed; the resulting segment would be assigned to class two. Hence, samples at $t - 1$ and $t$ become a segment, and all the segments whose end point is $t - 2$ can potentially transit to this segment. $\alpha[t][2][2]$ can be computed by marginalizing out these possibilities.

Finally, segment $x_j$ and its class are determined by backward sampling length $k$ and class $c$ of the segment, based on forward

probabilities in $\alpha$. From $t = T$, length $k_1$ and class $c_1$ are determined according to $k_1, c_1 \sim \alpha[T][k][c]$, and $s_{T-k_1:T}$ becomes a segment whose class is $c_1$. Then, length $k_2$ and class $c_2$ of the next segment are determined according to $k_2, c_2 \sim \alpha[T - k_1][k][c]$. By iterating this procedure until $t = 0$, the observed sequence can be divided into segments and their classes can be determined.

# 4. EXPERIMENTS

We conducted experiments to confirm the validity of the proposed method. We used two types of motion capture data: (1) data from the CMU motion capture dataset (CMU, 2009), and (2) data containing karate motions.

## 4.1. Segmentation of Exercise Motions

We first applied our proposed method to CMU motion capture data containing several exercise routines. The CMU motion capture data was captured using a Vicon motion capture system, and positions and angles of 31 body parts are available. The dataset contains 2605 trials in six categories and 23 subcategories, and motions in each subcategory were performed by one or a few subjects. In this experiment, three sequences from subject 14 in the general exercise and stretching category were used, and include running, jumping, squats, knee raises, reach out stretches, side stretches, body twists, up and down movements, and toe touches. To reduce computational cost, we downsampled from 120 frames per second to 4 frames per second. **Figure 7** shows the coordinate system of motion capture data used in this experiment; two-dimensional frontal views of the left hand $(x_{lh}, y_{lh})$, right hand $(x_{rh}, y_{rh})$, left foot $(x_{lf}, y_{lf})$, and right foot $(x_{rf}, y_{rf})$ were used. Therefore, each frame was represented by eight dimensional vectors:
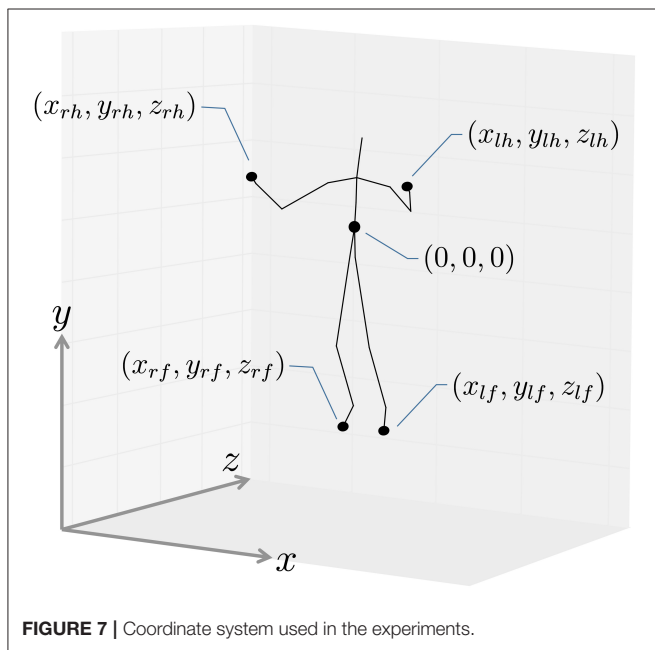


**FIGURE 7 |** Coordinate system used in the experiments.

**TABLE 1 |** Segmentation accuracy of CMU motion capture data.

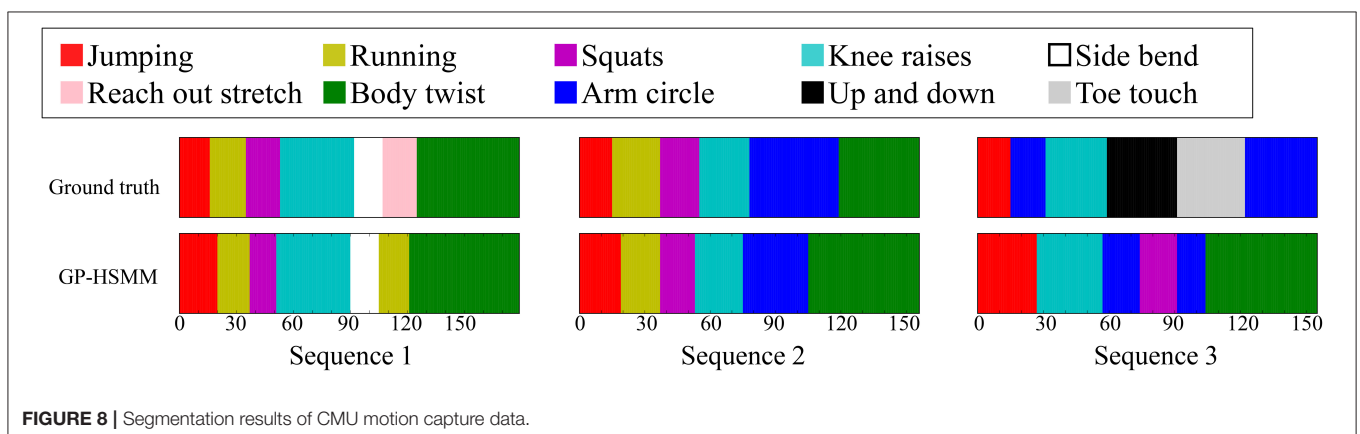| Hamming distance | Precision | Recall | F-measure |
|---|---|---|---|
| 0.33 | 0.81 | 0.81 | 0.81 |



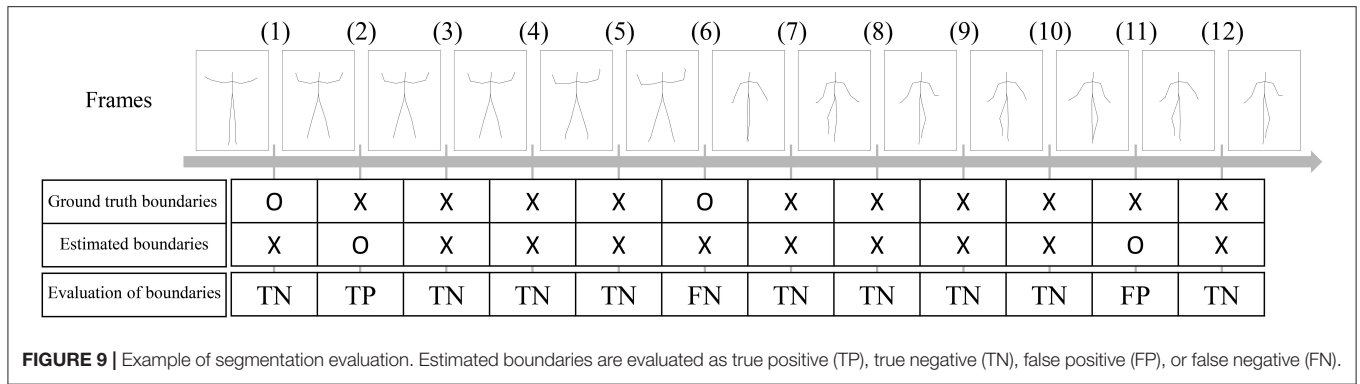**FIGURE 8 |** Segmentation results of CMU motion capture data.

**FIGURE 9 |** Example of segmentation evaluation. Estimated boundaries are evaluated as true positive (TP), true negative (TN), false positive (FP), or false negative (FN).
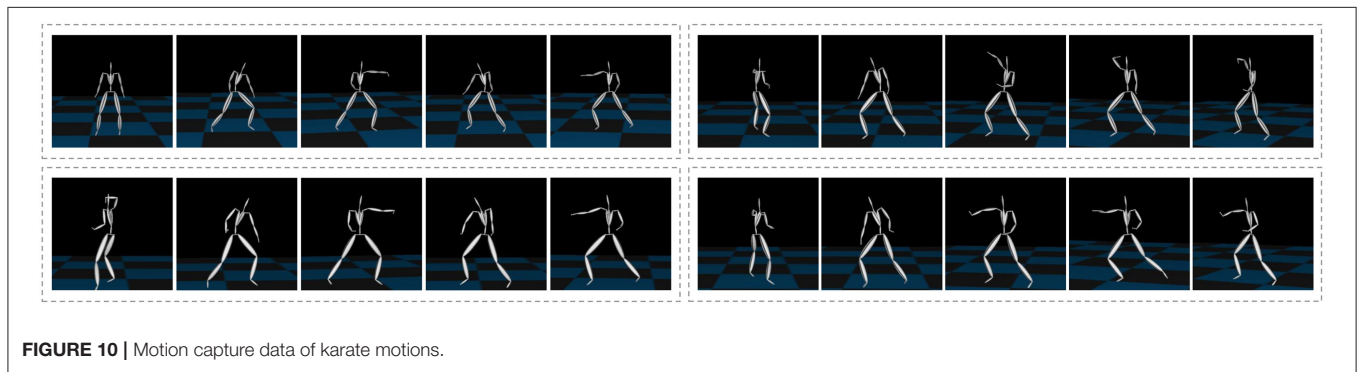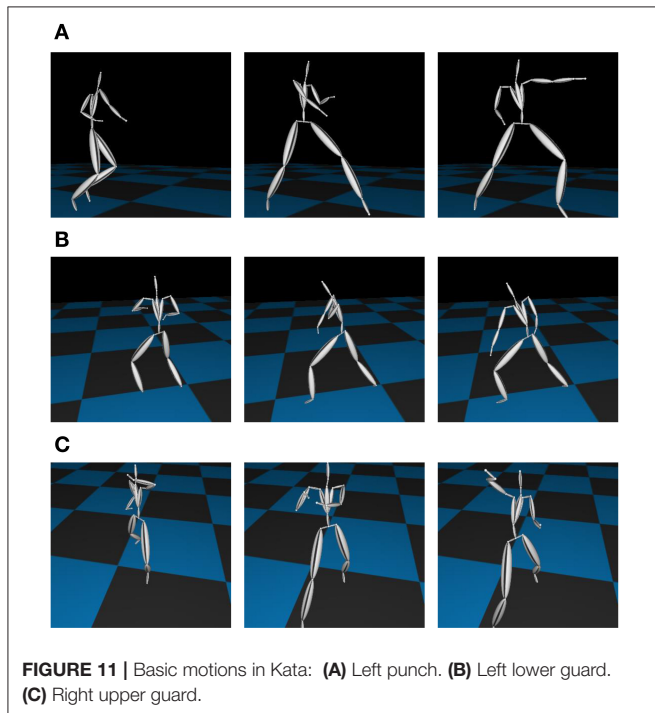


**FIGURE 10 |** Motion capture data of karate motions.



**FIGURE 11 |** Basic motions in Kata: **(A)** Left punch. **(B)** Left lower guard. **(C)** Right upper guard.

$(x_{lh}, y_{lh}, x_{rh}, y_{rh}, x_{lf}, y_{lf}, x_{rf}, y_{rf})$. Because GP-HSMM requires the number of classes to be specified in advance, we set it to eight.

**Figure 8** shows the results of the segmentation. The horizontal axis represents the frame number, and the colors represent motion classes into which each segment was classified. The segments were classified into seven classes out of eight. **Table 1** shows the accuracy of the segmentation. We computed the following normalized Hamming distance between the unsupervised segmentation and the ground truth:

$$ND(\boldsymbol{c}, \bar{\boldsymbol{c}}) = \frac{D(\boldsymbol{c}, \bar{\boldsymbol{c}})}{|\bar{\boldsymbol{c}}|}, \qquad (11)$$

where $\boldsymbol{c}$ and $\bar{\boldsymbol{c}}$ represent sequences of estimated motion classes and true motion classes, $D(\boldsymbol{c}, \bar{\boldsymbol{c}})$ is the Hamming distance between two sequences, and $|\bar{\boldsymbol{c}}|$ represents the length of the sequence. Therefore, the normalized Hamming distance ranges from 0 to 1; lower Hamming distances indicate more accurate segmentation. In this experiment, the Hamming distance was 0.33, which is comparable with the BP-HMM reported in (Fox et al., 2011). However, they also reported that some segments were split into two or more categories, and that those shorter segments were bridged. In contrast, we performed no such modifications, and **Figure 8** shows that there are no shorter segments. We also computed the precision, recall, and F-measure of the segmentation. To compute them, estimated boundaries of segments are evaluated as true positive (TP), true negative (TN), false positive (FP), or false negative (FN). **Figure 9** shows an example of segmentation evaluation. We considered the estimated boundary to be TP if it was within true boundary ± four frames, as shown in **Figure 9**(2). If
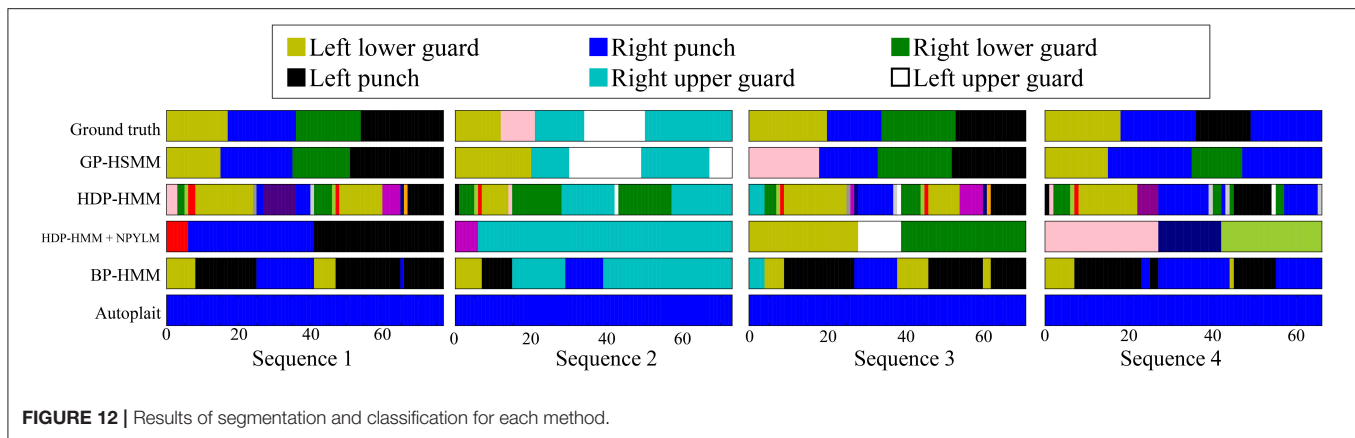
**FIGURE 12 |** Results of segmentation and classification for each method.

**TABLE 2 |** Segmentation accuracy of karate motions.

|  | Hamming distance | Precision | Recall | F-measure |
|---|---|---|---|---|
| GP-HSMM | 0.21 | 0.92 | 0.92 | 0.92 |
| HDP-HMM | 0.47 | 0.12 | 0.54 | 0.19 |
| HDP-HMM + NPYLM | 0.61 | 0.00 | 0.00 | 0.00 |
| BP-HMM | 0.49 | 0.13 | 0.23 | 0.16 |
| AutoPlait | 0.76 | 0.00 | 0.00 | 0.00 |

the ground truth boundary has no corresponding estimated boundary as shown in **Figure 9**(6), it was considered as FN. Conversely, if the estimated boundary has no corresponding ground truth boundary as shown in **Figure 9**(11), it was considered as FP. From these evaluations, the precision, recall, and F-measure of the segmentation are computed as follows:

$$P = \frac{N_{TP}}{N_{TP} + N_{FP}}, \qquad (12)$$

$$R = \frac{N_{TP}}{N_{TP} + N_{FN}}, \qquad (13)$$

$$F = \frac{2PR}{P + R}, \qquad (14)$$

where $N_{TP}$, $N_{FP}$, and $N_{TN}$ represent the number of points assessed as TP, FP, and FN. The F-measure of the segmentation was 0.81, and this fact indicates that GP-HSMM can estimate boundaries reasonably. This is because GP-HSMM estimates the length of segments as well as the classes of segments.

Moreover, **Figure 8** shows that most false segmentations are in sequence 3. This is because "up and down" and "toe touch" motions are included only in sequence 3, and GP-HSMM was not able to extract patterns that occur infrequently. However, this problem is not limited to GP-HSMM, and it is generally difficult for any learning method to extract infrequent patterns. The Hamming distance, which was computed only from sequence 1 and sequence 2, was 0.15. This result shows that GP-HSMM

can accurately estimate segments that appear multiple times in a sequence.

## 4.2. Segmentation of Karate Motion

We then applied our proposed method to more complex motion capture data, which consisted of the basic motions of karate (called kata in Japanese)[1] as shown in **Figure 10** from the motion capture library Mocapdata.com[2]. There are fixed motion patterns (punches or guards) in kata, and it is easy to form a ground truth for the segmentation. However, there might be shorter motion patterns, and GP-HSMM might be able to find those motion patterns if the number of classes is set to a larger number. Moreover, it is possible for GP-HSMM to discover patterns that cannot be labeled by humans, and GP-HSMM has the potential to analyze unlabeled time series data. However, in this experiment, we must evaluate the proposed method quantitatively, and fixed motion patterns (punches or guards) labeled by a human expert are used as ground truth. The type of kata we used was called heian 1, which is the most basic form of kata consisting of punches, lower guard, and upper guard (Tsuki, Gedanbarai, and Joudanuke in Japanese). **Figure 11** shows the basic movements used in heian 1. We divided this motion sequence into four parts, for use as four motion sequences to apply the blocked Gibbs sampler. Each motion sequence consisted of the following actions:

1. Left lower guard, right punch, right lower guard, and left punch.
2. Left lower guard, right upper guard, left upper guard, and right upper guard.
3. Left lower guard, right punch, right lower guard, and left punch.
4. Left lower guard, right punch, left punch, and right punch

By way of its preprocessing, as shown in **Figure 7**, the motion capture data was converted into motions with the body facing forward with a center of (0,0,0). To reduce computational cost, we downsampled the motion capture data from 30 frames per
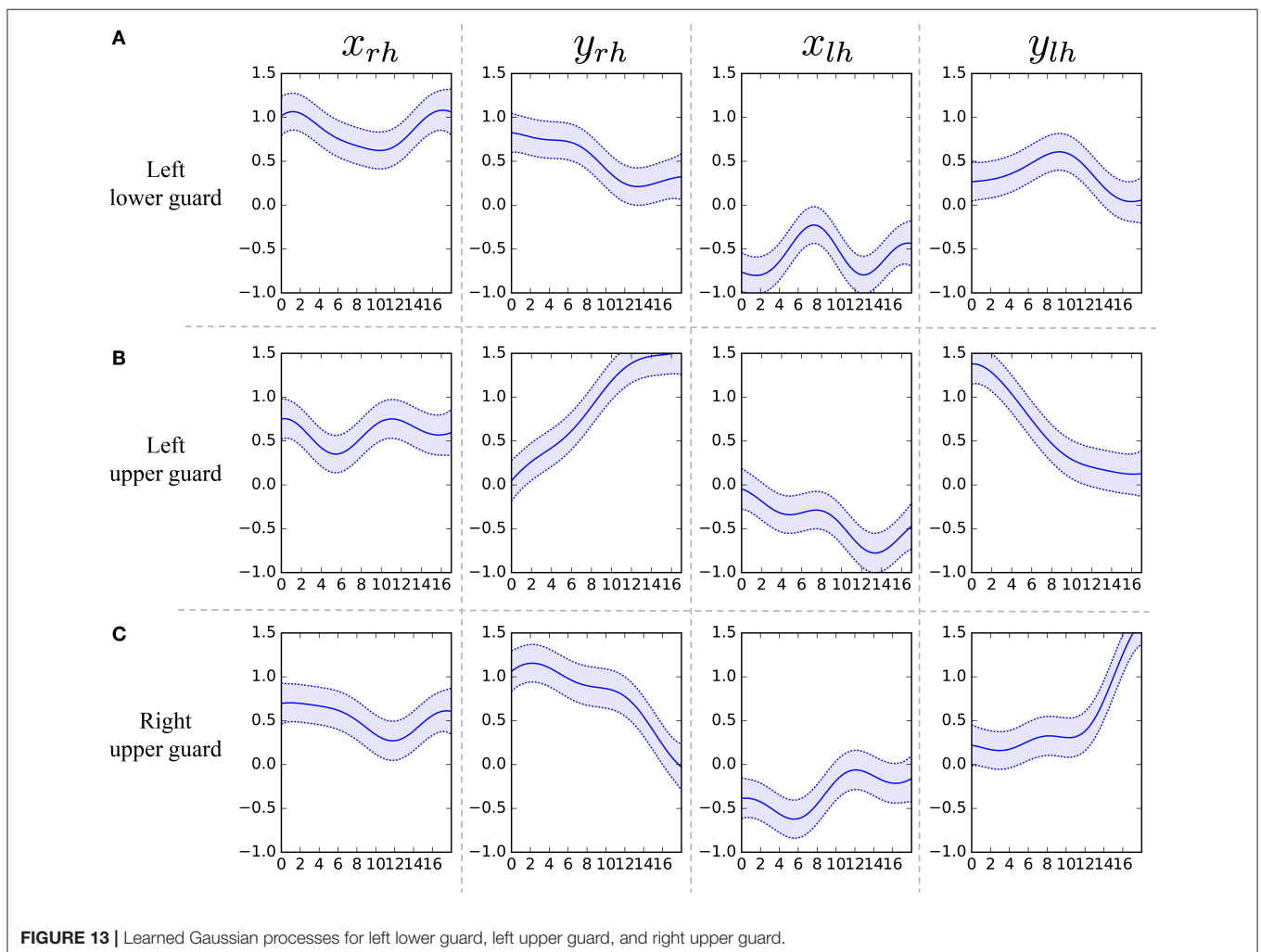
---

[1]https://mocapdata.blob.core.windows.net/freemotions/karate.zip
[2]http://www.mocapdata.com/

second to 15 frames per second, and used two-dimensional left-hand positions ($x_{lh}$, $y_{lh}$) and right-hand positions ($x_{rh}$, $y_{rh}$) in the frontal view, as shown in **Figure 7**. To compare our method with others, we used segmentation based on HDP-HMM (Beal et al., 2001) and segmentation based on NPYLM and HDP-HMM (Taniguchi and Nagasaka, 2011), where NPYLM (Mochihashi et al., 2009) divides sequences discretized by HDP-HMM. In addition, we compared our method with BP-HMM (Fox et al., 2011) and AutoPlait (Matsubara et al., 2014).

**Figure 12** shows the segmentation results. The horizontal axis represents the frame number, and the colors represent motion classes into which each segment was classified. The figure shows that HDP-HMM estimated shorter segments than the ground truth. This occurred because the emission distribution of HDP-HMM is a Gaussian distribution, which cannot represent continuous trajectories. Moreover, the result produced by segmentation, in which NPYLM divided sequences discretized by HDP-HMM, yielded longer segments. Moreover, NPYLM cannot extract fixed patterns of sequences. This is because the sequences discretized by HDP-HMM included noise and, therefore, NPYLM was unable to find a pattern in them.

It was also difficult for BP-HMM to estimate correct segments, and some shorter segments were present. Further, AutoPlait could not find any segments in the karate motion sequences. We believe this occurred because HMMs are too simple to model complex motions. On the contrary, we use Gaussian processes that make it possible to model complex sequences. **Table 2** shows the segmentation accuracy of each method. We considered the estimated boundary to be correct if it was within true boundary ± five frames. The F-measure of the proposed method was 0.92, which indicates that GP-HSMM can estimate boundaries

**TABLE 3 |** Computational time of each method.

|                | Time (s) |
|----------------|----------|
| GP-HSMM        | 248      |
| HDP-HMM        | 1.99     |
| HDP-HMM + NPYLM | 18.2    |
| BP-HMM         | 3.37     |
| AutoPlait      | 0.31     |



**FIGURE 13 |** Learned Gaussian processes for left lower guard, left upper guard, and right upper guard.

accurately. The results show that GP-HSMM outperforms the other methods. **Figure 13** shows the learned Gaussian process. $y_{rh}$ in **Figure 13A**, which represents the height of the left hand, is decreased, which indicates the motion where the left hand is dropped for the lower guard. In contrast, $y_{rh}$ in **Figure 13B** is increased, which indicates the motion where the left hand is raised for the upper guard. Conversely, $y_{lh}$ in **Figure 13C** is increased for the right upper guard. From this result, we can see that characteristics of motions can be learned by Gaussian processes.

Moreover, the motions were classified into seven classes, although we set the number of classes to eight. This result indicates that the number of classes can be estimated to a certain extent, if a number closer to the correct number is given. However, a smaller number leads to under-segmentation and misclassification, and a much larger number leads to over-segmentation. This is a limitation of the current GP-HSMM, and we believe it can be solved by introducing a non-parametric Bayesian model.

Computational cost is another limitation of GP-HSMM. **Table 3** shows the computational time required to segment karate motion. HMM-based methods such as HDP-HMM, BP-HMM, and AutoPlait are relatively faster. In particular, AutoPlait is the fastest because it uses a single scan algorithm proposed in (Matsubara et al., 2014) to find boundaries, and it has been demonstrated that AutoPlait can detect meaningful patterns from large datasets. In contrast, our proposed GP-HSMM is much slower than other methods, and cannot process such large datasets. This is another limitation of the proposed method.

## 5. CONCLUSION

In this paper, we proposed a method for motion segmentation based on a hidden semi-Markov model (HSMM) with a Gaussian process (GP) emission distribution. By employing HSMM, segment classes and their lengths can be estimated. Moreover, a forward filtering-backward sampling algorithm is used to estimate the parameters of GP-HSMM; this makes it possible to efficiently search for all possible segment lengths and classes. The experimental results showed that the proposed method can accurately segment motion capture data. Although motions that occurred in the sequences a single time were difficult to segment correctly, motions that occurred a few times could be segmented with higher accuracy.

However, some issues remain in the current GP-HSMM. The most significant problem is that GP-HSMM requires the number of classes to be specified in advance. We believe this value can be estimated by utilizing a non-parametric Bayesian model. We are planning to introduce a stick-breaking process as a prior distribution of the transition matrix, and beam sampling for parameter estimation; these techniques are utilized in Beal et al. (2001). Another problem is computational cost. The computational cost to learn a Gaussian process is $O(n^3)$, where $n$ denotes the number of data points classified in the GP. To overcome this problem, efficient computation methods have been proposed (Nguyen-Tuong et al., 2009; Okadome et al., 2014), and we will consider introducing these methods into GP-HSMM.

## AUTHOR CONTRIBUTIONS

ToN, TaN, DM, IK, and HA conceived of the presented idea. ToN, TaN, and DM developed the theory and performed the computations. IK and HA verified the theory and the analytical methods. ToN wrote the manuscript with support from TaN and MK. IK and HA supervised the project. All authors discussed the results and contributed to the final manuscript.

## ACKNOWLEDGMENTS

## REFERENCES

Argall, B. D., Chernova, S., Veloso, M., and Browning, B. (2009). A survey of robot learning from demonstration. *Robot. Auton. Sys.* 57, 469–483. doi: 10.1016/j.robot.2008.10.024

Beal, M. J., Ghahramani, Z., and Rasmussen, C. E. (2001). "The infinite hidden markov model," in *Advances in Neural Information Processing Systems* (Vancouver, BC), 577–584.

CMU (2009). *CMU Graphics Lab Motion Capture Database.* Available online at: http://mocap.cs.cmu.edu/

Fod, A., Matarić, M. J., and Jenkins, O. C. (2002). Automated derivation of primitives for movement classification. *Auton. Rob.* 12, 39–54. doi: 10.1023/A:1013254724861

Fox, E. B., Sudderth, E. B., Jordan, M. I., and Willsky, A. S. (2007). *The Sticky hdp-hmm: Bayesian Nonparametric Hidden Markov Models with Persistent States.* Technical Report, MIT Laboratory for Information and Decision Systems.

Fox, E. B., Sudderth, E. B., Jordan, M. I., and Willsky, A. S. (2011). Joint modeling of multiple related time series via the beta process. *arXiv preprint arXiv:1111.4226.*

Goldwater, S. (2006). *Nonparametric Bayesian Models of Lexical Acquisition.* Ph.D. thesis: Brown University, Providence, RI.

Gräve, K., and Behnke, S. (2012). "Incremental action recognition and generalizing motion generation based on goal-directed features," in *IEEE/RSJ International Conference on Intelligent Robots and Systems* (Vilamoura), 751–757.

Lin, J. F.-S., Karg, M., and Kulić, D. (2016). Movement primitive segmentation for human motion modeling: A framework for analysis. *IEEE Trans. Hum. Mach. Sys.* 46, 325–339. doi: 10.1109/THMS.2015.2493536

Lin, J. F.-S., and Kulić, D. (2012). "Segmenting human motion for automated rehabilitation exercise analysis," in *Annual International Conference of the IEEE Engineering in Medicine and Biology Society* (San Diego, CA), 2881–2884.

Lioutikov, R., Neumann, G., Maeda, G., and Peters, J. (2015). "Probabilistic segmentation applied to an assembly task," in *IEEE-RAS International Conference on Humanoid Robots* (Seoul), 533–540.

Manschitz, S., Kober, J., Gienger, M., and Peters, J. (2015). Learning movement primitive attractor goals and sequential skills from kinesthetic demonstrations. *Robot. Auton. Sys.* 74, 97–107. doi: 10.1016/j.robot.2015.07.005

Matsubara, Y., Sakurai, Y., and Faloutsos, C. (2014). "Autoplait: utomatic mining of co-evolving time sequences," in *ACM SIGMOD International Conference on Management of Data* (Snowbird, UT), 193–204.

Mochihashi, D., Yamada, T., and Ueda, N. (2009). "Bayesian Unsupervised Word Segmentation with Nested Pitman-Yor Language Modeling," in *Joint*

Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing Vol. 1 (Singapore), 100–108.

Nguyen-Tuong, D., Peters, J. R., and Seeger, M. (2009). "Local gaussian process regression for real time online model learning," in *Advances in Neural Information Processing Systems* (Vancouver, BC), 1193–1200.

Okadome, Y., Urai, K., Nakamura, Y., Yomo, T., and Ishiguro, H. (2014). Adaptive lsh based on the particle swarm method with the attractor selection model for fast approximation of gaussian process regression. *Art. Life Robot.* 19, 220–226. doi: 10.1007/s10015-014-0161-1

Shiratori, T., Nakazawa, A., and Ikeuchi, K. (2004). "Detecting dance motion structure through music analysis," in *IEEE International Conference on Automatic Face and Gesture Recognition* (Seoul), 857–862.

Takano, W., and Nakamura, Y. (2016). Real-time unsupervised segmentation of human whole-body motion and its application to humanoid robot acquisition of motion symbols. *Robot. Auton. Sys.* 75, 260–272. doi: 10.1016/j.robot.2015.09.021

Taniguchi, T. and Nagasaka, S. (2011). "Double articulation analyzer for unsegmented human motion using pitman-yor language model and infinite hidden markov model," in *IEEE/SICE International Symposium on System Integration* (Kyoto), 250–255.

Uchiumi, K., Hiroshi, T., and Mochihashi, D. (2015). "Inducing Word and Part-of-Speech with Pitman-Yor Hidden Semi-Markov Models," in *Joint Conference of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing* (Beijing), 1774–1782.

Wachter, M., and Asfour, T. (2015). "Hierarchical segmentation of manipulation actions based on object relations and motion characteristics," in *International Conference on Advanced Robotics* (Istanbul), 549–556.

**Conflict of Interest Statement:** The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.