

Segmenting Continuous Motions with Hidden Semi-Markov Models and Gaussian Processes

Tomoaki Nakamura^{1,*}, Takayuki Nagai¹, Daichi Mochihashi²,
Ichiro Kobayashi³, Hideki Asoh⁴, and Masahide Kaneko¹

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*

Correspondence*:
Tomoaki Nakamura
tnakmaura@uec.ac.jp

2 ABSTRACT

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.

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

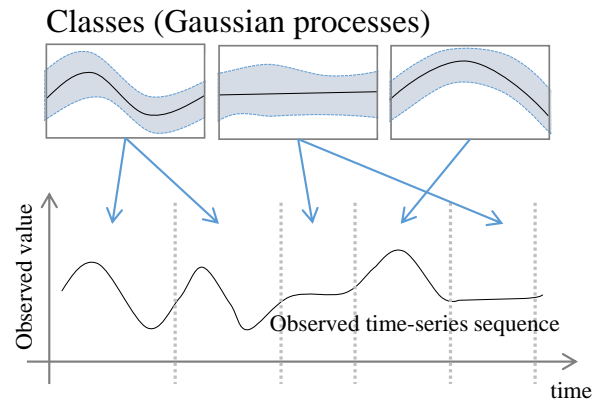


Figure 1. Overview of the proposed GP-HSMM.

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. Fig. 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 (Manschitz et al., 2015; Gräve and Behnke, 2012). 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

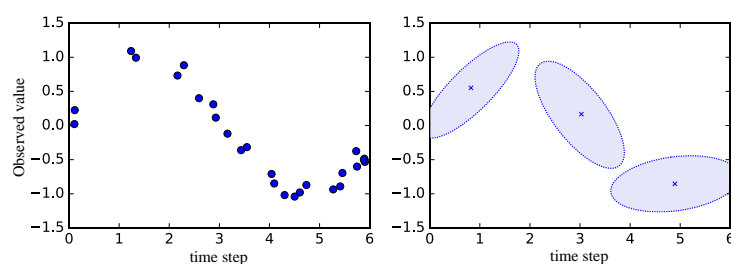


Figure 2. Example of representation of time series data by HMM. Left: Data points for learning HMM. Right: Mean and standard deviation learned by HMM.

48 motions using motion primitives that are learned by a supervised hidden Markov model. In these studies,
 49 the motions are segmented and labeled in advance. However, we consider that it is difficult to segment
 50 and label all possible motion primitives.

51 Additionally, some studies have proposed unsupervised motion segmentation. However, these studies
 52 rely on heuristics. For instance, Wächter et al. have proposed a method to segment human manipulation
 53 motions based on contact relations between the end-effectors and objects in a scene (Wächter and Asfour,
 54 2015); in their method, the points at which the end-effectors make contact with an object are determined
 55 as boundaries of motions. We believe this method works well in limited scenes; however, there are
 56 many motions, such as gestures and dances, in which objects are not manipulated. Lioutikov et al.
 57 proposed unsupervised segmentation; however, to reduce computational costs, this technique requires the
 58 possible boundary candidates between motion primitives to be specified in advance (Lioutikov et al.,
 59 2015). Therefore, the segmentation depends on those candidates, and motions cannot be segmented
 60 correctly if the correct candidates are not selected. In contrast, our proposed method does not require
 61 such candidates; all possible cutting points are considered by use of forward filtering-backward sampling,
 62 which uses the principles of dynamic programming. In some methods (Fod et al., 2002; Shiratori et al.,
 63 2004; Fod et al., 2002; Lin and Kulić, 2012), motion features (such as the zero velocity of joint angles) are
 64 used for motion segmentation. However, these features cannot be applied to all motions. Takano et al. use
 65 the error between actual movements and predicted movements as the criteria for specifying boundaries
 66 (Takano and Nakamura, 2016). However, the threshold must be manually tuned according to the motions
 67 to be segmented. Moreover, they used an HMM that is a stochastic model. We consider such an assumption
 68 to be unnatural from the viewpoint of stochastic models, and boundaries should be determined based on
 69 a stochastic model. In our proposed method, we do not use such heuristics and assumptions, and instead
 70 formulate the segmentation based on a stochastic model.

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

81 Matsubara et al. proposed an unsupervised segmentation method called AutoPlait (Matsubara et al.,
 82 2014). This method uses multiple HMMs, each of which represents a fixed pattern; moreover, transitions
 83 between the HMMs are allowed. Therefore, time series data is segmented at points at which the state is
 84 changed to another HMM's state. However, we believe that HMMs are too simple to represent complicated
 85 sequences such as motions. Fig. 2 illustrates an example of representation of time series data by HMM.
 86 The graph on the right in Fig. 2 represents the mean and standard deviation learned by HMM from data
 87 points shown in the graph on the left. HMM represents time series data using only the mean and standard
 88 deviation; therefore, details of time series data can be lost. Therefore, we use Gaussian processes, which
 89 are non-parametric methods that can represent complex time series data.

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

105 Taniguchi et al. proposed a method to divide continuous sequences into segments by utilizing NPYLM
 106 (Taniguchi and Nagasaka, 2011). In their method, continuous sequences are discretized and converted into
 107 discrete-valued sequences using the infinite hidden Markov model (Fox et al., 2007). The discrete-valued
 108 sequences are then divided into segments by using NPYLM. In this method, motions can be recognized by
 109 the learned model, but cannot be generated naively because they are discretized. Moreover, segmentation
 110 based on NPYLM does not work well if errors occur in the discretization step.

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

3 GAUSSIAN PROCESS-HIDDEN SEMI-MARKOV MODEL

120 Fig. 3 shows a graphical representation of the proposed GP-HSMM. In this figure, $c_j (j = 1, 2, \dots, J)$
 121 denotes classes of segments, and each segment x_j is generated by a Gaussian process, with parameters

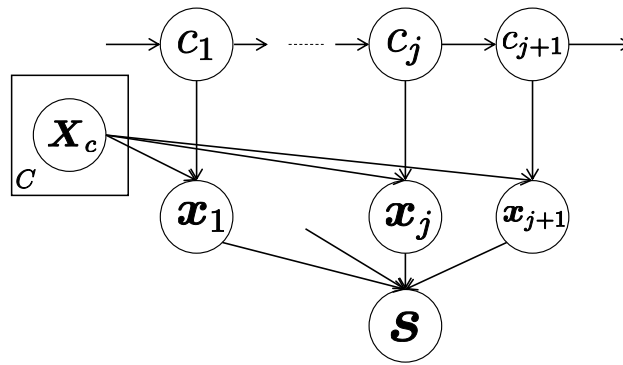


Figure 3. Graphical representation of the proposed GP-HSMM.

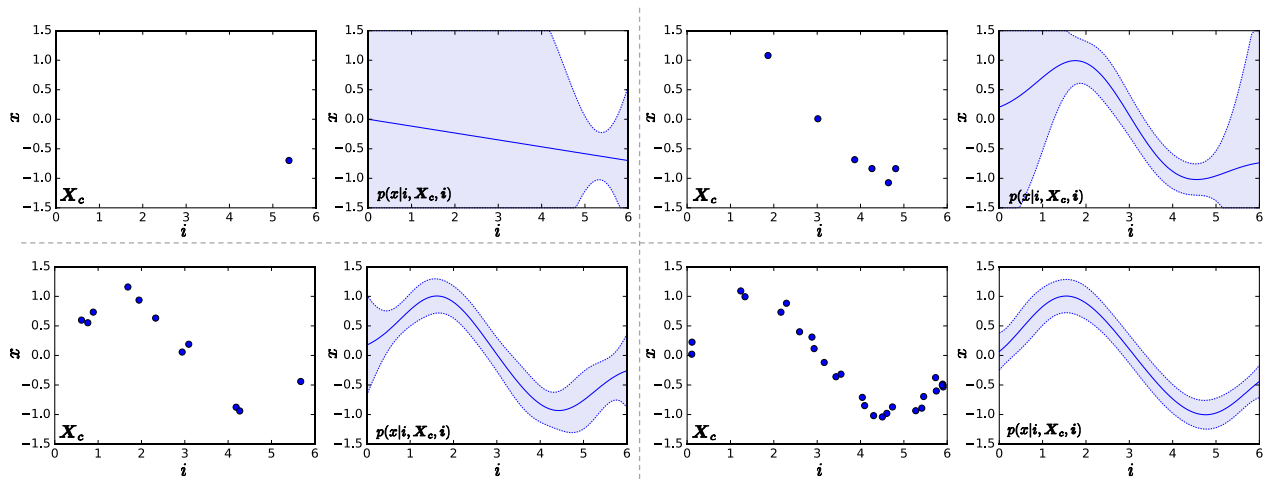


Figure 4. Examples of Gaussian processes. Left graph in each pair of graphs represents learning data points (i, \mathbf{X}_c) . Right graph shows the learned probabilistic distribution $p(x|i, \mathbf{X}_c, i)$; the solid line represents the mean, and the blue region represents the range of standard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

denoted by \mathbf{X}_c and given by the following generative process:

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

$$x_j \sim \mathcal{GP}(x|\mathbf{X}_c), \quad (2)$$

where \mathbf{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, \mathbf{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.

$$p(x|i, \mathbf{X}_c, i) \propto \mathcal{N}(\mathbf{k}^T \mathbf{C}^{-1} \mathbf{i}, c - \mathbf{k}^T \mathbf{C}^{-1} \mathbf{k}), \quad (3)$$

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

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

β is a hyperparameter that represents noise in the observation. In Eq. (3), \mathbf{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:

$$k(i_p, i_q) = \theta_0 \exp(-\frac{1}{2} \theta_1 \|i_p - i_q\|^2) + \theta_2 + \theta_3 i_p i_q, \quad (5)$$

where θ_* represents parameters of the kernel. Fig. 4 shows examples of Gaussian processes. The left graph in each pair of graphs represents learning data points (i, \mathbf{X}_c) , and the right graph shows the learned probabilistic distribution $p(x|i, \mathbf{X}_c, i)$. 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 $\mathbf{x} = (x_0, x_1, \dots)$, we assume that each dimension is generated independently, and a predictive distribution $\mathcal{GP}(\mathbf{x}|\mathbf{X}_c)$ is computed as follows:

$$\begin{aligned} \mathcal{GP}(\mathbf{x}|\mathbf{X}_c) &= p(x_0|i, \mathbf{X}_{c,0}, i_c) \\ &\quad \times p(x_1|i, \mathbf{X}_{c,1}, i_c) \\ &\quad \times p(x_2|i, \mathbf{X}_{c,2}, i_c) \cdots \end{aligned} \quad (6)$$

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 $\mathbf{x}_{nj}(j = 1, 2, \dots, J_n)$ in observed sequence \mathbf{s}_n are then removed from the learning data, and parameter \mathbf{X}_c of the Gaussian process and transition probability $P(c|c')$ of HSMM are updated. Segments $\mathbf{x}_{nj}(j = 1, 2, \dots, J_n)$ and their classes $c_{nj}(j = 1, 2, \dots, J_n)$ are then estimated as follows:

$$(\mathbf{x}_{n1}, \dots, \mathbf{x}_{nJ_n}), (c_{n1}, \dots, c_{nJ_n}) \sim P(\mathbf{X}, \mathbf{c}|\mathbf{s}_n), \quad (7)$$

where \mathbf{X} is a set of segments into which \mathbf{s}_n is divided, and \mathbf{c} denotes classes of the segments. To carry out this sampling efficiently, the probability of all possible segments \mathbf{X} and classes \mathbf{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 \mathbf{x}_{nj} and c_{nj} , parameter \mathbf{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 Eq. (10).

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:     Delete segments  $\mathbf{x}_{nj}$  from  $\mathbf{X}_{c_{nj}}$ 
7:   end for
8:
9:   // Sample segments and their classes
10:   $(\mathbf{x}_{n1}, \dots, \mathbf{x}_{nJ_n}), (c_{n1}, \dots, c_{nJ_n}) \sim P(\mathbf{X}, \mathbf{c} | \mathbf{s}_n)$ 
11:
12:  for  $j = 1$  to  $J_n$  do
13:     $N_{c_{nj}} + +$ 
14:     $N_{c_{nj}, c_{n,j+1}} + +$ 
15:    Add segments  $\mathbf{x}_{nj}$  into  $\mathbf{X}_{c_{nj}}$ 
16:  end for
17: end for

```

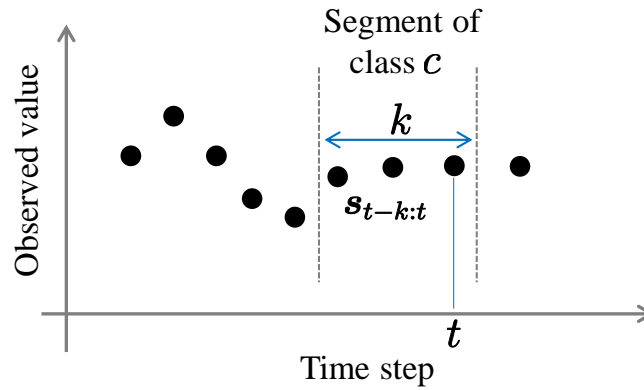


Figure 5. A segment whose probability is computed during forward filtering.

160 3.2.2 Forward filtering-backward sampling

161 In this study, we regard segments and their classes as latent variables that are sampled by forward
 162 filtering-backward sampling. In forward filtering, as shown in Fig. 5, the probability that k samples $\mathbf{s}_{t-k:t}$
 163 prior to time step t in observed sequence \mathbf{s} form a segment, and that the resulting segment belongs to class
 164 c , is computed as follows:

$$\begin{aligned}
 \alpha[t][k][c] &= P(\mathbf{s}_{t-k:t} | \mathbf{X}_c) \\
 &\times \sum_{k'=1}^K \sum_{c'=1}^C p(c|c') \alpha[t-k][k'][c'],
 \end{aligned} \tag{8}$$

165 where C and K denote the number of classes and the maximum length of segments, respectively.
 166 $P(\mathbf{s}_{t-k:t} | \mathbf{X}_c)$ represents the probability that $\mathbf{s}_{t-k:t}$ is generated from a class c ; this is computed as follows:

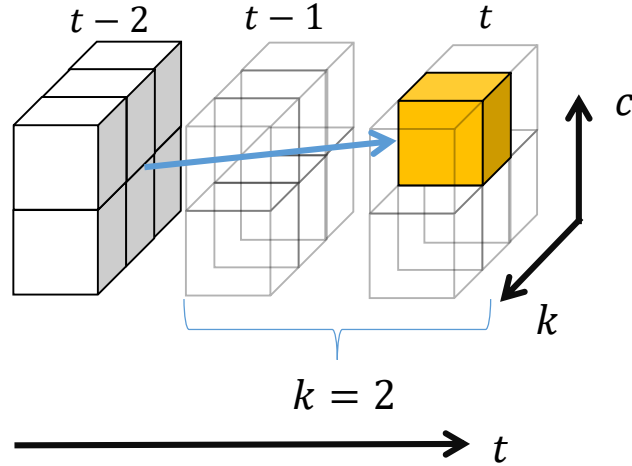
$$P(\mathbf{s}_{t-k:t} | \mathbf{X}_c) = \mathcal{GP}(\mathbf{s}_{t-k:t} | \mathbf{X}_c) P_{len}(k | \lambda). \tag{9}$$

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  $\alpha[t][k][c]$ 
6:     end for
7:   end for
8: end for
9:
10: // Backward sampling
11:  $t = T, j = 1$ 
12: while  $t > 0$  do
13:    $k, c \sim \alpha[t][k][c]$ 
14:    $\mathbf{x}_j = \mathbf{s}_{t-k:t}$ 
15:    $c_j = c$ 
16:    $t = t - k$ 
17:    $j = j + 1$ 
18: end while
19: return  $(\mathbf{x}_{J_n}, \mathbf{x}_{J_n-1}, \dots, \mathbf{x}_1), (c_{J_n}, c_{J_n-1}, \dots, c_1)$ 

```

**Figure 6.** Recursive computation in forward filtering.

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

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

170 where $N_{c'}$ and $N_{c'c}$ denote the number of segments whose classes are c' and the number of transitions from
 171 c' to c , respectively, and k' and c' respectively denote the length and class of the segment preceding $\mathbf{s}_{t-k:t}$;
 172 these are marginalized out in Eq. (8). Moreover, $\alpha[t][k][*] = 0$ if $t - k < 0$, and $\alpha[0][0][*] = 1.0$. All

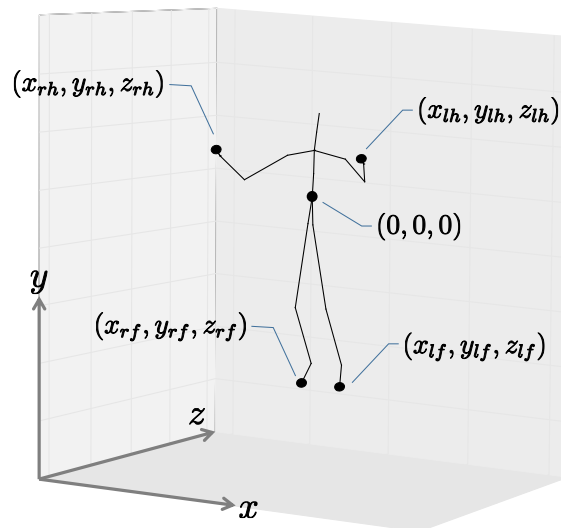


Figure 7. Coordinate system used in the experiments.

elements of $\alpha[*][*][*]$ in Eq. (8) can be recursively computed from $\alpha[1][1][*]$ by dynamic programming. Fig. 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 α . 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. Fig. 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

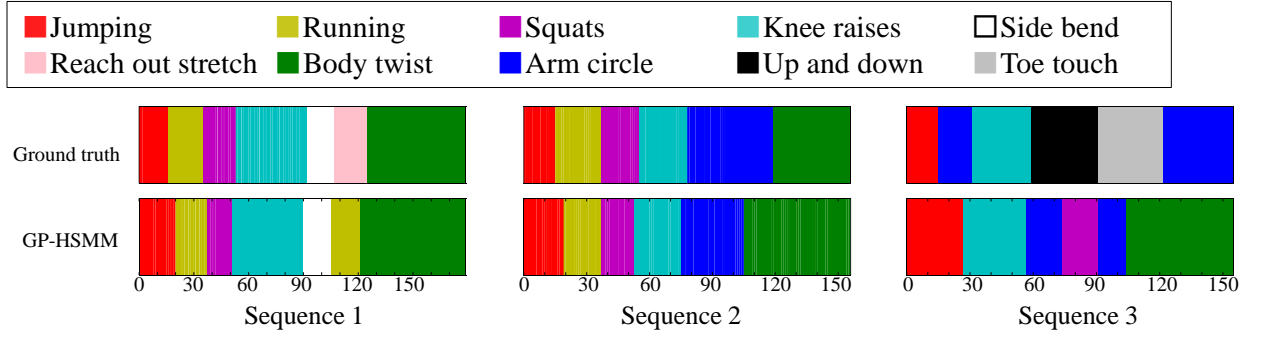


Figure 8. Segmentation results of CMU motion capture data.

Table 1. Segmentation accuracy of CMU motion capture data.

Hamming distance	Precision	Recall	F-measure
0.33	0.81	0.81	0.81

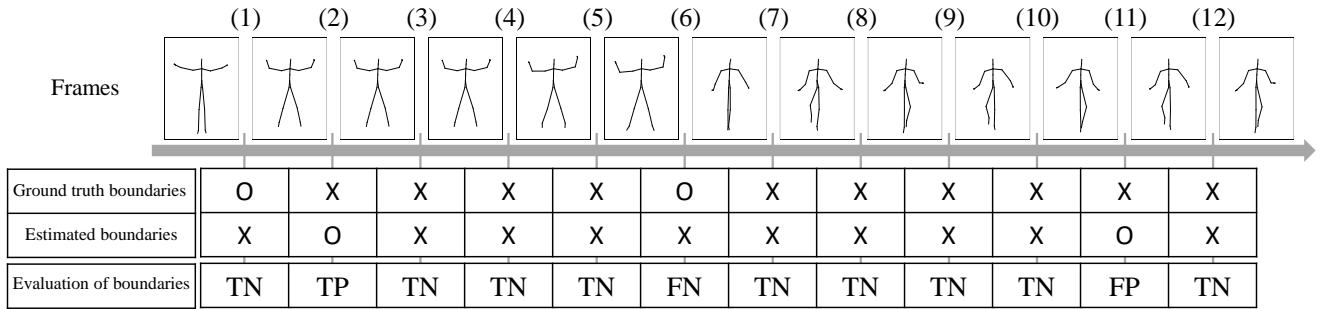


Figure 9. Example of segmentation evaluation. Estimated boundaries are evaluated as true positive (TP), true negative (TN), false positive (FP), or false negative (FN).

dimensional vectors: $(x_{lh}, y_{lh}, x_{rh}, y_{rh}, x_{lf}, y_{lf}, x_{rf}, y_{rf})$. Because GP-HSMM requires the number of classes to be specified in advance, we set it to 8.

Fig. 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:

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

where c and \bar{c} represent sequences of estimated motion classes and true motion classes, $D(c, \bar{c})$ is the Hamming distance between two sequences, and $|\bar{c}|$ 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 Fig. 8 shows that there are no shorter segments. We also computed the precision, recall,

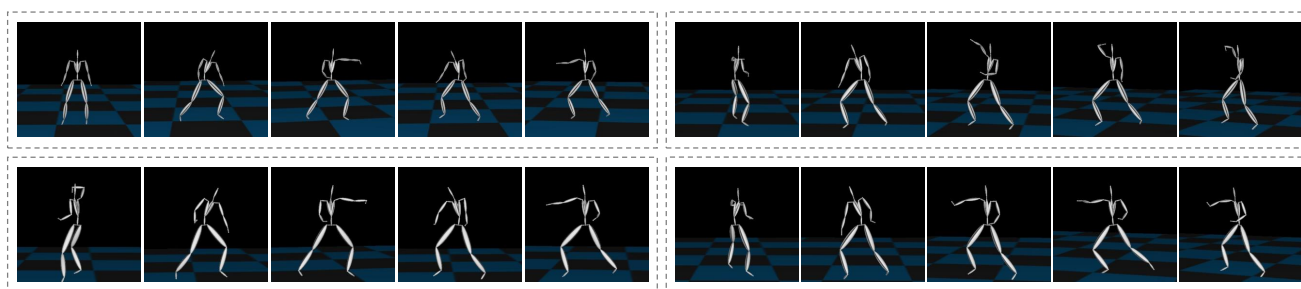


Figure 10. Motion capture data of karate motions.

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

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). Fig. 9 shows an example of segmentation evaluation. We considered the estimated boundary to be TP if it was within true boundary \pm four frames, as shown in Fig. 9(2). If the ground truth boundary has no corresponding estimated boundary as shown in Fig. 9(6), it was considered as FN. Conversely, if the estimated boundary has no corresponding ground truth boundary as shown in Fig. 9(11), it was considered as FP. From these evaluations, the precision, recall, and F-measure of the segmentation are computed as follows:

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

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

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

where N_{TP} , N_{FP} , and N_{FN} 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, Fig. 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.

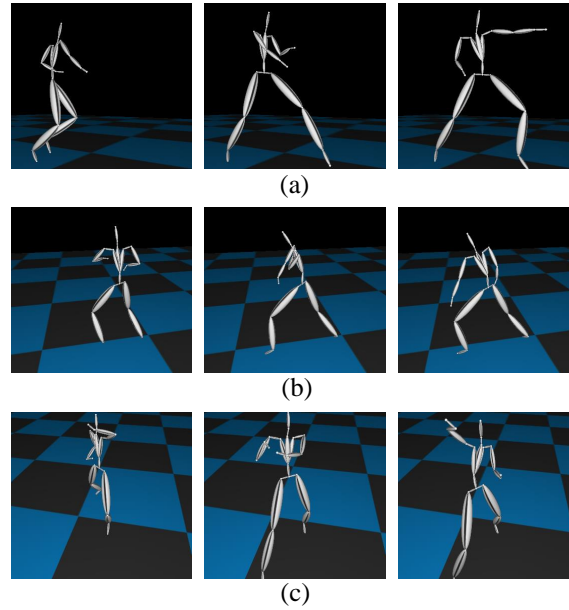


Figure 11. Basic motions in Kata: (a) Left punch. (b) Left lower guard. (c) Right upper guard.

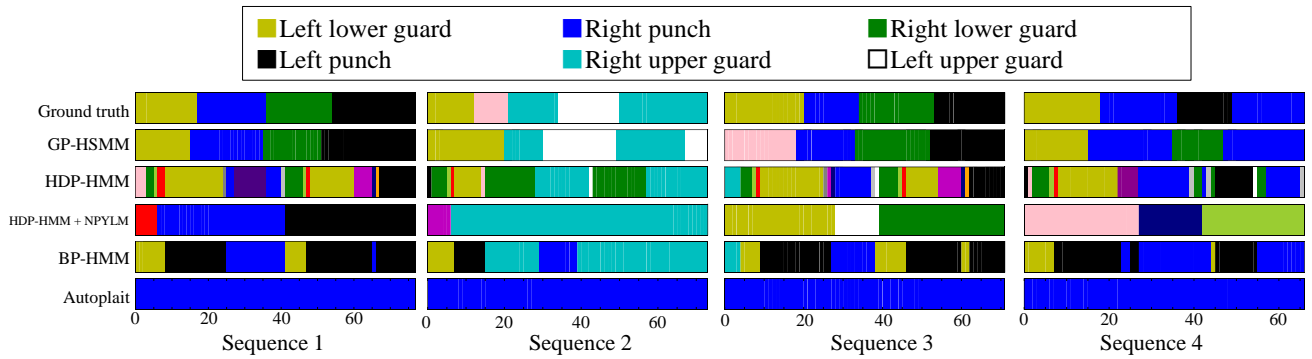


Figure 12. Results of segmentation and classification for each method.

212 4.2 Segmentation of karate motion

213 We then applied our proposed method to more complex motion capture data, which consisted of the basic
 214 motions of karate (called kata in Japanese)¹ (Fig. 10) from the motion capture library Mocapdata.com².
 215 There are fixed motion patterns (punches or guards) in kata, and it is easy to form a ground truth for the
 216 segmentation. However, there might be shorter motion patterns, and GP-HSMM might be able to find
 217 those motion patterns if the number of classes is set to a larger number. Moreover, it is possible for
 218 GP-HSMM to discover patterns that cannot be labeled by humans, and GP-HSMM has the potential to
 219 analyze unlabeled time series data. However, in this experiment, we must evaluate the proposed method
 220 quantitatively, and fixed motion patterns (punches or guards) labeled by a human expert are used as ground
 221 truth. The type of kata we used was called heian 1, which is the most basic form of kata consisting of
 222 punches, lower guard, and upper guard (Tsuki, Gedanbarai, and Joudanuke in Japanese). Fig. 11 shows

¹ http://mocapdata.com/product.cgi?product_id=10019

² <http://www.mocapdata.com/>

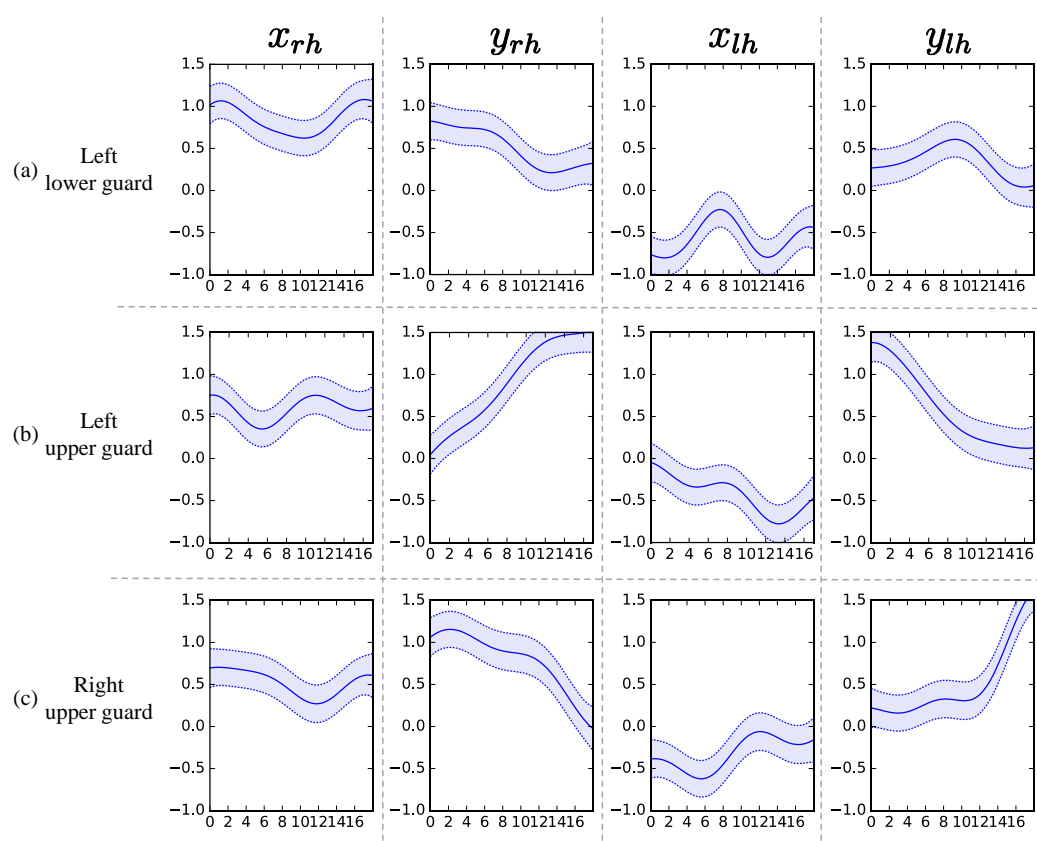


Figure 13. Learned Gaussian processes for left lower guard, left upper guard, and right upper guard.

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 Fig. 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 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 Fig. 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).

Fig. 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

Table 3. Computational time of each method.

	time (sec)
GP-HSMM	248
HDP-HMM	1.99
HDP-HMM + NPYLM	18.2
BP-HMM	3.37
AutoPlait	0.31

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 \pm five frames. The F-measure of the proposed method was 0.92, which indicates that GP-HSMM can estimate boundaries accurately. The results show that GP-HSMM outperforms the other methods. Fig. 13 shows the learned Gaussian process. y_{rh} in Fig. 13(a), 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 Fig. 13(b) is increased, which indicates the motion where the left hand is raised for the upper guard. Conversely, y_{lh} in Fig. 13(c) 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.

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

ACKNOWLEDGEMENT

283 This work was supported by JST CREST Grant Number JPMJCR15E3 and JSPS KAKENHI Grant
284 Number JP17K12758.

REFERENCES

- 285 Argall, B. D., Chernova, S., Veloso, M., and Browning, B. (2009). A survey of robot learning from
286 demonstration. *Robotics and autonomous systems* 57, 469–483
- 287 Beal, M. J., Ghahramani, Z., and Rasmussen, C. E. (2001). The infinite hidden markov model. In
288 *Advances in neural information processing systems*. 577–584
- 289 CMU (2009). CMU graphics lab motion capture database. <http://mocap.cs.cmu.edu/>
- 290 Fod, A., Matarić, M. J., and Jenkins, O. C. (2002). Automated derivation of primitives for movement
291 classification. *Autonomous Robots* 12, 39–54
- 292 Fox, E. B., Sudderth, E. B., Jordan, M. I., and Willsky, A. S. (2007). The sticky hdp-hmm: Bayesian
293 nonparametric hidden markov models with persistent states. *MIT Laboratory for Information and*
294 *Decision Systems Technical Report 2777*
- 295 Fox, E. B., Sudderth, E. B., Jordan, M. I., and Willsky, A. S. (2011). Joint modeling of multiple related
296 time series via the beta process. *arXiv preprint arXiv:1111.4226*
- 297 Goldwater, S. (2006). *Nonparametric Bayesian Models of Lexical Acquisition*. Ph.D. thesis, Brown
298 University
- 299 Gräve, K. and Behnke, S. (2012). Incremental action recognition and generalizing motion generation
300 based on goal-directed features. In *IEEE/RSJ International Conference on Intelligent Robots and*
301 *Systems*. 751–757
- 302 Lin, J. F.-S., Karg, M., and Kulić, D. (2016). Movement primitive segmentation for human motion
303 modeling: A framework for analysis. *IEEE Transactions on Human-Machine Systems* 46, 325–339
- 304 Lin, J. F.-S. and Kulić, D. (2012). Segmenting human motion for automated rehabilitation exercise
305 analysis. In *Annual International Conference of the IEEE Engineering in Medicine and Biology Society*.
306 2881–2884
- 307 Lioutikov, R., Neumann, G., Maeda, G., and Peters, J. (2015). Probabilistic segmentation applied to an
308 assembly task. In *IEEE-RAS International Conference on Humanoid Robots*. 533–540
- 309 Manschitz, S., Kober, J., Gienger, M., and Peters, J. (2015). Learning movement primitive attractor goals
310 and sequential skills from kinesthetic demonstrations. *Robotics and Autonomous Systems* 74, 97–107

- 311 Matsubara, Y., Sakurai, Y., and Faloutsos, C. (2014). Autoplait: Automatic mining of co-evolving time
312 sequences. In *ACM SIGMOD International Conference on Management of Data*. 193–204
- 313 Mochihashi, D., Yamada, T., and Ueda, N. (2009). Bayesian Unsupervised Word Segmentation with
314 Nested Pitman-Yor Language Modeling. In *Joint Conference of the 47th Annual Meeting of the ACL
315 and the 4th International Joint Conference on Natural Language Processing*. vol. 1, 100–108
- 316 Nguyen-Tuong, D., Peters, J. R., and Seeger, M. (2009). Local gaussian process regression for real time
317 online model learning. In *Advances in Neural Information Processing Systems*. 1193–1200
- 318 Okadome, Y., Urai, K., Nakamura, Y., Yomo, T., and Ishiguro, H. (2014). Adaptive lsh based on the
319 particle swarm method with the attractor selection model for fast approximation of gaussian process
320 regression. *Artificial Life and Robotics* 19, 220–226
- 321 Shiratori, T., Nakazawa, A., and Ikeuchi, K. (2004). Detecting dance motion structure through music
322 analysis. In *IEEE International Conference on Automatic Face and Gesture Recognition*. 857–862
- 323 Takano, W. and Nakamura, Y. (2016). Real-time unsupervised segmentation of human whole-body
324 motion and its application to humanoid robot acquisition of motion symbols. *Robotics and Autonomous
325 Systems* 75, 260–272
- 326 Taniguchi, T. and Nagasaka, S. (2011). Double articulation analyzer for unsegmented human motion using
327 pitman-yor language model and infinite hidden markov model. In *IEEE/SICE International Symposium
328 on System Integration*. 250–255
- 329 Uchiumi, K., Hiroshi, T., and Mochihashi, D. (2015). Inducing Word and Part-of-Speech with Pitman-
330 Yor Hidden Semi-Markov Models. In *Joint Conference of the 53rd Annual Meeting of the Association
331 for Computational Linguistics and the 7th International Joint Conference on Natural Language
332 Processing*. 1774–1782
- 333 Wächter, M. and Asfour, T. (2015). Hierarchical segmentation of manipulation actions based on object
334 relations and motion characteristics. In *International Conference on Advanced Robotics*. 549–556