

---

# Assisting Movement Training and Execution with Visual and Haptic Feedback

Marco Ewerton<sup>1,\*</sup>, David Rother<sup>1</sup>, Jakob Weimar<sup>1</sup>, Gerrit Kollegger<sup>2</sup>, Josef Wiemeyer<sup>2</sup>, Jan Peters<sup>1,3</sup> and Guilherme Maeda<sup>1,4</sup>

<sup>1</sup>*Intelligent Autonomous Systems group, department of Computer Science, Technische Universität Darmstadt, Darmstadt, Germany*

<sup>2</sup>*Institute of Sport Science, Technische Universität Darmstadt, Darmstadt, Germany*

<sup>3</sup>*Max Planck Institute for Intelligent Systems, Tübingen, Germany*

<sup>4</sup>*ATR Computational Neuroscience Labs, department of Brain-Robot Interface, Kyoto, Japan*

Correspondence\*:

Intelligent Autonomous Systems group, department of Computer Science, Technische Universität Darmstadt, Hochschulstr. 10, 64289 Darmstadt, Germany  
ewerton@ias.tu-darmstadt.de

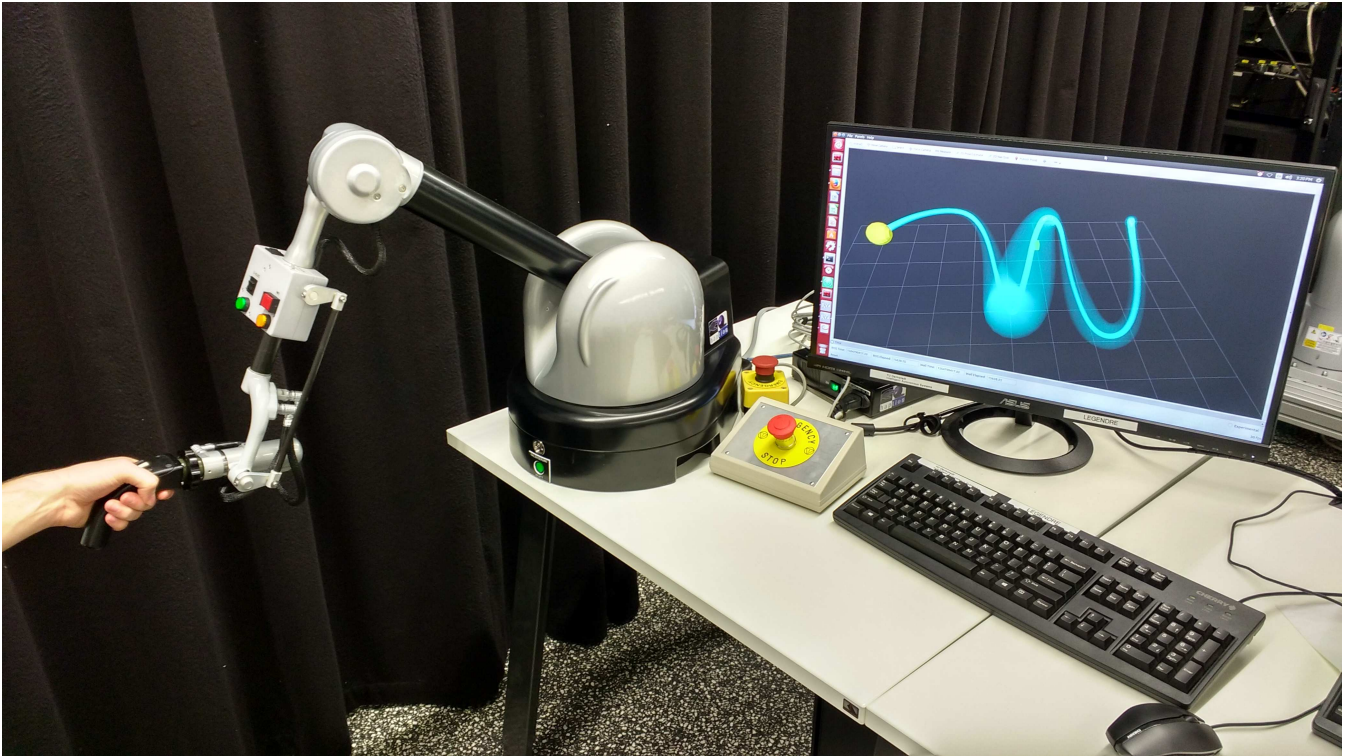
## 2 ABSTRACT

3 In the practice of motor skills in general, errors in the execution of movements may go unnoticed  
4 when a human instructor is not available. In this case, a computer system or robotic device able  
5 to detect movement errors and propose corrections would be of great help. This paper addresses  
6 the problem of how to detect such execution errors and how to provide feedback to the human to  
7 correct his/her motor skill using a general, principled methodology based on imitation learning.  
8 The core idea is to compare the observed skill with a probabilistic model learned from expert  
9 demonstrations. The intensity of the feedback is regulated by the likelihood of the model given the  
10 observed skill. Based on demonstrations, our system can, for example, detect errors in the writing  
11 of Japanese characters with multiple strokes. Moreover, by using a haptic device, the Haption  
12 Virtuose 6D, we demonstrate a method to generate haptic feedback based on a distribution  
13 over trajectories, which could be used as an auxiliary means of communication between an  
14 instructor and an apprentice. Additionally, given a performance measurement, the haptic device  
15 can help the human discover and perform better movements to solve a given task. In this case,  
16 the human first tries a few times to solve the task without assistance. Our framework, in turn,  
17 uses a reinforcement learning algorithm to compute haptic feedback, which guides the human  
18 towards better solutions.

19 **Keywords:** shared autonomy, HRI, movement primitives, reinforcement learning, policy search, cooperation, robotics, interaction

## 1 INTRODUCTION

20 In the absence of an instructor, errors in the execution of movements by a person trying to learn a new motor  
21 skill, such as calligraphy, for example, may go unnoticed. To counter this problem, we propose recording  
22 demonstrations of a motor skill provided by an instructor and processing them such that someone practicing



**Figure 1.** Human manipulating a haptic device, the Haption Virtuose 6D. In our experiments, the haptic device assists the movements of the human by providing force feedback which is inversely proportional to the standard deviation of a distribution over trajectories (example is shown on the computer screen).

23 that motor skill in the absence of the instructor can have the correctness of his/her trials automatically  
24 assessed and receive feedback based on the demonstrations.

25 More precisely, our system aligns demonstrated trajectories in space and time and computes a probability  
26 distribution over them. Often, demonstrations may have been executed at different speeds. In order to  
27 extract the underlying shape of the movement from multiple trajectories, it is thus necessary to time-align  
28 these trajectories. In some cases, such as writing characters, the scale and the absolute position of the  
29 movements are not as relevant as their shape, justifying the necessity of addressing space-alignment in our  
30 framework as well.

31 When a new trajectory is executed, our system aligns the observations in space and time with the post-  
32 processed demonstrations and computes the probability of each of the positions of this new trajectory  
33 under the distribution over the demonstrations. The computed probabilities provide a way of assessing the  
34 correctness of each position of the new trajectory.

35 Based on this assessment, our system can generate visual or haptic feedback. We demonstrate the  
36 generation of visual feedback with the task of assisting the practice of writing Japanese characters on a  
37 monitor with a computer mouse. The generation of haptic feedback is demonstrated in an experiment with  
38 a haptic device, the Haption Virtuose 6D (see Figure 1). Our system gives haptic feedback to the user in the  
39 form of forces that constrain his/her movements when manipulating the haptic device, which can be seen  
40 as a form of guiding virtual fixtures (Rosenberg, 1992). The produced force is perpendicular to the mean  
41 trajectory of the distribution and its intensity is inversely proportional to the standard deviation along the  
42 distribution, as detailed in Section 4.

43    There are situations where the initial demonstrations are not enough to successfully accomplish a task,  
44 but it is possible to define performance measurements accounting for certain objectives. Examples of such  
45 a situation could be found in a teleoperation task where the user perception and motor capabilities do  
46 not enable him/her to succeed. Such a task could be for instance telemanipulating a robot arm to move  
47 an object from a start position to an end position while avoiding obstacles. In such a task, a user can  
48 easily hit obstacles or fail to reach objects of interest. However, it may be possible to define performance  
49 measurements based on the positions of objects in the environment of the teleoperated robot. These  
50 positions could be computed from information provided by sensors in that environment. The framework  
51 presented in this paper deals with these situations by applying reinforcement learning to adapt the original  
52 distribution over trajectories. The adapted distribution is then used to guide the user towards a better  
53 solution to the task.

54    In general, the problem of finding a distribution over trajectories that avoid obstacles and pass through  
55 positions of interest involves multiple optimization subproblems. Tuning the hyperparameters of the reward  
56 function to satisfy all the objectives may be time-consuming and may not produce the desired results. For  
57 this reason, our proposed framework includes a novel reinforcement learning algorithm that makes use  
58 of a parametric representation of trajectories and identifies how relevant each policy parameter is to each  
59 of the objectives of the task. By identifying how relevant each policy parameter is to each objective, it  
60 is possible to achieve effective policies with simpler reward functions, one for each objective, instead of  
61 a single reward function with different user-defined weights for each objective. Moreover, it is possible  
62 to optimize each objective sequentially, exploring different values of the parameters that matter for that  
63 objective and preserving the uncertainty about the other parameters.

64    In summary, this paper presents a new framework to assist humans in training and executing movements  
65 by providing visual and haptic feedback to the user. This feedback can be given based on a probability  
66 distribution over expert demonstrations or based on an optimized distribution learned from a few non-  
67 expert demonstrations and performance criteria. By including methods for time and space-alignment of  
68 trajectories, this framework can potentially be applied to a large range of motor skills as long as the shape  
69 of the movement is critical, not its speed. In this work, our framework has been applied to the learning  
70 of Japanese characters and to teleoperation. As a secondary contribution, this paper presents a novel  
71 reinforcement learning algorithm for problems involving multiple objectives, which are often encountered  
72 in teleoperation scenarios.

## 2 RELATED WORK

73 This section primarily describes related work on techniques to assess the correctness of human motion  
74 and provide feedback to the user. It briefly introduces related work on the required components used for  
75 modeling the human demonstrations.

### 76 2.1 Human Motion Assessment and Feedback to the User

77    With similar goals as in our work, Solis et al. (2002) presented a method to teach users how to write  
78 characters using a haptic interface. In their method, characters are modeled with Hidden Markov Models  
79 (HMMs) with discrete hidden states and discrete observations. The system recognizes online what character  
80 the user intends to write and applies a proportional derivative (PD) controller with fixed gains to restrict the  
81 user to move along the trajectory that corresponds to the recognized character. Differently, in our work,  
82 the gains of the haptic device are adapted as a function of the user's deviation with respect to the model  
83 learned from expert demonstrations or through reinforcement learning. Adaptive gains allow for practicing

84 motor skills with multiple correct possibilities of execution, in case there is not a single correct trajectory.  
85 Also, it allows for regulating the stiffness of the robot to impose different levels of precision at different  
86 parts of the movement.

87 Parisi et al. (2016) proposed a “multilayer learning architecture with incremental self-organizing networks”  
88 to give the user real-time visual feedback during the execution of movements, e.g. powerlifting exercises. In  
89 our work, we have not addressed real-time visual feedback so far, although we do address real-time haptic  
90 feedback. On the other hand, our framework can deal with movements with different absolute positions and  
91 scales when producing visual feedback. By disabling this preprocessing, it would be possible to generate  
92 real-time visual feedback as well.

93 Kowsar et al. (2016) presented a workflow to detect anomalies in weight training exercises. In their  
94 work, movement repetitions are segmented based on the acceleration along an axis in space. A probability  
95 distribution over a number of time-aligned repetitions is built. Then, based on this distribution, movement  
96 segments can be deemed correct or incorrect. Our approach focuses rather on correcting movements with  
97 respect to their shape or position in space, not on correcting acceleration patterns.

98 A variable impedance controller based on an estimation of the stiffness of the human arm was proposed  
99 by Tsumugiwa et al. (2002). This controller enabled a robot to assist humans in calligraphic tasks. In the  
100 cited work, the tracked trajectories were not learned from demonstrations.

101 Our work is in line with approaches that aim to assist learning with demonstrations. Raiola et al. (2015),  
102 for instance, used probabilistic virtual guides learned from demonstrations to help humans manipulate  
103 a robot arm. In another related work, Soh and Demiris (2015) presented a system that learns from  
104 demonstrations how to assist humans using a smart wheelchair.

105 Visual, auditory and haptic feedback modalities have been successfully used for motor learning in the  
106 fields of sport and rehabilitation (Sigrist et al., 2013). Our method to provide visual feedback to the user,  
107 detailed in Section 3.4, is, for instance, similar in principle to bandwidth feedback. This sort of feedback  
108 means that the user only receives feedback when the movement error exceeds a certain threshold and it has  
109 been shown to be effective in rehabilitation (Timmermans et al., 2009). The work here presented relates and  
110 can potentially complement previous research on bandwidth feedback in the sense that our threshold is not  
111 constant, but depends on a probability distribution over trajectories. Our approach may find applications in  
112 tasks where it is desirable to give the user more freedom of movement around a certain position and less  
113 freedom around a different position or where multiple variations of movements are considered correct.

114 Ernst and Banks (2002) have demonstrated that maximum-likelihood estimation describes the way  
115 humans combine visual and haptic perception. The estimation of a certain environmental property that  
116 results from the combination of visual and haptic stimuli presents lower variance than estimations based  
117 only on visual or haptic stimuli. When the visual stimulus is noise-free, users tend to rely more on vision  
118 to perform their estimation. On the other hand, when the visual stimulus is noisy, users tend to rely more  
119 on haptics. Therefore, users may profit from multimodal feedback to learn a new motor skill. In our  
120 experimental section, we provide haptic feedback to users to help them perform a teleoperation task in a  
121 virtual environment. The findings in Ernst and Banks (2002) indicate that haptic feedback also helps users  
122 perceive some aspects of the task that they could not perceive only from visual stimuli, which could help  
123 them learn how to better solve the task without assistance next time. The usefulness of haptic feedback  
124 to learn motor skills is also demonstrated in Kümmel et al. (2014), where robotic haptic guidance has  
125 been shown to induce long-lasting changes in golf swing movements. The work here presented offers an

126 algorithmic solution to the acquisition of policies and control of a robotic device that could be applied to  
127 help humans learn and retain motor skills.

128 In contrast to most of the work on haptic feedback for human motor learning, our method modulates the  
129 stiffness of the haptic device according to demonstrations and uses reinforcement learning to improve upon  
130 the demonstrated movements. Those features may be interesting as a means of communication between an  
131 expert and an apprentice or patient and to enable improvement of initial demonstrations.

## 132 **2.2 Learning and Adapting Models from Demonstrations**

133 An essential component of this work is to construct a model from expert demonstrations, which is then  
134 queried at runtime to evaluate the performance of the user. One recurrent issue when building models  
135 from demonstration is the problem of handling the variability of phases (i.e. the speed of the execution)  
136 of different movements. Listgarten et al. (2004) proposed the Continuous Profile Model (CPM), which  
137 can align multiple continuous time series. It assumes that each continuous time series is a non-uniformly  
138 subsampled, noisy and locally rescaled version of a single latent trace. The model is similar to a Hidden  
139 Markov Model (HMM). The hidden states encode the corresponding time step of the latent trace and  
140 a rescaling factor. The CPM has been successfully applied to align speech data and data sets from an  
141 experimental biology laboratory.

142 Coates et al. (2008) augmented the model of Listgarten et al. (2004) by additionally learning the  
143 dynamics of the controlled system in the vicinity of the intended trajectory. With this modification, their  
144 model generates an ideal trajectory that not only is similar to the demonstrations but also obeys the  
145 system's dynamics. Moreover, differently from Listgarten et al. (2004), their algorithm to time-align the  
146 demonstrations and to determine an ideal trajectory relies both on an EM algorithm and on Dynamic Time  
147 Warping (Sakoe and Chiba, 1978). With this approach, they were able to achieve autonomous helicopter  
148 aerobatics after training with sub-optimal human expert demonstrations.

149 The same method was used by Van Den Berg et al. (2010) to extract an ideal trajectory from multiple  
150 demonstrations. The demonstrations were, in this case, movements of a surgical robot operated by a human  
151 expert.

152 Similarly to Coates et al. (2008); Van Den Berg et al. (2010), our system uses Dynamic Time Warping  
153 (DTW) to time-align trajectories. While DTW usually aligns pairs of temporal sequences, in Section 3.2  
154 we present a solution for aligning multiple trajectories. An alternative solution was presented by Sanguansat  
155 (2012), however, it suffers from scalability issues because distances need to be computed between every  
156 point of every temporal sequence.

157 Differences in the scale and shape of movements must also be addressed to account for the variability  
158 in human demonstrations. In practice, for tasks such as writing, we want our system to be invariant to  
159 the scale of the movements of different demonstrations. The analysis of the difference between shapes  
160 is usually addressed by Procrustes Analysis (Goodall, 1991). The output of this analysis is the affine  
161 transformation that maps one of the inputs to best match the other input, while the residual is quantified  
162 as the effective distance (deformation) between the shapes. As the analysis consists of computing such  
163 transformations in relation to the centroid, Procrustes Analysis provides a global, average assessment  
164 and has found applications in tasks of trajectory and transfer learning (Bocsi et al., 2013; Holladay and  
165 Srinivasa, 2016; Makondo et al., 2015) and manipulation (Collet et al., 2009). While this seems the most  
166 natural solution to our problem of aligning shapes, we noticed that it is not suitable for detecting anomalies.  
167 In fact, in the writing task, we are interested in finding the “outliers” that can be indicated to the human as

168 erroneous strokes. However, Procrustes Analysis aligns the shapes globally such that the positions of the  
 169 centroids are inappropriately biased towards such outliers. In Sections 3.1.1 and 3.1.2 we describe our own  
 170 alignment method that is suited for detecting particular errors with the introduction of a few heuristics.

### 3 PROCESSING DEMONSTRATIONS AND ASSESSING THE CORRECTNESS OF OBSERVED TRAJECTORIES

171 Assuming the availability of expert demonstrations, the workflow of our proposed method is the following:  
 172 First, the expert demonstrations are aligned in space and time and a probability distribution over these  
 173 demonstrations is computed. Afterward, a user tries to perform the motor task. The movements of the  
 174 user are also aligned in space and time with the demonstrations. Based on the probability distribution  
 175 over the demonstrations, our system highlights which parts of the user's movements need improvement.  
 176 A way of translating a distribution over trajectories into haptic feedback is presented later in Section 4.  
 177 A novel reinforcement learning algorithm to help the user achieve good movements according to certain  
 178 performance criteria without good demonstrations available is presented in Section 5.

#### 179 3.1 Rescaling and Repositioning

180 In assessing the correctness of individual executions of a motor skill, it is often not important what the  
 181 absolute position of the sequence of movements is, e.g. in weightlifting or gymnastics. In some situations,  
 182 it is also not of crucial importance what the scale of the movements is as long as they keep their relative  
 183 proportions, e.g. in drawing or calligraphy. Therefore, our system rescales all trajectories, both the ones  
 184 demonstrated by a human expert and the ones performed by a user practicing a motor skill. Moreover, all  
 185 trajectories are repositioned in such a way that the first position of the reference stroke is at the origin of  
 186 the coordinate system. In practice, each stroke composing a motor skill is used once as the reference for  
 187 rescaling and repositioning. For each reference stroke, a different score and visual feedback are computed.  
 188 The best score and the respective feedback are presented to the user. This procedure enables our algorithm  
 189 to present meaningful feedback to the user regardless the location of his/her errors. In this section, our  
 190 method for rescaling and repositioning is explained for two dimensions ( $x$  and  $y$ ) and exemplified with the  
 191 task of writing Japanese characters. This method can nevertheless be extended in a straightforward manner  
 192 for more than two dimensions.

##### 193 3.1.1 Rescaling

First, the system computes

$$\Delta x_{\text{ref}} = \max_t x_{\text{ref}}(t) - \min_t x_{\text{ref}}(t), \quad (1)$$

$$\Delta y_{\text{ref}} = \max_t y_{\text{ref}}(t) - \min_t y_{\text{ref}}(t), \quad (2)$$

194 where  $t$  indexes each time step,  $\max_t x_{\text{ref}}(t)$  is the maximum  $x$  coordinate of the reference stroke,  
 195  $\min_t x_{\text{ref}}(t)$  is the minimum  $x$  coordinate of the reference stroke, and similarly for  $\max_t y_{\text{ref}}(t)$  and  
 196  $\min_t y_{\text{ref}}(t)$ .

197 Subsequently, a rescaling factor  $\alpha$  is given by

$$\alpha = \begin{cases} \frac{1}{\Delta x_{\text{ref}}} & \text{if } \Delta x_{\text{ref}} \geq \Delta y_{\text{ref}}, \\ \frac{1}{\Delta y_{\text{ref}}} & \text{otherwise.} \end{cases} \quad (3)$$

198 The characters are written on a square window with side equal to 1. The rescaling factor  $\alpha$  expresses the  
 199 ratio between the constant 1 and the width  $\Delta x_{\text{ref}}$  or height  $\Delta y_{\text{ref}}$  of the reference stroke. If  $\Delta x_{\text{ref}} \geq \Delta y_{\text{ref}}$ ,  
 200 the width is used to compute  $\alpha$ . Otherwise, the height is used. Some strokes are much larger in width than  
 201 in height or vice versa. Therefore, this way of computing the rescaling factor selects the width or the height  
 202 of the reference stroke according to which one will lead to the smallest amount of rescaling. For example,  
 203 the characters depicted in Figure 2(a) will be rescaled according to the width of the first stroke of each of  
 204 them respectively, resulting in characters whose first stroke has width equal to 1.

205 The rescaling factor can also be written as

$$\alpha = \frac{x_{i,\text{rescaled}}(t) - \min_{\{j,k\}} x_j(k)}{x_i(t) - \min_{\{j,k\}} x_j(k)} = \frac{y_{i,\text{rescaled}}(t) - \min_{\{j,k\}} y_j(k)}{y_i(t) - \min_{\{j,k\}} y_j(k)}, \quad (4)$$

206 where both  $t$  and  $k$  are time step indexes, while the indexes  $i$  and  $j$  represent the strokes of a character. Here,  
 207  $x_{i,\text{rescaled}}(t) - \min_{\{j,k\}} x_j(k)$  is the difference between the  $x$  coordinate at time step  $t$  of stroke  $i$  after  
 208 rescaling and the minimum  $x$  coordinate of the character. The term  $x_i(t) - \min_{\{j,k\}} x_j(k)$  represents the  
 209 corresponding difference before rescaling. Equation (4) also includes similar terms for the  $y$  coordinates.  
 210 Therefore, after rescaling, the difference between the  $x$  coordinate of the position at time step  $t$  and the  
 211 minimum  $x$  coordinate is  $\alpha$  times this difference before rescaling, and similarly for the  $y$  coordinate. Thus  
 212 this rescaling keeps the proportion between the width and the height of the character.

Rearranging the terms of (4) leads to

$$x_{i,\text{rescaled}}(t) = \min_{\{j,k\}} x_j(k) + \left( x_i(t) - \min_{\{j,k\}} x_j(k) \right) \alpha, \quad (5)$$

$$y_{i,\text{rescaled}}(t) = \min_{\{j,k\}} y_j(k) + \left( y_i(t) - \min_{\{j,k\}} y_j(k) \right) \alpha, \quad (6)$$

213 which is how the coordinates of the rescaled version of a character are computed. Figure 2(a) shows two  
 214 demonstrations of the same character and Figure 2(b) shows the result of rescaling these characters.

### 215 3.1.2 Repositioning

In order to reposition a character such that the first position of the reference stroke is  $(x = 0, y = 0)$ , our system simply computes

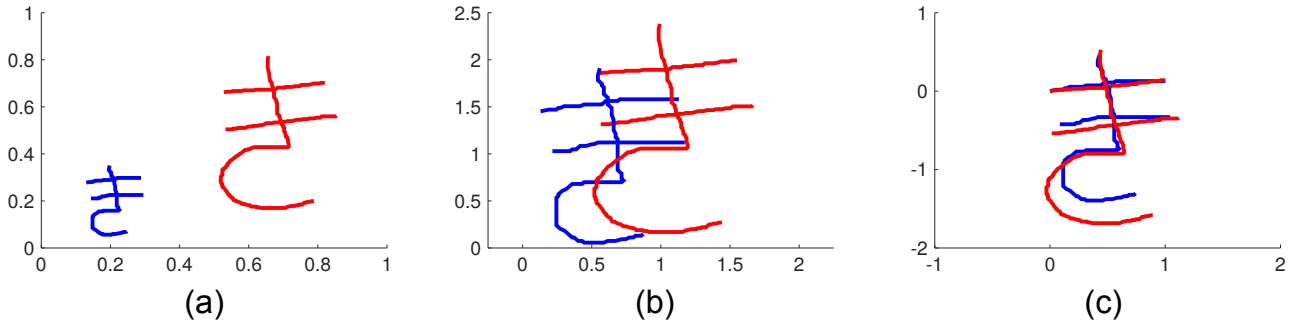
$$x_{i,\text{repositioned}}(t) = x_i(t) - x_{\text{ref}}(t = 1), \quad (7)$$

$$y_{i,\text{repositioned}}(t) = y_i(t) - y_{\text{ref}}(t = 1), \quad (8)$$

216 where  $x_i(t)$  and  $y_i(t)$  are the original coordinates of stroke  $i$  at time step  $t$ ,  $x_{i,\text{repositioned}}(t)$  and  
 217  $y_{i,\text{repositioned}}(t)$  are the coordinates of stroke  $i$  at time step  $t$  of the character after repositioning,  $x_{\text{ref}}(t = 1)$   
 218 and  $y_{\text{ref}}(t = 1)$  are the coordinates of the reference stroke at the first time step. Figure 2(c) shows two  
 219 demonstrations of the same character after rescaling and repositioning.

## 220 3.2 Time Alignment

221 The time alignment of all the demonstrations and of the user's movements is achieved in our system by  
 222 using Dynamic Time Warping (Sakoe and Chiba, 1978). Each stroke of an execution of a motor skill is  
 223 time-aligned with respect to the corresponding stroke of other executions of that same motor skill.



**Figure 2.** Rescaling and repositioning different executions of a motor skill. In this example, the motor skill is writing a Japanese character. **(a)** Two demonstrations of a Japanese character. **(b)** After rescaling both characters. **(c)** After repositioning the characters such that the first position of the first stroke is  $(x = 0, y = 0)$ . The first stroke is the reference stroke in this case.

224 Suppose two corresponding strokes need to be time-aligned. Let us represent these strokes by  $\tau_1$  and  
 225  $\tau_2$ , which are sequences of Cartesian coordinates from time step  $t = 1$  until time step  $t = T_1$  and  $t = T_2$ ,  
 226 respectively. Here,  $T_1$  and  $T_2$  represent the last time step of  $\tau_1$  and  $\tau_2$ , respectively.

First, the Euclidean distance  $D(i, j)$  between position at  $t = i$  of  $\tau_1$  and position at  $t = j$  of  $\tau_2$  is computed for all time steps of both strokes, i.e.

$$D(i, j) = \|\tau_1(i) - \tau_2(j)\|, \quad (9)$$

$$\forall i \in \{1, 2, \dots, T_1\}, \forall j \in \{1, 2, \dots, T_2\}.$$

Subsequently, assuming that the first position of  $\tau_1$  corresponds to the first position of  $\tau_2$ , the accumulated cost  $A(i, j)$  of associating  $\tau_1(i)$  with  $\tau_2(j)$  is computed according to

$$A(1, 1) = D(1, 1), \quad (10)$$

$$A(i, 1) = D(i, 1) + A(i - 1, 1), \quad (11)$$

$$A(1, j) = D(1, j) + A(1, j - 1), \quad (12)$$

$$A(i, j) = D(i, j) + \min \{A(i - 1, j), A(i - 1, j - 1), A(i, j - 1)\}. \quad (13)$$

227 Once the matrix of accumulated costs  $A$  has been determined, a path  $p$  can be computed that indicates  
 228 how each trajectory should progress in time such that the minimum total cost is achieved. This path is  
 229 computed backward in time in a dynamic programming fashion, as detailed in Algorithm 1.

230 The time warped versions of trajectories  $\tau_1$  and  $\tau_2$ , denoted by  $\tau'_1$  and  $\tau'_2$ , are computed with Algorithm 2.  
 231

232 Algorithms 1 and 2 represent a common form of DTW which aligns pairs of temporal sequences.  
 233 Algorithm 3 shows our proposed extension of DTW for time-aligning multiple temporal sequences. It  
 234 works as follows: Trajectories  $\tau_1$  and  $\tau_2$  are time-aligned with DTW, resulting in  $\tau'_1$  and  $\tau'_2$ . Then  $\tau'_2$  and  
 235  $\tau_3$  are time-aligned. Subsequently, the same warping applied to  $\tau'_2$  is also applied to  $\tau'_1$ . The algorithm  
 236 proceeds like that until  $\tau_n$ , always warping previous trajectories as well. For  $n$  trajectories, DTW needs to  
 237 be computed  $n - 1$  times and the computation of the distance matrix  $D$  remains the same as in the original  
 238 DTW. Figure 3 exemplifies the time-alignment of multiple trajectories.



---

**Algorithm 1** Path Search

---

```

1: procedure PATH( $T_1, T_2, \mathbf{A}$ )
2:    $k \leftarrow 1$ 
3:    $i \leftarrow T_1$ 
4:    $j \leftarrow T_2$ 
5:    $\mathbf{p}(k) \leftarrow (i, j)$ 
6:   while  $i \neq 1$  or  $j \neq 1$  do
7:     if  $i = 1$  then
8:        $j \leftarrow j - 1$ 
9:     else if  $j = 1$  then
10:       $i \leftarrow i - 1$ 
11:     else
12:       if  $\mathbf{A}(i - 1, j) = \min \{ \mathbf{A}(i - 1, j), \mathbf{A}(i - 1, j - 1), \mathbf{A}(i, j - 1) \}$  then
13:          $i \leftarrow i - 1$ 
14:       else if  $\mathbf{A}(i, j - 1) = \min \{ \mathbf{A}(i - 1, j), \mathbf{A}(i - 1, j - 1), \mathbf{A}(i, j - 1) \}$  then
15:          $j \leftarrow j - 1$ 
16:       else
17:          $i \leftarrow i - 1$ 
18:          $j \leftarrow j - 1$ 
19:       end if
20:     end if
21:      $k \leftarrow k + 1$ 
22:      $\mathbf{p}(k) \leftarrow (i, j)$ 
23:   end while
24:   return  $\mathbf{p}$ 
25: end procedure

```

---



---

**Algorithm 2** Warping for a Pair of Trajectories

---

```

1: procedure PAIRWARP( $\mathbf{p}, \tau_1, \tau_2$ )
2:    $t \leftarrow 0$ 
3:   for  $k = \mathbf{p}.\text{Length} \rightarrow 1$  do
4:      $t \leftarrow t + 1$ 
5:      $(i, j) \leftarrow \mathbf{p}(k)$ 
6:      $\tau'_1(t) \leftarrow \tau_1(i)$ 
7:      $\tau'_2(t) \leftarrow \tau_2(j)$ 
8:   end for
9:   return  $\tau'_1, \tau'_2$ 
10: end procedure

```

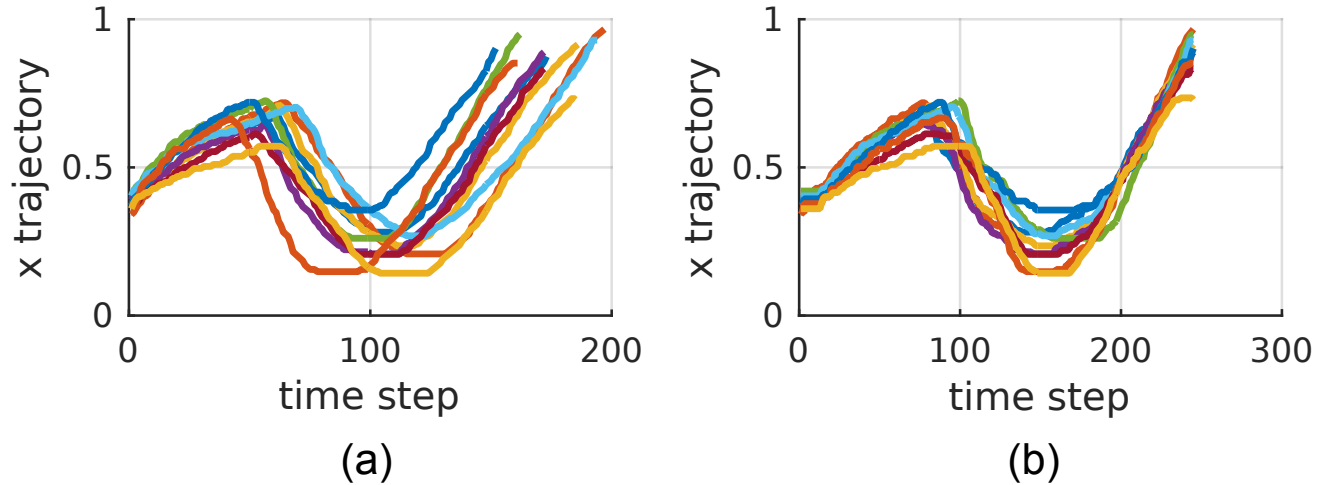
---

239 **3.3 Distribution over Trajectories**

240 In order to create a distribution over trajectories, we use the framework of Probabilistic Movement  
241 Primitives (Paraschos et al., 2013). Probabilistic Movement Primitives (ProMPs) allow for representing  
242 each trajectory with a relatively small number of parameters. A distribution over trajectories can then be  
243 computed by integrating out those parameters.

244 More precisely, in this framework, each trajectory  $\tau$  with a certain duration  $T$  is approximated by a  
245 weighted sum of  $N$  normalized Gaussian basis functions evenly distributed along the time axis. This  
246 approximation can be represented by

$$\tau = \Psi w + \epsilon, \tag{14}$$



**Figure 3.**  $x$  trajectories of corresponding strokes of multiple instances of a Japanese character. (a) Before time alignment. (b) After time alignment using DTW and our extension to deal with multiple trajectories.

---

**Algorithm 3** Warping for Multiple Trajectories

---

```

1: procedure MULTIPLEWARP( $\tau_1, \tau_2, \dots, \tau_n$ )
2:   for  $l = 1 \rightarrow n - 1$  do
3:      $(\mathbf{p}, \tau_l, \tau_{l+1}) \leftarrow \text{DTW}(\tau_l, \tau_{l+1})$ 
4:     for  $m = 1 \rightarrow l - 1$  do
5:        $t \leftarrow 0$ 
6:       for  $k = \mathbf{p}.\text{Length} \rightarrow 1$  do
7:          $t \leftarrow t + 1$ 
8:          $(i, j) \leftarrow \mathbf{p}(k)$ 
9:          $\tau_m(t) \leftarrow \tau_m(i)$ 
10:      end for
11:    end for
12:  end for
13:  return  $\tau_1, \tau_2, \dots, \tau_n$ 
14: end procedure
    
```

---

where  $\mathbf{w}$  is a weight vector,  $\epsilon$  is a zero-mean i.i.d. Gaussian noise, i.e.  $\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_{TxT})$ , and

$$\Psi = \begin{bmatrix} \psi_1(1) & \psi_2(1) & \cdots & \psi_N(1) \\ \psi_1(2) & \psi_2(2) & \cdots & \psi_N(2) \\ \vdots & \vdots & \ddots & \vdots \\ \psi_1(T) & \psi_2(T) & \cdots & \psi_N(T) \end{bmatrix}. \quad (15)$$

247 A term  $\psi_i(t)$  in this matrix represents the normalized Gaussian basis function with index  $i$  evaluated at  
 248 time step  $t$ .

249 Given a trajectory  $\tau$ , a pre-defined matrix of basis functions  $\Psi$  and a regularizing factor  $\lambda$ , the weight  
 250 vector  $\mathbf{w}$  can be computed with linear ridge regression as follows:

$$\mathbf{w} = \left( \Psi^T \Psi + \lambda \mathbf{I}_{N \times N} \right)^{-1} \Psi^T \tau. \quad (16)$$

251 Once the weight vectors  $w$  corresponding to a set of trajectories  $\tau$  have been computed, a Gaussian  
 252 distribution  $\mathcal{N}(\mu_w, \Sigma_w)$  over these vectors is determined using maximum likelihood estimation. The  
 253 distribution over trajectories  $\tau$  can be expressed as the marginal distribution

$$p(\tau) = \int p(\tau|w) p(w) dw, \quad (17)$$

254 where  $p(w) = \mathcal{N}(w|\mu_w, \Sigma_w)$ . Assuming that a Gaussian is a good approximation for the distribution  
 255 over  $w$ , this integral can be solved in closed-form, yielding

$$p(\tau) = \mathcal{N}(\tau|\mu_\tau, \Sigma_\tau), \quad (18)$$

with

$$\begin{aligned} \mu_\tau &= \Psi \mu_w, \\ \Sigma_\tau &= \sigma^2 I_{TxT} + \Psi \Sigma_w \Psi^T. \end{aligned} \quad (19)$$

256 To deal with not only one stroke and a single degree of freedom (DoF) but with multiple strokes and  
 257 multiple DoFs, one can think of  $\tau$  as a concatenation of trajectories. The matrix  $\Psi$  becomes, in this case, a  
 258 block diagonal matrix and  $w$  a concatenation of weight vectors. For further details about this formulation,  
 259 the interested reader is referred to our previous work (Maeda et al., 2016) in which ProMPs were used to  
 260 coordinate the movements of a human and a robot in collaborative scenarios.

261 The variance  $\sigma^2$  defining the Gaussian noise  $\epsilon$  determines how sensitive our system is to deviations from  
 262 the distribution over demonstrations because  $\sigma^2$  directly influences the variance along this distribution, as  
 263 expressed by (19). A small  $\sigma^2$  results in assessing positions as incorrect more often, while a high  $\sigma^2$  results  
 264 in a less strict evaluation.

### 265 **3.4 Assessing the Correctness of New Trajectories**

266 The correctness of each position of a new trajectory is assessed by comparing the probability density  
 267 function evaluated at that position with the probability density function evaluated at the corresponding  
 268 position along the mean trajectory, which is considered by our system the best achievable trajectory, since  
 269 it is the one with the highest probability under the Gaussian distribution over demonstrations.

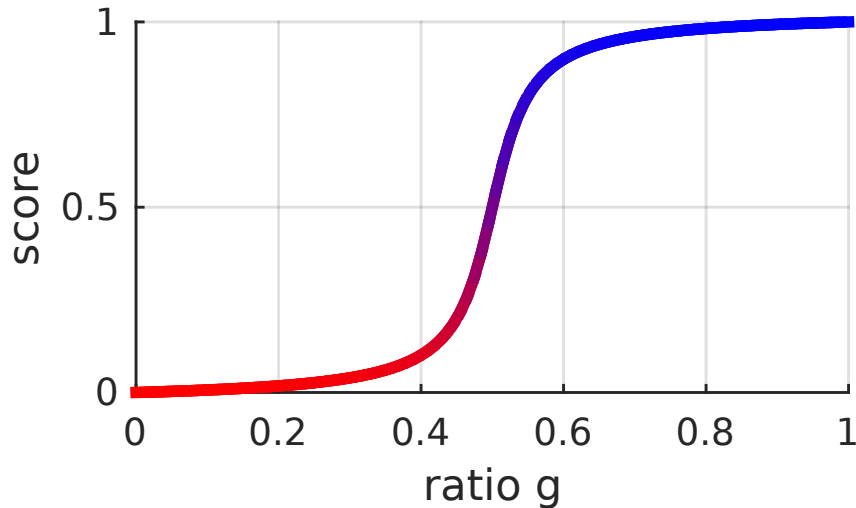
270 First, the ratio

$$g(t) = \frac{p(\tau(t))}{p(\mu_\tau(t))} \quad (20)$$

271 is computed for each time step  $t$ , where  $p(\tau(t))$  is the probability of position  $\tau(t)$  at time step  $t$  and  
 272  $p(\mu_\tau(t))$  is the probability of position  $\mu_\tau(t)$  at time step  $t$ . Since the highest achievable value of the  
 273 Gaussian probability density function at each time step is the one achieved by the mean trajectory,  $g$  is a  
 274 function with values between 0 and 1.

275 Subsequently a score

$$s(g(t)) = \frac{\arctan((g(t) + a) b)}{2c} + 0.5 \quad (21)$$



**Figure 4.** Score function relating the ratio  $g$  between the probability of a certain position and the probability of the corresponding position along the mean trajectory. This function is determined by (21) and was designed to be 0 when  $g = 0$ , 1 when  $g = 1$  and to monotonically increase with  $g$ . It is possible to change the steepness of this function by changing its hyperparameter  $b$ . The same color code as in this figure is used to give visual feedback to the user.

276 for each time step  $t$  is computed, where

$$c = \arctan((1 + a)b). \quad (22)$$

277 The score function  $s$  was designed with a few desired properties in mind. With  $a = -0.5$ ,  $s$  is equal  
 278 to 0 when the ratio  $g$  is equal to 0, it is 0.5 when  $g$  is 0.5 and it is 1 when  $g$  is 1. The score function  $s$   
 279 monotonically increases with  $g$ . Its steepness is regulated by the parameter  $b$ . We have been using  $a = -0.5$   
 280 and  $b = 25$ . One could consider using other score functions, depending on the preferences of the users.  
 281 The score function depicted in Figure 4 leads to a sharp distinction between right and wrong positions.  
 282 One might prefer a more gradual distinction. In this work, we did not investigate what score function the  
 283 users prefer nor whether certain score functions make the users learn faster. These considerations could be  
 284 subject of extensive user studies.

#### 4 METHOD TO PROVIDE HAPTIC FEEDBACK

285 Up to now, it has been solely discussed in this paper how to provide offline visual feedback to the user  
 286 assessing the correctness of his/her movements. Here, it is presented how our framework provides online  
 287 haptic feedback to the user, guiding him/her towards correct movements.

288 The Haption Virtuose 6D can provide force feedback to the user by simulating a virtual object attached to  
 289 its end effector constituting a mass-spring-damper system. Given the mass and the inertia of the virtual  
 290 object, the Virtuose API computes stiffness and damping coefficients that guarantee the stability of the  
 291 system. The intensity of the force produced by this system can be rescaled by a factor denoted in this paper  
 292 by  $\zeta$ .

293 In this work, we are interested in providing feedback to the user according to a probability distribution  
 294 over trajectories, which is computed as in Section 3.3. If the standard deviation at a certain part of the

295 distribution is high, the haptic device should become compliant in that region, while if the standard  
 296 deviation is low, the haptic device should become stiff. The virtual object always lies along the mean  
 297 trajectory of the distribution. The factor  $\zeta$  can be derived from

$$\frac{\zeta - \zeta_{\min}}{\zeta_{\max} - \zeta_{\min}} = \frac{\sigma - \sigma_{\max}}{\sigma_{\min} - \sigma_{\max}}, \quad (23)$$

298 where  $\zeta_{\min}$  and  $\zeta_{\max}$  are respectively the minimum and the maximum force scaling factor. These values  
 299 have been empirically defined in our experiments. The variable  $\sigma$  stands for the standard deviation that  
 300 corresponds to the current position of the virtual object. The variables  $\sigma_{\min}$  and  $\sigma_{\max}$  stand for the minimum  
 301 and maximum standard deviations of the distribution over trajectories. These values can be determined  
 302 from a set of demonstrated trajectories. The underlying assumption behind (23) is that the stiffness is the  
 303 highest when the standard deviation is the minimum and the lowest when the standard deviation is the  
 304 maximum. Moreover, we assume a linear dependence between  $\zeta - \zeta_{\min}$  and  $\sigma - \sigma_{\max}$ . Rearranging (23), we  
 305 get

$$\zeta = \zeta_{\min} + (\zeta_{\max} - \zeta_{\min}) \left( \frac{\sigma - \sigma_{\max}}{\sigma_{\min} - \sigma_{\max}} \right). \quad (24)$$

306 The closest point along the mean trajectory that is not further away from the previous position of the  
 307 virtual object than a certain threshold becomes the new position of the virtual object. This threshold is  
 308 especially necessary when dealing with convoluted trajectories to avoid large sudden variations in the  
 309 position of the virtual object.

310 In situations where there are no good demonstrations available, but there is a performance measurement  
 311 of the trajectories, it is possible to use reinforcement learning to improve the distribution over trajectories.  
 312 Such a situation could be found in a teleoperation scenario, where an optimization problem with multiple  
 313 objectives may have to be solved, accounting for distances to via points, distances to obstacles and other  
 314 performance measurements. In the next section, a novel reinforcement learning algorithm is presented to  
 315 address such problems.

## 5 RELEVANCE WEIGHTED POLICY OPTIMIZATION

316 We are interested in enabling a haptic device to assist a human in a task also when good demonstrations are  
 317 not available. As it will be presented in Section 6.2, our particular task is to move an object in a virtual  
 318 environment from a start position to an end position through a window in a wall. We have defined three  
 319 objectives to determine the performance of solutions to this task: distance to the start position, distance to  
 320 the center of the window and distance to the end position. An optimal policy for this task is a trajectory  
 321 that begins at the start position, passes through the center of the window and reaches the end position.  
 322 This problem can be decomposed into three subproblems w.r.t. which a policy parameter can be more  
 323 or less relevant. Therefore, in this section, a new policy search method is explained, which identifies the  
 324 relevance of each policy parameter to each subproblem in order to improve the learning of the global task.  
 325 This method makes use of Reward-weighted Regression (Peters and Schaal, 2007). The basic idea of this  
 326 method is to first find out how much each policy parameter influences each objective. Subsequently, this  
 327 information is used to optimize the policy with respect to the objectives. In our particular application, the  
 328 policy parameters are the elements of the weight vector  $w$  as in (14).

329 **5.1 Learning Relevance Functions**

330 Our approach to answering how much each policy parameter influences each objective consists of  
 331 learning a relevance function  $f_o$  for each objective  $o$ . The argument of this function is an index identifying  
 332 a policy parameter. In other words, a relevance function  $f_o(n)$  evaluated for policy parameter indexed by  $n$   
 333 represents how relevant this parameter is to the objective indexed by  $o$ . In order to learn this function, in  
 334 this paper, it is assumed that a relevance function can be represented by a weighted sum of basis functions  
 335 with lower bound 0 and upper bound 1 as follows:

$$f_o(n) = \begin{cases} 0, & \text{if } \boldsymbol{\rho}^T \boldsymbol{\phi}(n) \leq 0 \\ 1, & \text{if } \boldsymbol{\rho}^T \boldsymbol{\phi}(n) \geq 1, \\ \boldsymbol{\rho}^T \boldsymbol{\phi}(n), & \text{otherwise,} \end{cases} \quad (25)$$

336 where  $\boldsymbol{\rho}$  is a vector of weights  $\rho_i$  for the basis functions  $\phi_i$  and  $\boldsymbol{\phi}(n) = [\phi_1(n), \phi_2(n), \dots, \phi_I(n)]^T$ . It  
 337 will become clear in the remainder of this section why the lower bound of a relevance function is 0 and its  
 338 upper bound is 1.

The basis functions are

$$\phi_i(n) = \frac{1}{\exp(-k(n - m_i))}, \quad (26)$$

$$\phi_I = 1, \quad (27)$$

339 with  $i \in \{1, 2, \dots, I\}$ , where  $I$  is the total number of basis functions for the relevance function,  $n$  is an  
 340 index representing one of the policy parameters,  $k$  is a scalar determining steepness and  $m_i$  is a scalar  
 341 determining the midpoint of the logistic basis function with index  $i$ .

These basis functions have been chosen because weighted combinations of them lead to reasonable  
 relevance functions. For example, three relevance functions that can be constructed with the proposed basis  
 functions are depicted in Figure 5. The depicted relevance functions determine how each of the parameters  
 determining a movement influences objectives at the beginning of the movement, in the middle or in the  
 end. These relevance functions are

$$f_{\text{start}}(n) = \phi_3(n) - \phi_2(n), \quad (28)$$

$$f_{\text{middle}}(n) = \frac{1}{\max_n(\phi_1(n) - \phi_2(n))} \phi_1(n) - \frac{1}{\max_n(\phi_1(n) - \phi_2(n))} \phi_2(n), \quad (29)$$

$$f_{\text{end}}(n) = \phi_1(n), \quad (30)$$

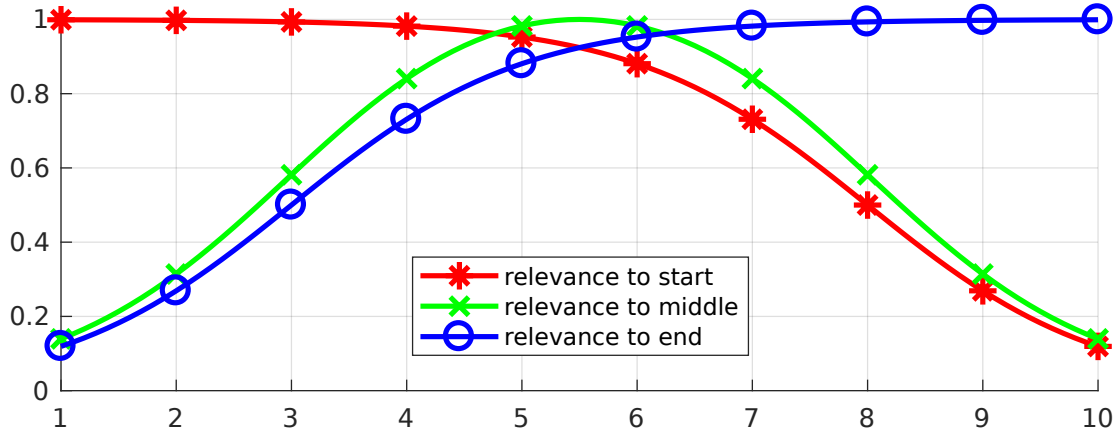
where the basis functions are

$$\phi_1(n) = \frac{1}{\exp(-(n - 3))}, \quad (31)$$

$$\phi_2(n) = \frac{1}{\exp(-(n - 8))}, \quad (32)$$

$$\phi_3(n) = 1. \quad (33)$$

342



**Figure 5.** Three examples of relevance functions. Let us assume that our goal is to optimize a certain movement with respect to an objective at the beginning of the movement, an objective in the middle and an objective in the end. Let us further assume that the movement to be optimized can be determined by 10 parameters and that the first parameters (close to 1) have higher influence over the beginning of the movement, while the last ones (close to 10) have higher influence over the end. The image depicts potentially suitable relevance functions for each of the objectives in this problem.

343 In this framework, learning a relevance function with respect to a certain objective means finding a vector  
 344  $\rho$  that leads to a high variability in the value of that objective and to a low variability in the values of other  
 345 objectives. How a relevance function influences the variability in the values of an objective will be made  
 346 explicit in the following.

347 First, a Gaussian distribution  $\mathcal{N}(\mu_\rho, \Sigma_\rho)$  over  $\rho$  is initialized with a certain mean  $\mu_\rho$  and a certain  
 348 covariance matrix  $\Sigma_\rho$ . Subsequently, parameter vectors  $\rho$  are sampled from this distribution and, for each  
 349 sample, the relevance function  $f_o$  is computed using (25).

350 Let us now assume that there is an initial Gaussian probability distribution  $\mathcal{N}(\mu_w, \Sigma_w)$  over the policy  
 351 parameters  $w$ . The mean  $\mu_w$  and the covariance matrix  $\Sigma_w$  can be computed from an initial set of  
 352 demonstrations or determined by the user.

For each  $f_o$  computed with the sampled vectors  $\rho$ , our algorithm generates samples of the policy parameters  $w$  from the distribution  $\mathcal{N}(\mu_w, \Sigma_w^{f_o})$ , where

$$\Sigma_w^{f_o} = \begin{bmatrix} \sigma_{w_1}^2 f_o(1) & 0 & \dots & 0 \\ 0 & \sigma_{w_2}^2 f_o(2) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_{w_N}^2 f_o(N) \end{bmatrix} \tag{34}$$

353 and  $\sigma_{w_n}^2, \forall n \in \{1, 2, \dots, N\}$ , are the variances in the diagonal of the matrix  $\Sigma_w$ . In other words, the  
 354 policy parameters are sampled in such a way that their original variance is weighted with a relevance  
 355 coefficient. The higher the relevance of a parameter, the larger the range of values for that parameter among  
 356 the samples.

357 Each sampled vector of policy parameters  $w$  determines a policy with a corresponding value for each  
 358 objective. In our teleoperation scenario, for example, each policy parameter vector  $w$  determines a

359 trajectory, which has a certain distance to the start position, a certain distance to the center of the window  
 360 and a certain distance to the end position. Given these objective values, our algorithm computes a reward  
 361 function

$$R_{\rho,o} = \exp \left( \beta_{\text{relevance}} \left( \sigma_o - \sum_{i \neq o} \sigma_i \right) \right), \quad (35)$$

362 where  $\sigma_o$  is the standard deviation of the values for objective  $o$  and  $\sigma_i$  with  $i \neq o$  is the standard deviation  
 363 of the values for the other objectives. The scalar  $\beta_{\text{relevance}}$  can be determined with line search.

364 Parameters  $\rho$  determining suitable relevance functions  $f_o$  result in higher reward  $R_{\rho,o}$  because the range  
 365 of values for the parameters that mainly affect objective  $o$  will be high, producing a high standard deviation  
 366 of the values for that objective. Moreover, the range of values for the parameters that mainly affect the  
 367 other objectives will be low, producing a low standard deviation of the values for the other objectives.

368 Finally, Reward-weighted Regression (RWR) is used to learn the relevance parameters  $\rho$ . RWR is an  
 369 iterative algorithm that finds the best Gaussian distribution over parameters of interest (in the particular  
 370 case of optimizing the relevance functions, the parameters of interest are given by  $\rho$ ) to maximize the  
 371 expected reward, given samples from the Gaussian distribution of the previous iteration. In order to do so,  
 372 RWR solves the optimization problem

$$\{\mu_{\rho}^{k+1}, \Sigma_{\rho}^{k+1}\} = \arg \max_{\{\mu_{\rho}, \Sigma_{\rho}\}} \sum_{i=1}^S R_{\rho,o,i} \mathcal{N}(\rho_i; \mu_{\rho}, \Sigma_{\rho}) \quad (36)$$

373 at each iteration, where  $S$  is the number of sampled parameter vectors  $\rho_i$  from the previous distribution  
 374  $\mathcal{N}(\mu_{\rho}^k, \Sigma_{\rho}^k)$ . The solution to this optimization problem is

$$\mu_{\rho}^{k+1} = \frac{\sum_{i=1}^S R_{\rho,o,i} \rho_i}{\sum_{i=1}^S R_{\rho,o,i}}, \quad (37)$$

375

$$\Sigma_{\rho}^{k+1} = \frac{\sum_{i=1}^S R_{\rho,o,i} (\rho_i - \mu_{\rho}^k) (\rho_i - \mu_{\rho}^k)^T}{\sum_{i=1}^S R_{\rho,o,i}}. \quad (38)$$

376 This procedure is repeated until convergence of  $R_{\rho,o}$  has been reached to learn a relevance function  
 377  $f_o$  for each objective  $o$ . The parameters determining the relevance functions  $f_o$  are given by the vector  
 378  $\mu_{\rho}$  computed in the last iteration. After this iterative procedure is finished, our algorithm computes  
 379  $f_o(n) / \max_n f_o(n)$ ,  $\forall n \in \{1, 2, \dots, N\}$ , and assigns this value to  $f_o(n)$ . This last step makes the  
 380 maximum value of  $f_o$  be not less than 1 and helps the exploration in the policy optimization phase, which  
 381 will be discussed in the next section. Algorithm 4 presents an informal description of the relevance learning  
 382 algorithm.

## 383 5.2 Policy Optimization using Relevance Functions

384 Now that a relevance function for each objective has been learned, our algorithm uses this information to  
 385 optimize a policy with respect to each objective. As in Section 5.1, it is assumed here that there is an initial  
 386 Gaussian probability distribution  $\mathcal{N}(\mu_w, \Sigma_w)$  over the policy parameters  $w$ .



**Algorithm 4** Learning Relevance Functions

```

1: Inputs: mean  $\mu_w$  and covariance  $\Sigma_w$  of the policy parameter vectors  $w$ 
2: Initialize the mean  $\mu_\rho$  and the covariance  $\Sigma_\rho$  of the Gaussian distribution over the parameter vectors
    $\rho$  that determine the relevance functions  $f_o$ 
3: repeat
4:   Sample parameter vectors  $\rho$  from  $\mathcal{N}(\mu_\rho, \Sigma_\rho)$ 
5:   for each sample vector  $\rho$  do
6:     for each objective  $o$  do
7:       Compute the relevance functions  $f_o$  (Equation 25)
8:       Compute matrix  $\Sigma_w^{f_o}$  (Equation 34)
9:       Sample policy parameter vectors  $w$  from  $\mathcal{N}(\mu_w, \Sigma_w^{f_o})$ 
10:      for each sample vector  $w$  do
11:        Compute value achieved by policy for objective  $o$ 
12:      end for
13:      Compute standard deviation  $\sigma_o$  of the values achieved for  $o$  with the different samples
14:    end for
15:    for each objective  $o$  do
16:      Compute  $R_{\rho,o}$  (Equation 35)
17:    end for
18:  end for
19:  Update  $\mu_\rho$  and  $\Sigma_\rho$  (Equations 37 and 38)
20: until convergence of the rewards  $R_{\rho,o}$ 
21:  $\rho^* = \mu_\rho$ 
22: for each objective  $o$  do
23:   Compute  $f_o^*$  using  $\rho^*$  (Equation 25)
24:   Normalize  $f_o^*$  by computing  $\frac{f_o^*(n)}{\max_n f_o^*(n)}$ 
25: end for
26: return the relevance functions  $f_o^*$ 

```

For each objective  $o$ , our algorithm samples policy parameters  $w$  from the distribution  $\mathcal{N}(\mu_w, \Sigma_w^{f_o^*})$ , where

$$\Sigma_w^{f_o^*} = \begin{bmatrix} \sigma_{w_1}^2 f_o^*(1) & 0 & \dots & 0 \\ 0 & \sigma_{w_2}^2 f_o^*(2) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_{w_N}^2 f_o^*(N) \end{bmatrix} \quad (39)$$

387 and  $f_o^*$  is the learned relevance function with respect to the objective  $o$ . Therefore, the policy parameters  $w$   
388 are sampled from a Gaussian distribution where the original variances  $\sigma_{w_n}^2$  are weighted with the learned  
389 relevance function. This procedure means that a larger range of values will be sampled for the policy  
390 parameters  $w_n$  that are more relevant to the objective  $o$  and a smaller range of values will be sampled for  
391 the policy parameters  $w_n$  that are less relevant to this objective.

392 For each sampled vector of policy parameters  $w_i$ , the reward  $R_{w,o,i}$  associated with the objective  $o$  is  
393 computed. These objectives and rewards depend on the problem. An objective might be for instance to  
394 achieve a certain goal position, in which case the reward could be a non-negative function monotonically  
395 decreasing with the distance to the goal position. In our particular teleoperation problem, the reward  
396 associated with the objective of being close to the start position is given by  $R = \exp(-\beta_{\text{policy}} d_{\text{start}})$ , where

397  $d_{\text{start}}$  is the distance between the first position of the trajectory and the position where the trajectories should  
 398 start.

399 Our algorithm uses once again RWR. This time, RWR is used to to maximize the expected reward with  
 400 respect to  $\mu_w$  and  $\Sigma_w^{f_o^*}$ . This maximization is done iteratively according to

$$\{\mu_w^{k+1}, C^{k+1}\} = \arg \max_{\{\mu_w, \Sigma_w^{f_o^*}\}} \sum_{i=1}^S R_{w,o,i} \mathcal{N}(w_i; \mu_w, \Sigma_w^{f_o^*}), \quad (40)$$

401 where  $S$  is the number of sampled policy parameter vectors  $w_i$  from the previous distribution  
 402  $\mathcal{N}(w; \mu_w^k, \Sigma_w^{f_o^*k})$ .

403 The solution to (40) is given by

$$\mu_w^{k+1} = \frac{\sum_{i=1}^S R_{w,o,i} w_i}{\sum_{i=1}^S R_{w,o,i}}, \quad (41)$$

404

$$C^{k+1} = \frac{\sum_{i=1}^S R_{w,o,i} (w_i - \mu_w^k) (w_i - \mu_w^k)^T}{\sum_{i=1}^S R_{w,o,i}}. \quad (42)$$

405 In each iteration, after computations (41) and (42), our algorithm updates the variances of each policy  
 406 parameter  $\sigma_{w_n}^2$  with

$$\sigma_{w_n,k+1}^2 = (1 - f_o(n)) \sigma_{w_n,k}^2 + f_o(n) C_{nn}^{k+1}, \quad (43)$$

407 where  $\sigma_{w_n,k}^2$  is the previous variance of policy parameter  $w_n$  and  $C_{nn}^{k+1}$  is the  $n^{\text{th}}$  element along the main  
 408 diagonal of the matrix  $C^{k+1}$ . This equation has the effect of keeping the previous variance of the parameters  
 409 that are less relevant to the objective  $o$  while updating the variance of the parameters that are more relevant  
 410 to this objective. The algorithm then uses  $\sigma_{w_n,k+1}^2$  to compute  $\Sigma_w^{f_o^*k+1}$  as in (39).

411 Finally, Equation 43 justifies the lower bound of 0 and the upper bound of 1 for the relevance function.  
 412 The closer the relevance of policy parameter  $w_n$  is to 0, the closer the updated variance of this parameter  
 413 is to the previous variance  $\sigma_{w_n,k}^2$ . The closer the relevance of policy parameter  $w_n$  is to 1, the closer the  
 414 updated variance of this parameter is to  $C_{nn}^{k+1}$ . In other words, the previous variance of irrelevant policy  
 415 parameters is preserved, while the variance of relevant policy parameters is updated. Algorithm 5 presents  
 416 an informal description of the algorithm for policy optimization using relevance functions.

### 417 5.3 Example of Policy Optimization with Relevance Weighting

418 In order to make the proposed Relevance Weighted Policy Optimization algorithm more clear, we present  
 419 an example using the 2D scenario depicted in Figure 6(a). This scenario is composed of a start position,  
 420 a wall with a window and an end position. Given the initial trajectories depicted in Figure 6(a), the goal  
 421 of our algorithm is to find a distribution over trajectories that begin at the start position, pass through the  
 422 center of the window and reach the end position.

423 First, the algorithm aligns the initial trajectories in time and computes the parameters  $w$  for each of them  
 424 using (16). Subsequently, the relevance functions for start position, center and end position are learned  
 425 as in Section 5.1. An example of learned relevance functions is depicted in Figure 6(b). After learning  
 426 the relevance functions, the algorithm uses the procedure explained in Section 5.2 to learn a policy that  
 427 satisfies the three above-stated objectives. Figure 7 shows how the distribution over trajectories changes

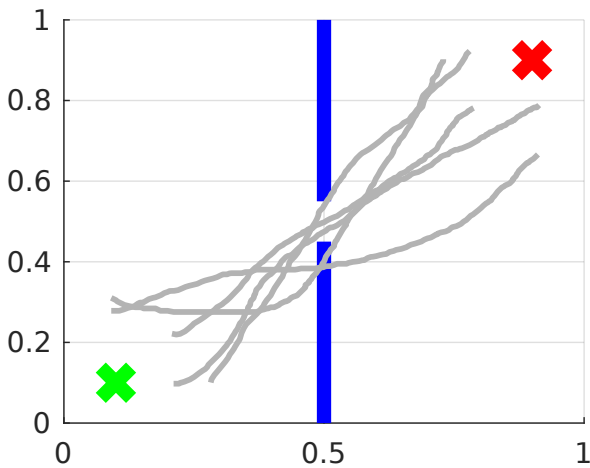
**Algorithm 5** Policy Optimization using Relevance Functions

```

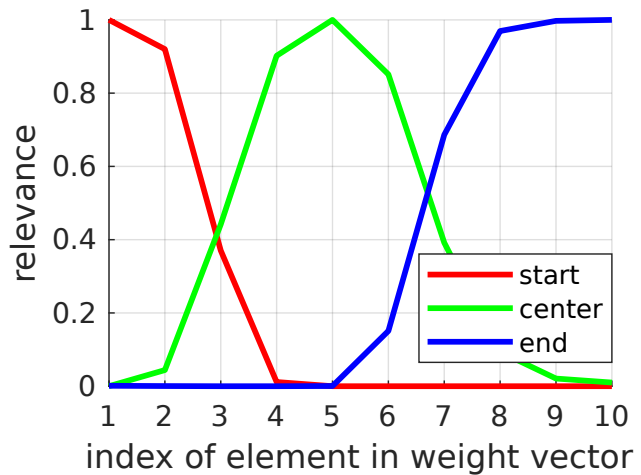
1: Inputs: mean  $\mu_w$  and covariance  $\Sigma_w$  of the policy parameter vectors  $w$ , learned relevance functions  $f_o^*$ 
2: repeat
3:   for each objective  $o$  do
4:     Compute matrix  $\Sigma_w^{f_o^*}$  (Equation 39)
5:     Sample policy parameter vectors  $w$  from  $\mathcal{N}(\mu_w, \Sigma_w^{f_o^*})$ 
6:     for each sample vector  $w$  do
7:       Compute the reward  $R_{w,o}$  of the policy with parameters given by vector  $w$  associated with objective  $o$ 
8:     end for
9:     Update  $\mu_w$  and compute  $C$  (Equations 41 and 42)
10:    Update the variances of the policy parameters  $\sigma_{w_n}^2$  (Equation 43)
11:  end for
12: until convergence of the rewards  $R_{w,o}$ 
13: return the mean  $\mu_w$  and the variances  $\sigma_{w_n}^2$ 

```

**(a) Demonstrated trajectories**



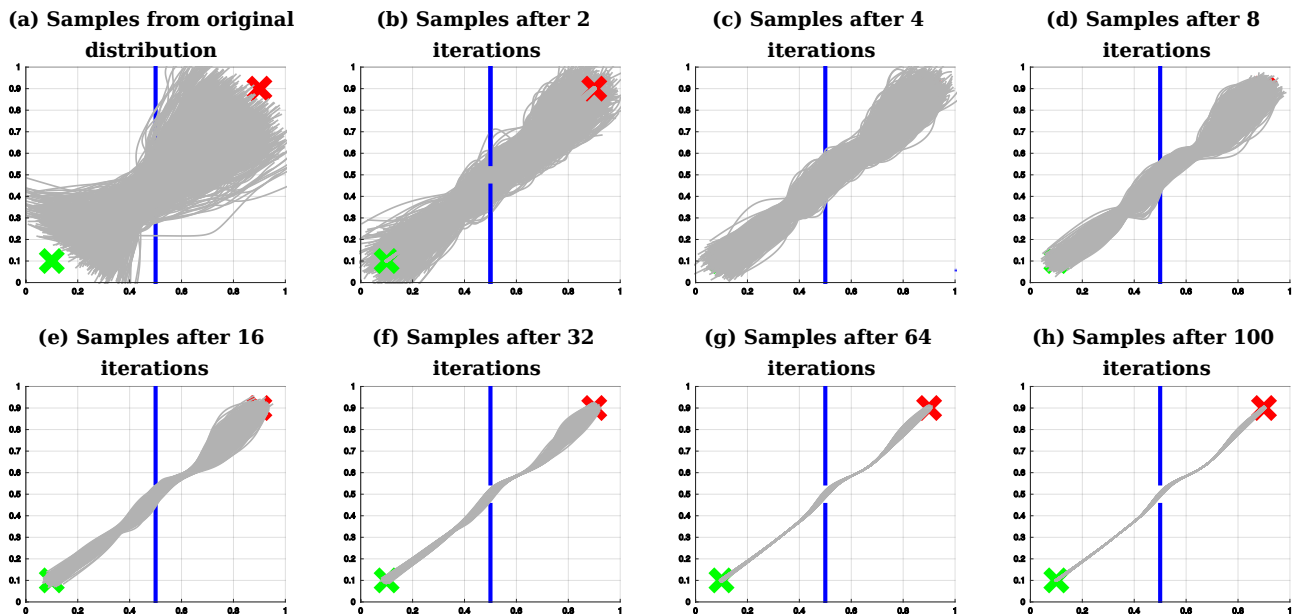
**(b) Learned relevance functions**



**Figure 6.** (a) 2D problem used to explain the proposed Relevance Weighted Policy Optimization (RWPO) algorithm. The green x at the lower left corner of the image represents the start position. The blue lines in the middle represent a wall with a window in the center. The red x at the upper-right corner represents the end position. The goal of our algorithm is, given a few initial trajectories (depicted in light gray), to find a distribution over trajectories that begin at the start position, pass through the center of the window and reach the end position. (b) Learned relevance functions for the 2D problem. The learned relevance functions show that policy parameters close to  $w_1$  are more important for beginning at the start position, policy parameters around  $w_5$  are more important to pass through the center of the window and policy parameters close to  $w_{10}$  are more important to reach the end position.

428 with the number of iterations of the algorithm. The distances to start, center and end positions decrease  
429 with the number of iterations and the return  $\exp(-\beta_{\text{policy}}(d_{\text{start}} + d_{\text{center}} + d_{\text{end}}))$  increases, as depicted  
430 by Figure 8. Here,  $\beta_{\text{policy}}$  is a parameter which can be determined with line search,  $d_{\text{start}}$  is the distance to  
431 the start,  $d_{\text{center}}$  is the distance to the center and  $d_{\text{end}}$  is the distance to the end.

432 Relevance Weighted Policy Optimization implements policy search for each objective sequentially.  
433 For each objective, the algorithm samples a larger range of values for the parameters that are more



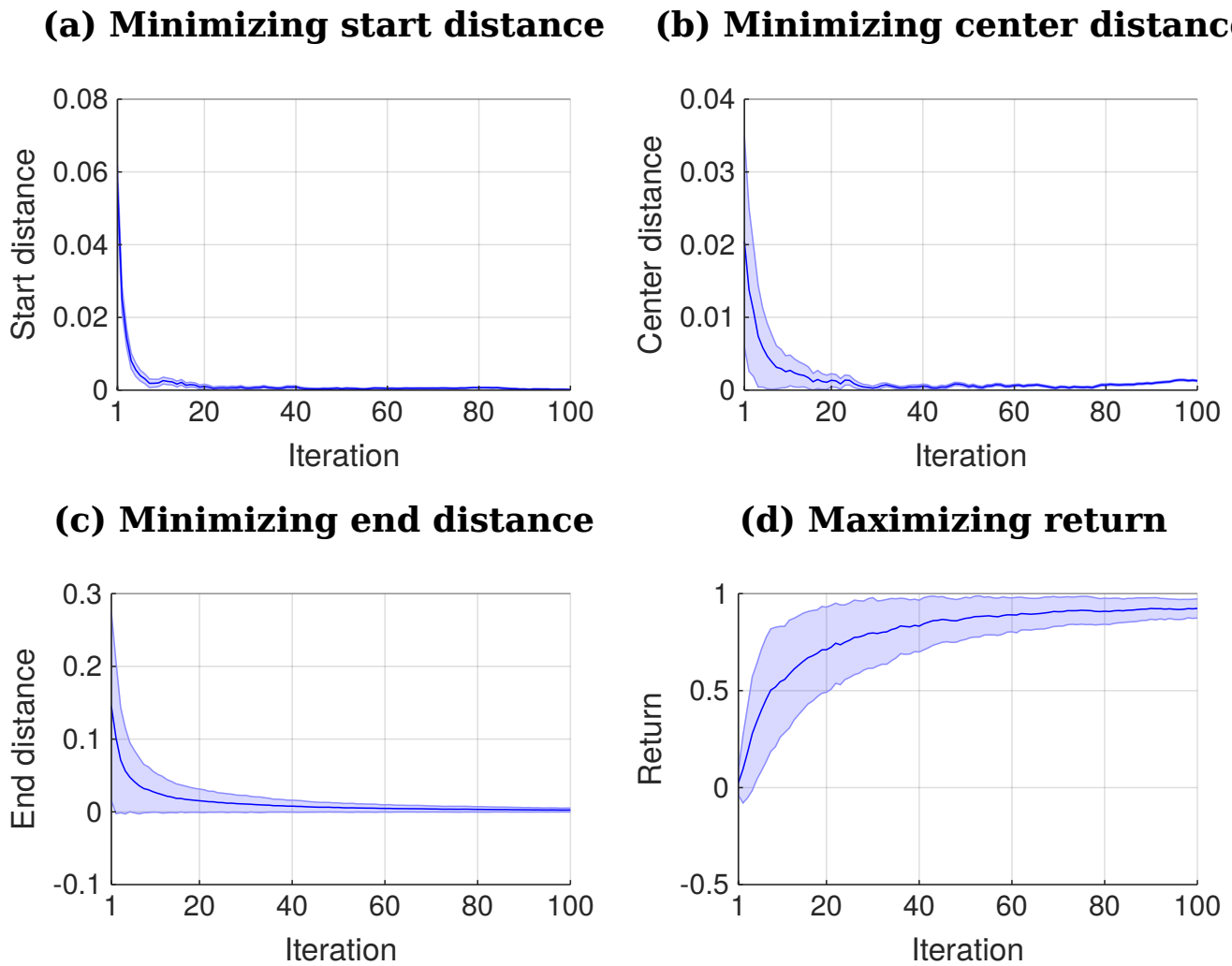
**Figure 7.** Example of policy search with relevance weighting. The proposed algorithm finds a distribution over trajectories that start and end at the correct positions (represented by the green x and by the red x, respectively) and do not hit the wall (represented by the blue lines).

434 relevant to that objective, while sampling values close to the mean for the parameters that are less relevant.  
 435 Subsequently, the algorithm optimizes the mean and the variances of the policy parameters given the  
 436 samples. After optimization, the mean and the variance of the parameters that matter more to that objective  
 437 are updated, while the mean and the variance of parameters that matter less remain similar to the previous  
 438 distribution. The algorithm does not require defining a reward function with different weights for the  
 439 different objectives, which can be time-consuming and ineffective. Moreover, at each iteration, when  
 440 optimizing the distribution over policy parameters with respect to a certain objective, the algorithm does  
 441 not accidentally find solutions that are good according to this objective, but bad according to the other  
 442 objectives because only the mean and the variance of the parameters that matter change substantially. The  
 443 mean and the variance of the other parameters remain close to the mean and the variance of the previous  
 444 distribution.

445 Figure 9 exemplifies how the algorithm samples trajectories in the 2D teleoperation problem. Figure  
 446 9(a) shows samples from the original distribution. Figure 9(b) shows samples of the first iteration of the  
 447 algorithm right before optimizing for beginning at the start position. Figure 9(c) depicts the next step, still  
 448 in the first iteration, after the first optimization for starting at the start position and before optimizing for  
 449 passing through the center of the window. Finally, Figure 9(d) shows samples at the first iteration of the  
 450 algorithm, right before optimizing for reaching the end position.

451 Figure 10(a) shows a distribution over trajectories learned by Reward-weighted Regression (RWR)  
 452 optimizing only for passing through the center of the window. Figure 10(b) shows the solution of Relevance  
 453 Weighted Policy Optimization (RWPO) for this same optimization problem. RWPO's solution achieves  
 454 the objective with higher accuracy and preserves a large variance for parts of the trajectory that do not  
 455 influence the objective.

Finally, Figure 11 shows a comparison between Comparison between Reward-weighted Regression (RWR), sequential Reward-weighted Regression (sRWR) and Relevance Weighted Policy Optimization



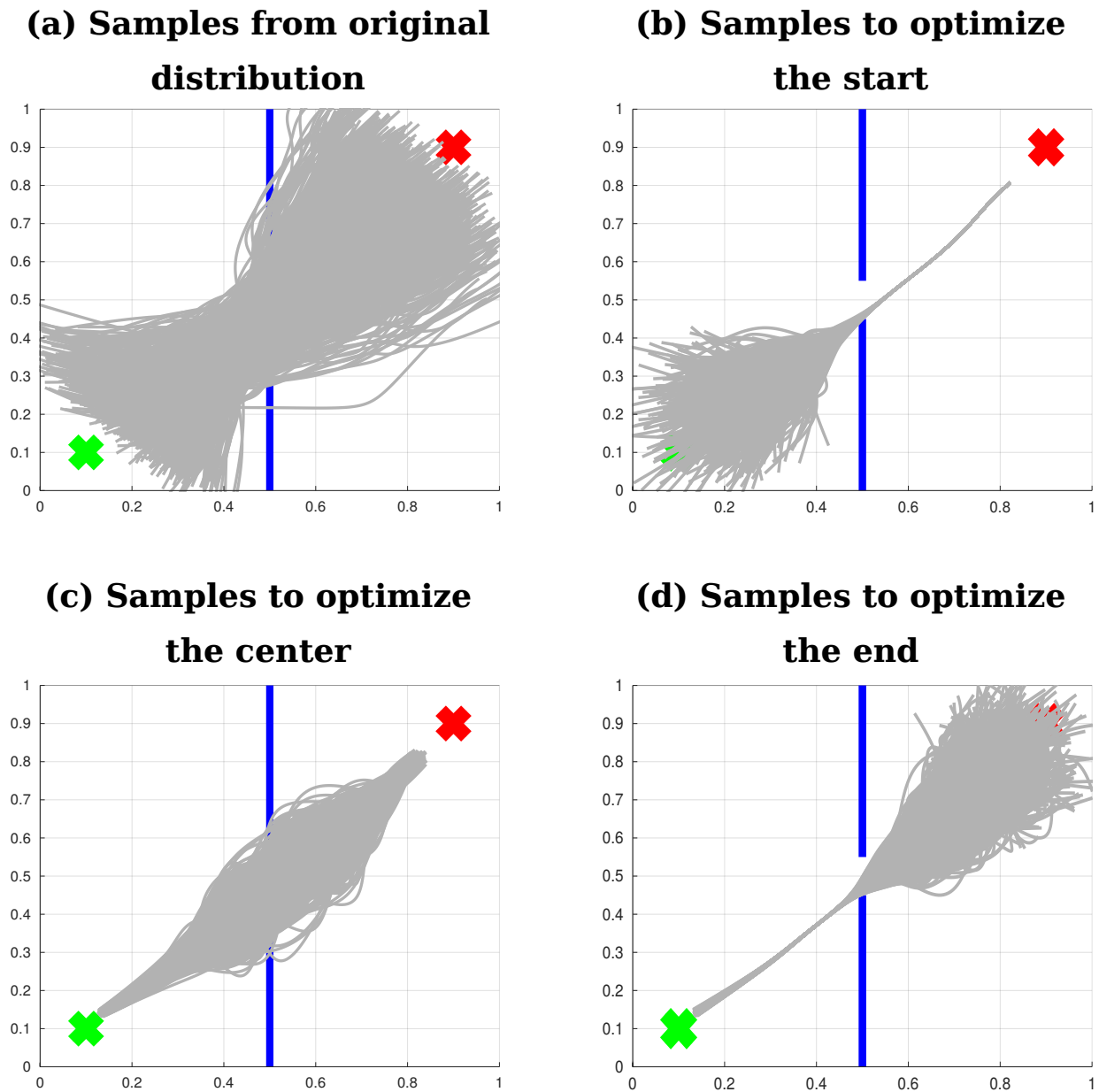
**Figure 8.** Iteration versus distances and iteration versus returns. The plots represent mean and two times the standard deviation. All the distances to the points of interest decrease to 0 or close to it with the number of iterations. A return function given by  $\exp(-\beta_{\text{policy}}(d_{\text{start}} + d_{\text{center}} + d_{\text{end}}))$  increases with the number of iterations. Here,  $\beta_{\text{policy}}$  is a parameter which can be determined with line search,  $d_{\text{start}}$  is the distance to the start,  $d_{\text{center}}$  is the distance to the center and  $d_{\text{end}}$  is the distance to the end.

(RWPO). RWR used here a reward function of the form  $R = \exp(-\beta_{\text{policy}}(d_{\text{start}} + d_{\text{center}} + d_{\text{end}}))$ , while sRWR and RWPO used one reward function for each objective:

$$R_{\text{start}} = \exp(-\beta_{\text{policy}}(d_{\text{start}})), \tag{44}$$

$$R_{\text{center}} = \exp(-\beta_{\text{policy}}(d_{\text{center}})), \tag{45}$$

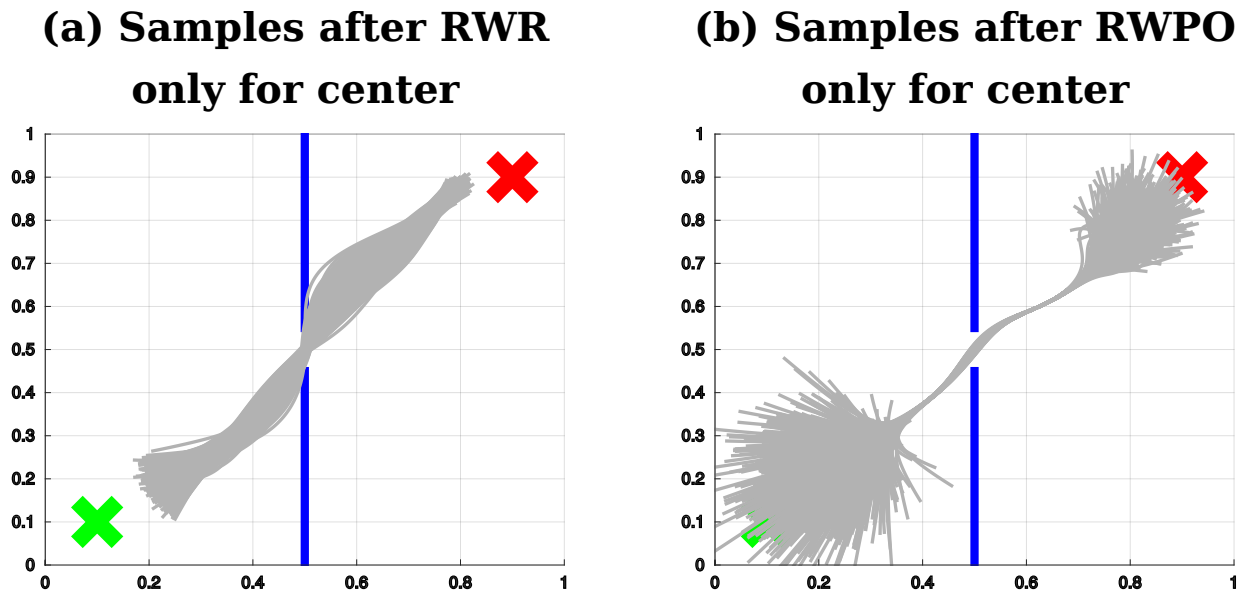
$$R_{\text{end}} = \exp(-\beta_{\text{policy}}(d_{\text{end}})). \tag{46}$$



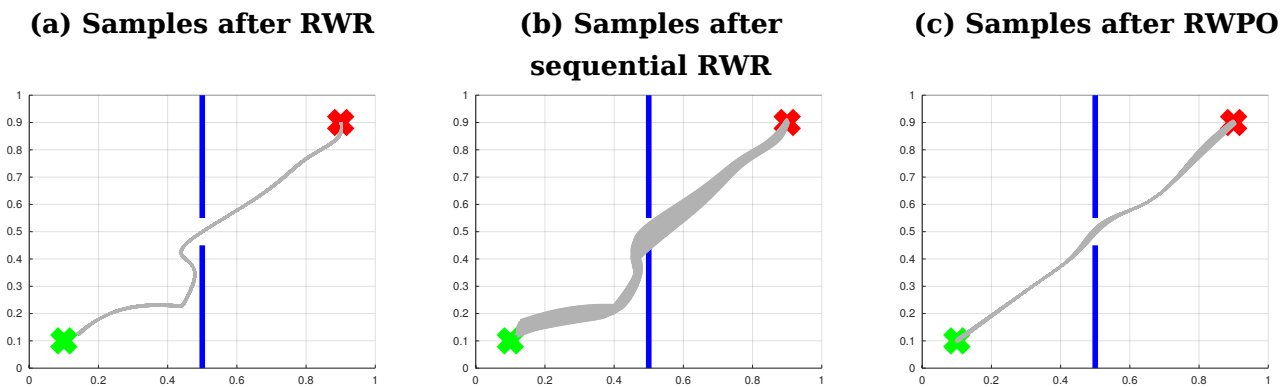
**Figure 9.** Sampling with relevance weighting. **(a)** Samples from the original distribution. **(b)** Samples to optimize the distribution over trajectories with respect to beginning at the start position. **(c)** Samples to optimize the distribution over trajectories with respect to passing through the center of the window. **(d)** Samples to optimize the distribution over trajectories with respect to reaching the end position. The proposed algorithm explores for each objective a large range of values for the policy parameters that are relevant to that objective, while sampling values close to the mean for the other policy parameters. The variance of the irrelevant parameters is recovered according to Equation 43. Therefore, after optimizing for each objective, the distribution over the relevant parameters is updated, while the distribution over the irrelevant parameters is preserved.

## 6 EXPERIMENTS

456 We demonstrate our method to assist the practice of motor skills by humans with the task of writing Japanese  
 457 characters. Moreover, an experiment involving a haptic device, the Haption Virtuose 6D, demonstrates how

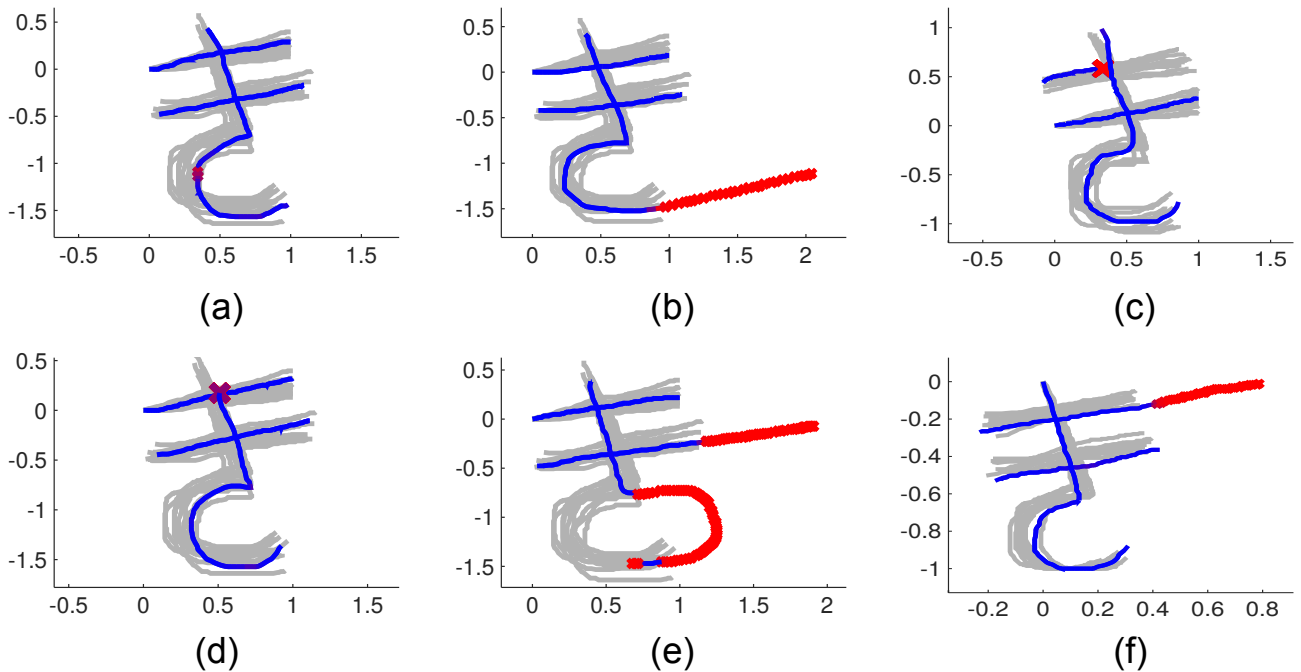


**Figure 10.** (a) Sample trajectories after using Reward-weighted Regression (RWR) to optimize the distribution over trajectories with respect to passing through the center of the window. (b) Sample trajectories after using Relevance Weighted Policy Optimization (RWPO) to optimize the distribution over trajectories with respect to the same objective, using the same reward function. In contrast to RWR, RWPO finds a better policy to avoid hitting the wall and does not shrink the variance of parts of the trajectories that are far away from the region of interest.



**Figure 11.** Comparison between Reward-weighted Regression (RWR), sequential Reward-weighted Regression (sRWR) and Relevance Weighted Policy Optimization (RWPO). This time, the three algorithms optimize the distribution over trajectories with respect to all objectives. (a) Distribution after optimization with RWR, which uses a single reward function with a term for each objective. (b) Distribution after optimization using sRWR, which optimizes for each objective sequentially and has a reward function for each objective. (c) Distribution after optimization using RWPO, which uses the concept of relevance functions and optimizes for each objective sequentially with a reward function for each objective.

458 our method can be used to give haptic feedback to the user, guiding him/her towards correct movements  
 459 according to certain performance criteria even in the absence of expert demonstrations.



**Figure 12.** The demonstrations after rescaling, repositioning and time-alignment are depicted in light gray. Parts of a new trajectory that are considered correct are depicted in blue. Parts of a new trajectory that are considered wrong are marked with red x's . **(a)** Instance with a small mistake in the third stroke. **(b)** Third stroke goes further than it should. **(c)** First stroke is too short. **(d)** Third stroke starts too low. **(e)** Second stroke is too long and third stroke has its arch in the wrong direction. **(f)** First stroke is too long.

## 460 6.1 Teaching Japanese Characters

461 In these experiments, first, a human provided with a computer mouse 10 demonstrations of a certain  
 462 Japanese character composed of multiple strokes. Our system aligned these demonstrations in space and  
 463 time. Afterward, a human provided a new trajectory. This new trajectory was also aligned in space and  
 464 time by our system with respect to the demonstrations. Once all the demonstrations and the new trajectory  
 465 had been time-aligned, our system computed a probability distribution over the demonstrations. Based  
 466 on the probability density function evaluated at each position of the new trajectory in comparison to the  
 467 probability density function evaluated at corresponding positions along the mean trajectory, our system  
 468 computed a score. This score was then used to highlight parts of the new trajectory that do not fit the  
 469 distribution over demonstrations with a high probability.

470 Figure 12 shows some examples of feedbacks provided by our system. The new trajectory provided by  
 471 the user is also aligned in space and time. Therefore the absolute position of his/her character and its scale  
 472 are not relevant. The speed profile of the new trajectory can also be different from the speed profile of the  
 473 demonstrations. Figure 12 shows the new trajectories already after alignment in space and time.

## 474 6.2 Haptic Feedback

475 When learning complex movements in a 3D environment or perhaps when manipulating objects, haptic  
 476 feedback may give the human information about how to adapt his/her movements that would be difficult  
 477 to extract only from visual feedback. Therefore, we investigated how to give haptic feedback based on a  
 478 probability distribution over trajectories possibly provided by an instructor or resulting from a reinforcement  
 479 learning algorithm. This study was carried out in accordance with the recommendations of the Declaration



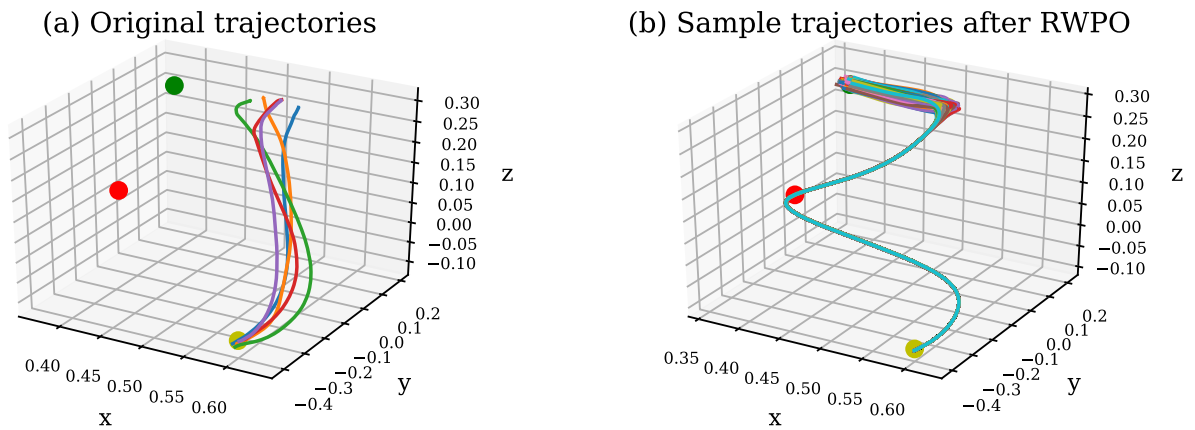
480 of Helsinki in its latest version. The protocol was approved by the Ethical Committee of the Technische  
481 Universität Darmstadt. All participants provided written informed consent before participation.

482 In this user experiment, users had to use the Haption Virtuose 6D device to teleoperate a small cube in  
483 a 3D environment (See Figure 14(a)). The users were instructed to begin at the position marked by the  
484 yellow sphere, pass through the center of the window in the wall and end at the position marked by the blue  
485 sphere. Moreover, it was allowed, at any time, to rotate the virtual environment, zoom in and zoom out  
486 using the computer mouse. Five users took part in our experiments: two females and three males, between  
487 27 and 29 years old. Three users had not used the Virtuose 6D before, while two users did have some  
488 experience with it.

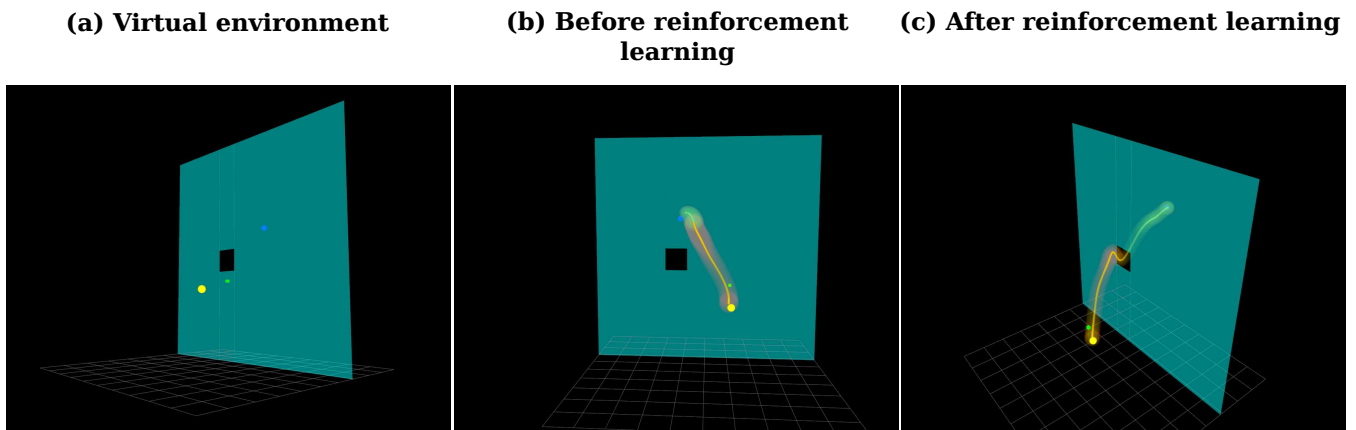
489 In the first phase of the experiment, users tried to perform the task ten times without force feedback.  
490 Right before each trial, the user pressed a button on the handle of the haptic device, indicating to our  
491 system when to start recording the trajectory. By pressing this same button another time, by the end of  
492 the trajectory, the user indicated to our system when to finish recording. The users were then instructed  
493 to move the cube back to the start position and perform another trial. This procedure would be repeated  
494 until ten trajectories had been recorded. Afterward, our system would align them in time and compute a  
495 probability distribution over them. Figure 14(b) shows a visualization of the distribution over trajectories  
496 of one user after this phase. Subsequently, our system optimized this distribution over trajectories using the  
497 proposed Relevance Weighted Policy Optimization (RWPO) algorithm. An example of trajectories before  
498 and after RWPO is depicted in Figure 13. Figure 14(c) shows the optimized distribution over trajectories  
499 given the initial distribution shown in Figure 14(b). After optimizing the distribution over trajectories, our  
500 system used it to give force feedback to the user according to the method explained in Section 4. The users  
501 were requested to try to perform the task with force feedback ten times using the aforementioned procedure  
502 to record the trajectories.

503 Results showing the performance of the users with and without force feedback are presented in Figure 15.  
504 The use of force feedback did not greatly influence the distance to the start because the force feedback  
505 was activated only when the user pressed a button, right before starting to move. The start distance of the  
506 third trial of user 2 with force feedback is an outlier. This outlier was due to the user starting far away  
507 from the start position. The use of force feedback decreased the distance to the center of the window for  
508 all users and the distance to the end for three out of five users. The plots of trial versus distances indicate  
509 that the users did not achieve a better performance with the force feedback only due to training through  
510 repetition because there is a clear difference between the performance in trials with force feedback and the  
511 performance in trials without force feedback.

512 A  $2$  (feedback)  $\times$   $3$  (distance measures) repeated-measures ANOVA was conducted to test the  
513 performance differences between the conditions with and without force feedback. The results reveal  
514 significant main effects of feedback ( $F(1, 4) = 16.31$ ;  $p < 0.05$ ;  $\eta_p^2 = 0.80$ ) and distance measure  
515 ( $F(1, 5) = 12.93$ ;  $p < 0.05$ ;  $\eta_p^2 = 0.76$ ; after  $\epsilon$  correction for lack of sphericity) as well as a significant  
516 interaction of feedback  $\times$  distance measure ( $F(1, 5) = 10.10$ ;  $p < 0.05$ ;  $\eta_p^2 = 0.72$ ; after  $\epsilon$  correction for  
517 lack of sphericity). Follow-up one-factor (feedback) repeated-measures ANOVA revealed a significant  
518 difference for distance to the center ( $F(1, 4) = 57.32$ ;  $p < 0.05$ ;  $\eta_p^2 = 0.94$ ), but not for distance to the  
519 start ( $F(1, 4) = 0.11$ ;  $p = 0.76$ ) and end ( $F(1, 4) = 3.61$ ;  $p = 0.13$ ), respectively. Therefore, feedback  
520 had only a significant and strong effect on the distance to the center. However, due to the small sample,  
521 the distance to the end test was slightly underpowered ( $1 - \beta = 0.798$ ; corresponding to  $\eta_p^2 = 0.474$ ).  
522 Thus, we conclude that force feedback has a differential influence on performance. Whereas force feedback



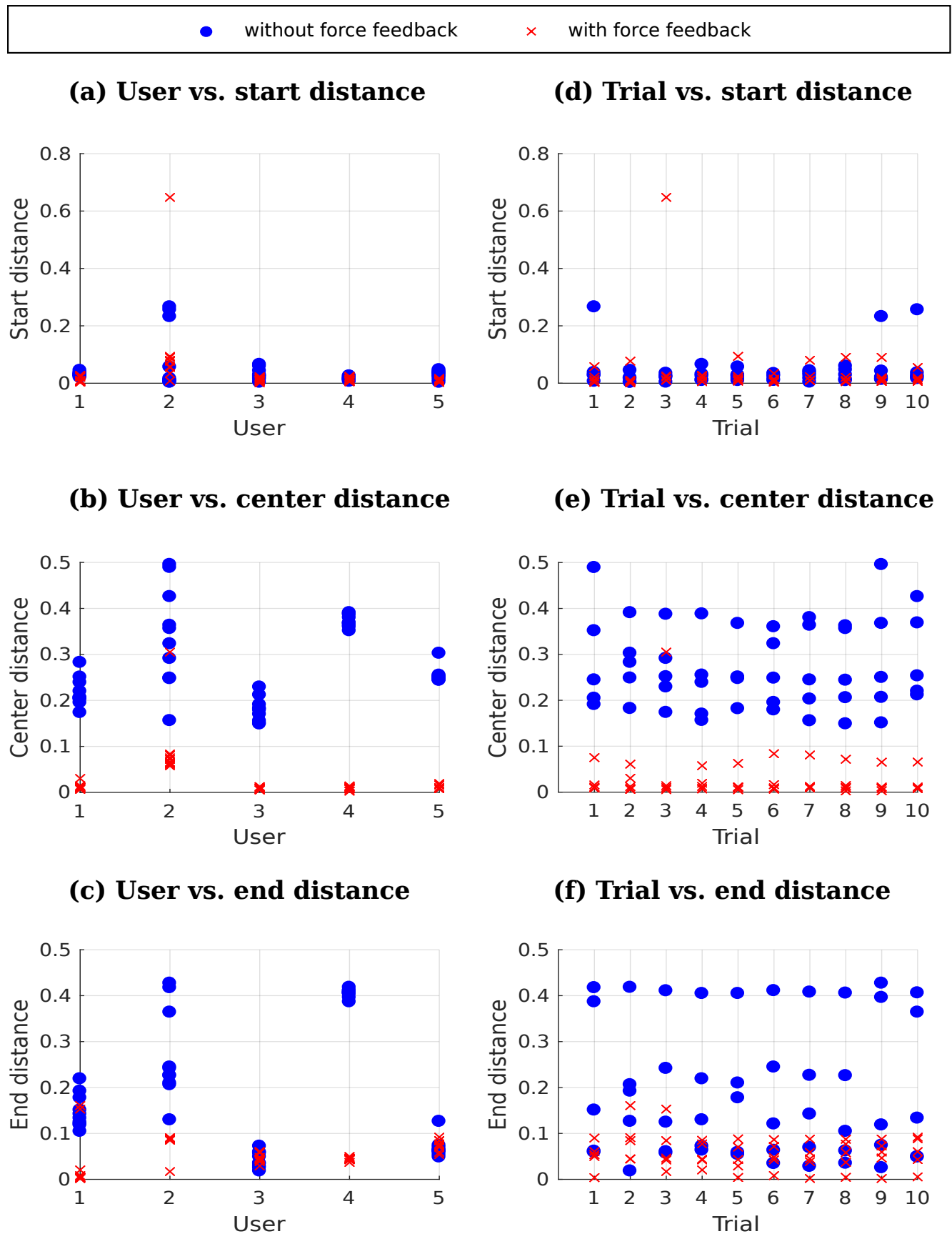
**Figure 13.** (a) Original trajectories. (b) Sample trajectories after Relevance Weighted Policy Optimization (RWPO).



**Figure 14.** (a) Virtual environment. The goal of the user is to teleoperate the green cube to move it from the position marked by the yellow sphere to the position marked by the blue sphere through the window to not hit the wall. (b) Distribution over trajectories before reinforcement learning. (c) Distribution over trajectories after reinforcement learning using the proposed Relevance Weighted Policy Optimization (RWPO) algorithm.

523 does not influence initial error, later errors are expected to be substantially influenced by force feedback.  
 524 However, further studies with bigger samples are required to confirm this conclusion.

525 As it can be seen in Figure 15(c), users 3 and 5 were able to reach the desired end position with  
 526 approximately the same accuracy with and without force feedback. Moreover, it has not been enforced  
 527 in our experiments that users really finish their trials at the end position. Users have been instructed to  
 528 finish their trials both with and without force feedback whenever they thought they have reached the end  
 529 position. We could instead, for the trials with force feedback, instruct users to stop only when they feel  
 530 force feedback contrary to the continuation of the trajectory, which could potentially help minimizing the  
 531 variance of the end distance with force feedback as observed for users 1 and 2.



**Figure 15.** Distances without force feedback and with force feedback. (a,b,c) User versus distances, where each data point corresponds to a different trial. (d,e,f) Trial versus distances, where each data point corresponds to a different user.

## 7 CONCLUSION AND FUTURE WORK

532 This paper presents a probabilistic approach for assisting the practice and the execution of motor tasks by  
533 humans. The method here presented addresses the alignment in space and time of trajectories representing  
534 different executions of a motor task, possibly composed of multiple strokes. Moreover, it addresses  
535 building a probability distribution over demonstrations provided by an expert, which can then be used  
536 to assess a new execution of a motor task by a user. When no expert demonstrations are available, our  
537 system uses a novel reinforcement learning algorithm to learn suitable distributions over trajectories given  
538 performance criteria. This novel algorithm, named Relevance Weighted Policy Optimization, is able to  
539 solve optimization problems with multiple objectives by introducing the concept of relevance functions of  
540 the policy parameters. The relevance functions determine how the policy parameters are sampled when  
541 optimizing the policy for each objective.

542 We evaluated our framework for providing visual feedback to a user practicing the writing of Japanese  
543 characters using a computer mouse. Moreover, we demonstrated how our framework can provide force  
544 feedback to a user, guiding him/her towards correct movements in a teleoperation task involving a haptic  
545 device and a 3D environment.

546 Our algorithm to give visual feedback to the user practicing Japanese characters has still some limitations  
547 that could possibly be addressed by introducing a few heuristics. For example, the current algorithm  
548 assumes that the orientation of the characters is approximately the same. A correct character written in  
549 a different orientation would be deemed wrong by our algorithm. Procrustes Analysis (Goodall, 1991)  
550 provides a solution to align objects with different orientations. Our algorithm could be extended in the  
551 future with a similar technique to give meaningful feedback to the user regardless the orientation of the  
552 characters.

553 In our system, the user has to enter the correct number of strokes to receive feedback. For example, if  
554 the user is practicing a character composed of three strokes, the system waits until the user has drawn  
555 three strokes. Furthermore, the user has to draw the characters in the right order to get meaningful  
556 feedback, otherwise, strokes that do not really correspond to each other are compared. These limitations  
557 can potentially be addressed by analyzing multiple possible alignments and multiple possible stroke orders,  
558 giving feedback to the user according to the alignment and order that result in the best score.

559 Our current framework can give the user feedback concerning the shape of a movement, but not concerning  
560 its speed. In previous work (Ewerton et al., 2016), we have demonstrated how to learn distributions over  
561 shape and phase parameters to represent multiple trajectories with multiple speed profiles. Instead of giving  
562 the user force feedback towards the closest position along the mean trajectory, the distribution over phase  
563 parameters could be used to determine the speed of the attractor along the mean trajectory and how much  
564 the user is allowed to deviate from that speed. This extension shall be made in future work.

565 The framework proposed here could be applied in a scenario where a human would hold a brush connected  
566 to a robot arm. The robot could give the user force feedback to help him/her learn both the position and the  
567 orientation of the brush when writing calligraphy.

568 Especially if our framework can be extended to give feedback to the user concerning the right speed of  
569 a movement, it could potentially be applied in sports. This work could, for example, help users perform  
570 correct movements in weight training, such as in Parisi et al. (2016); Kowsar et al. (2016). Another  
571 possibility would be to help users train golf swings given expert demonstrations or given optimized

572 probability distributions over swings. The training of golf swings could be based on haptic guidance and  
573 use a similar setup as in Kümmel et al. (2014).

574 Future work should also explore further applications of the proposed Relevance Weighted Policy  
575 Optimization algorithm. In particular, it should be verified whether this algorithm can help finding  
576 solutions in more complicated teleoperation scenarios with different performance criteria, favoring, for  
577 example, smooth movements.

## **CONFLICT OF INTEREST STATEMENT**

578 The authors declare that the research was conducted in the absence of any commercial or financial  
579 relationships that could be construed as a potential conflict of interest.

## **AUTHOR CONTRIBUTIONS**

580 ME contributed to this work by writing the manuscript, developing the proposed methods, coding, planning,  
581 preparing and executing the experiments. DR and JWe contributed to the code and to the preparation of  
582 the user experiments. JWi conducted the statistical analysis of our user experiments and contributed to the  
583 writing of the manuscript. GK, JP and GM contributed to the writing of the manuscript.

## **FUNDING**

584 The research leading to these results has received funding from the project BIMROB of the “Forum  
585 für interdisziplinäre Forschung” (FiF) of the TU Darmstadt, from the European Community’s Seventh  
586 Framework Programmes (FP7-ICT-2013-10) under Grant agreement No. 610878 (3rdHand), from the  
587 European Union’s Horizon 2020 research and innovation programme under grant agreement No. 645582  
588 (RoMaNS), from the European Union’s Horizon 2020 research and innovation programme under grant  
589 agreement No. 640554 (SKILLS4ROBOTS) and from a project commissioned by the Japanese New Energy  
590 and Industrial Technology Development Organization (NEDO).

## **REFERENCES**

- 591 Bocsi, B., Csató, L., and Peters, J. (2013). Alignment-based transfer learning for robot models. In *Neural*  
592 *Networks (IJCNN), The 2013 International Joint Conference on (IEEE)*, 1–7
- 593 Coates, A., Abbeel, P., and Ng, A. Y. (2008). Learning for control from multiple demonstrations. In  
594 *Proceedings of the 25th international conference on Machine learning (ACM)*, 144–151
- 595 Collet, A., Berenson, D., Srinivasa, S. S., and Ferguson, D. (2009). Object recognition and full pose  
596 registration from a single image for robotic manipulation. In *Robotics and Automation, 2009. ICRA’09.*  
597 *IEEE International Conference on (IEEE)*, 48–55
- 598 Ernst, M. O. and Banks, M. S. (2002). Humans integrate visual and haptic information in a statistically  
599 optimal fashion. *Nature* 415, 429
- 600 Ewerton, M., Maeda, G., Neumann, G., Kisner, V., Kollegger, G., Wiemeyer, J., et al. (2016). Movement  
601 primitives with multiple phase parameters. In *Robotics and Automation (ICRA), 2016 IEEE International*  
602 *Conference on (IEEE)*, 201–206
- 603 Goodall, C. (1991). Procrustes methods in the statistical analysis of shape. *Journal of the Royal Statistical*  
604 *Society. Series B (Methodological)* , 285–339

- 605 Holladay, R. M. and Srinivasa, S. S. (2016). Distance metrics and algorithms for task space path  
606 optimization. In *Intelligent Robots and Systems (IROS), 2016 IEEE/RSJ International Conference on*  
607 (IEEE), 5533–5540
- 608 Kowsar, Y., Moshtaghi, M., Velloso, E., Kulik, L., and Leckie, C. (2016). Detecting unseen anomalies  
609 in weight training exercises. In *Proceedings of the 28th Australian Conference on Computer-Human*  
610 *Interaction (ACM)*, 517–526
- 611 Kümmel, J., Kramer, A., and Gruber, M. (2014). Robotic guidance induces long-lasting changes in the  
612 movement pattern of a novel sport-specific motor task. *Human movement science* 38, 23–33
- 613 Listgarten, J., Neal, R. M., Roweis, S. T., and Emili, A. (2004). Multiple alignment of continuous time  
614 series. In *Advances in neural information processing systems*. 817–824
- 615 Maeda, G. J., Neumann, G., Ewerton, M., Lioutikov, R., Kroemer, O., and Peters, J. (2016). Probabilistic  
616 movement primitives for coordination of multiple human–robot collaborative tasks. *Autonomous Robots*  
617 , 1–20
- 618 Makondo, N., Rosman, B., and Hasegawa, O. (2015). Knowledge transfer for learning robot models  
619 via local procrustes analysis. In *Humanoid Robots (Humanoids), 2015 IEEE-RAS 15th International*  
620 *Conference on (IEEE)*, 1075–1082
- 621 Paraschos, A., Daniel, C., Peters, J., and Neumann, G. (2013). Probabilistic movement primitives. In  
622 *Advances in Neural Information Processing Systems (NIPS)*. 2616–2624
- 623 Parisi, G. I., Magg, S., and Wermter, S. (2016). Human motion assessment in real time using recurrent self-  
624 organization. In *Robot and Human Interactive Communication (RO-MAN), 2016 25th IEEE International*  
625 *Symposium on (IEEE)*, 71–76
- 626 Peters, J. and Schaal, S. (2007). Reinforcement learning by reward-weighted regression for operational  
627 space control. In *Proceedings of the 24th international conference on Machine learning (ACM)*,  
628 745–750
- 629 Raiola, G., Lamy, X., and Stulp, F. (2015). Co-manipulation with multiple probabilistic virtual guides. In  
630 *Intelligent Robots and Systems (IROS), 2015 IEEE/RSJ International Conference on (IEEE)*, 7–13
- 631 Rosenberg, L. B. (1992). *The Use of Virtual Fixtures as Perceptual Overlays to Enhance Operator*  
632 *Performance in Remote Environments*. Tech. rep., DTIC Document
- 633 Sakoe, H. and Chiba, S. (1978). Dynamic programming algorithm optimization for spoken word recognition.  
634 *Acoustics, Speech and Signal Processing, IEEE Transactions on* 26, 43–49
- 635 Sanguansat, P. (2012). Multiple multidimensional sequence alignment using generalized dynamic time  
636 warping. *WSEAS Transactions on Mathematics* 11, 668–678
- 637 Sigrist, R., Rauter, G., Riener, R., and Wolf, P. (2013). Augmented visual, auditory, haptic, and multimodal  
638 feedback in motor learning: a review. *Psychonomic bulletin & review* 20, 21–53
- 639 Soh, H. and Demiris, Y. (2015). Learning assistance by demonstration: Smart mobility with shared control  
640 and paired haptic controllers. *Journal of Human-Robot Interaction* 4, 76–100
- 641 Solis, J., Avizzano, C. A., and Bergamasco, M. (2002). Teaching to write japanese characters using a haptic  
642 interface. In *Haptic Interfaces for Virtual Environment and Teleoperator Systems, 2002. HAPTICS 2002.*  
643 *Proceedings. 10th Symposium on (IEEE)*, 255–262
- 644 Timmermans, A. A., Seelen, H. A., Willmann, R. D., and Kingma, H. (2009). Technology-assisted training  
645 of arm-hand skills in stroke: concepts on reacquisition of motor control and therapist guidelines for  
646 rehabilitation technology design. *Journal of neuroengineering and rehabilitation* 6, 1
- 647 Tsumugiwa, T., Yokogawa, R., and Hara, K. (2002). Variable impedance control based on estimation of  
648 human arm stiffness for human-robot cooperative calligraphic task. In *Robotics and Automation, 2002.*  
649 *Proceedings. ICRA'02. IEEE International Conference on (IEEE)*, vol. 1, 644–650

- 650 Van Den Berg, J., Miller, S., Duckworth, D., Hu, H., Wan, A., Fu, X.-Y., et al. (2010). Superhuman  
651 performance of surgical tasks by robots using iterative learning from human-guided demonstrations. In  
652 *Robotics and Automation (ICRA), 2010 IEEE International Conference on (IEEE)*, 2074–2081