

# Gaussian-Process-based Robot Learning from Demonstration

Miguel Arduengo<sup>1</sup>, Adrià Colomé<sup>1</sup>, Joan Lobo-Prat<sup>1</sup>, Luis Sentis<sup>2</sup> and Carme Torras<sup>1</sup>

**Abstract**—Endowed with higher levels of autonomy, robots are required to perform increasingly complex manipulation tasks. Learning from demonstration is arising as a promising paradigm for transferring skills to robots. It allows to implicitly learn task constraints from observing the motion executed by a human teacher, which can enable adaptive behavior. We present a novel Gaussian-Process-based learning from demonstration approach. This probabilistic representation allows to generalize over multiple demonstrations, and encode variability along the different phases of the task. In this paper, we address how Gaussian Processes can be used to effectively learn a policy from trajectories in task space. We also present a method to efficiently adapt the policy to fulfill new requirements, and to modulate the robot behavior as a function of task variability. This approach is illustrated through a real-world application using the TIAGo robot.

## I. INTRODUCTION

In the context of robotics, learning from demonstration (LfD) is the paradigm in which robots learn a task policy from examples provided by a human teacher. This facilitates non-expert robot programming, since task constraints and requirements are learned implicitly from the demonstrated motion, which can enable adaptive behavior [1].

Trajectory-based methods are a well-established approach to learn movement policies in robotics [2]. These methods encode skills by extracting trajectory patterns from demonstrations (Figure 1), using a variety of techniques to retrieve a generalized shape of the trajectory. Over the past two decades, it has been an intensive field of study. Among the most relevant contributions, the following methods can be highlighted: Dynamic Movement Primitives (DMP) [3,4], Probabilistic Movement Primitives (ProMP) [5], Gaussian Mixture Model-Gaussian Mixture Regression (GMM-GMR) [6,7], Kernelized Movement Primitives [8,9], and Gaussian Processes (GP) [10,11]. These representations have proved successful at learning and generalizing trajectories. However, each model presents its strengths and shortcomings.

The main advantage of probabilistic-based methods (GMM-GMR, ProMP, KMP and GP) is that they not only retrieve an estimate of the underlying trajectory across multiple demonstrations, but also encode its variability by means of a covariance matrix. This information, which can be inferred from the dispersion of the collected data, can be exploited for the execution of the task, i.e., specifying the robot tracking precision or switching the controller [12].

This work is partially funded by ERC Advanced Grant H2020-741930 (project CLOTHILDE).

<sup>1</sup>Institut de Robotica i Informàtica Industrial (IRI), Barcelona, Spain. {marduengo, acolome, jlobo, torras}@iri.upc.edu

<sup>2</sup>Human Centered Robotics Laboratory (HCRL), UT Austin, USA. lsentis@austin.utexas.edu

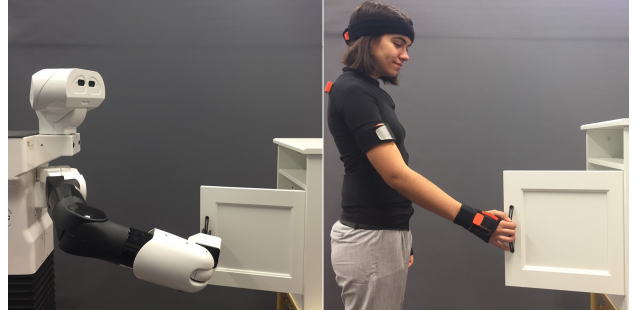


Fig. 1. The proposed Gaussian-Process-based learning from demonstration approach allows to teach the robot manipulation tasks such as opening doors.

Unlike probabilistic-based methods, at the cost of not encoding the variability of the task, DMP only requires a single demonstration. Generalization is achieved by assuming trajectories to be solutions of a deterministic dynamical system, achieving remarkable success. A drawback of DMP, and also ProMP, is that they rely on the manual specification of basis functions, which requires expert knowledge and makes the learning problem with high-dimensional inputs almost intractable. GMM-GMR, in contrast, has proven successful in handling this kind of demonstrations. KMP and GP, due to their kernel treatment, can be implemented for manipulation tasks where high-dimensional inputs and outputs are required.

In LfD, it is also interesting to transfer the learned motion to unseen scenarios while maintaining the general trajectory shape as in the demonstrations. By exploiting the properties of probability distributions, ProMP, KMP and GP allow for trajectory adaptations with via-points. On the other hand, despite GMM-GMR being formulated in terms of Gaussian distributions, the re-optimization of the learned policy requires to re-estimate the model parameters, which lie in a high-dimensional space. This makes the adaptation process very expensive, which prevents its use in unstructured environments, where the policy adjustment is essential.

Besides the generation of adaptive trajectories, another desired property in LfD is extrapolation. In this regard, there is an interesting duality between GMM-GMR and GP representations. The former covariance matrices, model the variability of the trajectories. Conversely, the latter provide a measure of the prediction uncertainty, the variance increasing with the absence of training data. This information is relevant when trying to generalize the learned motion outside of the demonstrated action space. The simultaneous exploitation of both measures is considered in KMP. Moreover, in a recent work [17] Jaquier et al. propose a GMM-based GP for encoding the trajectory.

TABLE I  
COMPARISON AMONG THE STATE-OF-THE-ART AND OUR APPROACH

	DMP [3,13]	ProMP [5]	GMM-GMR [6,14]	GP [15,16]	KMP [8,9]	GMM-GP [17]	Our approach
<i>High-dimensional Learning</i>	–	–	✓	✓	✓	✓	✓
<i>Via-point Adaptation</i>	–	✓	–	–	✓	✓	✓
<i>Task Variability</i>	–	✓	✓	✓	✓	✓	✓
<i>Prediction Uncertainty</i>	–	–	–	✓	✓	✓	✓
<i>Prior Information</i>	–	–	–	✓	–	✓	✓
<i>Task Space Rotations</i>	✓	✓	✓	–	✓	–	✓

In the recent years, there has been a growing interest in Gaussian Processes [18]. The main advantage of GP over the previously discussed methods, is their ability to encode prior beliefs through the mean and kernel functions. This allows the representation of more complex behaviors in the regions of the action space where demonstration data is sparse. A few works have studied the use of an entirely GP-based representation in the LfD context [10,11]. Among the most representative is the one presented by Schneider et al. [15]. They propose a representation of a pick-and-place task that effectively encodes the task variability. Similarly, Umlauf et al. [16] estimate the prediction uncertainty separately, using Wishart Processes. The learned trajectory is retrieved combining GP and DMP. Neither of these works consider the adaptation of the learned policy. Other works formulate the learning and motion planning problem within a single GP-based framework [19,20]. In these works the entire trajectory is retrieved from an optimization perspective. However, this becomes inefficient as the length of the trajectory and the dimensionality of the learning problem increase.

A drawback of GP is that they are usually only defined in Euclidean space, even though a formulation with non-Euclidean input space is possible in principle [21]. Thus, when it comes to the modeling of task space trajectories, representation of orientation imposes great challenges, since is accompanied with additional constraints. This is an aspect disregarded in the aforementioned GP-based methods, which is critical in LfD. Some works have successfully addressed this question with DMP [14], GMM-GMR [13] and KMP [9]. In a recent work, Lang et al. [22] proposed an efficient representation for GP, which we have adopted.

In this work, we present a general Gaussian-Process-based learning from demonstration approach. For the purpose of clear comparison, the main contributions of the state-of-the-art and our approach are summarized in Table I. We aim to unify in a single, entirely GP-based framework, the main features required for a state-of-the-art LfD approach. We show how to achieve an effective representation of the manipulation skill, inferred from the demonstrated trajectories. We unify both, the task variability and the prediction uncertainty, in a single concept we refer to as task uncertainty in the remainder of the paper. Furthermore, in order to achieve an effective generalization across demonstrations, we propose the novel Task Completion Index, for temporal alignment of task trajectories. Finally, we address the adap-

tation of the policy through via-points, and the modulation of the robot behavior depending on the task uncertainty through variable admittance control. The paper is structured as follows: in Section II we discuss the theoretical aspects of the considered GP models; in Section III we present the proposed learning from demonstration framework; in Section IV we illustrate the main aspects of the paper through a real-world application with the TIAGo robot; finally, in Section V, we summarize the final conclusions.

## II. GAUSSIAN PROCESS MODELS

In this section we discuss the theoretical background of the proposed LfD approach. First, we present the fundamentals of GP [23]. Then, we address the challenges of modeling rigid-body dynamics with them. Finally, we present how heteroscedastic GP allows to accurately represent the uncertainty of the taught manipulation task.

### A. Gaussian Process Fundamentals

Intuitively, one can think of a Gaussian process as defining a distribution over functions, and inference taking place directly in the space of functions. Formally, GP are a collection of random variables, any finite number of which have a joint Gaussian distribution [23]. It can be completely specified by its mean  $m(t)$  and covariance  $k(t, t')$  functions:

$$m(t) = \mathbb{E}[f(t)] \quad (1)$$

$$k(t, t') = \mathbb{E}[(f(t) - m(t))(f(t') - m(t'))] \quad (2)$$

where  $f(t)$  is the underlying process,  $m(t)$  depicts the prior knowledge of its mean, and  $k(t, t')$  is symmetric and positive semi-definite (usually referred to as kernel) that must be specified. We are interested in incorporating the knowledge that the training data  $\mathcal{D} = \{(t_i, y_i)\}_{i=1}^N$  provides about  $f(t)$ . We consider that we do not have available direct observations, but only noisy versions  $y$ . Let  $\mathbf{m}(t)$  be the vector of the mean function evaluated at all training points  $t$  and  $K(t, t^*)$  be the matrix of the covariances evaluated at all pairs of training and prediction points  $t^*$ . Assuming additive independent identically distributed Gaussian noise with variance  $\sigma_n^2$ , we can write the joint distribution of the observed target values  $\mathbf{y}$  and the function values at the test locations  $\mathbf{f}^*$  under the prior as:

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{f}^* \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \mathbf{m}(t) \\ \mathbf{m}(t^*) \end{bmatrix}, \begin{bmatrix} K(t, t) + \sigma_n^2 I & K(t, t^*) \\ K(t^*, t) & K(t^*, t^*) \end{bmatrix} \right) \quad (3)$$

The posterior distribution over functions can be computed by conditioning the joint Gaussian prior distribution on the observations  $p(\mathbf{f}^*|t, \mathbf{y}, t^*) \sim \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$  where:

$$\boldsymbol{\mu}^* = \mathbf{m}(t^*) + K(t^*, t) [K(t, t) + \sigma_n^2 I]^{-1} [\mathbf{y} - \mathbf{m}(t)] \quad (4)$$

$$\boldsymbol{\Sigma}^* = K(t^*, t^*) - K(t^*, t) [K(t, t) + \sigma_n^2 I]^{-1} K(t, t^*) \quad (5)$$

When we consider only the prediction of one output variable,  $k(t, t')$  is a scalar function. The previous concepts can be extended to multiple-output GP (MOGP) by taking a matrix covariance function  $\mathbf{k}(t, t')$ . Usual approaches to MOGP modelling are mostly formulated around the Linear Model of Coregionalization (LMC) [24]. For a  $d$ -dimensional output the kernel is expressed in the following form:

$$\mathbf{B} \otimes \mathbf{k}(t, t') = \begin{bmatrix} B_{11}k_{11}(t_1, t'_1) & \dots & B_{1d}k_{1d}(t_1, t'_d) \\ \vdots & \ddots & \vdots \\ B_{d1}k_{d1}(t_d, t'_1) & \dots & B_{dd}k_{dd}(t_d, t'_d) \end{bmatrix} \quad (6)$$

where  $\mathbf{B} \in \mathbb{R}^{d \times d}$  is regarded as the coregionalization matrix and  $t_i$  represents the input corresponding to the  $i$ -th output. Diagonal elements correspond to the single-output case, while the off-diagonal elements represent the prior assumption on the covariance of two different output dimensions [25]. If no a-priori assumption is made,  $B_{ij} = 0$  for  $i \neq j$  and the MOGP is equivalent to  $d$  independent GP.

Regarding the form of  $k(t, t')$ , typically kernel families have free hyperparameters  $\Theta$ . Such parameters can be determined by maximizing the log marginal likelihood:

$$\log p(\mathbf{y}|t, \Theta) = -\frac{1}{2} \mathbf{y}^T K_y^{-1} \mathbf{y} - \frac{1}{2} \log |K_y| - \frac{N}{2} \log 2\pi \quad (7)$$

where  $K_y = K(t, t) + \sigma_n^2 I$ . This optimization problem might suffer from multiple local optima.

### B. Rigid-Body Motion Representation

In the LfD context, representation of trajectories in task space is usually required. However, the modelling of rotations is not straightforward with GP, since the standard formulation is defined for an underlying Euclidean space. A common approach is to use the Euler angles, and exploit that locally the rotation group  $SO(3) \simeq \mathbb{R}^3$ , allowing distances to be computed as Euclidean. However, when this approximation is no longer valid (e.g. at low sampling frequency or if collected data is sparse) it might lead to inaccurate predictions. To overcome this issue, as proposed in [22], rotations can also be represented by a set of unit length Euler axes  $\mathbf{u}$  together with a rotation angle  $\theta$ :

$$SO(3) \subset \{ \boldsymbol{\theta} \mathbf{u} \in \mathbb{R}^3 / \|\mathbf{u}\| = 1 \wedge \theta \in [0, \pi] \} \quad (8)$$

This set defines the solid ball  $B_\pi(0)$  in  $\mathbb{R}^3$  with radius  $0 \leq r \leq \pi$  which is closed, dense and compact. Ambiguity in the representation occurs for  $\theta = \pi$ . To obtain an isomorphism between the rotation group  $SO(3)$  and the axis-angle representation, we can fix the axis representation for  $\theta = \pi$ :

$$\tilde{B}_\pi(0) = B_\pi(0) \setminus \{ \boldsymbol{\pi} \mathbf{u} / u_z < 0 \vee (u_z = 0 \wedge u_y < 0) \vee (u_z = u_y = 0 \wedge u_x < 0) \} \quad (9)$$

where  $\mathbf{u} = (u_x, u_y, u_z)$ . This parametrization is a minimal and unique  $SO(3) \simeq \tilde{B}_\pi(0)$ . Rigid motion dynamics is given by a mapping from time, to translation and rotation  $h: \mathbb{R} \rightarrow SE(3)$ . Let the translational components be defined by the Euclidean vector  $\mathbf{v} \in \mathbb{R}^3$ . Then  $SE(3)$  is defined isomorphically by  $SE(3) \simeq \mathbb{R}^3 \times \tilde{B}_\pi(0)$ . Thus, rigid body motion can be represented in MOGP with the 6-dimensional output vector structure  $(\mathbf{v}, \boldsymbol{\theta} \mathbf{u}) = (x, y, z, \theta u_x, \theta u_y, \theta u_z)$ .

Another possible, more accurate representation, can be achieved with dual quaternions [21]. However, as shown in [22], with the proposed parametrization, a good performance is attained and computations are more efficient.

### C. Heteroscedastic Gaussian Process

The standard Gaussian Process model assumes a constant noise level. This can be an important limitation when encoding a manipulation task. Consider the example shown in Figure 2: it is evident that while the initial and final positions are highly constrained, that is not the case for the path to follow between such positions. In graphs a) and b) we can see that with a standard approach we accurately represent the mean but not the variability of demonstrations.

Considering an independent normally distributed noise,  $\lambda \sim \mathcal{N}(0, r(t))$ , where the variance is input-dependent and modeled by  $r(t)$ . The mean and covariance of the predictive distribution can be modified to [26]:

$$\boldsymbol{\mu}^* = \mathbf{m}(t^*) + K(t^*, t) [K(t, t) + R(t)]^{-1} [\mathbf{y} - \mathbf{m}(t)] \quad (10)$$

$$\boldsymbol{\Sigma}^* = K(t^*, t^*) + R(t^*) - K(t^*, t) [K(t, t) + R(t)]^{-1} K(t, t^*) \quad (11)$$

where  $R(t)$  is a diagonal matrix, with elements  $r(t)$ .

Taking into account the input-dependent noise shown in Figure 2d) the uncertainty in the different phases of the manipulation task is effectively encoded by the uncertainty of Figure 2c). This approach is commonly referred to as heteroscedastic Gaussian Process.

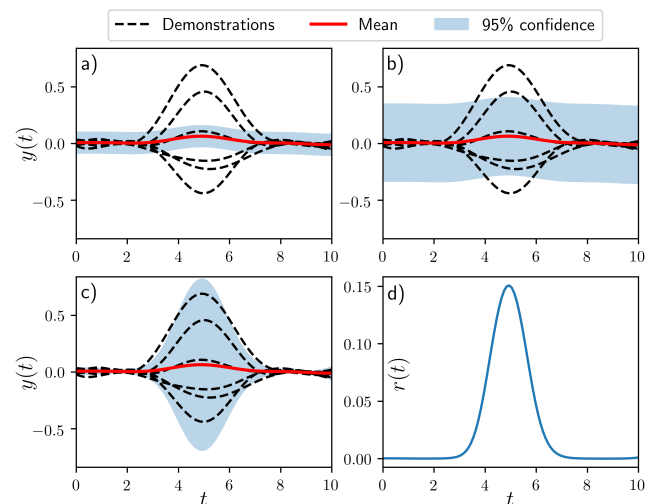


Fig. 2. Standard GP do not accurately model the uncertainty of the demonstrated task as can be seen in a) and b), where it is underestimated and overestimated, respectively. On the other hand, the heteroscedastic GP approach, in c), adequately encodes the uncertainty in the different phases of the task, considering the local noise in d).

### III. LEARNING FROM DEMONSTRATION FRAMEWORK

In this section, we present the proposed GP-based LfD framework. First, we formalize the problem of learning manipulation skills from demonstrated trajectories. Then, we propose an approach for encoding the learned policy with GP. Next, we discuss the temporal alignment of demonstrations. We also present a method that allows to adapt the learned policy through via-points. Finally, we study how the uncertainty model of GP can be exploited to stably modulate the robot behavior, varying end-effector virtual dynamics.

#### A. Problem Statement

In LfD we assume that a dataset of demonstrations is available. In the trajectory-learning case, the dataset consists of a set of trajectories  $\mathbf{s}$  together with a timestamp  $t \in \mathbb{R}$ ,  $\mathcal{D} = \{(t_i, \mathbf{s}_i)\}_{i=1}^N$ . Without loss of generality, we will consider  $\mathbf{s}_i \in SE(3)$ . The aim is to learn a policy  $\pi$  that infers, for a given time, the desired end-effector pose  $\mathbf{s}_i^d$  to perform the taught manipulation task:  $\mathbf{s}_i^d = \pi(t_i)$ . The policy must generate continuous and smooth paths, and generalize over multiple demonstrations.

#### B. Manipulation Task Representation with GP

Representing a manipulation task using heteroscedastic GP models requires the specification of  $m(t)$ ,  $k(t, t')$  and  $r(t)$ . As we have discussed in Section II-B, a suitable mapping for representing a trajectory is given by the following MOGP:

$$\pi(t) \sim \mathcal{GP}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*) : t \longrightarrow (x, y, z, \theta_{u_x}, \theta_{u_y}, \theta_{u_z}) \quad (12)$$

The prior mean function is commonly defined as  $m(t) = 0$ . Although not necessary in general, if no prior knowledge is available this is a simplifying assumption. The GP covariance function controls the policy function shape. The chosen kernel must generate continuous and smooth paths. Note also that the time parametrization of trajectories is invariant to translations in the time domain. Thus, the covariance function must be stationary. That is, it should be a function of  $\tau = t - t'$ . The Radial Basis Function (RBF) kernel fulfils all these requirements:

$$k(t, t') = \sigma_f^2 \exp\left(-\frac{[t - t']^2}{2l^2}\right) \quad (13)$$

with hyperparameters  $l$  and  $\sigma_f$ . Moreover, for multidimensional outputs, we have to consider the prior interaction. In the general case, we usually do not have any previous knowledge about how the different components of the demonstrated trajectories relate to each other. Thus, we can assume that the 6 components are independent a-priori. The matrix covariance function can then be written as:

$$\mathbf{k}(t, t') = \text{diag}\left(\sigma_{f_1}^2 e^{\left(\frac{[t-t']^2}{l_1^2}\right)}, \dots, \sigma_{f_6}^2 e^{\left(\frac{[t-t']^2}{l_6^2}\right)}\right) \quad (14)$$

where  $\text{diag}()$  refers to diagonal, and  $l_i$  and  $\sigma_{f_i}$  correspond to output dimension  $i$ . In Section II-C we discussed the convenience of specifying an input-dependent noise function  $r(t)$  for encoding the manipulation skill with GP. Usually, it

is not known a-priori and must be inferred from the demonstrations. As proposed in [27], first a standard GP can be fit to the data. Its predictions can be used to estimate the input-dependent noise empirically. Then, a second independent GP can be used to model  $z(t) = \log[r(t)]$ . Let  $\mathcal{Z}$  be the set of noise data  $\mathbf{z} = \{z_i\}_{i=1}^n$  and its predictions  $\mathbf{z}^*$ . The posterior predictive distribution can be approximated by:

$$p(\mathbf{f}^* | \mathcal{D}, t^*) = \iint p(\mathbf{f}^* | \mathcal{D}, \mathcal{Z}, t^*) p(\mathcal{Z} | \mathcal{D}, t^*) \simeq p(\mathbf{f}^* | \mathcal{D}, \mathcal{Z}, t^*) \quad (15)$$

where  $\mathcal{Z} = \arg \max_{\mathbf{z}, \mathbf{z}^*} p(\mathbf{z}, \mathbf{z}^* | \mathcal{D}, t^*)$ . Therefore, we have specified all the required functions of the model.

#### C. Temporal Alignment of Demonstrations

For inferring a time dependent policy, the correlation between the temporal and spatial coordinates of two demonstrations of the same task must remain constant. In general, it is very difficult for a human to repeat them at the same velocity. Thus, a time distortion appears (Figure 3a), and should be adequately corrected. Dynamic Time Warping (DTW) [28] is a well-known algorithm for finding the optimal match between two temporal sequences, which may vary in speed.

The algorithm finds a non-linear mapping of the demonstrated trajectories and a reference based on a similarity measure. A common measure in the LfD context is the Euclidean distance. This relies on the assumption that the manipulation task can be performed always following the same path. For instance, consider the case of a pick-and-place task where the objects have to be placed in shelves at different levels. Using the Euclidean distance as similarity measure will lead to an erroneous temporal alignment (Figure 3b), since intermediate points for placing the object at a higher level can be mapped to ending points of a lower level. We propose to use an index which considers the portion of the trajectory that has been covered for task completion as a similarity measure. We will refer to it as the Task Completion Index (TCI). We define it in discrete form as:

$$\zeta(t_k) = \frac{\sum_{j=1}^k d(\mathbf{s}_j, \mathbf{s}_{j-1})}{\sum_{j=1}^M d(\mathbf{s}_j, \mathbf{s}_{j-1})} \quad \forall k = 1, \dots, M \quad (16)$$

where  $\mathbf{s}_j \in SE(3)$  refers to the trajectory point at time instant  $t_j$ ,  $d(\cdot)$  to a scalar distance function and  $M$  to the total number of discrete points. Note that  $0 = \zeta(t_0) \leq \zeta(t_k) \leq \zeta(t_M) = 1$ . As a distance function on  $SE(3)$ , using the representation discussed in Section II-B, we define:

$$d(\mathbf{s}_i, \mathbf{s}_j) = \sqrt{\omega_1 [d_{\text{arc}}(\boldsymbol{\theta}_i \mathbf{u}_i, \boldsymbol{\theta}_j \mathbf{u}_j)]^2 + \omega_2 \|\mathbf{v}_i - \mathbf{v}_j\|^2} \quad (17)$$

where  $\omega_k$  are a convex combination of weights for application dependent scaling and  $d_{\text{arc}}(\cdot)$  is the length of the geodesic between rotations [22]:

$$d_{\text{arc}}(\boldsymbol{\theta}_i \mathbf{u}_i, \boldsymbol{\theta}_j \mathbf{u}_j) = 2 \arccos \left| \cos \frac{\theta_i}{2} \cos \frac{\theta_j}{2} + \sin \frac{\theta_i}{2} \sin \frac{\theta_j}{2} \mathbf{u}_i^T \mathbf{u}_j \right| \quad (18)$$

In Figure 3c we show that the trajectories are warped correctly, allowing then an effective encoding of the manipulation task, with the proposed TCI (Figure 3d).

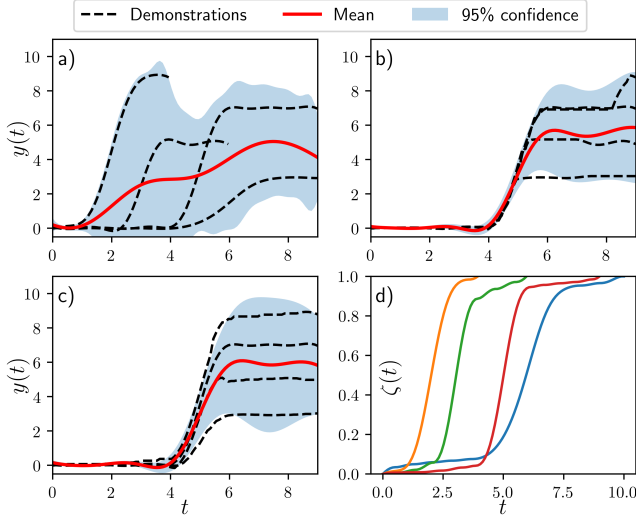


Fig. 3. In a) we observe that due to distortion in time, task constraints are not encoded correctly. In b) the trajectories are aligned with DTW using the Euclidean distance as similarity measure. In c) we show the resulting alignment using the proposed TCI d), as similarity measure.

#### D. Policy Adaptation through Via-points

The modulation of the learned policy through via-points is an important property to adapt to new situations. Let  $\mathcal{V} = \{(t_i, s_i^v)\}$  be the set of via-points  $s_i^v$  which are desired to be reached by the policy at time instant  $t_i$ . In the proposed probabilistic framework, generalization can be implemented by conditioning the policy on both  $\mathcal{D}$  and  $\mathcal{V}$ . Assuming that the predictive distribution of each set can be computed independently, the conditioned policy is [29]:

$$p(\mathbf{f}^*|\mathcal{D}, \mathcal{V}, t^*) = p(\mathbf{f}^*|\mathcal{D}, t^*) p(\mathbf{f}^*|\mathcal{V}, t^*) \quad (19)$$

If  $p(\mathbf{f}^*|\mathcal{D}, t^*) \sim \mathcal{N}(\boldsymbol{\mu}^d, \boldsymbol{\Sigma}^d)$  and  $p(\mathbf{f}^*|\mathcal{V}, t^*) \sim \mathcal{N}(\boldsymbol{\mu}^v, \boldsymbol{\Sigma}^v)$ , then, it holds that  $p(\mathbf{f}^*|\mathcal{D}, \mathcal{V}, t^*) \sim \mathcal{N}(\boldsymbol{\mu}^{**}, \boldsymbol{\Sigma}^{**})$  where:

$$\boldsymbol{\mu}^{**} = \boldsymbol{\Sigma}^v (\boldsymbol{\Sigma}^d + \boldsymbol{\Sigma}^v)^{-1} \boldsymbol{\mu}^d + \boldsymbol{\Sigma}^d (\boldsymbol{\Sigma}^d + \boldsymbol{\Sigma}^v)^{-1} \boldsymbol{\mu}^v \quad (20)$$

$$\boldsymbol{\Sigma}^{**} = \boldsymbol{\Sigma}^d (\boldsymbol{\Sigma}^d + \boldsymbol{\Sigma}^v)^{-1} \boldsymbol{\Sigma}^v \quad (21)$$

The resulting distribution is computed as a product of Gaussians, and is a compromise between the via-point constraints and the demonstrated trajectories, weighted inversely by their variances. Considering an heteroscedastic GP model for  $\mathcal{V}$  (equations 10 and 11), the strength of the via-point constraints can then be easily specified by means of the latent noise function. For instance, via-points with low noise will have a higher relative weight, modifying significantly the learned policy. On the other hand, via-points with a high noise level will produce a more subtle effect. In Figure 4 we illustrate how the distribution adapts to strong and weak defined via-points.

It should be remarked that the posterior predictive distribution of  $\mathcal{D}$  only needs to be computed once. Thus, adaptation of the policy just involves a computational cost of  $\mathcal{O}(m^3)$ , where  $m$  is the number of predicted outputs. Since  $m$  can be specified, the proposed approach is suitable for on-line applications (for further insight on GP complexity see [30]).

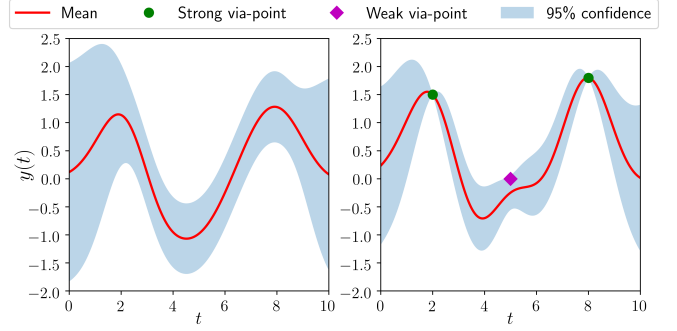


Fig. 4. On the left, a GP model based on the demonstrated trajectories. On the right, the policy adapted through via-points.

#### E. Modulation of the Robot Behavior

In LfD, it is often convenient to adapt the behavior of the robot as a function of the uncertainty in the different phases of the task. Let the robot end-effector be controlled through a virtual spring-mass-damper model dynamics:

$$\mathbf{M}(t)\ddot{\mathbf{e}}(t) + \mathbf{D}(t)\dot{\mathbf{e}}(t) + \mathbf{K}_p(t)\mathbf{e}(t) = \mathbf{F}_{\text{ext}}(t) \quad (22)$$

where  $\mathbf{M}(t), \mathbf{D}(t), \mathbf{K}_p(t) \in \mathbb{R}^{6 \times 6}$  refer to inertia, damping and stiffness, respectively, and  $\mathbf{e}(t) \in \mathbb{R}^{6 \times 1}$  is the tracking error, when subjected to an external force  $\mathbf{F}_{\text{ext}}(t) \in \mathbb{R}^{6 \times 1}$ .

It can be proved (see [31]) that for a constant, symmetric, positive definite  $\mathbf{M}$ , and  $\mathbf{D}(t), \mathbf{K}_p(t)$  continuously differentiable, the system is globally asymptotically stable if there exists a  $\gamma > 0$  such that:

- 1)  $\gamma\mathbf{M} - \mathbf{D}(t)$  is negative semidefinite
- 2)  $\dot{\mathbf{K}}_p(t) + \gamma\dot{\mathbf{D}}(t) - 2\gamma\mathbf{K}_p(t)$  is negative definite

Now consider  $\mathbf{M}, \mathbf{D}(t)$  and  $\mathbf{K}_p(t)$  diagonal matrices, and a constant damping ratio  $\delta$ . Substituting  $d(t) = 2\delta\sqrt{mk_p(t)}$  on the second stability condition, it yields the following upper bound for the stiffness derivative:

$$\dot{k}_p(t) < \frac{2\gamma\sqrt{k_p(t)}^3}{\sqrt{k_p(t)} + 2\delta\gamma\sqrt{m}} \quad (23)$$

where  $m$  and  $k_p(t)$  are an arbitrary diagonal element of  $\mathbf{M}$  and  $\mathbf{K}_p(t)$ , respectively. In order to modulate the robot behavior, we propose the following variable stiffness profile:

$$k_p(t) = k_p^{\max} - \frac{k_p^{\max} - k_p^{\min}}{1 + e^{-\alpha(\sigma(t) - \beta)}} \quad (24)$$

which increases the stiffness inversely to the uncertainty  $\sigma(t)$  and saturates at  $k_p^{\min}$  and  $k_p^{\max}$  for high and low values respectively. Differentiating we have:

$$\dot{k}_p(t) = \alpha k_p(t) \left( 1 - \frac{k_p(t)}{k_p^{\max} - k_p^{\min}} \right) \frac{d\sigma(t)}{dt} \quad (25)$$

For a constant  $d\sigma(t)/dt$ , the maximum value of the stiffness derivative  $\dot{k}_p(t)$  is obtained for  $k_p(t) = (k_p^{\max} - k_p^{\min})/2$ . Thus, substituting in (25), it yields the following upper bound:

$$\dot{k}_p(t) \leq \frac{\alpha}{4} (k_p^{\max} - k_p^{\min}) \frac{d\sigma(t)}{dt} \quad (26)$$



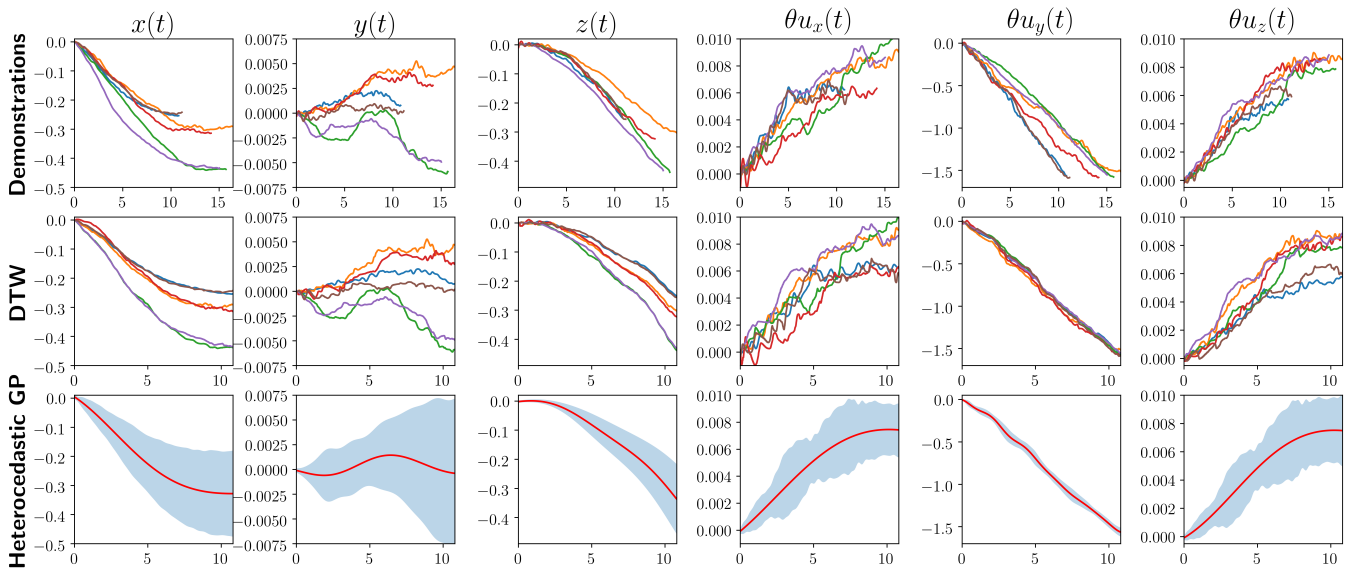


Fig. 6. In the first row, demonstrated trajectories. In the second row, alignment with DTW using TCI index. In the third row, heteroscedastic GP of each dimension of the MOGP encoding the learned policy. Note that significant variations are only observed in  $x$ ,  $z$  and  $\theta u_y$  (see vertical axis scale).

Then, from inspection of the first stability condition, we can see that  $\gamma$  defines a lower bound for the minimum allowed damping  $d(t)$ . Thus, given the variable stiffness profile in Equation 24, and assuming constant damping ratio, the most restrictive value is  $\gamma = 2\delta\sqrt{k_p^{min}/m}$ . Substituting in (23), we can obtain the following lower bound:

$$\dot{k}_p(t) < \frac{4\delta\sqrt{(k_p^{min})^3}}{(1+4\delta^2)\sqrt{m}} \leq \frac{2\gamma\sqrt{k_p(t)^3}}{\sqrt{k_p(t)} + 2\delta\gamma\sqrt{m}} \quad (27)$$

Then, from equations (26) and (27) the following sufficient stability condition can be derived:

$$\frac{d\sigma(t)}{dt} < \frac{16\delta}{\alpha} \frac{\sqrt{(k_p^{min})^3}}{(k_p^{max} - k_p^{min})(1+4\delta^2)\sqrt{m}} \quad (28)$$

The control parameters can then be tuned to ensure the satisfaction of this inequality. Note that sharper uncertainty profiles  $\sigma(t)$  are more restrictive with respect to variations of the stiffness. For instance, stability is favored by a smaller range  $(k_p^{max} - k_p^{min})$  or lower values of  $\alpha$ , i.e. slower transition between stiff and compliant behaviors. For the limit cases  $k_p^{max} \rightarrow k_p^{min}$  and  $\alpha \rightarrow 0$ , that is, constant stiffness, stability can be achieved regardless of  $\sigma(t)$ . It can also be observed, since the right-hand side of the inequality is always positive, that with the proposed variable stiffness profile, stability is ensured if the uncertainty decreases.

#### IV. AN EXAMPLE APPLICATION: DOOR OPENING TASK

In order to illustrate how the proposed GP-based LfD approach can be applied to real-world manipulation tasks, we address the problem of opening doors using a TIAGo robot. This is a relevant skill for robots operating in domestic environments [32], since they need to open doors when navigating, to pick up objects in fetch-and-carry applications or assist people in their mobility.

#### A. Policy Inference from Human Demonstrations

We performed human demonstrations using an Xsens MVN motion capture system. Right hand trajectories of the human teacher relative to the initial closed door position were recorded for three different doors (Figure 5). Coordinate axes were chosen such as the pulling direction is parallel to the  $x$  axis and the  $y$  axis is perpendicular to the floor. The demonstration dataset consisted in a total of 6 trajectories, two per each door.



Fig. 5. Demonstrations were recorded using an Xsens MVN motion capture system. The teacher opens three doors with different radius.

The recorded trajectories were then temporally aligned using DTW and the proposed TCI index. Next, the data was used to infer the door opening policy. In Figure 6 we show these steps for each output dimension. On the third row, we can see the resulting heteroscedastic MOGP representation. Note that in the door opening motion, significant variations are only observed in the  $x$ ,  $z$  and  $\theta u_y$  components. We can see that the trajectories are warped effectively using the proposed TCI similarity measure, since they are clearly clustered in three groups, one for each type of door. It can also be observed that the resulting heteroscedastic MOGP effectively encodes the skill. This is more evident in Figure 7, where position uncertainty has been projected onto the  $x-z$  plane. We can observe the uncertainty in the door radius is accurately captured from demonstrations.

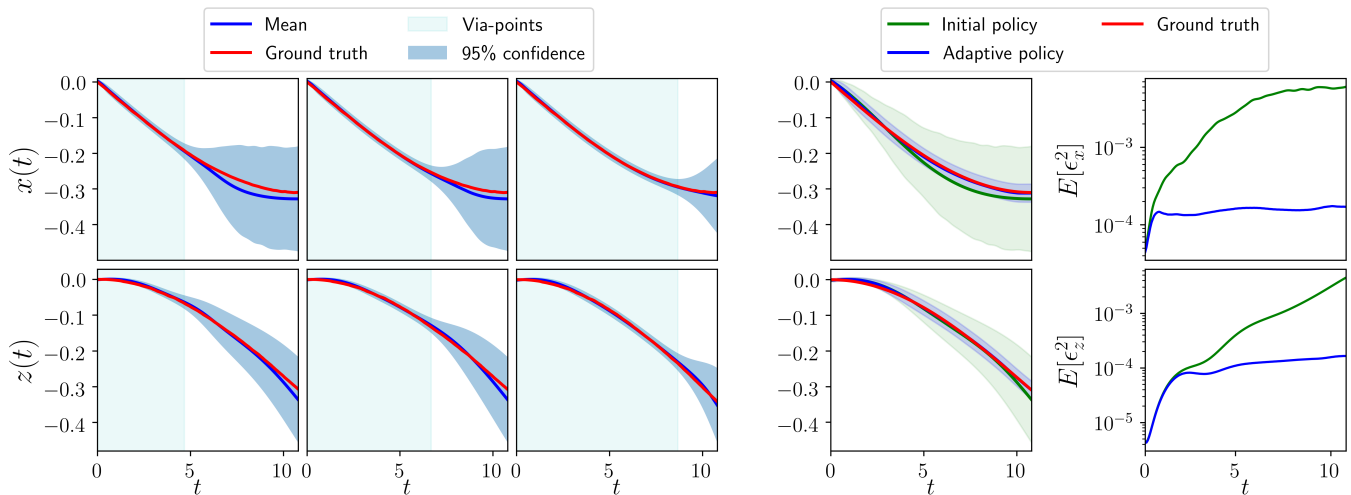


Fig. 9. The first three columns show how does the posterior distribution vary considering as via-points the observations of the door motion in the light-blue shaded area. The next column shows the comparison between the predictive distribution considering the adaptive policy or the policy based only on human demonstrations. The shaded areas and lines of the same color correspond to the 95% confidence interval and mean, respectively. The column on the right shows the mean squared prediction error of each policy.

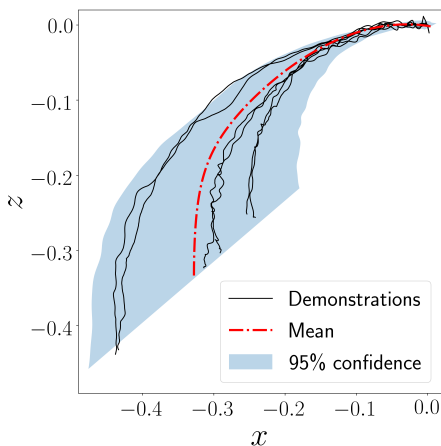


Fig. 7. Door opening policy projected on the  $x-z$  plane.

### B. Policy Adaptation and Modulation of the Robot Behavior

In the reproduction stage, we gather observations of the door motion solving the forward kinematics of the robot. These observations can then be defined as via-points to improve the policy prediction capabilities. Additionally, we can take advantage of the forces exerted by the door to correct small biases on the policy, adopting a variable admittance control scheme. The set-point of the controller is changed through the virtual dynamics discussed in Section III-E.



Fig. 8. TIAGo robot opening the door.

With the proposed approach we successfully performed the door opening task with TIAGo (Figure 8). We can see in Figure 9, on the first three columns on the left, how the posterior distribution varies for coordinates  $x$  and  $z$  as the door is opened. We show that the adaptive policy converges to the ground truth trajectory. On the next column, we show the comparison between the resulting distribution at step  $t_i$  considering via-points up to  $t_{i-1}$ , and the initial policy. In order to obtain a quantitative measure of the prediction performance, we evaluated the mean squared prediction error,  $E[\epsilon^2] = (E[f^*(t^*)] - f(t^*))^2 + \text{var}[f^*(t^*)]$ . We can observe that the adaptive policy clearly achieves a better performance.

The resulting variable stiffness profile is shown in Figure 10. We have tuned the parameters empirically, being the used values  $k_p^{max} = 500$ ,  $k_p^{min} = 100$ ,  $m = 1$ ,  $\delta = 1$ ,  $\alpha = 600$  and  $\beta = 0.01$ . For simplicity, we have considered the same law for the 6 degrees of freedom. We can observe that the robot behavior is modulated towards a more compliant behavior towards the final phases, where the policy is more uncertain. We can also see that the stability bound is not crossed, which is coherent with the behavior observed in the conducted experiments, where no instabilities occurred.

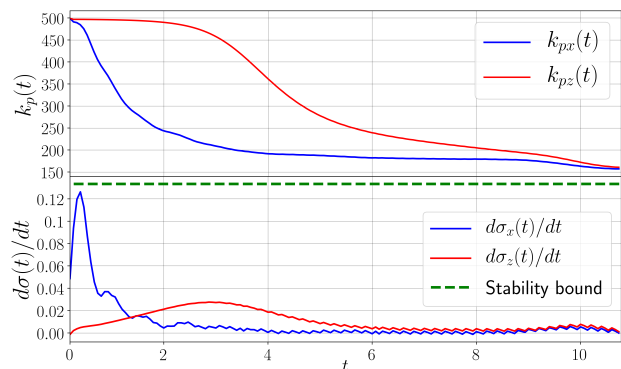


Fig. 10. On top, variable stiffness profile for  $x$  and  $z$ . Below, the evolution of the uncertainty derivative of the adaptive policy.

## V. CONCLUSION

Gaussian Processes (GP) are a promising paradigm for learning manipulation skills from human demonstrations. In this paper, we present a novel approach that takes advantage of the versatility and expressiveness of these models to encode task policies. We propose an heteroscedastic multi-output GP policy representation, inferred from demonstrations. This model considers a suitable parametrization of task space rotations for GP and ensures that only continuous and smooth paths are generated. Furthermore, the introduction of an input-dependent latent noise function allows an effective simultaneous encoding of the prediction uncertainty and the variability of demonstrated trajectories.

In order to effectively establish a correlation between temporal and spatial coordinates, demonstrations must be aligned. This operation can be performed with the Dynamic Time Warping algorithm. We introduce the novel Task Completion Index, a similarity measure that allows to achieve an effective warping when the learned task requires the consideration of different paths. Adaptation of the policy can be performed by conditioning it on a set of specified via-points. We also introduce a new computationally efficient method, where the relative importance of the constraints can also be defined. Additionally, we propose an innovative variable stiffness profile that takes advantage of the uncertainty measure provided by the GP model to stably modulate the robot end-effector dynamics.

We illustrated the proposed learning from demonstration framework through the door opening task and evaluated the performance of the learned policy through real-world experiments with the TIAGo robot. Results show that the manipulation skill is effectively encoded and a successful reproduction can be achieved taking advantage of the presented policy adaptation and robot behavior modulation approaches.

This work aims to push the state-of-the-art in learning from demonstration towards easily extending robot capabilities. Future research will be conducted focusing on its applicability on complex tasks, such as cloth manipulation.

## REFERENCES

- [1] H. Ravichandar, A. Polydoros, S. S. Chernova, and A. Billard, "Recent advances in robot learning from demonstration," *Annual Review of Control, Robotics, and Autonomous Systems*, vol. 3, May 2020.
- [2] A. Colomé and C. Torras, *Reinforcement Learning of Bimanual Robot Skills*. Springer Tracts in Advanced Robotics (STAR), vol. 134. Springer International Publishing, 2020.
- [3] A. J. Ijspeert, J. Nakanishi, and S. Schaal, "Trajectory formation for imitation with nonlinear dynamical systems," in *Proceedings 2001 IEEE/RSJ International Conference on Intelligent Robots and Systems*, vol. 2, Oct 2001, pp. 752–757.
- [4] P. Pastor, H. Hoffmann, T. Asfour, and S. Schaal, "Learning and generalization of motor skills by learning from demonstration," in *IEEE International Conference on Robotics and Automation*, May 2009, pp. 763–768.
- [5] A. Paraschos, C. Daniel, J. Peters, and G. Neumann, "Using probabilistic movement primitives in robotics," *Autonomous Robots*, vol. 42, no. 3, pp. 529–551, March 2018.
- [6] S. Calinon, "A tutorial on task-parameterized movement learning and retrieval," *Intelligent Service Robotics*, vol. 9, no. 1, Jan 2016.
- [7] E. Pignat and S. Calinon, "Bayesian Gaussian Mixture Model for robotic policy imitation," *IEEE Robotics and Automation Letters (RA-L)*, vol. 4, no. 4, pp. 4452–4458, 2019.
- [8] Y. Huang, L. Rozo, J. Silvério, and D.-G. Caldwell, "Non-parametric imitation learning of robot motor skills," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2019, pp. 5266–5272.
- [9] Y. Huang, F. J. Abu-Dakka, J. Silvério, and D.-G. Caldwell, "Generalized orientation learning in robot task space," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2019.
- [10] D. Nguyen-Tuong and J. Peters, "Local Gaussian Process regression for real-time model-based robot control," in *IEEE/RSJ International Conference on Intelligent Robots and Systems*, Sept 2008.
- [11] D. Forte, A. Ude, and A. Kos, "Robot learning by Gaussian Process regression," in *19th International Workshop on Robotics in Alpe-Adria-Danube Region (RAAD 2010)*, June 2010, pp. 303–308.
- [12] J. Silvério, Y. Huang, L. Rozo, S. Calinon, and D. G. Caldwell, "Probabilistic learning of torque controllers from kinematic and force constraints," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Oct 2018, pp. 1–8.
- [13] L. Koutras and Z. Doulgeri, "A correct formulation for the Orientation Dynamic Movement Primitives for robot control in the Cartesian space," in *3rd Conference on Robot Learning (CoRL)*, Osaka, 11 2019.
- [14] M.-J. Zeestraten, I. Havoutis, J. Silvério, S. Calinon, and D.-G. Caldwell, "An approach for imitation learning on riemannian manifolds," *IEEE Robotics and Automation Letters*, vol. 2, pp. 1240–1247, 2017.
- [15] M. Schneider and W. Ertel, "Robot learning by demonstration with local Gaussian Process regression," in *2010 IEEE/RSJ International Conference on Intelligent Robots and Systems*, Oct 2010.
- [16] J. Umlauf, Y. Fanger, and S. Hirche, "Bayesian uncertainty modeling for programming by demonstration," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2017, pp. 6428–6434.
- [17] N. Jaquier, D. Ginsbourger, and S. Calinon, "Learning from demonstration with model-based Gaussian Process," in *3rd Conference on Robot Learning (CoRL)*, Osaka, Japan, Oct 2019.
- [18] E. Schulz, M. Speekenbrink, and A. Krause, "A tutorial on Gaussian Process regression: Modelling, exploring, and exploiting functions," *Journal of Mathematical Psychology*, vol. 85, pp. 1–16, 2018.
- [19] M.-A. Rana, M. Mukadam, S.-R. Ahmadzadeh, S. Chernova, and B. Boots, "Towards robust skill generalization: unifying learning from demonstration and motion planning," in *1st Conference on Robot Learning (CoRL)*, CA, USA, 10 2017.
- [20] T. Osa, N. Sugita, and M. Mitsuishi, "Online trajectory planning and force control for automation of surgical tasks," *IEEE Transactions on Automation Science and Engineering*, vol. 15, pp. 675–691, 2018.
- [21] M. Lang, M. Kleinstueber, and S. Hirche, "Gaussian Process for 6-DoF rigid motions," *Autonomous Robots*, vol. 42, no. 6, 2018.
- [22] M. Lang and S. Hirche, "Computationally efficient rigid-body Gaussian Process for motion dynamics," *IEEE Robotics and Automation Letters*, vol. 2, no. 3, pp. 1601–1608, July 2017.
- [23] C.-E. Rasmussen and C.-K.-I. Williams, *Gaussian Processes for Machine Learning*. The MIT Press, 2006.
- [24] M.-A. Álvarez, L. Rosasco, and N.-D. Lawrence, "Kernels for vector-valued functions: A review," *Foundations and Trends in Machine Learning*, vol. 4, no. 3, pp. 195–266, Mar 2012.
- [25] H. Liu, J. Cai, and Y.-S. Ong, "Remarks on Multi-Output Gaussian Process Regression," *Knowledge-Based Systems*, vol. 144, pp. 102–121, March 2018.
- [26] P. Goldberg, C. Williams, and C. Bishop, "Regression with input-dependent noise: A Gaussian Process treatment," *Advances in Neural Information Processing Systems*, vol. 10, pp. 493–499, Jan 1998.
- [27] K. Kersting, C. Plagemann, P. Pfaff, and W. Burgard, "Most-likely Heteroscedastic Gaussian Process regression," in *ACM International Conference Proceeding Series*, vol. 227, Jan 2007, pp. 393–400.
- [28] P. Senin, "Dynamic Time Warping algorithm review," Information and Computer Science Department, University of Hawaii at Manoa, Honolulu, USA, Tech. Rep., Dec 2008.
- [29] M. Deisenroth and J.-W. Ng, "Distributed Gaussian Processes," in *32nd International Conference on Machine Learning (ICML)*, vol. 37, July 2015, pp. 1481–1490.
- [30] H.-L. Bilj, "LQG and Gaussian Process Techniques for Fixed-Structure Wind Turbine Control," PhD Dissertation, Delft University of Technology, The Netherlands, Tech. Rep., Oct 2018.
- [31] K. Kronander and A. Billard, "Stability considerations for variable impedance control," *IEEE Transactions on Robotics*, vol. 32, no. 5, pp. 1298–1305, Oct 2016.
- [32] D. Kim, J.-H. Kang, C.-S. H. CS., and G.-T. Park, "Mobile robot for door opening in a house," in *Knowledge-Based Intelligent Information and Engineering Systems (KES), Part II*, Sept 2004, pp. 596–602.