
Probabilistic Movement Primitives

Alexandros Paraschos¹, Christian Daniel¹, Jan Peters^{1,2}, and Gerhard Neumann¹

1: Intelligent Autonomous Systems, Technische Universität Darmstadt
Hochschulstr. 10, 64289 Darmstadt, Germany

{paraschos,daniel,peters,neumann}@ias.tu-darmstadt.de

2: Max Planck Institute for Intelligent Systems
Spemannstr. 41, 72076 Tübingen, Germany

Abstract

Movement Primitives (MP) are a well-established approach for representing modular and re-usable robot movement generators. Many state-of-the-art robot learning successes are based MP, due to their compact representation of the inherently continuous and high dimensional robot movements. A major goal in robot learning is to combine multiple MPs as building blocks in a modular control architecture to solve complex tasks. To this effect, a MP representation has to allow for blending between motions, adapting to altered task variables, and co-activating multiple MPs in parallel. We present a probabilistic formulation of the MP concept that maintains a distribution over trajectories. Our probabilistic approach allows for the derivation of new operations which are essential for implementing all aforementioned properties in one framework. In order to use such a trajectory distribution for robot movement control, we analytically derive a stochastic feedback controller which reproduces the given trajectory distribution. We evaluate and compare our approach to existing methods on several simulated as well as real robot scenarios.

1 Introduction

Movement Primitives (MPs) are commonly used for representing and learning basic movements in robotics, e.g., hitting and batting, grasping, etc. [1, 2, 3]. MP formulations are compact parameterizations of the robot’s control policy. Modulating their parameters permits imitation and reinforcement learning as well as adapting to different scenarios. MPs have been used to solve many complex tasks, including ‘Ball-in-the-Cup’ [4], Ball-Throwing [5, 6], Pancake-Flipping [7] and Tetherball [8].

The aim of MPs is to allow for composing complex robot skills out of elemental movements with a modular control architecture. Hence, we require a MP architecture that supports parallel activation and smooth blending of MPs for composing complex movements of sequentially [9] and simultaneously [10] activated primitives. Moreover, adaptation to a new task or a new situation requires modulation of the MP to an altered desired target position, target velocity or via-points [3]. Additionally, the execution speed of the movement needs to be adjustable to change the speed of, for example, a ball-hitting movement. As we want to learn the movement from data, another crucial requirement is that the parameters of the MPs should be straightforward to learn from demonstrations as well as through trial and error for reinforcement learning approaches. Ideally, the same architecture is applicable for both stroke-based and periodic movements, and capable of representing optimal behavior in deterministic and stochastic environments.

While many of these properties are implemented by one or more existing MP architectures [1, 11, 10, 2, 12, 13, 14, 15], no approach exists which exhibits all of these properties in *one* framework. For

example, [13] also offers a probabilistic interpretation of MPs by representing an MP as a learned graphical model. However, this approach heavily depends on the quality of the used planner and the movement can not be temporally scaled. Rozo et. al. [12, 16] use a combination of primitives, yet, their control policy of the MP is based on heuristics and it is unclear how the combination of MPs affects the resulting movements.

In this paper, we introduce the concept of probabilistic movement primitives (ProMPs) as a general probabilistic framework for representing and learning MPs. Such a ProMP is a distribution over trajectories. Working with distributions enables us to formulate the described properties by operations from probability theory. For example, modulation of a movement to a novel target can be realized by conditioning on the desired target’s positions or velocities. Similarly, consistent parallel activation of two elementary behaviors can be accomplished by a product of two independent trajectory probability distributions. Moreover, a trajectory distribution can also encode the variance of the movement, and, hence, a ProMP can often directly encode optimal behavior in stochastic systems [17]. Finally, a probabilistic framework allows us to model the covariance between trajectories of different degrees of freedom, that can be used to couple the joints of the robot.

Such properties of trajectory distributions have so far not been properly exploited for representing and learning MPs. The main reason for the absence of such an approach has been the difficulty of extracting a policy for controlling the robot from a trajectory distribution. We show how this step can be accomplished and derive a control policy that exactly reproduces a given trajectory distribution. To the best of our knowledge, we present the first principled MP approach that can exploit the power of operations from probability theory.

While the ProMPs’ representation introduces many novel components, it incorporates many advantages from well-known previous movement primitive representations [18, 10], such as phase variables for timing of the movement that enable temporal rescaling of movements, and the ability to represent both rhythmic and stroke based movements. However, since ProMPs incorporate the variance of demonstrations, the increased flexibility and advantageous properties of the representation come at the price of requiring multiple demonstrations to learn the primitives as opposed to past approaches [18, 3] that can clone movements from a single demonstration.

2 Probabilistic Movement Primitives (ProMPs)

A movement primitive representation should exhibit several desirable properties, such as co-activation, adaptability and optimality in order to be a powerful MP representation. The goal of this paper is to unify these properties in one framework. We accomplish this objective by using a probabilistic formulation for MPs. We summarized all the properties and how they are implemented in our framework in Table 1. In this section, we will sequentially explain the importance of each of these property and discuss the implementation in our framework. As crucial part of our objective, we will introduce conditioning and a product of ProMPs as new operations that can be applied on the ProMPs due to the probabilistic formulation. Finally, we show how to derive a controller which follows a given trajectory distribution.

Table 1: Desirable properties and their implementation in the ProMP

| Property | Implementation |
|--------------------|------------------|
| Co-Activation | Product |
| Modulation | Conditioning |
| Optimality | Encode variance |
| Coupling | Mean, Covariance |
| Learning | Max. Likelihood |
| Temporal Scaling | Modulate Phase |
| Rhythmic Movements | Periodic Basis |

2.1 Probabilistic Trajectory Representation

We model a single movement execution as a trajectory $\tau = \{q_t\}_{t=0\dots T}$, defined by the joint angles q_t over time. In our framework, a MP describes multiple ways to execute a movement, which naturally leads to a probability distribution over trajectories.

Encoding a Time-Varying Variance of Movements. Our movement primitive representation models the time-varying variance of the trajectories to be able to capture multiple demonstrations

with high-variability. Representing the variance information is crucial as it reflects the importance of single time points for the movement execution and it is often a requirement for representing optimal behavior in stochastic systems [17].

We use a weight vector \mathbf{w} to compactly represent a single trajectory. The probability of observing a trajectory τ given the underlying weight vector \mathbf{w} is given as a linear basis function model

$$\mathbf{y}_t = \begin{bmatrix} q_t \\ \dot{q}_t \end{bmatrix} = \Phi_t^T \mathbf{w} + \epsilon_y, \quad p(\tau|\mathbf{w}) = \prod_t \mathcal{N}(\mathbf{y}_t | \Phi_t^T \mathbf{w}, \Sigma_y), \quad (1)$$

where $\Phi_t = [\phi_t, \dot{\phi}_t]$ defines the $n \times 2$ dimensional time-dependent basis matrix for the joint positions q_t and velocities \dot{q}_t , n defines the number of basis functions and $\epsilon_y \sim \mathcal{N}(\mathbf{0}, \Sigma_y)$ is zero-mean i.i.d. Gaussian noise. By weighing the basis functions Φ_t with the parameter vector \mathbf{w} , we can represent the mean of a trajectory.

In order to capture the variance of the trajectories, we introduce a distribution $p(\mathbf{w}; \theta)$ over the weight vector \mathbf{w} , with parameters θ . The trajectory distribution $p(\tau; \theta)$ can now be computed by marginalizing out the weight vector \mathbf{w} , i.e., $p(\tau; \theta) = \int p(\tau|\mathbf{w})p(\mathbf{w}; \theta)d\mathbf{w}$. The distribution $p(\tau; \theta)$ defines a Hierarchical Bayesian Model (HBM) whose parameters are given by the observation noise variance Σ_y and the parameters θ of $p(\mathbf{w}; \theta)$.

Temporal Modulation. Temporal modulation is needed for a faster or slower execution of the movement. We introduce a phase variable z to decouple the movement from the time signal as for previous non-probabilistic approaches [18]. The phase can be any function monotonically increasing with time $z(t)$. By modifying the rate of the phase variable, we can modulate the speed of the movement. Without loss of generality, we define the phase as $z_0 = 0$ at the beginning of the movement and as $z_T = 1$ at the end. The basis functions ϕ_t now directly depend on the phase instead of time, such that $\phi_t = \phi(z_t)$ and the corresponding derivative becomes $\dot{\phi}_t = \phi'(z_t)\dot{z}_t$.

Rhythmic and Stroke-Based Movements. The choice of the basis functions depends on the type of movement, which can be either rhythmic or stroke-based. For stroke-based movements, we use Gaussian basis functions b_i^G , while for rhythmic movements we use Von-Mises basis functions b_i^{VM} to model periodicity in the phase variable z , i.e.,

$$b_i^G(z) = \exp\left(-\frac{(z_t - c_i)^2}{2h}\right), \quad b_i^{\text{VM}}(z) = \exp\left(\frac{\cos(2\pi(z_t - c_i))}{h}\right), \quad (2)$$

where h defines the width of the basis and c_i the center for the i th basis function. We normalize the basis functions with $\phi_i(z_t) = b_i(z)/\sum_j b_j(z)$.

Encoding Coupling between Joints. So far, we have considered each degree of freedom to be modeled independently. However, for many tasks we have to coordinate the movement of the joints. A common way to implement such coordination is via the phase variable z_t that couples the mean of the trajectory distribution [18]. Yet, it is often desirable to also encode higher-order moments of the coupling, such as the covariance of the joints at time point t . Hence, we extend our model to multiple dimensions. For each dimension i , we maintain a parameter vector \mathbf{w}_i , and we define the combined, weight vector \mathbf{w} as $\mathbf{w} = [\mathbf{w}_1^T, \dots, \mathbf{w}_n^T]^T$. The basis matrix Φ_t now extends to a block-diagonal matrix containing the basis functions and their derivatives for each dimension. The observation vector \mathbf{y}_t consists of the angles and velocities of all joints. The probability of an observation \mathbf{y} at time t is given by

$$p(\mathbf{y}_t|\mathbf{w}) = \mathcal{N}\left(\begin{bmatrix} \mathbf{y}_{1,t} \\ \vdots \\ \mathbf{y}_{d,t} \end{bmatrix} \middle| \begin{bmatrix} \Phi_t^T & \dots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ \mathbf{0} & \dots & \Phi_t^T \end{bmatrix} \mathbf{w}, \Sigma_y\right) = \mathcal{N}(\mathbf{y}_t | \Psi_t \mathbf{w}, \Sigma_y) \quad (3)$$

where $\mathbf{y}_{i,t} = [q_{i,t}, \dot{q}_{i,t}]^T$ denotes the joint angle and velocity for the i^{th} joint. We now maintain a distribution $p(\mathbf{w}; \theta)$ over the combined parameter vector \mathbf{w} . Using this distribution, we can also capture the covariance between joints.

Learning from Demonstrations. One crucial requirement of a MP representation is that the parameters of a single primitive are easy to acquire from demonstrations. To facilitate the estimation of the parameters, we will assume a Gaussian distribution for $p(\mathbf{w}; \boldsymbol{\theta}) = \mathcal{N}(\mathbf{w} | \boldsymbol{\mu}_w, \boldsymbol{\Sigma}_w)$ over the parameters \mathbf{w} . Consequently, the distribution of the state $p(\mathbf{y}_t | \boldsymbol{\theta})$ for time step t is given by

$$p(\mathbf{y}_t; \boldsymbol{\theta}) = \int \mathcal{N}(\mathbf{y}_t | \boldsymbol{\Psi}_t^T \mathbf{w}, \boldsymbol{\Sigma}_y) \mathcal{N}(\mathbf{w} | \boldsymbol{\mu}_w, \boldsymbol{\Sigma}_w) d\mathbf{w} = \mathcal{N}(\mathbf{y}_t | \boldsymbol{\Psi}_t^T \boldsymbol{\mu}_w, \boldsymbol{\Psi}_t^T \boldsymbol{\Sigma}_w \boldsymbol{\Psi}_t + \boldsymbol{\Sigma}_y), \quad (4)$$

and, thus, we can easily evaluate the mean and the variance for any time point t . As a ProMP represents multiple ways to execute an elemental movement, we also need multiple demonstrations to learn $p(\mathbf{w}; \boldsymbol{\theta})$. The parameters $\boldsymbol{\theta} = \{\boldsymbol{\mu}_w, \boldsymbol{\Sigma}_w\}$ can be learned from multiple demonstrations by maximum likelihood estimation, for example, by using the expectation maximization algorithm for HBMs with Gaussian distributions [19].

2.2 New Probabilistic Operators for Movement Primitives

The ProMPs allow for the formulation of new operators from probability theory, e.g., conditioning for modulating the trajectory and a product of distributions for co-activating MPs.

Modulation of Via-Points, Final Positions or Velocities by Conditioning. The modulation of via-points and final positions are important properties of any MP framework such that the MP can be adapted to new situations. In our probabilistic formulation, such operations can be described by conditioning the MP to reach a certain state \mathbf{y}_t^* at time t . Conditioning is performed by adding a desired observation $\mathbf{x}_t = [\mathbf{y}_t^*, \boldsymbol{\Sigma}_y^*]$ to our probabilistic model and applying Bayes theorem, i.e., $p(\mathbf{w} | \mathbf{x}_t^*) \propto \mathcal{N}(\mathbf{y}_t^* | \boldsymbol{\Psi}_t^T \mathbf{w}, \boldsymbol{\Sigma}_y^*) p(\mathbf{w})$. The state vector \mathbf{y}_t^* represents the desired position and velocity vector at time t and $\boldsymbol{\Sigma}_y^*$ describes the accuracy of the desired observation. We can also condition on any subset of \mathbf{y}_t^* . For example, by specifying a desired joint position q_1 for the first joint the trajectory distribution will automatically infer the most probable joint positions for the other joints.

For Gaussian trajectory distributions the conditional distribution $p(\mathbf{w} | \mathbf{x}_t^*)$ for \mathbf{w} is Gaussian with mean and variance

$$\boldsymbol{\mu}_w^{[\text{new}]} = \boldsymbol{\mu}_w + \boldsymbol{\Sigma}_w \boldsymbol{\Psi}_t \left(\boldsymbol{\Sigma}_y^* + \boldsymbol{\Psi}_t^T \boldsymbol{\Sigma}_w \boldsymbol{\Psi}_t \right)^{-1} \left(\mathbf{y}_t^* - \boldsymbol{\Psi}_t^T \boldsymbol{\mu}_w \right), \quad (5)$$

$$\boldsymbol{\Sigma}_w^{[\text{new}]} = \boldsymbol{\Sigma}_w - \boldsymbol{\Sigma}_w \boldsymbol{\Psi}_t \left(\boldsymbol{\Sigma}_y^* + \boldsymbol{\Psi}_t^T \boldsymbol{\Sigma}_w \boldsymbol{\Psi}_t \right)^{-1} \boldsymbol{\Psi}_t^T \boldsymbol{\Sigma}_w. \quad (6)$$

Conditioning a ProMP to different target states is also illustrated in Figure 1(a). We can see that, despite the modulation of the ProMP by conditioning, the ProMP stays within the original distribution, and, hence, the modulation is also learned from the original demonstrations. Modulation strategies in current approaches such as the DMPs do not show this beneficial effect [18].

Combination and Blending of Movement Primitives. Another beneficial probabilistic operation is to continuously combine and blend different MPs into a single movement. Suppose that we maintain a set of i different primitives that we want to combine. We can co-activate them by taking the products of distributions, i.e., $p_{\text{new}}(\boldsymbol{\tau}) \propto \prod_i p_i(\boldsymbol{\tau})^{\alpha^{[i]}}$ where the $\alpha^{[i]} \in [0, 1]$ factors denote the activation of the i^{th} primitive. This product captures the overlapping region of the active MPs, i.e., the part of the trajectory space where all MPs have high probability mass.

However, we also want to be able to modulate the activations of the primitives, for example, to continuously blend the movement execution from one primitive to the next. Hence, we decompose the trajectory into single time steps and use time-varying activation functions $\alpha_t^{[i]}$, i.e.,

$$p^*(\boldsymbol{\tau}) \propto \prod_t \prod_i p_i(\mathbf{y}_t)^{\alpha_t^{[i]}}, \quad p_i(\mathbf{y}_t) = \int p_i(\mathbf{y}_t | \mathbf{w}^{[i]}) p_i(\mathbf{w}^{[i]}) d\mathbf{w}^{[i]}. \quad (7)$$

For Gaussian distributions $p_i(\mathbf{y}_t) = \mathcal{N}(\mathbf{y}_t | \boldsymbol{\mu}_t^{[i]}, \boldsymbol{\Sigma}_t^{[i]})$, the resulting distribution $p^*(\mathbf{y}_t)$ is again Gaussian with variance and mean

$$\boldsymbol{\Sigma}_t^* = \left(\sum_i \left(\boldsymbol{\Sigma}_t^{[i]} / \alpha_t^{[i]} \right)^{-1} \right)^{-1}, \quad \boldsymbol{\mu}_t^* = (\boldsymbol{\Sigma}_t^*)^{-1} \left(\sum_i \left(\boldsymbol{\Sigma}_t^{[i]} / \alpha_t^{[i]} \right)^{-1} \boldsymbol{\mu}_t^{[i]} \right) \quad (8)$$

Both terms, and their derivatives, are required to obtain the stochastic feedback controller which is finally used to control the robot. We illustrated the co-activation of two ProMPs in Figure 1(b) and the blending of two ProMPs in Figure 1(c).

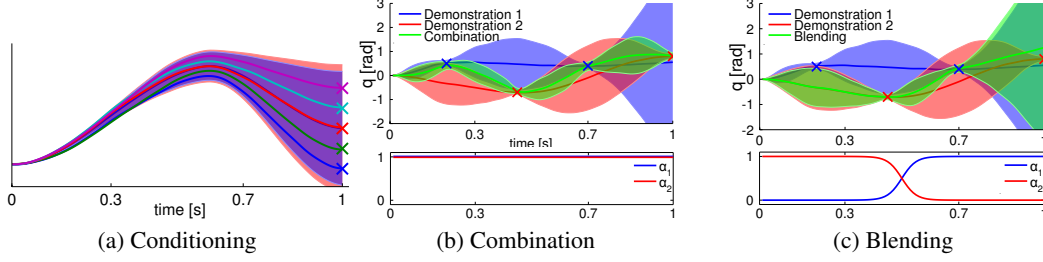


Figure 1: (a) *Conditioning on different target states.* The blue shaded area represents the learned trajectory distribution. We condition on different target positions, indicated by the ‘x’-markers. The produced trajectories exactly reach the desired targets while keeping the shape of the demonstrations. (b) *Combination of two ProMPs.* The trajectory distributions are indicated by the blue and red shaded areas. Both primitives have to reach via-points at different points in time, indicated by the ‘x’-markers. We co-activate both primitives with the same activation factor. The trajectory distribution generated by the resulting feedback controller now goes through all four via-points. (c) *Blending of two ProMPs.* We smoothly blend from the red primitive to the blue primitive. The activation factors are shown in the bottom. The resulting movement (green) first follows the red primitive and, subsequently, switches to following the blue primitive.

2.3 Using Trajectory Distributions for Robot Control

In order to fully exploit the properties of trajectory distributions, a policy for controlling the robot is needed that reproduces these distributions. To this effect, we analytically derivate a stochastic feedback controller that can accurately reproduce the mean vectors μ_t and the variances Σ_t for all t of a given trajectory distribution.

We follow a model-based approach. First, we approximate the continuous time dynamics of the system by a linearized discrete-time system with step duration dt ,

$$\mathbf{y}_{t+dt} = (\mathbf{I} + \mathbf{A}_t dt) \mathbf{y}_t + \mathbf{B}_t dt \mathbf{u} + \mathbf{c}_t dt, \quad (9)$$

where the system matrices \mathbf{A}_t , the input matrices \mathbf{B}_t and the drift vectors \mathbf{c}_t can be obtained by first order Taylor expansion of the dynamical system¹. We assume a stochastic linear feedback controller with time varying feedback gains is generating the control actions, i.e.,

$$\mathbf{u} = \mathbf{K}_t \mathbf{y}_t + \mathbf{k}_t + \epsilon_u, \quad \epsilon \sim \mathcal{N}(\epsilon_u | 0, \Sigma_u / dt), \quad (10)$$

where the matrix \mathbf{K}_t denotes a feedback gain matrix and \mathbf{k}_t a feed-forward component. We use a control noise which behaves like a Wiener process [21], and, hence, its variance grows linearly with the step duration² dt . By substituting Eq. (10) into Eq. (9), we rewrite the next state of the system as

$$\begin{aligned} \mathbf{y}_{t+dt} &= (\mathbf{I} + (\mathbf{A}_t + \mathbf{B}_t \mathbf{K}_t) dt) \mathbf{y}_t + \mathbf{B}_t dt (\mathbf{k}_t + \epsilon_u) + \mathbf{c}_t dt = \mathbf{F}_t \mathbf{y}_t + \mathbf{f}_t + \mathbf{B}_t dt \epsilon_u, \\ \text{with } \mathbf{F}_t &= (\mathbf{I} + (\mathbf{A}_t + \mathbf{B}_t \mathbf{K}_t) dt), \quad \mathbf{f}_t = \mathbf{B}_t \mathbf{k}_t dt + \mathbf{c}_t dt. \end{aligned} \quad (11)$$

For improved clarity, we will omit the time-index as subscript for most matrices in the remainder of the paper. From Eq. 4 we know that the distribution for our current state \mathbf{y}_t is Gaussian with mean $\mu_t = \Psi_t^T \mu_w$ and covariance³ $\Sigma_t = \Psi_t^T \Sigma_w \Psi_t$. As the system dynamics are modeled by a Gaussian linear model, we can obtain the distribution of the next state $p(\mathbf{y}_{t+dt})$ analytically from the forward model

$$\begin{aligned} p(\mathbf{y}_{t+dt}) &= \int \mathcal{N}(\mathbf{y}_{t+dt} | \mathbf{F} \mathbf{y}_t + \mathbf{f}, \Sigma_s dt) \mathcal{N}(\mathbf{y}_t | \mu_t, \Sigma_t) d\mathbf{y}_t \\ &= \mathcal{N}(\mathbf{y}_{t+dt} | \mathbf{F} \mu_t + \mathbf{f}, \mathbf{F} \Sigma_t \mathbf{F}^T + \Sigma_s dt), \end{aligned} \quad (12)$$

¹If inverse dynamics control [20] is used for the robot, the system reduces to a linear system where the terms \mathbf{A}_t , \mathbf{B}_t and \mathbf{c}_t are constant in time.

²As we multiply the noise by $\mathbf{B}dt$, we need to divide the covariance Σ_u of the control noise ϵ_u by dt to obtain this desired behavior.

³The observation noise is omitted as it represents independent noise which is not used for predicting the next state.

where $\text{dt}\Sigma_s = \text{dt}B\Sigma_u B^T$ represents the system noise matrix. Both sides of Eq. 12 are Gaussian distributions, where the left-hand side can also be computed by our desired trajectory distribution $p(\tau; \theta)$. We match the mean and the variances of both sides with our control law, i.e.,

$$\mu_{t+\text{dt}} = F\mu_t + (Bk + c)\text{dt}, \quad \Sigma_{t+\text{dt}} = F\Sigma_t F^T + \Sigma_s \text{dt}, \quad (13)$$

where F is given in Eq. (11) and contains the time varying feedback gains K . Using both constraints, we can now obtain the time dependent gains K and k .

Derivation of the Controller Gains. By rearranging terms, the covariance constraint becomes

$$\Sigma_{t+\text{dt}} - \Sigma_t = \Sigma_s \text{dt} + (A + BK)\Sigma_t \text{dt} + \Sigma_t(A + BK)^T \text{dt} + O(\text{dt}^2), \quad (14)$$

where $O(\text{dt}^2)$ denotes all second order terms in dt . After dividing by dt and taking the limit of $\text{dt} \rightarrow 0$, the second order terms disappear and we obtain the time derivative of the covariance

$$\dot{\Sigma}_t = \lim_{\text{dt} \rightarrow 0} \frac{\Sigma_{t+\text{dt}} - \Sigma_t}{\text{dt}} = (A + BK)\Sigma_t + \Sigma_t(A + BK)^T + \Sigma_s. \quad (15)$$

The matrix $\dot{\Sigma}_t$ can also be obtained from the trajectory distribution $\dot{\Sigma}_t = \dot{\Psi}_t^T \Sigma_w \Psi_t + \Psi_t^T \Sigma_w \dot{\Psi}_t$, which we substitute into Eq. (15). After rearranging terms, the equation reads

$$M + M^T = BK\Sigma_t + (BK\Sigma_t)^T, \text{ with } M = \dot{\Psi}_t^T \Sigma_w \Phi_t^T - A\Sigma_t - \Sigma_s/2. \quad (16)$$

Setting $M = BK\Sigma_t$ and solving for the gain matrix K

$$K = B^\dagger \left(\dot{\Psi}_t^T \Sigma_w \Psi_t - A\Sigma_t - \Sigma_s/2 \right) \Sigma_t^{-1}, \quad (17)$$

yields the solution, where B^\dagger denotes the pseudo-inverse of the control matrix B .

Derivation of the Feed-Forward Controls. Similarly, we obtain the feed-forward control signal k by matching the mean of the trajectory distribution $\mu_{t+\text{dt}}$ with the mean computed with the forward model. After rearranging terms, dividing by dt and taking the limit of $\text{dt} \rightarrow 0$, we arrive at the continuous time constraint for the vector k ,

$$\dot{\mu}_t = (A + BK)\mu_t + Bk + c. \quad (18)$$

We can again use the trajectory distribution $p(\tau; \theta)$ to obtain $\mu_t = \Psi_t \mu_w$ and $\dot{\mu}_t = \dot{\Psi}_t \mu_w$ and solve Eq. (18) for k ,

$$k = B^\dagger \left(\dot{\Psi}_t \mu_w - (A + BK) \Psi_t \mu_w - c \right) \quad (19)$$

Estimation of the Control Noise. In order to match a trajectory distribution, we also need to match the control noise matrix Σ_u which has been applied to generate the distribution. We first compute the system noise covariance $\Sigma_s = B\Sigma_u B^T$ by examining the cross-correlation between time steps of the trajectory distribution. To do so, we compute the joint distribution $p(y_t, y_{t+\text{dt}})$ of the current state y_t and the next state $y_{t+\text{dt}}$,

$$p(y_t, y_{t+\text{dt}}) = \mathcal{N} \left(\begin{bmatrix} y_t \\ y_{t+\text{dt}} \end{bmatrix} \middle| \begin{bmatrix} \mu_t \\ \mu_{t+\text{dt}} \end{bmatrix}, \begin{bmatrix} \Sigma_t & C_t \\ C_t^T & \Sigma_{t+\text{dt}} \end{bmatrix} \right), \quad (20)$$

where $C_t = \Psi_t \Sigma_w \Psi_{t+\text{dt}}^T$ is the cross-correlation. We can again use our model to match the cross correlation. The joint distribution for y_t and $y_{t+\text{dt}}$ is obtained by our system dynamics by $p(y_t, y_{t+\text{dt}}) = \mathcal{N}(y_t | \mu_t, \Sigma_t) \mathcal{N}(y_{t+\text{dt}} | Fy_t + f, \Sigma_u)$ which yields

$$p(y_t, y_{t+\text{dt}}) = \mathcal{N} \left(\begin{bmatrix} y_t \\ y_{t+\text{dt}} \end{bmatrix} \middle| \begin{bmatrix} \mu_t \\ F\mu_t + f \end{bmatrix}, \begin{bmatrix} \Sigma_t & \Sigma_t F^T \\ F\Sigma_t & F\Sigma_t F^T + \Sigma_s \text{dt} \end{bmatrix} \right). \quad (21)$$

The noise covariance Σ_s can be obtained by matching both covariance matrices given in Eq. (20) and (21),

$$\Sigma_s \text{dt} = \Sigma_{t+\text{dt}} - F\Sigma_t F^T = \Sigma_{t+\text{dt}} - F\Sigma_t \Sigma_t^{-1} \Sigma_t F^T = \Sigma_{t+\text{dt}} - C_t^T \Sigma_t^{-1} C_t \quad (22)$$

The variance Σ_u of the control noise is then given by $\Sigma_u = B^\dagger \Sigma_s B^{\dagger T}$. As we can see from Eq. (22) the variance of our stochastic feedback controller does not depend on the controller gains and can be pre-computed before estimating the controller gains.

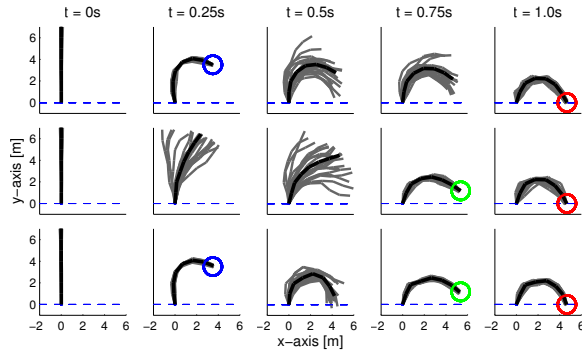


Figure 2: A 7-link planar robot has to reach a target position at $T = 1.0s$ with its end-effector while passing a via-point at $t_1 = 0.25s$ (top) or $t_2 = 0.75s$ (middle). The plot shows the mean posture of the robot at different time steps in black and samples generated by the ProMP in gray. The ProMP approach was able to exactly reproduce the demonstration which have been generated by an optimal control law. The combination of both learned ProMPs is shown in the bottom. The resulting movement reached both via-points with high accuracy.

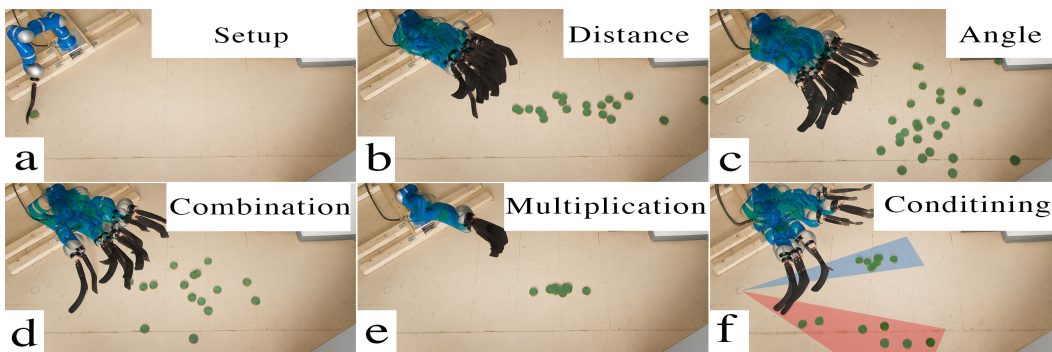


Figure 3: Robot Hockey. The robot shoots a hockey puck. We demonstrate ten straight shots for varying distances and ten shots for varying angles. The pictures show samples from the ProMP model for straight shots (b) and angled shots (c). Learning from combined data set yields a model that represents variance in both, distance and angle (d). Multiplying the individual models leads to a model that only reproduces shots where both models had probability mass, in the center at medium distance (e). The last picture shows the effect of conditioning on only left and right angles (f).

3 Experiments

We evaluated our approach on two different real robot tasks, one stroke based movement and one rhythmic movements. Additionally, we illustrate our approach on a 7-link simulated planar robot. For all real robot experiments we use a seven degrees of freedom KUKA lightweight robot arm. A more detailed description of the experiments is given in the supplementary material.

7-link Reaching Task. In this task, a seven link planar robot has to reach a target position in end-effector space. While doing so, it also has to reach a via-point at a certain time point. We generated the demonstrations for learning the MPs with an optimal control law [22]. In the first set of demonstrations, the robot has to reach the via-point at $t_1 = 0.25s$. The reproduced behavior with the ProMPs is illustrated in Figure 2(top). We learned the coupling of all seven joints with one ProMP. The ProMP exactly reproduced the via-points in task space while exhibiting a large variability in between the time points of the via-points. Moreover, the ProMP could also reproduce the coupling of the joints from the optimal control law which can be seen by the small variance of the end-effector in comparison to the rather large variance of the single joints at the via-points. The ProMP could achieve an average cost value of a similar quality as the optimal controller. We also used a second set of demonstrations where the first via-point was located at time step $t_2 = 0.75$, which is illustrated in Figure 2(middle). We combined the ProMPs learned from both demonstrations, which resulted in the movement illustrated in Figure 2(bottom). The combination of both MPs accurately reaches both via-points at $t_1 = 0.25$ and $t_2 = 0.75$.

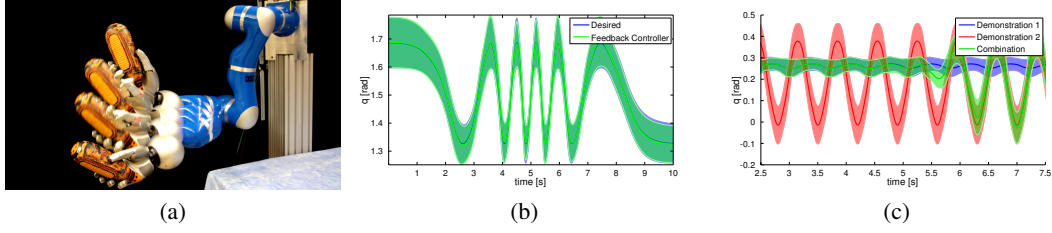


Figure 4: (a) The maracas task. We demonstrate a slow version of the rhythmic shaking movement required to play the maracas. After learning the ProMP model for the demonstration, we progressively increase the speed of the movement by modulating the phase, such that the robot successfully plays the instrument. (b) Trajectory distribution for playing maracas (joint number 4). By modulating the speed of the phase signal z_t , the speed of the movement can be adapted. The plot shows the desired distribution in blue and the generated distribution from the feedback controller in green. Both distributions match. (c) Blending between two rhythmic movements (blue and red shaded areas) for playing maracas. The green shaded is produced by continuously switching from the blue to the red movement.

Robot Hockey. In the hockey task, the robot has to shoot a hockey puck in different directions and distances. The task setup can be seen in Figure 3(a). We record two different sets of demonstrations, one that contains straight shots with varying distances while the second set contains shots with a varying shooting angle. Both data sets contain ten demonstrations each. Sampling from the two models generated by the different data sets yields shots that exhibit the demonstrated variance in either angle or distance, as shown in Figure 3(b) and 3(c). When combining the two individual primitives, the resulting model shoots only in the center at medium distance, i.e., the intersection of both MPs. We also learn a joint distribution over the final puck position and the weight vectors w and condition on the angle of the shot. The conditioning yields a model that shoots in different directions, depending on the conditioning, see Figure 3(f).

Robot Maracas. A maracas is a musical instrument containing grains, such that shaking it produces sounds. Demonstrating fast movements can be difficult on the robot arm, due to the inertia of the arm. Instead, we demonstrate a slower movement of ten periods to learn the motion. Since the ProMPs can change the phase of the motion, we use this slow demonstration and change the phase after learning the model to achieve a shaking movement of appropriate speed to generate the desired sound of the instrument. Using a variable phase also allows us to change the speed of the motion during one execution to achieve different sound patterns. We show an example movement of the robot in Figure 4(a). The desired trajectory distribution of the rhythmic movement and the resulting distribution generated from the feedback controller are shown in Figure 4(b). Both distributions match. We also demonstrated a second type of rhythmic shaking movement which we use to continuously blend between both movements to produce different sounds. One such transition between the two ProMPs is shown for one joint in Figure 4(c).

4 Conclusion

Probabilistic movement primitives are a promising approach for learning, modulating, and re-using movements in a modular control architecture. To effectively take advantage of such a control architecture, ProMPs support simultaneous activation, match the quality of the encoded behavior from the demonstrations, are able to adapt to different desired target positions, and efficiently learn by imitation. We parametrize the desired trajectory distribution of the primitive by a Hierarchical Bayesian Model with Gaussian distributions. The trajectory distribution can be easily obtained from demonstrations. Learning from a single or a limited set of demonstrations is feasible given prior knowledge of the task in hand, i.e. given an appropriate prior for Σ_w . It simultaneously defines a feedback controller which is used for movement execution. Our probabilistic formulation allows for new operations for movement primitives, including conditioning and combination of primitives. Future work will focus on using the ProMPs in a modular control architecture and improving upon imitation learning by reinforcement learning.

References

- [1] A. Ijspeert and S. Schaal. Learning Attractor Landscapes for Learning Motor Primitives. In *Advances in Neural Information Processing Systems 15*, (NIPS), MIT Press, Cambridge, MA, 2003.
- [2] M. Khansari-Zadeh and A. Billard. Learning Stable Non-Linear Dynamical Systems with Gaussian Mixture Models. *IEEE Transaction on Robotics*, 2011.
- [3] Jens Kober, Katharina Mülling, Oliver Kroemer, Christoph H. Lampert, Bernhard Schölkopf, and Jan Peters. Movement Templates for Learning of Hitting and Batting. In *International Conference on Robotics and Automation (ICRA)*, 2010.
- [4] J. Kober and J. Peters. Policy Search for Motor Primitives in Robotics. *Machine Learning*, pages 1–33, 2010.
- [5] A. Ude, A. Gams, T. Asfour, and J. Morimoto. Task-Specific Generalization of Discrete and Periodic Dynamic Movement Primitives. *Trans. Rob.*, (5), October 2010.
- [6] B. da Silva, G. Konidaris, and A. Barto. Learning Parameterized Skills. In *International Conference on Machine Learning*, 2012.
- [7] P. Kormushev, S. Calinon, and D. Caldwell. Robot Motor Skill Coordination with EM-based Reinforcement Learning. In *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2010.
- [8] C. Daniel, G. Neumann, and J. Peters. Learning Concurrent Motor Skills in Versatile Solution Spaces. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2012.
- [9] George Konidaris, Scott Kuindersma, Roderic Grupen, and Andrew Barto. Robot Learning from Demonstration by Constructing Skill Trees. *International Journal of Robotics Research*, 31(3):360–375, March 2012.
- [10] A. d’Avella and E. Bizzi. Shared and Specific Muscle Synergies in Natural Motor Behaviors. *Proceedings of the National Academy of Sciences (PNAS)*, 102(3):3076–3081, 2005.
- [11] M. Williams, B. and Toussaint and A. Storkey. Modelling Motion Primitives and their Timing in Biologically Executed Movements. In *Advances in Neural Information Processing Systems (NIPS)*, 2007.
- [12] L. Roza, S. Calinon, D. G. Caldwell, P. Jimenez, and C. Torras. Learning Collaborative Impedance-Based Robot Behaviors. In *AAAI Conference on Artificial Intelligence*, 2013.
- [13] E. Rueckert, G. Neumann, M. Toussaint, and W. Pr Maass. Learned Graphical Models for Probabilistic Planning provide a new Class of Movement Primitives. 2012.
- [14] L. Righetti and A. Ijspeert. Programmable central pattern generators: an application to biped locomotion control. In *Proceedings of the 2006 IEEE International Conference on Robotics and Automation*, 2006.
- [15] A. Paraschos, G. Neumann, and J. Peters. A probabilistic approach to robot trajectory generation. In *Proceedings of the International Conference on Humanoid Robots (HUMANOIDS)*, 2013.
- [16] S. Calinon, P. Kormushev, and D. G. Caldwell. Compliant skills acquisition and multi-optima policy search with em-based reinforcement learning. *Robotics and Autonomous Systems*, 61(4):369–379, April 2013.
- [17] E. Todorov and M. Jordan. Optimal Feedback Control as a Theory of Motor Coordination. *Nature Neuroscience*, 5:1226–1235, 2002.
- [18] S. Schaal, J. Peters, J. Nakanishi, and A. Ijspeert. Learning Movement Primitives. In *International Symposium on Robotics Research*, (ISRR), 2003.
- [19] A. Lazaric and M. Ghavamzadeh. Bayesian Multi-Task Reinforcement Learning. In *Proceedings of the 27th International Conference on Machine Learning (ICML)*, 2010.
- [20] J. Peters, M. Mistry, F. E. Udawadia, J. Nakanishi, and S. Schaal. A Unifying Methodology for Robot Control with Redundant DOFs. *Autonomous Robots*, (1):1–12, 2008.
- [21] H. Stark and J. Woods. *Probability and Random Processes with Applications to Signal Processing (3rd Edition)*. 3 edition, August 2001.
- [22] M. Toussaint. Robot Trajectory Optimization using Approximate Inference. In *Proceedings of the 26th International Conference on Machine Learning*, (ICML), 2009.