# Composition of Movement Primitives

Andrea Pierré

June 20, 2025

# Contents

# 1 ProMPs

## 1.1 Recap

From (Paraschos et al., 2013, 2018):

Table 1: Notation

| | |
|---|---|
| $q_t$ | joint angle over time |
| $\dot{q}_t$ | joint velocity over time |
| $\boldsymbol{\tau} = \{q_t\}_{t=0\ldots T}$ | trajectory |
| $\boldsymbol{w}$ | weight vector of a single trajectory $[n \times 1]$ |
| $\phi_t$ | basis function |
| $n$ | number of basis functions |
| $\boldsymbol{\Phi}_t = [\phi_t, \dot{\phi}_t]$ | $n \times 2$ dimensional time-dependent basis matrix |
| $z(t)$ | monotonically increasing phase variable |
| $\boldsymbol{\epsilon}_y \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{\Sigma}_y)$ | zero-mean i.i.d. Gaussian noise |

$$\mathbf{\Phi}_t = \begin{bmatrix} \phi_1 & \dot{\phi}_1 \\ \vdots & \vdots \\ \phi_n & \dot{\phi}_n \end{bmatrix} \tag{1}$$

$$\mathbf{y}_t = \begin{bmatrix} q_t \\ \dot{q}_t \end{bmatrix} = \mathbf{\Phi}_t^\top \mathbf{w} + \boldsymbol{\epsilon}_y \tag{2}$$

$$p(\boldsymbol{\tau}|\mathbf{w}) = \prod_t \mathcal{N}\left(\mathbf{y}_t | \mathbf{\Phi}_t^\top \mathbf{w}, \mathbf{\Sigma}_y\right) \tag{3}$$

$$p(\boldsymbol{\tau}; \boldsymbol{\theta}) = \int p(\boldsymbol{\tau}|\mathbf{w}) \cdot p(\mathbf{w}; \boldsymbol{\theta}) d\mathbf{w} \tag{4}$$

## 1.2 Coupling between joints

$$p(\mathbf{y}_t|\mathbf{w}) = \mathcal{N}\left( \begin{bmatrix} \mathbf{y}_{1,t} \\ \vdots \\ \mathbf{y}_{d,t} \end{bmatrix} \middle| \begin{bmatrix} \mathbf{\Phi}_t^\top & \cdots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ \mathbf{0} & \cdots & \mathbf{\Phi}_t^\top \end{bmatrix} \mathbf{w}, \mathbf{\Sigma}_y \right) = \mathcal{N}\left(\mathbf{y}_t | \mathbf{\Psi}_t \mathbf{w}, \mathbf{\Sigma}_y\right) \tag{5}$$

with:

Table 2: Notation

| | |
|---|---|
| $\mathbf{w} = [\mathbf{w}_1^\top, \ldots, \mathbf{w}_n^\top]^\top$ | combined weight vector $[n \times n]$ |
| $\mathbf{\Psi}_t$ | block-diagonal basis matrix containing the basis functions and their derivatives for each dimension |
| $\mathbf{y}_{i,t} = [q_{i,t}, \dot{q}_{i,t}]^\top$ | joint angle and velocity for the $i^{\text{th}}$ joint |

## 1.3 Hierarchical Bayesian Model

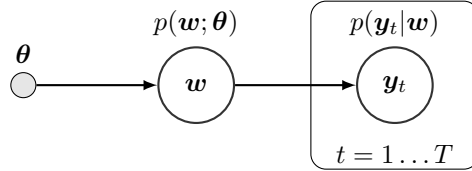The Hierarchical Bayesian Model used in ProMPs is illustrated in Fig. 1.



Figure 1: Hierarchical Bayesian Model used in ProMPs.

$$p(\mathbf{y}_t; \boldsymbol{\theta}) = \int \mathcal{N}\left(\mathbf{y}_t | \mathbf{\Psi}_t^\top \mathbf{w}, \mathbf{\Sigma}_y\right) \cdot p(\mathbf{w}; \boldsymbol{\theta}) \, d\mathbf{w} \tag{6}$$

$$= \int \mathcal{N}\left(\mathbf{y}_t | \mathbf{\Psi}_t^\top \mathbf{w}, \mathbf{\Sigma}_y\right) \cdot \mathcal{N}\left(\mathbf{w} | \boldsymbol{\mu}_w, \mathbf{\Sigma}_w\right) \, d\mathbf{w} \tag{7}$$

$$= \mathcal{N}\left(\mathbf{y}_t | \mathbf{\Psi}_t^\top \boldsymbol{\mu}_w, \mathbf{\Psi}_t^\top \mathbf{\Sigma}_w \mathbf{\Psi}_t + \mathbf{\Sigma}_y\right) \tag{8}$$

See Appendix A for the proof.

Table 3: Notation

| | |
|---|---|
| $\boldsymbol{\theta} = \{\boldsymbol{\mu}_w, \mathbf{\Sigma}_w\}$ | parameters |
| $p(\mathbf{w}; \boldsymbol{\theta}) = \mathcal{N}(\mathbf{w}|\boldsymbol{\mu}_w, \mathbf{\Sigma}_w)$ | prior over the weight vector $\mathbf{w}$, with parameters $\boldsymbol{\theta}$, assumed to be Gaussian |

Table 4: Notation

| | |
|---|---|
| $\boldsymbol{x}_t^\star = [\boldsymbol{y}_t^\star, \boldsymbol{\Sigma}_t^\star]$ | desired observation |
| $\boldsymbol{y}_t^\star$ | desired position and velocity vector at time $t$ |
| $\boldsymbol{\Sigma}_t^\star$ | accuracy of the desired observation |

## 1.4 Via-Points Modulation

Using Bayes rule:

$$p(\boldsymbol{w}|\boldsymbol{x}_t^\star) = \frac{p(\boldsymbol{x}_t^\star|\boldsymbol{w}) \cdot p(\boldsymbol{w})}{p(\boldsymbol{x}_t^\star)} \tag{9}$$

$$p(\boldsymbol{w}|\boldsymbol{x}_t^\star) \propto \mathcal{N}\left(\boldsymbol{y}_t^\star|\boldsymbol{\Psi}_t^\top \boldsymbol{w}, \boldsymbol{\Sigma}_t^\star\right) \cdot \mathcal{N}(\boldsymbol{w}|\boldsymbol{\mu}_w, \boldsymbol{\Sigma}_w) \tag{10}$$

$$\boldsymbol{\mu}_{\boldsymbol{w}}^{[new]} = \boldsymbol{\mu}_{\boldsymbol{w}} + \boldsymbol{\Sigma}_{\boldsymbol{w}}\boldsymbol{\Psi}_t\left(\boldsymbol{\Sigma}_y^\star + \boldsymbol{\Psi}_t^\top\boldsymbol{\Sigma}_{\boldsymbol{w}}\boldsymbol{\Psi}_t\right)^{-1}(\boldsymbol{y}_t^\star - \boldsymbol{\Psi}_t^\top\boldsymbol{\mu}_{\boldsymbol{w}}) \tag{11}$$

$$\boldsymbol{\Sigma}_{\boldsymbol{w}}^{[new]} = \boldsymbol{\Sigma}_{\boldsymbol{w}} - \boldsymbol{\Sigma}_{\boldsymbol{w}}\boldsymbol{\Psi}_t\left(\boldsymbol{\Sigma}_y^\star + \boldsymbol{\Psi}_t^\top\boldsymbol{\Sigma}_{\boldsymbol{w}}\boldsymbol{\Psi}_t\right)^{-1}\boldsymbol{\Psi}_t^\top\boldsymbol{\Sigma}_{\boldsymbol{w}} \tag{12}$$

See Appendix B for the proof.

### 1.4.1 Do we actually get the desired mean by applying the conditioning update?

*Proof that the posterior mean equals the observed mean.*

$$\mathbb{E}[\boldsymbol{y}_t|\boldsymbol{x}_t^\star] = \boldsymbol{\mu}_{\boldsymbol{y}_t|\boldsymbol{x}_t^\star} = \boldsymbol{\Psi}_t^\top\boldsymbol{\mu}_{\boldsymbol{w}|\boldsymbol{x}_t^\star} = \boldsymbol{\Psi}_t^\top\boldsymbol{\mu}_{\boldsymbol{w}} + \boldsymbol{\Psi}_t^\top\boldsymbol{\Sigma}_{\boldsymbol{w}}\boldsymbol{\Psi}_t\left(\boldsymbol{\Sigma}_t^\star + \boldsymbol{\Psi}_t^\top\boldsymbol{\Sigma}_{\boldsymbol{w}}\boldsymbol{\Psi}_t\right)^{-1}(\boldsymbol{y}_t^\star - \boldsymbol{\Psi}_t^\top\boldsymbol{\mu}_{\boldsymbol{w}}) \tag{13}$$

We set the observed covariance $\boldsymbol{\Sigma}_t^\star$ to 0 so as to have perfect accuracy around our observed position.

$$\boldsymbol{\Psi}_t^\top\boldsymbol{\mu}_{\boldsymbol{w}|\boldsymbol{x}_t^\star} = \boldsymbol{\Psi}_t^\top\boldsymbol{\mu}_{\boldsymbol{w}} + \cancel{\boldsymbol{\Psi}_t^\top\boldsymbol{\Sigma}_{\boldsymbol{w}}\boldsymbol{\Psi}_t}\left(\cancel{\boldsymbol{\Psi}_t^\top\boldsymbol{\Sigma}_{\boldsymbol{w}}\boldsymbol{\Psi}_t}\right)^{-1}(\boldsymbol{y}_t^\star - \boldsymbol{\Psi}_t^\top\boldsymbol{\mu}_{\boldsymbol{w}}) \tag{14}$$

$$= \cancel{\boldsymbol{\Psi}_t^\top\boldsymbol{\mu}_{\boldsymbol{w}}} + \boldsymbol{y}_t^\star - \cancel{\boldsymbol{\Psi}_t^\top\boldsymbol{\mu}_{\boldsymbol{w}}} \tag{15}$$

$$= \boldsymbol{y}_t^\star \tag{16}$$
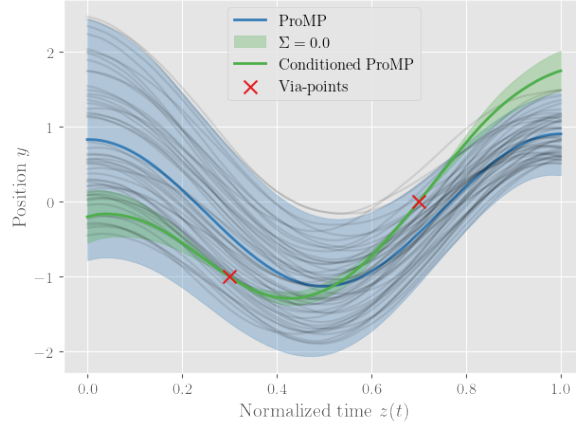
$\square$

### 1.4.2 Multi via-points



Figure 2: Example of ProMP with two via-points.

1. For the first via-point conditioning update with the observed via-point $\boldsymbol{x}_{t_1}^\star = [\boldsymbol{y}_{t_1}^\star, \boldsymbol{\Sigma}_{t_1}^\star]$, we can directly apply Eq. (11) and (12), with $\boldsymbol{\Psi}_{t_1}$ the observation matrix at time $t_1$:

$$\boldsymbol{\mu}_{\boldsymbol{w}|\boldsymbol{x}_{t_1}^\star} = \boldsymbol{\mu}_{\boldsymbol{w}} + \boldsymbol{\Sigma}_{\boldsymbol{w}}\boldsymbol{\Psi}_{t_1}\left(\boldsymbol{\Sigma}_{t_1}^\star + \boldsymbol{\Psi}_{t_1}^\top\boldsymbol{\Sigma}_{\boldsymbol{w}}\boldsymbol{\Psi}_{t_1}\right)^{-1}(\boldsymbol{y}_{t_1}^\star - \boldsymbol{\Psi}_{t_1}^\top\boldsymbol{\mu}_{\boldsymbol{w}}) \tag{17}$$

$$\boldsymbol{\Sigma}_{\boldsymbol{w}|\boldsymbol{x}_{t_1}^\star} = \boldsymbol{\Sigma}_{\boldsymbol{w}} - \boldsymbol{\Sigma}_{\boldsymbol{w}}\boldsymbol{\Psi}_{t_1}\left(\boldsymbol{\Sigma}_{t_1}^\star + \boldsymbol{\Psi}_{t_1}^\top\boldsymbol{\Sigma}_{\boldsymbol{w}}\boldsymbol{\Psi}_{t_1}\right)^{-1}\boldsymbol{\Psi}_{t_1}^\top\boldsymbol{\Sigma}_{\boldsymbol{w}} \tag{18}$$

2. For the second via-point conditioning update with the observed via-point $\boldsymbol{x}_{t_2}^\star = [\boldsymbol{y}_{t_2}^\star, \boldsymbol{\Sigma}_{t_2}^\star]$, the prior is the posterior from the first via-point, i.e., $\boldsymbol{w} \sim \mathcal{N}(\boldsymbol{\mu}_{\boldsymbol{w}|\boldsymbol{x}_{t_1}^\star}, \boldsymbol{\Sigma}_{\boldsymbol{w}|\boldsymbol{x}_{t_1}^\star})$, the likelihood is $\boldsymbol{y}_{t_2}^\star \sim \mathcal{N}(\boldsymbol{\Psi}_{t_2}^\top\boldsymbol{w}, \boldsymbol{\Sigma}_{t_2}^\star)$, with $\boldsymbol{\Psi}_{t_2}$ the observation matrix at time $t_2$, and the posterior update becomes:

$$\boldsymbol{\mu}_{\boldsymbol{w}|\boldsymbol{x}_{t_1}^\star, \boldsymbol{x}_{t_2}^\star} = \boldsymbol{\mu}_{\boldsymbol{w}|\boldsymbol{x}_{t_1}^\star} + \boldsymbol{\Sigma}_{\boldsymbol{w}|\boldsymbol{x}_{t_1}^\star}\boldsymbol{\Psi}_{t_2}\left(\boldsymbol{\Sigma}_{t_2}^\star + \boldsymbol{\Psi}_{t_2}^\top\boldsymbol{\Sigma}_{\boldsymbol{w}|\boldsymbol{x}_{t_1}^\star}\boldsymbol{\Psi}_{t_2}\right)^{-1}(\boldsymbol{y}_{t_2}^\star - \boldsymbol{\Psi}_{t_2}^\top\boldsymbol{\mu}_{\boldsymbol{w}|\boldsymbol{x}_{t_1}^\star}) \tag{19}$$

$$\boldsymbol{\Sigma}_{\boldsymbol{w}|\boldsymbol{x}_{t_1}^\star, \boldsymbol{x}_{t_2}^\star} = \boldsymbol{\Sigma}_{\boldsymbol{w}|\boldsymbol{x}_{t_1}^\star} - \boldsymbol{\Sigma}_{\boldsymbol{w}|\boldsymbol{x}_{t_1}^\star}\boldsymbol{\Psi}_{t_2}\left(\boldsymbol{\Sigma}_{t_2}^\star + \boldsymbol{\Psi}_{t_2}^\top\boldsymbol{\Sigma}_{\boldsymbol{w}|\boldsymbol{x}_{t_1}^\star}\boldsymbol{\Psi}_{t_2}\right)^{-1}\boldsymbol{\Psi}_{t_2}^\top\boldsymbol{\Sigma}_{\boldsymbol{w}|\boldsymbol{x}_{t_1}^\star} \tag{20}$$

3. For the $k^{\text{th}}$ via-point conditioning update with the observed via-point $\boldsymbol{x}_{t_k}^\star = [\boldsymbol{y}_{t_k}^\star, \boldsymbol{\Sigma}_{t_k}^\star]$, the prior is the posterior after conditioning on the previous $k-1$ via-points, i.e., $\boldsymbol{w} \sim \mathcal{N}(\boldsymbol{\mu}_{\boldsymbol{w}|\boldsymbol{x}_{t_1}^\star,...,\boldsymbol{x}_{t_{k-1}}^\star}, \boldsymbol{\Sigma}_{\boldsymbol{w}|\boldsymbol{x}_{t_1}^\star,...,\boldsymbol{x}_{t_{k-1}}^\star})$, the likelihood is $\boldsymbol{y}_{t_k}^\star \sim \mathcal{N}(\boldsymbol{\Psi}_{t_k}^\top\boldsymbol{w}, \boldsymbol{\Sigma}_{t_k}^\star)$, with $\boldsymbol{\Psi}_{t_k}$ the observation matrix at time $t_k$, and the posterior update becomes:

$$\boldsymbol{\mu}_{\boldsymbol{w}|\boldsymbol{x}_{t_1}^\star,...,\boldsymbol{x}_{t_k}^\star} = \boldsymbol{\mu}_{\boldsymbol{w}|\boldsymbol{x}_{t_1}^\star,...,\boldsymbol{x}_{t_{k-1}}^\star}$$
$$+ \boldsymbol{\Sigma}_{\boldsymbol{w}|\boldsymbol{x}_{t_1}^\star,...,\boldsymbol{x}_{t_{k-1}}^\star}\boldsymbol{\Psi}_{t_k}\left(\boldsymbol{\Sigma}_{t_k}^\star + \boldsymbol{\Psi}_{t_k}^\top\boldsymbol{\Sigma}_{\boldsymbol{w}|\boldsymbol{x}_{t_1}^\star,...,\boldsymbol{x}_{t_{k-1}}^\star}\boldsymbol{\Psi}_{t_k}\right)^{-1}(\boldsymbol{y}_{t_k}^\star - \boldsymbol{\Psi}_{t_k}^\top\boldsymbol{\mu}_{\boldsymbol{w}|\boldsymbol{x}_{t_1}^\star,...,\boldsymbol{x}_{t_{k-1}}^\star}) \tag{21}$$

$$\boldsymbol{\Sigma}_{\boldsymbol{w}|\boldsymbol{x}_{t_1}^\star,...,\boldsymbol{x}_{t_k}^\star} = \boldsymbol{\Sigma}_{\boldsymbol{w}|\boldsymbol{x}_{t_1}^\star,...,\boldsymbol{x}_{t_{k-1}}^\star} \tag{22}$$
$$- \boldsymbol{\Sigma}_{\boldsymbol{w}|\boldsymbol{x}_{t_1}^\star,...,\boldsymbol{x}_{t_{k-1}}^\star}\boldsymbol{\Psi}_{t_k}\left(\boldsymbol{\Sigma}_{t_k}^\star + \boldsymbol{\Psi}_{t_k}^\top\boldsymbol{\Sigma}_{\boldsymbol{w}|\boldsymbol{x}_{t_1}^\star,...,\boldsymbol{x}_{t_{k-1}}^\star}\boldsymbol{\Psi}_{t_k}\right)^{-1}\boldsymbol{\Psi}_{t_k}^\top\boldsymbol{\Sigma}_{\boldsymbol{w}|\boldsymbol{x}_{t_1}^\star,...,\boldsymbol{x}_{t_{k-1}}^\star}$$

**Alternative Batch Formulation** Instead of iterative updates, we could condition on all via-points simultaneously by stacking the observations:

$$\boldsymbol{y}^\star = \begin{bmatrix} \boldsymbol{y}_{t_1}^\star \\ \vdots \\ \boldsymbol{y}_{t_k}^\star \end{bmatrix}, \quad \boldsymbol{\Psi} = \begin{bmatrix} \boldsymbol{\Psi}_{t_1} \\ \vdots \\ \boldsymbol{\Psi}_{t_k} \end{bmatrix}, \quad \boldsymbol{\Sigma}^\star = \text{diag}(\boldsymbol{\Sigma}_{t_1}^\star, \ldots, \boldsymbol{\Sigma}_{t_k}^\star) \tag{23}$$

$$\boldsymbol{\mu}_{\boldsymbol{w}|\{\boldsymbol{x}_{t_k}^\star\}_{k=1}^K} = \boldsymbol{\mu}_{\boldsymbol{w}} + \boldsymbol{\Sigma}_{\boldsymbol{w}}\boldsymbol{\Psi}\left(\boldsymbol{\Sigma}^\star + \boldsymbol{\Psi}^\top\boldsymbol{\Sigma}_{\boldsymbol{w}}\boldsymbol{\Psi}\right)^{-1}(\boldsymbol{y}^\star - \boldsymbol{\Psi}^\top\boldsymbol{\mu}_{\boldsymbol{w}}) \tag{24}$$

$$\boldsymbol{\Sigma}_{\boldsymbol{w}|\{\boldsymbol{x}_{t_k}^\star\}_{k=1}^K} = \boldsymbol{\Sigma}_{\boldsymbol{w}} - \boldsymbol{\Sigma}_{\boldsymbol{w}}\boldsymbol{\Psi}\left(\boldsymbol{\Sigma}^\star + \boldsymbol{\Psi}^\top\boldsymbol{\Sigma}_{\boldsymbol{w}}\boldsymbol{\Psi}\right)^{-1}\boldsymbol{\Psi}^\top\boldsymbol{\Sigma}_{\boldsymbol{w}} \tag{25}$$

# 2 Gaussian mixture modeling (GMM)/Gaussian mixture regression (GMR) recap

## 2.1 Gaussian Mixture Modeling (GMM)

$$p(\boldsymbol{x}|\boldsymbol{\theta}) = \sum_{k=1}^K \pi_k \mathcal{N}(\boldsymbol{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \tag{26}$$

$$0 \le \pi_k \le 1, \quad \sum_{k=1}^K \pi_k = 1 \tag{27}$$

$$r_{nk} := \frac{\pi_k \mathcal{N}(\boldsymbol{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\boldsymbol{x}_n|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \tag{28}$$

| | |
|---|---|
| $\pi_k$ | mixture weights |
| $\boldsymbol{\theta} := \{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k : k = 1, \ldots, K\}$ | collection of all parameters of the model |
| $r_{nk}$ | responsibility of the $k^{\text{th}}$ mixture component for the $n^{\text{th}}$ data point |
| $N$ | number of data points |
| $N_k := \sum_{n=1}^{N} r_{nk}$ | total responsibility of the $k^{\text{th}}$ mixture component for the entire dataset |

Update of the GMM means:

$$\boldsymbol{\mu}_k^{new} = \frac{\sum_{n=1}^{N} r_{nk} \boldsymbol{x}_n}{\sum_{n=1}^{N} r_{nk}} \tag{29}$$

Update of the GMM covariances:

$$\boldsymbol{\Sigma}_k^{new} = \frac{1}{N_k} \sum_{n=1}^{N} r_{nk} (\boldsymbol{x}_n - \boldsymbol{\mu}_k)(\boldsymbol{x}_n - \boldsymbol{\mu}_k)^\top \tag{30}$$

Update of the GMM mixture weights:

$$\pi_k^{new} = \frac{N_k}{N}, \quad k = 1, \ldots, K \tag{31}$$

The parameters are updated using Algorithm 1.

## 2.2 Gaussian Mixture Regression (GMR)

At each iteration step $t$, the datapoint $\boldsymbol{x}_t$ can be decomposed as two subvectors $\boldsymbol{x}_t^I$ and $\boldsymbol{x}_t^O$ spanning for the input and output dimensions. For trajectory encoding in task space, $I$ corresponds to the time input dimension (*e.g.*, value of a decay term), and $O$ corresponds to the output dimensions describing a path (*e.g.*, end-effector position in task space).

During the training phase we learn the joint probability $p(\boldsymbol{x}_t^I, \boldsymbol{x}_t^O)$ with GMM through EM (Algorithm 1) with:

$$p(\boldsymbol{x}_t^I, \boldsymbol{x}_t^O) = \sum_{k=1}^{K} \pi_k \mathcal{N}_k(\boldsymbol{x}_t^I, \boldsymbol{x}_t^O | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \tag{32}$$

$$= \sum_{k=1}^{K} \pi_k p(\boldsymbol{x}_t^O | \boldsymbol{x}_t^I) \cdot p(\boldsymbol{x}_t^I) \tag{33}$$

$$= \sum_{k=1}^{K} \pi_k \mathcal{N}_k(\boldsymbol{x}_t^O | \hat{\boldsymbol{\mu}}_k^O(\boldsymbol{x}_t^I), \hat{\boldsymbol{\Sigma}}_k^O) \cdot \mathcal{N}_k(\boldsymbol{\mu}_k^I, \boldsymbol{\Sigma}_k^I) \tag{34}$$

$$\boldsymbol{x}_t = \begin{bmatrix} \boldsymbol{x}_t^I \\ \boldsymbol{x}_t^O \end{bmatrix}, \quad \boldsymbol{\mu}_k = \begin{bmatrix} \boldsymbol{\mu}_k^I \\ \boldsymbol{\mu}_k^O \end{bmatrix}, \quad \boldsymbol{\Sigma}_k = \begin{bmatrix} \boldsymbol{\Sigma}_k^I & \boldsymbol{\Sigma}_k^{IO} \\ \boldsymbol{\Sigma}_k^{OI} & \boldsymbol{\Sigma}_k^O \end{bmatrix} \tag{35}$$

$$\hat{\boldsymbol{\mu}}_k^O(\boldsymbol{x}_t^I) = \boldsymbol{\mu}_k^O + \boldsymbol{\Sigma}_k^{OI}(\boldsymbol{\Sigma}_k^I)^{-1}(\boldsymbol{x}_t^I - \boldsymbol{\mu}_k^I) \tag{36}$$

$$\hat{\boldsymbol{\Sigma}}_k^O = \boldsymbol{\Sigma}_k^O - \boldsymbol{\Sigma}_k^{OI}(\boldsymbol{\Sigma}_k^I)^{-1} \cdot \boldsymbol{\Sigma}_k^{IO} \tag{37}$$

The marginal probability $p(\boldsymbol{x}_t^I)$ is:

$$p(\boldsymbol{x}_t^I) = \int p(\boldsymbol{x}_t^I, \boldsymbol{x}_t^O) d\boldsymbol{x}_t^O = \sum_{k=1}^{K} \pi_k \cdot \mathcal{N}_k(\boldsymbol{x}_t^I | \boldsymbol{\mu}_k^I, \boldsymbol{\Sigma}_k^I) \tag{38}$$

**Algorithm 1:** EXPECTATION MAXIMIZATION (EM) algorithm for a Gaussian mixture model

**Input** : Initial model parameters $\{\boldsymbol{\mu}_k\}$, $\{\boldsymbol{\Sigma}_k\}$, $\{\pi_k\}$
**Input** : Data set $\{\mathbf{x}_1, \ldots, \mathbf{x}_N\}$
**Output:** Final model parameters $\{\boldsymbol{\mu}_k\}$, $\{\boldsymbol{\Sigma}_k\}$, $\{\pi_k\}$

**repeat**

    // E-step
    **for** $n \in \{1, \ldots, N\}$ **do**
        **for** $k \in \{1, \ldots, K\}$ **do**

$$r_{nk} \leftarrow \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^{K} \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$

        **end**
    **end**

    // M-step
    **for** $k \in \{1, \ldots, K\}$ **do**

$$N_k = \sum_{n=1}^{N} r_{nk}$$

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^{N} r_{nk} \mathbf{x}_n$$

$$\boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^{N} r_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^{\top}$$

$$\pi_k = \frac{N_k}{N}$$

    **end**

    // Log likelihood

$$\mathcal{L} \leftarrow \sum_{n=1}^{N} \ln \left[ \sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right]$$

**until** *convergence*
**return** $\{\boldsymbol{\mu}_k\}$, $\{\boldsymbol{\Sigma}_k\}$, $\{\pi_k\}$

Table 6: Notation

| | |
|---|---|
| $\boldsymbol{x}_t \in \mathbb{R}^D$ | datapoint at timestep $t$ |
| $\boldsymbol{\mu}_k$ | center of the $k^{\text{th}}$ Gaussian |
| $\boldsymbol{\Sigma}_k$ | covariance of the $k^{\text{th}}$ Gaussian |

To retrieve the multimodal conditional distribution $p(\boldsymbol{x}_t^O|\boldsymbol{x}_t^I)$ for each input/timestep we have:

$$p(\boldsymbol{x}_t^O|\boldsymbol{x}_t^I) = \frac{p(\boldsymbol{x}_t^O, \boldsymbol{x}_t^I)}{p(\boldsymbol{x}_t^I)} \tag{39}$$

$$= \frac{\sum_{k=1}^K \pi_k \mathcal{N}_k(\hat{\boldsymbol{\mu}}_k^O(\boldsymbol{x}_t^I), \hat{\boldsymbol{\Sigma}}_k^O) \cdot \mathcal{N}_k(\boldsymbol{x}_t^I|\boldsymbol{\mu}_k^I, \boldsymbol{\Sigma}_k^I)}{\sum_{k=1}^K \pi_k \cdot \mathcal{N}_k(\boldsymbol{x}_t^I|\boldsymbol{\mu}_k^I, \boldsymbol{\Sigma}_k^I)} \tag{40}$$

$$= \sum_{k=1}^K r_{nk} \cdot \mathcal{N}_k(\hat{\boldsymbol{\mu}}_k^O(\boldsymbol{x}_t^I), \hat{\boldsymbol{\Sigma}}_k^O) \tag{41}$$

When a unimodal output distribution is required, the law of total mean and variance can be used to approximate the distribution with the Gaussian (see Appendix C):

ToDo

$$m(x) = \mathbb{E}(x_j|t_j) = \sum_{i=1}^K r_{nk} \cdot m_i(t_j) \tag{42}$$

$$var(x) = \sum_{j=1}^K r_{nk} \cdot cov_i \tag{43}$$

# 3 Composition of MPs

## 3.1 Stitching

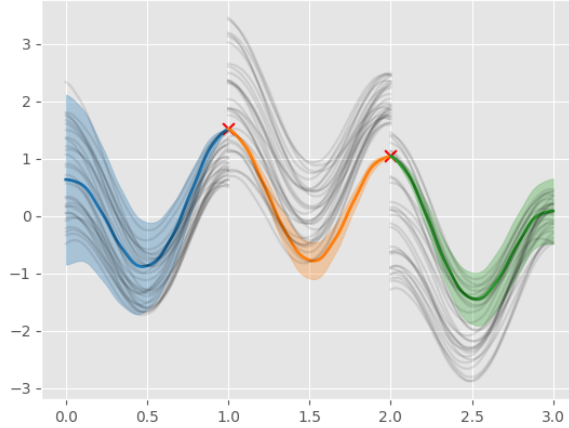The main issue with stitching is the smoothness of the mean and covariance between ProMPs, see Fig. 3.



Figure 3: Stitching three ProMPs.

7

| | |
|---|---|
| $\boldsymbol{t} = [t_1, \ldots, t_m]$ | timestep vector |
| $\boldsymbol{\mu}_{kt}$ | mean of $t_i$ for the $k^{\text{th}}$ Gaussian cluster |
| $h_k(t_i)$ | responsibility of the $k^{\text{th}}$ mixture component at timestep $t_i$ |

Table 7: Notation

## 3.2 Mixture of ProMPs

### 3.2.1 Hierarchical Bayesian Model

Learn the joint distribution $p(t_i, y_i)$ as a mixture of ProMPs, and retrieve the conditional distribution $p(y_i|t_i)$ for each timestep $t_i$.

$$p(t_i, y_i) = \sum_{k=1}^{K} \pi_k \cdot \mathcal{N}_k\big(y_i|t_i; \hat{\boldsymbol{\mu}}_k(t_i), \hat{\boldsymbol{\Sigma}}_k\big) \cdot \mathcal{N}_k\big(t_i|\mu_{kt}, \boldsymbol{\Sigma}_{kt}\big) \tag{44}$$

$$\boldsymbol{x}_t = \begin{bmatrix} \boldsymbol{t} \\ \boldsymbol{y}_t \end{bmatrix}, \quad \boldsymbol{\mu}_k = \begin{bmatrix} \boldsymbol{\mu}_{kt} \\ \boldsymbol{\mu}_{ky} \end{bmatrix} = \begin{bmatrix} \cdot \\ \boldsymbol{\Psi}_t^\top \boldsymbol{\mu}_w \end{bmatrix}, \quad \boldsymbol{\Sigma}_k = \begin{bmatrix} \boldsymbol{\Sigma}_{ktt} & \boldsymbol{\Sigma}_{kty} \\ \boldsymbol{\Sigma}_{kyt} & \boldsymbol{\Sigma}_{kyy} \end{bmatrix} = \begin{bmatrix} \cdot & \cdot \\ \cdot & \boldsymbol{\Psi}_t^\top \boldsymbol{\Sigma}_w \boldsymbol{\Psi}_t + \boldsymbol{\Sigma}_y \end{bmatrix} \tag{45}$$

$$\hat{\boldsymbol{\mu}}_k(t_i) = \boldsymbol{\mu}_{ky} + \boldsymbol{\Sigma}_{kyt}\boldsymbol{\Sigma}_{ktt}^{-1}(t_i - \boldsymbol{\mu}_{kt}) \tag{46}$$

$$\hat{\boldsymbol{\Sigma}}_k = \boldsymbol{\Sigma}_{kyy} - \boldsymbol{\Sigma}_{kyt}\boldsymbol{\Sigma}_{ktt}^{-1}\boldsymbol{\Sigma}_{kty} \tag{47}$$

Marginal probability $p(t_i)$:

$$p(t_i) = \int p(t_i, y_i) dy = \sum_{k=1}^{K} \pi_k \mathcal{N}_k(t_i|\boldsymbol{\mu}_{kt}, \boldsymbol{\Sigma}_{kt}) \tag{48}$$

Conditional probability $p\big(y_i|t_i; \hat{\boldsymbol{\mu}}_k(t_i), \hat{\boldsymbol{\Sigma}}_k\big)$:

$$p\big(y_i|t_i; \hat{\boldsymbol{\mu}}_k(t_i), \hat{\boldsymbol{\Sigma}}_k\big) = \frac{p(t_i, y_i)}{p(t_i)} \tag{49}$$

$$= \frac{\sum_{k=1}^{K} \pi_k \cdot \mathcal{N}_k\big(y_i|t_i; \hat{\boldsymbol{\mu}}_k(t_i), \hat{\boldsymbol{\Sigma}}_k\big) \cdot \mathcal{N}_k\big(t_i|\mu_{kt}, \boldsymbol{\Sigma}_{kt}\big)}{\sum_{k=1}^{K} \pi_k \mathcal{N}_k\big(t_i|\boldsymbol{\mu}_{kt}, \boldsymbol{\Sigma}_{kt}\big)} \tag{50}$$

$$= \sum_{k=1}^{K} h_k(t_i) \mathcal{N}_k\big(y_i|t_i; \hat{\boldsymbol{\mu}}_k(t_i), \hat{\boldsymbol{\Sigma}}_k\big) \tag{51}$$

with:

$$h_k(t_i) = \frac{\pi_k \mathcal{N}_k\big(t_i|\mu_{kt}, \boldsymbol{\Sigma}_{kt}\big)}{\sum_{k=1}^{K} \pi_k \mathcal{N}_k\big(t_i|\boldsymbol{\mu}_{kt}, \boldsymbol{\Sigma}_{kt}\big)} \tag{52}$$

### 3.2.2 Via-Points Modulation

ToDo

## 3.3 Piecewise Gaussian Process

# References

A. Paraschos, C. Daniel, J. R. Peters, and G. Neumann, "Probabilistic Movement Primitives," in *Advances in Neural Information Processing Systems*, vol. 26. Curran Associates, Inc., 2013. [Online]. Available: https://proceedings.neurips.cc/paper/2013/hash/e53a0a2978c28872a4505bdb51db06dc-Abstract.html

A. Paraschos, C. Daniel, J. Peters, and G. Neumann, "Using probabilistic movement primitives in robotics," *Autonomous Robots*, vol. 42, no. 3, pp. 529–551, Mar. 2018. [Online]. Available: https://doi.org/10.1007/s10514-017-9648-7

M. P. Deisenroth, A. A. Faisal, and C. S. Ong, *Mathematics for machine learning*. Cambridge University Press, 2020.

C. M. Bishop and H. Bishop, *Deep Learning: Foundations and Concepts*. Springer International Publishing, 2024. [Online]. Available: https://doi.org/10.1007/978-3-031-45468-4

S. Calinon, "A tutorial on task-parameterized movement learning and retrieval," *Intelligent Service Robotics*, vol. 9, no. 1, pp. 1–29, Jan. 2016. [Online]. Available: https://doi.org/10.1007/s11370-015-0187-9

# A    Hierarchical Bayesian Model proof

*Proof of Eq.* (8). From (Deisenroth et al., 2020), we have the joint distribution:

$$p(\mathbf{x}_a, \mathbf{x}_b) = \mathcal{N}\left( \begin{bmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{bmatrix} \right) \tag{53}$$

and the marginal distribution $p(\mathbf{x}_a)$ of a joint Gaussian distribution $p(\mathbf{x}_a, \mathbf{x}_b)$:

$$p(\mathbf{x}_a) = \int p(\mathbf{x}_a, \mathbf{x}_b) d\mathbf{x}_b = \mathcal{N}(\mathbf{x}_a | \boldsymbol{\mu}_a, \boldsymbol{\Sigma}_{aa}) \tag{54}$$

Since $\boldsymbol{y}_t$ and $\boldsymbol{w}$ are jointly Gaussian, we have:

$$\begin{bmatrix} \boldsymbol{y}_t \\ \boldsymbol{w} \end{bmatrix} \sim \mathcal{N}\left( \begin{bmatrix} \boldsymbol{\mu}_{\boldsymbol{y}_t} \\ \boldsymbol{\mu}_{\boldsymbol{w}} \end{bmatrix}, \begin{bmatrix} \text{Cov}[\boldsymbol{y}_t, \boldsymbol{y}_t] & \text{Cov}[\boldsymbol{y}_t, \boldsymbol{w}] \\ \text{Cov}[\boldsymbol{w}, \boldsymbol{y}_t] & \text{Cov}[\boldsymbol{w}, \boldsymbol{w}] \end{bmatrix} \right) \tag{55}$$

$$\boldsymbol{\mu}_{\boldsymbol{y}_t} = \mathbb{E}[\boldsymbol{y}_t] \tag{56}$$
$$= \mathbb{E}[\boldsymbol{\Psi}_t^\top \boldsymbol{w} + \boldsymbol{\epsilon}_y] \tag{57}$$
$$= \boldsymbol{\Psi}_t^\top \mathbb{E}[\boldsymbol{w}] + \mathbb{E}[\boldsymbol{\epsilon}_y] \tag{58}$$
$$= \boldsymbol{\Psi}_t^\top \boldsymbol{\mu}_{\boldsymbol{w}} + 0 \tag{59}$$
$$= \boldsymbol{\Psi}_t^\top \boldsymbol{\mu}_{\boldsymbol{w}} \tag{60}$$

$$\text{Cov}[\boldsymbol{y}_t, \boldsymbol{y}_t] = \text{Cov}[\boldsymbol{\Psi}_t^\top \boldsymbol{w} + \boldsymbol{\epsilon}_y] \tag{61}$$
$$= \text{Cov}[\boldsymbol{\Psi}_t^\top \boldsymbol{w}] + \text{Cov}[\boldsymbol{\epsilon}_y] \tag{62}$$
$$= \boldsymbol{\Psi}_t^\top \text{Cov}[\boldsymbol{w}] \boldsymbol{\Psi}_t + \boldsymbol{\Sigma}_y \tag{63}$$
$$= \boldsymbol{\Psi}_t^\top \boldsymbol{\Sigma}_{\boldsymbol{w}} \boldsymbol{\Psi}_t + \boldsymbol{\Sigma}_y \tag{64}$$

$$\begin{bmatrix} \boldsymbol{y}_t \\ \boldsymbol{w} \end{bmatrix} \sim \mathcal{N}\left( \begin{bmatrix} \boldsymbol{\Psi}_t^\top \boldsymbol{\mu}_{\boldsymbol{w}} \\ \boldsymbol{\mu}_{\boldsymbol{w}} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Psi}_t^\top \boldsymbol{\Sigma}_{\boldsymbol{w}} \boldsymbol{\Psi}_t + \boldsymbol{\Sigma}_y & \boldsymbol{\Psi}_t^\top \boldsymbol{\Sigma}_{\boldsymbol{w}} \\ \boldsymbol{\Sigma}_{\boldsymbol{w}} \boldsymbol{\Psi}_t & \boldsymbol{\Sigma}_{\boldsymbol{w}} \end{bmatrix} \right) \tag{65}$$

$$p(\boldsymbol{y}_t; \boldsymbol{\theta}) = \mathcal{N}\left( \boldsymbol{y}_t | \boldsymbol{\Psi}_t^\top \boldsymbol{\mu}_{\boldsymbol{w}}, \boldsymbol{\Psi}_t^\top \boldsymbol{\Sigma}_{\boldsymbol{w}} \boldsymbol{\Psi}_t + \boldsymbol{\Sigma}_y \right) \tag{66}$$

$\square$

# B  Via-Points conditioning proof

*Proof of Eq. (11) and Eq. (12).* With the joint distribution $p(\mathbf{x}_a, \mathbf{x}_b)$ in Eq. (53), and from (Bishop and Bishop, 2024), the parameters of a conditional multivariate Gaussian $p(\mathbf{x}_a|\mathbf{x}_b) = \mathcal{N}(\boldsymbol{\mu}_{a|b}, \boldsymbol{\Sigma}_{a|b})$ are the following:

$$\boldsymbol{\mu}_{a|b} = \boldsymbol{\mu}_a + \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}(\mathbf{x}_b - \boldsymbol{\mu}_b) \tag{67}$$

$$\boldsymbol{\Sigma}_{a|b} = \boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba} \tag{68}$$

We want the posterior $p(\boldsymbol{w}|\boldsymbol{x}_t^\star)$, knowing the likelihood $\boldsymbol{x}_t^\star|\boldsymbol{w} \sim \mathcal{N}\left(\boldsymbol{y}_t^\star|\boldsymbol{\Psi}_t^\top \boldsymbol{w}, \boldsymbol{\Sigma}_t^\star\right)$, and the prior $\boldsymbol{w} \sim \mathcal{N}(\boldsymbol{w}|\boldsymbol{\mu}_w, \boldsymbol{\Sigma}_w)$.

$$\begin{bmatrix} \boldsymbol{w} \\ \boldsymbol{x}_t^\star \end{bmatrix} \sim \mathcal{N}\left( \begin{bmatrix} \boldsymbol{\mu}_w \\ \boldsymbol{\Psi}_t^\top \boldsymbol{\mu}_w \end{bmatrix}, \begin{bmatrix} \mathrm{Cov}[\boldsymbol{w}, \boldsymbol{w}] & \mathrm{Cov}[\boldsymbol{w}, \boldsymbol{x}_t^\star] \\ \mathrm{Cov}[\boldsymbol{x}_t^\star, \boldsymbol{w}] & \mathrm{Cov}[\boldsymbol{x}_t^\star, \boldsymbol{x}_t^\star] \end{bmatrix} \right) \tag{69}$$

$\mathrm{Cov}[\boldsymbol{x}_t^\star, \boldsymbol{x}_t^\star]$ follows from Eq. (64).

$$\begin{aligned}
\mathrm{Cov}[\boldsymbol{w}, \boldsymbol{x}_t^\star] &= \mathrm{Cov}[\boldsymbol{w}, \boldsymbol{\Psi}_t^\top \boldsymbol{w} + \boldsymbol{\epsilon}_y] & &\tag{70} \\
&= \mathrm{Cov}[\boldsymbol{w}, \boldsymbol{\Psi}_t^\top \boldsymbol{w}] & (\mathrm{Cov}[\boldsymbol{w}, \boldsymbol{\epsilon}_y] = 0 \text{ since } \boldsymbol{\epsilon}_y \text{ is independent of } \boldsymbol{w}) &\tag{71} \\
&= \mathbb{E}[(\boldsymbol{w} - \boldsymbol{\mu}_w)(\boldsymbol{\Psi}_t^\top \boldsymbol{w} - \boldsymbol{\Psi}_t^\top \boldsymbol{\mu}_w)^\top] & &\tag{72} \\
&= \mathbb{E}[(\boldsymbol{w} - \boldsymbol{\mu}_w)(\boldsymbol{w} - \boldsymbol{\mu}_w)^\top \boldsymbol{\Psi}_t] & &\tag{73} \\
&= \mathrm{Cov}[\boldsymbol{w}, \boldsymbol{w}] \cdot \boldsymbol{\Psi}_t & &\tag{74} \\
&= \boldsymbol{\Sigma}_w \boldsymbol{\Psi}_t & &\tag{75}
\end{aligned}$$

$$\begin{bmatrix} \boldsymbol{w} \\ \boldsymbol{x}_t^\star \end{bmatrix} \sim \mathcal{N}\left( \begin{bmatrix} \boldsymbol{\mu}_w \\ \boldsymbol{\Psi}_t^\top \boldsymbol{\mu}_w \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_w & \boldsymbol{\Sigma}_w \boldsymbol{\Psi}_t \\ \boldsymbol{\Psi}_t^\top \boldsymbol{\Sigma}_w & \boldsymbol{\Psi}_t^\top \boldsymbol{\Sigma}_w \boldsymbol{\Psi}_t + \boldsymbol{\Sigma}_t^\star \end{bmatrix} \right) \tag{76}$$

Using Eq. (67) we get:

$$\boldsymbol{\mu}_{\boldsymbol{w}|\boldsymbol{x}_t^\star} = \boldsymbol{\mu}_w + \boldsymbol{\Sigma}_w \boldsymbol{\Psi}_t \left( \boldsymbol{\Sigma}_t^\star + \boldsymbol{\Psi}_t^\top \boldsymbol{\Sigma}_w \boldsymbol{\Psi}_t \right)^{-1} (\boldsymbol{y}_t^\star - \boldsymbol{\Psi}_t^\top \boldsymbol{\mu}_w) \tag{77}$$

Using Eq. (68) we get:

$$\boldsymbol{\Sigma}_{\boldsymbol{w}|\boldsymbol{x}_t^\star} = \boldsymbol{\Sigma}_w - \boldsymbol{\Sigma}_w \boldsymbol{\Psi}_t \left( \boldsymbol{\Sigma}_t^\star + \boldsymbol{\Psi}_t^\top \boldsymbol{\Sigma}_w \boldsymbol{\Psi}_t \right)^{-1} \boldsymbol{\Psi}_t^\top \boldsymbol{\Sigma}_w \tag{78}$$

$\square$

# C  Gaussian mixture regression approximated by a single normal distribution

*Proof of Eq. (11) and Eq. (12).* From (Calinon, 2016), we have:

ToDo

$\square$