
Data-Efficient Reinforcement Learning in Continuous State-Action POMDPs

Anonymous Author(s)

Affiliation

Address

email

Abstract

We present a data-efficient reinforcement learning method for continuous state-action systems under significant observation noise. Data-efficient solutions under small noise exist, such as PILCO which learns the cartpole swing-up task in 30s. PILCO evaluates policies by planning state-trajectories using a dynamics model. However, PILCO applies policies to the observed state, therefore planning in *observation*-space. We extend PILCO with filtering to instead plan in *belief*-space, consistent partially observable Markov decisions process (POMDP) planning. This enables data-efficient learning under significant observation noise, outperforming more naive methods such as *post-hoc* application of a filter to policies optimised by the original (unfiltered) PILCO algorithm. We test our method on the cartpole swing-up task, which involves nonlinear dynamics and requires nonlinear control.

1 Introduction

The Probabilistic Inference and Learning for Control (PILCO) [3] framework is a reinforcement learning algorithm, which uses Gaussian Processes (GPs) to learn the dynamics in continuous state spaces. The method has shown to be highly *efficient* in the sense that it can learn with only very few interactions with the real system. However, a serious limitation of PILCO is that it assumes that the observation noise level is small. There are two main reasons which make this assumption necessary. Firstly, the dynamics are learnt from the noisy observations, but learning the transition model in this way doesn't correctly account of the noise in the observations. If the noise is assumed small, then this will be a good approximation to the real transition function. Secondly, PILCO uses the noisy observation directly to calculate the action, which is problematic if the observation noise is substantial. Imagine a policy controlling an unstable system, where high gain feed-back is necessary for good performance. Observation noise is *amplified* when the noisy input is fed directly to the high gain controller, which in turn injects noise back into the state, creating cycles of increasing variance and instability.

In this paper we extend PILCO to address these two shortcomings, enabling PILCO to be used in situations with substantial observation noise. The first issue is addressed using the so-called *direct* method for training the transition model, see section 3.3. The second problem can be tackled by *filtering* the observations. One way to look at this is that PILCO does planning in observation space, rather than in belief space. In this paper we extend PILCO to allow filtering of the state, by combining the previous state distribution with the dynamics model and the observation using Bayes rule. Note, that this is easily done when the controller is being applied, but to gain the full benefit, we have to also take the filter into account when optimising the policy.

PILCO trains its policy through minimising the expected predicted loss when *simulating* the the system and controller actions. Since the dynamics are not known exactly, the simulation in PILCO had to simulate *distributions* of possible trajectories of the physical state of the system. This was achieved using an analytical approximation based on moment-matching and Gaussian state distributions. In

38 this paper we thus need to augment the simulation over physical states to include the state of the
 39 filter, an *information state* or *belief state*. This is a bit more complicated as the belief state itself is a
 40 probability distribution, we will now have to simulate distributions over distributions. This will allow
 41 the algorithm both to apply filtering during control but also to anticipate the effect of filtering during
 42 training, thereby learning a better policy.

43 We will first give a brief outline of related work in section 2 and the original PILCO algorithm
 44 in section 3, including the proposed use of the ‘direct method’ for training dynamics from noisy
 45 observations in section 3.3. In section 4 will derive the algorithm for POMDP training or planning
 46 in belief space. Note an assumption is that we observe noisy versions of the state variables. We
 47 do not handle more general cases where other unobserved states are also learnt nor learn any other
 48 mapping from the state space to observations other than additive Gaussian noise. In the final sections
 49 we show experimental results showing that the proposed algorithm handles observation noise better
 50 than competing algorithms, and close with conclusions and discussion.

51 2 Related Work

52 Implementing a filter is straightforward when the system dynamics are *known* and *linear*, referred
 53 to as Kalman filtering. For known nonlinear systems, the extended Kalman filter (EKF) is often
 54 adequate (e.g. [10]), as long as the dynamics are *locally linear*, meaning approximately linear within
 55 the region covered by the belief distribution. Otherwise, the EKF’s first order Taylor expansion
 56 approximation breaks down. Greater nonlinearities usually warrant the unscented Kalman filter (UKF)
 57 or particle methods [5, 8]. The UKF uses a deterministic sampling technique to estimate moments.
 58 However, if moments can be computed analytically and exactly, moment-matching methods are
 59 preferred. Moment-matching using distributions from the exponential family (e.g. Gaussians) is
 60 equivalent to optimising the Kullback-Leibler divergence $KL(p||q)$ between the true distribution p
 61 and an approximate distribution q . In such cases, moment-matching is less susceptible to model bias
 62 than the EKF due to its conservative predictions [2]. Unfortunately, the literature does not provide a
 63 continuous state-action method that is both data efficient and resistant to noise when the dynamics are
 64 *unknown* and *locally nonlinear*. Sometimes the dynamics are partially-known, with known functional
 65 form yet unknown parameters. Such ‘grey-box’ problems have the aesthetic solution of incorporating
 66 the unknown dynamics parameters into state, reducing the learning task into a POMDP planning task
 67 [4, 11, 9]. Finite state-action space tasks can be similarly solved, e.g. only two Dirichlet parameters
 68 are required to model the probability of each of the finitely-many state-action-state transitions [7].
 69 However, such solutions are not suitable for continuous-state ‘black-box’ problems with no prior
 70 dynamics knowledge. The original PILCO does not assume any prior dynamics knowledge, yet
 71 assumes full state observability and fails under moderate sensor noise. One proposed solution is to
 72 filter observations during policy execution [2]. However, without also predicting system trajectories
 73 w.r.t. the filtering process, the above method merely optimises policies for unfiltered control, not
 74 for filtered control. The mismatch between unfiltered-prediction and filtered-execution restricts
 75 PILCO’s ability to take full advantage of filtering. Dallaire et al. [1] optimise a policy using a more
 76 realistic filtered-prediction. However, the method neglects model uncertainty by only using the
 77 maximum a posteriori (MAP) model. Unlike the method of Deisenroth and Peters [2] which gives a
 78 full probabilistic treatment of the dynamics predictions, work by Dallaire et al. [1] is therefore highly
 79 susceptible to model error, hampering data-efficiency.

80 We instead predict system trajectories using closed loop filtered control precisely because we execute
 81 closed loop filtered control. The resulting policies are thus optimised for the specific case in which
 82 they are used. Doing so, our method retains the same data-efficiency properties of PILCO whilst
 83 applicable to tasks with high observation noise. To evaluate our method, we use the benchmark
 84 cartpole swing-up task with noisy sensors. We show realistic and probabilistic prediction (to consider
 85 uncertainty) helps our method outperform aforementioned methods.

86 3 The PILCO Algorithm

87 PILCO is a model-based policy-search RL algorithm. It applies to continuous-state, continuous-
 88 action, continuous-observation and discrete-time control tasks. A probabilistic dynamics model is
 89 used to predict one-step system dynamics (from one timestep to the next). This allows PILCO to
 90 probabilistically predict multi-step system trajectories over arbitrary time horizon T , by repeatedly
 91 using the predictive dynamics model’s output at one timestep, as the (uncertain) input in the following

92 timestep. For tractability PILCO uses moment-matching to keep the latent state distribution Gaussian.
 93 The result is an analytic distribution of state-trajectories, approximated as a joint Gaussian distribution
 94 over T states. The policy is evaluated as the expected total cost of the trajectories. Next, the policy is
 95 improved using local gradient-based optimisation, searching over policy-parameter space. A distinct
 96 advantage of moment-matched prediction for policy search instead of particle methods is smoother
 97 policy gradients and less local optima [6]. Finally, the policy is executed, generating additional data
 98 to re-train the dynamics model. The whole process then repeats until policy convergence. For the
 99 remainder of this section we discuss, step by step, PILCO summarised by Algorithm 1. The user first
 100 defines a parametric policy π function (Algorithm 1, line 1) and then initialises the policy parameters
 101 ψ randomly (line 2) since we begin without any data.

102 3.1 System Execution Phase

103 With a policy now defined, PILCO is ready to *execute* the system (Algorithm 1, line 4). Let the
 104 latent state of the system at time t be $x_t \in \mathbb{R}^D$, which is noisily observed as $z_t = x_t + \epsilon_t$, where
 105 $\epsilon_t \stackrel{iid}{\sim} \mathcal{N}(0, \Sigma^\epsilon)$. The policy π , parameterised by ψ , takes observation z_t as input, and outputs a
 106 control action $u_t = \pi(z_t, \psi) \in \mathbb{R}^F$. Applying action u_t to the dynamical system in state x_t , results
 107 in a new system state x_{t+1} . Repeating until horizon T results in a new single state-trajectory of data.

108 3.2 Learning Dynamics

109 To learn the unknown dynamics (Algorithm 1, line 5), any probabilistic model flexible enough
 110 to capture the complexity of the dynamics can be used. Bayesian nonparametric models are par-
 111 ticularly suited given their resistance to overfitting and underfitting respectively. Overfitting oth-
 112 erwise leads to model bias, and underfitting limits the complexity of the system this method can
 113 learn to control. In a nonparametric model no prior dynamics knowledge is required, not even
 114 knowledge of *how complex* the unknown dynamics might be since the model’s complexity grows
 115 with the available data. We define the latent dynamics $f : \tilde{x}_t \rightarrow x_{t+1}$, where $\tilde{x}_t \doteq [x_t^\top, u_t^\top]^\top$.
 116 PILCO models the dynamics with D independent Gaussian process (GP) priors, one for each
 117 dynamics output variable: $f^a : \tilde{x}_t \rightarrow x_{t+1}^a$, where $a \in [1, D]$ is the a ’th dynamics output, and
 118 $f^a \sim \mathcal{GP}(\phi_a^\top \tilde{x}, k(\tilde{x}_i, \tilde{x}_j))$. Note we implement PILCO with a linear mean function¹ $\phi_a^\top \tilde{x}$. The
 119 covariance function k is squared exponential, with length scales $\Lambda_a = \text{diag}([l_{a,1}^2, \dots, l_{a,D+F}^2])$, and
 120 signal variance s_a^2 : $k(\tilde{x}_i, \tilde{x}_j) = s_a^2 \exp(-\frac{1}{2}(\tilde{x}_i - \tilde{x}_j)^\top \Lambda_a^{-1}(\tilde{x}_i - \tilde{x}_j))$.

121 3.3 Learning Dynamics from Noisy Observations

122 The original PILCO algorithm ignored sensor noise when training each GP by assuming each
 123 observation z_t to be the latent state x_t . However, this approximation breaks down under significant
 124 noise. More complex training schemes are required for each GP that correctly treat each training
 125 datum x_t as latent, yet noisily-observed as z_t . We resort to GP state space model methods, specifically
 126 the ‘direct method’ [6, section 3.5]. The direct method infers the marginal likelihood $p(z_{1:N})$
 127 approximately using moment-matching and a single forward-pass. Doing so, it specifically exploits
 128 the time series structure that generated observations $z_{1:N}$. We use the direct method to set the
 129 GP’s training data $\{x_{1:N}, u_{1:N}\}$ and observation noise variance Σ^ϵ to the inducing point parameters
 130 and noise parameters that optimise the marginal likelihood. In this paper we use the superior
 131 Direct method to train GPs, both in our extended version of PILCO presented section 4, and in our
 132 implementation of the original PILCO algorithm for fair comparison in the experiments.

133 3.4 System Prediction Phase

134 In contrast to executions, PILCO also *predicts* analytic distributions of state-trajectories (Algorithm 1,
 135 line 6) for policy evaluation. PILCO does this offline, between the online system executions. Predicted
 136 control is identical to executed control except each aforementioned quantity is instead now a random
 137 variable, distinguished with capitals: $X_t, Z_t, U_t, \tilde{X}_t$ and X_{t+1} , all approximated as jointly Gaussian.
 138 These variables interact both in execution and prediction according to Figure 1. To predict X_{t+1} now
 139 that \tilde{X}_t is uncertain PILCO uses the iterated law of expectation and variance:

$$p(X_{t+1} | \tilde{X}_t) = \mathcal{N}(\mu_{t+1}^x = \mathbb{E}_{\tilde{X}}[\mathbb{E}_f[f(\tilde{X}_t)]], \Sigma_{t+1}^x = \mathbb{V}_{\tilde{X}}[\mathbb{E}_f[f(\tilde{X}_t)]] + \mathbb{E}_{\tilde{X}}[\mathbb{V}_f[f(\tilde{X}_t)]]). \quad (1)$$

140 After a one-step prediction from X_0 to X_1 , PILCO repeats the process from X_1 to X_2 , and up to X_T ,
 141 resulting in a multi-step prediction whose joint we refer to as a distribution over state-trajectories.

¹The original PILCO instead uses a zero mean function, and instead predicts relative changes in state.

Algorithm 1 PILCO

```
1: Define policy’s functional form:  $\pi : z_t \times \psi \rightarrow u_t$ .
2: Initialise policy parameters  $\psi$  randomly.
3: repeat
4:   Execute system, record data.
5:   Learn dynamics model  $p(f)$ .
6:   Predict state trajectories from  $p(X_0)$  to  $p(X_T)$ .
7:   Evaluate policy:  $J(\psi) = \sum_{t=0}^T \gamma^t \mathcal{E}_t$ ,  $\mathcal{E}_t = \mathbb{E}_X[\text{cost}(X_t)|\psi]$ .
8:   Improve policy:  $\psi \leftarrow \text{argmin}_{\psi} J(\psi)$ .
9: until policy parameters  $\psi$  converge
```

142 3.5 Policy Evaluation and Improvement

143 To evaluate the policy π (or more specifically, the policy parameters ψ), PILCO computes the loss
144 $J(\psi)$ by applying a cost function to the marginal state distribution at each timestep (see Algorithm 1,
145 line 7). After policy evaluation, PILCO optimises the policy using the analytic gradients of J . A
146 BFGS optimisation method searches for the set of policy parameters ψ that minimise the total cost
147 $J(\psi)$ using gradients information $dJ/d\psi$ (Algorithm 1, line 8). To compute $dJ/d\psi$ we require
148 derivatives $d\mathcal{E}_t/d\psi$ at each time t to chain together, detailed in [3].

149 4 Our Method: PILCO Extended with Bayesian Filtering

150 Here we describe the novel aspects of our method. Our method uses the same high-level algorithm
151 as PILCO (Algorithm 1). However, we modify² two subroutines to extend PILCO from MDPs to
152 a special-case of POMDPs (specifically where the partial observability has the form of Gaussian
153 noise on the latent state). First, we filter observations during system execution (Algorithm 1, line 4),
154 detailed section 4.1. Second, we predict *belief*-trajectories (instead of state-trajectories) through
155 the filter in addition to PILCO’s dynamics model (line 6), detailed section 4.2. Filtering maintains
156 a belief posterior of the latent system state. The belief is conditioned on, not just the most recent
157 observation, but all previous observations (Figure 2). Such additional conditioning has the benefit of
158 providing a less-noisy and more-informed input to the policy: the filtered belief-mean instead of the
159 raw observation z_t . Our implementation continues PILCO’s distinction between *executing* the system
160 (resulting in a single real belief-trajectory) and *predicting* the system’s responses (which in our case
161 yields an analytic distribution of multiple possible future belief-trajectories). During the execution
162 phase, the system reads specific observations z_t . Our method additionally maintains a belief state
163 $b \sim \mathcal{N}(m, V)$ by filtering observations. This belief state b can be treated as a random variable
164 with a distribution parameterised by belief-mean m and belief-certainty V . Note both m and V are
165 functions of previous observations $z_{1:t}$. Now, during the (probabilistic) system prediction phase, future
166 observations are instead *random* variables (since they have not been observed yet), distinguished
167 as Z . Since the belief parameters m and V are functions of the now-random observations, the
168 belief parameters must be random also, distinguished as M and V' . Given the belief’s distribution
169 parameters are now random, the belief is *hierarchically*-random, denoted $B \sim \mathcal{N}(M, V')$. Our
170 framework allows us to consider multiple possible future belief-states *analytically* during policy
171 evaluation. Intuitively, this framework is the analytical analogue of POMDP policy evaluation using
172 particle methods. In particle methods, each particle is associated with a distinct belief, due to each
173 conditioning on independent samples of future observations. A particle distribution thus defines a
174 distribution over beliefs. Our method is the analytical analogue of this particle distribution. And by
175 restricting our beliefs as (parametric) Gaussian, we can tractably encode a distribution over beliefs by
176 a distribution over belief-parameters.

177 4.1 Filtered-System Execution Phase

178 When an actual filter is applied, it starts with three pieces of information: $m_{t|t-1}$, $V_{t|t-1}$ and a noisy
179 observation of the system z_t (the dual subscript means belief of the latent physical state x at time t
180 given all observations up until time $t - 1$ inclusive). The filtering ‘update step’ combines prior belief
181 $b_{t|t-1} = p(x_t|z_{1:t-1}, u_{1:t-1}) \sim \mathcal{N}(m_{t|t-1}, V_{t|t-1})$ with observational likelihood $p(x_t) = \mathcal{N}(z_t, \Sigma^\epsilon)$

²We implement our method by modifying PILCO’s source code: <http://mlg.eng.cam.ac.uk/pilco/>.

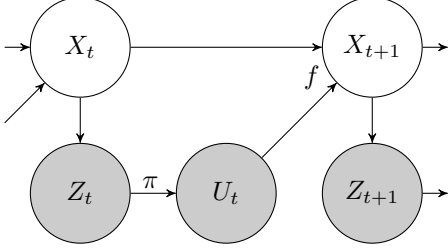


Figure 1: **The original (unfiltered) PILCO, as a probabilistic graphical model.** At each timestep, the latent system X_t is observed noisily as Z_t which is inputted directly into policy function π to decide action U_t . Finally, the latent system will evolve to X_{t+1} , according to the unknown, nonlinear dynamics function f of the previous state X_t and action U_t .

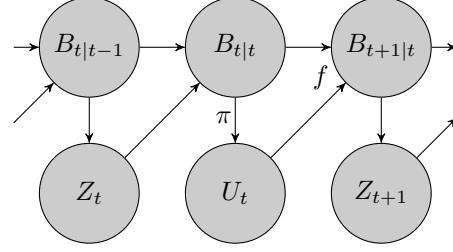


Figure 2: **Our method (PILCO extended with Bayesian filtering).** Our prior belief $B_{t|t-1}$ (over latent system X_t), generates observation Z_t . The prior belief $B_{t|t-1}$ then combines with observation Z_t resulting in posterior belief $B_{t|t}$ (the update step). Then, the mean posterior belief $\mathbb{E}[B_{t|t}]$ is inputted into policy function π to decide action U_t . Finally, the next timestep’s prior belief $B_{t+1|t}$ is predicted using dynamics model f (the prediction step).

182 using Bayes rule to yield posterior belief $b_{t|t} = p(x_t | z_{1:t}, u_{1:t-1})$:

$$b_{t|t} \sim \mathcal{N}(m_{t|t}, V_{t|t}), \quad m_{t|t} = W_m m_{t|t-1} + W_z z_t, \quad V_{t|t} = W_m V_{t|t-1}, \quad (2)$$

183 with weight matrices $W_m = \Sigma^\epsilon (V_{t|t-1} + \Sigma^\epsilon)^{-1}$ and $W_z = V_{t|t-1} (V_{t|t-1} + \Sigma^\epsilon)^{-1}$ computed from
 184 the standard result of a product of two Gaussians. The policy π instead uses updated belief-mean
 185 $m_{t|t}$ (smoother and better-informed than z_t) to decide the action: $u_t = \pi(m_{t|t}, \psi)$. Thus, the joint
 186 distribution over the updated (random) belief and the (non-random) action is

$$\tilde{b}_{t|t} \doteq \begin{bmatrix} b_{t|t} \\ u_t \end{bmatrix} \sim \mathcal{N} \left(\tilde{m}_{t|t} \doteq \begin{bmatrix} m_{t|t} \\ u_t \end{bmatrix}, \tilde{V}_{t|t} \doteq \begin{bmatrix} V_{t|t} & 0 \\ 0 & 0 \end{bmatrix} \right). \quad (3)$$

187 Next, the filtering ‘prediction step’ computes the predictive-distribution of $b_{t+1|t} = p(x_{t+1} | z_{1:t}, u_{1:t})$
 188 from the output of dynamics model f given random input $\tilde{b}_{t|t}$. The distribution $f(\tilde{b}_{t|t})$ is non-
 189 Gaussian yet has analytically computable moments [3]. For tractability, we approximate $b_{t+1|t}$ as
 190 Gaussian-distributed using moment-matching:

$$b_{t+1|t} \sim \mathcal{N}(m_{t+1|t}, V_{t+1|t}), \quad m_{t+1|t}^a = \mathbb{E}_{\tilde{b}_{t|t}}[f^a(\tilde{b}_{t|t})], \quad V_{t+1|t}^{ab} = \mathbb{C}_{\tilde{b}_{t|t}}[f^a(\tilde{b}_{t|t}), f^b(\tilde{b}_{t|t})], \quad (4)$$

191 where a and b refer to the a ’th and b ’th dynamics output. Both $m_{t+1|t}^a$ and $V_{t+1|t}^{ab}$ are derived in
 192 Appendix C. The process then repeats using the predictive belief (Eq. 4) as the prior belief in the
 193 following timestep. This completes the specification of the system in execution.

194 4.2 Filtered-System Prediction Phase

195 In *system prediction*, we compute the probabilistic behaviour of the filtered system via an analytic
 196 distribution of possible beliefs. A distribution over beliefs b is in principle a distribution over its
 197 parameters m and V . To distinguish m , V and b as now being *random* and *hierarchically-random*
 198 respectively, we capitalise them: M , V' and B . As an approximation we do not consider random-
 199 variance V' , but instead consider the non-random $V = \mathbb{E}[V']$ (which has a fixed value, for a given
 200 timestep). Restricting M to being Gaussian distributed then we begin system prediction with the
 201 joint:

$$\begin{bmatrix} M_{t|t-1} \\ Z_t \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mu_{t|t-1}^m \\ \mu_t^z \end{bmatrix}, \begin{bmatrix} \Sigma_{t|t-1}^m & 0 \\ 0 & \Sigma_t^z \end{bmatrix} \right), \quad (5)$$

202 where $\Sigma_t^z = \Sigma_{t|t-1}^m + V_{t|t-1} + \Sigma^\epsilon$. The updated belief posterior is also Gaussian,

$$M_{t|t} \sim \mathcal{N} \left(\mu_{t|t}^m = \mu_{t|t-1}^m, \Sigma_{t|t}^m = W_m \Sigma_{t|t-1}^m W_m^\top + W_z \Sigma_t^z W_z^\top \right). \quad (6)$$

203 The policy now has a random input $M_{t|t}$, thus the control output must also be random (even though π is
 204 a deterministic function): $U_t = \pi(M_{t|t}, \psi)$, which we implement by overloading the policy function:
 205 $(\mu_t^u, \Sigma_t^u, C_t^{mu}) = \pi(\mu_{t|t}^m, \Sigma_{t|t}^m, \psi)$, where μ_t^u is the output mean, Σ_t^u the output variance and C_t^{mu}
 206 input-output covariance with premultiplied inverse input variance, $C_t^{mu} \doteq (\Sigma_{t|t}^m)^{-1} C_M[M_{t|t}, U_t]$.

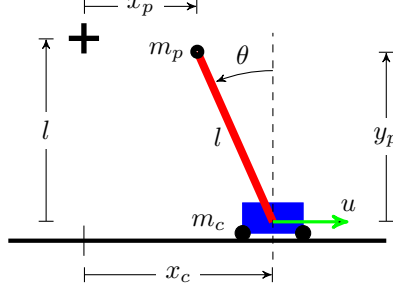


Figure 3: **The cartpole swing-up task.** A pendulum of length l is attached to a cart by a frictionless pivot. The cart has mass m_c and position x_c . The pendulum’s endpoint has mass m_p and position (x_p, y_p) , with angle θ from vertical. The cart begins at position $x_c = 0$ and pendulum hanging down: $\theta = \pi$. The goal is to accelerate the cart by applying horizontal force u_t at each timestep t to invert then stabilise the pendulum’s endpoint at the goal (black cross), i.e. to maintain $x_c = 0$ and $\theta = 0$.

207 Making a moment-matched approximation yields a joint Gaussian:

$$\tilde{M}_{t|t} \doteq \begin{bmatrix} M_{t|t} \\ U_t \end{bmatrix} \sim \mathcal{N} \left(\tilde{\mu}_{t|t} \doteq \begin{bmatrix} \mu_{t|t}^m \\ \mu_{t|t}^u \end{bmatrix}, \tilde{\Sigma}_{t|t} \doteq \begin{bmatrix} \Sigma_{t|t}^m & \Sigma_{t|t}^{mu} \\ (C_t^{mu})^\top \Sigma_{t|t}^m & \Sigma_t^u \end{bmatrix} \right) \quad (7)$$

208 Finally, we probabilistically predict the belief-mean distribution $p(M_{t+1|t})$ and the expected belief-
 209 variance $\tilde{V}_{t+1|t} = \mathbb{E}_{\tilde{M}_{t|t}}[V_{t+1|t}]$, both detailed in Appendix D. We have now discussed the one-step
 210 prediction of the filtered system, from $B_{t|t-1}$ to $B_{t+1|t}$. Using this process repeatedly, from initial
 211 belief $B_{0|0}$ we one-step predict to $B_{1|1}$, then to $B_{2|2}$ etc., up to $B_{T|T}$.

212 4.3 Policy Evaluation and Improvement

213 To evaluate a policy we again apply the loss function J (Algorithm 1, line 7) to the multi-step
 214 prediction (section 4.2). The policy is again optimised using the analytic gradients of J . Since
 215 J now is a function of beliefs, we additionally consider the gradients of $B_{t|t-1}$ w.r.t. ψ . As the
 216 belief is distributed by $B_{t|t-1} \sim \mathcal{N}(M_{t|t-1}, V_{t|t-1}) \sim \mathcal{N}(\mathcal{N}(\mu_{t|t-1}^m, \Sigma_{t|t-1}^m), V_{t|t-1})$, we use partial
 217 derivatives of $\mu_{t|t-1}^m$, $\Sigma_{t|t-1}^m$ and $V_{t|t-1}$ w.r.t. each other and w.r.t ψ , detailed in Appendix A.

218 5 Experiments

219 We test our algorithm on the cartpole swing-up problem (Figure 3), a benchmark for comparing
 220 controllers of nonlinear dynamical systems. We experiment using a physics simulator by solving the
 221 differential equations of the system. The pendulum begins each episode hanging downwards with
 222 the goal of swinging it up and stabilising it. The use a cart mass of $m_c = 0.5\text{kg}$. A zero-order hold
 223 controller applies horizontal forces to the cart within range $[-10, 10]\text{N}$. The policy is a mixture of
 224 100 radial basis functions. Friction resists the cart’s motion with damping coefficient $b = 0.1\text{Ns/m}$.
 225 Connected to the cart is a pole of length $l = 0.2\text{m}$ and mass $m_p = 0.5\text{kg}$ located at its endpoint,
 226 which swings due to gravity’s acceleration $g = 9.82\text{m/s}^2$. An inexpensive camera observes the
 227 system. Frame rates of \$10 webcams are typically 30Hz at maximum resolution, thus the time
 228 discretisation is $\Delta t = 1/30\text{s}$. The state x comprises the cart position, pendulum angle, and their time
 229 derivatives $x = [x_c, \theta, \dot{x}_c, \dot{\theta}]^\top$. The cartpole’s motion is described with the differential equation:

$$\dot{x}^\top = \left[\dot{x}_c, \dot{\theta}, \frac{-2m_p l \dot{\theta}^2 s + 3m_p g s c + 4u - 4b \dot{x}_c}{4(m_c + m_p) - 3m_p c^2}, \frac{-3m_p l \dot{\theta}^2 s c + 6(m_c + m_p) g s + 6(u - b \dot{x}_c c)}{4l(m_c + m_p) - 3m_p l c^2} \right], \quad (8)$$

230 using shorthand $s = \sin \theta$ and $c = \cos \theta$. We both randomly-initialise the system and set the
 231 initial belief of the system according to $B_{0|0} \sim \mathcal{N}(M_{0|0}, V_{0|0})$ where $M_{0|0} \sim \delta([0, \pi, 0, 0]^\top)$
 232 and $V_{0|0}^{1/2} = \text{diag}([0.2\text{m}, 0.2\text{rad}, 0.2\text{m/s}, 0.2\text{rad/s}])$. The camera’s noise standard deviation is:
 233 $(\Sigma^\epsilon)^{1/2} = \text{diag}([0.03\text{m}, 0.03\text{rad}, \frac{0.03}{\Delta t}\text{m/s}, \frac{0.03}{\Delta t}\text{rad/s}])$, noting $0.03\text{rad} \approx 1.7^\circ$. We use the $\frac{0.03}{\Delta t}$ terms
 234 since using a camera we cannot observe velocities directly but can estimate them with finite differ-
 235 ences. Each episode has a two second time horizon (60 timesteps). The cost function we impose is
 236 $1 - \exp(-\frac{1}{2}d^2/\sigma_c^2)$ where $\sigma_c = 0.25\text{m}$ and d^2 is the squared Euclidean distance between the pendu-
 237 lum’s end point (x_p, y_p) and its goal $(0, l)$. I.e. $d^2 = x_p^2 + (l - y_p)^2 = (x_c - l \sin \theta)^2 + (l - l \cos \theta)^2$.

We compare four algorithms: 1) PILCO by Deisenroth and Rasmussen [3] as a baseline (unfiltered execution, and unfiltered full-prediction); 2) the method by Dallaire et al. [1] (filtered execution, and filtered MAP-prediction); 3) the method by Deisenroth and Peters [2] (filtered execution, and unfiltered full-prediction); and lastly 4) our method (filtered execution, and filtered full-prediction). For clear comparison we opted for a tightly controlled experiment. We control for data and dynamics models, i.e. each algorithm has access to the exact same data and exact same dynamics model. The reason is to eliminate variance in performance caused by different algorithms choosing different actions. We generate a single dataset by running the baseline PILCO algorithm for 11 episodes (totalling 22 seconds of system interaction). The independent variables of our experiment are 1) the method of system prediction and 2) the method of system execution. We then optimise each policy from the same initialisation using their respective prediction methods. Finally, we measure and compare their performances in both prediction and execution.

6 Results and Analysis

6.1 Predictive Performance

We now compare algorithm performance, both predictive (Figure 4) and empirical (Figure 5). First, we analyse predictive costs per timestep (Figure 4). Since predictions are probabilistic, the costs have distributions, with the exception of Dallaire et al. [1] which predicts MAP trajectories and therefore has deterministic cost. Even though we plot distributed costs, policies are optimised w.r.t. expected total cost only. Using the same dynamics, the different prediction methods optimise different policies (with the exception of Deisenroth and Rasmussen [3] and Deisenroth and Peters [2], whose prediction methods are identical). During the first 10 timesteps, we note identical performance with maximum cost due to the non-zero time required physically swing the pendulum up near the goal. Performances thereafter diverge. Since we predict w.r.t. a filtering process, less noise is predicted to be injected into the policy, and the optimiser can thus afford higher gain parameters w.r.t. the pole at balance point. If we linearise our policy around the goal point, our policy has a gain of -81.7N/rad w.r.t. pendulum angle, a larger-magnitude than both Deisenroth method gains of -39.1N/rad (negative values refer to *left* forces in Figure 3). Being afforded higher gains our policy is more reactive and more likely to catch a falling pendulum. Finally, we note Dallaire et al. [1] predict very high performance. Without balancing the costs across multiple possible trajectories, the method instead optimises a sequence of deterministic states to near perfection.

6.2 Empirical Performance

To compare the predictive results against the empirical, we used 100 executions of each algorithm (Figure 5). First, we notice a stark difference between predictive and executed performances from Dallaire et al. [1], due to neglecting model uncertainty, suffering model bias. In contrast, the other methods consider uncertainty and have relatively unbiased predictions, judging by the similarity between predictive-vs-empirical performances. Deisenroth’s methods, which differ only in execution, illustrate that filtering during execution-only can be better than no filtering at all. However, the major benefit comes when the policy is evaluated from multi-step predictions of a filtered system. Opposed to Deisenroth and Peters [2], our method’s predictions reflect reality closer because we both predict and execute system trajectories using closed loop filtering control.

To test statistical significance of empirical cost differences given 100 executions, we use a Wilcoxon rank-sum test at each time step. Excluding time steps ranging $t = [0, 29]$ (whose costs are similar), the minimum z -score over timesteps $t = [30, 60]$ that our method has superior average-cost than each other methods follows: Deisenroth 2011 $\min(z) = 4.99$, Dallaire 2009’s $\min(z) = 8.08$, Deisenroth 2012’s $\min(z) = 3.51$. Since the minimum $\min(z) = 3.51$, we have $p > 99.9\%$ certainty our method’s average empirical cost is superior than each other method.

6.3 Training Time Complexity

Training the GP dynamics model involved $N = 660$ data, $M = 50$ inducing points under a sparse GP FITC $P = 100$ policy RBF centroids, $D = 4$ state dimensions, $F = 1$ action dimensions, and $T = 60$ timestep horizon. To train the dynamics model scales $\mathcal{O}(DNM^2)$. Policy optimisation (with 300 steps, each of which require trajectory prediction with gradients) is the most intense part: our method and both Deisenroth’s methods scale $\mathcal{O}(M^2D^2(D+F)^2T + P^2D^2F^2T)$, whilst Dallaire’s only scales $\mathcal{O}(MD(D+F)T + PDFT)$. Worst case we require $M = \mathcal{O}(\exp(D+F))$ inducing

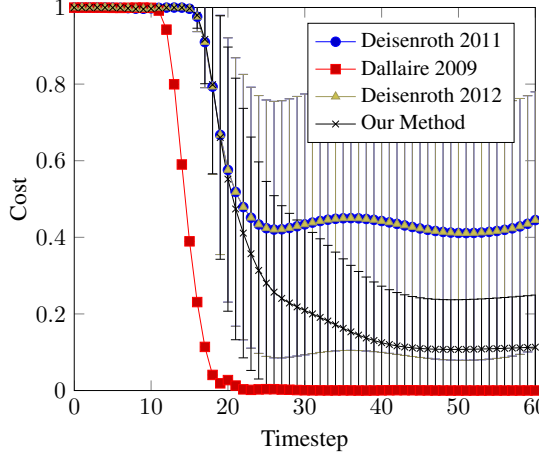


Figure 4: **Predictive costs per timestep.** The error bars show ± 1 standard deviation. Each algorithm has access to the same data set (generated by baseline Deisenroth 2011) and dynamics model. Algorithms differ in their multi-step prediction methods (except Deisenroth’s algorithms whose predictions overlap).

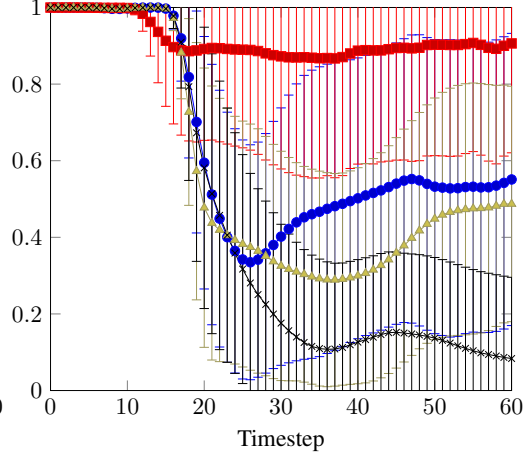


Figure 5: **Empirical costs per timestep.** We generate empirical cost distributions from 100 executions per algorithm. Error bars show ± 1 standard deviation. The plot colours and shapes correspond to the legend in Figure 4.

points to capture dynamics, the average case is unknown. Total training time was four hours to train the original PILCO method with an additional one hour to re-optimize the policy.

7 Conclusion and Future Work

In this paper, we extended the original PILCO algorithm [3] to filter observations, both during system execution and multi-step probabilistic prediction required for policy evaluation. The extended framework enables learning in *partially-observed* environments (POMDPs) whilst retaining PILCO’s data-efficiency property. We demonstrated successful application to a benchmark control problem, the noisily-observed cartpole swing-up. Our algorithm learned a good policy under significant observation noise in less than 30 seconds of system interaction. Importantly, our algorithm evaluates policies with predictions that are faithful to reality. We predict w.r.t. closed loop filtered control precisely because we execute closed loop filtered control, unlike other RL algorithms.

We showed experimentally that *faithful* and *probabilistic* predictions give greater performance gains than otherwise. For clear comparison we constrained each algorithm to use the same dynamics dataset rather than each interacting with the system to generate their own. Thus, we cannot currently claim superior data-collection abilities of our method, only superior data-usage. In future work we wish to relax this experimental constraint, to test a difference in data-collection abilities. However, the extra variance in empirical performance (caused by selection of different data) means a much larger number of experiments will be required to detect if such *additional* performance gains exist.

Several more challenges remain for future work. Firstly the assumption of zero variance of the belief-variance could be relaxed. A relaxation allows distributed trajectories to more accurately consider belief states having various degrees of certainty (belief-variance). E.g. system trajectories have larger belief-variance when passing through data-sparse regions of state-space, and smaller belief-variance in data-dense regions. Secondly, the policy could be a function of the full belief distribution (mean and variance) rather than just the mean. Such flexibility could help the policy make more ‘cautious’ actions when more uncertain about the state. Thirdly, the framework could be extended to active learning. Currently, the framework is a passive learner, greedily optimising the total cost-means and ignoring cost-variance information which could otherwise better inform exploration, increasing data-efficiency further.

References

- [1] Dallaire, P., Besse, C., Ross, S., and Chaib-Draa, B. (2009). Bayesian reinforcement learning in continuous POMDPs with Gaussian processes. In *International Conference on Intelligent Robots and Systems (IROS 2009)*, pages 2604–2609. IEEE.
- [2] Deisenroth, M. and Peters, J. (2012). Solving nonlinear continuous state-action-observation POMDPs for mechanical systems with Gaussian noise. In *European Workshop on Reinforcement Learning (EWRL 2012)*.
- [3] Deisenroth, M. and Rasmussen, C. (2011). PILCO: A model-based and data-efficient approach to policy search. In *International Conference on Machine Learning (ICML 2011)*, pages 465–472, New York, NY, USA.
- [4] Duff, M. (2002). *Optimal Learning: Computational procedures for Bayes-adaptive Markov decision processes*. PhD thesis, Department of Computer Science, University of Massachusetts Amherst.
- [5] Ko, J. and Fox, D. (2009). GP-BayesFilters: Bayesian filtering using Gaussian process prediction and observation models. *Autonomous Robots*, 27(1):75–90.
- [6] McHutchon, A. (2014). *Nonlinear modelling and control using Gaussian processes*. PhD thesis, Department of Engineering, University of Cambridge.
- [7] Poupart, P., Vlassis, N., Hoey, J., and Regan, K. (2006). An analytic solution to discrete Bayesian reinforcement learning. *Proceedings of the 23rd international conference on Machine learning*, pages 697–704.
- [8] Ross, S., Chaib-Draa, B., and Pineau, J. (2008a). Bayesian reinforcement learning in continuous POMDPs with application to robot navigation. In *International Conference on Robotics and Automation (ICRA 2008)*, pages 2845–2851. IEEE.
- [9] Ross, S., Chaib-Draa, B., and Pineau, J. (2008b). Bayesian reinforcement learning in continuous POMDPs with application to robot navigation. In *Robotics and Automation, 2008. ICRA 2008. IEEE International Conference on*, pages 2845–2851. IEEE.
- [10] Van Den Berg, J., Patil, S., and Alterovitz, R. (2012). Efficient approximate value iteration for continuous gaussian pomdps. In *AAAI*.
- [11] Webb, D. J., Crandall, K. L., and van den Berg, J. (2014). Online parameter estimation via real-time replanning of continuous gaussian pomdps. In *Robotics and Automation (ICRA), 2014 IEEE International Conference on*, pages 5998–6005. IEEE.

350 A Gradients for Policy Improvement

351 Let $\text{vec}(\cdot)$ be the ‘unwrap operator’ that reshapes matrices columnwise into vectors. We define a
 352 Markov filtered-system from the belief’s parameters: $S_t = [M_{t|t-1}^\top, \text{vec}(V_{t|t-1})^\top]^\top$. To predict
 353 system evolution, the state distribution is defined:

$$p(S_t) \sim \mathcal{N}\left(\mu_t^s = \begin{bmatrix} \mu_{t|t-1}^m \\ \text{vec}(V_{t+1|t}) \end{bmatrix}, \Sigma_t^s = \begin{bmatrix} \Sigma_{t|t-1}^m & 0 \\ 0 & 0 \end{bmatrix}\right). \quad (9)$$

354 To compute policy gradient $dJ/d\psi$ we first require $d\mathcal{E}_t/d\psi$:

$$\begin{aligned} \frac{d\mathcal{E}_t}{d\theta} &= \frac{d\mathcal{E}_t}{dp(S_t)} \frac{dp(S_t)}{d\theta} \\ &= \frac{d\mathcal{E}_t}{\partial \mu_t^s} \frac{\partial \mu_t^s}{d\theta} + \frac{d\mathcal{E}_t}{\partial \Sigma_t^s} \frac{\partial \Sigma_t^s}{d\theta}, \quad \text{and} \end{aligned} \quad (10)$$

$$\frac{dp(S_{t+1})}{d\theta} = \frac{dp(S_{t+1})}{dp(S_t)} \frac{dp(S_t)}{d\theta} + \frac{\partial p(S_{t+1})}{\partial \theta}. \quad (11)$$

355 Application of the chain rule backwards from the state distribution at the horizon S_T , to S_t at arbitrary
 356 time t , is analogous to that detailed in PILCO [3], where we use S_t , μ_t^s and Σ_t^s in the place of x_t , μ_t
 357 and Σ_t respectively.

358 B Identities for Gaussian Process Prediction with Hierarchical Uncertain In- 359 puts

360 The two functions

$$\begin{aligned} q(x, x', \Lambda, V) &\triangleq |\Lambda^{-1}V + I|^{-1/2} \exp\left(-\frac{1}{2}(x - x')[\Lambda + V]^{-1}(x - x')\right), \\ Q(x, x', \Lambda_a, \Lambda_b, V, \mu, \Sigma) &\triangleq c_1 \exp\left(-\frac{1}{2}(x - x')^\top [\Lambda_a + \Lambda_b + 2V]^{-1}(x - x')\right) \\ &\quad \times \exp\left(-\frac{1}{2}(z - \mu)^\top [((\Lambda_a + V)^{-1} + (\Lambda_b + V)^{-1})^{-1} + \Sigma]^{-1}(z - \mu)\right), \\ &= c_2 q(x, \mu, \Lambda_a, V) q(\mu, x', \Lambda_b, V) \\ &\quad \times \exp\left(\frac{1}{2}\mathbf{r}^\top [(\Lambda_a + V)^{-1} + (\Lambda_b + V)^{-1} + \Sigma^{-1}]^{-1}\mathbf{r}\right), \end{aligned}$$

$$\text{where } \begin{cases} z = (\Lambda_b + V)(\Lambda_a + \Lambda_b + 2V)^{-1}x + (\Lambda_a + V)(\Lambda_a + \Lambda_b + 2V)^{-1}x' \\ \mathbf{r} = (\Lambda_a + V)^{-1}(x - \mu) + (\Lambda_b + V)^{-1}(x' - \mu) \\ c_1 = |(\Lambda_a + V)(\Lambda_b + V) + (\Lambda_a + \Lambda_b + 2V)\Sigma|^{-1/2} |\Lambda_a \Lambda_b|^{1/2} \\ c_2 = |((\Lambda_a + V)^{-1} + (\Lambda_b + V)^{-1})\Sigma + I|^{-1/2}, \end{cases} \quad (12)$$

361 have the following Gaussian integrals

$$\begin{aligned} \int q(x, t, \Lambda, V) \mathcal{N}(t|\mu, \Sigma) dt &= q(x, \mu, \Lambda, \Sigma + V), \\ \int q(x, t, \Lambda_a, V) q(t, x', \Lambda_b, V) \mathcal{N}(t|\mu, \Sigma) dt &= Q(x, x', \Lambda_a, \Lambda_b, V, \mu, \Sigma), \\ \int Q(x, x', \Lambda_a, \Lambda_b, 0, \mu, V) \mathcal{N}(\mu|\mathbf{m}, \Sigma) d\mu &= Q(x, x', \Lambda_a, \Lambda_b, 0, \mathbf{m}, \Sigma + V). \end{aligned} \quad (13)$$

362 We want to model data with E output coordinates, and use separate combinations of linear models
 363 and GPs to make predictions, $a = 1, \dots, E$:

$$f_a(x^*) = f_a^* \sim \mathcal{N}(\theta_a^\top x^* + k_a(x^*, \mathbf{x})\beta_a, k_a(x^*, x^*) - k_a(x^*, \mathbf{x})(K_a + \Sigma_\varepsilon^a)^{-1}k_a(\mathbf{x}, x^*)),$$

364 where the E squared exponential covariance functions are

$$k_a(x, x') = s_a^2 q(x, x', \Lambda_a, 0), \quad \text{where } a = 1, \dots, E, \quad (14)$$

365 and s_a^2 are the signal variances and Λ_a is a diagonal matrix of squared length scales for GP number
 366 a . The noise variances are Σ_ε^a . The inputs are \mathbf{x} and the outputs y_a and we define $\beta_a = (K_a +$
 367 $\Sigma_\varepsilon^a)^{-1}(y_a - \theta_a^\top \mathbf{x})$, where K_a is the Gram matrix.

368 **B.1 Derivatives**

369 For symmetric Λ and V and Σ :

$$\begin{aligned}\frac{\partial \ln q(x, x', \Lambda, V)}{\partial x} &= -(\Lambda + V)^{-1}(x - x') = -(\Lambda^{-1}V + I)^{-1}\Lambda^{-1}(x - x') \\ \frac{\partial \ln q(x, x', \Lambda, V)}{\partial x'} &= (\Lambda + V)^{-1}(x - x') \\ \frac{\partial \ln q(x, x', \Lambda, V)}{\partial V} &= -\frac{1}{2}(\Lambda + V)^{-1} + \frac{1}{2}(\Lambda + V)^{-1}(x - x')(x - x')^\top (\Lambda + V)^{-1}\end{aligned}\tag{15}$$

370 Let $L = (\Lambda_a + V)^{-1} + (\Lambda_b + V)^{-1}$, $R = \Sigma L + I$, $Y = R^{-1}\Sigma = [L + \Sigma^{-1}]^{-1}$, $T : X \rightarrow XX^\top$:

$$\partial Q(x, x', \Lambda_a, \Lambda_b, V, \mu, \Sigma) = Q \circ \partial \left(\ln c_2 + \ln q(x, \mu, \Lambda_a, V) + \ln q(\mu, x', \Lambda_b, V) + \frac{1}{2}y^\top Y y \right)$$

$$\frac{1}{2} \frac{\partial y^\top Y y}{\partial \mu} = y^\top Y \frac{\partial y}{\partial \mu} = -y^\top Y L$$

$$\frac{\partial \ln c_2}{\partial \Sigma} = -\frac{1}{2} \frac{\partial \ln |L\Sigma + I|}{\partial \Sigma} = -\frac{1}{2} L^\top (L\Sigma + I)^{-\top} = -\frac{1}{2} L R^{-1}$$

$$\frac{\partial y^\top Y y}{\partial \Sigma} = \Sigma^{-\top} Y^\top y y^\top Y^\top \Sigma^{-\top} = T(R^{-\top} y)$$

$$\frac{\partial \ln c_2}{\partial V} = -\frac{1}{2} \frac{\partial \ln |L\Sigma + I|}{\partial V} = -\frac{1}{2} \frac{\partial \ln |\sum_i [(\Lambda_i + V)^{-1}] \Sigma + I|}{\partial V}$$

$$= \frac{1}{2} \sum_i \left[(\Lambda_i + V)^{-\top} \left(\sum_j [(\Lambda_j + V)^{-1}] \Sigma + I \right)^{-\top} \Sigma^\top (\Lambda_i + V)^{-\top} \right]$$

$$= \frac{1}{2} \sum_i \left[(\Lambda_i + V)^{-1} Y (\Lambda_i + V)^{-1} \right]$$

$$\frac{\partial y^\top Y y}{\partial V} = y^\top \frac{\partial Y}{\partial V} y + \frac{\partial y^\top}{\partial V} Y y + y^\top Y \frac{\partial y}{\partial V} = \sum_i \left[(\Lambda_i + V)^{-1} Y^\top y y^\top Y^\top (\Lambda_i + V)^{-1} \right]$$

$$- \sum_i \left[(\Lambda_i + V)^{-1} (x_{n_i} - \mu) (Y y)^\top (\Lambda_i + V)^{-1} \right]$$

$$- \sum_i \left[(\Lambda_i + V)^{-1} (y^\top Y)^\top (x_{n_i} - \mu)^\top (\Lambda_i + V)^{-1} \right]$$

$$= \sum_i \left[T \left((\Lambda_i + V)^{-1} (Y y - (x_{n_i} - \mu)) \right) - T \left((\Lambda_i + V)^{-1} (x_{n_i} - \mu) \right) \right]$$

(16)

C Dynamics Predictions in System-Execution

Here we specify the predictive distribution $p(b_{t+1|t})$, whose moments are equal to the moments from dynamics model output f with uncertain input $\tilde{b}_{t|t} \sim \mathcal{N}(\tilde{m}_{t|t}, \tilde{V}_{t|t})$ similar to Deisenroth and Rasmussen [3]. Consider making predictions from $a = 1, \dots, E$ GPs at $\tilde{b}_{t|t}$ with specification $p(\tilde{b}_{t|t}) \sim \mathcal{N}(\tilde{m}_{t|t}, \tilde{V}_{t|t})$. We have the following expressions for the predictive mean, variances and input-output covariances using the law of iterated expectations and variances:

$$b_{t+1|t} \sim \mathcal{N}(m_{t+1|t}, V_{t+1|t}), \quad (17)$$

$$\begin{aligned} m_{t+1|t}^a &= \mathbb{E}_{\tilde{b}_{t|t}}[f^a(\tilde{b}_{t|t})] \\ &= \int (s_a^2 \beta_a^\top q(x_i, \tilde{b}_{t|t}, \Lambda_a, 0) + \phi_a^\top \tilde{b}_{t|t}) \mathcal{N}(\tilde{b}_{t|t}; \tilde{m}_{t|t}, \tilde{V}_{t|t}) d\tilde{b}_{t|t} \\ &= s_a^2 \beta_a^\top q^a + \phi_a^\top \tilde{m}_{t|t}, \end{aligned} \quad (18)$$

$$\begin{aligned} C_a &\doteq \tilde{V}_{t|t}^{-1} \mathbb{C}_{\tilde{b}_{t|t}}[\tilde{b}_{t|t}, f^a(\tilde{b}_{t|t}) - \phi_a^\top \tilde{b}_{t|t}], \\ &= \tilde{V}_{t|t}^{-1} \int (\tilde{b}_{t|t} - \tilde{m}_{t|t}) s_a^2 \beta_a^\top q(x, \tilde{b}_{t|t}, \Lambda_a, 0) \mathcal{N}(\tilde{b}_{t|t}; \tilde{m}_{t|t}, \tilde{V}_{t|t}) d\tilde{b}_{t|t} \\ &= s_a^2 (\Lambda_a + \tilde{V}_{t|t})^{-1} (x - \tilde{m}_{t|t}) \beta_a q^a, \end{aligned} \quad (19)$$

$$\begin{aligned} V_{t+1|t}^{ab} &= \mathbb{C}_{\tilde{b}_{t|t}}[f^a(\tilde{b}_{t|t}), f^b(\tilde{b}_{t|t})] \\ &= \mathbb{C}_{\tilde{b}_{t|t}}[\mathbb{E}_f[f^a(\tilde{b}_{t|t})], \mathbb{E}_f[f^b(\tilde{b}_{t|t})]] + \mathbb{E}_{\tilde{b}_{t|t}}[\mathbb{C}_f[f^a(\tilde{b}_{t|t}), f^b(\tilde{b}_{t|t})]] \\ &= \mathbb{C}_{\tilde{b}_{t|t}}[s_a^2 \beta_a^\top q(x, \tilde{b}_{t|t}, \Lambda_a, 0) + \phi_a^\top \tilde{b}_{t|t}, s_b^2 \beta_b^\top q(x, \tilde{b}_{t|t}, \Lambda_b, 0) + \phi_b^\top \tilde{b}_{t|t}] + \\ &\quad \delta_{ab} \mathbb{E}[s_a^2 - k_a(\tilde{b}_{t|t}, x)(K_a + \Sigma_\varepsilon^a)^{-1} k_a(x, \tilde{b}_{t|t})] \\ &= s_a^2 s_b^2 [\beta_a^\top (Q^{ab} - q^a q^{b\top}) \beta_b + \\ &\quad \delta_{ab} (s_a^{-2} - \text{tr}((K_a + \Sigma_\varepsilon^a)^{-1} Q^{aa}))] + C_a^\top \tilde{V}_{t|t} \phi_b + \phi_a^\top \tilde{V}_{t|t} C_b + \phi_a^\top \tilde{V}_{t|t} \phi_b, \end{aligned} \quad (20)$$

where

$$\begin{aligned} q_i^a &= q(x_i, \tilde{m}_{t|t}, \Lambda_a, \tilde{V}_{t|t}), \\ Q_{ij}^{ab} &= Q(x_i, x_j, \Lambda_a, \Lambda_b, 0, \tilde{m}_{t|t}, \tilde{V}_{t|t}), \\ \beta_a &= (K_a + \Sigma^{\varepsilon, a})^{-1} (y_a - \phi_a^\top x), \end{aligned}$$

and training inputs are x , outputs are y_a (determined by the ‘Direct method’), K_a is a Gram matrix.

D Dynamics Predictions in System-Prediction

Here we describe the prediction formulae for the random belief state in system-prediction. We again note, during execution, our belief distribution is specified by certain parameters, $b_{t|t} \sim \mathcal{N}(m_{t|t}, V_{t|t})$. By contrast, during system prediction, our belief distribution is specified by an uncertain belief-mean and certain belief-variance: $B_{t|t} \sim \mathcal{N}(M_{t|t}, V_{t|t}) \sim \mathcal{N}(\mathcal{N}(\mu_{t|t}^m, \Sigma_{t|t}^m), \tilde{V}_{t|t})$, where we assumed a delta distribution on $\tilde{V}_{t|t}$ for mathematical simplicity, i.e. $\text{vec}(\tilde{V}_{t|t}) \sim \mathcal{N}(\text{vec}(\tilde{V}_{t|t}), 0)$. Therefore we conduct GP prediction given hierarchically-uncertain inputs, outlining each output moment below. I.e. consider making predictions from $a = 1, \dots, E$ GPs at $B_{t|t}$ with *hierarchical* specification

$$p(B_{t|t}) \sim \mathcal{N}(M_{t|t}, V_{t|t}), \text{ and } M_{t|t} \sim \mathcal{N}(\mu_{t|t}^m, \Sigma_{t|t}^m), \quad (21)$$

or equivalently the joint

$$p\left(\begin{bmatrix} B_{t|t} \\ M_{t|t} \end{bmatrix}\right) \sim \mathcal{N}\left(\begin{bmatrix} \mu_{t|t}^m \\ \mu_{t|t}^m \end{bmatrix}, \begin{bmatrix} \Sigma_{t|t}^m + \tilde{V}_{t|t} & \Sigma_{t|t}^m \\ \Sigma_{t|t}^m & \Sigma_{t|t}^m \end{bmatrix}\right). \quad (22)$$

388 **Mean of the Belief-Mean:** dynamics prediction uses input $\tilde{M}_{t|t} \sim \mathcal{N}(\mu_{t|t}^{\tilde{m}}, \Sigma_{t|t}^{\tilde{m}})$, which is jointly
 389 distributed according to Eq. 7. Using the belief-mean $m_{t+1|t}^a$ definition (Eq. 18),

$$\begin{aligned}\mu_{t+1|t}^{m,a} &= \mathbb{E}_{\tilde{M}_{t|t}}[M_{t+1|t}^a] \\ &= \int M_{t+1|t}^a \mathcal{N}(\tilde{M}_{t|t}; \mu_{t|t}^{\tilde{m}}, \Sigma_{t|t}^{\tilde{m}}) d\tilde{M}_{t|t}, \\ &= s_a^2 \beta_a^\top \int q(x, \tilde{M}_{t|t}, \Lambda_a, V) \mathcal{N}(\tilde{M}_{t|t}; \mu_{t|t}^{\tilde{m}}, \Sigma_{t|t}^{\tilde{m}}) d\tilde{M}_{t|t} + \phi_a^\top \mu_{t|t}^{\tilde{m}} \\ &= s_a^2 \beta_a^\top \hat{q}^a + \phi_a^\top \mu_{t|t}^{\tilde{m}},\end{aligned}\tag{23}$$

$$\hat{q}_i^a = q(x_i, \mu_{t|t}^{\tilde{m}}, \Lambda_a, \Sigma_{t|t}^{\tilde{m}} + \tilde{V}_{t|t}).\tag{24}$$

390 **Input-Output Covariance:** the expected input-output covariance belief term (Eq. 19) (equivalent
 391 to the input-output covariance of the belief-mean) is:

$$\begin{aligned}\hat{C}_a &\doteq \tilde{V}_{t|t}^{-1} \mathbb{E}_{\tilde{M}_{t|t}}[\mathbb{C}_{B_{t|t}}[\tilde{B}_{t|t}, f(\tilde{B}_{t|t}) - \phi_a^\top \tilde{M}_{t|t}]], \text{ and similarly defined} \\ &\doteq (\Sigma_{t|t}^{\tilde{m}})^{-1} \mathbb{C}_{\tilde{M}_{t|t}}[\tilde{M}_{t|t}, \mathbb{E}_{B_{t|t}}[f(\tilde{B}_{t|t}) - \phi_a^\top \tilde{M}_{t|t}]], \\ &= (\Sigma_{t|t}^{\tilde{m}})^{-1} \int (\tilde{M}_{t|t} - \mu_{t|t}^{\tilde{m}}) \mathbb{E}_{B_{t|t}}[f(\tilde{B}_{t|t})] \mathcal{N}(\tilde{M}_{t|t}; \mu_{t|t}^{\tilde{m}}, \Sigma_{t|t}^{\tilde{m}}) d\tilde{M}_{t|t} \\ &= (\Sigma_{t|t}^{\tilde{m}})^{-1} \int (\tilde{M}_{t|t} - \mu_{t|t}^{\tilde{m}}) (s_a^2 \beta_a^\top q(x_i, \tilde{M}_{t|t}, \Lambda_a, \tilde{V}_{t|t})) \mathcal{N}(\tilde{M}_{t|t}; \mu_{t|t}^{\tilde{m}}, \Sigma_{t|t}^{\tilde{m}}) d\tilde{M}_{t|t} \\ &= s_a^2 (\Lambda_a + \Sigma_{t|t}^{\tilde{m}} + \tilde{V}_{t|t})^{-1} (x - \mu_{t|t}^{\tilde{m}}) \beta_a \hat{q}_i^a.\end{aligned}\tag{25}$$

392 **Variance of the Belief-Mean:** the variance of randomised belief-mean (Eq 18) is:

$$\begin{aligned}\Sigma_{t+1|t}^{m,ab} &= \mathbb{C}_{\tilde{M}_{t|t}}[M_{t+1|t}^a, M_{t+1|t}^b], \\ &= \int M_{t+1|t}^a M_{t+1|t}^b \mathcal{N}(\tilde{M}_{t|t}; \mu_{t|t}^{\tilde{m}}, \Sigma_{t|t}^{\tilde{m}}) d\tilde{M}_{t|t} - \mu_{m_{t+1|t}}^a \mu_{m_{t+1|t}}^b, \\ &= s_a^2 s_b^2 \beta_a^\top (\hat{Q}^{ab} - \hat{q}^a \hat{q}^{b\top}) \beta_b + \hat{C}_a^\top \Sigma_{t|t}^{\tilde{m}} \phi_b + \phi_a^\top \Sigma_{t|t}^{\tilde{m}} \hat{C}_b + \phi_a^\top \Sigma_{t|t}^{\tilde{m}} \phi_b,\end{aligned}\tag{26}$$

$$\hat{Q}_{ij}^{ab} = Q(x_i, x_j, \Lambda_a, \Lambda_b, \tilde{V}_{t|t}, \mu_{t|t}^{\tilde{m}}, \Sigma_{t|t}^{\tilde{m}}).\tag{27}$$

393 **Mean of the Belief-Variance:** using the belief-variance $V_{t+1|t}^{ab}$ definition (Eq. 20),

$$\begin{aligned}\tilde{V}_{t+1|t}^{ab} &= \mathbb{E}_{\tilde{M}_{t|t}}[V_{t+1|t}^{ab}] \\ &= \int V_{t+1|t}^{ab} \mathcal{N}(\tilde{M}_{t|t}; \mu_{t|t}^{\tilde{m}}, \Sigma_{t|t}^{\tilde{m}}) d\tilde{M}_{t|t} \\ &= s_a^2 s_b^2 [\beta_a^\top (\tilde{Q}^{ab} - \hat{Q}^{ab}) \beta_b + \delta_{ab} (s_a^{-2} - \text{tr}((K_a + \Sigma_\varepsilon^a)^{-1} \tilde{Q}^{aa}))] \\ &\quad + \hat{C}_a^\top \tilde{V}_{t|t} \phi_b + \phi_a^\top \tilde{V}_{t|t} \hat{C}_b + \phi_a^\top \tilde{V}_{t|t} \phi_b,\end{aligned}\tag{28}$$

$$\tilde{Q}_{ij}^{ab} = Q(x_i, x_j, \Lambda_a, \Lambda_b, 0, \mu_{t|t}^{\tilde{m}}, \Sigma_{t|t}^{\tilde{m}} + \tilde{V}_{t|t}).\tag{29}$$