# Bayesian Meta-Reinforcement Learning with Laplace Variational Recurrent Networks

**Anonymous authors**
Paper under double-blind review

**Keywords:** Variational Inference, Bayesian Reinforcement Learning, Meta-Reinforcement Learning, Uncertainty Estimation

## Summary

Meta-reinforcement learning trains a single reinforcement learning agent on a distribution of tasks to quickly generalize to new tasks outside of the training set at test time. From a Bayesian perspective, one can interpret this as performing amortized variational inference on the posterior distribution over training tasks. Among the various meta-reinforcement learning approaches, a common method is to represent this distribution with a point-estimate using a recurrent neural network. We show how one can augment this point estimate to give full distributions through the Laplace approximation, either at the start of, during, or after learning, without modifying the base model architecture. With our approximation, we are able to estimate distribution statistics (e.g., the entropy) of non-Bayesian agents and observe that point-estimate based methods produce overconfident estimators while not satisfying consistency. Furthermore, when comparing our approach to full-distribution based learning of the task posterior, our method performs on par with variational baselines while having much fewer parameters.

## Contribution(s)

1. We formulate a probabilistic graphical model to match the practical design of memory-based meta-reinforcement learning agents, in order to perform uncertainty quantification through the Laplace approximation *without* retraining or architecture modifications.
   **Context:** Ours is an extension of the variational recurrent neural networks by Chung et al. (2015), *maximum a posteriori* policy optimization by Abdolmaleki et al. (2018), and adopts the control-as-inference framework (Levine, 2018).

2. We investigate how different assumptions on the posterior model over Markov decision processes interact with representation learning and uncertainty quantification of the recurrent neural network.
   **Context:** The agents trained with a recurrent neural network are non-Bayesian agents to which we try to apply a Bayesian approximation. Although we obtain a method for quantifying uncertainty in their learned representation, there is still a degree of misspecification.

3. When used as an alternative to the baseline variational recurrent network, we show that our method either matches or improves performance.
   **Context:** This shows that our probabilistic formulation provides an alternative approximation for variational online learning while using fewer learnable parameters and without needing architecture modifications.

4. Our results show that the recurrent neural network representations learned by non-Bayesian meta-reinforcement learning agents, judging over multiple assumptions on the graphical model, produces overconfident estimators.
   **Context:** This extends prior insight by Xiong et al. (2021) that the representations of memory-based meta-reinforcement learning agents learn *inconsistent* estimators.

# Bayesian Meta-Reinforcement Learning with Laplace Variational Recurrent Networks

**Anonymous authors**
Paper under double-blind review

## Abstract

Meta-reinforcement learning trains a single reinforcement learning agent on a distribution of tasks to quickly generalize to new tasks outside of the training set at test time. From a Bayesian perspective, one can interpret this as performing amortized variational inference on the posterior distribution over training tasks. Among the various meta-reinforcement learning approaches, a common method is to represent this distribution with a point-estimate using a recurrent neural network. We show how one can augment this point estimate to give full distributions through the Laplace approximation, either at the start of, during, or after learning, without modifying the base model architecture. With our approximation, we are able to estimate distribution statistics (e.g., the entropy) of non-Bayesian agents and observe that point-estimate based methods produce overconfident estimators while not satisfying consistency. Furthermore, when comparing our approach to full-distribution based learning of the task posterior, our method performs on par with variational baselines while having much fewer parameters.

## 1   Introduction

Reinforcement Learning (RL) concerns itself with making optimal decisions from data (Sutton & Barto, 2018). This is typically achieved by letting an agent generate data in an environment and then optimizing a cost function of the agent's parameters given this data. In meta-RL, an agent is trained to optimize an expected cost over a prior distribution of environments (Finn et al., 2017; Chen et al., 2017; Beck et al., 2023). The idea is then that, given a trajectory of new data, an agent can infer latent environment parameters and successfully adapt its action policy online. This is known as zero-shot or few-shot adaptation or learning (Beck et al., 2023). In recent years, this paradigm has shown impressive results, for example, by the Capture the Flag agent (Jaderberg et al., 2019) or the Adaptive Agent (Bauer et al., 2023).

In any meta-RL algorithm, an accurate approximation of the latent parameter distribution given data, also known as the task posterior distribution, is useful to quantify the agent's uncertainty (Grant et al., 2018). Accurate quantification of uncertainty enables agents to detect distribution shifts (Daxberger et al., 2021) or guide exploration through novelty signals (Osband et al., 2013; Sekar et al., 2020). Importantly, on deployment, distribution shift detection is essential for timely human intervention or retraining. This ultimately improves the robustness and efficacy of our algorithms and allows us to more reliably inspect failure cases.

A common approach in meta-RL is to model the task posterior with point estimates, e.g., using the hidden state of recurrent neural networks (RNN) (Chen et al., 2017). However, this prevents us from exploiting useful distributional statistics. Another downside of using point-estimates is the increased risk of overconfidence unless the true posterior is sharply peaked at that particular point. This would imply that there exists almost no uncertainty about the environment, which is typically a strong and unrealistic assumption. As a consequence, point-estimate based meta-RL has been known to overfit

to its training distribution, leading to brittle downstream performance (Xiong et al., 2021; Greenberg et al., 2023).

Arguably, a better approach would be to explicitly parameterize some full distribution (e.g., a Gaussian) (Chung et al., 2015; Zintgraf et al., 2021). However, approximate Bayesian methods are slower to train and are still often outperformed by simple point-estimate methods in terms of expected returns (Greenberg et al., 2023) even though they model the posterior more accurately. This could be explained by the fact that non-Bayesian methods enjoy reduced sampling noise, easier numerical representation, and improved model capacity by not having to learn a complex posterior model (Goyal et al., 2017; Hafner et al., 2019).

To get benefits from Bayesian methods when using non-Bayesian models, we introduce the *Laplace variational recurrent neural network* (Laplace VRNN) which utilizes the Laplace approximation to extend RNN-based meta-reinforcement learning. Our method can perform uncertainty quantification for non-Bayesian meta-RL agents without modifying the model architecture or loss function, and without needing to retrain any parameters. In other words, the consequence of the Laplace approximation is that we can apply it at any point during model training. When applied after training, this is often referred to as a *post-hoc* posterior (Daxberger et al., 2021). This allows us to make use of deterministic pre-training schedules and benefit from their aforementioned advantages while also enjoying the benefits of Bayesian methods. Although the Laplace approximation has already been explored in meta-RL (Grant et al., 2018; Finn et al., 2018), it has not been applied in memory-based methods (Duan et al., 2016; Beck et al., 2023) which is what we explore.

The Laplace approximation is a simple method that only requires the *curvature* of a distribution's log-likelihood at a local maximum (Daxberger et al., 2021). For a Gaussian mean-field assumption on our posterior model (Bishop, 2007), we only require the Jacobian matrix of the RNN output with respect to its hidden state. This gives us a Gaussian distribution for the task posterior distribution centered at the RNN hidden state with inverse covariance equal to the sum of Jacobian outer products. This is a comparatively cheap approximation compared to typical methods that apply the Laplace approximation to the much higher dimensional neural network parameters (Grant et al., 2018; Daxberger et al., 2021; Martens & Grosse, 2015).

We empirically validate that our method can *reliably estimate posterior statistics of our non-Bayesian baselines without degrading performance* on supervised and reinforcement learning domains. Similarly to Xiong et al. (2021), our results show that non-Bayesian meta-RL agents do not learn consistent estimators, however, we also find that the learned representations are overconfident. This could be seen by inspecting the *post-hoc* posterior provided by the Laplace approximation, which showed low entropy while not converging to a stable distribution. Furthermore, when comparing our method against variational inference baselines, we find that our Laplace method performs on par in terms of mean returns. Ultimately, this shows that the Laplace approximation can complement (or serve as an alternative to) variational inference methods for uncertainty quantification, since we do not *learn* our local uncertainty but estimate this based on the model's fitted parameters.

## 2   Related Work

**Meta-Reinforcement Learning**   Meta-learning has been described from various viewpoints, ranging from contexts (Sodhani et al., 2022), latent-variable models (Garnelo et al., 2018; Wu et al., 2020; Gordon et al., 2019), amortized inference (Gershman & Goodman, 2014; Wu et al., 2020), and "learning to learn" (Beck et al., 2023; Wang et al., 2017; Hospedales et al., 2021). Applying these ideas to reinforcement learning has been gaining traction within the field recently, for example, learning maximum likelihood estimation algorithms (like our work) (Andrychowicz et al., 2016; Garnelo et al., 2018), probability density functions (Lu et al., 2022; Bechtle et al., 2021), or model exploration strategies (Gupta et al., 2018).

Related to our work is the neural process by Garnelo et al. (2018) which formalizes using meta-learning to infer a set of (global) latent variables for a generative distribution. Since we test on

reinforcement learning problems, one could view our model as a type of non-stationary stochastic process or sequential neural process (Øksendal, 2003; Singh et al., 2019). This makes our model and optimization objective similar to the PlaNet model (Hafner et al., 2019), however, we do not condition our recurrent model on samples from the task posterior so we can obtain an analytical solution. This choice does reduce the non-linearity of our model, the topic of linear vs. non-linear state space models is still active research (Gu & Dao, 2024).

**Bayesian Reinforcement Learning**  Learning an optimal control policy conditional on a task-posterior amounts to approximating the Bayes-adaptive optimal policy (Duff, 2002). In this framework an agent is conditioned on its current state and a history of observations, the history can then be used to produce a belief distribution over latent variables. The Bayes-adaptive optimal policy maximizes the environment returns in expectation over this belief (Duff, 2002; Ghavamzadeh et al., 2015; Zintgraf et al., 2021; Mikulik et al., 2020). Although meta-learning induces uncertainty only over the reward and transition function, many have also successfully tackled the problem as a general partially observable Markov decision process (Chen et al., 2017; Bauer et al., 2023). Doing this is orthogonal to our method, however, we focus on the Bayes-adaptive framework.

**Laplace Approximation**  The Laplace approximation has been explored for meta-learning in, for example, the model agnostic meta-learning algorithm (Finn et al., 2017; 2018) to achieve more accurate inference, or for continual learning (Kirkpatrick et al., 2017) as a regularizer for weight-updates. The main obstacle to using the Laplace approximation in practice is the computation of the inverse Hessian (MacKay, 1992), which alone has quadratic memory scaling in the number of model parameters. For this reason, in Bayesian neural networks, the block-diagonal factorization has become quite popular (Martens & Grosse, 2015), as used by TRPO (Schulman et al., 2015) or second order optimizers (Botev et al., 2017). Our method bypasses the costly Hessian problem by modeling a distribution on a small subset of the full parameter set. In contrast to doing Bayesian linear regression on the last layer of a neural network, our method can express multimodal distributions.

# 3 Preliminaries

We want to find an optimal policy $\pi$ for a sequential decision-making problem, which we formalize as an episodic Markov decision process (Sutton & Barto, 2018). We define states $S \in \mathcal{S}$, actions $A \in \mathcal{A}$, and rewards $R \in \mathbb{R}$ as random variables that we sample in sequences. We write $H^i = \{S_t, A_t, R_t\}_{t=1}^T$ to abbreviate the joint random variable of episode $i \in \mathbb{N}$,

$$p(H^i) = \prod_{t=1}^T p(R_t|S_t, A_t)\pi(A_t|S_t)p(S_t|S_{t-1}, A_{t-1}),$$

where $p(S_1|A_0, S_0) \stackrel{\Delta}{=} p(S_1)$ is the initial state distribution, $\pi(A_t|S_t)$ is the policy, $p(S_{t+1}|S_t, A_t)$ is the transition model, and $p(R_t|S_t, A_t)$ is the reward model. To avoid confusion, we denote episodes in the *superscript* from $i = 1, \ldots, n$ and time in the *subscripts* from $t = 1, \ldots, T$. For convenience, we subsume the common discount factors $\gamma \in [0, 1]$ into the transition probabilities as a global termination probability of $p_{\text{term}} = 1 - \gamma$ (Levine, 2018) assuming that the MDP will end in an absorbing state with zero rewards. The objective is to find $\pi^*$ such that $\mathbb{E}_{p(H)} \sum_t R_t$ is maximized.

## 3.1 Inference in Meta-RL

In contrast to single-task Reinforcement Learning (RL), in meta-RL we want to find the optimal policy $\pi^*$ to a distribution over different environments. The agent typically does not know which environment it is currently being deployed in and needs to adaptively switch strategies based on online feedback. We assume that the agent can adapt over *multiple* episodes $H^{1:n}$, as opposed to only one episode $H^1$, which is also known as zero-shot or few-shot adaptation (Beck et al., 2023). This approach can be formalized using the concept of global latent variables $Z$. For a fixed $\pi$,

129  each trajectory $H$ that we sample depends on a sampled latent variable $Z \sim p(Z)$, which could be
130  interpreted as a unique identifier of the current environment. Our agent does not directly observe $Z$,
131  but this variable influences the reward and transition models of the environment.

132  With full generality, we define the generative process,

$$p(H^{1:n}, Z^{1:n}) = \prod_{i=1}^{n} p(H^i | Z^i) p(Z^i | Z^{<i}, H^{<i}), \tag{1}$$

133  where the term $p(H^i | Z^i)$ indicates the sampling distribution of the environment under our current
134  model for $Z$, and the term $p(Z^i | Z^{<i}, H^{<i})$ denotes the posterior distribution over latent variables $Z$
135  given all the data we have observed so far. For brevity, we do not expand the sampling distribution
136  $p(H^i | Z^i)$ here (see Appendix A.1), however, this expression hides that the posterior model is also
137  updated inter-episodically at every $S_t^i, A_t^i, R_t^i \in H_t^i$. $Z^i$. Note that this model is fully general, and
138  is perhaps more common in continual-RL settings (Khetarpal et al., 2022), yet it reflects the model
139  factorization and inference capabilities captured by most memory-based meta-RL methods (Duan
140  et al., 2016). This generality can be both a feature and a downside; the agent can capture broad
141  environment settings, but combined with function approximation it obfuscates what the agent learns
142  and how posterior uncertainty is represented.

**Amortized Inference**   The posterior $p(Z^i | Z^{<i}, H^{<i})$ from Eq. (1) is usually intractable; a com-
144  mon approach to deal with this is to use variational inference (Bishop, 2007). This approach rep-
145  resents the posterior with another distribution $q \in \mathcal{Q}$ within some simpler model class, and then
146  chooses $q$ to maximize a lower-bound to the data marginal (see Appendix A.1), i.e.,

$$\ln p(H^{1:n}) \geq \max_{q \in \mathcal{Q}} \mathbb{E}_{q(Z^{1:n}|H^{1:n})} \sum_{i=1}^{n} \ln p(H^i | Z^i) - KL\left(q(Z^{\leq i} | H^{<i}) \| p(Z^{\leq i} | H^{<i})\right). \tag{2}$$

147  This involves a functional optimization problem to be repeatedly solved at every timestep. There-
148  fore, a more desirable approach is to amortize this with a learned neural network $f_\theta$ that maps past
149  observations and latents directly to a distribution $q \in \mathcal{Q}$ (Gershman & Goodman, 2014). In practice,
150  this can be achieved by using a parametric family for $q_\phi$ and using $f_\theta$ to predict the parameters
151  $\phi \in \Phi$.

152  To find the parameters $\theta$ that maximize the evidence lower-bound, we can derive an amortized
153  learning objective. If we abbreviate the maximand from Eq. (2) as $\mathcal{L}(q, H^{1:n})$ and define $f_\theta$ :
154  $\mathcal{H}^n \to \Phi$ (omitting $Z$ for brevity), we always have the inequality

$$\max_\theta \mathbb{E}_{p(H^{1:n})} \mathcal{L}(q_{\phi=f_\theta(H^{1:n})}, H^{1:n}) \leq \mathbb{E}_{p(H^{1:n})} [\max_{\phi \in \Phi} \mathcal{L}(q_\phi, H^{1:n})]. \tag{3}$$

155  This shows that **1)** with a sufficient function class for $f_\theta$ and optimal parameters $\theta^*$, we obtain
156  equality when $f_{\theta^*}(H^{1:n}) = \arg\max_{\phi \in \Phi} \mathcal{L}(q_\phi, H^{1:n}), \forall H^{1:n} \in \mathcal{H}^n$, and **2)** $\theta^*$ can be obtained
157  through a straightforward training procedure. The l.h.s. requires us to be able to sample from
158  $p(H^{1:n})$, and that $\mathcal{L}$ is end-to-end differentiable with respect to $\theta$, which enables the use of stochastic
159  gradient methods (Kingma & Ba, 2017).

**From Inference to Control**   So far, we have mostly discussed how meta-RL agents can perform
161  inference to latent variables of the environment for a fixed policy $\pi$. To define optimal behavior in
162  the generative process we extend the probabilistic model of Eq. (1) using the control as inference
163  framework (Levine, 2018). This enables us to apply our Bayesian tools directly to our RL-agent in
164  a theoretically sound manner and, for specific design choices, recovers the typical meta-RL training
165  objective (Duan et al., 2016; Wang et al., 2017).

166  Control as inference reformulates classical RL as an inference problem by conditioning the distribu-
167  tion over trajectories $p(H^i)$ on a desired outcome $\mathcal{O}$. This outcome $\mathcal{O} \in \{0, 1\}$ is a binary variable
168  indicating whether a trajectory achieves the outcome (1) or not (0). The likelihood of this outcome

169  given a trajectory $H^i$ follows the exponentiated sum of rewards, $p(\mathcal{O} = 1|H^i) \propto \exp(\sum_t R_t^i)$.
170  Using Bayes rule, we can then infer the desired policy as $p(H^i|\mathcal{O} = 1)$.

171  Similarly to the latent variable posterior $q_\phi$, we can estimate $p(H^{1:n}|\mathcal{O} = 1)$ through variational
172  inference by defining a lower bound to the log-likelihood of the outcome variable $\ln p(\mathcal{O} = 1|H^{1:n})$
173  using a variational policy $q_\pi(A_t^i|S_t^i, Z_t^i)$. Given our choice for the outcome likelihood, this recovers
174  a regularized RL objective (Geist et al., 2019) (see the derivation in Appendix A.1.2),

$$\mathcal{L}(q_\phi, q_\pi) = \mathbb{E}_{q_\phi(H^{1:n}, Z^{1:n})} \sum_{i=1}^{n} \sum_{t=1}^{T_i} R_t^i - KL(q_\pi \| \pi) \tag{4}$$
$$- KL\left(q_\phi(Z_t^i|Z_{<t}^{\leq i}, H_{<t}^{\leq i}) \| p(Z_t^i|Z_{<t}^{\leq i}, H_{<t}^{\leq i})\right)$$

175  which is an extension of the MPO objective by Abdolmaleki et al. (2018) to include the latent-
176  variable posterior and its KL-term. The indexing of the conditional is slightly overloaded for brevity,
177  it indicates that the task posterior is conditioned on all prior data.

178  Unfortunately, this lower-bound does not give a practical training objective for both inference or
179  amortization as shown in Eq. (3), nor does it find the optimal policy given a fixed $\pi$. Therefore,
180  practitioners often include the following design choices.

181  1. The true posterior $p(Z_t^i|Z_{<t}^{\leq i}, H_{<t}^{\leq i})$ is substituted by the previous variational posterior
182     $q(Z_{t-1}^i|Z_{<t-1}^{\leq i}, H_{<t-1}^{\leq i})$, using only a "true" prior $p(Z_0)$ at time and episode 0 (Hafner et al.,
183     2019).

184  2. The policy $\pi$ is iteratively updated to the variational optimum $q_\pi$ (Abdolmaleki et al., 2018).

185  3. KL-penalties are scaled by parameters $\beta_\pi, \beta_q$.

186  Finally, observe that the amortized objective for Eq. (4) (i.e., substituted into Eq. (3)) recov-
187  ers the typical memory-based meta-RL objective when using a point-estimate (Dirac posterior)
188  $q_\phi(Z|\dots) = \delta(\phi - Z)$ and ignoring the KL-penalties (Duan et al., 2016).

189  ## 4   Laplace Variational RNNs

190  We introduce the Laplace variational recurrent neural network (Laplace VRNN) to make a relatively
191  simple approximation to the variational task posterior $q_\theta \approx \hat{q}_\theta$, to be used in the lower-bounds
192  of Eq. (2) and Eq. (4), using the Laplace approximation (MacKay, 1992; Daxberger et al., 2021).
193  This enables the construction of proper distributions over the latent-variables without introducing
194  additional variational parameters. The idea is that we use this to extend point-estimate methods
195  (i.e., base RNNs) to use distributions at any point during training. For exposition, we introduce our
196  approximation starting from a simpler variational distribution $q_\phi(Z_t|H_{<t})$, dropping superscripts.
197  The full derivation is given in Appendix A.3.

198  We use the predicted distribution parameters $\phi_t = f_\theta(H_{<t})$ as a helper variable, and interchange the
199  notation $q_{\phi_t}(Z_t|H_{<t}) = q_\theta(Z_t|H_{<t}, \phi_t)$. Most importantly, the statistical amortization by Eq. (3)
200  allows us to interpret the mapping $f_\theta : H \mapsto \phi$ as a learned summary statistic for the distribution
201  of $Z$; i.e., a *maximum a posteriori* estimate. We assume that $\phi_t$ is computed autoregressively with a
202  recurrent neural network (RNN) such that $\phi_{t+1} = f_\theta(S_t, A_t, R_t; \phi_t)$.

203  We then factorize using a mean-field assumption,

$$q_\theta(Z_t|H_{<t}, \phi_t) = \frac{1}{q_\theta(Z_t|\phi_t)^{t-2}} \prod_{i=1}^{t-1} q_\theta(Z_t|S_i, R_i, A_i, \phi_t),$$
$$= \exp[(2 - t)\ln q_\theta(Z_t|\phi_t) + \sum_{i=1}^{t-1} \ln q_\theta(Z_t|S_i, R_i, A_i, \phi_t)]$$
$$= \exp h_\theta(Z_t; H_{<t}, \phi_t), \tag{5}$$

204 which gives us the target function $h_\theta$ that we wish to approximate (for the first step, see Lemma 1;
205 Appendix A.2). For a given $\theta$ and data $H_{<t}$, we use the second order Taylor expansion of $h_\theta \approx \hat{h}_\theta$
206 linearized at $\phi = \phi_t$, we then exponentiate $\hat{h}_\theta$ and renormalize. Assume that $\phi_t$ is *maximum a*
207 *posteriori* to $q_\theta(Z_t|H_{<t}, \phi_t)$, then we obtain the Laplace approximation,

$$
q_\theta(Z_t|H_{<t}) \approx \frac{\exp \hat{h}_\theta(Z_t; H_{<t}, \phi_t)}{\int \exp \hat{h}_\theta(z; H_{<t}, \phi_t) dz}
$$

$$
= \mathcal{N}(\phi_t, \nabla_\phi^2 \ln q_\theta(Z_t|H_{<t}, \phi)|_{\phi=\phi_t}), \tag{6}
$$

208 where $\nabla_\phi^2 \ln q_\theta$ is the Hessian of our log-posterior.

209 To complete our model from Eq. (2) and Eq. (4), we can solve the expectation over $\mathbb{E}_{q_\theta(Z_{<t}|H_{<t})}$
210 for the posterior $q_\theta(Z_t|H_{<t}, Z_{<t})$ at each time-step $t$, by assuming a convolution of Gaussian den-
211 sities. The result of this convolution is well-known to be another Gaussian with summed parameters
212 (Bromiley, 2003),

$$
\mu_t = \phi_t + \sum_{i=1}^{t-1} \mu_i, \tag{7}
$$

$$
\Lambda_t = -\nabla_\phi^2 \ln q_\theta(Z|H_{<t}, \phi)|_{\phi=\phi_t} + \sum_{i=1}^{t-1} \Lambda_i. \tag{8}
$$

213 In practice, we use a smaller window $H_{k:t-1}$ and $Z_{k:t-1}$ inside the conditional for efficiency. In
214 summary, this implements an RNN where we sum the last $k$ hidden states and covariances for the
215 output-Gaussian, and where the covariances are produced by the Hessian of the log-posterior with
216 respect to the hidden state.

217 The assumption that $\phi_t$ is *maximum a posteriori* is quite strict and assumes equality in the amor-
218 tization objective Eq. (3). For a sub-optimal $\theta$ or insufficient function class $f_\theta$, this can induce a
219 first-order error in the Taylor expansion of the variational log-posterior, which results in a worse
220 approximation. Furthermore, most RNN methods do not sum their hidden states for the predictive
221 model, which mismatches our formulation. It is also not obvious what representations RNN-based
222 meta-RL agents learn and, thus, which model factorization would best suit the agent for construc-
223 tion of our Laplace approximation. We investigate these technicalities and assumptions in the next
224 section.

225 **Special Case**   Our method obtains a particularly nice form if we choose $q_\theta(Z_t|S_i, A_i, R_i, \phi_t)$ to
226 be standard Gaussian and use an uninformative prior for $q_\theta(Z_t|\phi_t)$. In that case, the Hessian of the
227 log-posterior $q_\theta(Z_t|H_{<t}, \phi_t)$ becomes a sum of outer products of our RNN state Jacobians w.r.t. $\phi$.
228 Let $x_i = (S_i, A_i, R_i)$, this gives the inverse covariance,

$$
\Lambda_t = \sum_{i=1}^{t-1} (\nabla_\phi f_\theta(x_i; \phi)|_{\phi=\phi_t})(\nabla_\phi f_\theta(x_i; \phi)|_{\phi=\phi_t})^\top, \tag{9}
$$

229 which is cheap to compute with forward accumulation (Bradbury et al., 2018), see Prop. 2 in Ap-
230 pendix A.3.

231 **Posterior Predictive**   If we now choose a policy $\pi(A_t|S_t, Z_t)$ that is linear in $Z_t$, then our full
232 model would recover a type of Gaussian process (Immer et al., 2021; Rasmussen & Williams, 2005).
233 However, we model this term with another neural network $\pi_\psi(A_t|S_t, Z_t)$ to improve expressiveness.
234 Our policy is then defined by the *posterior predictive* $\pi_\psi(A_t|S_t, H_{<t})$, which we compute using
235 Monte-Carlo,

$$
\int \pi_\psi(A_t|S_t, z_t) q_\theta(z_t|H_{<t}) dz_t \approx \frac{1}{k} \sum_{i=1}^{k} \pi_\psi(A_t|S_t, Z_t = z^{(i)}), \quad z^{(i)} \sim q_\theta(Z_t|H_{<t}), \tag{10}
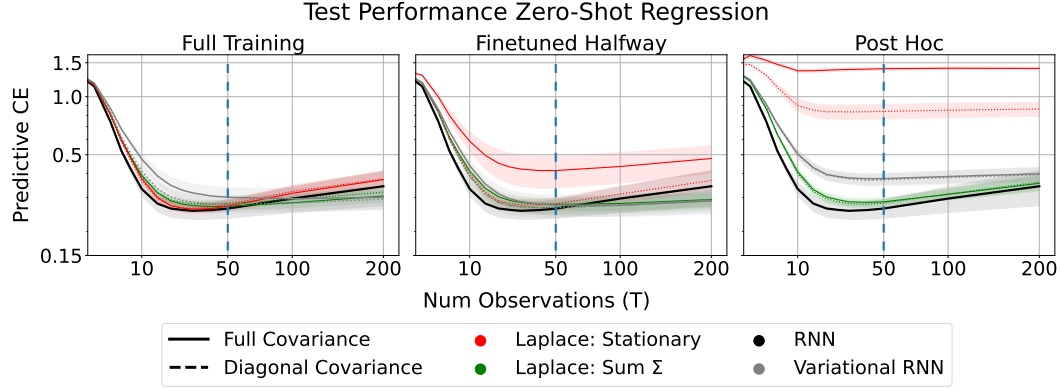$$

Figure 1: Final performance on the zero-shot regression task in terms of predictive cross-entropy, this should decrease over time. The left column shows results for complete model training, the middle and right columns perform model pre-training with the RNN (black). The middle column includes parameter finetuning, the right column does not. The blue dashed line indicates the training cut-off ($T = 50$).

overloading superscripts to index Monte-Carlo samples. This induces a finite mixture for the policy where $k = 1$ corresponds to posterior sampling (Osband et al., 2013). The parameters $\theta$ and $\psi$ were trained jointly in an end-to-end manner (as defined in the l.h.s. of Eq.(3)).

Interestingly, during training we found output aggregation to train more stably in the loss when $k > 1$. Thus, during training we chose to average the predicted logits of $\pi_\psi$ over samples $z^{(i)}$, or in the continuous case, we averaged over the parameters of a parametric distribution (Wang et al., 2020), e.g., the mean and variance of a Gaussian (see Appendix B.3 for a discussion).

## 5   Experimental Validation

In order to apply our method to memory-based meta-RL agents in a manner that didn't alter the network architecture, we required a few extra simplifications to our more general model from Section 4. Crucially, for the *non-stationary* assumption on the log-posterior only, this required us to omit the mean aggregation of Eq. (7). This wasn't needed for the stationary factorizations, see Appendix A.3.1 for a more detailed discussion. To evaluate our factorizations and design choices, we performed experiments to answer the following:

1. **Utility:** Does our method give useful posterior statistics for a non-Bayesian baseline?

2. **Sensitivity:** What model assumptions for the Laplace VRNN are empirically effective?

3. **Performance:** When used as an alternative to variational inference, does the Laplace VRNN perform at least on par with existing methods?

Our point-estimate (RNN) baseline was implemented with a long-short term memory architecture (Hochreiter & Schmidhuber, 1997). The VRNN baseline (Chung et al., 2015) extends the RNN by predicting the mean and covariance for a Gaussian distribution as a transformation of the RNN output. For the RNN and VRNN we assumed a stationary factorization of the posterior $q_\theta(Z_t|H_{t-1})$, which is a simplification of the fully general posterior shown in Eq. (1) and most accurate to the true generative process. We intermittently created model snapshots of the point-estimate baseline (RNN) and finetuned these snapshots over our parameter grid for the Laplace VRNN and VRNN.

All experiments were repeated over $r = 30$ seeds (number of network initializations), we tested intermediate model parameters by measuring their in-distribution performance and model statistics for $B = 128$ samples (number of test-tasks). We report 2-sided confidence intervals with a confi-
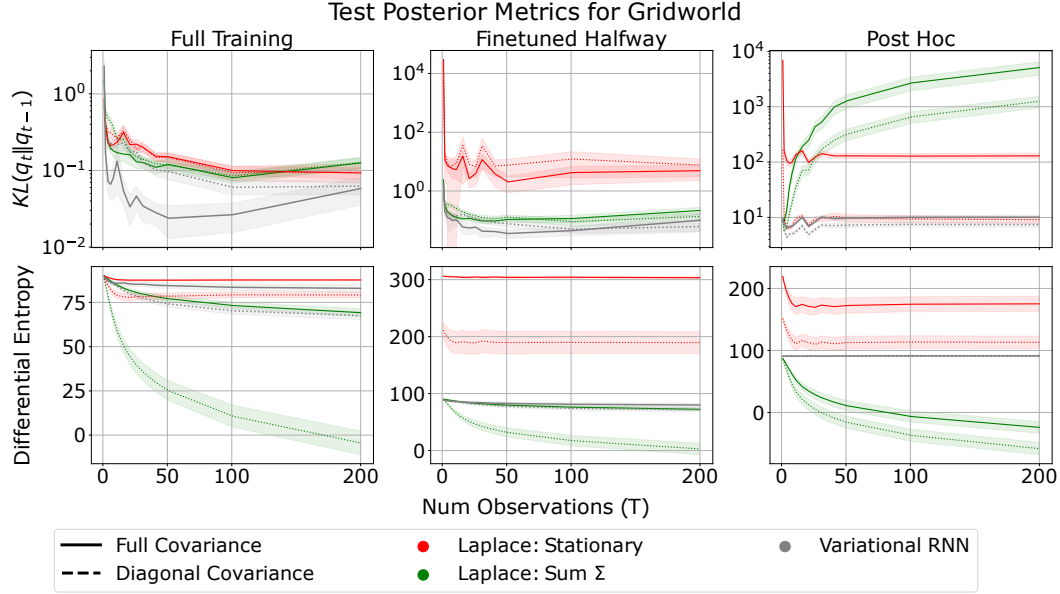
Figure 2: Evolution of summary statistics for the posterior model during testing. The top row shows the KL-divergences between consecutive posteriors $q_t$ and $q_{t+1}$, and the bottom row shows the model entropy over time. In principle, we expect all lines to decrease gradually with more observations. When applying our Laplace approximation with summed covariances (green) after deterministic pre-training (right-column), we see that the posterior becomes more and more confident but does not converge to a stable distribution.

dence level $\alpha = 0.99$ for each metric $X$ aggregated over the seeds $r$ and the test-tasks $B$. For the full details on the experiment and baseline setup, see Appendix B (code available upon publication).

## 5.1 Supervised Learning

As a didactic test-setup, we evaluated our method on noiseless 1D regression tasks. We generated data by sampling parameters to a Fourier expansion and then sampling datasets $\{\{(X_i, f_j(X_i))\}_{i=1}^T\}_{j=1}^n$ where each $X_i \sim \text{Unif}(-1, 1)$, $f_j \sim p_{\text{Fourier}}(f)$, and $n = 256, T = 50$. During training, we optimized a lower bound for a supervised domain using a weight for the KL-term of $\beta = 10^{-2}$ (see Eq. (2); Appendix A.1.1). During testing, we computed the predictive cross-entropy (CE) with the true data-generating distribution and our model. So, at each step $t$, we used $H_t = \{X_i, f(X_i)\}_{i=1}^t$ to estimate the posterior predictive distribution $\mathbb{E}_{q_\theta(Z_t|H_t)} p_\theta(Y_i|X_i, Z_t)$ with Monte-Carlo using $m = 30$ samples. The predictive CE was estimated using Monte-Carlo over a large test dataset.

**Variations** Our Laplace VRNN used a stationary $q_\theta(Z_t|H_{<t})$ assumption (Laplace: Stationary; red) and a Markovian $q_\theta(Z_t|H_{t-1}, Z_{t-1})$ assumption (Laplace: Sum $\Sigma$; green) on the graphical model from Eq. (1). In practice, the stationary model computes the covariance for each datapoint in $H_{<t}$ at each $t$, whereas the Markovian model sums the covariances for each pair $(X_i, Y_i)$. To reduce clutter, we only show the Laplace VRNN ablation that sums the covariance, which also performed best among our variations (c.f., Appendix C.1). We tested both diagonal and full covariance matrices.

**Results** As shown in Figure 1, across our comparisons the predictive CE goes down initially (except for the post-hoc stationary Laplace VRNN), however slightly increases again after the size of the dataset exceeds that seen during training $T > 50$. Notably, we see that our stationary Laplace VRNN strongly degrades performance when not used at the beginning (red), most notably the full-
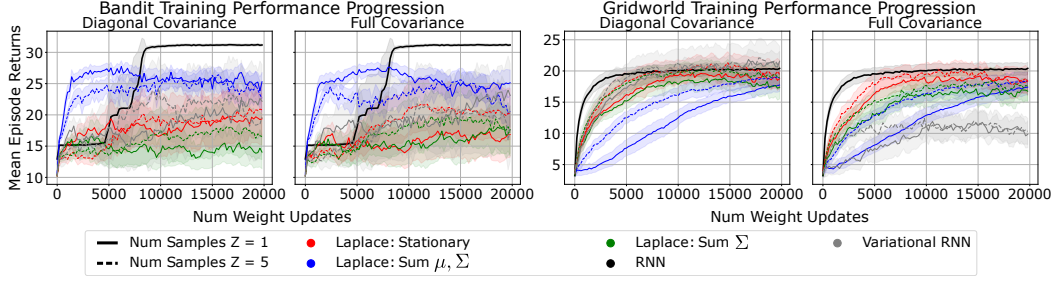
Figure 3: Average return curves during training for the Reinforcement Learning experiments. The dashed and solid lines (Num $Z$) indicate the number of Monte-Carlo samples used for the posterior model during training inside the modified lower bound of Eq. (4), to validate off-policy robustness in our loss. As expected, the deterministic RNN performs best, but our Laplace VRNN also outperforms the baseline VRNN.

covariance variation. This could be an indication that the Laplace approximated posterior is too wide, whereas the point-estimate is extremely sharp, causing samples from our method to be out-of-distribution for the predictive model. In contrast, this result also shows that our method with the Markovian assumption (green) performs at least as well as the baselines in all cases while also providing a Bayesian posterior for the RNN after training.

## 5.2 Reinforcement Learning

To show that our method can perform uncertainty quantification while maintaining strong performance in reinforcement learning problems, we evaluated our method on a stochastic 5-armed bandit and a deterministic $5 \times 5$ gridworld with sparse rewards. We tested all models using a variant of recurrent PPO (Schulman et al., 2017) as a simple approximation for Eq. (4). During training, we used a batch size of $B = 256$ task samples and a sequence length of $T = 50$ interactions with the bandit and $T = 100$ for the gridworld task.

For the bandit, we generated training tasks by sampling reward probabilities from a Dirichlet prior using $\vec{\alpha} = 0.2$. In this domain we only condition our policy on the sampled model hypotheses $z_t \sim q_\theta(Z_t|H_t)$, $a_t \sim \pi_\theta(Z = z_t)$ as is typical in Thompson sampling (Osband et al., 2013). This experiment also aimed to investigate robustness to model sampling noise. For the gridworld (Zintgraf et al., 2021) we sampled tasks by generating the agent's start- and goal-tile uniformly randomly across the grid. In contrast to the bandit problem, the gridworld agent modeled the task as a Bayes-adaptive Markov decision process (Duff, 2002). Meaning that the policy conditions on both the model samples $Z$ and the current state, $a_t \sim \pi_\theta(Z = z_t, S = s_t)$.

**Variations** As before, we test different assumptions for our Laplace VRNN agent's model from Eq. (1). In this instance, we used a *windowed* version of the stationary Laplace VRNN $q_\theta(Z_t|H_{t-w-1:t-1})$ for $w = 10$ (red). I.e., this truncates the history up to a certain timestep to improve the runtime of the covariance computation, which otherwise scales in $\mathcal{O}(t)$-time. We also tested two variations of the Markovian $q_\theta(Z_t|H_{t-1}, Z_{t-1})$ factorization, which scaled in $\mathcal{O}(1)$-time. The proper-Markovian method (blue) sums the mean and covariance computed at each state-action $(S_t, A_t)$ whereas the second variant only sums the covariance (green).

**Results** We visualize the evolution of estimated posterior statistics during testing in figure 2, where we removed the ablation that sums both the mean and covariance (blue) to reduce clutter (this ablation performed in between the other two, see Appendix C.3). We plotted the differential entropy of the posteriors $q_\theta$ and the consecutive KL-divergences $KL(q_\theta(Z_t|\ldots)\|q_\theta(Z_{t-1}|\ldots))$ between posteriors over time, to see whether their behavior matches that of the true posterior. For the true posterior, we expect the entropy to decrease gradually with more observations $T$, which indicates

that our model concentrates around some true value. Furthermore, the KL-divergences should converge to zero.

As expected, we see that the posterior entropy of the Bayesian methods reliably goes down and the KL-divergences gets close to zero (left-column). We see a similar pattern when doing finetuning (middle-columns) except for our stationary Laplace variation (red). Most importantly, we see a strong effect of our accumulating covariances variation (green) when using a post-hoc posterior approximation. We see that the entropy steadily decreases while the KL-divergences between consecutive posteriors grows larger and larger. In contrast, the post-hoc VRNN (grey) and stationary Laplace (red) stay relatively constant, and are therefore non-informative. This result shows that the deterministic RNN does not converge to a stable hidden state when not explicitly regularized during training. This means that the learned estimator becomes more and more confident while not being consistent (Xiong et al., 2021).

The average training returns for our model ablations are shown in Figure 3. As argued in the introduction, we find that the deterministic method (black) has strong performance while also being the least noisy in the mean episode returns and being the fastest to train in terms of algorithm runtime. Interestingly, the proper-Markovian factorization (blue) of our Laplace VRNN showed faster learning in the Bandit up to a certain point, whereas it degraded training performance for the gridworld. All Bayesian methods tested on the bandit task were significantly noisy and only achieved sub-linear cumulative regret during test-time about $50\%$ over all experiment repetitions. On the grid task, all methods achieved sub-linear cumulative regret.

In summary, none of the ablations for our Bayesian methods degraded performance when applied after deterministic pre-training without finetuning (post-hoc). Our Laplace VRNN also typically performed on par with the VRNN in terms of returns. However, only our Markovian Laplace variation that summed the covariances (green) could produce insightful posterior statistics of the pre-trained model. This confirms all our research questions of whether our proposed Laplace VRNN, and what model assumptions, can give useful posterior statistics while not degrading performance.

# 6   Conclusions

We have described how the Laplace approximation can be applied to recurrent neural network models in a zero-shot meta-reinforcement learning context. Our method is a cheap transformation of an existing recurrent network to a Bayesian model. This enables trained agents to more accurately model their task-uncertainty which can be used to create better or more robust methods.

We tested our method on supervised and reinforcement learning tasks to investigate the utility of our approximation (the quality of posterior statistics), how it depends on model assumptions (ablations), and how it compares against variational inference or point-estimate baselines (no degradation in performance). Our results show that the proposed *Laplace variational recurrent neural network* can reliably transform existing non-Bayesian models to produce a Bayesian posterior, at any point during training without modifying the model or training procedure. In contrast, variational inference requires altering the model architecture and training setup despite matching (or underperforming) compared to our method.

One limitation of our method is the computation of the Jacobians and possible restrictiveness of the Gaussian distribution. Furthermore, the RNN based agents required a variety of simplifications to our probabilistic model formulation (at least, for some of the tested configurations). Although the results matched expected behavior, and is consistent with prior work (Xiong et al., 2021; Mikulik et al., 2020), future work should investigate what the induced biases entail for our approximated posterior. Our method also does not fix the overconfidence of the non-Bayesian agents, but enables us to observe this effect. Extending our approach to enable statistically sound representation learning or to improve exploration through e.g., distribution-shift detection in meta-RL are a promising directions for further study (Daxberger et al., 2021).

# References

Abbas Abdolmaleki, Jost Tobias Springenberg, Yuval Tassa, Remi Munos, Nicolas Heess, and Martin Riedmiller. Maximum a Posteriori Policy Optimisation. In *International Conference on Learning Representations*, 2018.

Marcin Andrychowicz, Misha Denil, Sergio Gómez, Matthew W Hoffman, David Pfau, Tom Schaul, Brendan Shillingford, and Nando de Freitas. Learning to learn by gradient descent by gradient descent. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016.

Søren Asmussen and Peter W. Glynn. *Stochastic Simulation: Algorithms and Analysis*, volume 57 of *Stochastic Modelling and Applied Probability*. Springer, New York, NY, 2007. DOI: 10.1007/978-0-387-69033-9.

Jakob Bauer, Kate Baumli, Feryal Behbahani, Avishkar Bhoopchand, Nathalie Bradley-Schmieg, Michael Chang, Natalie Clay, Adrian Collister, Vibhavari Dasagi, Lucy Gonzalez, Karol Gregor, Edward Hughes, Sheleem Kashem, Maria Loks-Thompson, Hannah Openshaw, Jack Parker-Holder, Shreya Pathak, Nicolas Perez-Nieves, Nemanja Rakicevic, Tim Rocktäschel, Yannick Schroecker, Satinder Singh, Jakub Sygnowski, Karl Tuyls, Sarah York, Alexander Zacherl, and Lei M Zhang. Human-Timescale Adaptation in an Open-Ended Task Space. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202, pp. 1887–1935. PMLR, 2023.

Sarah Bechtle, Artem Molchanov, Yevgen Chebotar, Edward Grefenstette, Ludovic Righetti, Gaurav Sukhatme, and Franziska Meier. Meta Learning via Learned Loss. In *2020 25th International Conference on Pattern Recognition (ICPR)*, pp. 4161–4168, 2021. DOI: 10.1109/ICPR48806.2021.9412010.

Jacob Beck, Risto Vuorio, Evan Zheran Liu, Zheng Xiong, Luisa Zintgraf, Chelsea Finn, and Shimon Whiteson. A survey of meta-reinforcement learning. arXiv:2301.08028, 2023.

Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer, 1st edition, 2007.

Aleksandar Botev, Hippolyt Ritter, and David Barber. Practical Gauss-Newton Optimisation for Deep Learning. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70, pp. 557–565. PMLR, 2017.

James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018.

Paul Bromiley. Products and convolutions of Gaussian probability density functions. *Tina-Vision Memo*, 3(4):1, 2003.

Yutian Chen, Matthew W. Hoffman, Sergio Gómez Colmenarejo, Misha Denil, Timothy P. Lillicrap, Matt Botvinick, and Nando de Freitas. Learning to Learn without Gradient Descent by Gradient Descent. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70, pp. 748–756. PMLR, 2017.

Junyoung Chung, Kyle Kastner, Laurent Dinh, Kratarth Goel, Aaron C Courville, and Yoshua Bengio. A Recurrent Latent Variable Model for Sequential Data. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015.

George Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of control, signals and systems*, 2(4):303–314, 1989.

Erik Daxberger, Agustinus Kristiadi, Alexander Immer, Runa Eschenhagen, Matthias Bauer, and Philipp Hennig. Laplace Redux - Effortless Bayesian Deep Learning. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, volume 34, pp. 20089–20103. Curran Associates, Inc., 2021.

DeepMind, Igor Babuschkin, Kate Baumli, Alison Bell, Surya Bhupatiraju, Jake Bruce, Peter Buchlovsky, David Budden, Trevor Cai, Aidan Clark, Ivo Danihelka, Antoine Dedieu, Claudio Fantacci, Jonathan Godwin, Chris Jones, Ross Hemsley, Tom Hennigan, Matteo Hessel, Shaobo Hou, Steven Kapturowski, Thomas Keck, Iurii Kemaev, Michael King, Markus Kunesch, Lena Martens, Hamza Merzic, Vladimir Mikulik, Tamara Norman, George Papamakarios, John Quan, Roman Ring, Francisco Ruiz, Alvaro Sanchez, Laurent Sartran, Rosalia Schneider, Eren Sezener, Stephen Spencer, Srivatsan Srinivasan, Miloš Stanojević, Wojciech Stokowiec, Luyu Wang, Guangyao Zhou, and Fabio Viola. The DeepMind JAX Ecosystem, 2020.

Yan Duan, John Schulman, Xi Chen, Peter L. Bartlett, Ilya Sutskever, and Pieter Abbeel. $RL^2$: Fast reinforcement learning via slow reinforcement learning. arXiv:1611.02779, 2016.

Michael O'Gordon Duff. *Optimal learning: Computational procedures for Bayes-adaptive Markov decision processes*. PhD thesis, The University of Massachusetts Amherst, 2002.

Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70, pp. 1126–1135. PMLR, 2017.

Chelsea Finn, Kelvin Xu, and Sergey Levine. Probabilistic Model-Agnostic Meta-Learning. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.

Marta Garnelo, Jonathan Schwarz, Dan Rosenbaum, Fabio Viola, Danilo J. Rezende, S. M. Ali Eslami, and Yee Whye Teh. Neural Processes. arXiv:1807.01622, 2018.

Matthieu Geist, Bruno Scherrer, and Olivier Pietquin. A Theory of Regularized Markov Decision Processes. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97, pp. 2160–2169. PMLR, 2019.

Samuel Gershman and Noah Goodman. Amortized inference in probabilistic reasoning. In *Proceedings of the annual meeting of the cognitive science society*, volume 36, 2014.

Mohammad Ghavamzadeh, Shie Mannor, Joelle Pineau, and Aviv Tamar. Bayesian Reinforcement Learning: A Survey. *Found. Trends Mach. Learn.*, 8(5–6):359–483, 2015. DOI: 10.1561/2200000049.

Jonathan Gordon, John Bronskill, Matthias Bauer, Sebastian Nowozin, and Richard Turner. Meta-Learning Probabilistic Inference for Prediction. In *International Conference on Learning Representations*, 2019.

Anirudh Goyal, Alessandro Sordoni, Marc-Alexandre Côté, Nan Rosemary Ke, and Yoshua Bengio. Z-Forcing: Training Stochastic Recurrent Networks. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.

Erin Grant, Chelsea Finn, Sergey Levine, Trevor Darrell, and Thomas Griffiths. Recasting Gradient-Based Meta-Learning as Hierarchical Bayes. In *International Conference on Learning Representations*, 2018.

Ido Greenberg, Shie Mannor, Gal Chechik, and Eli Meirom. Train Hard, Fight Easy: Robust Meta Reinforcement Learning. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 68276–68299. Curran Associates, Inc., 2023.

Albert Gu and Tri Dao. Mamba: Linear-Time Sequence Modeling with Selective State Spaces. In *First Conference on Language Modeling*, 2024.

Abhishek Gupta, Russell Mendonca, YuXuan Liu, Pieter Abbeel, and Sergey Levine. Meta-Reinforcement Learning of Structured Exploration Strategies. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.

Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and James Davidson. Learning Latent Dynamics for Planning from Pixels. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97, pp. 2555–2565. PMLR, 2019.

Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to Control: Learning Behaviors by Latent Imagination. In *International Conference on Learning Representations*, 2020.

Sepp Hochreiter and Jürgen Schmidhuber. Long Short-Term Memory. *Neural Computation MIT-Press*, 1997.

Timothy M. Hospedales, Antreas Antoniou, Paul Micaelli, and Amos J. Storkey. Meta-Learning in Neural Networks: A Survey. *IEEE transactions on pattern analysis and machine intelligence*, 44 (9):5149–5169, 2021.

Alexander Immer, Maciej Korzepa, and Matthias Bauer. Improving predictions of Bayesian neural nets via local linearization. In Arindam Banerjee and Kenji Fukumizu (eds.), *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, volume 130, pp. 703–711. PMLR, 2021.

Max Jaderberg, Wojciech M. Czarnecki, Iain Dunning, Luke Marris, Guy Lever, Antonio Garcia Castañeda, Charles Beattie, Neil C. Rabinowitz, Ari S. Morcos, Avraham Ruderman, Nicolas Sonnerat, Tim Green, Louise Deason, Joel Z. Leibo, David Silver, Demis Hassabis, Koray Kavukcuoglu, and Thore Graepel. Human-level performance in 3D multiplayer games with population-based reinforcement learning. *Science*, 364(6443):859–865, 2019. DOI: 10.1126/science.aau6249.

Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are RNNs: Fast Autoregressive Transformers with Linear Attention. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119, pp. 5156–5165. PMLR, 2020.

Khimya Khetarpal, Matthew Riemer, Irina Rish, and Doina Precup. Towards Continual Reinforcement Learning: A Review and Perspectives. *Journal of Artificial Intelligence Research*, 75:1401–1476, 2022.

Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. arXiv:1412.6980, 2017.

James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A. Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, Demis Hassabis, Claudia Clopath, Dharshan Kumaran, and Raia Hadsell. Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences*, 114(13):3521–3526, 2017. DOI: 10.1073/pnas.1611835114.

Sergey Levine. Reinforcement Learning and Control as Probabilistic Inference: Tutorial and Review. arXiv:1805.00909, 2018.

Ilya Loshchilov and Frank Hutter. Decoupled Weight Decay Regularization. In *International Conference on Learning Representations*, 2019.

Chris Lu, Jakub Kuba, Alistair Letcher, Luke Metz, Christian Schroeder de Witt, and Jakob Foerster. Discovered Policy Optimisation. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 16455–16468. Curran Associates, Inc., 2022.

David J. C. MacKay. Bayesian Interpolation. *Neural Computation*, 4(3):415–447, 1992. DOI: 10.1162/neco.1992.4.3.415.

James Martens and Roger Grosse. Optimizing Neural Networks with Kronecker-factored Approximate Curvature. In Francis Bach and David Blei (eds.), *Proceedings of the 32nd International Conference on Machine Learning*, volume 37, pp. 2408–2417, Lille, France, 2015. PMLR.

Vladimir Mikulik, Grégoire Delétang, Tom McGrath, Tim Genewein, Miljan Martic, Shane Legg, and Pedro Ortega. Meta-trained agents implement Bayes-optimal agents. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 18691–18703. Curran Associates, Inc., 2020.

Ian Osband, Daniel Russo, and Benjamin Van Roy. (More) Efficient Reinforcement Learning via Posterior Sampling. In C.J. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger (eds.), *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc., 2013.

Carl Edward Rasmussen and Christopher K. I. Williams. *Gaussian Processes for Machine Learning*, volume 2. The MIT Press, Cambridge, MA, 2005. DOI: 10.7551/mitpress/3206.001.0001.

Hippolyt Ritter, Aleksandar Botev, and David Barber. Online Structured Laplace Approximations for Overcoming Catastrophic Forgetting. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.

John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust Region Policy Optimization. In Francis Bach and David Blei (eds.), *Proceedings of the 32nd International Conference on Machine Learning*, volume 37, pp. 1889–1897, Lille, France, 2015. PMLR.

John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal Policy Optimization Algorithms. arXiv:1707.06347, 2017.

Ramanan Sekar, Oleh Rybkin, Kostas Daniilidis, Pieter Abbeel, Danijar Hafner, and Deepak Pathak. Planning to Explore via Self-Supervised World Models. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119, pp. 8583–8592. PMLR, 2020.

Ilia Shumailov, Zakhar Shumaylov, Yiren Zhao, Yarin Gal, Nicolas Papernot, and Ross Anderson. The Curse of Recursion: Training on Generated Data Makes Models Forget. arXiv:2306.17493, 2023.

Gautam Singh, Jaesik Yoon, Youngsung Son, and Sungjin Ahn. Sequential Neural Processes. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.

Shagun Sodhani, Franziska Meier, Joelle Pineau, and Amy Zhang. Block Contextual MDPs for Continual Learning. In Roya Firoozi, Negar Mehr, Esen Yel, Rika Antonova, Jeannette Bohg, Mac Schwager, and Mykel Kochenderfer (eds.), *Proceedings of The 4th Annual Learning for Dynamics and Control Conference*, volume 168, pp. 608–623. PMLR, 2022.

Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. A Bradford Book, 2nd edition, 2018.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. Attention is All you Need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.

Che Wang, Yanqiu Wu, Quan Vuong, and Keith Ross. Striving for Simplicity and Performance in Off-Policy DRL: Output Normalization and Non-Uniform Sampling. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119, pp. 10070–10080. PMLR, 2020.

Jane X Wang, Zeb Kurth-Nelson, Dhruva Tirumala, Hubert Soyer, Joel Z Leibo, Remi Munos, Charles Blundell, Dharshan Kumaran, and Matt Botvinick. Learning to reinforcement learn. arXiv:1611.05763, 2017.

Mike Wu, Kristy Choi, Noah Goodman, and Stefano Ermon. Meta-Amortized Variational Inference and Learning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(04):6404–6412, 2020. DOI: 10.1609/aaai.v34i04.6111.

Zheng Xiong, Luisa M Zintgraf, Jacob Austin Beck, Risto Vuorio, and Shimon Whiteson. On the Practical Consistency of Meta-Reinforcement Learning Algorithms. In *Fifth Workshop on Meta-Learning at the Conference on Neural Information Processing Systems*, 2021.

Luisa Zintgraf, Sebastian Schulze, Cong Lu, Leo Feng, Maximilian Igl, Kyriacos Shiarlis, Yarin Gal, Katja Hofmann, and Shimon Whiteson. VariBAD: Variational Bayes-Adaptive Deep RL via Meta-Learning. *Journal of Machine Learning Research*, 22(289):1–39, 2021.

Bernt Øksendal. *Stochastic Differential Equations*. Universitext. Springer, 2003. DOI: 10.1007/978-3-642-14394-6.

# Supplementary Materials

*The following content was not necessarily subject to peer review.*

## A   Derivations

### A.1   Lower Bounds

In this section, we derive the two lower bounds used for model training in the main paper. These lower bounds are not particularly new or special, they only show how one can derive a learning objective from a probabilistic graphical modeling perspective.

#### A.1.1   Supervised Learning

For the supervised learning domain we can derive an evidence lower bound on the data-marginal as a training objective for our neural network parameters in the following way. For all permutations of $H^{1:n} = \{X^i, Y^i\}_{i=1}^n$, where $X^i \in \mathcal{X}, Y^i \in \mathcal{Y}$, we have,

$$p(H^{1:n}) = \int p(H^{1:n}|z)p(z)dz \tag{11}$$

$$= \int p(X^n, Y^n|z, H^{<n})p(H^{<n}|z)p(z)dz \tag{12}$$

$$= \int p(X^n, Y^n|z)p(H^{<n}|z)p(z)dz \tag{13}$$

$$= p(H^{<n}) \int p(X^n, Y^n|z)p(z|H^{<n})dz \tag{14}$$

if we complete the recursion for $p(H^{<n})$ and do importance sampling on the posterior with $q$, we get,

$$\ln p(H^{1:n}) = \ln \prod_{i=1}^n \int p(X^i, Y^i|z^i)p(z^i|H^{1:i})dz^{1:n} \tag{15}$$

$$= \sum_{i=1}^n \ln \int p(X^i, Y^i|z^i)\frac{q(z^i|H^{1:i})}{q(z^i|H^{1:i})}p(z^i|H^{1:i})dz^{1:n} \tag{16}$$

$$\geq \sum_{i=1}^n \int q(z^i|H^{1:i})\left[\ln p(X^i, Y^i|z^i) + \ln \frac{p(z^i|H^{1:i})}{q(z^i|H^{1:i})}\right]dz^{1:n} \tag{17}$$

$$= \sum_{i=1}^n \mathbb{E}_{q(Z^i|H^{1:i})} \ln p(X^i, Y^i|Z^i) - KL(q(Z^i|H^{1:i})\|p(Z^i|H^{1:i})), \tag{18}$$

which gives us the lower-bound for our approximate inference model when we use neural network parameters $\theta$ for the predictive and posterior models (overloading notation for $q_\phi$ in Eq.(3)),

$$\mathcal{L}(\theta, H^{1:n}) = \sum_{i=1}^n \mathbb{E}_{q_\theta(Z^i|H^{1:i})} \underbrace{\ln p_\theta(X^i, Y^i|Z^i)}_{\text{Prediction Loss}}$$
$$- \beta \cdot \underbrace{KL(q_\theta(Z^i|H^{1:i})\|\texttt{stop\_grad}[q_\theta(Z^{i-1}|H^{1:i-1}))}_{\text{Complexity Penalty}}, \tag{19}$$

where $\texttt{stop\_grad}[\cdot]$ indicates a stop-gradient operation and $\mathcal{L}$ should be *maximized* with respect to $\theta$. The hyperparameter $\beta \in \mathbb{R}_+$ accounts for differences in scaling. The stop-gradient is necessary so that the posterior at time $t$ does not depend on the future. During generation and training, we also assume a uniform prior over the inputs $p(X^i|Z) = \text{Unif}$. Of course, this is just one lower bound,

16

592 the one we use in the main paper for the reinforcement learning tasks also assumes that each $z^i$ is
593 sequentially dependent. In this case, the product would appear inside the integral in the first line for
594 $\ln p(H^{1:n})$, this is only relevant for the Laplace VRNN that accumulates the mean and covariances.

595 For simplicity, we only perform training on a single permutation of $H^n$ (i.e., canonical order), as in
596 expectation all permutations are covered anyway and this provides training batches with more di-
597 verse examples. Unfortunately, when amortizing the computation of this lower bound with recurrent
598 models it can be difficult to properly distill this permutation invariance of the data into the model.
599 Using a recurrent model that linearly transforms the state, like a transformer (Vaswani et al., 2017;
600 Katharopoulos et al., 2020) or general state space model (Bishop, 2007), would prevent this. We
601 leave this open for future work.

### A.1.2 Reinforcement Learning

603 Consider the joint distribution over environment traces $H^i = \{S_t, R_t, A_t\}_{t=1}^T$ and latent variables
604 $Z$, we'll write episode indices (*extra*-episodic) $i = 1 \ldots, n$, in the superscript and time indices
605 (*inter*-episodic) $t = 1, \ldots, T$, in the subscript,

$$p(H^{1:n}, Z^{1:n}) = \prod_{i=1}^{n} p(H^i, Z^i | H^{<i}, Z^{<i}) \tag{20}$$

$$= \prod_{i=1}^{n} \prod_{t=1}^{T_i} p(S_t^i, R_t^i, A_t^i, Z_t^i | S_{<t}^i, R_{<t}^i, A_{<t}^i, Z_{<t}^i, Z^{<i}, H^{<i}) \tag{21}$$

$$= \prod_{i=1}^{n} \prod_{t=1}^{T_i} p(S_t^i, R_t^i, A_t^i | Z_t^i, H_{<t}^i) p(Z_t^i | \underbrace{Z_{<t}^i, H_{<t}^i}_{inter}, \underbrace{Z^{<i}, H^{<i}}_{extra}) \tag{22}$$

$$= \prod_{i=1}^{n} \prod_{t=1}^{T_i} \underbrace{p(R_t^i | S_t^i, A_t^i, Z_t^i)}_{\text{Reward Model}} \underbrace{\pi(A_t^i | S_t^i, Z_t^i)}_{\text{Action Model}}$$
$$\underbrace{p(S_t^i | S_{t-1}^i, A_{t-1}^i, Z_t^i)}_{\text{Transition Model}} \underbrace{p(Z_t^i | Z_{<t}^i, H_{<t}^i, Z^{<i}, H^{<i})}_{\text{Posterior Model}}, \tag{23}$$

606 the lower-bound in Eq. 2 can then be easily derived by doing importance sampling on the posterior
607 model with $q$, marginalizing out the latent variables, and assuming that $H^i$ is independent of all
608 other variables given the latent-variable $Z^i$,

$$\ln p(H^{1:n}) = \ln \int \prod_{i=1}^{n} p(H^i, z^i | H^{<i}, z^{<i}) dz^{1:n} \tag{24}$$

$$= \ln \int \prod_{i=1}^{n} p(H^i | z^i) \frac{q(z^{1:i} | H^{<i})}{q(z^{1:i} | H^{<i})} p(z^i | H^{<i}, z^{<i}) dz^{1:n} \tag{25}$$

$$\geq \int \sum_{i=1}^{n} q(z^{1:i} | H^{1:i}) \left( \ln p(H^i | z^i) - \ln \frac{q(z^{1:i} | H^{<i})}{p(z^i | H^{<i}, z^{<i})} \right) dz^{1:n} \tag{26}$$

$$\propto \mathbb{E}_{q(Z^{1:n} | H^{1:n})} \sum_{i=1}^{n} \ln p(H^i | Z^i) - KL(q(Z^i | Z^{<i}, H^{<i}) \| p(Z^i | Z^{<i}, H^{<i})). \tag{27}$$

609 As stated in the paper, this lower bound only reproduces the data but does not maximize the rewards
610 per se. So, using the control as inference framework (Levine, 2018), if we write the conditional
611 that a given trajectory $H$ is desirable as $p(\mathcal{O} = 1 | H) \propto \exp(\sum_{t=1}^{T} R_t)$, then we can derive a lower
612 bound for the sampling distribution for a reinforcement learning agent as (again simplifying notation

613 for $\mathcal{L}$ compared to the main text),

$$\ln p(\mathcal{O} = 1) = \ln \mathbb{E}_{q(H^{1:n}, Z^{1:n})} p(\mathcal{O} = 1 | H^{1:n}, Z^{1:n}) \frac{p(H^{1:n}, Z^{1:n})}{q(H^{1:n}, Z^{1:n})} \tag{28}$$

$$\geq \mathbb{E}_{q(H^{1:n}, Z^{1:n})} \ln p(\mathcal{O} = 1 | H^{1:n}) - KL(q(H^{1:n}, Z^{1:n}) \| p(H^{1:n}, Z^{1:n})), \tag{29}$$

$$= \mathcal{L}(q) \tag{30}$$

614 where we define the variational distribution $q(H^{1:n}, Z^{1:n})$ to factorize in exactly the same way
615 as $p(H^{1:n}, Z^{1:n})$ where we fix the reward and transition models and then modify the action and
616 posterior models. This choice of factorization cancels out the fixed terms in the KL-divergence,
617 giving us the lower bound (Eq.(4) in the main-text),

$$\mathcal{L}(q) = \mathbb{E}_{q(H^{1:n}, Z^{1:n})} \sum_{i=1}^{n} \sum_{t=1}^{T_i} R_t^i - KL\left(q(A_t^i | S_t^i, Z_t^i) \| \pi(A_t^i | S_t^i, Z_t^i)\right) \tag{31}$$

$$- KL\left(q(Z_t^i | Z_{<t}^i, H_{<t}^i, Z^{<i}, H^{<i}) \| p(Z_t^i | Z_{<t}^i, H_{<t}^i, Z^{<i}, H^{<i})\right). \tag{32}$$

618 To amortize computation of this lower bound and make this practical to compute, we parametrize
619 the variational posterior $q_\theta(Z | \dots)$ and action model $\pi_\psi(A | \dots)$. To then finally give us a practical
620 optimization objective for the parameters $\theta, \psi$, we substitute for the action model $\pi(A | \dots) = \pi_{\psi_{old}}$
621 and for the true posterior we simply use $q_\theta(Z | \dots)$ with a stop-gradient $\square$. We scale the KL-penalty
622 with a hyperparameter $\beta \in \mathbb{R}_+$ to account for differences in scaling. This gives us our final lower-
623 bound,

$$\mathcal{L}(\theta, \psi) = \mathbb{E}_{q_\theta(H^{1:n}, Z^{1:n})} \sum_{i=1}^{n} \sum_{t=1}^{T_i} R_t^i - KL\left(\pi_\psi(A_t^i | S_t^i, Z_t^i) \| \pi_{\psi_{old}}(A_t^i | S_t^i, Z_t^i)\right) \tag{33}$$

$$- \beta \cdot KL\left(q_\theta(Z_t^i | Z_{<t}^i, H_{<t}^i, Z^{<i}, H^{<i}) \| \square q_\theta(Z_{t-1}^i | Z_{<t-1}^i, H_{<t-1}^i, Z^{<i}, H^{<i})\right),$$

624 which we can plug into the l.h.s. of the amortization objective in Eq. (3) enabling us to optimize our
625 model parameters through sampling and end-to-end differentiation from the policy to the posterior.
626 Although this is the objective we desire, we make further heuristic approximations through the
627 use of the Proximal Policy Optimization algorithm (Schulman et al., 2017). This could roughly
628 be interpreted as doing expectation-propagation (Bishop, 2007) on the policy (i.e., swapping the
629 KL-arguments for the policies).

630 Our lower bound is an extension of the one by Abdolmaleki et al. (2018), for standard Markov deci-
631 sion processes, to include the latent variable posterior for use in memory-based meta-reinforcement
632 learning (Duan et al., 2016). When using an RNN to approximate the posterior, the KL-penalty for
633 $q_\theta(Z | \dots)$ is typically ignored since this is undefined for point-estimates. Doing this would recover
634 the RL$^2$ objective in combination with MPO (Abdolmaleki et al., 2018; Duan et al., 2016).

### A.2 Posterior Factorization

636 To choose an efficient factorization for our variational model we need the following result,

637 **Lemma 1.** *We can write $p(Z | \{X_i\}_{i=1}^n) = \frac{1}{p(Z)^{n-1}} \prod_{i=1}^n p(Z | X_i)$ iff $X_i \perp\!\!\!\perp X_j, \forall j \neq i$.*

638 *Proof.* This result can be shown by applying Bayes rule then factorizing each $X_i$ to be independent
639 of $X_j, \forall j \neq i$ and then applying Bayes rule again,

$$p(Z|\{X_i\}_{i=1}^n) = \frac{p(X_1, X_2, \ldots, X_n|Z)p(Z)}{p(X_1, X_2, \ldots, X_n)} \tag{34}$$

$$= \frac{p(Z)}{\prod_{i=1}^n p(X_i)} \prod_{i=1}^n p(X_i|Z) \qquad \text{(Independence)}$$

$$= \frac{p(Z)}{\prod_{i=1}^n p(X_i)} \left[ \prod_{i=1}^n p(Z|X_i) \frac{p(X_i)}{p(Z)} \right] \tag{35}$$

$$= \frac{p(Z)}{\prod_{i=1}^n p(X_i)} \left[ \frac{\prod_{i=1}^n p(X_i)}{p(Z)^n} \prod_{i=1}^n p(Z|X_i) \right] \tag{36}$$

$$= \frac{1}{p(Z)^{n-1}} \prod_{i=1}^n p(Z|X_i) \tag{37}$$

640 $\square$

## A.3 Laplace Variational Recurrent Model

642 **Proposition 1.** *Given a mean-field assumption on the data for our posterior $q_\theta$ (Lemma 1). The*
643 *second order Taylor Expansion of $\ln q_\theta(Z_t|H_{<t}, \phi_t)$ linearized at $\phi_t$, where $\phi_t = \phi^*$ is a local*
644 *maximizer of $q_\theta$ and occupies the same space as $Z_t$, yields the following Gaussian distribution,*

$$q_\theta(Z_t|H_{<t}, \phi_t) = \mathcal{N}\left(Z_t; \mu = \phi_t, \Sigma = (-\nabla_\phi^2 \ln q_\theta(Z_t|H_{<t}, \phi)|_{\phi=\phi_t})^{-1}\right). \tag{38}$$

645 *Proof.* Reiterating the results from the main paper, we choose to factorize our model as,

$$q_\theta(Z_t|H_{<t}, \phi_t) = \frac{1}{q_\theta(Z_t|\phi_t)^{t-2}} \prod_{i=1}^{t-1} q_\theta(Z_t|S_i, R_i, A_i, \phi_t), \qquad \text{(Lemma 1)}$$

$$= \exp\left[ (2-t)\ln q_\theta(Z_t|\phi_t) + \sum_{i=1}^{t-1} \ln q_\theta(Z_t|S_i, R_i, A_i, \phi_t) \right] \tag{39}$$

$$= \exp h_\theta(Z_t; H_{<t}, \phi_t), \tag{40}$$

646 where our aim is to make a local approximation to $h_\theta$, the rest of the proof follows Appendix A from
647 Daxberger et al. (2021).

648 The second order Taylor expansion of $h_\theta(Z_t; H_{<t}, \phi_t)$ where $\phi_t$ (locally) maximizes $h_\theta$ keeping all
649 other arguments fixed, gives us,

$$\hat{h}_\theta(Z_t; H_{<t}, \phi) = h_\theta(Z_t; H_{<t}, \phi_t) + \underbrace{\nabla_\phi h_\theta|_{\phi=\phi_t}(\phi - \phi_t)}_{= 0} + \frac{1}{2}(\phi - \phi_t)^\top \nabla_\phi^2 h_\theta|_{\phi=\phi_t}(\phi - \phi_t),$$

$$= h_\theta(Z_t; H_{<t}, \phi_t) - \frac{1}{2}(\phi - \phi_t)^\top \nabla_\phi^2 h_\theta|_{\phi=\phi_t}(\phi - \phi_t), \tag{41}$$

650  dropping the function arguments to $h_\theta$ for the higher-order terms for brevity. When exponentiating
651  $\hat{h}_\theta$ and renormalizing it to integrate to 1, it is easy to show that we recover a Gaussian,

$$h_\theta(Z_t; H_{<t}, \phi_t) \approx \frac{1}{\int_{\mathbb{R}} \exp \hat{h}_\theta(Z_t; H_{<t}, \phi')d\phi'} \exp \hat{h}_\theta(Z_t; H_{<t}, \phi) \tag{42}$$

$$= \frac{\exp\{h_\theta(Z_t; H_{<t}, \phi_t) - \frac{1}{2}(\phi - \phi_t)^\top (-\nabla_\phi^2 h_\theta|_{\phi=\phi_t})(\phi - \phi_t)\}}{\int_{\mathbb{R}} \exp\{h_\theta(Z_t; H_{<t}, \phi_t) - \frac{1}{2}(\phi' - \phi_t)^\top (-\nabla_\phi^2 h_\theta|_{\phi=\phi_t})(\phi' - \phi_t)\}d\phi'} \tag{43}$$

$$= \frac{\exp\{-\frac{1}{2}(\phi - \phi_t)^\top (-\nabla_\phi^2 h_\theta|_{\phi=\phi_t})(\phi - \phi_t)\}}{\int_{\mathbb{R}} \exp\{-\frac{1}{2}(\phi' - \phi_t)^\top (-\nabla_\phi^2 h_\theta|_{\phi=\phi_t})(\phi' - \phi_t)\}d\phi'} \tag{44}$$

$$= \mathcal{N}\left(Z; \mu = \phi_t, \Sigma = (-\nabla_\phi^2 \ln q_\theta(Z_t|H_{<t}, \phi)|_{\phi=\phi_t})^{-1}\right). \tag{45}$$

652  $\square$

653  Observe that the above result technically gives us a distribution over $\phi$, and misleadingly not $Z$.
654  However, this is simply a consequence of our notation and distributional assumptions for $q$. In
655  practice, $\phi$ is computed by our representation model (recurrent neural network), and our probabilistic
656  model conflates the learned representation $\phi$ as the mode of the distribution for the latent distribution
657  $Z$.

658  **Proposition 2.** *If we choose $q_\theta(Z_t|S_i, R_i, A_i, \phi) = \mathcal{N}(Z_t; \mu = f_\theta(S_i, R_i, A_i; \phi), \Sigma = I_n)$ and*
659  *$q_\theta(Z_t|\phi) = \mathcal{N}(Z_t; \mu = \phi, \Sigma = \sigma_\phi^2 I_n)$ where we take the limit for $\sigma_\phi^2$ to infinity, then our Laplace*
660  *approximated posterior (Proposition 1) has an inverse covariance that is computed as,*

$$\Sigma_t^{-1} = \sum_{i=1}^{t-1} (\nabla_\phi f_\theta(S_i, R_i, A_i; \phi)|_{\phi=\phi_t})(\nabla_\phi f_\theta(S_i, R_i, A_i; \phi)|_{\phi=\phi_t})^\top. \tag{46}$$

661  *Proof.* To see this we only need to write down the Hessian under a local maximum assumption of
662  $\phi = \phi^*$ (Laplace approximation) and substitute the chosen Gaussian distributions in for all terms.

$$\nabla_\phi^2 \ln q_\theta(Z_t|H_{<t}, \phi) = \nabla_\phi^2 \left[ (2-t) \ln q_\theta(Z_t|\phi) + \sum_{i=1}^{t-1} \ln q_\theta(Z_t|S_i, R_i, A_i, \phi) \right] \tag{47}$$

$$= \nabla_\phi^2 \left[ (2-t) \ln \mathcal{N}(Z_t; \phi, \sigma_\phi^2 I_n) + \sum_{i=1}^{t-1} \ln \mathcal{N}(Z_t; f_\theta(S_i, R_i, A_i; \phi), I_n) \right] \tag{48}$$

$$= \frac{t-2}{\sigma_\phi^2} I_n + \sum_{i=1}^{t-1} (J_\phi)_i \underbrace{\left( \nabla_\mu^2 \ln \mathcal{N}(Z_t; (f_\theta)_i, I_n) \right)}_{=-1} (J_\phi)_i^\top$$

$$+ \underbrace{\sum_{i=1}^{t-1} \nabla_\phi^2 (f_\theta)_i|_{\phi=\phi_t} \nabla_\mu \ln \mathcal{N}(Z_t; (f_\theta)_i, I_n)}_{=0, \text{ when } \phi_t=\phi^* \text{ (Laplace approx.)}} \tag{49}$$

$$= \frac{t-2}{\sigma_\phi^2} I_n - \sum_{i=1}^{t-1} (J_\phi)_i (J_\phi)_i^\top, \tag{50}$$

20

where we abbreviate $(J_\phi)_i = \nabla_\phi f_\theta(S_i, R_i, A_i; \phi)$ and $(f_\theta)_i = f_\theta(S_i, R_i, A_i; \phi)$. Then, in the case of using an infinite variance Gaussian for the prior, $\lim_{\sigma_\phi^2 \to \infty} \nabla_\phi^2 \ln q_\theta(Z_t | H_{<t}, \phi)$, we get,

$$\Sigma_t^{-1} = \left( -\nabla_\phi^2 \ln q_\theta(Z_t | H_{<t}, \phi)|_{\phi=\phi_t} \right) = \sum_{i=1}^{t-1} (J_\phi)_i (J_\phi)_i^\top \tag{51}$$

$$= \sum_{i=1}^{t-1} (\nabla_\phi f_\theta(S_i, R_i, A_i; \phi)|_{\phi=\phi_t})(\nabla_\phi f_\theta(S_i, R_i, A_i; \phi)|_{\phi=\phi_t})^\top. \tag{52}$$

$\square$

### A.3.1 Final Model

To complete our fully general Laplace approximated variational recurrent neural network, we need to define the posterior over *all* previous latent variables, $q_\theta(Z_t | Z_{<t}, H_{<t})$. We can easily plug this dependency in for our Laplace approximation from Prop. 1 by assuming that each consecutive posterior has an additive effect on all future posteriors (i.e., a Gaussian convolution),

$$q_\theta(Z_t | Z_{<t}, H_{<t}) = \mathcal{N} \left( Z_t; \mu_t = \phi_t + \sum_{i=1}^{t-1} Z_i, \Lambda_t = -\nabla_\phi^2 \ln q_\theta(Z_t | Z_{<t}, H_{<t}, \phi_t)|_{\phi=\phi_t} \right). \tag{53}$$

As long as we do not condition $\phi_t$ on $Z_{<t}$, the dependency on past latent variables becomes a constant w.r.t. $\nabla_\phi^2$, making the covariance independent of these terms. This is in contrast to a recurrent state-space model architecture which does condition on these values (Hafner et al., 2020), however, our choice permits an analytical solution. It is a known result that the expected posterior then becomes another Gaussian with the means and inverse covariances summed (Bromiley, 2003),

$$\mathbb{E}_{q_\theta(Z_{<t}|H_{<t})} q_\theta(Z_t | Z_{<t}, H_{<t})$$
$$= \mathcal{N} \left( \mu_t = \phi_t + \sum_{i=1}^{t-1} \mu_i, \Lambda_t = -\nabla_\phi^2 \ln q_\theta(Z | H_{<t}, \phi)|_{\phi=\phi_t} + \sum_{i=1}^{t-1} \Lambda_i \right). \tag{54}$$

This particular form has also been described by Ritter et al. (2018) in the context of continual learning. It is easy to accumulate the mean and covariances terms sequentially over $t$. Depending on the assumptions one makes on the data-generating distribution, one can sum over fewer terms to make the calculation more efficient. The ones we ran experiments for in the main paper include:

1. **Stationary Posterior**: $q_\theta(Z_t | H_{<t})$. Full summation over $H_{<t}$ in the calculation of $\Lambda_t$. No summation over previous posteriors $Z_{<t}$.

2. **Markov Chain Posterior**: $q_\theta(Z_t | H_{t-1}, Z_{t-1})$. The inverse covariance is calculated only using the most recent observation $t-1$. Only the previous posterior mean and precision are summed with the current mean and precision.

3. **Windowed Markov Chain Posterior**: $q_\theta(Z_t | H_{k:t-1}, Z_{t-1})$. The inverse covariance is calculated using the $k$ most recent observations. Only the previous posterior mean and precision are summed with the current mean and precision.

However, these are all simplified models, whereas the variational recurrent model (Chung et al., 2015) we discuss in the main paper is fully general.

**On summation of the means** The current formulation for the probabilistic model mismatches with typical recurrent neural network (RNN) architectures. Typically, RNN models do not aggregate their past hidden states for the outputs. However, as discussed in the main paper, we strictly required this simplification for the Markov chain (non-stationary) factorizations, since this was the only approach that would leave the base RNN architecture untouched. Despite that, we can make two heuristic arguments that can partially explain the effect of summing all the previous covariances but omit summation of the mean:

697 • **Representation Learning:** When using covariance summation throughout training (no post-hoc
698 posterior, or finetuning of deterministic baselines), an RNN can learn to represent the hidden state
699 aggregation implicitly within the state-update.

700 • **Exponential Tilting:** Omitting the mean summation can be interpreted as a form of exponential
701 tilting of Gaussian distributions. This is an importance-sampling technique for rare-event simu-
702 lation (Asmussen & Glynn, 2007). For the correctly chosen tilting parameter, this can have the
703 same effect as subtracting all previous means (at the cost of some bias).

## B  Implementation Details

### B.1  Model Architecture and Optimization

706 Following the main text we can define our model according to the following components,

$$
\begin{aligned}
\text{Embedding} \qquad & S_t^g, A_t^g, R_t^g = g_\theta(S_t), h_\theta(A_t), w_\theta(R_t), \\
\text{Recurrent Model} \qquad & \phi_{t+1} = f_\theta(S_t^g, A_t^g, R_t^g; \phi_t), \\
\text{Posterior Model} \qquad & Z_t \sim q_\theta(Z_t | H_{<t}, \phi_t), \\
\text{Reward Model} \qquad & \hat{R} \sim p_\theta(\hat{R} | Z_t, A^g, S_t^g), \\
\text{Action Model} \qquad & A_t \sim \pi_\theta(A_t | Z_t, S_t^g),
\end{aligned}
$$

707 note that we do not train an action model for the supervised experiments, and we do not use the
708 reward model in the reinforcement learning experiments (i.e., we do not use it to select actions or to
709 do planning). In the supervised case, the reward model simply learns a direct function prediction $\hat{R}$
710 where the state can be considered stationary. For brevity, we denoted the full set of parameters as
711 $\theta$, in practice each component has its parameters but we jointly optimize for these using end-to-end
712 differentiation.

713 For the embedding model we used a multi-layered perceptron (MLP) (Cybenko, 1989) of two hidden
714 layers of width 256 nodes, we used leaky-ReLU for the activation. The action and predictive model
715 used a three hidden layer MLP with sizes (256, 256, 64), also with leaky-ReLU. For the recurrent
716 model, we use a long-short-term memory module (Hochreiter & Schmidhuber, 1997) with $n = 128$
717 hidden nodes. We projected the outputs (not the carried state) of the LSTM to a smaller $n = 64$
718 feature output vector with a learned affine transform (i.e., a 1-layer MLP without an activation).

719 For discrete environments, the action model predicted the $n$ logits for the full action space. For the
720 regression task, the reward model outputs a Gaussian with a learned mean and input-independent
721 variance.

722 As noted in the main paper, we apply stop-gradients on the prior term appearing in the KL-
723 divergences (Eq. 2 and Eq. 4). Doing this prevents past posteriors from fitting to data beyond their
724 respective timestep, while still constraining our current posterior to not deviate too much from these
725 terms. We also apply stop-gradients to the *past* mean and covariance terms when *accumulating* these
726 terms. This is only relevant for the Laplace VRNN ablations. The reasoning for this is the same as
727 for the stop-gradient in the KL term, we want to use the past means and covariances as constants
728 at timestep $t$, and not as another learnable parameter. As a side note, we also found that doing this
729 sped up training orders of magnitude.

730 We developed everything discussed in this paper in Jax v.0.4.23 (Bradbury et al., 2018). For neural
731 network design, we used the Flax library v0.8.1 (DeepMind et al., 2020). For optimization, we
732 used the Adamw optimizer (Loshchilov & Hutter, 2019) implemented in Optax with learning-rate
733 $= 10^{-3}$, weight-decay $= 10^{-6}$, and the rest on default settings at version v0.1.7. We used gradient-
734 clipping to have a max global norm of 1.0 and a maximum individual gradient of [-5.0, 5.0]. We
735 used our implementation for the PPO algorithm with help from the RLax library to compute the
736 generalized advantage estimators. For a full list of dependencies, requirements, and program flags,
737 our code will be made available upon publication.

### B.2  Variational Posterior Baseline

As discussed in the main paper, typically in meta-reinforcement learning we see that the posterior $q_\theta(Z_t|H_{<t})$ is modeled with a point-estimate and we propose to use the Laplace approximation to convert this point-estimate into a Gaussian distribution. Instead of computing the covariance matrix using the Laplace approximation, we can also directly predict this covariance with another neural network.

So, our variational recurrent neural network (VRNN) baseline predicted the $m(m-1)$ lower triangular elements of the $m \times m$ dimensional covariance matrix (or just $m$ for the diagonal ablations) and the $m$ dimensional mean. We opted for a spectral decomposition $\Sigma = USV$ where $S$ is diagonal and $U$ was predicted through the Cayley map starting from a skew-symmetric matrix. First, we compute $\phi_t$ with our RNN, then we project $\phi_t$ to $\mu_t, S_t, L_t$ with a linear layer such that $\mu_t \in \mathbb{R}^m$, $S_t \in \mathbb{R}^{m \times m}$ is diagonal and $L_t \in \mathbb{R}^{m \times m}$ is lower-triangular, we then compute $U_t = (I - A_t)(I + A_t)^{-1}$ where $A_t = L_t - L_t^\top$ to get an orthogonal eigenbasis. We then construct the Gaussian as,

$$q_\theta(Z_t|H_{<t}, \phi_t) = \mathcal{N}(Z_t; \mu = \mu_t, \Sigma = U_t(\exp S_t)U_t^\top), \tag{55}$$

this representation also made matrix inversion incredibly easy as $(Ue^S V)^{-1} = Ue^{-S}V$ since $U = V^\top$ are orthogonal.

The reason for using the spectral decomposition was that we did not achieve stable training using any variant of the Cholesky factorization for the covariance matrix ($LU$ or $LDU$ decomposition), even if explicitly transforming the eigenvalues to be positive. We suspect that the spectral parameterization trained more stably than the LU or LDL parameterization as the basis matrices $U$ or $V$ are constrained to the group of orthogonal matrices. Therefore, each element in the predicted parameters $L_n$ is also constrained. In contrast, the LDL parameterization leaves the triangular parameters in an unconstrained representation which might enable an unstable representation. However, we did not investigate this problem beyond what was needed to get our method working.

We also did not accumulate the mean or covariances for the VRNN, unlike for the Laplace VRNN, as this model parameterization could simply learn this function instead (or learn to undo this parameterization). The point of our experiments was not to squeeze performance out of our baseline but to have a strong reference for comparison. We also found that the VRNN was highly competitive with the Laplace VRNN.

### B.3  Predictive Ensemble Averaging

Since we sample latent variables $Z$ from our posterior $q_\theta$, when we pass multiple samples through our predictive model or the action model, we obtain multiple distributions for the output modalities. Typically this induces a mixture distribution, however, we simply averaged out either the logits or the means and variances (Wang et al., 2020). This can be interpreted as a normalized Gaussian convolution over the output parameters or simply as a kind of bagging strategy over ensemble members $Z = z^{(i)}, i = 1, \ldots, k$.

The reasoning for opting for ensemble averaging instead of mixture distributions is that this simply achieved more stable training in combination with PPO (Schulman et al., 2017). The main problem was caused by the penalty term of the policy entropy, our implementation did not achieve stable training when approximating this term in any way (so neither with mixtures nor with singular distributions), we only achieved stable training when the entropy term could be computed *exactly*.

We did not investigate in depth why an approximate entropy loss caused training divergence, but we speculate this is due to the problems of training on generated data as discussed in the context of language models by Shumailov et al. (2023). In this case, the tails of the action distribution slowly shrink over time as Monte-Carlo estimation of the entropy might not cover low-likelihood events sufficiently with small sample sizes. This phenomenon is referred to as *model collapse*, not to be confused with *posterior collapse* (Goyal et al., 2017). We suspect that this problem could be

784 reduced by using a proper (approximate) Bayesian inference algorithm, like maximum a posteriori
785 optimization (Abdolmaleki et al., 2018), which essentially removes the instability of reinforcement
786 learning losses by casting policy learning as a supervised learning problem.

## B.4 Environment Design

788 For our testing environments we implemented three problem domains, a 1D function regression
789 problem (Finn et al., 2017), a discrete $n$-armed bandit problem (Duan et al., 2016), and an $n \times n$
790 open grid problem (Zintgraf et al., 2021), as shown in Figure 4. We generated many variations of
791 these environments to learn from by sampling their dependent task parameters.



Figure 4: Visualization of sampled tasks we evaluated our method on. 1) Zero-shot learning of
a function (left), 2) learning a stochastic best-arm selection algorithm (middle), and 3) learning a
deterministic grid exploration agent (right).

792 **Supervised.** For the supervised problem we generated noiseless 1D test-functions on the bounded
793 domain $[-1, 1]$. We did this through a Fourier expansion of $n = 4$ components where we randomly
794 generate amplitudes, phase shifts, and input shifts. We manually tuned the ranges for these param-
795 eters and sampled uniformly random within these ranges. In other words, we can define the joint
796 distribution over a function dataset as,

$$p(X, Y) = \delta(Y = \text{FourierSum}(Y; X, \varphi_{1:n}, c_{\text{shift}}, A_{0:n})) \cdot \text{Unif}(X; -1.0, 1.0), \quad (56)$$

$$c_{\text{shift}} \sim \text{Unif}(0.0, \pi), \varphi_i \sim \text{Unif}(0.0, \pi), A_i \sim \text{Unif}(-1.0, 1.0) \quad (57)$$

797 where the 'FourierSum' is computed in its amplitude-phase form. The training was performed with
798 50 examples per sampled function.

799 **Bandit** To generate bandit problems we used a Dirichlet distribution with $\alpha = 0.2$ during training
800 and $\alpha = 0.3$ during testing (higher $\alpha$ makes the problem more difficult), this gave us a normalized
801 vector of probabilities, $p_n \sim \text{Dir}(\alpha = 1_n \cdot 0.2)$ for which the agent needed to find $\max_i p_i$. Ob-
802 servations (which are equivalent to the rewards) were generated by sampling Bernoulli outcomes
803 given $p_i$. Since each interaction of the agent with the bandit is seen as an episode, the environment
804 returned discount factors of $\gamma = 0$. The training was performed over 50 total interactions.

805 **Gridworld** For the discrete grid environment we closely match the implementation of (Zintgraf
806 et al., 2021). We reimplemented this environment in Jax to benefit from GPU acceleration for the
807 data-generation process. The environment constructs a $n \times n$ open grid for the agent and uniformly
808 randomly initializes the start and goal state (such that they don't overlap). The agent can choose
809 between moving up, down, left, or right, the environment transitions deterministically but does not
810 move the agent if it moves outside the bounds. The agent is rewarded with +1 if it encounters the
811 goal and 0 otherwise, the observations were constructed as two one-hot-encoded vectors for the row
812 and column index. The goal tile is *not* observed. If the agent did not find the goal within $T = 15$
813 steps it would be reset to its starting state and the discount factor would be set to $\gamma = 0$ at that
814 transition. The training was performed over 100 total interactions.

## B.5 Experimental Design and Hyperparameters

The supervised experiments did not provide additional hyperparameters to set, all necessary parameters were learned using an empirical cross-entropy with the data (see the main paper). For the reinforcement learning experiments we needed to set several hyperparameters for PPO (Schulman et al., 2017), these are given in Table 1. We did not use mini batching for our version of recurrent PPO and accumulated the loss over the full trajectories and batches. We found that larger batch sizes gave us faster and more stable learning.

Table 1: Proximal Policy Optimization loss parameters. Note that we use our Recurrent implementation for this algorithm.

| Name | Symbol | Value |
|---|---|---|
| Minibatches | | Full-Batch |
| Batch-Size | | 256 |
| TD-Lambda | $\lambda$ | 0.9 |
| Discount | $\gamma$ | 0.9 |
| Policy-Ratio clipping | $\epsilon$ | 0.2 |
| Standardize Advantages | | False |
| Exact Policy Entropy | | True |
| Value Loss Scale | | 1.0 |
| Policy Loss Scale | | 1.0 |
| Entropy Loss Scale | | 0.1 |

Then our experimental design implied running an exhaustive parameter grid over the domain presented in Table 2. This grid was adjusted over successive experiments to reduce the computational footprint, this was manually tuned to select the parameter values that performed the best for both the baseline and our method. The parameter grid was applied equivalently on all problem domains for $r = 30$ distinct seeds. Although this experiment design is still quite modest, running this full grid induces 144 distinct configurations times 30 repetitions for the Laplace VRNN alone. One single run took on average $1\frac{1}{2}$ hour to complete for the gridworld environment on an A100 80GB NVIDIA GPU. Although, a better learning algorithm other than PPO (e.g., MPO), and minibatch optimizations could probably get the wallclock time down drastically while still achieving similar performance.

For the finetuning experiments we essentially made model snapshots of the deterministic baselines (RNNs) and reran the ablations as shown in Table 2 using the snapshots as starting weights. For each experiment, these snapshots were taken halfway, and three-quarters way during training in terms of the number of weight-updates.

To make some final informal notes on the choices for the parameters,

- We found that scaling the KL-penalties in the lower bound with hyperparameters that were slightly larger than $\beta > 10^{-2}$ caused posterior collapse (Goyal et al., 2017). Meaning, our posterior simply fitted its parameters to always match the prior despite accumulating more data.

- For the regression problem, using multiple samples for $z^{(i)}$ inside the lower bound decreased the predictive loss a lot. Showing that integration over the predictive is effective. Although, this did not result in better test-time performance.

- Using a small buffer size for the history window in the posterior $q_\theta(Z_t|H_{t-k:t-1})$ is more practical since the computation of the Jacobians is the main bottleneck of our method. We found that

Table 2: Proximal Policy Optimization loss parameters. Note that we use our recurrent implementation for this algorithm. All configurations were repeated for $r = 30$ repetitions (distinct random seeds). We also drop configurations that do not induce a valid model, e.g., $k_Z = 0$ and accumulation of $\Sigma$ is not a valid configuration since there is no window to accumulate over.

| Name | Symbol | Value |
|------|--------|-------|
| **Deterministic RNN** | | |
| Posterior Dimensionality | $n$ | $\{32, 64\}$ |
| **Variational RNN** | | |
| Posterior Dimensionality | $n$ | $\{32, 64\}$ |
| Covariance Parameterization | | $\{$Full, Diagonal$\}$ |
| Posterior KL-Penalty | $\beta$ | $\{1.0, 10^{-2}, 10^{-4}\}$ |
| Number of Posterior Samples | $n_Z$ | $\{1, 5\}$ |
| **Laplace VRNN** | | |
| Posterior Dimensionality | $n$ | $\{32, 64\}$ |
| Covariance Parameterization | | $\{$Full, Diagonal$\}$ |
| Posterior KL-Penalty | $\beta$ | $\{1.0, 10^{-2}, 10^{-4}\}$ |
| Number of Posterior Samples | $n_Z$ | $\{1, 5\}$ |
| History Buffer Window | $k_H$ | $\{1, 10\}$ |
| Latent Variable Window | $k_Z$ | $\{0, 1\}$ |
| Accumulation | | $\{(\mu, \Sigma), \Sigma\}$ |

accumulation of only the covariance where each covariance is computed with a window of just 1, $H_{t-1}$ results in the fastest method (in wallclock time) while being on par with many of the ablations.

# C   Supplementary Results
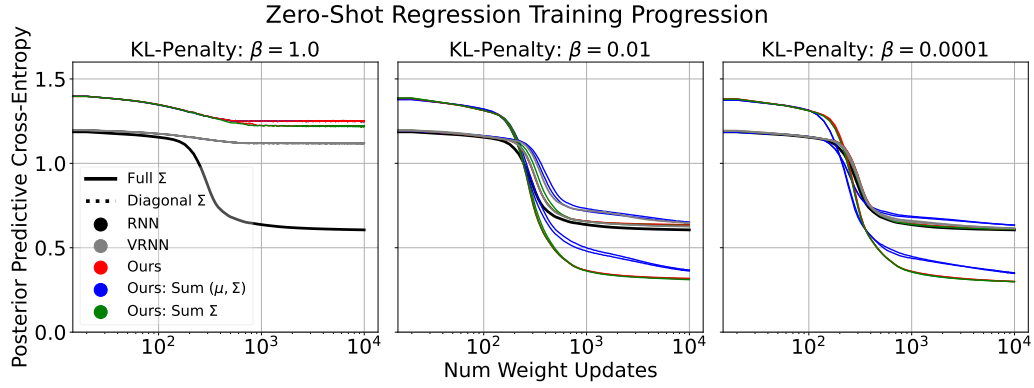
## C.1   Supplementary Supervised Results



Figure 5: Progression of the predictive error during supervised model training of all ablations. Ablations are averaged over parameter groups as indicated by the legend. This figure does not show the finetuning results. To reduce computation time for subsequent results, we picked $\beta = 0.01$ for the reinforcement learning task ablations.



Figure 6: Zoom-in of Figure 5 for three finetune runs (left), for most ablations the predictive error quickly goes down to their full variational training error.

Figure 7: Posterior statistics during testing on the supervised task. The consecutive KL divergences (top) and entropy (bottom) should go down over time. The Laplace VRNN where we sum the covariances (green) is the only method that performs as expected for the entropy estimation but seems to grow more unstable for the KL-divergences. Results use $\beta = 10^{-2}$.
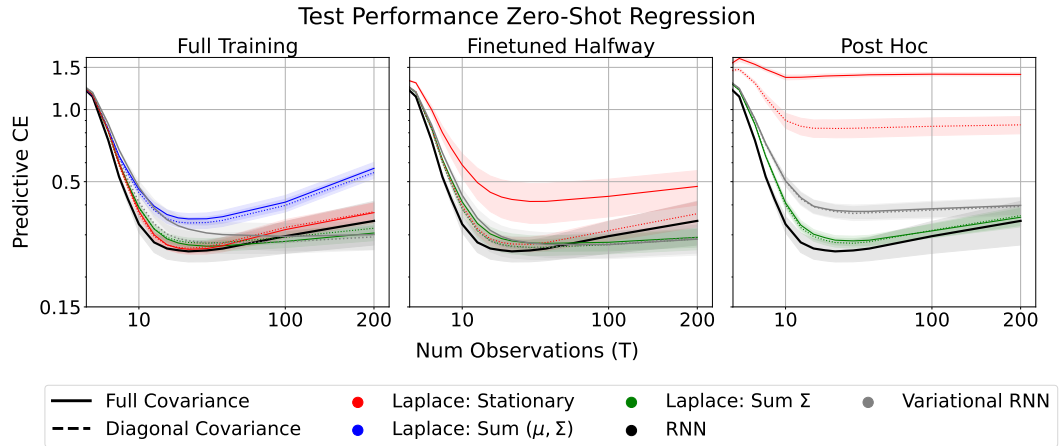


Figure 8: Complete plot of Figure 1 from the main paper to include the accumulation over means and covariances (blue) in the left plot. This was left out of the main paper to improve visibility. Results use $\beta = 10^{-2}$.

**C.2    Supplementary Bandit Results**



Figure 9: Zoom-in of Figure 3 for the bandit task when finetuning the deterministic model weights intermittently with a diagonal variational model. In this domain, it seems that finetuning slightly actually helps the expected training performance. Results use $\beta = 10^{-2}$.
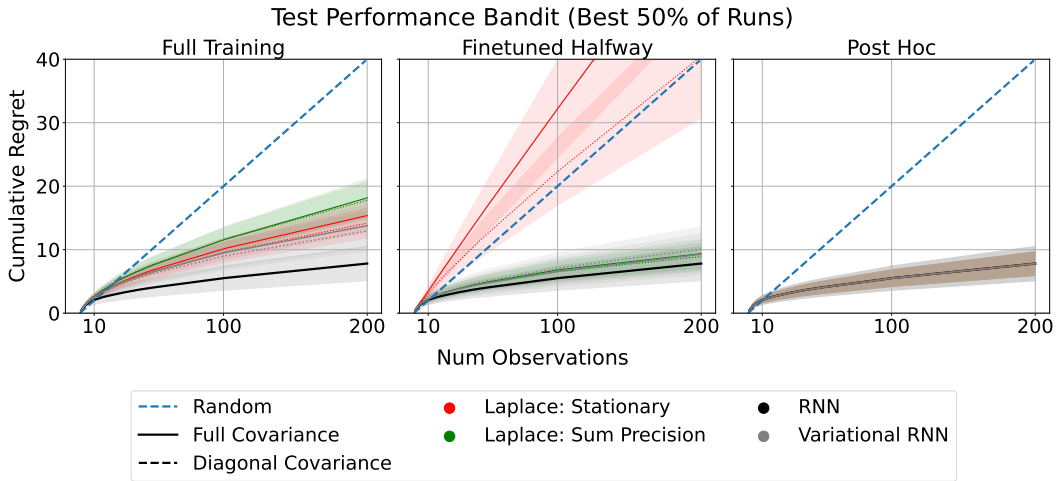


Figure 10: Cumulative regret of trained agents on the bandit task, lower is better. Since the training was quite unstable, about half of the repetitions did not find model weights with strong final performance. Only when we filtered out the better-than-median agents, did we find stronger than random performance. Results use $\beta = 10^{-2}$.

**C.3   Supplementary Gridworld Results**
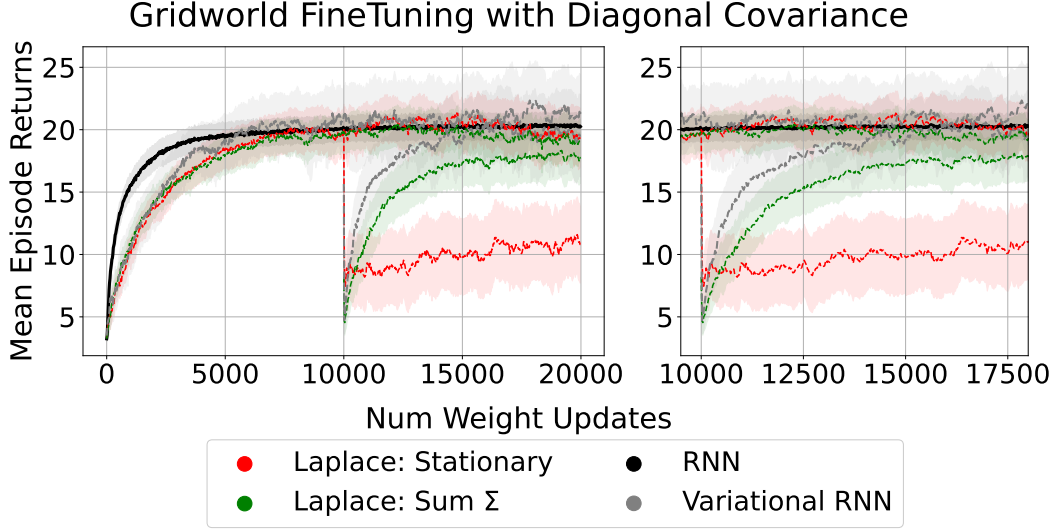
## Gridworld FineTuning with Diagonal Covariance



Figure 11: Zoom-in of Figure 3 for the gridworld task when finetuning the deterministic model weights intermittently with a diagonal variational model. In contrast to the bandit task, the agent needs to recover from this sudden change of additional model noise. Results use $\beta = 10^{-2}$.
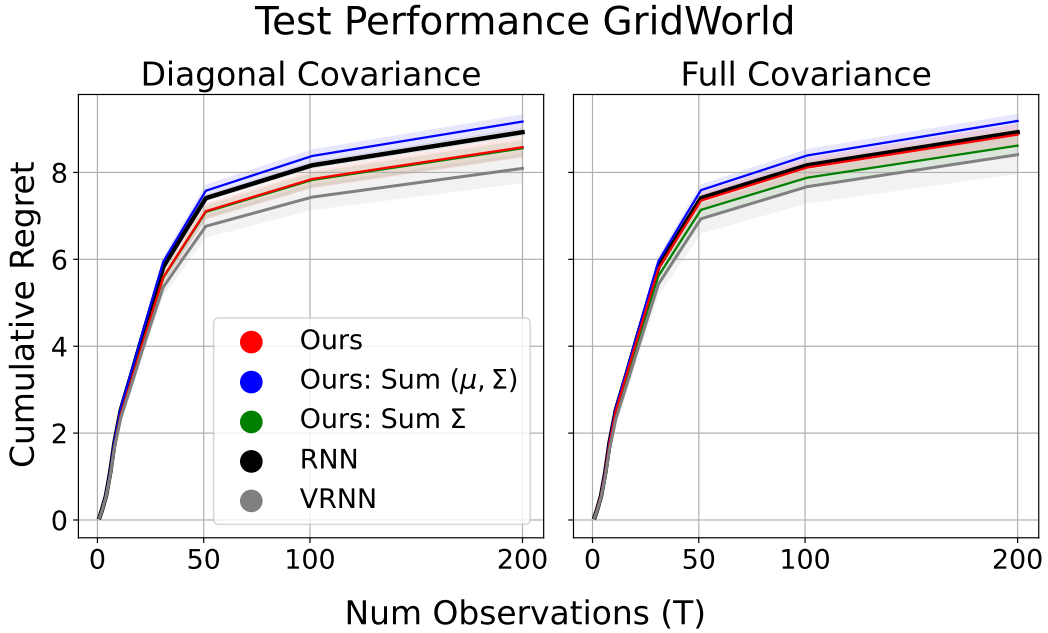
## Test Performance GridWorld



Figure 12: Cumulative regret of trained agents on the gridworld task, lower is better. All agents perform well on this task, the diagonal Variational RNN slightly outperforms all other agents.
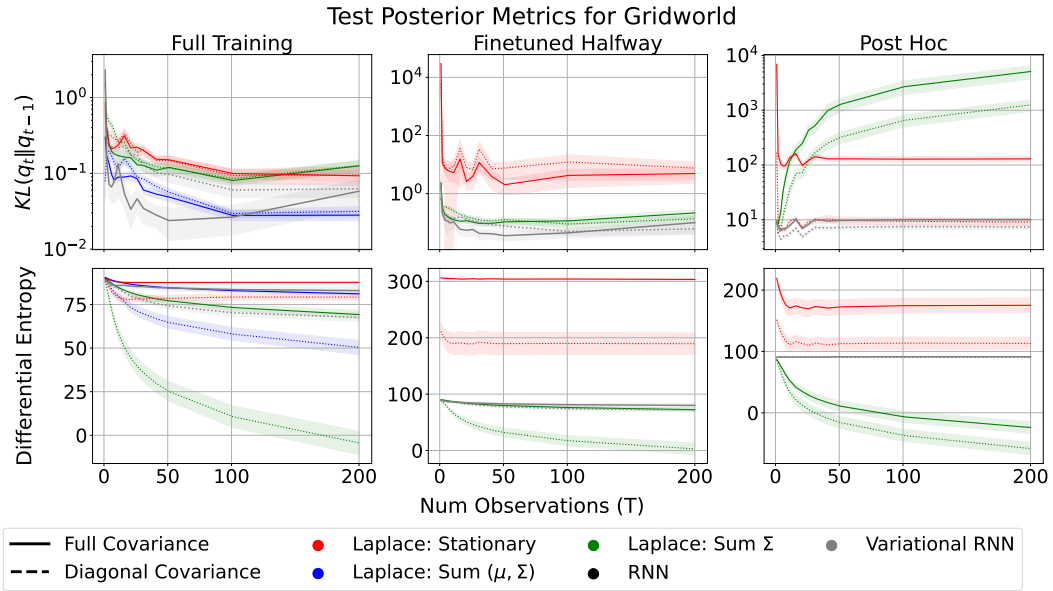
Figure 13: Complete plot of Figure 2 from the main paper to include the accumulation over mans and covariances (blue) in the left plot. Posterior statistics during testing on the gridworld task. The consecutive KL divergences (top) and entropy (bottom) should go down over time. Results use $\beta = 10^{-2}$.