
DynaMITE-RL: A Dynamic Model for Improved Temporal Meta-Reinforcement Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We introduce *DynaMITE-RL*, a meta-reinforcement learning (meta-RL) approach
2 to approximate inference in environments where the latent state evolves at varying
3 rates. We model episode sessions—parts of the episode where the latent state
4 is fixed—and propose three key modifications to existing meta-RL methods: (i)
5 consistency of latent information within sessions, (ii) session masking, and (iii)
6 prior latent conditioning. We demonstrate the importance of these modifications
7 in various domains, ranging from discrete Gridworld environments to continuous-
8 control and simulated robot assistive tasks, illustrating the efficacy of DynaMITE-
9 RL over state-of-the-art baselines in both online and offline RL settings.

10 1 Introduction

11 Markov decision processes (MDPs) [4] provide a general framework in reinforcement learning (RL),
12 and can be used to model sequential decision problems in a variety of domains, e.g., recommender
13 systems (RSs), robot and autonomous vehicle control, and healthcare [22, 21, 7, 46, 31, 5]. MDPs
14 assume a static environment with fixed transition probabilities and rewards [3]. In many real-world
15 systems, however, the dynamics of the environment are intrinsically tied to latent factors subject
16 to temporal variation. While non-stationary MDPs are special instances of partially observable
17 MDPs (POMDPs) [24], in many applications these latent variables change infrequently, i.e. the latent
18 variable remains fixed for some duration before changing. One class of problems exhibiting this latent
19 transition structure is recommender systems, where a user’s preferences are a latent variable which
20 gradually evolves over time [23, 26]. For instance, a user may initially have a strong affinity for a
21 particular genre (e.g., action movies), but their viewing habits could change over time, influenced by
22 external factors such as trending movies, mood, etc. A robust system should adapt to these evolving
23 tastes to provide suitable recommendations. Another example is in manufacturing settings, where
24 industrial robots may experience unobserved gradual deterioration of their mechanical components
25 affecting the overall functionality of the system. Accurately modelling such latent transitions caused
26 by hardware degradation can help manufacturers optimize performance, cost, and equipment lifespan.

27 Our goal in this work is to leverage such a temporal structure to obviate the need to solve a fully general
28 POMDP. To this end, we propose **Dynamic Model for Improved Temporal Meta Reinforcement**
29 **Learning (DynaMITE-RL)**, a method designed to exploit the temporal structure of sessions, i.e.,
30 sub-trajectories within the history of observations in which the latent state is fixed. We formulate our
31 problem as a *dynamic latent contextual MDP (DLCMDP)*, and identify three crucial elements needed
32 to enable tractable and efficient policy learning in environments with the latent dynamics captured by
33 a DLCMDP. First, we consider consistency of latent information, by exploiting time steps for which
34 we have high confidence that the latent variable is constant. To do so, we introduce a consistency loss
35 to regularize the posterior update model, providing better posterior estimates of the latent variable.
36 Second, we enforce the posterior update model to learn the dynamics of the latent variable. This

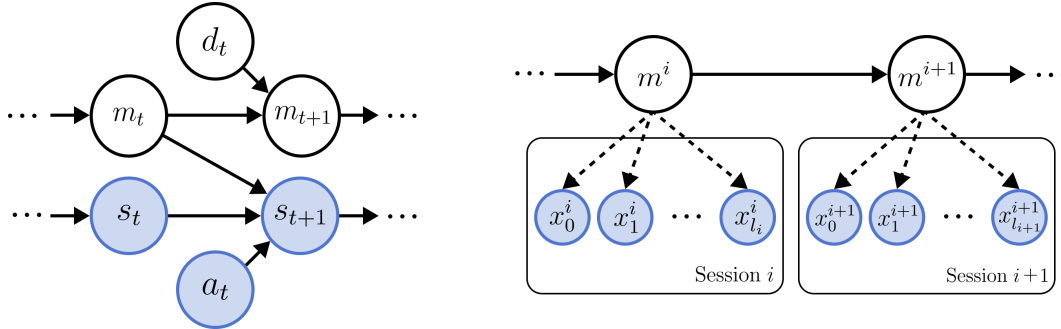


Figure 1: **(Left)** The graphical model for a DLCMDP. The transition dynamics of the environment follows $T(s_{t+1}, m_{t+1} \mid s_t, a_t, m_t)$. At every timestep t , an i.i.d. Bernoulli random variable, d_t , denotes the change in the latent context, m_t . Blue shaded variables are observed, whereas white shaded variables are latent. **(Right)** A realization of a DLCMDP episode. Each session i is governed by a latent variable m^i which is changing between sessions according to a fixed transition function, $T_m(m' \mid m)$. We denote l_i as the length of session i . The state-action pair (s_t^i, a_t^i) at timestep t in session i is summarized into a single observed variable, x_t^i . We emphasize that session terminations are not explicitly observed.

37 allows the trained policy to better infer, and adapt to, temporal shifts in latent context in unknown
 38 environments. Finally, we show that the variational objective in meta-RL algorithms, which attempts
 39 to reconstruct the entire trajectory, can hurt performance when the latent context is nonstationary. We
 40 modify this objective to reconstruct only the transitions that share the same latent context.

41 Closest to our work is VariBAD [47], a meta-RL [1] approach for learning a Bayes-optimal policy,
 42 enabling an agent to quickly adapt to a new environment with unknown dynamics and reward
 43 functions. VariBAD uses variational inference to learn a posterior update model that approximates
 44 the belief over the distribution of transition and reward functions. It augments the state space with
 45 this belief to encode the agent’s uncertainty during decision-making. Nevertheless, VariBAD and the
 46 Bayes-Adaptive MDP framework [35] assume the latent context is static *across an episode* and do
 47 not address settings with latent state dynamics. In this work, we focus on the dynamic latent state
 48 formulation of the meta-RL problem.

49 Our core contributions are as follows: (1) We introduce DynaMITE-RL, a meta-RL approach to
 50 handle environments with evolving latent context variables. (2) We introduce three key elements
 51 for learning an improved posterior update model: session consistency, modeling dynamics of latent
 52 context, and session reconstruction masking. (3) We validate our approach on a diverse set of
 53 challenging simulation environments and demonstrate significantly improved results over multiple
 54 state-of-the-art baselines in both online and offline-RL settings.

55 2 Background

56 We begin by reviewing relevant background including meta-RL and Bayesian RL. We also briefly
 57 summarize the VariBAD [47] algorithm for learning Bayes-adaptive policies.

58 **Meta-RL.** The goal of meta-RL [1] is to quickly adapt an RL agent to an unseen test environment.
 59 Meta-RL assumes a distribution $p(\mathcal{T})$ over possible environments or *tasks*, and learns this distribution
 60 by repeatedly sampling batches of tasks during meta-training. Each task $\mathcal{T}_i \sim p(\mathcal{T})$ is described by
 61 an MDP $\mathcal{M}_i = (\mathcal{S}, \mathcal{A}, R_i, T_i, \gamma)$, where the state space \mathcal{S} , action space \mathcal{A} , and discount factor γ are
 62 shared across tasks, while R_i and T_i are task-specific reward and transition functions, respectively.
 63 The objective of meta-RL is to learn a policy that efficiently maximizes reward given a new task
 64 $\mathcal{T}_i \sim p(\mathcal{T})$ sampled from the task distribution at meta-test time. Meta-RL is a special case of
 65 a POMDP in which the unobserved variables are R and T , which are assumed to be stationary
 66 throughout an episode.

67 **Bayesian Reinforcement Learning (BRL).** BRL [18] utilizes Bayesian inference to model the
 68 uncertainty of agent and environment in sequential decision making problems. In BRL, R
 69 and T are unknown a priori and treated as random variables with associated prior distributions.

70 At time t , the *observed history* of states, actions and re-
71 wards is $\tau_{:t} = \{s_0, a_0, r_1, \dots, r_t, s_t\}$, and the belief b_t
72 represents the posterior over task parameters R and T
73 given the transition history, i.e. $b_t \triangleq p(R, T \mid \tau_{:t})$. Given
74 the initial belief $b_0(R, T)$, the belief can be updated iteratively
75 using Bayes’ rule: $b_{t+1} = p(R, T \mid \tau_{:t+1}) \propto$
76 $p(s_{t+1}, r_{t+1} \mid \tau_{:t}, R, T) \cdot b_t$. This Bayesian approach to
77 RL can be formalized as a *Bayes-adaptive MDP (BAMDP)*
78 [14]. A BAMDP is an MDP over the *augmented state*
79 *space* $S^+ = S \times \mathcal{B}$, where \mathcal{B} denotes the belief space. Given
80 the augmented state $s_t^+ = (s_t, b_t)$, the transition function is
81 given by $T^+(s_{t+1}^+ \mid s_t^+, a_t) = \mathbb{E}_{b_t}[T(s_{t+1} \mid s_t, a_t) \cdot \delta(b_{t+1} =$
82 $p(R, T \mid \tau_{:t+1})]$, and reward function is the expected re-
83 ward given the belief, $R^+(s_t^+, a_t) = \mathbb{E}_{b_t}[R(s_t, a_t)]$. The
84 BAMDP formulation naturally resolves the exploration-
85 exploitation tradeoff. A Bayes-optimal RL agent takes
86 information-gathering actions to reduce its uncertainty in
87 the MDP parameters while simultaneously maximizing its
88 returns. However, for most interesting problems, solving
89 the BAMDP—and even computing posterior updates—is
90 intractable given the continuous and typically high-
91 dimensional nature of its state space.

92 **VariBAD.** Zintgraf et al. [47] approximates the Bayes-optimal solution by modeling uncertainty over
93 the MDP parameters. These parameters are represented by a latent vector $m \in \mathbb{R}^d$, the posterior over
94 which is $p(m \mid \tau_{:H})$, where H is the BAMDP horizon. VariBAD uses a variational approximation
95 $q_\phi(m \mid \tau_{:t})$ parameterized by ϕ and conditioned on the observed history up to time t . Zintgraf
96 et al. [47] show that $q_\phi(m \mid \tau_{:t})$ approximates the belief b_t . In practice, $q_\phi(m \mid \tau_{:t})$ is represented
97 by a Gaussian distribution $q_\phi(m \mid \tau_{:t}) = \mathcal{N}(\mu(\tau_{:t}), \Sigma(\tau_{:t}))$, where μ and Σ are sequence models
98 (e.g., recurrent neural networks or transformers [42]) that encode trajectories to latent statistics. The
99 variational lower bound at time t is $\mathbb{E}_{q_\phi(m \mid \tau_{:t})}[\log p_\theta(\tau_{:H} \mid m)] - D_{KL}(q_\phi(m \mid \tau_{:t}) \parallel p_\theta(m))$, where
100 the first term reconstructs the trajectory likelihood $p_\theta(\tau_{:H} \mid m)$ and the second term regularizes
101 the variational posterior to a prior distribution over the latent space, typically modeled with a
102 standard Gaussian distribution. Importantly, the trajectory up to time t , i.e., $\tau_{:t}$, is used in the
103 ELBO equation to infer the posterior belief at time t , which then decodes the entire trajectory $\tau_{:H}$,
104 including future transitions. Given the belief state distribution q_ϕ of a BAMDP, the policy maps
105 both the state and belief to actions, i.e., $\pi(a_t \mid s_t, q_\phi(m \mid \tau_{:t}))$. The BAMDP solution policy π^*
106 is trained, e.g., via policy gradient methods, to maximize the expected cumulative return of meta-RL:
107 $J(\pi) = \mathbb{E}_{R, T} \left[\mathbb{E}_\pi \left[\sum_{t=0}^{H-1} \gamma^t r(s_t, a_t) \right] \right]$, where the first expectation is averaged over environments.
108 The RL agent is trained jointly with the variational belief distribution q_ϕ .

109 3 Dynamic Latent Contextual MDPs

110 As a special case of a BAMDP, where the belief state is parameterized with a latent context vector
111 (analogous to the problem formulation of VariBAD), the *dynamic latent contextual MDP (DLCMDP)*
112 is denoted by $\langle \mathcal{S}, \mathcal{A}, \mathcal{M}, R, T, \nu_0, H \rangle$, where \mathcal{S} is the state space, \mathcal{A} is the action space, \mathcal{M} is the
113 *latent* context space, $R : \mathcal{S} \times \mathcal{A} \times \mathcal{M} \mapsto \Delta_{[0,1]}$ is a reward function, $T : \mathcal{S} \times \mathcal{A} \times \mathcal{M} \mapsto \Delta_{\mathcal{S} \times \mathcal{M}}$
114 is a transition function, $\nu_0 \in \Delta_{\mathcal{S} \times \mathcal{M}}$ is an initial state distribution, $\gamma \in (0, 1)$ is a discount factor, and
115 H is the (possibly infinite) horizon.

116 We assume an episodic setting in which each episode begins in a state-context pair $(s_0, m_0) \sim \nu_0$. At
117 time t , the agent is at state s_t and context m_t , and has observed history $\tau_{:t} = \{s_0, a_0, r_1, \dots, r_t, s_t\}$.
118 Given the history, the agent selects an action $a_t \in \mathcal{A}$, after which the state and latent context
119 transitions according to $T(s_{t+1}, m_{t+1} \mid s_t, a_t, m_t)$, and the agent receives a reward sampled from
120 $R(s_t, a_t, m_t)$. Throughout this process, the context m_t is latent (i.e., *not observed* by the agent).

121 DLCMDPs embody the causal independence depicted by the graphical model in Figure 1. Particularly,
122 DLCMDPs impose a structure on changes of the latent variable m , allowing the latent context m to
123 change less or more frequently. We denote by d_t the random variable at which a transition occurs in

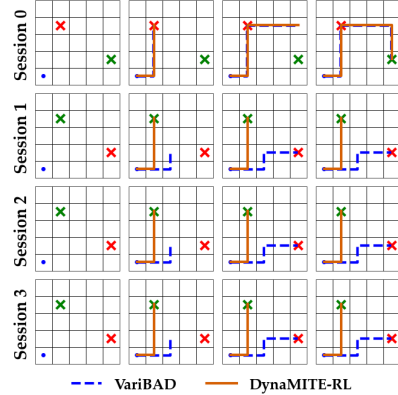


Figure 2: A DLCMDP rollout. VariBAD does not model the transition dynamics of the latent context and fails to adapt to the changing goal location. By contrast, DynaMITE-RL correctly infers the transition and consistently reaches the rewarding cell (green cross).

DynaMITE-RL Training

- 1: **Input:** env, policy, critic, belief model
 - 2: **for** iter = 1 to num_rl_updates **do**
 - 3: Collect DLCMDP episodes
 - 4: Train posterior belief model by maximizing ELBO (Eq. (2))
 - 5: Train policy and critic with any online RL algorithm
 - 6: **end for**
-

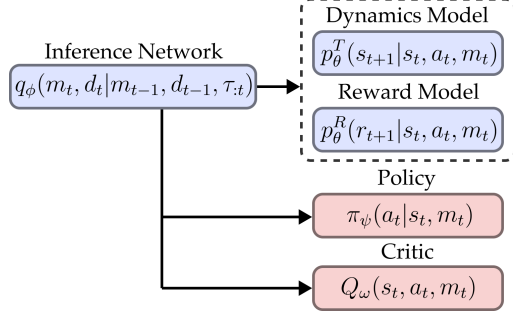


Figure 3: Pseudo-code (online RL training) and model architecture of DynaMITE-RL.

124 m_t . Let $\Omega = \{d_t\}_{t=0}^{H-1}$ denote a sequence of i.i.d. Bernoulli random variables, according to Figure 1,
 125 the transition function T is represented by the following factored distribution:

$$\begin{aligned}
 T(s_{t+1} = s', m_{t+1} = m' \mid s_t = s, a_t = a, m_t = m) \\
 = T_s(s' \mid s, a, m) \mathbb{1}\{m' = m, d_t = 0\} T_d(d_t = 0) + \nu_0(s' \mid m') T_m(m' \mid m) \mathbb{1}\{d_t = 1\} T_d(d_t = 1),
 \end{aligned}$$

126 where $T_m : \mathcal{M} \mapsto \mathcal{M}$ is the latent dynamics function, T_s is the context-dependent state transition
 127 function, and T_d is the termination probability distribution. We refer to sub-trajectories between
 128 changes in the latent context as *sessions*, which may vary in length. At the start of a new session,
 129 a new state and a new latent context are sampled based on the distribution ν_0 . Each session itself
 130 is governed by an MDP parameterized with a latent context $m \in \mathcal{M}$, which changes stochastically
 131 between sessions according to the latent transition function $T_m(m' \mid m)$. For notational simplicity
 132 we use index i to denote the i^{th} session in a trajectory, and m^i the respective latent context of that
 133 session. We emphasize that sessions switching times are latent random variables.

134 Notice that DLCMDPs are more general than latent MDPs [38, 29], in which the latent context is
 135 fixed throughout the entire episode; this corresponds to $d_t \equiv 0$. Moreover, DLCMDPs are closely
 136 related to POMDPs; letting $d_t \equiv 1$, a DLCMDP reduces to a general POMDP with state space \mathcal{M} ,
 137 observation space \mathcal{S} , and observation function ν_0 . As a consequence DLCMDPs are as general as
 138 POMDPs, rendering them very expressive. Moreover, the specific temporal structure of DLCMDPs
 139 allows us to devise efficient learning algorithms that exploit the transition dynamics of the latent
 140 context, improving learning efficiency. DLCMDPs are related to DCMDPs [40], LSMDPs [8], and
 141 DP-MDP [45]. However, DCMDPs assume contexts are observed, and focus on aggregated context
 142 dynamics, LSMDPs assume that the latent contexts across sessions are i.i.d (i.e., there is no latent
 143 dynamics) and DP-MDPs assume that sessions are fixed length.

144 We aim to learn a policy $\pi(a_t \mid s_t, m_t)$ which maximizes the expected return $J(\pi)$ over unseen test
 145 environments. As in BAMDPs, the optimal DLCMDP Q-function satisfies the Bellman equation;
 146 $\forall s^+ \in \mathcal{S}^+, a \in \mathcal{A} : Q(s^+, a) = R^+(s^+, a) + \gamma \sum_{s^{+'} \in \mathcal{S}^+} T^+(s^{+'} \mid s^+, a) \max_{a'} Q(s^{+'}, a)$. In the
 147 following section, we present DynaMITE-RL for learning a Bayes-optimal agent in a DLCMDP.

148 4 DynaMITE-RL

149 We detail DynaMITE-RL, first deriving a variational lower bound for learning a DLCMDP posterior
 150 model, then outlining three principles for training DLCMDPs, and finally integrating them into our
 151 training objective.

152 **Variational Inference for Dynamic Latent Contexts.** Given that we do not have direct access to
 153 the transition and reward functions of the DLCMDP, following Zintgraf et al. [47], we infer the
 154 posterior $p(m \mid \tau_{:t})$, and reason about the latent context vector m instead. Since exact posterior
 155 computation over m is computationally infeasible, given the need to marginalize over task space, we
 156 introduce the variational posterior $q_\phi(m \mid \tau_{:t})$, parameterized by $\phi \in \mathbb{R}^d$, to enable fast inference at
 157 every step. Our learning objective maximizes the log-likelihood $\mathbb{E}_\pi[\log p(\tau)]$ of observed trajectories.
 158 In general, the true posterior over the latent context is intractable, as is the empirical estimate of the

159 log-likelihood. To circumvent this, we derive the *evidence lower bound (ELBO)* [27] to approximate
 160 the posterior over m under the variational inference framework.

161 Let $\mathcal{Z} = \{m^i\}_{i=0}^{K-1}$ be the sequence of latent context vectors for K sessions in an episode (note that K
 162 is inherently a random variable—the exact number of sessions in an episode is not known). As defined
 163 previously, Ω is the collection of the session terminations. We use a parametric generative distribution
 164 model for the state-reward trajectory, conditioned on the action sequence: $p_\theta(s_0, r_1, s_1, \dots, r_H, s_H \mid$
 165 $a_0, \dots, a_{H-1})$. In what follows, we drop the conditioning on $a_{:H-1}$ for the sake of brevity.

166 The variational lower bound can be expressed as:

$$\log p_\theta(\tau) \geq \underbrace{\mathbb{E}_{q_\phi(\mathcal{Z}, \Omega | \tau_{:t})} [\log p_\theta(\tau | \mathcal{Z}, \Omega)]}_{\text{reconstruction}} - \underbrace{D_{KL}(q_\phi(\mathcal{Z}, \Omega | \tau_{:t}) \parallel p_\theta(\mathcal{Z}, \Omega))}_{\text{regularization}} = \mathcal{L}_{\text{ELBO}, t}, \quad (1)$$

167 which can be estimated via Monte Carlo sampling over a learnable approximate posterior q_ϕ . In
 168 optimizing the reconstruction loss of session transitions and rewards, the learned latent variables
 169 should capture the unobserved MDP parameters. The full derivation of the ELBO for a DLCMDP is
 170 provided in Appendix A.1.

171 Figure 2 depicts a (qualitative) didactic GridWorld example with two possible rewarding goals that
 172 alternate between sessions. The VariBAD agent does not account for latent goal dynamics and gets
 173 stuck after reaching the goal in the first session. By contrast, DynaMITE-RL employs the latent
 174 context dynamics model to capture goal changes, and adapts to the context changes across sessions.

175 **Consistency of Latent Information.** In the DLCMDP formulation, each session is itself an MDP
 176 with a latent context fixed across the session. This within-context stationarity means new observations
 177 can only increase the information the agent has about this context. In other words, the agent’s
 178 posterior over latent contexts gradually hone in on the true latent distribution. Although this true
 179 distribution remain unknown, this insight suggest the use of a *session-based consistency loss*, which
 180 penalizes an increase in KL-divergence between the current and final posterior belief within a session.
 181 Let $d_{H-1} = 1$ and $t_i \in \{0, \dots, H\}$ be a random variable denoting the last timestep of session
 182 $i \in \{0, \dots, K-1\}$, i.e., $t_i = \min\{t' \in \mathbb{Z}_{\geq 0} : \sum_{t=0}^{t'} d_t = i + 1\}$. At each time t in session i , we
 183 define the temporal, session-based consistency loss as

$$\mathcal{L}_{\text{consistency}, t} = \max\{D_{KL}(q_\phi(m^i | \tau_{:t+1}) \parallel q_\phi(m^i | \tau_{:t_i})) - D_{KL}(q_\phi(m^i | \tau_{:t}) \parallel q_\phi(m^i | \tau_{:t_i})), 0\},$$

184 where $q_\phi(m^i | \tau_{:t_i})$ is the final posterior in session i . Using temporal consistency to regularize
 185 inference introduces an explicit inductive bias that allows for better posterior estimation.

186 *Remark 4.1.* We introduce session-based consistency for DLCMDPs, though it is also relevant in
 187 single-session settings with non-dynamic latent context. Indeed, as we discuss below, while VariBAD
 188 focuses on single sessions, it does not constrain the latent’s posterior to be identical to final posterior
 189 belief. Consistency may be useful in settings where the underlying latent variable is stationary, but
 190 may hurt performance when this variable is indeed changing. Since our modeling approach allows
 191 latent context changes across sessions, incorporating consistency regularization does not generally
 192 hurt performance.

193 **Latent Belief Conditioning.** Unlike the usual BAMDP framework, DLCMDPs allow one to model
 194 temporal changes of latent contexts via dynamics $T_m(m' | m)$ across sessions. To incorporate this
 195 model into belief estimation, in addition to the history $(\tau_{:t}, d_{:t})$, we condition the posterior on the final
 196 latent belief $q_\phi(m', d' | m, d, \tau_{:t})$ from the previous session, and impose KL-divergence matching
 197 between this belief and the prior distribution $p_\theta(m' | m)$.

198 **Reconstruction Masking.** When the agent is at time t , Zintgraf et al. [47] encode past interactions to
 199 obtain the current posterior $q_\phi(m | \tau_{:t})$ since this is all the information available for inference about
 200 the current task (see Eq. (1)). They use this posterior to decode the entire trajectory—including future
 201 transitions—from different sessions to optimize the lower bound during training. The insight is that
 202 decoding both the past and future allows the posterior model to perform inference about unseen states.
 203 However, we observe that when the latent context is stochastic, reconstruction over the full sequence
 204 is detrimental to training efficiency. The model is attempting to reconstruct transitions outside of the
 205 current session that may be irrelevant or biased given the latent-state dynamics, rendering it a more
 206 difficult learning problem. Instead we reconstruct only the transitions within the session defined by
 207 the predicted termination indicators, i.e., at any arbitrary time t within session i , the session-based
 208 reconstruction loss is given by

$$\mathcal{L}_{\text{session-ELBO}, t} = \mathbb{E}_{q_\phi(\mathcal{Z}, \Omega | \tau_{:t})} [\log p_\theta(\tau_{t_{i-1}+1:t_i} | \mathcal{Z}, \Omega)] - D_{KL}(q_\phi(\mathcal{Z}, \Omega | \tau_{:t}) \parallel p_\theta(\mathcal{Z}, \Omega)).$$

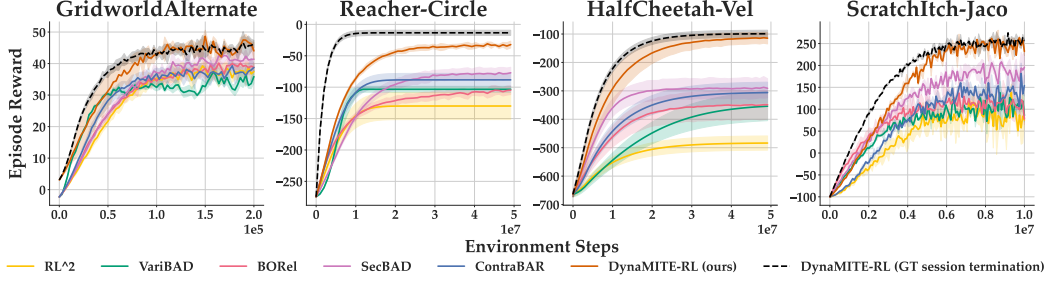


Figure 4: Learning curves for **DynaMITE-RL** and state-of-the-art baseline methods. Shaded areas represent standard deviation over 5 different random seeds for each method and 3 for ScratchItch. In each of the evaluation environments, we observe that **DynaMITE-RL** exhibits better sample efficiency and converges to a policy with better environment returns than the baseline methods.

209 **DynaMITE-RL.** By incorporating the three modifications above, we obtain at the following training
 210 objective for our variational meta-RL approach:

$$\mathcal{L}_{\text{DynaMITE-RL}}(\theta, \phi) = \sum_{t=0}^{H-1} \left[\mathcal{L}_{\text{session-ELBO},t}(\theta, \phi) + \beta \cdot \mathcal{L}_{\text{consistency},t}(\phi) \right], \quad (2)$$

211 where $\beta > 0$ is a hyper-parameter that regularizes the consistency loss. We present a simplified
 212 pseudocode for online training of DynaMITE-RL in Algorithm 3a and a detailed algorithm in
 213 Appendix A.2.

214 **Implementation Details.** We use proximal policy optimization (PPO) [37] for online RL training.
 215 We introduce a posterior inference network that outputs a Gaussian over the latent context for
 216 the i -th session and the session termination indicators, $q_{\phi}(m^i, d_{:t} \mid \tau_{:t}, m^{i-1})$, conditioned on the
 217 history and posterior belief from the previous session. We parameterize the inference network
 218 as a sequence model, with e.g., an RNN [9] or a Transformer [42], with different multi-layer
 219 perceptron (MLP) output heads for predicting the logits for session termination and the posterior
 220 belief. In practice, the posterior MLP outputs the parameters of a Gaussian belief distribution
 221 $q_{\phi_m}(m^i \mid \tau_{:t}, m^{i-1}) = \mathcal{N}(\mu(\tau_{:t}), \Sigma(\tau_{:t}))$. The session termination network applies a sigmoid
 222 activation function $\sigma(x) = \frac{1}{1+e^{-x}}$ to the MLP output. Following PPO [37], the actor loss \mathcal{J}_{π}
 223 and critic loss \mathcal{J}_{ω} are respectively given by $\mathcal{J}_{\pi} = \mathbb{E}_{\tau \sim \pi_{\psi}}[\log \pi_{\psi}(a \mid s, m) \hat{A}(s, a, m)]$ and $\mathcal{J}_{\omega} =$
 224 $\mathbb{E}_{\tau \sim \pi_{\psi}}[(Q_{\omega}(s, a, m) - (r + V_{\omega}(s', m)))^2]$, where V is the target network, and \hat{A} is the advantage
 225 function. We also add an entropy bonus to ensure sufficient exploration in more complex domains.
 226 A decoder network, also parameterized using MLPs, reconstructs transitions and rewards given
 227 the session’s latent context m^i , current state s_t , and action a_t , i.e., $p_{\theta}^T(s_{t+1} \mid s_t, a_t, m_t)$ and
 228 $p_{\theta}^R(r_{t+1} \mid s_t, a_t, m_t)$. Figure 3b depicts the implemented model architecture. The final objective
 229 of DLCMDP is to jointly learn the policy π_{ψ} , the variational posterior model q_{ϕ} , and the factored
 230 likelihood model p_{θ} that minimizes the following loss:

$$\mathcal{L}(\theta, \phi, \psi) = \mathbb{E} \left[\mathcal{J}_{\pi}(\psi) + \lambda \cdot \mathcal{L}_{\text{DynaMITE-RL}}(\phi, \theta) \right], \quad (3)$$

231 where \mathcal{J} is the expected return, and $\lambda > 0$ is a hyper-parameter trades off this return with DynaMITE-
 232 RL’s variational inference objective. We also evaluate DynaMITE-RL in an offline RL setting, in
 233 which we collect an offline dataset of trajectories following an oracle goal-conditioned policy and
 234 subsequently approximate the optimal value function and RL agent using offline RL methods, e.g.,
 235 IQL [28]. The value function and the policy are parameterized with the same architecture as in the
 236 online setting and will be detailed in Appendix A.5.

237 5 Experiments

238 We present experiments that demonstrate, while VariBAD and other meta-RL methods struggle to
 239 learn good policies given nonstationary latent contexts, DynaMITE-RL exploits the causal structure

240 of a DLCMDP to more efficiently learn performant policies. We compare our approach to several
 241 state-of-the-art meta-RL baselines, showing its significantly better evaluation returns.

242 **Environments.** We test DynaMITE-RL on a suite of standard meta-RL benchmark tasks including a
 243 didactic gridworld navigation, continuous control, and human-in-the-loop robot assistance as shown
 244 in Figure 8. Gridworld navigation and MuJoCo [41] locomotion tasks are considered by Zintgraf et al.
 245 [47], Dorfman et al. [12], and Choshen and Tamar [10]. We modify these environments to incorporate
 246 temporal shifts in the reward and/or environment dynamics. To achieve good performance under
 247 these conditions, a learned policy must adapt to the latent state dynamics. More details about the
 248 environments and hyperparameters can be found in Appendix A.4 and A.5.

249 *Gridworld.* We modify the Gridworld environment used by Zintgraf et al. [47]. In a 5×5 gridworld,
 250 two possible goals are sampled uniformly at random in each episode. One of the two goals has a
 251 $+1$ reward while the other has 0 reward. The rewarding goal location changes after each session
 252 according to a predefined transition function. Goal locations are provided to the agent in the state—the
 253 only latent information is which goal has positive reward.

254 *Continuous Control.* We experiment with two tasks from OpenAI Gym [6]: Reacher and HalfCheetah.
 255 Reacher is a two-jointed robot arm tasked with reaching a 2D goal location that moves along a
 256 circular path according to some unknown transition function. HalfCheetah is a locomotion task which
 257 we modify to incorporate changing latent contexts w.r.t. the target direction (HalfCheetah-Dir), target
 258 velocity (HalfCheetah-Vel), and target velocity with opposing wind forces (HalfCheetah-Wind+Vel).

259 *Assistive Itch Scratching.* Assistive Itch Scratch is part of the Assistive-Gym benchmark [15]
 260 consisting of a human and a wheelchair-mounted 7-degree-of-freedom (DOF) Jaco robot arm. The
 261 human has limited-mobility and requires robot assistance to scratch an itch. We simulate stochastic
 262 latent context by moving the itch location—unobserved by the agent—along the human’s right arm.

263 **Meta-RL Baselines.** We compare DynaMITE-
 264 RL to several state-of-the-art (approximately)
 265 Bayes-optimal meta-RL methods including RL²
 266 [13], VariBAD [47], BOREl [12], SecBAD [8],
 267 and ContraBAR [10]. RL² [13] is an RNN-
 268 based policy gradient method which encodes
 269 environment transitions in the hidden state and
 270 maintains them across episodes. VariBAD re-
 271 duces to RL² without the decoder and the vari-
 272 ational reconstruction objective for environment
 273 transitions. BOREl primarily investigates offline
 274 meta-RL (OMRL) and proposes a few modifica-
 275 tions such as reward relabelling to address the
 276 identifiability issue in OMRL. Chen et al. [8]
 277 proposes the latent situational MDP (LS-MDP),
 278 in which there is non-stationary latent contexts
 279 that are sampled i.i.d., and SecBAD, an algo-
 280 rithm for learning in an LS-MDP. However, they
 281 do not consider latent dynamics which a crucial
 282 aspect in many applications. ContraBAR em-
 283 ploys a contrastive learning objective to discrim-
 284 inate future observations from negative samples
 285 to learn an *approximate* sufficient statistic of the
 286 history. As Zintgraf et al. [47] already demonstrate better performance by VariBAD than posterior
 287 sampling methods (e.g., PEARL [34]) we exclude such methods from our comparison.

288 **DynaMITE-RL outperforms prior meta-RL methods in a DLCMDP in both online and offline
 289 RL settings.** In Figure 4, we show the learning curves for DynaMITE-RL and baseline methods.
 290 We first observe that **DynaMITE-RL** significantly outperforms the baselines across all domains in
 291 sample efficiency and average environment returns. **RL²**, **VariBAD**, **BOReL**, **SecBAD**, and **ContraBAR**
 292 all perform poorly in the DLCMDP, converging to a suboptimal policy. By contrast, **DynaMITE-RL**
 293 accurately models the latent dynamics and consistently achieves high rewards despite the nonstation-
 294 ary latent context. We also evaluate an oracle with access to ground-truth session terminations and
 295 find that **DynaMITE-RL** with learned session terminations effectively recovers session boundaries and

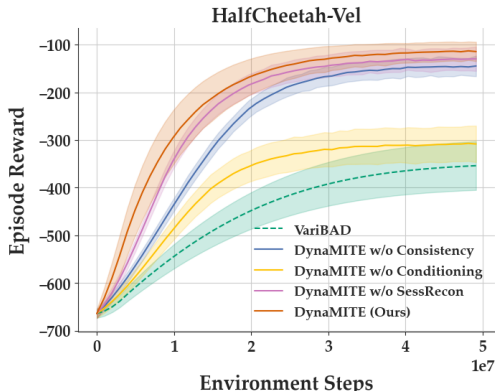


Figure 5: Ablating components of **DynaMITE-RL**. We observe that modelling latent dynamics is crucial in achieving good performance in a DLCMDP. Additionally, consistency regularization and session reconstruction improve the sample efficiency and convergence to a better performing policy.

Table 1: Average single episode returns for DynaMITE-RL and other state-of-the-art meta-RL algorithms across different environments. Results for all environments are averaged across 5 seeds beside ScratchItch which has 3 seeds. DynaMITE-RL, in bold, achieves the highest return on all of the evaluation environments and is the only method able to recover an optimal policy.

	Gridworld	Reacher	HC-Dir	HC-Vel	Wind+Vel	ScratchItch
RL ²	33.4±1.6	-150.6±1.2	-420.0±8.4	-513.2±8.7	-493.5±1.8	50.4±16.8
VariBAD	31.8±1.9	-102.4±4.2	-242.5±4.8	-363.5±3.2	-188.5±4.4	81.8±6.9
BOrel	32.4±2.4	-103.5±4.6	-240.6±4.3	-343.4±3.6	-167.8±5.4	82.5±6.0
SecBAD	38.5±3.1	-96.2±4.8	-202.4±10.4	-323.5±3.4	-155.3±5.4	101.4±9.2
ContraBAR	34.5±0.9	-101.6±3.2	-256.5±3.6	-312.3±4.8	-243.4±2.6	114.6±24.4
DynaMITE-RL	42.9±0.5	-8.4±5.1	-68.5±2.3	-146.0±8.1	-42.8±6.9	231.2±23.3

Table 2: Average single episode returns with Offline RL. Results are averaged across 5 random seeds. Algorithm with the highest average return are shown in bold. We present results for an oracle agent trained with goal information for reference.

	Gridworld	Reacher	HC-Dir	HC-Vel	HC-Dir+Vel	ScratchItch
BOrel	31.4±3.5	-102.0±5.8	-245.0±12.4	-354.0±8.3	-170.0±5.4	72.5±4.6
w/o Consistency	38.2±1.2	-33.2±2.7	-206.0±5.6	-212.0±6.4	-120.0±12.4	105.8±8.5
w/o Sess. Dynamics	33.4±1.3	-95.0±5.2	-244.0±6.0	-342.0±8.6	-166.0±9.5	74.1±2.3
DynaMITE-RL	41.8±0.6	-15.5±3.2	-154.0±8.6	-156.0±4.8	-48.0±8.6	225.5±10.6
w/ Transformer	43.8±0.6	-8.4±2.8	-132.0±7.4	-144.0±6.5	-33.0±5.8	242.5±7.4
Oracle (w/ goal)	44.6	-4.8	-112.0	-132.2	-24.4	245.3

296 matches oracle performance with sufficient training. Our empirical results validate that **DynaMITE-RL**
 297 learns a policy robust to changing latent contexts at inference time, while the baseline methods fail to
 298 adapt and get stuck in suboptimal behavior. We also demonstrate that **DynaMITE-RL** outperforms
 299 **BOrel** in an offline RL setting in Table 2 in all environments. This highlights the importance of
 300 **DynaMITE-RL** training objectives in learning a more accurate posterior belief model even without
 301 online environment interactions. We also experimented with a Transformer encoder to parameterize
 302 our belief model and find that a more powerful model further improves the evaluation performance.

303 **Each component of DynaMITE-RL contributes**

304 **to efficient learning in a DLCMDP:** We ablate the
 305 three key components of **DynaMITE-RL** to under-
 306 stand their impact on the resulting policy. We com-
 307 pare full **DynaMITE-RL** to: (i) DynaMITE-RL w/o
 308 Consistency, which does not include consistency reg-
 309 ularization; (ii) DynaMITE-RL w/o Conditioning,
 310 which does not include latent conditioning; and (iii)
 311 DynaMITE-RL w/o SessRecon, which does not in-
 312 clude session reconstruction. In Figure 5, we re-
 313 port the performance for each of these ablations and
 314 vanilla VariBAD for comparisons. First, without prior
 315 latent belief conditioning, the model converges to a
 316 suboptimal policy slightly better than **VariBAD**, con-
 317 firming the importance of modeling the latent transi-
 318 tion dynamics of a DLCMDP. Second, we find that
 319 session consistency regularization reinforces the in-
 320 ductive bias of changing dynamics and improves the
 321 sample efficiency of learning an accurate posterior
 322 model in DLCMDPs. Finally, session reconstruction
 323 masking also improves the sample efficiency by
 324 neglecting terms that are irrelevant and potentially bi-
 325 ased. Similar ablation studies in the offline RL setting
 326 can be found in Table 2, reinforcing the importance
 327 of our proposed training objectives.

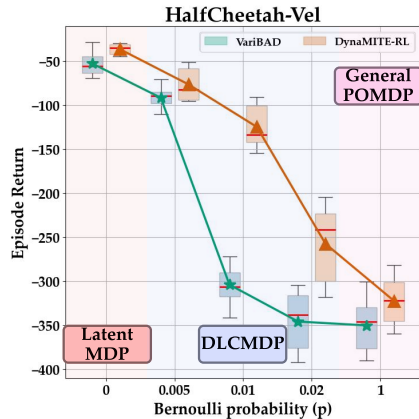


Figure 6: Ablation studies on various frequencies of latent context switches within an episode in the HalfCheetah-Vel environment. The boxplot shows the distribution over evaluation returns for 25 rollouts of trained policies with **VariBAD** and **DynaMITE-RL**. When $p = 0$, we have a latent MDP and when $p = 1$ this is equivalent to a general POMDP.

328 **DynaMITE-RL is robust to varying levels of latent stochasticity.** We study the effect of varying
 329 the number of latent context switches over an episode of fixed time horizon. For the HalfCheetah-Vel
 330 environment, we fix the episode horizon $H = 400$ to create multiple problems. We introduce a
 331 Bernoulli random variable, e.g. $d_t \sim \text{Bernoulli}(p)$ where p is a hyperparameter we set to determine
 332 the probability that the latent context changes at timestep t . If $p = 0$, the latent context remains
 333 unchanged throughout the entire episode, corresponding to a latent MDP. If $p = 1$, the latent
 334 context changes at every timestep, which is equivalent to a general POMDP. As shown in Figure 6,
 335 DynaMITE-RL performs better, on average, than VariBAD, with lower variance in a latent MDP. We
 336 hypothesize that, in the case of latent MDP, consistency regularization helps learn a more accurate
 337 posterior model by enforcing the inductive bias that the latent is static. Otherwise, there is no inherent
 338 advantage in modeling the latent dynamics if it is stationary. As we gradually increase the number
 339 of context switches, the problem becomes more difficult and closer to a general POMDP. VariBAD
 340 performance decreases drastically because it is unable to model the changing latent dynamics while
 341 DynaMITE-RL is less affected, highlighting the robustness of our approach. When we set the number
 342 of contexts equal to the episode horizon length, we recreate a fully general POMDP and again the
 343 performance between VariBAD and DynaMITE-RL converges.

344 6 Related Work

345 POMDPs provide a general framework modeling non-stationality and partial observability in sequen-
 346 tial decision problems. Many model variants have been introduced, defining a rich spectrum between
 347 episodic MDPs and POMDPs. The Bayes-adaptive MDP (BAMDP) [14] and hidden parameter MDP
 348 (HiP-MDP) [25] are both special cases of POMDPs in which environment parameters are unknown
 349 and the goal is to infer these parameters online during an episode. However, neither framework
 350 addresses the dynamics of the latent parameters across sessions, but rather assumes it is constant
 351 throughout an episode. LSMDP [8] and DP-MDP [44] do investigate nonstationary latent contexts
 352 but LSMDP samples them i.i.d., not considering the dynamics, while DP-MDP assumes fixed session
 353 lengths. By contrast, DLCMDPs models the dynamics of the latent state and simultaneously infers
 354 when the transition occurs, allowing better posterior updates at inference time.

355 DynaMITE-RL shares conceptual similarities with other meta-RL algorithms. Firstly, optimization-
 356 based techniques [16, 11, 36] learn neural network policies that can quickly adapt to new tasks at
 357 test time using policy gradient updates. However, these methods do not optimize for Bayes-optimal
 358 behavior and generally exhibit suboptimal test-time adaptation. Context-based meta-RL techniques
 359 aim to learn policies that directly infer task parameters at test time, conditioning the policy on
 360 the posterior belief. Such methods include recurrent memory-based architectures [13, 43, 30, 2]
 361 and variational approaches [20, 47, 12]. VariBAD, closest to our work, uses variational inference
 362 to approximate Bayes-optimal policies. However, we have demonstrated above the limitations of
 363 VariBAD in DLCMDPs, and have developed several crucial modifications to drive effective learning
 364 a highly performant policies in our setting.

365 7 Conclusion

366 We developed DynaMITE-RL, a meta-RL method to approximate Bayes-optimal behavior using
 367 a latent variable model. We presented the dynamic latent contextual Markov Decision Process
 368 (DLCMDP), a model in which latent context information changes according to an unknown transition
 369 function, that captures many natural settings. We derived a graphical model for this problem setting
 370 and formalized it as an instance of a POMDP. DynaMITE-RL is designed to exploit the causal
 371 structure of this model, and in a didactic GridWorld environment and several challenging continuous
 372 control tasks, we demonstrated that it outperforms existing meta-RL methods w.r.t. both learning
 373 efficiency and test-time adaptation in both online and offline-RL settings.

374 There are a number of exciting directions for future research building on the DLCMDP model. While
 375 we only consider Markovian latent dynamics in this work (i.e. future latent states are independent of
 376 prior latent states given the current latent state), we plan to investigate richer non-Markovian latent
 377 dynamics. We hope to extend DynaMITE-RL to other real-world applications including recommender
 378 systems (RS), autonomous driving, multi-agent collaborative systems, etc. DLCMDPs are a good
 379 model for RS as recommender agents often interact with users over long periods of time during which
 380 the user’s latent context changes irregularly, directly influencing their preferences.

References

- 381
- 382 [1] J. Beck, R. Vuorio, E. Z. Liu, Z. Xiong, L. M. Zintgraf, C. Finn, and S. Whiteson. A survey of
383 meta-reinforcement learning. *arXiv preprint arXiv:2301.08028*, 2023.
- 384 [2] J. Beck, R. Vuorio, Z. Xiong, and S. Whiteson. Recurrent hypernetworks are
385 surprisingly strong in meta-rl. In *Advances in Neural Information Processing*
386 *Systems*, 2023. URL [http://papers.nips.cc/paper_files/paper/2023/hash/
387 c3fa3a7d50b34732c6d08f6f66380d75-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/c3fa3a7d50b34732c6d08f6f66380d75-Abstract-Conference.html).
- 388 [3] R. Bellman. A markovian decision process. *Journal of Mathematics and Mechanics*, 6(5):
389 679—84, 1957.
- 390 [4] D. Bertsekas. *Dynamic programming and optimal control: Volume I*, volume 4. Athena
391 scientific, 2012.
- 392 [5] E. Biyik, J. Margoliash, S. R. Alimo, and D. Sadigh. Efficient and safe exploration in deter-
393 ministic markov decision processes with unknown transition models. In *American Control*
394 *Conference*, pages 1792–1799. IEEE, 2019.
- 395 [6] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba.
396 OpenAI gym. *arXiv preprint arXiv:1606.01540*, 2016.
- 397 [7] Z. Cao, E. Biyik, W. Z. Wang, A. Raventos, A. Gaidon, G. Rosman, and D. Sadigh. Reinforce-
398 ment learning based control of imitative policies for near-accident driving. *Robotics: Science*
399 *and Systems*, 2020.
- 400 [8] X. Chen, X. Zhu, Y. Zheng, P. Zhang, L. Zhao, W. Cheng, P. Cheng, Y. Xiong, T. Qin, J. Chen,
401 et al. An adaptive deep rl method for non-stationary environments with piecewise stable context.
402 *Neural Information Processing Systems*, 35:35449–35461, 2022.
- 403 [9] K. Cho, B. Van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and
404 Y. Bengio. Learning phrase representations using RNN encoder-decoder for statistical machine
405 translation. In *Conference on Empirical Methods in Natural Language Processing*, pages
406 1724–1734, 2014.
- 407 [10] E. Choshen and A. Tamar. Contrabar: Contrastive bayes-adaptive deep rl. In *International*
408 *Conference on Machine Learning*, volume 202, pages 6005–6027, 2023.
- 409 [11] I. Clavera, J. Rothfuss, J. Schulman, Y. Fujita, T. Asfour, and P. Abbeel. Model-based reinforce-
410 ment learning via meta-policy optimization. In *Conference on Robot Learning*, pages 617–629.
411 PMLR, 2018.
- 412 [12] R. Dorfman, I. Shenfeld, and A. Tamar. Offline meta reinforcement learning—identifiability
413 challenges and effective data collection strategies. *Neural Information Processing Systems*, 34:
414 4607–4618, 2021.
- 415 [13] Y. Duan, J. Schulman, X. Chen, P. L. Bartlett, I. Sutskever, and P. Abbeel. RI^2 : Fast reinforce-
416 ment learning via slow reinforcement learning. *arXiv preprint arXiv:1611.02779*, 2016.
- 417 [14] M. O. Duff. *Optimal learning: computational procedures for bayes-adaptive markov decision*
418 *processes*. PhD thesis, University of Massachusetts Amherst, 2002.
- 419 [15] Z. Erickson, V. Gangaram, A. Kapusta, C. K. Liu, and C. C. Kemp. Assistive gym: A physics
420 simulation framework for assistive robotics. In *IEEE International Conference on Robotics and*
421 *Automation*. IEEE, 2020.
- 422 [16] C. Finn, P. Abbeel, and S. Levine. Model-agnostic meta-learning for fast adaptation of deep
423 networks. In *International Conference on Machine Learning*, pages 1126–1135. PMLR, 2017.
- 424 [17] C. D. Freeman, E. Frey, A. Raichuk, S. Girgin, I. Mordatch, and O. Bachem. Brax - a differen-
425 tiable physics engine for large scale rigid body simulation. *arXiv preprint arXiv:2106.13281*,
426 2021. URL <http://github.com/google/brax>.

- 427 [18] M. Ghavamzadeh, S. Mannor, J. Pineau, and A. Tamar. Bayesian reinforcement learning: A
428 survey. *Foundations and Trends in Machine Learning*, 8(5-6):359–483, 2015.
- 429 [19] S. Huang, R. F. J. Dossa, A. Raffin, A. Kanervisto, and W. Wang. The 37 implementa-
430 tion details of proximal policy optimization. In *ICLR Blog Track*, 2022. URL <https://iclr-blog-track.github.io/2022/03/25/ppo-implementation-details/>.
431
- 432 [20] J. Humprik, A. Galashov, L. Hasenclever, P. A. Ortega, Y. W. Teh, and N. Heess. Meta
433 reinforcement learning as task inference. *arXiv preprint arXiv:1905.06424*, 2019.
- 434 [21] E. Ie, C. Hsu, M. Mladenov, V. Jain, S. Narvekar, J. Wang, R. Wu, and C. Boutilier. RecSim: A
435 configurable simulation platform for recommender systems. *arXiv preprint arXiv:1909.04847*,
436 2019.
- 437 [22] D. Jannach, A. Manzoor, W. Cai, and L. Chen. A survey on conversational recommender
438 systems. *ACM Computing Surveys (CSUR)*, 54(5):1–36, 2021.
- 439 [23] G. Jawaheer, P. Weller, and P. Kostkova. Modeling user preferences in recommender systems:
440 A classification framework for explicit and implicit user feedback. *ACM Transactions on*
441 *Interactive Intelligent Systems*, 4(2):1–26, 2014.
- 442 [24] L. P. Kaelbling, M. L. Littman, and A. R. Cassandra. Planning and acting in partially observable
443 stochastic domains. *Artificial intelligence*, 101(1-2):99–134, 1998.
- 444 [25] T. W. Killian, S. Daulton, G. Konidaris, and F. Doshi-Velez. Robust and efficient transfer
445 learning with hidden parameter markov decision processes. *Neural Information Processing*
446 *Systems*, 2017.
- 447 [26] C. Kim, J. Park, J. Shin, H. Lee, P. Abbeel, and K. Lee. Preference transformer: Modeling human
448 preferences using transformers for rl. *International Conference of Learning Representations*,
449 2023.
- 450 [27] D. P. Kingma and M. Welling. Auto-encoding variational bayes. In *International Conference*
451 *on Learning Representations*, 2014.
- 452 [28] I. Kostrikov, A. Nair, and S. Levine. Offline reinforcement learning with implicit q-learning. In
453 *International Conference on Learning Representations*, 2021.
- 454 [29] J. Kwon, Y. Efroni, C. Caramanis, and S. Mannor. RL for latent mdps: Regret guarantees and a
455 lower bound. *Neural Information Processing Systems*, 34:24523–24534, 2021.
- 456 [30] G. Lee, B. Hou, A. Mandalika, J. Lee, S. Choudhury, and S. S. Srinivasa. Bayesian policy
457 optimization for model uncertainty. *International Conference on Learning Representations*,
458 2018.
- 459 [31] S. Liu, K. C. See, K. Y. Ngiam, L. A. Celi, X. Sun, and M. Feng. Reinforcement learning for
460 clinical decision support in critical care: comprehensive review. *Journal of Medical Internet*
461 *Research*, 22(7):e18477, 2020.
- 462 [32] A. v. d. Oord, Y. Li, and O. Vinyals. Representation learning with contrastive predictive coding.
463 *arXiv preprint arXiv:1807.03748*, 2018.
- 464 [33] X. B. Peng, A. Kumar, G. Zhang, and S. Levine. Advantage-weighted regression: Simple and
465 scalable off-policy reinforcement learning. *arXiv preprint arXiv:1910.00177*, 2019.
- 466 [34] K. Rakelly, A. Zhou, C. Finn, S. Levine, and D. Quillen. Efficient off-policy meta-reinforcement
467 learning via probabilistic context variables. In *International Conference on Machine Learning*,
468 pages 5331–5340. PMLR, 2019.
- 469 [35] S. Ross, B. Chaib-draa, and J. Pineau. Bayes-adaptive pomdps. *Neural Information Processing*
470 *Systems*, 2007.
- 471 [36] J. Rothfuss, D. Lee, I. Clavera, T. Asfour, and P. Abbeel. Promp: Proximal meta-policy search.
472 *International Conference on Learning Representations*, 2018.

- 473 [37] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization
474 algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 475 [38] L. N. Steimle, D. L. Kaufman, and B. T. Denton. Multi-model markov decision processes. *IJSE*
476 *Transactions*, 53(10):1124–1139, 2021.
- 477 [39] G. Tennenholtz, A. Hallak, G. Dalal, S. Mannor, G. Chechik, and U. Shalit. On covariate shift
478 of latent confounders in imitation and reinforcement learning. *International Conference of*
479 *Learning Representations*, 2022.
- 480 [40] G. Tennenholtz, N. Merlis, L. Shani, M. Mladenov, and C. Boutilier. Reinforcement learning
481 with history dependent dynamic contexts. In A. Krause, E. Brunskill, K. Cho, B. Engelhardt,
482 S. Sabato, and J. Scarlett, editors, *Proceedings of the 40th International Conference on Machine*
483 *Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 34011–34053.
484 PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/tennenholtz23a.html>.
485
- 486 [41] E. Todorov, T. Erez, and Y. Tassa. Mujoco: A physics engine for model-based control. In *2012*
487 *IEEE/RSJ international conference on intelligent robots and systems*, pages 5026–5033. IEEE,
488 2012.
- 489 [42] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and
490 I. Polosukhin. Attention is all you need. *Neural Information Processing Systems*, 30, 2017.
- 491 [43] J. X. Wang, Z. Kurth-Nelson, D. Tirumala, H. Soyer, J. Z. Leibo, R. Munos, C. Blundell, D. Ku-
492 maran, and M. Botvinick. Learning to reinforcement learn. *arXiv preprint arXiv:1611.05763*,
493 2016.
- 494 [44] A. Xie and C. Finn. Lifelong robotic reinforcement learning by retaining experiences. In
495 *Conference on Lifelong Learning Agents, CoLLAs 2022*, volume 199 of *Proceedings of Machine*
496 *Learning Research*, pages 838–855, 2022.
- 497 [45] A. Xie, J. Harrison, and C. Finn. Deep reinforcement learning amidst continual structured
498 non-stationarity. In *Proceedings of the 38th International Conference on Machine Learning*,
499 volume 139 of *Proceedings of Machine Learning Research*, pages 11393–11403, 2021.
- 500 [46] C. Yu, J. Liu, S. Nemati, and G. Yin. Reinforcement learning in healthcare: A survey. *ACM*
501 *Computing Surveys (CSUR)*, 55(1):1–36, 2021.
- 502 [47] L. Zintgraf, K. Shiarlis, M. Igl, S. Schulze, Y. Gal, K. Hofmann, and S. Whiteson. VariBAD: A
503 very good method for bayes-adaptive deep rl via meta-learning. *International Conference of*
504 *Learning Representations*, 2020.

505 **NeurIPS Paper Checklist**

506 **1. Claims**

507 Question: Do the main claims made in the abstract and introduction accurately reflect the
508 paper's contributions and scope?

509 Answer: [Yes]

510 Justification: Section 5 demonstrates, while VariBAD and other meta-RL methods struggle
511 to learn good policies given nonstationary latent contexts, DynaMITE-RL exploits the causal
512 structure of a DLCMDP to more efficiently learn performant policies in both online and
513 offline-RL settings.

514 Guidelines:

- 515 • The answer NA means that the abstract and introduction do not include the claims
516 made in the paper.
- 517 • The abstract and/or introduction should clearly state the claims made, including the
518 contributions made in the paper and important assumptions and limitations. A No or
519 NA answer to this question will not be perceived well by the reviewers.
- 520 • The claims made should match theoretical and experimental results, and reflect how
521 much the results can be expected to generalize to other settings.
- 522 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
523 are not attained by the paper.

524 **2. Limitations**

525 Question: Does the paper discuss the limitations of the work performed by the authors?

526 Answer: [Yes]

527 Justification: We only consider Markovian latent dynamics here (i.e. future latent states are
528 independent of prior latent states given the current latent state). It would be interesting to
529 explore complex non-Markovian latent dynamics.

530 Guidelines:

- 531 • The answer NA means that the paper has no limitation while the answer No means that
532 the paper has limitations, but those are not discussed in the paper.
- 533 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 534 • The paper should point out any strong assumptions and how robust the results are to
535 violations of these assumptions (e.g., independence assumptions, noiseless settings,
536 model well-specification, asymptotic approximations only holding locally). The authors
537 should reflect on how these assumptions might be violated in practice and what the
538 implications would be.
- 539 • The authors should reflect on the scope of the claims made, e.g., if the approach was
540 only tested on a few datasets or with a few runs. In general, empirical results often
541 depend on implicit assumptions, which should be articulated.
- 542 • The authors should reflect on the factors that influence the performance of the approach.
543 For example, a facial recognition algorithm may perform poorly when image resolution
544 is low or images are taken in low lighting. Or a speech-to-text system might not be
545 used reliably to provide closed captions for online lectures because it fails to handle
546 technical jargon.
- 547 • The authors should discuss the computational efficiency of the proposed algorithms
548 and how they scale with dataset size.
- 549 • If applicable, the authors should discuss possible limitations of their approach to
550 address problems of privacy and fairness.
- 551 • While the authors might fear that complete honesty about limitations might be used by
552 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
553 limitations that aren't acknowledged in the paper. The authors should use their best
554 judgment and recognize that individual actions in favor of transparency play an impor-
555 tant role in developing norms that preserve the integrity of the community. Reviewers
556 will be specifically instructed to not penalize honesty concerning limitations.

557 **3. Theory Assumptions and Proofs**

558 Question: For each theoretical result, does the paper provide the full set of assumptions and
559 a complete (and correct) proof?

560 Answer: [NA]

561 Justification: We do not derive new theoretical results in this work.

562 Guidelines:

- 563 • The answer NA means that the paper does not include theoretical results.
- 564 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
565 referenced.
- 566 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 567 • The proofs can either appear in the main paper or the supplemental material, but if
568 they appear in the supplemental material, the authors are encouraged to provide a short
569 proof sketch to provide intuition.
- 570 • Inversely, any informal proof provided in the core of the paper should be complemented
571 by formal proofs provided in appendix or supplemental material.
- 572 • Theorems and Lemmas that the proof relies upon should be properly referenced.

573 4. Experimental Result Reproducibility

574 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
575 perimental results of the paper to the extent that it affects the main claims and/or conclusions
576 of the paper (regardless of whether the code and data are provided or not)?

577 Answer: [Yes]

578 Justification: We present all the information needed in the Appendix.

579 Guidelines:

- 580 • The answer NA means that the paper does not include experiments.
- 581 • If the paper includes experiments, a No answer to this question will not be perceived
582 well by the reviewers: Making the paper reproducible is important, regardless of
583 whether the code and data are provided or not.
- 584 • If the contribution is a dataset and/or model, the authors should describe the steps taken
585 to make their results reproducible or verifiable.
- 586 • Depending on the contribution, reproducibility can be accomplished in various ways.
587 For example, if the contribution is a novel architecture, describing the architecture fully
588 might suffice, or if the contribution is a specific model and empirical evaluation, it may
589 be necessary to either make it possible for others to replicate the model with the same
590 dataset, or provide access to the model. In general, releasing code and data is often
591 one good way to accomplish this, but reproducibility can also be provided via detailed
592 instructions for how to replicate the results, access to a hosted model (e.g., in the case
593 of a large language model), releasing of a model checkpoint, or other means that are
594 appropriate to the research performed.
- 595 • While NeurIPS does not require releasing code, the conference does require all submis-
596 sions to provide some reasonable avenue for reproducibility, which may depend on the
597 nature of the contribution. For example
 - 598 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
599 to reproduce that algorithm.
 - 600 (b) If the contribution is primarily a new model architecture, the paper should describe
601 the architecture clearly and fully.
 - 602 (c) If the contribution is a new model (e.g., a large language model), then there should
603 either be a way to access this model for reproducing the results or a way to reproduce
604 the model (e.g., with an open-source dataset or instructions for how to construct
605 the dataset).
 - 606 (d) We recognize that reproducibility may be tricky in some cases, in which case
607 authors are welcome to describe the particular way they provide for reproducibility.
608 In the case of closed-source models, it may be that access to the model is limited in
609 some way (e.g., to registered users), but it should be possible for other researchers
610 to have some path to reproducing or verifying the results.

611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [No]

Justification: Not at this point but we will release the code along with the camera ready version of the paper. We will integrate several other meta-RL environments in addition to the ones discussed in the paper.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: As said, we present all the information needed in the Appendix. We disclose hyperparameters in Appendix A.5.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: Tables 1 and 2 have error bars. Figures 5 and 6 also have error bars.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.

- 662 • The factors of variability that the error bars are capturing should be clearly stated (for
663 example, train/test split, initialization, random drawing of some parameter, or overall
664 run with given experimental conditions).
- 665 • The method for calculating the error bars should be explained (closed form formula,
666 call to a library function, bootstrap, etc.)
- 667 • The assumptions made should be given (e.g., Normally distributed errors).
- 668 • It should be clear whether the error bar is the standard deviation or the standard error
669 of the mean.
- 670 • It is OK to report 1-sigma error bars, but one should state it. The authors should
671 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
672 of Normality of errors is not verified.
- 673 • For asymmetric distributions, the authors should be careful not to show in tables or
674 figures symmetric error bars that would yield results that are out of range (e.g. negative
675 error rates).
- 676 • If error bars are reported in tables or plots, The authors should explain in the text how
677 they were calculated and reference the corresponding figures or tables in the text.

678 8. Experiments Compute Resources

679 Question: For each experiment, does the paper provide sufficient information on the com-
680 puter resources (type of compute workers, memory, time of execution) needed to reproduce
681 the experiments?

682 Answer: [Yes]

683 Justification: Section A.5.2 provides information on the computer resources.

684 Guidelines:

- 685 • The answer NA means that the paper does not include experiments.
- 686 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
687 or cloud provider, including relevant memory and storage.
- 688 • The paper should provide the amount of compute required for each of the individual
689 experimental runs as well as estimate the total compute.
- 690 • The paper should disclose whether the full research project required more compute
691 than the experiments reported in the paper (e.g., preliminary or failed experiments that
692 didn't make it into the paper).

693 9. Code Of Ethics

694 Question: Does the research conducted in the paper conform, in every respect, with the
695 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

696 Answer: [Yes]

697 Justification: We confirm that this paper conforms the NeurIPS Code of Ethics.

698 Guidelines:

- 699 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 700 • If the authors answer No, they should explain the special circumstances that require a
701 deviation from the Code of Ethics.
- 702 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
703 eration due to laws or regulations in their jurisdiction).

704 10. Broader Impacts

705 Question: Does the paper discuss both potential positive societal impacts and negative
706 societal impacts of the work performed?

707 Answer: [NA]

708 Justification: This paper is about foundational research and not tied to particular applications
709 currently. In the future, DynaMITE-RL can be used in assistive robots to improve healthcare
710 delivery and patient satisfaction as we demonstrate in the experiments with Assistive Itch
711 Scratch.

712 Guidelines:

- 713 • The answer NA means that there is no societal impact of the work performed.
- 714 • If the authors answer NA or No, they should explain why their work has no societal
- 715 impact or why the paper does not address societal impact.
- 716 • Examples of negative societal impacts include potential malicious or unintended uses
- 717 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
- 718 (e.g., deployment of technologies that could make decisions that unfairly impact specific
- 719 groups), privacy considerations, and security considerations.
- 720 • The conference expects that many papers will be foundational research and not tied
- 721 to particular applications, let alone deployments. However, if there is a direct path to
- 722 any negative applications, the authors should point it out. For example, it is legitimate
- 723 to point out that an improvement in the quality of generative models could be used to
- 724 generate deepfakes for disinformation. On the other hand, it is not needed to point out
- 725 that a generic algorithm for optimizing neural networks could enable people to train
- 726 models that generate Deepfakes faster.
- 727 • The authors should consider possible harms that could arise when the technology is
- 728 being used as intended and functioning correctly, harms that could arise when the
- 729 technology is being used as intended but gives incorrect results, and harms following
- 730 from (intentional or unintentional) misuse of the technology.
- 731 • If there are negative societal impacts, the authors could also discuss possible mitigation
- 732 strategies (e.g., gated release of models, providing defenses in addition to attacks,
- 733 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
- 734 feedback over time, improving the efficiency and accessibility of ML).

735 11. Safeguards

736 Question: Does the paper describe safeguards that have been put in place for responsible
 737 release of data or models that have a high risk for misuse (e.g., pretrained language models,
 738 image generators, or scraped datasets)?

739 Answer: [NA]

740 Justification: This paper poses no such risks.

741 Guidelines:

- 742 • The answer NA means that the paper poses no such risks.
- 743 • Released models that have a high risk for misuse or dual-use should be released with
- 744 necessary safeguards to allow for controlled use of the model, for example by requiring
- 745 that users adhere to usage guidelines or restrictions to access the model or implementing
- 746 safety filters.
- 747 • Datasets that have been scraped from the Internet could pose safety risks. The authors
- 748 should describe how they avoided releasing unsafe images.
- 749 • We recognize that providing effective safeguards is challenging, and many papers do
- 750 not require this, but we encourage authors to take this into account and make a best
- 751 faith effort.

752 12. Licenses for existing assets

753 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
 754 the paper, properly credited and are the license and terms of use explicitly mentioned and
 755 properly respected?

756 Answer: [Yes]

757 Justification: All the creators of code used in the paper are credited and VariBAD, RL²,
 758 BOREl, SecBAD, and ContraBAR are under MIT License.

759 Guidelines:

- 760 • The answer NA means that the paper does not use existing assets.
- 761 • The authors should cite the original paper that produced the code package or dataset.
- 762 • The authors should state which version of the asset is used and, if possible, include a
- 763 URL.
- 764 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.

- 765
- 766
- 767
- 768
- 769
- 770
- 771
- 772
- 773
- 774
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
 - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
 - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
 - If this information is not available online, the authors are encouraged to reach out to the asset's creators.

775 **13. New Assets**

776 Question: Are new assets introduced in the paper well documented and is the documentation
777 provided alongside the assets?

778 Answer: [NA]

779 Justification: This paper does not release new assets.

780 Guidelines:

- 781
- 782
- 783
- 784
- 785
- 786
- 787
- 788
- The answer NA means that the paper does not release new assets.
 - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
 - The paper should discuss whether and how consent was obtained from people whose asset is used.
 - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

789 **14. Crowdsourcing and Research with Human Subjects**

790 Question: For crowdsourcing experiments and research with human subjects, does the paper
791 include the full text of instructions given to participants and screenshots, if applicable, as
792 well as details about compensation (if any)?

793 Answer: [NA]

794 Justification: This paper does not involve crowdsourcing nor research with human subjects.

795 Guidelines:

- 796
- 797
- 798
- 799
- 800
- 801
- 802
- 803
- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
 - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
 - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

804 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human
805 Subjects**

806 Question: Does the paper describe potential risks incurred by study participants, whether
807 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
808 approvals (or an equivalent approval/review based on the requirements of your country or
809 institution) were obtained?

810 Answer: [NA]

811 Justification: This paper does not involve crowdsourcing nor research with human subjects.

812 Guidelines:

- 813
- 814
- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.

- 815
- 816
- 817
- 818
- 819
- 820
- 821
- 822
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
 - We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
 - For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.