# Deep Reinforcement Learning amidst Continual Structured Non-Stationarity

Annie Xie<sup>1</sup> James Harrison<sup>1</sup> Chelsea Finn<sup>1</sup>

# Abstract

As humans, our goals and our environment are persistently changing throughout our lifetime based on our experiences, actions, and internal and external drives. In contrast, typical reinforcement learning problem set-ups consider decision processes that are stationary across episodes. Can we develop reinforcement learning algorithms that can cope with the persistent change in the former, more realistic problem settings? While on-policy algorithms such as policy gradients in principle can be extended to non-stationary settings, the same cannot be said for more efficient off-policy algorithms that replay past experiences when learning. In this work, we formalize this problem setting, and draw upon ideas from the online learning and probabilistic inference literature to derive an off-policy RL algorithm that can reason about and tackle such lifelong nonstationarity. Our method leverages latent variable models to learn a representation of the environment from current and past experiences, and performs off-policy RL with this representation. We further introduce several simulation environments that exhibit lifelong non-stationarity, and empirically find that our approach substantially outperforms approaches that do not reason about environment shift.

# 1. Introduction

In the standard reinforcement learning (RL) set-up, the agent is assumed to operate in a stationary environment, i.e., under fixed dynamics and reward. However, the assumption of stationarity rarely holds in more realistic settings, such as in the context of lifelong learning systems (Thrun, 1998). That is, over the course of its lifetime, an agent may be subjected to environment dynamics and rewards that vary with time. In robotics applications, for example, this nonstationarity manifests itself in changing terrains and weather conditions. In some situations, not even the objective is necessarily fixed: consider an assistive robot helping a human whose preferences gradually change over time. And, because stationarity is a core assumption in many existing RL algorithms, they are unlikely to perform well in these environments.

Crucially, in each of the above scenarios, the environment is specified by unknown, time-varying parameters. These latent parameters are also not i.i.d., and in fact have associated but unobserved dynamics. For example, outdoor robots experience weather conditions that are determined by the season; a user-facing robot's task depends on the user's preferences which can vary based on their day-to-day routine. We formalize this problem setting with the dynamic parameter Markov decision process (DP-MDP). The DP-MDP corresponds to a sequence of stationary MDPs, related through a set of latent parameters governed by an autonomous dynamical system. While non-stationary MDPs are special instances of the partially observable Markov decision process (POMDP) (Kaelbling et al., 1998), in this setting, we can leverage structure available in the dynamics of the hidden parameters and avoid solving POMDPs in the general case.

On-policy RL algorithms can in principle cope with such non-stationarity (Sutton et al., 2007). However, in highly dynamic environments, only a limited amount of interaction is permitted before the environment changes, and on-policy methods may fail to adapt rapidly enough in this low-shot setting (Al-Shedivat et al., 2017). Instead, we desire an off-policy RL algorithm that can use past experience both to improve sample efficiency and to reason about the environment dynamics. In order to adapt, the agent needs the ability to predict how the MDP parameters will shift. We thus require a representation of the MDP as well as a model of how parameters evolve in this space, both of which can be learned from off-policy experience.

To this end, our core contribution is an off-policy RL algorithm that can operate under non-stationarity by jointly learning (1) a latent variable model, which lends a compact representation of the MDP, and (2) a maximum entropy policy with this representation. We validate our approach, which we call Lifelong Latent Actor-Critic (LILAC), on a

<sup>&</sup>lt;sup>1</sup>Stanford University. Correspondence to: Annie Xie <anniexie@stanford.edu>.

Proceedings of the 38<sup>th</sup> International Conference on Machine Learning, PMLR 139, 2021. Copyright 2021 by the author(s).

set of simulated environments that demonstrate persistent non-stationarity. In our experimental evaluation, we find that our method far outperforms RL algorithms that do not account for environment non-stationary, handles extrapolating environment shifts, and retains strong performance in *stationary* settings.

# 2. Dynamic Parameter Markov Decision Processes

The standard RL setting assumes episodic interaction with a fixed MDP (Sutton & Barto, 2018). In the real world, the assumption of episodic interaction with identical MDPs is limiting as it does not capture the wide variety of exogenous factors that may effect the decision-making problem. A common model to avoid the strict assumption of Markovian observations is the partially observed MDP (POMDP) formulation (Kaelbling et al., 1998). While the POMDP is highly general, we focus on leveraging known structure of the non-stationary MDP in this work to improve performance. In particular, we consider an episodic environment, which we call the *dynamic parameter MDP* (DP-MDP), where a new MDP (we also refer to MDPs as tasks) is presented in each episode. In reflection of the regularity of real-world non-stationarity, the tasks are sequentially related through a set of continuous parameters.

Formally, the DP-MDP is equipped with state space S, action space  $\mathcal{A}$ , and initial state distribution  $\rho_{s}(s_{1})$ . Following the formulation of the Hidden Parameter MDP (HiP-MDP) (Doshi-Velez & Konidaris, 2016), a set of unobserved task parameters  $\mathbf{z} \in \mathcal{Z}$  defines the dynamics  $p_{\mathbf{s}}(\mathbf{s}_{t+1}|\mathbf{s}_t,\mathbf{a}_t;\mathbf{z})$  and reward function  $r(\mathbf{s}_t,\mathbf{a}_t;\mathbf{z})$  for each task. In contrast to the HiP-MDP, the task parameters z in the DP-MDP are not sampled i.i.d. but instead shift stochastically according to  $p_z(z^{i+1}|z^{1:i})$ , with initial distribution  $\rho_{\mathbf{z}}(\mathbf{z}^1)$ . In other words, the DP-MDP is a sequence of tasks with parameters determined by the transition function  $p_z$ . If the task parameters z for each episode were known, the augmented state space  $S \times Z$  would define a fully observable MDP for which we can use standard RL algorithms to solve. Hence, in our approach, we aim to infer the hidden task parameters and learn their transition function, allowing us to leverage existing RL algorithms by augmenting the observations with the inferred task parameters.

Approximate model of continuously varying environments. Some environments may not exhibit shifts only at episode boundaries, and instead change more smoothly at every timestep. Formally, continuously varying environments have a set of task parameters  $z_t^i$  for each timestep tin each episodic interaction i. While these environments do not explicitly fall under the setting of DP-MDPs, the DP-MDP can exactly represent these environments when the intra-episode timestep t is either provided as part of the



Figure 1. The graphical model for the RL-as-Inference framework consists of states  $s_t$ , actions  $a_t$ , and optimality variables  $\mathcal{O}_t$ . By incorporating rewards through the optimality variables, learning an RL policy amounts to performing inference in this model.

state s or can be inferred. One way to see this mapping is to define our DP-MDP such that the task parameters for episodic interaction *i* is the concatenation of all parameters of the episode, i.e.,  $\tilde{\mathbf{z}}^i = [\mathbf{z}_t^i]_{t=1}^T$ . Then, if the continuously varying environment has dynamics  $p'_{\mathbf{s}}(\mathbf{s}_{t+1}^i | \mathbf{s}_t^i, \mathbf{a}_t^i; \mathbf{z}_t^i)$  and reward function  $r'(\mathbf{s}_t^i, \mathbf{a}_t^i; \mathbf{z}_t^i)$ , the equivalent DP-MDP, with state  $\tilde{\mathbf{s}} = [\mathbf{s}, t]$ , is defined by:

$$p_{\mathbf{s}}(\tilde{\mathbf{s}}_{t+1}^{i}|\tilde{\mathbf{s}}_{t}^{i}, \mathbf{a}_{t}^{i}; \tilde{\mathbf{z}}^{i}) = p_{\mathbf{s}}'(\tilde{\mathbf{s}}_{t+1}^{i}|\mathbf{s}_{t}^{i}, \mathbf{a}_{t}^{i}; \tilde{\mathbf{z}}^{i}[t])$$
$$r(\tilde{\mathbf{s}}_{t}^{i}, \mathbf{a}_{t}^{i}; \tilde{\mathbf{z}}^{i}) = r'(\mathbf{s}_{t}^{i}, \mathbf{a}_{t}^{i}; \tilde{\mathbf{z}}^{i}[t]).$$

Furthermore, even when the timestep is not provided, the DP-MDP can still be viewed as quantized model of these forms of environment shifts, and using this quantization can be significantly more efficient in computation than modeling small changes at every single timestep. Under this interpretation, algorithms for solving DP-MDPs are not necessarily limited to environments with inter-episode shifts, and can be applied to fairly general non-stationary environments. We validate this claim in the experiments, and indeed, find that the algorithm proposed in the next section can solve instances of continuously varying environments.

### 3. Preliminaries: RL as Inference

We first discuss an established connection between probabilistic inference and reinforcement learning (Toussaint, 2009; Levine, 2018) to provide some context for our approach. At a high level, this framework casts sequential decision-making as a probabilistic graphical model, and from this perspective, the maximum-entropy RL objective can be derived as an inference procedure in this model.

### 3.1. A Probabilistic Graphical Model for RL

As depicted in Figure 1, the proposed model consists of states  $\mathbf{s}_t$ , actions  $\mathbf{a}_t$ , and per-timestep optimality variables  $\mathcal{O}_t$ , which are related to rewards by  $p(\mathcal{O}_t = 1 | \mathbf{s}_t, \mathbf{a}_t) = \exp(r(\mathbf{s}_t, \mathbf{a}_t))$  and denote whether the action  $\mathbf{a}_t$  taken from state  $\mathbf{s}_t$  is optimal. While rewards are required to be non-positive through this relation, so long the rewards are bounded, they can be scaled and centered to be no greater than 0. A trajectory is the sequence of states and actions,  $(\mathbf{s}_1, \mathbf{a}_1, \mathbf{s}_2, \dots, \mathbf{s}_T, \mathbf{a}_T)$ , and we aim to infer the posterior



*Figure 2.* The graphical model for the DP-MDP. Each episode presents a new task, or MDP, determined by latent variables  $\mathbf{z}$ . The MDPs are related through a transition function  $p_{\mathbf{z}}(\mathbf{z}^{i+1}|\mathbf{z}^{1:i})$ .

distribution  $p(\mathbf{s}_{1:T}, \mathbf{a}_{1:T} | \mathcal{O}_{1:T} = 1)$ , i.e., the trajectory distribution that is optimal for all timesteps.

### 3.2. Variational Inference

Among existing inference tools, structured variational inference is particularly appealing for its scalability and efficiency to approximate the distribution of interest. In the variational inference framework, a variational distribution qis optimized through the variational lower bound to approximate another distribution p. Assuming a uniform prior over actions, the optimal trajectory distribution is:

$$p(\mathbf{s}_{1:T}, \mathbf{a}_{1:T} | \mathcal{O}_{1:T} = 1) \propto p(\mathbf{s}_{1:T}, \mathbf{a}_{1:T}, \mathcal{O}_{1:T} = 1)$$
$$= p(\mathbf{s}_1) \prod_{t=1}^{T} \exp(r(\mathbf{s}_t, \mathbf{a}_t)) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t).$$

For our approximating distribution, we can choose the form  $q(\mathbf{s}_{1:T}, \mathbf{a}_{1:T}) = p(\mathbf{s}_1) \prod_{t=1}^T p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) q(\mathbf{a}_t | \mathbf{s}_t)$ , where  $p(\mathbf{s}_1)$  and  $p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$  are fixed and given by the environment. We now rename  $q(\mathbf{a}_t | \mathbf{s}_t)$  to  $\pi(\mathbf{a}_t | \mathbf{s}_t)$  since this represents the desired policy. By Jensen's inequality, the variational lower bound for the evidence  $\mathcal{O}_{1:T} = 1$  is

$$\log p(\mathcal{O}_{1:T} = 1) = \log \mathbb{E}_q \left[ \frac{p(\mathbf{s}_{1:T}, \mathbf{a}_{1:T}, \mathcal{O}_{1:T} = 1)}{q(\mathbf{s}_{1:T}, \mathbf{a}_{1:T})} \right]$$
$$\geq \mathbb{E}_\pi \left[ \sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) - \log \pi(\mathbf{a}_t | \mathbf{s}_t) \right].$$

which is the maximum entropy RL objective (Ziebart et al., 2008; Toussaint, 2009; Rawlik et al., 2013; Fox et al., 2015; Haarnoja et al., 2017). This objective adds a conditional entropy term and thus maximizes both returns and the entropy of the policy. This formulation is known for its improvements in exploration, robustness, and stability over other RL algorithms, thus we build upon it in our method to inherit these qualities. We capture non-stationarity by augmenting the RL-as-inference model with latent variables  $z^i$  for each task *i*. As we will see in the next section, by viewing non-stationarity from this probabilistic perspective, our algorithm can be derived as an inference procedure in a unified model.

# 4. Off-Policy Reinforcement Learning in Non-Stationary Environments

Building upon the RL-as-inference framework, in this section, we offer a probabilistic graphical model that underlies the dynamic parameter MDP setting introduced in Section 2. Then, using tools from variational inference, we derive a variational lower bound that performs joint RL and representation learning. Finally, we present our RL algorithm, which we call Lifelong Latent Actor-Critic (LILAC), that optimizes this objective and builds upon on soft actorcritic (Haarnoja et al., 2018), an off-policy maximum entropy RL algorithm.

### 4.1. Non-stationarity as a Probabilistic Model

We can cast the dynamic parameter MDP as a probabilistic hierarchical model, where non-stationarity occurs at the episodic level, and within each episode is an instance of a stationary MDP. To do so, we construct a two-tiered model: on the first level, we have the sequence of latent variables  $z^i$  as a Markov chain, and on the second level, a Markov decision process corresponding to each  $z^i$ . The graphical model formulation of the DP-MDP is illustrated in Figure 2.

Within this formulation, the trajectories gathered from each episode are modeled individually, rather than amortized as in Subsection 3.2. Let  $\mathbf{u}^i$  represent the sequence of actions  $\mathbf{a}_{1:T}^i$  taken in trajectory *i*. Then, the probability distribution  $p(\mathbf{z}^{1:N}, \tau^{1:N} | \mathbf{u}^{1:N})$  is defined as follows:

$$p(\mathbf{z}^1)p(\tau^1|\mathbf{z}^1,\mathbf{u}^1)\prod_{i=1}^N p(\mathbf{z}^i|\mathbf{z}^{1:i-1})p(\tau^i|\mathbf{z}^i,\mathbf{u}^i)$$

where the probability of each trajectory  $\tau$  conditioned on z and action sequence u is

$$p(\tau | \mathbf{z}, \mathbf{u}) = p(\mathbf{s}_1) \prod_{t=1}^T p(\mathcal{O}_t = 1 | \mathbf{s}_t, \mathbf{a}_t; \mathbf{z}) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t; \mathbf{z})$$
$$= p(\mathbf{s}_1) \prod_{t=1}^T \exp(r(\mathbf{s}_t, \mathbf{a}_t; \mathbf{z})) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t; \mathbf{z}).$$

With this factorization, the non-stationary elements of the environment are captured by the latent variables z, and within a task, the dynamics and reward functions are necessarily stationary. This suggests that learning to infer z, which amounts to representing the non-stationarity elements of the environment with z, will reduce this RL setting to a stationary one. Taking this type of approach is appealing since there already exists a rich body of algorithms for the standard RL setting. In the next subsection, we describe how we can approximate the posterior over z, by deriving the evidence lower bound for this model under the variational inference framework.

# 4.2. Joint Representation and Reinforcement Learning via Variational Inference

Recall the agent is operating in an online learning setting. That is, it must continuously adapt to a stream of tasks and leverage experience gathered from previous tasks for learning. Thus, at any episode i > 1, the agent has observed all of the trajectories collected from episodes 1 through i - 1,  $\tau^{1:i-1} = \{\tau^1, \dots, \tau^{i-1}\}$ , where  $\tau = \{\mathbf{s}_1, \mathbf{a}_1, r_1, \dots, \mathbf{s}_T, \mathbf{a}_T, r_T\}$ .

We aim to infer, at every episode *i*, the posterior distribution over actions, given the evidence  $\mathcal{O}_{1:T}^i = 1$  and the experience from the previous episodes  $\tau^{1:i-1}$ . Following Subsection 3.2, we can leverage variational inference to optimize a variational lower bound to the logprobability of this set of evidence conditioned on the actions taken,  $\log p(\tau^{1:i-1}, \mathcal{O}_{1:T}^i = 1 | \mathbf{u}^{1:i-1})$ , where  $\mathbf{u}^i$  represents  $\mathbf{a}_{1:T}^i$ . Since  $p(\tau^{1:i-1}, \mathcal{O}_{1:T}^i = 1, \mathbf{u}^{1:i-1})$  factorizes as  $p(\tau^{1:i-1} | \mathbf{u}^{1:i-1}) p(\mathcal{O}_{1:T}^i = 1 | \tau^{1:i-1})$ , the log-probability of the evidence can be decomposed into  $\log p(\tau^{1:i-1} | \mathbf{u}^{1:i-1} | + \log p(\mathcal{O}_{1:T}^i = 1 | \tau^{1:i-1})$ . These two terms can be separately lower bounded and summed to form a single objective.

The variational lower bound of the first term follows from that of a variational auto-encoder (Kingma & Welling, 2014) with evidence  $\tau^{1:i-1}$  and latent variables  $z^{1:i-1}$ :

$$\log p(\tau^{1:i-1} | \mathbf{u}^{1:i-1}) = \log \mathbb{E}_q \left[ \frac{p(\tau^{1:i-1}, \mathbf{z}^{1:i-1} | \mathbf{u}^{1:i-1})}{q(\mathbf{z}^{1:i-1})} \right].$$

We choose our approximating distribution over the latent variables  $\mathbf{z}^i$  to be conditioned on the trajectory from episode *i*, i.e.  $q(\mathbf{z}^i | \tau^i)$ . Then, the variational lower bound is:

$$\mathcal{L}_{\text{rep}} = \mathbb{E}_q \left[ \sum_{j=1}^{i} \sum_{t=1}^{T} \log p(\mathbf{s}_{t+1}^j, r_t^j | \mathbf{s}_t^j, \mathbf{a}_t^j, \mathbf{z}^j) - D_{\text{KL}}(q(\mathbf{z}^j | \tau^j)) || p(\mathbf{z}^j | \mathbf{z}^{1:j-1})) \right].$$

The lower bound  $\mathcal{L}_{rep}$  corresponds to an objective for unsupervised representation learning in a sequential latent variable model. By optimizing the reconstruction loss of



Figure 3. An overview of our network architecture. Our method consists of the actor  $\pi$ , the critic Q, an inference network q, a decoder network, and a learned prior over latent embeddings. Each component is implemented with a neural network.

the transitions and rewards for each episode, the learned latent variables should encode the varying parameters of the MDP. Further, by imposing the prior  $p(\mathbf{z}^i | \mathbf{z}^{1:i-1})$  on the approximated distribution q through the KL divergence, the latent variables are encouraged to be sequentially consistent across time. This prior corresponds to a model of the environment's latent dynamics and gives the agent a predictive estimate of future conditions of the environment (to the extent to which the DP-MDP is predictable).

For the second term,

$$\log p(\mathcal{O}_{1:T}^{i} = 1 | \tau^{1:i-1}) = \log \int p(\mathcal{O}_{1:T}^{i} = 1, \mathbf{z}^{i} | \tau^{1:i-1}) d\mathbf{z}^{i}$$
$$= \log \int p(\mathcal{O}_{1:T}^{i} = 1 | \mathbf{z}^{i}) p(\mathbf{z}^{i} | \tau^{1:i-1}) d\mathbf{z}^{i}$$
$$\geq \mathbb{E}_{p(\mathbf{z}^{i} | \tau^{1:i-1})} \left[ \log p(\mathcal{O}_{1:T}^{i} = 1 | \mathbf{z}^{i}) \right]$$
$$\geq \mathbb{E}_{\substack{p(\mathbf{z}^{i} | \tau^{1:i-1}) \\ \pi(\mathbf{a}_{t} | \mathbf{s}_{t}, \mathbf{z}^{i})}} \left[ \sum_{i=1}^{T} r(\mathbf{s}_{t}, \mathbf{a}_{t}; \mathbf{z}^{i}) - \log \pi(\mathbf{a}_{t} | \mathbf{s}_{t}, \mathbf{z}^{i}) \right]$$
$$= \mathcal{L}_{\text{RL}}.$$

The final inequality is given by steps from Subsection 3.2. The bound  $\mathcal{L}_{RL}$  optimizes for both policy returns and policy entropy, as in the maximum entropy RL objective, but here the policy is also conditioned on the inferred latent embeddings of the MDP. This objective essentially performs task-conditioned reinforcement learning where the task variables at episode *i* are given by  $p(\mathbf{z}^i|\tau^{1:i-1})$ . Learning a multi-task RL policy is appealing, especially over a policy that adapts between episodes. That is, if the shifts in the environment are similar to those seen previously, we do not expect its performance to degrade even if the environment is shifting quickly, whereas a single-task policy would likely struggle to adapt quickly enough.

Our proposed objective is the sum of the above two terms

 $\mathcal{L} = \mathcal{L}_{rep} + \mathcal{L}_{RL}$ , which is also a variational lower bound for our entire model. Hence, while our objective was derived from and can be understood as an inference procedure in our probabilistic model, it also decomposes into two very intuitive objectives, with the first corresponding to unsupervised representation learning and the second corresponding to reinforcement learning.

### 4.3. Implementation

We introduce an inference network that outputs a distribution over latent variables,  $q(\mathbf{z}^i | \tau^i)$ , conditioned on the trajectory from the *i*-th episode. The network outputs parameters of a Gaussian distribution, and we use the reparameterization trick (Kingma & Welling, 2014) to sample  $\mathbf{z}^i$ . The weights of the inference network are trained with gradients from both  $\mathcal{L}_{rep}$  and  $\mathcal{L}_{RL}$ , which we detail below.

**Optimizing**  $\mathcal{L}_{rep}$ . Like in the standard VAE objective, the lower bound is  $\mathcal{L}_{rep} = -(\mathcal{J}_{dec} + \mathcal{J}_{KL})$  where

$$\mathcal{J}_{dec} = -\mathbb{E}_q \left[ \sum_{j=1}^{i} \sum_{t=1}^{T} \log p(\mathbf{s}_{t+1}^j, r_t^j | \mathbf{s}_t^j, \mathbf{a}_t^j, \mathbf{z}^j) \right]$$
$$\mathcal{J}_{KL} = \mathbb{E}_q \left[ \sum_{j=1}^{i} D_{KL}(q(\mathbf{z}^j | \tau^j)) \mid\mid p(\mathbf{z}^j | \mathbf{z}^{1:j-1})) \right]$$

A decoder neural network reconstructs transitions and rewards given the latent embedding  $\mathbf{z}^i$ , current state  $\mathbf{s}_t$ , and action taken  $\mathbf{a}_t$ . Finally, we approximate  $p(\mathbf{z}^i|\mathbf{z}^{1:i-1})$  and  $p(\mathbf{z}^i|\tau^{1:i-1})$  with a shared long short-term memory (LSTM) network (Hochreiter & Schmidhuber, 1997), which receives  $\mathbf{z}^{i-1}$  from  $q(\mathbf{z}^{i-1}|\tau^{i-1})$  and hidden state  $h_{i-1}$ , and produces  $\mathbf{z}^i$  and the next hidden state  $h_i$ .

**Optimizing**  $\mathcal{L}_{RL}$ . To optimize  $\mathcal{L}_{RL}$ , we extend soft actorcritic (SAC) (Haarnoja et al., 2018), which implements maximum entropy off-policy RL. As depicted in Figure 3, the policy and critic are conditioned on the environment state and the latent variables z. During training, z is sampled from  $q(\mathbf{z}|\tau)$  outputted by the inference network. At execution time, the latent variables z the policy conditions on are given by the LSTM network, based on the inferred latent variables from the previous episode. Following SAC (Haarnoja et al., 2018), the actor loss  $\mathcal{J}_{\pi}$  and critic loss  $\mathcal{J}_{Q}$  are

$$\begin{aligned} \mathcal{J}_{\pi} &= \mathop{\mathbb{E}}_{\substack{\tau \sim \mathcal{D}, \\ \mathbf{z} \sim q(\cdot | \tau)}} \left[ D_{\mathrm{KL}} \left( \pi(\mathbf{a} | \mathbf{s}, \mathbf{z}) \middle| \middle| \frac{\exp(Q(\mathbf{s}, \mathbf{a}, \mathbf{z}))}{Z(\mathbf{s}_t)} \right) \right], \\ \mathcal{J}_{Q} &= \mathop{\mathbb{E}}_{\substack{\tau \sim \mathcal{D}, \\ \mathbf{z} \sim q(\cdot | \tau)}} \left[ (Q(\mathbf{s}, \mathbf{a}, \mathbf{z}) - (r + V(\mathbf{s}', \mathbf{z})))^2 \right], \end{aligned}$$

where V denotes the target network. Our complete algorithm, Lifelong Latent Actor-Critic (LILAC), is summarized in Algorithm 1.

#### Algorithm 1 Lifelong Latent Actor-Critic (LILAC)

**Input:** env,  $\alpha_Q$ ,  $\alpha_{\pi}$ ,  $\alpha_{enc}$ ,  $\alpha_{dec}$ ,  $\alpha_{\psi}$ Randomly initialize  $\theta_Q$ ,  $\theta_\pi$ ,  $\phi_{enc}$ ,  $\phi_{dec}$ , and  $\psi$ Initialize empty replay buffer  $\mathcal{D}$ Assign  $\mathbf{z}^1 \leftarrow \vec{0}$ for i = 1, 2, ... do Sample  $\mathbf{z}^i \sim p_{\psi}(\mathbf{z}^i | \mathbf{z}^{1:i-1})$ Collect trajectory  $\tau^i$  from env with  $\pi_{\theta}(\mathbf{a}|\mathbf{s}, \mathbf{z})$ Update replay buffer  $\mathcal{D}[i] \leftarrow \tau^i$ for j = 1, 2, ..., N do Sample a batch of episodes E from  $\mathcal{D}$ > Update actor and critic  $\begin{array}{l} \theta_Q \leftarrow \theta_Q - \alpha_Q \nabla_{\theta_Q} \mathcal{J}_Q \\ \theta_\pi \leftarrow \theta_\pi - \alpha_\pi \nabla_{\theta_\pi} \mathcal{J}_\pi \end{array}$ ▷ Update inference network  $\phi_{\text{enc}} \leftarrow \phi_{\text{enc}} - \alpha_{\text{enc}} \nabla_{\phi_{\text{enc}}} \left( \mathcal{J}_{\text{dec}} + \mathcal{J}_{\text{KL}} + \mathcal{J}_Q \right)$ ⊳ Update model  $\phi_{\text{dec}} \leftarrow \phi_{\text{dec}} - \alpha_{\text{dec}} \nabla_{\phi_{\text{dec}}} \mathcal{J}_{\text{dec}}$  $\psi \leftarrow \psi - \alpha_{\psi} \nabla_{\psi} \mathcal{J}_{\mathrm{KL}}$ end for end for

# 5. Related Work

Partial observability in RL. The POMDP is a general, flexible framework capturing non-stationarity and partial observability in sequential decision-making problems. While exact solution methods are tractable only for tiny state and actions spaces (Kaelbling et al., 1998), methods based (primarily) on approximate Bayesian inference have enabled scaling to larger problems over the course of the past two decades (Kurniawati et al., 2008; Roy et al., 2005). In recent years, representation learning, and especially deep learning paired with amortized variational inference, has enabled scaling to a larger class of problems, including continuous state and action spaces (Igl et al., 2018; Han et al., 2020; Lee et al., 2019a; Hafner et al., 2019) and image observations (Lee et al., 2019a; Kapturowski et al., 2019). However, the generality of the POMDP formulation both ignores possible performance improvements that may be realized by exploiting the structure of the DP-MDP, and does not explicitly consider between-episode non-stationarity.

A variety of intermediate problem statements between episodic MDPs and POMDPs have been proposed. The Bayes-adaptive MDP formulation (BAMDP) (Duff, 2002; Ross et al., 2008), as well as the hidden parameter MDP (HiP-MDP) (Doshi-Velez & Konidaris, 2016) consider an MDP with unknown parameters governing the reward and dynamics, which we aim to infer online over the course of one episode. In this formulation, the explorationexploitation dilemma is resolved by augmenting the state space with a representation of posterior belief over the latent parameters. As noted by Duff (2002) in the RL literature and Feldbaum (1960); Bar-Shalom & Tse (1974) in control theory, this representation rapidly becomes intractable due to exploding state dimensionality. Recent work has developed effective methods for policy optimization in BAMDPs via, primarily, amortized inference (Zintgraf et al., 2020; Rakelly et al., 2019; Lee et al., 2019b). However, the BAMDP framework does not address the dynamics of the latent parameter between episodes, assuming a temporallyfixed structure. In contrast, we are capable of modeling the evolution of the latent variable over the course of episodes, leading to better priors for online inference.

A strongly related setting is the hidden-mode MDP (Choi et al., 2000), which augments the MDP with a latent parameter that evolves via a hidden Markov model with a discrete number of states. Algorithms that study the HM-MDP setting aim to quickly detect changes in the environment (Da Silva et al., 2006; Hadoux et al., 2014; Banerjee et al., 2017; Padakandla et al., 2020), while LILAC aims to anticipate future changes and adapt as they happen. In both the HM-MDP and the DP-MDP, the latent variable evolves infrequently, as opposed to at every time step as in the POMDP. The HM-MDP is limited to a fixed number of latent variable states due to the use of standard HMM inference algorithms. In contrast, our approach allows continuous latent variables, thus widely extending the range of applicability.

Non-stationarity in learning. LILAC also shares conceptual similarities with methods from online learning and lifelong learning (Shalev-Shwartz, 2012; Gama et al., 2014), which aim to capture non-stationarity in supervised learning, as well as meta-learning and meta-reinforcement learning algorithms, which aim to rapidly adapt to new settings. Within meta-reinforcement learning, two dominant techniques exist: optimization-based (Finn et al., 2017; Rothfuss et al., 2019; Zintgraf et al., 2019; Stadie et al., 2018) and contextbased, which includes both recurrent architectures (Duan et al., 2016; Wang et al., 2016; Mishra et al., 2018) and architectures based on latent variable inference (Rakelly et al., 2019; Lee et al., 2019a; Zintgraf et al., 2020). LILAC fits into this last category within this taxonomy, but extends previous methods by considering inter-episode latent variable dynamics. Previous embedding-based meta-RL algorithmswhile able to perform online inference of latent variables and incorporate this posterior belief into action selectiondo not consider how these latent variables evolve over the lifetime of the agent, as in the DP-MDP setting. The inner latent variable inference component of LILAC possesses strong similarities to the continual and lifelong learning setting (Gama et al., 2014). Many continual and lifelong learning aim to learn a variety of tasks without forgetting previous tasks (Kirkpatrick et al., 2017; Zenke et al., 2017; Lopez-Paz et al., 2017; Aljundi et al., 2019; Parisi et al., 2019; Rusu et al., 2016; Shmelkov et al., 2017; Rebuffi et al.,

2017; Shin et al., 2017). We consider a setting where it is practical to store past experiences in a replay buffer (Rolnick et al., 2019; Finn et al., 2019). Unlike these prior works, LILAC aims to learn the dynamics associated with latent factors, and perform online inference.

Within RL, non-stationarity Chandak et al. (2020) study a setting similar to ours, where the reward and transition dynamics change smoothly across episodes, and propose to use curve-fitting to estimate performance on future MDPs and learn a single policy that optimizes for future performance. This need for continual policy adaptation can result in performance lag in quickly changing environments; in contrast, LILAC learns a latent variable-conditioned policy, where different MDPs map to different values for these latent variables, and thus should be less sensitive to the rate of non-stationarity.

## 6. Experiments

In our experiments, we aim to address our central hypothesis: that existing off-policy RL algorithms struggle under persistent non-stationarity and that, by leveraging our latent variable model, LILAC can make learning in such settings both effective and efficient. To do so, we evaluate the agent's learning performance in various non-stationary environments, including environments with varying rates of change, intra-episodic shifts, and task parameters that exhibit extrapolating shifts.

Environments. We construct four continuous control environments with varying sources of change in the reward and/or dynamics. These environments are designed such that the policy needs to change in order to achieve good performance. The first is derived from the simulated Sawyer reaching task in the Meta-World benchmark (Yu et al., 2019). in which the target position is not observed and moves between episodes. In the second environment based on Half-Cheetah from OpenAI Gym (Brockman et al., 2016), we consider changes in the direction and magnitude of wind forces on the agent, and changes in the target velocity. We next consider the 8-DoF minitaur environment (Tan et al., 2018) and vary the mass of the agent between episodes, representative of a varying payload. Finally, we construct a 2D navigation task in an infinite, non-episodic environment with non-stationary dynamics which we call 2D Open World. The agent's goal is to collect food pellets and to avoid other objects and obstacles, whilst subjected to unknown perturbations that vary on an episodic schedule. These environments are illustrated in Figure 4. For full environment details, see Appendix A.

**Comparisons.** We compare our approach to standard softactor critic (SAC) (Haarnoja et al., 2018), which corresponds to our method without any latent variables, allowing



*Figure 4.* The environments in our evaluation. Each environment changes over the course of learning, including a changing target reaching position (left), variable wind and goal velocities (middle left), and variable payloads (middle right). We also introduce a 2D open world environment with non-stationary dynamics and visualize a partial snapshot of the LILAC agent's lifetime in purple (right).



*Figure 5.* Learning curves across our experimental domains. In all settings, our approach is substantially more stable and successful than SAC, SLAC, and PPO. As demonstrated in Half-Cheetah with varying target velocities and wind forces, our method can cope with non-stationarity in *both* dynamics and rewards. Error bars reflect 95% confidence intervals.

us to evaluate the performance of off-policy algorithms amid non-stationarity. We also compare to stochastic latent actorcritic (SLAC) (Lee et al., 2019a), which learns to model partially observed environments with a latent variable model but does not address inter-episode non-stationarity. This comparison allows us to evaluate the importance of modeling non-stationarity between episodes. Finally, we include proximal policy optimization (PPO) (Schulman et al., 2017) as a comparison to on-policy RL. Since the tasks in the Sawyer and Half-Cheetah domains involve goal reaching, we can obtain an oracle by training a goal-conditioned SAC policy, i.e. with the true goal concatenated to the observation. We provide this comparison to help contextualize the performance of our method against other algorithms. We tune the hyperparameters for all approaches, and run each with the best hyperparameter setting with 3 random seeds<sup>1</sup>. For all hyperparameter details, see Appendix B.

**Results.** Our experimental results are shown in Figure 5. Since on-policy algorithms tend to have worse sample complexity, we run PPO for 10 million environment steps and plot only the asymptotic returns. In all domains, LILAC attains higher and more stable returns compared to SAC, SLAC, and PPO. Since SAC amortizes experience collected across episodes into a single replay buffer, we observe that the algorithm converges to an averaged behavior. Meanwhile, SLAC does not have the mechanism to model non-stationarity across episodes, and has to infer the unknown

dynamics and reward from the initial steps taken during each episode, which the algorithm is not very successful at. Due to the cyclical nature of the tasks, the learned behavior of SLAC results in oscillating returns across tasks. Similarly, PPO cannot adapt to per-episode changes in the environment and ultimately converges to learning an average policy. In contrast to these methods, LILAC infers how the environment changes in future episodes and steadily maintains high rewards over the training procedure, despite experiencing persistent shifts in the environment in each episode. Further, LILAC can learn under simultaneous shifts in *both* dynamics and rewards, verified by the HC WindVel results. LILAC can also adeptly handle shifts in the 2D Open World environment without episodic resets. A partial snapshot of the agent's lifetime from this task is visualized in Figure 4.

Varying rates of environment shift. We next evaluate LILAC under varying rates of non-stationarity. To do so, we use the Sawyer reaching domain, where the goal moves along a fixed-radius circle, and vary the step size along the circle (0.2, 0.4, 0.6, and 0.8 radians/step) to generate environments that shift at different speeds. As depicted in Figure 6a, LILAC's performance is largely independent of the environment's rate of change. We also evaluate LILAC under stationary conditions, i.e. with a fixed goal, and find it achieves the same performance as SAC, thus retaining the ability to learn as effectively as SAC in a fixed environment. These results demonstrate LILAC's efficacy under a range of rates of non-stationarity, including the stationary case.

The gap in LILAC's performance between the nonstationary and stationary cases can likely be explained by the estimation error of future environment conditions given

<sup>&</sup>lt;sup>1</sup>We ran SAC in the Minitaur task with additional seeds for a total of 5 seeds, as recommended by a significance test of our results. The analysis, presented in Appendix C, suggests that no additional seeds are necessary for each of the other algorithms or environments.



*Figure 6.* (a) LILAC and SAC evaluated in the Sawyer task with varying rates of non-stationarity (0.2, 0.4, 0.6, and 0.8 radians/step). (b) We introduce a continuously varying variant of the Sawyer task with intra-episodic shifts, and evaluate LILAC with and without the timestep *t* included in the state s, finding that both are robust to shifts at every timestep. (c) The task parameters in this setting exhibit extrapolating shift: the target moves along a never-ending line between episodic trials. The LILAC agent can continually reach new goals, while the performance of SAC degrades over time.

by the prior  $p_{\phi}(\mathbf{z}^i|\mathbf{z}^{1:i-1})$ . Currently, the executed policy uses a fixed  $\mathbf{z}$  given by the prior for the entire duration of the episode, but a natural extension that may improve performance is updating  $\mathbf{z}$  *during* each episode. In particular, we could encode the collected partial trajectory with the inference network and combine the inferred values with the prior to form an updated estimate, akin to Bayesian filtering.

Intra-episodic environment shifts. As described in Section 2, the DP-MDP can exactly represent environments that change at every timestep, when the timestep t is provided as part of the state s. Even when the timestep t is not given, however, the DP-MDP can still be viewed as a quantization of these environments. To empirically investigate this proposition, we evaluate LILAC in a modified Sawyer reaching task: the target now moves after every time-step instead of every episode, thereby introducing intra-episode shifts. Here, the target moves at the same rate per episode as the original setting, but moves smoothly over time. We evaluate LILAC with and without the timestep t given in the state, and as shown in Figure 6b, LILAC is robust to shifts in both scenarios and significantly outperforms SAC. Hence, our approach can handle a wide subset of non-stationary environments, including those that change at every time-step.

**Extrapolating environment shifts.** To understand whether LILAC can cope within other open-world environments, we study a setting in which the dynamic parameters of the environment exhibit *extrapolating* shift. Deep RL algorithms generally struggle to generalize to out-of-distribution environment conditions (Kumar et al., 2020; Mendonca et al., 2020; Agarwal et al., 2021). In this experiment, we study an instance of extrapolating variations in the task. Specifically, we construct a Sawyer reaching task in which the goal gradually moves along a never-ending line between trials. Our results, presented in Figure 6c, indicate that LILAC can indeed learn to model as well as reach extrapolating goal positions, especially when compared to the SAC agent whose performance degrades over time.

### 7. Conclusion

We considered the problem of reinforcement learning with persistent but structured non-stationarity, a problem which we believe is a step towards reinforcement learning systems operating in the real world. This problem is at the intersection of reinforcement learning under partial observability (i.e. POMDPs) and online learning; hence we formalized the problem as a special case of a POMDP that is also significantly more tractable. We derive a graphical model underlying this problem setting, and utilize it to derive our approach under the formalism of reinforcement learning as probabilistic inference (Levine, 2018). Our method leverages this latent variable model to model the change in the environment, and conditions the policy and critic on the inferred values of these latent variables. On several challenging continuous control tasks with significant non-stationarity, we observe that our approach leads to substantial improvement compared to state-of-the-art RL methods.

While the DP-MDP formulation represents a strict generalization of the commonly-considered meta-reinforcement learning settings (typically, a BAMDP (Zintgraf et al., 2020)), it is still somewhat limited in its generality. In particular, the assumption of task parameters shifting between episodes presents a possibly unrealistic limitation, but can be relaxed when the timestep is given, or can be inferred. Otherwise, the DP-MDP can still be viewed as an approximation of environments with intra-episodic shifts. For highly *infrequent* shifts however, we may need to leverage alternative tools; in particular, this notion of infrequent, discrete shifts underlies the changepoint detection literature (Adams & MacKay, 2007; Fearnhead & Liu, 2007). Previous work within sequential decision making in changing environments (Da Silva et al., 2006; Hadoux et al., 2014; Banerjee et al., 2017) and meta-learning within changing data streams (Harrison et al., 2019) may enable a version of LILAC capable of handling unobserved changepoints.

# Acknowledgements

The authors would also like to thank Allan Zhou, Evan Liu, and Laura Smith for helpful feedback on an early version of this paper. This work was supported by an NSF GRFP, ONR grant N00014-20-1-2675, and JPMorgan Chase & Co. Any views or opinions expressed herein are solely those of the authors listed, and may differ from the views and opinions expressed by JPMorgan Chase & Co. or its affiliates. This material is not a product of the Research Department of J.P. Morgan Securities LLC. This material should not be construed as an individual recommendation for any particular client and is not intended as a recommendation of particular securities, financial instruments or strategies for a particular client. This material does not constitute a solicitation or offer in any jurisdiction.

### References

- Adams, R. P. and MacKay, D. J. Bayesian online changepoint detection. arXiv:0710.3742, 2007.
- Agarwal, R., Machado, M. C., Castro, P. S., and Bellemare, M. G. Contrastive behavioral similarity embeddings for generalization in reinforcement learning. *International Conference on Learning Representations (ICLR)*, 2021.
- Al-Shedivat, M., Bansal, T., Burda, Y., Sutskever, I., Mordatch, I., and Abbeel, P. Continuous adaptation via metalearning in nonstationary and competitive environments. *International Conference on Learning Representations* (*ICLR*), 2017.
- Aljundi, R., Kelchtermans, K., and Tuytelaars, T. Task-free continual learning. *IEEE Conference on Computer Vision* and Pattern Recognition (CVPR), 2019.
- Banerjee, T., Liu, M., and How, J. P. Quickest change detection approach to optimal control in markov decision processes with model changes. *American Control Conference (ACC)*, 2017.
- Bar-Shalom, Y. and Tse, E. Dual effect, certainty equivalence, and separation in stochastic control. *IEEE Transactions on Automatic Control*, 19(5):494–500, 1974.
- Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., and Zaremba, W. Openai gym. arXiv preprint arXiv:1606.01540, 2016.
- Chandak, Y., Theocharous, G., Shankar, S., Mahadevan, S., White, M., and Thomas, P. S. Optimizing for the future in non-stationary mdps. *ICML*, 2020.
- Choi, S. P., Yeung, D.-Y., and Zhang, N. L. Hidden-mode markov decision processes for nonstationary sequential decision making. In *Sequence Learning*, pp. 264–287. Springer, 2000.

- Colas, C., Sigaud, O., and Oudeyer, P.-Y. How many random seeds? statistical power analysis in deep reinforcement learning experiments. *arXiv preprint arXiv:1806.08295*, 2018.
- Da Silva, B. C., Basso, E. W., Bazzan, A. L., and Engel, P. M. Dealing with non-stationary environments using context detection. *ICML*, 2006.
- Doshi-Velez, F. and Konidaris, G. Hidden parameter markov decision processes: A semiparametric regression approach for discovering latent task parametrizations. *International Joint Conference on Artificial Intelligence* (*IJCAI*), 2016.
- Duan, Y., Schulman, J., Chen, X., Bartlett, P. L., Sutskever, I., and Abbeel, P. RL2: Fast reinforcement learning via slow reinforcement learning. arXiv preprint arXiv:1611.02779, 2016.
- Duff, M. O. *Optimal Learning: Computational procedures* for Bayes-adaptive Markov decision processes. PhD thesis, University of Massachusetts at Amherst, 2002.
- Fearnhead, P. and Liu, Z. On-line inference for multiple changepoint problems. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 2007.
- Feldbaum, A. Dual control theory. Avtomatika i Telemekhanika, 1960.
- Finn, C., Abbeel, P., and Levine, S. Model-agnostic metalearning for fast adaptation of deep networks. *ICML*, 2017.
- Finn, C., Rajeswaran, A., Kakade, S., and Levine, S. Online meta-learning. *ICML*, 2019.
- Fox, R., Pakman, A., and Tishby, N. Taming the noise in reinforcement learning via soft updates. *arXiv preprint arXiv:1512.08562*, 2015.
- Gama, J., Žliobaitė, I., Bifet, A., Pechenizkiy, M., and Bouchachia, A. A survey on concept drift adaptation. *ACM computing surveys (CSUR)*, 2014.
- Haarnoja, T., Tang, H., Abbeel, P., and Levine, S. Reinforcement learning with deep energy-based policies. *ICML*, 2017.
- Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. Soft actorcritic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. *ICML*, 2018.
- Hadoux, E., Beynier, A., and Weng, P. Sequential decisionmaking under non-stationary environments via sequential change-point detection. 2014.

- Hafner, D., Lillicrap, T., Fischer, I., Villegas, R., Ha, D., Lee, H., and Davidson, J. Learning latent dynamics for planning from pixels. *ICML*, 2019.
- Han, D., Doya, K., and Tani, J. Variational recurrent models for solving partially observable control tasks. *International Conference on Learning Representations (ICLR)*, 2020.
- Harrison, J., Sharma, A., Finn, C., and Pavone, M. Continuous meta-learning without tasks. arXiv preprint arXiv:1912.08866, 2019.
- Henderson, P., Islam, R., Bachman, P., Pineau, J., Precup, D., and Meger, D. Deep reinforcement learning that matters. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- Hochreiter, S. and Schmidhuber, J. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- Igl, M., Zintgraf, L., Le, T. A., Wood, F., and Whiteson, S. Deep variational reinforcement learning for pomdps. *ICML*, 2018.
- Kaelbling, L. P., Littman, M. L., and Cassandra, A. R. Planning and acting in partially observable stochastic domains. *Artificial intelligence*, 1998.
- Kapturowski, S., Ostrovski, G., Dabney, W., Quan, J., and Munos, R. Recurrent experience replay in distributed reinforcement learning. *International Conference on Learning Representations (ICLR)*, 2019.
- Kingma, D. P. and Welling, M. Auto-encoding variational Bayes. International Conference on Learning Representations (ICLR), 2014.
- Kirkpatrick, J., Pascanu, R., Rabinowitz, N., Veness, J., Desjardins, G., Rusu, A. A., Milan, K., Quan, J., Ramalho, T., Grabska-Barwinska, A., et al. Overcoming catastrophic forgetting in neural networks. *Proceedings* of the National Academy of Sciences, 2017.
- Kumar, S., Kumar, A., Levine, S., and Finn, C. One solution is not all you need: Few-shot extrapolation via structured maxent rl. *NeurIPS*, 33, 2020.
- Kurniawati, H., Hsu, D., and Lee, W. S. Sarsop: Efficient point-based pomdp planning by approximating optimally reachable belief spaces. *Robotics: Science and Systems* (*RSS*), 2008.
- Lee, A. X., Nagabandi, A., Abbeel, P., and Levine, S. Stochastic latent actor-critic: Deep reinforcement learning with a latent variable model. *arXiv preprint arXiv:1907.00953*, 2019a.

- Lee, G., Hou, B., Mandalika, A., Lee, J., Choudhury, S., and Srinivasa, S. S. Bayesian policy optimization for model uncertainty. *International Conference on Learning Representations (ICLR)*, 2019b.
- Levine, S. Reinforcement learning and control as probabilistic inference: Tutorial and review. *arXiv preprint arXiv:1805.00909*, 2018.
- Lopez-Paz, D. et al. Gradient episodic memory for continual learning. *NeurIPS*, 2017.
- Mendonca, R., Geng, X., Finn, C., and Levine, S. Metareinforcement learning robust to distributional shift via model identification and experience relabeling. *arXiv preprint arXiv:2006.07178*, 2020.
- Mishra, N., Rohaninejad, M., Chen, X., and Abbeel, P. A simple neural attentive meta-learner. *International Conference on Learning Representations (ICLR)*, 2018.
- Padakandla, S., Prabuchandran, K., and Bhatnagar, S. Reinforcement learning algorithm for non-stationary environments. *Applied Intelligence*, 50(11):3590–3606, 2020.
- Parisi, G. I., Kemker, R., Part, J. L., Kanan, C., and Wermter, S. Continual lifelong learning with neural networks: A review. *Neural Networks*, 2019.
- Rakelly, K., Zhou, A., Quillen, D., Finn, C., and Levine, S. Efficient off-policy meta-reinforcement learning via probabilistic context variables. *ICML*, 2019.
- Rawlik, K., Toussaint, M., and Vijayakumar, S. On stochastic optimal control and reinforcement learning by approximate inference. *International Joint Conference on Artificial Intelligence (IJCAI)*, 2013.
- Rebuffi, S.-A., Kolesnikov, A., and Lampert, C. H. icarl: Incremental classifier and representation learning. *IEEE Conference on Computer Vision and Pattern Recognition* (*CVPR*), 2017.
- Rolnick, D., Ahuja, A., Schwarz, J., Lillicrap, T., and Wayne, G. Experience replay for continual learning. *NeurIPS*, 2019.
- Ross, S., Chaib-draa, B., and Pineau, J. Bayes-adaptive pomdps. *NeurIPS*, 2008.
- Rothfuss, J., Lee, D., Clavera, I., Asfour, T., and Abbeel,
  P. Promp: Proximal meta-policy search. *International Conference on Learning Representations (ICLR)*, 2019.
- Roy, N., Gordon, G., and Thrun, S. Finding approximate pomdp solutions through belief compression. *Journal of Artificial Intelligence Research*, 23:1–40, 2005.

- Rusu, A. A., Rabinowitz, N. C., Desjardins, G., Soyer, H., Kirkpatrick, J., Kavukcuoglu, K., Pascanu, R., and Hadsell, R. Progressive neural networks. *arXiv:1606.04671*, 2016.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization algorithms. arXiv preprint arXiv:1707.06347, 2017.
- Shalev-Shwartz, S. Online learning and online convex optimization. "Foundations and Trends in Machine Learning", 2012.
- Shin, H., Lee, J. K., Kim, J., and Kim, J. Continual learning with deep generative replay. *NeurIPS*, 2017.
- Shmelkov, K., Schmid, C., and Alahari, K. Incremental learning of object detectors without catastrophic forgetting. arXiv:1708.06977, 2017.
- Stadie, B. C., Yang, G., Houthooft, R., Chen, X., Duan, Y., Wu, Y., Abbeel, P., and Sutskever, I. Some considerations on learning to explore via meta-reinforcement learning. *arXiv*:1803.01118, 2018.
- Sutton, R. S. and Barto, A. G. *Reinforcement learning: An introduction*. MIT press, 2018.
- Sutton, R. S., Koop, A., and Silver, D. On the role of tracking in stationary environments. *ICML*, 2007.
- Tan, J., Zhang, T., Coumans, E., Iscen, A., Bai, Y., Hafner, D., Bohez, S., and Vanhoucke, V. Sim-to-real: Learning agile locomotion for quadruped robots. *Robotics: Science* and Systems (RSS), 2018.
- Thrun, S. Lifelong learning algorithms. In *Learning to learn*, pp. 181–209. Springer, 1998.
- Toussaint, M. Robot trajectory optimization using approximate inference. *ICML*, 2009.
- Wang, J. X., Kurth-Nelson, Z., Tirumala, D., Soyer, H., Leibo, J. Z., Munos, R., Blundell, C., Kumaran, D., and Botvinick, M. Learning to reinforcement learn. *arXiv* preprint arXiv:1611.05763, 2016.
- Yu, T., Quillen, D., He, Z., Julian, R., Hausman, K., Finn, C., and Levine, S. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning. *Conference on Robot Learning (CoRL)*, 2019.
- Zenke, F., Poole, B., and Ganguli, S. Continual learning through synaptic intelligence. *ICML*, 2017.
- Ziebart, B. D., Maas, A. L., Bagnell, J. A., and Dey, A. K. Maximum entropy inverse reinforcement learning. 2008.

- Zintgraf, L., Shiarlis, K., Igl, M., Schulze, S., Gal, Y., Hofmann, K., and Whiteson, S. Varibad: A very good method for bayes-adaptive deep rl via meta-learning. *International Conference on Learning Representations (ICLR)*, 2020.
- Zintgraf, L. M., Shiarlis, K., Kurin, V., Hofmann, K., and Whiteson, S. Fast context adaptation via meta-learning. *ICML*, 2019.