

---

# Factored Adaptation for Non-stationary Reinforcement Learning

---

Fan Feng<sup>1</sup>, Biwei Huang<sup>2</sup>, Kun Zhang<sup>2,3</sup>, Sara Magliacane<sup>4,5</sup>

<sup>1</sup>City University of Hong Kong <sup>2</sup>Carnegie Mellon University

<sup>3</sup>Mohamed bin Zayed University of Artificial Intelligence

<sup>4</sup>University of Amsterdam <sup>5</sup>MIT-IBM Watson AI Lab

{ffeng1017, sara.magliacane}@gmail.com, biweih@andrew.cmu.edu, kunz1@cmu.edu

## Abstract

Dealing with non-stationarity in environments (e.g., in the transition dynamics) and objectives (e.g., in the reward functions) is a challenging problem that is crucial in real-world applications of reinforcement learning (RL). While most current approaches model the changes as a single shared embedding vector, we leverage insights from the recent causality literature to model non-stationarity in terms of individual latent change factors, and causal graphs across different environments. In particular, we propose Factored Adaptation for Non-Stationary RL (FANS-RL), a factored adaption approach that learns jointly both the causal structure in terms of a factored MDP, and a factored representation of the individual time-varying change factors. We prove that under standard assumptions, we can completely recover the causal graph representing the factored transition and reward function, as well as a partial structure between the individual change factors and the state components. Through our general framework, we can consider general non-stationary scenarios with different function types and changing frequency, including changes across episodes and within episodes. Experimental results demonstrate that FANS-RL outperforms existing approaches in terms of return, compactness of the latent state representation, and robustness to varying degrees of non-stationarity.

## 1 Introduction

Learning a stable policy under non-stationary environments is a long-standing challenge in Reinforcement learning (RL) [1, 2, 3]. While most RL approaches assume stationarity, in many real-world applications of RL there can be changes in the dynamics or the reward function, both *across* different episodes and *within* each episode. Recently, several works adapted Meta-RL methods to learn sequences of non-stationary tasks [4, 5]. However, the continuous MAML [6] adaptation for non-stationary RL [4] does not explicitly model temporal changing components, while TRIO [5] needs to meta-train the model on a set of non-stationary tasks. LILAC [7] and ZeUS [8] leverage latent variable models to directly model the change factors in the environment in a shared embedding space. In particular, they consider families of MDPs indexed by a single latent parameter. In this paper, we argue that disentangling the changes as separate latent parameters and modeling the process with a factored representation improves the efficiency of adapting to non-stationarity.

In particular, we leverage insights from the causality literature [9, 10] that model non-stationarity in terms of individual latent change factors and causal graphs across different environments. We propose Factored Adaptation for Non-Stationary RL (FANS-RL), a factored adaptation framework that jointly learns the causal structure of the MDP and a factored representation of the individual change factors, allowing for changes at discrete timepoints and continuously varying environments. While we provide a specific architecture (FN-VAE), the theoretical framework of FANS-RL can be implemented with different architectures and combined with various RL algorithms. We formalize our setting as a

*Factored Non-stationary MDP* (FN-MDP), which combines a Factored-MDP [11, 12, 13] with latent change factors that evolve in time following a Markov process.

We build upon the AdaRL framework [14], a recently proposed fast adaptation approach. AdaRL learns a factored representation that explicitly models changes (i.e., domain-specific components) in observation, dynamics and reward functions across a set of source domains. An optimal policy learnt on the source domains can then be adapted to a new target domain simply by identifying a low-dimensional change factor, without any additional finetuning. FANS-RL extends AdaRL from the stationary case with constant change factors to a general non-stationary framework. Specifically, FANS-RL learns the low-dimensional and time-evolving representations  $\theta_t^s$  and  $\theta_t^r$  that fully capture the non-stationarity of dynamics and rewards, allowing for continuous and discrete changing functions, both *within-episode* and *across-episode*. Our main contributions can be summarized as:

- We formalize FN-MDPs, a unified factored framework that can handle many non-stationary settings, including discrete and continuous changes, both within and across episodes. We prove that, under standard assumptions, the causal graph of the transition and reward function is identifiable, while we can recover a partial structure for the change factors.
- We introduce Factored Adaptation for Non-Stationary RL (FANS-RL), a general non-stationary RL approach that interleaves model estimation of an FN-MDP and policy optimization. We also describe FN-VAE, an example architecture for learning FN-MDPs.
- We evaluate FANS-RL on simulated benchmarks for continuous control and robotic manipulation tasks and show it outperforms the state of the art on the return, compactness of the latent space representation and robustness to varying degrees of non-stationarity.

## 2 Factored Non-stationary MDPs

To model different types of non-stationarity in a unified and factored way, we propose Factored Non-stationary Markov Decision Processes (FN-MDPs). FN-MDPs are an augmented form of a factored MDPs [11, 12, 13] with latent change factors that evolve over time following a Markov process. Since the change factors are latent, FN-MDPs are partially observed. We define them as:

**Definition 1.** A *Factored Non-stationary Markov Decision Process (FN-MDP)* is a tuple  $(\mathcal{S}, \mathcal{A}, \Theta^s, \Theta^r, \gamma, \mathcal{G}, \mathbb{P}_s, \mathcal{R}, \cdot, \mathbb{P}_{\theta^r}, \mathbb{P}_{\theta^s})$ , where  $\mathcal{S}$  is the state space,  $\mathcal{A}$  the action space,  $\Theta^s$  the space of the change factors for the dynamics,  $\Theta^r$  the space of the reward change factors and  $\gamma$  the discount factor. We assume  $\mathcal{G}$  is a Dynamic Bayesian Network over  $\{s_{1,t}, \dots, s_{d,t}, a_{1,t}, \dots, a_{m,t}, r_t, \theta_{1,t}^s, \dots, \theta_{p,t}^s, \theta_{1,t}^r, \dots, \theta_{q,t}^r\}$ , where  $d, m, p$ , and  $q$  are the dimensions of states, action, change factors on dynamics and reward, respectively. We define the factored state transition distribution  $\mathbb{P}_s$  as:

$$\mathbb{P}_s(\mathbf{s}_t | \mathbf{s}_{t-1}, \mathbf{a}_{t-1}, \boldsymbol{\theta}_t^s) = \prod_{i=1}^d \mathbb{P}_s(s_{i,t} | pa(s_{i,t}))$$

where  $pa(s_{i,t})$  denotes the causal parents of  $s_{i,t}$  in  $\mathcal{G}$ , which are a subset of the dimensions of  $\mathbf{s}_{t-1}$ ,  $\mathbf{a}_{t-1}$  and  $\boldsymbol{\theta}_t^s$ . Note that the action  $\mathbf{a}_{t-1}$  is a vector of  $m$  dimensions in our setting. We assume a given initial state distribution  $\mathbb{P}_s(\mathbf{s}_0)$ . Similarly, we define the reward function  $\mathcal{R}$  as a function of the parents of  $r_t$  in  $\mathcal{G}$ , i.e.,  $\mathcal{R}(\mathbf{s}_t, \mathbf{a}_t, \boldsymbol{\theta}_t^r) = \mathcal{R}(pa(r_t))$ , where  $pa(r_t)$  are a subset of dimensions of  $\mathbf{s}_t, \mathbf{a}_t$ , and  $\boldsymbol{\theta}_t^r$ . We define the factored latent change factors transition distributions  $\mathbb{P}_{\theta^s}$  and  $\mathbb{P}_{\theta^r}$  as:

$$\mathbb{P}_{\theta^s}(\boldsymbol{\theta}_t^s | \boldsymbol{\theta}_{t-1}^s) = \prod_{j=1}^p \mathbb{P}_{\theta^s}(\theta_{j,t}^s | pa(\theta_{j,t}^s)), \quad \mathbb{P}_{\theta^r}(\boldsymbol{\theta}_t^r | \boldsymbol{\theta}_{t-1}^r) = \prod_{k=1}^q \mathbb{P}_{\theta^r}(\theta_{k,t}^r | pa(\theta_{k,t}^r))$$

where  $pa(\theta_{j,t}^s)$  are a subset of the dimensions of  $\boldsymbol{\theta}_{t-1}^s$ , while  $pa(\theta_{k,t}^r)$  are a subset of dimensions of  $\boldsymbol{\theta}_{t-1}^r$ . We assume the initial distributions  $\mathbb{P}_{\theta^s}(\boldsymbol{\theta}_0^s)$  and  $\mathbb{P}_{\theta^r}(\boldsymbol{\theta}_0^r)$  are given.

We show an example DBN representing the graph  $\mathcal{G}$  of an FN-MDP in Fig. 1(a). Since we are interested in learning the graphical structure of the FN-MDP, as well as identifying the values of the latent change factors from data, we describe a generative process of an FN-MDP environment. In particular, we assume that the graph  $\mathcal{G}$  is time-invariant throughout the non-stationarity, and there are no unobserved confounders and instantaneous causal effects in the system. We will learn a set of binary masks  $\mathbf{c}^{\rightarrow}$  and  $\mathbf{C}^{\rightarrow}$  that are the indicators for edges in  $\mathcal{G}$ .

**Generative environment model.** We adapt the generative model in AdaRL [14] across  $k$  different domains to a time-varying setting on a single domain. We assume the generative process of the environment at timestep  $t$  in terms of the transition function for each dimension  $i = 1, \dots, d$  of  $\mathbf{s}_t$  is:

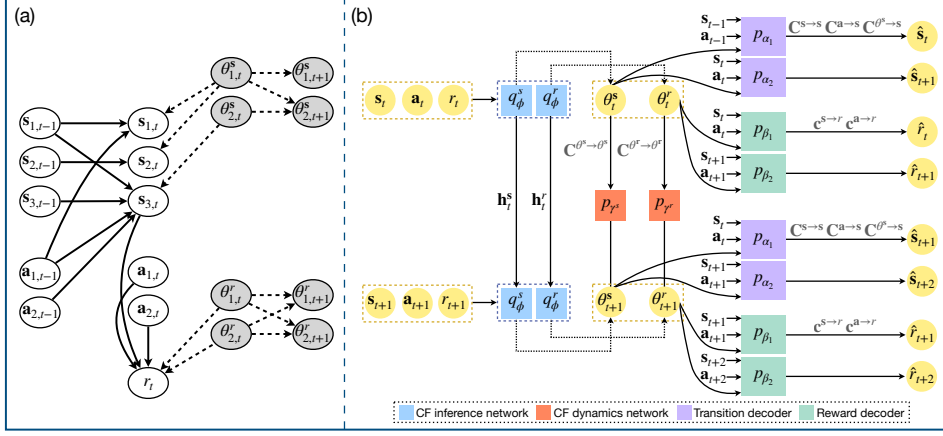


Figure 1: (a). A graphical representation of an FN-MDP. For readability, we only illustrate a subsection of dimensions of states, actions, and latent change factors. The shaded variables are unobserved; (b). The architecture of FN-VAE, which learns the generative model, explained in Sec. 3.

$$s_{i,t} = f_i(\mathbf{c}_i^{s \rightarrow s} \odot \mathbf{s}_{t-1}, \mathbf{c}_i^{a \rightarrow s} \odot \mathbf{a}_{t-1}, \mathbf{c}_i^{\theta_t^s} \odot \boldsymbol{\theta}_t^s, \epsilon_{i,t}^s) \quad (1)$$

where  $\odot$  is the element-wise product,  $f_i$  are non-linear functions and  $\epsilon_{i,t}^s$  is an i.i.d. random noise. The binary mask  $\mathbf{c}_i^{s \rightarrow s} \in \{0, 1\}^d$  represents which state components  $s_{j,t-1}$  are used in the transition function of  $s_{i,t}$ . Similarly,  $\mathbf{c}_i^{a \rightarrow s} \in \{0, 1\}^m$  indicates whether the action directly affects  $s_{i,t}$ . The change factor  $\boldsymbol{\theta}_t^s \in \mathbb{R}^p$  encodes any change in the dynamics. The binary mask  $\mathbf{c}_i^{\theta_t^s} \in \{0, 1\}^p$  represents which of the components of  $\boldsymbol{\theta}_t^s$  influence  $s_{i,t}$ . We model the reward function as:

$$r_t = h(\mathbf{c}^{s \rightarrow r} \odot \mathbf{s}_t, \mathbf{c}^{a \rightarrow r} \odot \mathbf{a}_t, \boldsymbol{\theta}_t^r, \epsilon_t^r) \quad (2)$$

where  $\mathbf{c}^{s \rightarrow r} \in \{0, 1\}^d$ ,  $\mathbf{c}^{a \rightarrow r} \in \{0, 1\}^m$ , and  $\epsilon_t^r$  is an i.i.d. random noise. The change factor  $\boldsymbol{\theta}_t^r \in \mathbb{R}^q$  encodes any change in the reward function. The binary masks  $\mathbf{c}^{\cdot \rightarrow \cdot}$  can be seen as indicators of edges in the DBN  $\mathcal{G}$ . In AdaRL, all change factors are assumed to be constant in each domain. Since in this paper we allow the change parameters to evolve in time, we introduce two additional equations:

$$\begin{aligned} \theta_{j,t}^s &= g^s(\mathbf{c}_j^{\theta_t^s \rightarrow \theta^s} \odot \boldsymbol{\theta}_{t-1}^s, \epsilon_t^{\theta_t^s}) \\ \theta_{k,t}^r &= g^r(\mathbf{c}_k^{\theta_t^r \rightarrow \theta^r} \odot \boldsymbol{\theta}_{t-1}^r, \epsilon_t^{\theta_t^r}) \end{aligned} \quad (3)$$

for  $i = 1, \dots, d$ ,  $j = 1, \dots, p$ ,  $k = 1, \dots, q$ , and  $g^s$ , and  $g^r$  are non-linear functions. We assume the binary masks  $\mathbf{c}^{\cdot \rightarrow \cdot}$  are stationary across timesteps and so are the  $\epsilon_{i,t}^s$ ,  $\epsilon_t^r$ ,  $\epsilon_t^{\theta_t^s}$  and  $\epsilon_t^{\theta_t^r}$ , the i.i.d. random noises. Although  $\mathbf{c}^{\cdot \rightarrow \cdot}$  and  $\epsilon$  are stationary, we model the changes in the functions and some changes in the graph structure through  $\boldsymbol{\theta}$ . For example a certain value of  $\boldsymbol{\theta}_t^r$  can switch off the contribution of some of the state or action dimensions in the reward function, or in other words nullify the effect of some edges in  $\mathcal{G}$ . Similarly the contribution of the noise distribution to each function can be modulated via the change factors. On the other hand, this setup does not allow adding edges that are not captured by the binary masks  $\mathbf{c}^{\cdot \rightarrow \cdot}$ . We group the binary masks in the matrices  $\mathbf{C}^{s \rightarrow s} := [\mathbf{c}_i^{s \rightarrow s}]_{i=1}^d$ ,  $\mathbf{C}^{\theta^s \rightarrow \theta^s} := [\mathbf{c}_j^{\theta^s \rightarrow \theta^s}]_{j=1}^p$ , and  $\mathbf{C}^{a \rightarrow s} := [\mathbf{c}_i^{a \rightarrow s}]_{i=1}^d$ . Similarly, we also group the binary vectors in the dynamics of the latent change factors in matrices  $\mathbf{C}^{\theta^s \rightarrow \theta^s} := [\mathbf{c}_j^{\theta^s \rightarrow \theta^s}]_{j=1}^p$  and  $\mathbf{C}^{\theta^r \rightarrow \theta^r} := [\mathbf{c}_k^{\theta^r \rightarrow \theta^r}]_{k=1}^q$ . Since latent change factors  $\boldsymbol{\theta}^s$  and  $\boldsymbol{\theta}^r$  follow a Markov process based on  $g^s$  and  $g^r$ , we can consider different types of changes by varying the form of  $g^s$  and  $g^r$ , generalizing the approaches in literature. We can also model concurrent changes in dynamics and reward, including different types of changes, e.g. a continuous  $g^s$  and a piecewise-constant  $g^r$ .

**Compact representations.** Huang et al. [14] show that the only dimensions of the state and change factors useful for policy learning are those that eventually affect the reward. These dimensions are called *compact representations* and are defined as the dimensions of the state and change factors with

a path (i.e. a sequence of edges  $\rightarrow$ ) to the present or future reward  $r_{t+\tau}$  for  $\tau \geq 0$  in the DBN  $\mathcal{G}$ :

$$s_{i,t} \in \mathbf{s}^{\min} \iff s_{i,t} \rightarrow \dots \rightarrow r_{t+\tau} \text{ for } \tau \geq 0, \text{ and } \theta_i \in \boldsymbol{\theta}^{\min} \iff \theta_i \rightarrow \dots \rightarrow r_{t+\tau} \text{ for } \tau \geq 0$$

**Continuous changes.** If  $g^s$  and  $g^r$  are continuous, then they can model smooth changes in the environment, including across episodes. While the functions in Eq. 1-3 allow us to model *within-episode* changes, i.e. changes that can happen only before  $t = H$  where  $H$  is the horizon, we also want to model *across-episode* changes. We use a separate time index  $\tilde{t}$  that models the agent’s lifetime. Initially  $\tilde{t} = t$  for  $t \leq H$ , but while we reset  $t = 0$  afterwards,  $\tilde{t}$  continues to grow indefinitely.

**Discrete changes.** We assume discrete changes happen at specific timesteps and can be represented with a piecewise constant function. In particular we denote change timesteps as  $\tilde{\mathbf{t}} = (\tilde{t}_1, \dots, \tilde{t}_M)$  where  $\tilde{t}_i$  describes a specific timepoint in the agent’s lifetime time index  $\tilde{t}$ . This allows us to model *within-episode changes*. In this case, we assume that the change happens always at the same steps  $(t_1, \dots, t_m)$  in each episode, i.e., we assume that  $\tilde{\mathbf{t}} = (t_1, \dots, t_m, H + t_1, \dots, H + t_m, 2H + t_1, \dots, 2H + t_m, \dots)$ , where  $H$  is the horizon. We can also model *across-episode changes*, when we assume the change points only occur at the end of each episode, i.e.,  $\tilde{\mathbf{t}} = (H, 2H, 3H, \dots)$ .

We extend the result on identifiability of factored MDPs in AdaRL [14] to non-stationary environments. We first assume that we observe the change factors and show we can identify the true causal graph  $\mathcal{G}$ :

**Proposition 1** (Full identifiability with observed change factors). *Suppose the generative process follows Eq. 1-3 and all change factors  $\boldsymbol{\theta}_t^s$  and  $\boldsymbol{\theta}_t^r$  are observed, i.e., Eq. 1-3 is an MDP. Under the Markov and faithfulness assumptions, i.e. conditional independences correspond exactly to  $d$ -separations, all binary masks  $\mathbf{C}^{\rightarrow\cdot}$  are identifiable, i.e., we can fully recover the causal graph  $\mathcal{G}$ .*

We provide all proofs and a detailed explanation in Appendix B. If we do not observe the change factors  $\boldsymbol{\theta}_t^s$  and  $\boldsymbol{\theta}_t^r$ , we cannot identify their dimensions, and we cannot fully recover the causal graph  $\mathcal{G}$ . On the other hand, we can still identify the partial causal graph over the state variables  $\mathbf{s}_t$ , reward variable  $r_t$ , and action variable  $\mathbf{a}_t$ . We can also identify which dimensions in  $s_{i,t}$  have changes, i.e., we can identify  $\mathbf{C}^{\theta^s \rightarrow s}$ . We formalize this idea in the following (proof in Appendix B):

**Proposition 2** (Partial Identifiability with latent change factors). *Suppose the generative process follows Eq. 1-3 and the change factors  $\boldsymbol{\theta}_t^s$  and  $\boldsymbol{\theta}_t^r$  are unobserved. Under the Markov and faithfulness assumptions, the binary masks  $\mathbf{C}^{s \rightarrow s}$ ,  $\mathbf{C}^{a \rightarrow s}$ ,  $\mathbf{c}^{s \rightarrow r}$  and  $\mathbf{c}^{a \rightarrow r}$  are identifiable. Moreover, we can identify which state dimensions are affected by  $\boldsymbol{\theta}_t^s$  and whether the reward function changes.*

This means that even in the most general case, we can learn most of the true causal graph  $G$  in an FN-MDP, with the exception of the transition structure of the latent change factors. In the following, we show a variational autoencoder setup to learn the generative process in FN-MDPs.

### 3 Learning the Generative Process in FN-MDPs

There are many possible architectures to learn FN-MDPs through Eq. 1-3. We propose *FN-VAE*, a variational autoencoder architecture described in Fig. 1(b). In FN-VAE, we jointly learn the structural relationships, state transition function, reward function, and transition function of the latent change factors, as described in detail in Appendix Alg. A1. An FN-VAE has four types of components: *change factor (CF) inference networks* that reconstruct the latent change factors, *change factor (CF) dynamics networks* that model their dynamics with an LSTM [15], *transition decoders* that reconstruct the state dynamics at the time  $t$  and predict one step further at  $t + 1$ , and *reward decoders* that reconstruct the reward at  $t$  and predict the future reward at  $t + 1$ . We now describe them in detail.

**CF inference networks (blue boxes in Fig. 1(b)).** The two inference models for latent change factors  $q_{\phi^s}(\boldsymbol{\theta}_t^s | \mathbf{s}_t, \mathbf{a}_t)$  and  $q_{\phi^r}(\boldsymbol{\theta}_t^r | \mathbf{s}_t, \mathbf{a}_t, r_t)$  are parameterised by  $\phi^s$  and  $\phi^r$ , respectively. To model the time-dependence of  $\boldsymbol{\theta}_t^s$  and  $\boldsymbol{\theta}_t^r$ , we use LSTMs [15] as inference networks. At timestep  $t$ , the dynamics change factor LSTM infers  $q_{\phi^s}(\boldsymbol{\theta}_t^s | \mathbf{s}_t, \mathbf{a}_t, r_t, \mathbf{h}_{t-1}^s)$ , where  $\mathbf{h}_{t-1}^s \in \mathbb{R}^L$  is the hidden state in the LSTM. Thus we can obtain  $\mu_{\phi^s}(\boldsymbol{\tau}_{0:t})$  and  $\sigma_{\phi^s}^2(\boldsymbol{\tau}_{0:t})$  using  $q_{\phi^s}$ , and sample the latent changing factor  $\boldsymbol{\theta}_t^s \sim \mathcal{N}(\mu_{\phi^s}(\boldsymbol{\tau}_{0:t}), \sigma_{\phi^s}^2(\boldsymbol{\tau}_{0:t}))$ , where  $\boldsymbol{\tau}_{0:t} = (\mathbf{s}_0, \mathbf{a}_0, r_0, \mathbf{s}_1, \mathbf{a}_1, r_1, \dots, \mathbf{s}_t, \mathbf{a}_t, r_t)$ . Similarly, the reward change factor LSTM infers  $q_{\phi^r}(\boldsymbol{\theta}_t^r | \mathbf{s}_t, \mathbf{a}_t, r_t, \mathbf{h}_{t-1}^r)$ , where  $\mathbf{h}_{t-1}^r \in \mathbb{R}^L$  is the hidden state, such that we can sample  $\boldsymbol{\theta}_t^r \sim \mathcal{N}(\mu_{\phi^r}(\boldsymbol{\tau}_{0:t}), \sigma_{\phi^r}^2(\boldsymbol{\tau}_{0:t}))$ .

**CF dynamics network (orange boxes in Fig. 1(b)).** We model the dynamics of latent change factors with  $p_{\gamma^s}(\boldsymbol{\theta}_{t+1}^s | \boldsymbol{\theta}_t^s, \mathbf{C}^{\theta^s \rightarrow \theta^s})$  and  $p_{\gamma^r}(\boldsymbol{\theta}_{t+1}^r | \boldsymbol{\theta}_t^r, \mathbf{C}^{\theta^r \rightarrow \theta^r})$ . To ensure the Markovianity of  $\boldsymbol{\theta}_t^s$  and

$\theta_t^r$ , we define a loss  $\mathcal{L}_{\text{KL}}$  that helps minimize the KL-divergence between  $q_\phi$  and  $p_\gamma$ .

$$\begin{aligned} \mathcal{L}_{\text{KL}} = & \sum_{t=2}^T \text{KL}(q_{\phi^s}(\theta_t^s | \mathbf{s}_t, \mathbf{a}_t, r_t, \mathbf{h}_{t-1}^s) \| p_{\gamma^s}(\theta_t^s | \theta_{t-1}^s; \mathbf{C}^{\theta^s \rightarrow \theta^s})) \\ & + \text{KL}(q_{\phi^r}(\theta_t^r | \mathbf{s}_t, \mathbf{a}_t, r_t, \mathbf{h}_{t-1}^r) \| p_{\gamma^r}(\theta_t^r | \theta_{t-1}^r; \mathbf{C}^{\theta^r \rightarrow \theta^r})) \end{aligned} \quad (4)$$

If we assume that the change between  $\theta_t^s$  and  $\theta_{t+1}^s$ , and similarly  $\theta_t^r$ , is smooth, we can add a smoothness loss  $\mathcal{L}_{\text{smooth}}$ . We provide a smooth loss for discrete changes in Appendix D.2.

$$\mathcal{L}_{\text{smooth}} = \sum_{t=2}^T (\|\theta_t^s - \theta_{t-1}^s\|_1 + \|\theta_t^r - \theta_{t-1}^r\|_1) \quad (5)$$

**Transition decoders (purple boxes in Fig. 1(b)).** We learn an approximation of the transition dynamics in Eq. 1 by learning a reconstruction, parameterized by  $\alpha_1$ , and a prediction encoder, parameterized by  $\alpha_2$ . To simplify the formulas, we define  $\mathbf{C}^{\cdot \rightarrow \cdot} := (\mathbf{C}^{s \rightarrow s}, \mathbf{C}^{a \rightarrow s}, \mathbf{C}^{\theta^s \rightarrow s})$ . At timestep  $t$ , the reconstruction encoder  $p_{\alpha_1}(\mathbf{s}_t | \mathbf{s}_{t-1}, \mathbf{a}_{t-1}, \theta_t^s; \mathbf{C}^{\cdot \rightarrow \cdot})$  reconstructs the state from current state  $\mathbf{s}_t$  with sampled  $\theta_t^s$ . The one-step prediction encoder  $p_{\alpha_2}(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t, \theta_t^s)$  instead tries to approximate the next state  $\mathbf{s}_{t+1}$ . We do not use the prediction loss, when the one-step prediction is not smooth. In particular, we do not use it for the last time-step in episode  $i$  if there is a change happening at the first step in episode  $(i+1)$ , since the states in new episodes will be randomly initiated. We also do not use it in the case of discrete changes at the timesteps  $(\tilde{t}_1 - 1, \dots, \tilde{t}_M - 1)$ . The loss functions are:

$$\begin{aligned} \mathcal{L}_{\text{rec-dyn}} &= \sum_{t=1}^{T-2} \mathbb{E}_{\theta_t^s \sim q_\phi} \log p_{\alpha_1}(\mathbf{s}_t | \mathbf{s}_{t-1}, \mathbf{a}_{t-1}, \theta_t^s; \mathbf{C}^{\cdot \rightarrow \cdot}) \\ \mathcal{L}_{\text{pred-dyn}} &= \sum_{t=1}^{T-2} \mathbb{E}_{\theta_t^s \sim q_\phi} \log p_{\alpha_2}(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t, \theta_t^s) \end{aligned} \quad (6)$$

**Reward decoders (green boxes in Fig. 1(b)).** Similarly, we use a reconstruction encoder  $p_{\beta_1}(r_t | \mathbf{s}_t, \mathbf{a}_t, \theta_t^r, \mathbf{c}^{s \rightarrow r}, \mathbf{c}^{a \rightarrow r})$ , parameterized by  $\beta_1$ , and a one-step prediction encoder  $p_{\beta_2}(r_{t+1} | \mathbf{s}_{t+1}, \mathbf{a}_{t+1}, \theta_t^r)$ , parameterized by  $\beta_2$ , to approximate the reward function. Similarly to transition decoders, we do not use the one-step prediction loss, if it is not smooth. The losses are:

$$\begin{aligned} \mathcal{L}_{\text{rec-rw}} &= \sum_{t=1}^{T-2} \mathbb{E}_{\theta_t^r \sim q_\phi} \log p_{\beta_1}(r_t | \mathbf{s}_t, \mathbf{a}_t, \theta_t^r; \mathbf{c}^{s \rightarrow r}, \mathbf{c}^{a \rightarrow r}) \\ \mathcal{L}_{\text{pred-rw}} &= \sum_{t=1}^{T-2} \mathbb{E}_{\theta_t^r \sim q_\phi} \log p_{\beta_2}(r_{t+1} | \mathbf{s}_{t+1}, \mathbf{a}_{t+1}, \theta_t^r) \end{aligned} \quad (7)$$

**Sparsity loss.** We encourage sparsity in the binary masks  $\mathbf{C}^{\cdot \rightarrow \cdot}$  to improve identifiability, by using following loss with adjustable hyperparameters  $(w_1, \dots, w_7)$ , which we learn through grid search.

$$\begin{aligned} \mathcal{L}_{\text{sparse}} = & w_1 \|\mathbf{C}^{s \rightarrow s}\|_1 + w_2 \|\mathbf{C}^{a \rightarrow s}\|_1 + w_3 \|\mathbf{C}^{\theta^s \rightarrow s}\|_1 + w_4 \|\mathbf{c}^{s \rightarrow r}\|_1 + w_5 \|\mathbf{c}^{a \rightarrow r}\|_1 \\ & + w_6 \|\mathbf{C}^{\theta^s \rightarrow \theta^s}\|_1 + w_7 \|\mathbf{C}^{\theta^r \rightarrow \theta^r}\|_1 \end{aligned} \quad (8)$$

The total loss is  $\mathcal{L}_{\text{vae}} = k_1(\mathcal{L}_{\text{rec-dyn}} + \mathcal{L}_{\text{rec-rw}}) + k_2(\mathcal{L}_{\text{pred-dyn}} + \mathcal{L}_{\text{pred-rw}}) - k_3 \mathcal{L}_{\text{KL}} - k_4 \mathcal{L}_{\text{sparse}} - k_5 \mathcal{L}_{\text{smooth}}$ , where  $(k_1, \dots, k_5)$ : hyper-parameters, which we learn with an automatic weighting method [16].

**Learning from raw pixels.** Our framework can be easily extended to image inputs by adding an encoder  $\phi^o$  to learn the latent state variables from pixels, similar to other works [17, 18, 14]. In this case, our identifiability results do not hold anymore, since we cannot guarantee that we identify the true causal variables. We describe this component in Appendix D.3.

## 4 FANS-RL: Online Model Estimation and Policy Optimization

We propose Factored Adaptation for Non-Stationary RL (FANS-RL), a general algorithm that interleaves model estimation and policy optimization, as shown in Alg. 1 and Appendix Fig. A9.

---

**Algorithm 1:** Factored Adaptation for non-stationary RL

---

```
1: Init: Env; VAE parameters:  $\phi = (\phi^s, \phi^r)$ ,  $\alpha = (\alpha_1, \alpha_2)$ ,  $\beta = (\beta_1, \beta_2)$ ,  $\gamma$ ; Binary masks:  $C^{**}$ ;
   Policy parameters:  $\psi$ ; replay buffer:  $\mathcal{D}$ ; Number of episodes:  $N$ ; Episode horizon:  $H$ ; Initial  $\theta$ :
    $\theta_{\text{old}}^s$  and  $\theta_{\text{old}}^r$ ; Length of collected trajectory:  $k$ .
2: Output: VAE parameters:  $\phi$ ,  $\alpha$ ,  $\beta$ ,  $\gamma$ ; Policy parameters:  $\psi$ 
3: Collect multiple trajectories of length  $k$ :  $\tau = \{\tau_{0:k}^1, \tau_{0:k}^2, \dots\}$  with policy  $\pi_\psi$  from Env;
4: Learn FN-VAE from  $\tau$ , including masks  $C^{**}$  that represent the graph  $\mathcal{G}$  (Appendix Alg. A1)
5: Identify the compact representations  $s^{\text{min}}$  and change factors  $\theta^{\text{min}}$  based on  $C^{**}$ 
6: for  $n = 0, \dots, N - 1$  do
7:   for  $t = 0, \dots, H - 1$  do
8:     Observe  $s_t$  from Env;
9:     if  $t = 0$  then
10:       $\theta^s \leftarrow \theta_{\text{old}}^s$  and  $\theta^r \leftarrow \theta_{\text{old}}^r$ 
11:     else
12:       $\theta^s \leftarrow \theta_{t-1}^s$  and  $\theta^r \leftarrow \theta_{t-1}^r$ 
13:     end if
14:     for  $j = s, r$  do
15:       Infer mean  $\mu_{\gamma^j}(\theta^j)$  and variance  $\sigma_{\gamma^j}^2(\theta^j)$  of the change parameter  $\theta_t^j$  via  $p_{\gamma^j}$ 
16:       Sample  $\theta_t^j \sim \mathcal{N}(\mu_{\gamma^j}(\theta^j), \sigma_{\gamma^j}^2(\theta^j))$ 
17:     end for
18:     if  $t = H - 1$  then
19:        $\theta_{\text{old}}^s \leftarrow \theta_t^s$  and  $\theta_{\text{old}}^r \leftarrow \theta_t^r$ ;
20:     end if
21:     Generate  $\mathbf{a}_t \sim \pi_\psi(\mathbf{a}_t | s_t^{\text{min}}, \theta_t^{\text{min}})$  and receive reward  $r_{n,t}$  from Env;
22:     Add  $(s_t, \mathbf{a}_t, r_t, \theta_t^s, \theta_t^r)$  to replay buffer  $\mathcal{D}$ ;
23:     Extract a trajectory with length  $k$  from replay buffer  $\mathcal{D}$ ;
24:     Learn FN-VAE (Appendix Alg. A1) with updateG=False (i.e. with fixed masks  $C^{**}$ );
25:     Sample a batch of data from replay buffer  $\mathcal{D}$  and update policy network parameters  $\psi$ ;
26:   end for
27: end for
```

---

After we estimate the initial FN-MDP with the FN-VAE, we can identify compact representations  $s^{\text{min}}$  and  $\theta^{\text{min}}$  following AdaRL [14]. In particular, the only dimensions of the state and change factors that are useful for policy learning are those that have a directed path to the reward in the graph  $\mathcal{G}$ . The online policy  $\pi_\psi(\mathbf{a}_t | s_t^{\text{min}}, q_\phi(\theta_t^{\text{min}} | \tau_{0:t}))$  can be learned end-to-end, including learning the FN-VAE, as shown in Alg. 1. We use SAC [19] as our policy learning model, so the policy parameters are  $\psi = (\pi, Q)$ .

**Continuous changes.** In Alg. 1 we describe our framework in case of continuous changes that can span across episodes. We start by collecting a few trajectories  $\tau$  and then learn our initial FN-VAE (Lines 3-4). We can use the graphical structure of the initial FN-VAE to identify the compact representations (Line 5). During the online model estimation and policy learning stage, we estimate the latent change factors  $\theta_t^s$  and  $\theta_t^r$  using the CF dynamics networks  $\gamma^s$  and  $\gamma^r$  (Lines 8-20). Since in this case, we assume the dynamics of the change factors are smooth across episodes, at time  $t = 0$  we will use the last timestep ( $H - 1$ ) of the previous episode as a prior on the change factors (Line 10). Otherwise, we will estimate the change factors using their values in the previous timestep  $t - 1$ . We use the estimated latent factors  $\theta_t$  and observed state  $s_t$  to generate  $\mathbf{a}_t$  using  $\pi_\psi$  and receive a reward  $r_t$  (Line 21). We add  $(s_t, \mathbf{a}_t, r_t, \theta_t^s, \theta_t^r)$  to the replay buffer (Line 22). We now update our estimation of the FN-VAE, but we keep the graph  $\mathcal{G}$  fixed (Lines 23-24). Finally, we sample a batch of trajectories in the replay buffer and update the policy network  $\psi$  (Line 25).

**Discrete changes.** Since we assume discrete changes happen at specific timestep  $\tilde{t} = (\tilde{t}_1, \dots, \tilde{t}_M)$ , we can easily modify Alg. 1 for discrete changes, both within-episode and across-episode, by changing Lines 9-20 to only update the change parameters at the timesteps in  $\tilde{t}$ , as shown in Appendix Alg. A2.

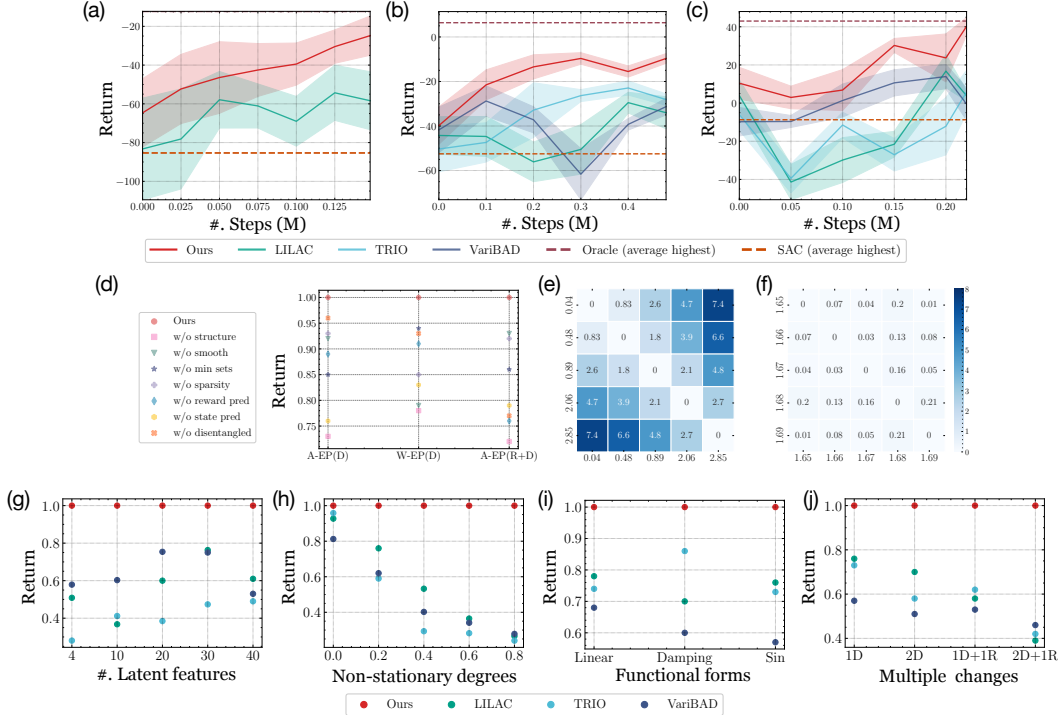


Figure 2: Summary of experimental results. (a)-(c). Average return (smoothed) across 10 runs. We only indicate the average of the highest result of all times for oracle and SAC. The shaded region is  $(\mu - \sigma, \mu + \sigma)$ , where  $\mu$  is the mean and  $\sigma$  is the standard deviation. (a) Half-Cheetah-V3 with continuous (sine) changes on  $f_w$ ; (b) Sawyer-Reaching with discrete across-episode changes on  $s^g$ ; and (c) Minitaur with discrete across-episode changes on  $m$  and  $s_{t,v}$  concurrently. (d) Ablation studies on Half-Cheetah with across & within episode changes on dynamics and across episode changes on both dynamics and rewards. (e)-(f): Pairwise distance on learned  $\theta$  between different time steps in Half-Cheetah experiment with across-episode changes on rewards. (g)-(j): Average and normalized final return on 10 runs on Half-Cheetah (g) with within-in episode changes on wind forces using a different number of dimensions in latent representation space; (h) with different non-stationary degrees on across-episode and multi-factor changes; (i) with different functional forms (across episode changes on dynamics); and (j) with different combinations of across-episode changes.

## 5 Evaluation

We evaluate our approach on four well-established benchmarks, including Half-Cheetah-V3 from MuJoCo [20, 21], Sawyer-Reaching and Sawyer-Peg from Sawyer [22, 18], and Minitaur [23]. We modified these tasks to test several non-stationary RL scenarios with continuous and discrete changes. The results suggest that FANS-RL can (1) obtain high rewards, (2) learn meaningful mappings to the change factors with compact latent vectors, and (3) be robust to different non-stationary levels, different functional forms (e.g. piecewise linear, sine, damped sine) and multiple concurrent changes. For space limits, we only highlight a subset of results, and report the full results in Appendix C.4. The implementation will be open-sourced at <https://bit.ly/3erKoWm>.

**Half-Cheetah-v3 in MuJoCo.** In this task, the agent is moving forward using the joint legs and the objective is to achieve the target velocity  $v^g$ . We consider both the changes in the dynamics (change of the wind forces  $f^w$ , change of gravity) and reward functions (change of target velocity  $v^g$ ). We also consider changes on the agent’s mechanism, where one random joint is disabled. The reward function is  $r_t = -\|v_t^o - v_t^g\|_2 - 0.05\|a_t\|_2$ , where  $v_t^o$  and  $a_t$  are the agent’s velocity and action, respectively, at timestep  $t$ . The number of time steps in each episode is 50. For dynamics, we change the wind forces  $f^w$  in the environment. Moreover, in terms of the reward functions, we change the target velocity  $v^g$  to be a time-dependent variable. The change function in the dynamics  $f^w$  can be either continuous or discrete, and discrete changes can happen both within and across episodes.

Methods	ZeUS	Meld	CaDM	Hyper-Dynamics	Ours
Best Avg. Return	12.45	6.38	4.04	3.91	<b>18.01</b>

Table 1: Average highest return on Sawyer-Peg. The number of random trails is 10.

Similarly to LILAC [7], we choose different functions (piecewise linear, sine and damped sine), besides allowing the change at specified intervals, we also allow it to change smoothly. The change in the reward function  $v^g$  is not generally stable in the continuous case, so we only consider the discrete and across episode change functions  $v^g$  for the reward. We also design a scenario where dynamic and reward functions change concurrently. We report all equations for  $g^s$  and  $g^r$  in Appendix C.1.

**Sawyer.** We consider two robotic manipulation tasks, Sawyer-Reaching and Sawyer-Peg. We describe the non-stationary settings in Appendix C.2.

In Sawyer-Reaching, the sawyer arm is trained to reach a target position  $s^g$ . The reward  $r_t$  is the difference between the current position  $s_t$  and the target position  $r_t = -\|s_t - s^g\|_2$ . In this task, we cannot directly modify the dynamics in the simulator, so consider a reward-varying scenario, where the target location changes across each episode following a periodic function.

In Sawyer-Peg, the robot arm is trained to insert a peg into a designed target location  $s^g$ . In this task, following [8], we consider a reward-varying scenario, where the target location changes across each episode following a periodic function. In order to compare with similar approaches, e.g., ZeUS [8], CADM [24], Hyperdynamics [25] and Meld [18], we evaluate our method on raw pixels. Following [8], we consider discrete across-episode changes, where the target location can change in each episode, and is randomly sampled from a small interval.

**Minitaur.** A minitaur robot is a simulated quadruped robot with eight direct-drive actuators. The minitaur is trained to move at the target speed  $s_{t,v}$ . The reward is  $r_t = 0.3 - |0.3 - s_{t,v}| - 0.01 \cdot \|a_t - 2a_{t-1} + a_{t-2}\|_1$ . We modify (1) the mass  $m$  (dynamics) and (2) target speed  $s_{t,v}$  (reward) of the minitaur. We consider continuous, discrete across-episode and discrete within-episode changes for the dynamics, and across-episode changes for the reward. We describe the settings in Appendix C.3.

**Baselines.** We compare our approach with a meta-RL approach for stationary RL, **VariBAD** [26], a meta-RL approach for non-stationary RL, **TRIO** [5], as well as with two representative task embedding approaches, **LILAC** [7] and **ZeUS** [8]. The details on the meta-learning setups are given in Appendix D.4. We also compare with stationary RL method, **SAC** [19], which will be our lower-bound, and compare with an **oracle** agent that has access to the full information of non-stationarity (e.g., the wind forces) and can use it to learn a policy, which will be our upper-bound. For all baselines, we use SAC for policy learning. We compare with the baselines on *average return* and the *compactness of the latent space* in varying degrees and functional forms of non-stationarity.

**Experimental results.** Fig. 2(a)-(c) shows the smoothed curves of average return across timesteps in a subsection of Half-Cheetah, Sawyer-Reaching, and Minitaur experiments. Smoothed curves of other experiments are given in Appendix Fig. A2 and A3. For continuous changes, we only compare with LILAC [7] since other approaches are not applicable. We smooth the learning curves by uniformly selecting a few data points for readability. Table 1 shows the results on Sawyer-Peg experiments using raw pixels as input, based on the reported results in [8]. The learning curve is given in Appendix Fig. A11. For a fair comparison, we indicate the best and final average return for each baseline in Fig. A11. Full results of all experiments are given in Appendix C.4, including significance tests for Wilcoxon signed-rank test at  $\alpha = 0.05$  showing that FANS-RL is significantly better than baselines. We also visualize the learned graph in Sawyer in Appendix Sec C.7, showing that FN-VAE recovers reasonable causal relations in this domain, where the true graph is unknown.

**Ablation studies.** We conduct ablation studies on each component of FANS-RL and report some key results in Fig. 2(d). The ablation studies verify the effectiveness of all components, including binary masks/structure, smoothness loss, sparsity loss, reward prediction or state prediction. The results show that *the largest gain is provided by the factored representation*, validating our original hypothesis, followed by state prediction. As expected, reward prediction is also important when there is a nonstationary reward, while smoothness is important for within-episode changes. The disentangled design of CF inference networks is valuable when there are changes on both dynamics



and reward functions. Full results are in Appendix C.5, showing that learning which state components are affected by change, i.e. estimating the binary masks  $C^{\theta^s \rightarrow s}$  and  $C^{\theta^r \rightarrow r}$ , provides the largest gains.

**Visualization on the learned  $\theta$ .** To verify that the learned  $\theta$  can capture the true change factors, we compute the pairwise distance between learned  $\theta^r$  at different time steps. We randomly sample 10 time steps from the Half-Cheetah experiment with across-episode changes on reward functions. Fig. 2(e) gives the pairwise distance of  $\theta^r$  among 5 time steps (from episode 148, 279, 155, 159, 230) with different target speed values (0.04, 0.48, 0.89, 2.06, 2.85), respectively. We can find that there is a positive correlation between the distance of learned  $\theta^r$  and values of change factors. Meanwhile, we also sample 5 time steps from episodes 31, 188, 234, 408 with target values around 1.67. Fig. 2(f) shows that the distance of  $\theta^r$  among these 5 time steps is very small, indicating that the learned  $\theta^r$  are almost the same for the similar values of change factors at different time steps. The visualization suggests that the learned  $\theta$  can capture meaningful mappings from the time-varying factors.

**Varying latent dimensions, non-stationary levels and functional forms.** Fig. 2(g) shows the normalized averaged return versus the number of two latent features on Half-Cheetah. Our framework can learn a better policy with relatively smaller feature dimensions in the latent space than other approaches. As we show in Appendix C.4, Saywer and Minitaur have the similar trend, where we learn a better policy than the baselines with fewer latent features. We also vary the non-stationary levels in Half-Cheetah with discrete across-episode changes on both dynamics and rewards. A higher non-stationary degree indicates a faster change rate. Fig. 2(h) shows that FANS-RL achieves the highest return across all tested non-stationary degrees and that the gap increases with the non-stationarity. We also test FANS-RL together with all baselines on different non-stationary functions, including piecewise linear, damping-like and sinusoid waves. Fig. 2(i) displays the results, which indicate that FANS-RL can generally outperform the baselines on diverse non-stationary function forms. Detailed function equations and experimental setups can be referred to Appendix C.

**Multiple change factors** We consider different numbers and types of changes to verify the benefits from the factored structure in FANS-RL. We conduct experiments with 1) only change wind forces (1D); 2) change wind forces and gravity concurrently (2D); 3) change wind force and target speed (1D+1R); and 4) change wind force, gravity, and target speed together (2D+1R) in an across-episode way in Half-Cheetah. From Fig. 2(j), we find that, thanks to the factored representation, FANS-RL performs better in those more complicated scenarios with multiple numbers and types of changes.

## 6 Related Work

**Non-stationary and transfer RL.** Early works in non-stationary and transfer RL [27, 28] only detect changes that have already happened instead of anticipating them. If the evolution of non-stationary environments is a (Semi)-Markov chain, one can deal with non-stationarity with HM-MDPs [29] or HS3MDPs [30]. Several methods learn to anticipate changes in non-stationary deep RL. Chandak et al. [31] propose to maximize future rewards without explicitly modeling non-stationary environments. MBCD [32] uses change-point detection to decide if the agent should learn a novel policy or reuse previously trained policies. Al-Shedivat et al. [4] extend MAML [6] for the non-stationary setting, but do not explicitly model the temporal changes. TRIO [5] tracks the non-stationarity by inferring the evolution of latent parameters, which captures the temporal change factors during the meta-testing phase. ReBAL and GrBAL [33] meta-train the dynamic prior, which adapts to the local contexts efficiently. However, these methods have to meta-train the model on a set of non-stationary tasks, which may not be accessible in real-world applications. Another line of research directly learns the latent representation to capture the non-stationary components. In particular, LILAC [7] and ZeUS [8] leverage latent variable models to directly model the change factors in environments, and Guo et al. [34] estimate latent vectors that describe the non-stationary or variable part of the dynamics. Perez et al. [35, 36] estimate low-dimensional latent vectors that capture the changes in dynamics and rewards and leverage these inferred change factors to facilitate transfer RL across tasks. Similarly, François-Lavet et al. [37] use a compact state abstraction for recovering the sufficient low-dimensional representation of the environment, and only retrain or update some of its abstracted states for transferring to a novel but related task. Our approach fits in this line of work; however, as opposed to these methods, which model changes using a shared embedding space or a mixture model, and cannot distinguish which state components are affected, we use a factored representation.

**Factored MDPs.** Among works on factored representations in stationary settings, Hallak et al. [38] learn factored MDPs [11, 12, 13] to improve sample efficiency in model-based off-policy RL. Balaji

et al. [39] employ known factored MDP to improve both model-free and model-based RL algorithms. Working memory graphs [40] learn the factored observation space using Transformers. NeverNet [41] factorizes the state-action space through graph neural networks. Zholus et al. [42] factorize the visual states into actor-related, object of manipulation, and the latent influence factor between these two states to achieve sparse interaction in robotic control tasks. Differently from methods modeling the factored dynamics and rewards only, Zhou et al. [43] and Tang et al. [44] explore factored entities and actions, respectively. Zhou et al. [43] extend the factored MDP to the multi-entity environments, learning the compositional structure in tasks with multiple entities involved. Tang et al. [44] leverage the factored action space to improve the sample efficiency in healthcare applications. However, the factored structures in these two works are derived from inductive bias or domain experts.

**Factored MDPs, causality and multiple environments.** Several works leverage factored MDPs to improve the sample efficiency and generalization of RL. Most of these works focus on learning an invariant (e.g. causal) representation that fits all environments and do not support learning latent change factors. For example, AFaR [45] learns factored value functions via attention modules to improve sample efficiency and generalization across tasks. Mutti et al. [46] learn a causal structure that can generalize across a family of MDPs under different environments, assuming that there are no latent causal factors. Similarly, Wang et al. [47] propose to learn the factored and causal dynamics in model-based RL, in which the learned causal structure is assumed to generalize to unseen states. By deriving the state abstraction based on the causal graph, it can improve both the sample efficiency and generalizability of policy learning of MBRL. While most of the previous works focus on learning a causal structure that is time-invariant, Pitis et al. [48] learn a locally causal dynamics that can vary at each timestep and use it to generate counterfactual dynamics transitions in RL. While the previously described methods focus on a domain or time-invariant representations, in our work we also focus on modelling domain or time-specific factors in the form of latent change factors. A related work, AdaRL [14] learns the factored representation and the model change factors under heterogeneous domains with varying dynamics or reward functions. However, AdaRL is designed only for the domain adaptation setting and constant change factors without considering non-stationarity.

**Independent causal mechanisms.** Another related line of work is based on independent causal mechanisms [49, 50]. Recurrent independent mechanisms (RIMs) [51] learn the independent transition dynamics in RL with sparse communication among the latent states. Meta-RIMs [52] leverage meta-learning and soft attention to learn a set of RIMs with competition and communication. As opposed to these works, we do not assume that the mechanisms are independent and we learn the factored structure among all components in MDPs with a DBN.

## 7 Conclusions, Limitations and Future Work

We describe Factored Adaptation for Non-Stationary RL (FANS-RL), a framework that learns a factored representation for non-stationarity that can be combined with any RL algorithm. We formalize our problem as a Factored Non-stationary MDP (FN-MDP), augmenting a factored MDP with latent change factors evolving as a Markov process. FN-MDPs do not model a family of MDPs, but instead include the dynamics of change factors analogously to the dynamics of the states. This allows us to capture different non-stationarities, e.g., continuous and discrete changes, both within and across different episodes. To learn FN-MDPs we propose FN-VAEs, which we integrated in FANS-RL, an online model estimation and policy evaluation approach. We evaluate FANS-RL on benchmarks for continuous control and robotic manipulation, also with pixel inputs, and show it outperforms the state of the art on rewards and robustness to varying degrees of non-stationarity. Learning the graph in model estimation is computationally expensive, which limits the scalability of our approach. In future work, we plan to meta-learn the graphs among different tasks to improve the scalability of our approach and its applicability to complex RL problems, e.g., multi-agent RL.

### Acknowledgments

FF would like to acknowledge the CityU High-Performance Computing (HPC) resources in Hong Kong SAR and LISA HPC from the SURF.nl. SM was supported by the MIT-IBM Watson AI Lab and the Air Force Office of Scientific Research under award number FA8655-22-1-7155. BH would like to acknowledge the support of Apple Scholarship. KZ was partially supported by the National Institutes of Health (NIH) under Contract R01HL159805, by the NSF-Convergence Accelerator Track-D award #2134901, by a grant from Apple Inc., and by a grant from KDDI Research Inc.

## References

- [1] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [2] Gabriel Dulac-Arnold, Nir Levine, Daniel J Mankowitz, Jerry Li, Cosmin Paduraru, Sven Gowal, and Todd Hester. Challenges of real-world reinforcement learning: definitions, benchmarks and analysis. *Machine Learning*, pages 1–50, 2021.
- [3] Sindhu Padakandla. A survey of reinforcement learning algorithms for dynamically varying environments. *ACM Computing Surveys (CSUR)*, 54(6):1–25, 2021.
- [4] Maruan Al-Shedivat, Trapit Bansal, Yura Burda, Ilya Sutskever, Igor Mordatch, and Pieter Abbeel. Continuous adaptation via meta-learning in nonstationary and competitive environments. In *International Conference on Learning Representations (ICLR)*, 2018.
- [5] Riccardo Poiani, Andrea Tirinzoni, and Marcello Restelli. Meta-reinforcement learning by tracking task non-stationarity. In *International Joint Conference on Artificial Intelligence (IJCAI)*, pages 2899–2905, 2021.
- [6] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International Conference on Machine Learning (ICML)*, pages 1126–1135. PMLR, 2017.
- [7] Annie Xie, James Harrison, and Chelsea Finn. Deep reinforcement learning amidst lifelong non-stationarity. *International Conference on Machine Learning (ICML)*, 2021.
- [8] Shagun Sodhani, Franziska Meier, Joelle Pineau, and Amy Zhang. Block contextual mdps for continual learning. In *Learning for Dynamics and Control Conference (LADC)*, pages 608–623. PMLR, 2022.
- [9] Biwei Huang, Kun Zhang, Jiji Zhang, Joseph D Ramsey, Ruben Sanchez-Romero, Clark Glymour, and Bernhard Schölkopf. Causal discovery from heterogeneous/nonstationary data. *Journal of Machine Learning Research*, 21(89):612–634, 2020.
- [10] Kun Zhang, Mingming Gong, Petar Stojanov, Biwei Huang, Qingsong Liu, and Clark Glymour. Domain adaptation as a problem of inference on graphical models. *Advances in Neural Information Processing Systems (NeurIPS)*, 33:4965–4976, 2020.
- [11] Craig Boutilier, Richard Dearden, and Moisés Goldszmidt. Stochastic dynamic programming with factored representations. *Artificial intelligence*, 121(1-2):49–107, 2000.
- [12] Michael Kearns and Daphne Koller. Efficient reinforcement learning in factored mdps. In *International Joint Conference on Artificial Intelligence (IJCAI)*, volume 16, pages 740–747, 1999.
- [13] Ian Osband and Benjamin Van Roy. Near-optimal reinforcement learning in factored mdps. *Advances in Neural Information Processing Systems (NeurIPS)*, 27:604–612, 2014.
- [14] Biwei Huang, Fan Feng, Chaochao Lu, Sara Magliacane, and Kun Zhang. Adarl: What, where, and how to adapt in transfer reinforcement learning. *International Conference on Learning Representations (ICLR)*, 2022.
- [15] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8): 1735–1780, 1997.
- [16] Lukas Liebel and Marco Körner. Auxiliary tasks in multi-task learning. *arXiv preprint arXiv:1805.06334*, 2018.
- [17] David Ha and Jürgen Schmidhuber. Recurrent world models facilitate policy evolution. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 2455–2467, 2018.

- [18] Zihao Zhao, Anusha Nagabandi, Kate Rakelly, Chelsea Finn, and Sergey Levine. MELD: meta-reinforcement learning from images via latent state models. In *Conference on Robot Learning (CORL)*, pages 1246–1261. PMLR, 2020.
- [19] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International Conference on Machine Learning (ICML)*, pages 1861–1870. PMLR, 2018.
- [20] Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 5026–5033. IEEE, 2012.
- [21] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- [22] Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Karol Hausman, Chelsea Finn, and Sergey Levine. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning. In *Conference on Robot Learning (CORL)*, pages 1094–1100. PMLR, 2020.
- [23] Jie Tan, Tingnan Zhang, Erwin Coumans, Atıl İscen, Yunfei Bai, Danijar Hafner, Steven Bohez, and Vincent Vanhoucke. Sim-to-real: Learning agile locomotion for quadruped robots. In *Robotics: Science and Systems*, 2018.
- [24] Kimin Lee, Younggyo Seo, Seunghyun Lee, Honglak Lee, and Jinwoo Shin. Context-aware dynamics model for generalization in model-based reinforcement learning. In *International Conference on Machine Learning (ICML)*, pages 5757–5766. PMLR, 2020.
- [25] Zhou Xian, Shamit Lal, Hsiao-Yu Tung, Emmanouil Antonios Platanios, and Katerina Fragkiadaki. Hyperdynamics: Meta-learning object and agent dynamics with hypernetworks. In *International Conference on Learning Representations (ICLR)*, 2021.
- [26] Luisa Zintgraf, Sebastian Schulze, Cong Lu, Leo Feng, Maximilian Igl, Kyriacos Shiarlis, Yarin Gal, Katja Hofmann, and Shimon Whiteson. Varibad: Variational bayes-adaptive deep rl via meta-learning. *Journal of Machine Learning Research*, 22(289):1–39, 2021.
- [27] Bruno C Da Silva, Eduardo W Basso, Ana LC Bazzan, and Paulo M Engel. Dealing with non-stationary environments using context detection. In *International Conference on Machine Learning (ICML)*, pages 217–224, 2006.
- [28] Richard S Sutton, Anna Koop, and David Silver. On the role of tracking in stationary environments. In *International Conference on Machine Learning (ICML)*, pages 871–878, 2007.
- [29] Samuel Choi, Dit-Yan Yeung, and Nevin Zhang. An environment model for nonstationary reinforcement learning. *Advances in neural information processing systems (NeurIPS)*, 12, 1999.
- [30] Emmanuel Hadoux, Aurélie Beynier, and Paul Weng. Solving hidden-semi-markov-mode markov decision problems. In *International Conference on Scalable Uncertainty Management*, pages 176–189. Springer, 2014.
- [31] Yash Chandak, Georgios Theodorou, Shiv Shankar, Martha White, Sridhar Mahadevan, and Philip Thomas. Optimizing for the future in non-stationary mdps. In *International Conference on Machine Learning (ICML)*, pages 1414–1425. PMLR, 2020.
- [32] Lucas N. Alegre, Ana L. C. Bazzan, and Bruno C. da Silva. Minimum-delay adaptation in non-stationary reinforcement learning via online high-confidence change-point detection. In *Proceedings of the 20th International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, pages 97–105, 2021.
- [33] Ignasi Clavera, Anusha Nagabandi, Simin Liu, Ronald S. Fearing, Pieter Abbeel, Sergey Levine, and Chelsea Finn. Learning to adapt in dynamic, real-world environments through meta-reinforcement learning. In *International Conference on Learning Representations (ICLR)*, 2019.

- [34] Jiaxian Guo, Mingming Gong, and Dacheng Tao. A relational intervention approach for unsupervised dynamics generalization in model-based reinforcement learnings. *International Conference on Learning Representations (ICLR)*, 2022.
- [35] Christian F Perez, Felipe Petroski Such, and Theofanis Karaletsos. Efficient transfer learning and online adaptation with latent variable models for continuous control. In *Continual Learning Workshop at NIPS*, 2018.
- [36] Christian Perez, Felipe Petroski Such, and Theofanis Karaletsos. Generalized hidden parameter mdps: Transferable model-based rl in a handful of trials. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, volume 34, pages 5403–5411, 2020.
- [37] Vincent François-Lavet, Yoshua Bengio, Doina Precup, and Joelle Pineau. Combined reinforcement learning via abstract representations. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, volume 33, pages 3582–3589, 2019.
- [38] Assaf Hallak, François Schnitzler, Timothy Mann, and Shie Mannor. Off-policy model-based learning under unknown factored dynamics. In *International Conference on Machine Learning (ICML)*, pages 711–719. PMLR, 2015.
- [39] Bharathan Balaji, Petros Christodoulou, Xiaoyu Lu, Byungsoo Jeon, and Jordan Bell-Masterson. Factoredrl: Leveraging factored graphs for deep reinforcement learning. *NeurIPS Deep Reinforcement Learning Workshop*, 2020.
- [40] Ricky Loynd, Roland Fernandez, Asli Celikyilmaz, Adith Swaminathan, and Matthew Hausknecht. Working memory graphs. In *International Conference on Machine Learning (ICML)*, pages 6404–6414. PMLR, 2020.
- [41] Tingwu Wang, Renjie Liao, Jimmy Ba, and Sanja Fidler. Nervenet: Learning structured policy with graph neural networks. In *International conference on learning representations (ICLR)*, 2018.
- [42] Artem Zhohus, Yaroslav Ivchenkov, and Aleksandr Panov. Factorized world models for learning causal relationships. In *ICLR Workshop on the Elements of Reasoning: Objects, Structure and Causality*, 2022.
- [43] Allan Zhou, Vikash Kumar, Chelsea Finn, and Aravind Rajeswaran. Policy architectures for compositional generalization in control. In *ICML Workshop on Spurious Correlations, Invariance and Stability*, 2022.
- [44] Shengpu Tang, Maggie Makar, Michael Sjoding, Finale Doshi-Velez, and Jenna Wiens. Leveraging factored action spaces for efficient offline reinforcement learning in healthcare. In *Decision Awareness in Reinforcement Learning Workshop at ICML*, 2022.
- [45] Shagun Sodhani, Sergey Levine, and Amy Zhang. Improving generalization with approximate factored value functions. In *ICLR Workshop on the Elements of Reasoning: Objects, Structure and Causality*, 2022.
- [46] Mirco Mutti, Riccardo De Santi, Emanuele Rossi, Juan Felipe Calderon, Michael Bronstein, and Marcello Restelli. Provably efficient causal model-based reinforcement learning for systematic generalization. *arXiv preprint arXiv:2202.06545*, 2022.
- [47] Zizhao Wang, Xuesu Xiao, Zifan Xu, Yuke Zhu, and Peter Stone. Causal dynamics learning for task-independent state abstraction. In *International Conference on Machine Learning (ICML)*, pages 23151–23180. PMLR, 2022.
- [48] Silviu Pitis, Elliot Creager, Ajay Mandlekar, and Animesh Garg. Mocoda: Model-based counterfactual data augmentation. In *Decision Awareness in Reinforcement Learning Workshop at ICML*, 2022.
- [49] Bernhard Schölkopf, Francesco Locatello, Stefan Bauer, Nan Rosemary Ke, Nal Kalchbrenner, Anirudh Goyal, and Yoshua Bengio. Toward causal representation learning. *Proceedings of the IEEE*, 109(5):612–634, 2021.

- [50] Giambattista Parascandolo, Niki Kilbertus, Mateo Rojas-Carulla, and Bernhard Schölkopf. Learning independent causal mechanisms. In *International Conference on Machine Learning (ICML)*, pages 4036–4044. PMLR, 2018.
- [51] Anirudh Goyal, Alex Lamb, Jordan Hoffmann, Shagun Sodhani, Sergey Levine, Yoshua Bengio, and Bernhard Schölkopf. Recurrent independent mechanisms. In *International Conference on Learning Representations (ICLR)*, 2021.
- [52] Kanika Madan, Nan Rosemary Ke, Anirudh Goyal, Bernhard Schölkopf, and Yoshua Bengio. Fast and slow learning of recurrent independent mechanisms. In *International Conference on Learning Representations (ICLR)*, 2021.
- [53] Peter Spirtes, Clark N Glymour, and Richard Scheines. *Causation, Prediction, and Search*. Springer-Verlag Lectures in Statistics, 1993.
- [54] Kevin Patrick Murphy. *Dynamic bayesian networks: representation, inference and learning*. University of California, Berkeley, 2002.
- [55] Judea Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge University Press, Cambridge, 2000.
- [56] Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. Causal inference on time series using restricted structural equation models. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 154–162, 2013.

## Checklist

1. For all authors...
  - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? **[Yes]** There are theoretical claims backed by two theoretical results in Prop. 1 and Prop. 2. The empirical claims are backed by an extensive evaluation in Section 5.
  - (b) Did you describe the limitations of your work? **[Yes]** We did provide a discussion on the limitations and potential future work in Section 7.
  - (c) Did you discuss any potential negative societal impacts of your work? **[Yes]** We discussed the broader societal impact in Appendix A.
  - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? **[Yes]**
2. If you are including theoretical results...
  - (a) Did you state the full set of assumptions of all theoretical results? **[Yes]** The assumptions are stated in Prop. 1 and Prop. 2.
  - (b) Did you include complete proofs of all theoretical results? **[Yes]** The complete proofs are in the Appendix B
3. If you ran experiments...
  - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplementary material or as a URL)? **[Yes]** We will include the code and instructions in the supplementary files.
  - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? **[Yes]** The details are in the Appendix D.4.
  - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? **[Yes]** The learning curves (Fig. 2 (a-c), Appendix Fig. A2-A3) include error bars with 10 runs.
  - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? **[Yes]** The information of computational resource is given in Appendix E.1.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
  - (a) If your work uses existing assets, did you cite the creators? **[Yes]** Yes, we have cited the OpenAI Gym and Mujoco libraries in our paper.
  - (b) Did you mention the license of the assets? **[Yes]** Yes, the licenses are given in Appendix E.2.
  - (c) Did you include any new assets either in the supplemental material or as a URL? **[N/A]**
  - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? **[N/A]**
  - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? **[N/A]**
5. If you used crowdsourcing or conducted research with human subjects...
  - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? **[N/A]**
  - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? **[N/A]**
  - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? **[N/A]**