



Deposited via The University of Sheffield.

White Rose Research Online URL for this paper:

<https://eprints.whiterose.ac.uk/id/eprint/222340/>

Version: Preprint

---

**Preprint:**

Han, C., Basu, D., Mangan, M. et al. (Submitted: 2024) Dynamical-VAE-based hindsight to learn the causal dynamics of factored-POMDPs. [Preprint - arXiv] (Submitted)

<https://doi.org/10.48550/arXiv.2411.07832>

---

© 2024 The Author(s). This preprint is made available under a Creative Commons Attribution 4.0 International License. (<https://creativecommons.org/licenses/by/4.0/>)

**Reuse**

This article is distributed under the terms of the Creative Commons Attribution (CC BY) licence. This licence allows you to distribute, remix, tweak, and build upon the work, even commercially, as long as you credit the authors for the original work. More information and the full terms of the licence here:

<https://creativecommons.org/licenses/>

**Takedown**

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing [eprints@whiterose.ac.uk](mailto:eprints@whiterose.ac.uk) including the URL of the record and the reason for the withdrawal request.

# Dynamical-VAE-based Hindsight to Learn the Causal Dynamics of Factored-POMDPs

**Chao Han**

*School of Computer Science, The University of Sheffield, UK*

C.HAN@SHEFFIELD.AC.UK

**Debabrota Basu**

*Equipe Scool, Univ. Lille, Inria, CNRS, Centrale Lille, UMR 9189 CRIStal, France*

DEBABROTA.BASU@INRIA.FR

**Michael Mangan**

*School of Computer Science, The University of Sheffield, UK*

M.MANGAN@SHEFFIELD.AC.UK

**Eleni Vasilaki**

*School of Computer Science, The University of Sheffield, UK*

E.VASILAKI@SHEFFIELD.AC.UK

**Aditya Gilra**

*Machine Learning group, Centrum Wiskunde & Informatica, Amsterdam, Netherlands*

ADITYA.GILRA@CWI.NL

*School of Computer Science, The University of Sheffield, UK*

## Abstract

Learning representations of underlying environmental dynamics from partial observations is a critical challenge in machine learning. In the context of Partially Observable Markov Decision Processes (POMDPs), state representations are often inferred from the history of past observations and actions. We demonstrate that incorporating future information is essential to accurately capture causal dynamics and enhance state representations. To address this, we introduce a Dynamical Variational Auto-Encoder (DVAE) designed to learn causal Markovian dynamics from offline trajectories in a POMDP. Our method employs an extended hindsight framework that integrates past, current, and multi-step future information within a factored-POMDP setting. Empirical results reveal that this approach uncovers the causal graph governing hidden state transitions more effectively than history-based and typical hindsight-based models.

**Keywords:** Causal Discovery, Representation Learning, POMDP, Variational Autoencoders

## 1. Introduction

Accurately learning the underlying dynamics of an environment is essential for developing models that can reliably predict future states, particularly in partially observable settings (Wang et al., 2019; Moerland et al., 2023). Existing self-predictive approaches to state representation aim to learn a Markovian transition model (Ni et al., 2024). However, in partially observable contexts, the true underlying state remains hidden, making it necessary to construct an approximate belief state from prior state-action histories as a proxy for the latent state. This approach effectively reformulates the Partially Observable Markov Decision Process (POMDP) as a Markov Decision Process (MDP) that depends solely on past observations and actions to approximate the full state information (Åström, 1965; Subramanian et al., 2022). Such an approach may, in general, only lead to an approximation of the true MDP.

In online settings, the agent is limited to past information alone, but in offline RL or model learning, both past and future data around each time step are accessible. This availability raises the question of whether combining both past and future information can improve our ability to identify

the generating MDP. By maximizing the log-likelihood of complete trajectories of observations and actions, we leverage the formalism of Dynamical Variational Auto-Encoders (DVAE) (Girin et al., 2020) to determine which elements of the past and future are essential for decoding unobservable variables at each time step. We separate unobservable variables into deterministic hidden ones, and using the Reparameterization Lemma (Buesing et al., 2018), into exogenous stochastic ones. We find that the 1-step past (including bootstrapped hidden), present, and future observables and actions are needed to accurately reconstruct deterministic unobserved hidden variables. We term our approach “DVAE-based hindsight” to contrast it with prior hindsight methods for latent identification that utilized only the present and 1-step future (Jarrett et al., 2023).

We utilize Causal Dynamical Learning (CDL) (Wang et al., 2022), employing Conditional Mutual Information (CMI), to learn a causal transition graph of the environment. The stationary Markovian transition model can be represented as a Directed Acyclic Graph (DAG), mapping the Markovian states and action at time step  $t$  to the Markovian states at  $t + 1$ . We extend CDL to a partially observable setting by learning to identify deterministic hidden variables and constructing the causal transition graph, combining the DVAE and CDL approaches in an end-to-end framework. We demonstrate the effectiveness of our approach against history-based (Littman and Sutton, 2001; Baisero and Amato, 2020; Ni et al., 2024) and earlier hindsight-based methods (Jarrett et al., 2023), in a factored-POMDP setting (Oliehoek et al., 2021) which highlights the advantages of our method.

## 2. Preliminaries and Problem Formulation

### 2.1. Partially Observable Markov Decision Processes (POMDPs)

A Markov Decision Process (MDP) in the context of reinforcement learning is defined by a tuple  $(S, A, T_a, R_a)$ , where  $S$  is the set of states,  $A$  the set of actions,  $T_a(s'|s)$  the probability of transitioning from state  $s$  to  $s'$  under action  $a$ , and  $R_a(s', s)$  the reward received for this transition. However, many real-world systems or environments are only partially observable. It is typically assumed that there exists an underlying or generating MDP that gives rise to a Partially Observable Markov Decision Process (POMDP)  $(S, A, T_a, R_a, \Omega, O)$ , where the states are not directly observable. Instead, we observe elements  $o$  from a set  $\Omega$ , governed by conditional probabilities  $O(o|s)$ . A POMDP can be converted into an MDP by relying solely on the history of observations and actions (Åström, 1965). This approach forms the basis for using a sequence of past observations (or their representation) and actions as a proxy, or belief state, for the environment’s current state (Subramanian et al., 2022).

### 2.2. Problem formulation: Learning the causal dynamics underlying a factored-POMDP

Our objective is to learn the underlying state transitions and associated causal graph (represented in Figure 1) from offline data in a factored-POMDP environment. A factored-POMDP (Oliehoek et al., 2021) allows us to focus on learning the underlying transition function and graph, without additional details of representation learning.

**Definition 1 (Factored-POMDP)** (Oliehoek et al., 2021). A factored partially observable Markov decision process is defined as a tuple  $\langle S, O, H, A, T, R, \bar{O} \rangle$  where:

- the state space  $S$  is spanned as  $S = S^1 \times \dots \times S^{d_S}$  (each state variable  $S^k$  is called a factor), such that every state  $s \in S$  is a  $d_S$ -dimension vector  $s = (s^1, \dots, s^{d_S})$ .
- the space of observed states  $O \subseteq S$  is denoted as  $O = O^1 \times \dots \times O^{d_O}$  with  $d_O \leq d_S$ .

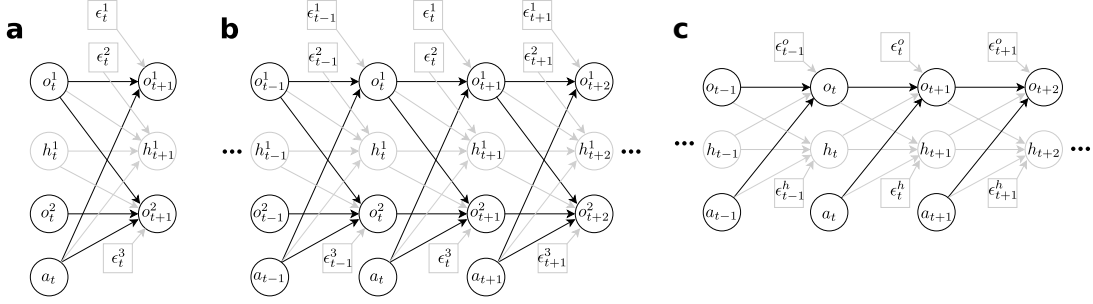


Figure 1: **(a)** The stationary transition model of a factored-POMDP is shown as a Structural Causal Model (SCM) from time step  $t$  to  $t + 1$ . The factored states are represented as circle nodes, which are deterministic as per Eq. (1). They can be either observed (black) or hidden (gray). Gray squares represent unobserved exogenous (i.e. no parents) stochastic nodes. The arrows connecting nodes represent directed causal edges from parents to children. The connectivity of the deterministic nodes is only an example. **(b)** The stationary transition model can be unrolled over time, by repeating the graph in panel (a). over multiple time steps, to obtain a SCM for a full trajectory. **(c)** We collect hidden factored states into vector  $h$ , and observable factored states into vector  $o$  while maintaining the underlying causal model. This is the general SCM for any factored-POMDP.

- the space of hidden states  $H \subseteq S$  is spanned as  $H = H^1 \times \dots \times H^{d_H}$  with  $d_H \leq d_S$ .
- $O \cup H = S, O \cap H = \emptyset$ , such that  $s = (o, h)$ .
- $A$  is the set of actions  $a$ .
- $T(s_{t+1} | s_t, a_t) = \prod_{j=1}^{d_H} \prod_{i=1}^{d_O} p(h_{t+1}^j | s_t, a_t) p(o_{t+1}^i | s_t, a_t)$  is the transition probability function.
- $R(s_t, a_t, s_{t+1})$  is the reward function
- $\bar{O}(o_t | s_t)$  is the observation probability function that outputs 1 if  $o_t \in O$  is subvector of  $s_t \in S$  and 0 otherwise.

In this factored-POMDP setting, the state  $s$  is represented as a concatenation of observed and hidden states, denoted by  $s = (o, h)$ . The state transition probability distribution  $T$  can be factorised as  $T(s_{t+1} | s_t, a_t) = \prod_{j=1}^{d_S} p(s_{t+1}^j | s_t, a_t)$ . Consequently, our goal reduces to learning the factored transitions  $p(o_{t+1}^j | \{s_t^i\}_{i=1}^{d_S}, a_t)$  for  $j = 1, \dots, d_O$  and  $p(h_{t+1}^j | \{s_t^i\}_{i=1}^{d_S}, a_t)$  for  $j = 1, \dots, d_H$ .

**Representing stochasticity in transitions as independent exogenous noise.** Via the Reparameterization Lemma (Appendix B of Buesing et al. (2018)), we can always reparameterize the stochasticity to be exogenous, and write the probabilistic MDP transition of factored state variables as a Structural Causal Model (SCM)

$$s_{t+1}^i := f_i(\mathbf{PA}_{s_{t+1}^i}, a_t, \epsilon_t^i), \quad i = 1, \dots, d_S \quad (1)$$

where each  $f_i$  represents an arbitrary deterministic function.  $\mathbf{PA}_{s_{t+1}^i}$  denotes the set of parent state factors at time  $t$ , of  $s_{t+1}^i$ , such that there exists an edge from each element  $s_t^j \in \mathbf{PA}_{s_{t+1}^i}$  to  $s_{t+1}^i$  in the transition graph  $\mathcal{G}$ . Action  $a_t$  is represented separately for clarity. The exogenous noise variable  $\epsilon_t^i$  for each factor  $i$  is jointly independent at each time step  $t$ , that is  $p_{\epsilon_t^1, \dots, \epsilon_t^{d_S}} = \prod_{i=1}^{d_S} p_{\epsilon_t^i}$ . This noise variable can be seen as introducing stochasticity in the transitions, such that every  $s_{t+1}^i = f_i(\mathbf{PA}_{s_{t+1}^i}, a_t^i, \epsilon_t^i)$  is a sample drawn from  $p(s_{t+1}^i | \mathbf{PA}_{s_{t+1}^i}, a_t^i)$ , for every  $\epsilon_t^i$ , consistent with the reparameterization lemma (Buesing et al., 2018). Thus, in Fig. 1, we can represent all stochasticity in transitions with independent exogenous noise nodes.

Furthermore, any stochastic factored-MDP can be converted to the factored-POMDP setting by hiding a set of factors  $h = (h^1, \dots, h^{d_H})$  from the agent. From the perspective of an agent, the uncertainty in predicting the next observables  $o_{t+1}^i$  from the current observables and action, in such a setting, arises from two sources: the effect of current values of hidden factors  $h_t$  and the unobservable stochasticity in the transition encapsulated by the current noise  $\epsilon_t^i$ . Therefore, if we somehow had access to the current hidden states  $h_t$  and the noise  $\epsilon_t^i$ , then each next state  $s_{t+1}^i$  would be deterministically predictable given the current observed states  $o_t$  and action  $a_t^i$ . For our factored-POMDP, similar to examples in real life, both  $h_t$  and  $\epsilon_t^i$  are not observable.

### 3. Deriving the algorithm for learning the transition dynamics of factored-POMDPs

In subsection 3.1, we derive the DVAE-based framework for identifying the transition model using our extended hindsight encoder for hidden factors. In subsection A.3, we outline how we estimate the transition graph. In subsection 3.3, we outline our Modulo environment, an example factored-POMDP to illustrate our results.

#### 3.1. DVAE for Factored-POMDP

We aim to maximize the conditional marginal log-likelihood of the observations  $o_{1:T}$  given the actions  $a_{1:T}$ , parameterized by  $\theta$ , under the true data distribution  $p(o_{1:T}|a_{1:T})$ :

$$\max_{\theta} \mathbb{E}_{p(o_{1:T}|a_{1:T})} [\log p_{\theta}(o_{1:T}|a_{1:T})] \quad (2)$$

By introducing a variational distribution  $q_{\phi}(h_{1:T}|o_{1:T}, a_{1:T})$ , parameterized by  $\phi$ , we can decompose the objective in Eq. (2) as follows (see Appendix A.1 for derivation):

$$\max_{\theta, \phi} \mathbb{E}_{p(o_{1:T}|a_{1:T})} [\ell_{\text{VLB}}(\theta, \phi; o_{1:T}, a_{1:T}) + D_{\text{KL}}(q_{\phi}(h_{1:T}|o_{1:T}, a_{1:T}) \parallel p_{\theta}(h_{1:T}|o_{1:T}, a_{1:T}))] \quad (3)$$

Here,  $\ell_{\text{VLB}}(\theta, \phi; o_{1:T}, a_{1:T})$  is the variational lower bound (VLB) on the marginal log-likelihood, serving as a lower bound due to the non-negativity of the KL divergence term. VLB is defined as:

$$\ell_{\text{VLB}}(\theta, \phi; o_{1:T}, a_{1:T}) = \mathbb{E}_{q_{\phi}(h_{1:T}|o_{1:T}, a_{1:T})} [\log p_{\theta}(o_{1:T}, h_{1:T}|a_{1:T}) - \log q_{\phi}(h_{1:T}|o_{1:T}, a_{1:T})] \quad (4)$$

Thus, optimizing Eq. (2) reduces to maximizing the expected VLB. In practice, we approximate the expectation of the data distribution  $p(o_{1:T}|a_{1:T})$ , using observed data sequences. We employ independent and identically distributed (i.i.d.) sampled trajectories from the collected dataset  $\mathcal{D}$  to construct a Monte Carlo estimate of the expected VLB, defined as follows:

$$\mathcal{L}_{\text{VLB}}(\theta, \phi; o_{1:T}, a_{1:T}) = \mathbb{E}_{(o_{1:T}, a_{1:T}) \sim \mathcal{D}} [\ell_{\text{VLB}}(\theta, \phi; o_{1:T}, a_{1:T})] \quad (5)$$

Leveraging the Markov property in the transition dynamics, the generative model  $p_{\theta}$  and inference model  $q_{\phi}$  in VLB of Eq. (4) can be further decomposed along time indices and state factors as follows (see Appendix A.2 for derivation):

$$p_{\theta}(o_{1:T}, h_{1:T}|a_{1:T}) = \prod_{t=0}^{T-1} \prod_{j=1}^{d_H} \prod_{i=1}^{d_O} p_{\theta_h}(h_{t+1}^j | s_t, a_t) p_{\theta_o}(o_{t+1}^i | s_t, a_t), \quad (6)$$

$$q_\phi(h_{1:T}|o_{1:T}, a_{1:T}) = \prod_{t=0}^{T-1} \prod_{j=1}^{d_H} q_\phi(h_{t+1}^j | h_t, o_{t:T}, a_{t:T}) \quad (7)$$

where  $o_t = (o_t^1, \dots, o_t^{d_O})$ ,  $h_t = (h_t^1, \dots, h_t^{d_H})$  and  $s_t = (o_t, h_t)$ . Note that the decomposed terms  $p_{\theta_h}(h_{t+1}^j | s_t, a_t)$  and  $p_{\theta_o}(o_{t+1}^i | s_t, a_t)$  can be interpreted as the transition probabilities of the  $j$ -th hidden state and  $i$ -th observed state, respectively. Similarly,  $q_\phi(h_{t+1}^j | h_t, o_t, a_t)$  serves as the encoder for the  $j$ -th hidden state, which we refer to as the *DVAE-based hindsight encoder*.

**Remark 1 (DVAE-based hindsight encoder for inferring the current hidden)** Eq. (7) shows that the joint conditional of the hidden states can be decomposed into  $T$  conditionals, each conditioned on 1-step past, current and all future observations and actions, as well as the 1-step past hidden states. The previous hidden states are recursively chained across the  $T$  time steps, effectively incorporating the entire past. Thus, the hidden encoder  $q_\phi(h_{t+1}^j | h_t, o_{t:T}, a_{t:T})$  systematically leverages all available information to infer the distribution of hidden states.

**Remark 2 (History-based encoder vs. DVAE-based hindsight encoder)** A history-based encoder, which conditions only on past and current observables and actions, cannot fully infer the current hidden, as it depends on an exogenous noise variable that is independent of past and current observations and actions (Fig. 1). In contrast, future observations, which depend on the current hidden state and carry this noise information, are utilized by our DVAE-based hindsight encoder.

**Remark 3 (Current and 1-step hindsight encoder vs. DVAE-based hindsight encoder)** Specifically, Eq. (1) can be rewritten as  $o_{t+1}^i := f_i(o_t, h_t^j, a_t, \epsilon_t^i)$ , for every  $i$  and  $j$ . An encoder conditioned on  $o_t$ ,  $a_t$ , and  $o_{t+1}^i$ , as in Jarrett et al. (2023), would infer  $h_t^j$  by inverting the transition function of the parent of  $h_t^j$ , i.e.,  $o_{t+1}^i$ . However, the inferred  $h_t^j$  would be contaminated with  $\epsilon_t^i$ . In fact, they do not include any hidden states in their environment, encoding only the exogenous noise using current and 1-step future observations and actions. In our DVAE-based hindsight encoder, the additional information from the 1-step past, along with further future observations and actions and the bootstrapped 1-step past hidden state, better disentangles the current hidden state from the exogenous noise.

Substituting Eqs. (6) and (7) into Eq. (4) yields:

$$\begin{aligned} \ell_{\text{VLB}}(\theta, \phi, \bar{\phi}; o_{1:T}, a_{1:T}) = & \sum_{t=0}^{T-1} \mathbb{E}_{q_\phi(h_{1:t}|o_{1:T}, a_{1:T})} \left[ \sum_{j=1}^{d_O} \log p_{\theta_o}(o_{t+1}^j | s_t, a_t) \right. \\ & \left. - \sum_{j=1}^{d_H} D_{\text{KL}}(q_{\bar{\phi}}(h_{t+1}^j | h_t, o_{t:T}, a_{t:T}) \parallel p_{\theta_h}(h_{t+1}^j | s_t, a_t)) \right] \quad (8) \end{aligned}$$

Here,  $q_{\bar{\phi}}(h_{t+1}^j | h_t, o_{t:T}, a_{t:T})$  serves as the target distribution of the next encoded hidden state in the KL divergence term, comparing it to the distribution of the next predicted hidden state  $p_{\theta_h}(h_{t+1}^j | s_t, a_t)$ . The notation  $\bar{\phi}$  denotes the stop-gradient version of  $\phi$ , which is detached from the computation graph and replaced by a copy of  $\phi$  from the previous training step. Using a stop-gradient target in self-predictive representations is common in practice (Zhang et al., 2020; Ghugare et al., 2022), as this technique helps avoid representational collapse (Ni et al., 2024).

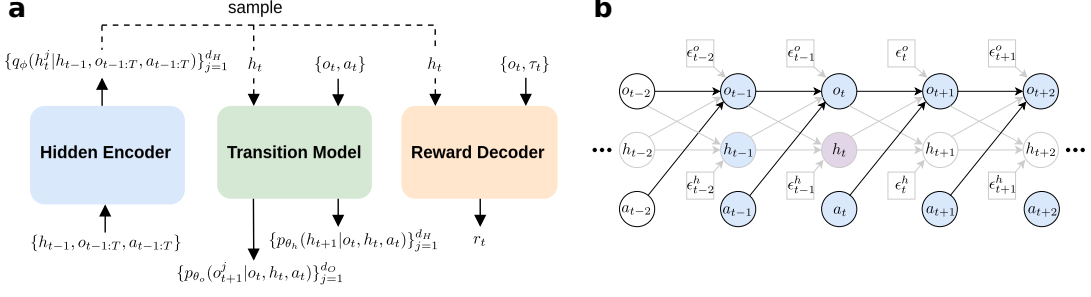


Figure 2: **(a)** Model architecture for computing the objective function in Eq. (11). **(b)** The current hidden state (light purple) is inferred from the hidden states, observables, and action of the previous step, along with the current and all future observables and actions (light blue) within the DVAE-based hindsight encoder.

Building upon this, to infer the factor-wise transition graph, we employ the approach from Wang et al. (2022). This involves modifying Eq. (8) to include 3 types of losses — full, masked, and causal. Thus, in addition to the full transition distribution (without masking any input factor), we compute a masked transition distribution by masking a randomly chosen state factor  $s_t^i$  or action  $a_t$  from each input  $\{s_t, a_t\}$  to the transition model, and also a causal transition distribution by masking out all input factors except for causal parents of  $s_{t+1}^j$  identified using the transition graph learned so far (see section 3.2). These are used in the 3 loss types, for both the Negative Log-Likelihood (NLL) of observed states and the KL-divergence (KL-div) of hidden states, to yield 6 loss terms:

$$\begin{aligned}
 \ell_{\text{VLB}}(\theta, \phi, \bar{\phi}; o_{1:T}, a_{1:T}) = & \sum_{t=0}^{T-1} \mathbb{E}_{q_{\phi}(h_{1:t} | o_{1:T}, a_{1:T})} \left[ - \sum_{j=1}^{d_O} \underbrace{[-\log p_{\theta_o}(o_{t+1}^j | s_t, a_t)]}_{\text{Full NLL Loss}} - \log p_{\theta_o}(o_{t+1}^j | s_t \setminus s_t^i, a_t) \right. \\
 & \underbrace{- \log p_{\theta_o}(o_{t+1}^j | \mathbf{PA}_{o_{t+1}^j})}_{\text{Causal NLL Loss}} - \sum_{j=1}^{d_H} \underbrace{[D_{\text{KL}}(q_{\bar{\phi}}(h_{t+1}^j | h_t, o_{t:T}, a_{t:T}) \| p_{\theta_h}(h_{t+1}^j | s_t, a_t))]}_{\text{Full KL-Div Loss}} \\
 & \underbrace{+ D_{\text{KL}}(q_{\bar{\phi}}(h_{t+1}^j | h_t, o_{t:T}, a_{t:T}) \| p_{\theta_h}(h_{t+1}^j | s_t \setminus s_t^i, a_t))}_{\text{Masked KL-Div Loss}} + \underbrace{D_{\text{KL}}(q_{\bar{\phi}}(h_{t+1}^j | h_t, o_{t:T}, a_{t:T}) \| p_{\theta_h}(h_{t+1}^j | \mathbf{PA}_{h_{t+1}^j}))}_{\text{Causal KL-Div Loss}} \left. \right] \quad (9)
 \end{aligned}$$

Here,  $s_t \setminus s_t^i = \{s_t^1, \dots, s_t^{i-1}, s_t^{i+1}, \dots, s_t^{d_S}\}$  denotes the set of all state factors at time  $t$  except for the  $i$ -th factor  $s_t^i$ . For each  $j$ , the index  $i$  is uniformly sampled from  $\{1, \dots, d_S\}$ . The term  $\mathbf{PA}_{s_{t+1}^j}$  are inferred from the learned transition graph so far using the conditional mutual information between each pair of factors, as discussed in Section 3.2.

Finally, the KL-div term in Eq. (9) which matches the distributions between the encoded next hidden state and the predicted next hidden state, can lead to convergence to a trivial constant representation of the hidden state (Ni et al., 2024). To prevent such degeneration, additional constraints on the hidden representation need to be applied alongside the VLB. Here, we employ a reward predictor parameterized by  $\psi$  to condition the encoded hidden representations  $h_{1:T}$ , trained by minimizing the prediction error:

$$\mathcal{L}_{\text{rew}}(\phi, \psi; o_{1:T}, \tau_{1:T}, r_{1:T}) = \mathbb{E}_{\substack{(o_{1:T}, \tau_{1:T}, r_{1:T}) \sim \mathcal{D} \\ h_{1:T} \sim q_{\phi}(h_{1:T} | o_{1:T}, a_{1:T})}} [\ell_{\text{rew}}(\psi; o_{1:T}, h_{1:T}, \tau_{1:T}, r_{1:T})] \quad (10)$$

**Algorithm 1:** Causal Dynamics Learning with Hindsight

**Input:** Initial hidden encoder  $q_\phi$ , initial transition models  $p_{\theta_o}$  and  $p_{\theta_h}$ , initial reward predictor  $R_\psi$ , and replay buffer  $\mathcal{D}$  containing pre-collected data.

**Parameters:** Learning rate  $\alpha > 0$ , CMI threshold  $\delta > 0$ , training steps  $M$ , CMI eval. period  $N$ .

**Output:** Converged hidden encoder  $q_{\phi^*}$ , transition models  $p_{\theta_o^*}$ ,  $p_{\theta_h^*}$ , and graph  $\mathcal{G}^*$ , reward predictor  $R_{\psi^*}$ .

**for**  $k = 1$  to  $M$  training steps **do**

Update  $\mathcal{D}$  and randomly sample a minibatch of  $m$  episodes  $\left\{ o_{1:T}^{(e)}, a_{1:T}^{(e)}, \tau_{1:T}^{(e)}, r_{1:T}^{(e)} \right\}_{e=1}^m$

Compute the mean objective  $\mathcal{L}_{\text{obj}} \left( \theta, \phi, \bar{\phi}, \psi; o_{1:T}^{(1:m)}, a_{1:T}^{(1:m)}, \tau_{1:T}^{(1:m)}, r_{1:T}^{(1:m)} \right)$  using Eq. (9)

Update the model parameters;

$$\begin{aligned} [\theta_o, \theta_h, \phi, \psi] &\leftarrow [\theta_o, \theta_h, \phi, \psi] + \alpha \nabla \mathcal{L}_{\text{obj}} \left( \theta, \phi, \bar{\phi}, \psi; o_{1:T}^{(1:m)}, a_{1:T}^{(1:m)}, \tau_{1:T}^{(1:m)}, r_{1:T}^{(1:m)} \right) \\ \bar{\phi} &\leftarrow \phi \end{aligned}$$

**if**  $k \bmod N = 0$  **then**

Evaluate  $\text{CMI}^{i,j}$  using Eqs. (12) and (13), and update it with an exponential moving average;

Binarize  $\text{CMI}^{i,j}$  to construct  $\mathcal{G}$  by checking if  $\text{CMI}^{i,j} \geq \delta$ ;

**end**

**end**

Here,  $\tau_t$  denotes any reward-related variables (e.g., a time-dependent/episodic target) used to predict the reward accurately.  $\ell_{\text{rew}}$  can be any supervised loss function; in our experiments, we use cross-entropy loss for categorical rewards.

Combining all components, we obtain the final objective to be minimized. This objective is a weighted sum of the mean VLB from Eqs. (5) and (9), and mean reward loss from Eq. (10), with a weight coefficient  $\lambda > 0$ :

$$\mathcal{L}_{\text{obj}} \left( \theta, \phi, \bar{\phi}, \psi; o_{1:T}, a_{1:T}, \tau_{1:T}, r_{1:T} \right) = -\mathcal{L}_{\text{VLB}} \left( \theta, \phi, \bar{\phi}; o_{1:T}, a_{1:T} \right) + \lambda \mathcal{L}_{\text{rew}} \left( \phi, \psi; o_{1:T}, \tau_{1:T}, r_{1:T} \right) \quad (11)$$

The model architecture depicted in Fig. 2a illustrates that every hidden states  $h_t^j$  is obtained through temporally recursive sampling from  $q_\phi(h_\tau^j | h_{\tau-1}, o_{\tau-1:T}, a_{\tau-1:T})$  for  $\tau = 1$  to  $t$ . Then, the hidden sample at each time step  $t$  is fed into the transition model and reward decoder to predict next states and reward. The unrolled probabilistic transition graph in Fig. 2b highlights the temporal data used as inputs to the DVAE-based hindsight encoder for the hidden states.

### 3.2. Transition Graph Estimation

The causal dependency of each transition pair  $s_t^i \rightarrow s_{t+1}^j$  or  $a_t \rightarrow s_{t+1}^j$  is estimated through conditional mutual information (CMI) (Wang et al., 2022). During evaluation, the CMI is computed based on two learned transition distributions: the full transition model  $p_\theta(s_{t+1}^j | s_t, a_t)$ , which leverages all state variables and the action to predict the next state of the  $j$ -th factor, and the masked transition model  $p_\theta(s_{t+1}^j | s_t \setminus s_t^i, a_t)$ , which relies on all state factors except for  $s_t^i$  for prediction.

Specifically, when the next state  $s_{t+1}^j$  is observable (denoted as  $o_{t+1}^j$ ), the CMI<sup>*i,j*</sup> between  $s_t^i$  and  $o_{t+1}^j$  given  $\{s_t \setminus s_t^i, a_t\}$  is formulated as:

$$I(s_t^i; o_{t+1}^j | s_t \setminus s_t^i, a_t) = \mathbb{E}_{s_t, a_t, o_{t+1}^j \sim \mathcal{D}, q_\phi} \left[ \log \frac{p_{\theta_o}(o_{t+1}^j | s_t, a_t)}{p_{\theta_o}(o_{t+1}^j | s_t \setminus s_t^i, a_t)} \right] \quad (12)$$

Here,  $s_t^i$  can be either an observed state or a hidden state sampled from the hidden encoder. The expectation in the CMI is approximated by aggregating transitions from all episodes in a mini-batch.

When the next state  $s_{t+1}^j$  is hidden (denoted as  $h_{t+1}^j$ ), the CMI<sup>*i,j*</sup> between  $s_t^i$  and  $h_{t+1}^j$  conditioned on  $\{s_t \setminus s_t^i, a_t\}$  is given by:

$$I(s_t^i; h_{t+1}^j | s_t \setminus s_t^i, a_t) = \mathbb{E}_{s_t, a_t \sim \mathcal{D}, q_\phi} \left[ D_{\text{KL}}(p_{\theta_h}(h_{t+1}^j | s_t, a_t) \parallel p_{\theta_h}(h_{t+1}^j | s_t \setminus s_t^i, a_t)) \right] \quad (13)$$

The derivations of Eqs. (12) and (13) are provided in Appendix A.3. Note that for causal dependency between the action and the next state,  $a_t \rightarrow s_{t+1}^j$ , the same CMI formula applies by replacing  $s_t^i$  with  $a_t$  in the conditioning set, which then becomes  $\{s_t\}$ .

In practice, the existence of an edge in the transition graph, i.e.,  $s_t^i \rightarrow s_{t+1}^j$  or  $a_t \rightarrow s_{t+1}^j$ , is determined by whether the corresponding CMI value CMI<sup>*i,j*</sup> exceeds a predefined threshold  $\delta$ . The binarized CMI matrix is then applied to select the parents of each next state in the causal transition losses in Eq. (9), and thus, refines learning of the causal transition dynamics  $p_\theta(s_{t+1}^j | \mathbf{PA}_{s_{t+1}^j})$ .

### 3.3. Modulo environment: a stochastic, discrete state-action, factored-POMDP

Modified from Ke et al. (2021), we construct a probabilistic discrete Factored-POMDP environment, to examine the performance of our model on inferring the hidden states and underlying transition graph. We called this environment modulo environment as the modulo operator is involved in its transition dynamics defined as  $s_{t+1} := (As_t + a_t + \epsilon_{t+1}) \bmod l$ , where  $l$  denotes the number of possible discrete values and  $A$  is the adjacency matrix of the transition graph  $\mathcal{G}$ . At time step  $t$ , each discrete factor  $s_t^i$  of the state vector  $s_t = (s_t^1, \dots, s_t^{d_S})^\top$  has values within  $\{0, \dots, l-1\}$ , the binary element  $a_t^i$  of the action vector  $a_t = (a_t^1, \dots, a_t^{d_S})^\top$  represents if the  $i$ -th factor is intervened or not by setting  $a_t^i = 1$  or 0 respectively, and the noise vector  $\epsilon_t = (\epsilon_t^1, \dots, \epsilon_t^{d_S})^\top \in E$  is sampled from a jointly independent distribution  $p_{\epsilon_t} = \prod_{i=1}^{d_S} p_{\epsilon_t^i}$ . Fig. 3 depicts noise-free transition dynamics with different underlying transition graph structures.

Our environment satisfies two properties. **(P1)** For every hidden factor  $h_t^i$ , there exists at least one observable  $o_{t+1}^j$ , such that  $h_t^i \in \mathbf{PA}_{o_{t+1}^j}$ . **(P2)** The transition map  $f \equiv \{f_i\}_{i=1}^{d_S}$  in Eq. (1), for every  $a \in A$  and every  $\epsilon \in E$ , i.e.  $(f)_{a,\epsilon} : S \rightarrow S$  from any  $s_t \in \mathcal{S}$  to  $s_{t+1} \in \mathcal{S}$ , is bijective, where  $s$  is the full state with all observable and hidden factors.

Indeed, by assuming a version of **(P2)**, in an environment with only exogenous noise but no hidden factors, we can deterministically infer these exogenous noise variables at  $t$ , by using a current and 1-step hindsight encoder for the hidden similar to the latent generator in Jarrett et al. (2023), which learns to invert  $f$  using observables and action at current  $t$  and observables at 1-step future  $t+1$ . However, *with both hidden factors and exogenous noise, despite these simplifying properties, history-based, and current and 1-step hindsight-based approaches are unable to learn the hidden factor and the graph*, as shown by the following experiments (see also Remark 3).

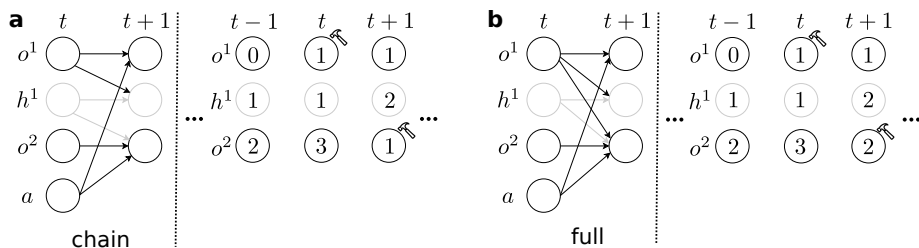


Figure 3: Illustration of Modulo environment with different types of transition graphs that have  $d_S = 3$  and  $l = 4$ . **(a)** Chain structure. left: ground truth transition graph, right: next states depend on current states and action. **(b)** Same demonstration for the full structured (lower-triangular adjacency matrix) transition graph. The hammer symbol denotes the action intervened on any of the observed states at each time step.

#### 4. Experiments demonstrate the effectiveness of DVAE-based hindsight encoder

**Environment setting.** We consider a straightforward yet non-trivial setup using the modulo environment: a chain-structured transition graph with  $d_S = 3$  and  $l = 4$ , with 3 factors: an observable  $o^1$ , a middle hidden state  $h^1$ , then an observable  $o^2$ . The environment includes a stationary discrete noise distribution defined as  $p(\epsilon_t^i = -1) = p(\epsilon_t^i = 1) = 0.05$  and  $p(\epsilon_t^i = 0) = 0.9$  for  $i = 1, 2, 3$ . The initial hidden state remains fixed across episodes. The principles outlined here can be extended to other graph structures and larger values of  $d_S$ , as empirically demonstrated later. Specifically, the transition dynamics in this setup are defined as  $o_{t+1}^1 := (o_t^1 + a_t^1 + \epsilon_t^1) \bmod 4$ ,  $h_{t+1}^1 := (o_t^1 + h_t^1 + \epsilon_t^2) \bmod 4$ , and  $o_{t+1}^2 := (h_t^1 + o_t^2 + a_t^3 + \epsilon_t^3) \bmod 4$ .

**Baselines and our DVAE encoders.** We compare the performance of 5 different hidden encoders, each learned end-to-end with the same transition model and reward predictor architecture: (i) history-based encoder, using complete past and current observations and actions (**History Enc.**):  $q_\phi(h_t | o_{1:t}, a_{1:t})$  parameterized by a forward RNN; (ii) Current and 1-step hindsight encoder (**Jarrett et al., 2023**), using current observations and action, and next step future observations (**Current & 1-Step Hindsight Enc.**):  $q_\phi(h_t | o_{t:t+1}, a_{t:t+1})$  parameterized by an MLP; (iii) Current and full hindsight encoder, using current and all future observations and actions (**Current & Full Hindsight Enc.**):  $q_\phi(h_t | o_{t:T}, a_{t:T})$  parameterized by a backward RNN; (iv) DVAE-based encoder with 1-step hindsight, using 1-step past (including sampled hidden), current, and 1-step future observations and actions (**DVAE 1-Step Hindsight Enc.**):  $q_\phi(h_t | h_{t-1}, o_{t-1:t+1}, a_{t-1:t+1})$ ; and (v) DVAE-based encoder with full hindsight, using 1-step past (including sampled hidden), current, and all future observations and actions (**DVAE Full Hindsight Enc.**):  $q_\phi(h_t | h_{t-1}, o_{t-1:T}, a_{t-1:T})$ .

**Implementation details.** We use the Adam optimizer with a step-decayed learning rate  $\alpha$ . The prediction horizon  $T$  is 5, and CMI threshold  $\delta$  is 0.03. Details on the neural network parameterization of the hidden encoder, transition model, and reward predictor are provided in Appendix A.4.

#### Results: DVAE Hindsight Encoders outperform History or Current & Hindsight Encoders.

In Fig. 4, we empirically compare the training performances and evaluated CMI matrices across 5 types of encoders under 2 settings with exogenous noise  $\epsilon_t$  applied to either the hidden state transition (noisy hidden setting) or the observed state transition (noisy observation setting).

In the noisy hidden setting (Fig. 4a and b), encoders with hindsight information converge to zero loss for both observed state predictions (the full NLL term of the VLB in Eq. (9)) and reward

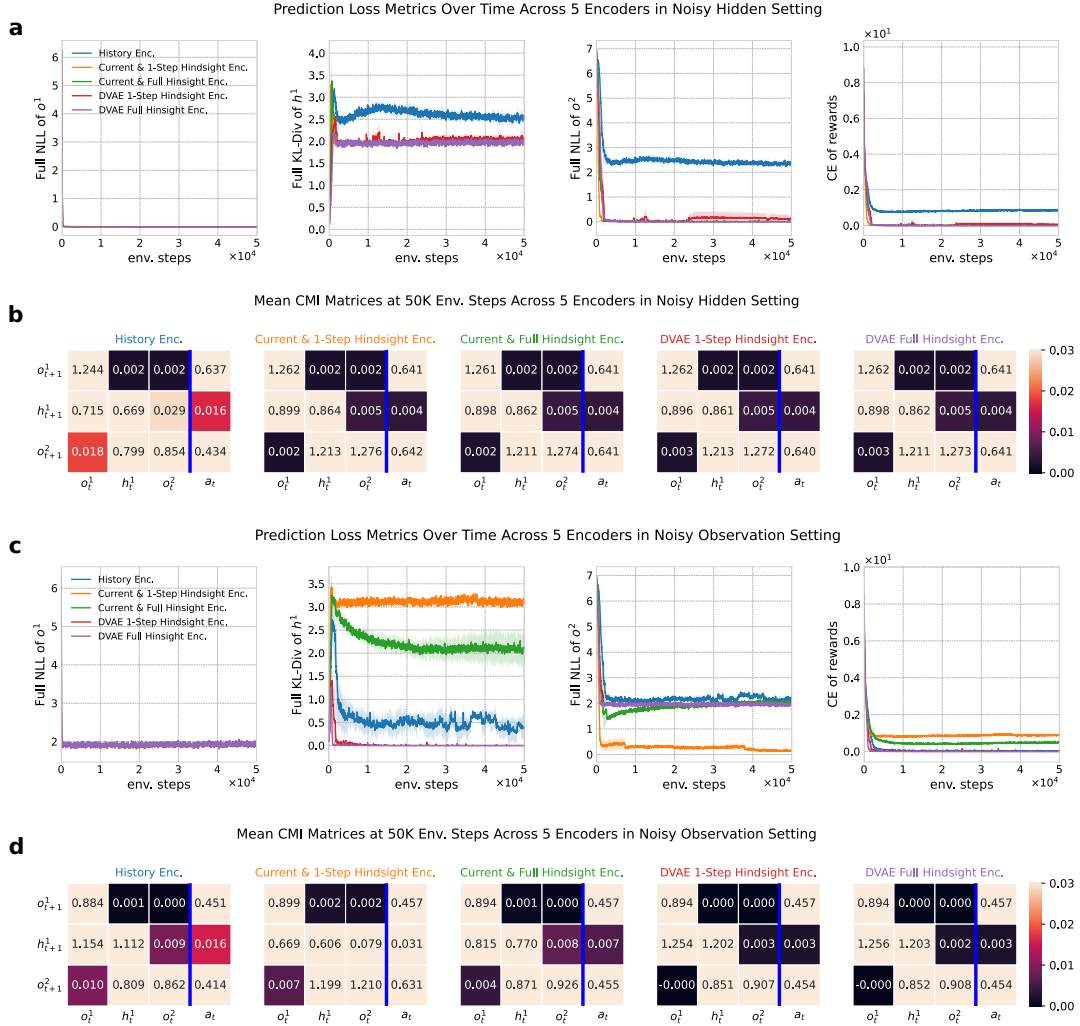


Figure 4: **(a)** Comparison of 5 encoder types, showing training profiles of state and reward prediction losses with mean and standard deviation (std) when noise  $\epsilon_t^h$  is applied to the hidden state transition. **(b)** Corresponding mean CMI matrices evaluated at the end of training across the 5 encoders, displayed as heatmaps under the same conditions. Similarly, **(c)** and **(d)** present training performance and evaluated CMI matrices, respectively, when noise  $\epsilon_t^h$  is applied to the observed state transitions. In all experiment results, each loss metric and CMI calculation for each encoder is run over 3 seeds. The color bar range is capped at the CMI threshold  $\delta$ , so that light color denotes an edge, and dark color no edge. The DVAE encoders produce CMI matrices whose binarized values match ground-truth.

prediction (the cross-entropy loss in Eq. (10)). These encoders also infer correct transition graphs after binarizing their evaluated CMI matrices using the threshold  $\delta$ . In contrast, the history-based encoder struggles to train effectively, resulting in a CMI matrix with values close to  $\delta$ , which reflects less statistical confidence in the existence of corresponding causal edges. Without access to the next observation  $o_{t+1}$ , the history-based encoder cannot deterministically infer the current hidden state  $h_t$ , given the unknown noise  $\epsilon_{t-1}$  affecting the transition to  $h_t$ . However, hindsight-based

		Evaluation Accuracy in Noisy Hidden / Observation Setting				
# Past	# Future	Graph	$h^1$ Decoding	$o^1, \sigma^2$ Prediction	$h^1$ Prediction	Reward Prediction
<b>History-Based Encoder</b>						
all	0	0.944 <sub>(0.039)</sub> / 1.000 <sub>(0.000)</sub>	0.865 <sub>(0.019)</sub> / 0.971 <sub>(0.021)</sub>	1.000 <sub>(0.000)</sub> , 0.872 <sub>(0.023)</sub> / 0.915 <sub>(0.017)</sub> , 0.876 <sub>(0.012)</sub>	0.850 <sub>(0.029)</sub> / 0.964 <sub>(0.026)</sub>	0.927 <sub>(0.020)</sub> / 0.997 <sub>(0.002)</sub>
<b>Current and Hindsight-Based Encoder</b>						
0	1	1.000 <sub>(0.000)</sub> / 0.944 <sub>(0.079)</sub>	1.000 <sub>(0.000)</sub> / 0.866 <sub>(0.044)</sub>	1.000 <sub>(0.000)</sub> , 1.000 <sub>(0.000)</sub> / 0.915 <sub>(0.017)</sub> , 0.996 <sub>(0.003)</sub>	0.899 <sub>(0.009)</sub> / 0.814 <sub>(0.009)</sub>	1.000 <sub>(0.000)</sub> / 0.949 <sub>(0.005)</sub>
0	all	1.000 <sub>(0.000)</sub> / 1.000 <sub>(0.000)</sub>	1.000 <sub>(0.000)</sub> / 0.914 <sub>(0.017)</sub>	1.000 <sub>(0.000)</sub> , 1.000 <sub>(0.000)</sub> / 0.915 <sub>(0.017)</sub> , 0.897 <sub>(0.007)</sub>	0.899 <sub>(0.009)</sub> / 0.870 <sub>(0.031)</sub>	1.000 <sub>(0.000)</sub> / 0.959 <sub>(0.009)</sub>
<b>DVAE-Based Hindsight Encoder</b>						
all	1	1.000 <sub>(0.000)</sub> / 1.000 <sub>(0.000)</sub>	1.000 <sub>(0.000)</sub> / 1.000 <sub>(0.000)</sub>	1.000 <sub>(0.000)</sub> , 1.000 <sub>(0.000)</sub> / 0.915 <sub>(0.017)</sub> , 0.905 <sub>(0.009)</sub>	0.895 <sub>(0.009)</sub> / 1.000 <sub>(0.000)</sub>	1.000 <sub>(0.000)</sub> / 1.000 <sub>(0.000)</sub>
all	all	1.000 <sub>(0.000)</sub> / 1.000 <sub>(0.000)</sub>	1.000 <sub>(0.000)</sub> / 1.000 <sub>(0.000)</sub>	1.000 <sub>(0.000)</sub> , 1.000 <sub>(0.000)</sub> / 0.915 <sub>(0.017)</sub> , 0.905 <sub>(0.009)</sub>	0.899 <sub>(0.009)</sub> / 1.000 <sub>(0.000)</sub>	1.000 <sub>(0.000)</sub> / 1.000 <sub>(0.000)</sub>

Table 1: Evaluation accuracies across various metrics, including transition graph accuracy (measured by the match between inferred and ground truth edges), hidden state decoding accuracy (linear decoding accuracy of encoded hidden states to ground truth hidden states), observation prediction accuracy, hidden state prediction accuracy (measured by the match between predicted and encoded next hidden states), and reward prediction accuracy. These metrics are reported for 5 types of encoders utilizing different steps of past and future observables in both noisy hidden and noisy observation settings. Each accuracy value is presented as  $\text{mean}_{\text{std}}$  over 3 runs. Lavender and beige highlights indicate suboptimal accuracy values for certain encoders in the noisy hidden and observation settings, respectively. Note that the DVAE-based encoder is labeled as using all past observables, as it estimates the 1-step past hidden state based on recursive hidden samples from the beginning of an episode, which requires all past observables.

approaches can learn  $h_t$  by utilizing information from observed states, which serve as children of the hidden state in the transition graph, thereby enabling accurate learning of the transition graph. Due to the unobserved exogenous noise injected into the hidden state transition, the transition model can only predict the next hidden state in distribution. As a result, prediction losses for the hidden state (measured by the full KL divergence between the encoded and predicted next hidden states in Eq. (9)) do not decrease to zero for all encoders.

In the noisy observation setting (Fig. 4c and d), the DVAE-based encoder successfully learns hidden representations, allowing it to accurately predict the next hidden states and rewards, while the history-based encoder and current hindsight-based encoder fail to achieve similar performance (as seen in the second and fourth panels of Fig. 4c). In the third panel of Fig. 4c), it appears that Current and Hindsight Encoders achieve lower loss, but this is due to learning to copy  $o_{t+1}^2$  to the hidden, as described for Table 1. Encoders other than DVAE-based encoders produce CMI matrices with values closer to threshold, or even infer spurious edges. We hypothesize that the DVAE-based model’s ability to identify the current hidden state  $h_t$  stems from its recursive structure, which combines sample-based past (Markovian) information with future information. In contrast, other encoders, which rely on a single directional view of observables along the trajectory, lack sufficient information to identify the current hidden state in the noisy observation setting. Similar to the prediction of next state hidden in the noisy hidden setting, the transition model can only predict noisy observations in distribution.

We also tabulate the accuracy of graph edges, decoding of encoded hidden, and state transitions, after convergence, of the five encoder architectures across both noise settings, in Table 1. In the noisy hidden setting, the lower accuracies of the history-based encoder, highlighted in lavender, indicate its inability to learn the hidden state and accurately perform the corresponding transition and reward predictions. Ideally, the encoded hidden state should be linearly decodable to its ground truth value and deterministically predictive of the reward, as reflected by perfect  $h^1$  decoding and reward prediction accuracy in all other encoders. Additionally, the expected  $h^1$  prediction accuracy

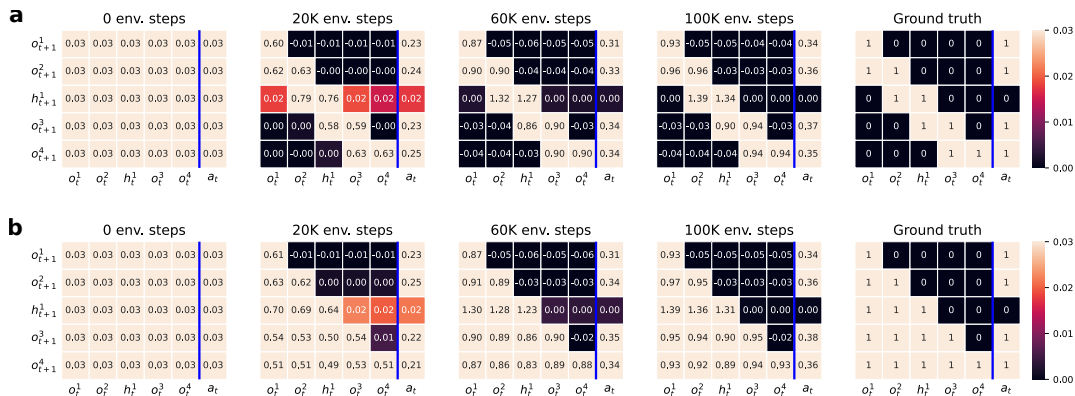


Figure 5: Evolution of CMI matrices for the (a) chain and (b) full graph structures. The ground truth graphs are shown on the far right.

should be approximately 0.9, accounting for the 10% noise in the hidden transition, assuming both the encoded and predicted hidden states are optimally learned. Indeed, the mean  $h^1$  prediction accuracy for all encoders, except the history-based one, is very close to 0.9. Similarly, in the noisy observation setting, the accuracies highlighted in beige indicate suboptimal encoding and prediction of the hidden states for the history-based and current hindsight-based encoders. Interestingly, for the current and hindsight-based encoder, the mean  $o^2$  prediction accuracy exceeds the expected value of 0.9 and approaches 1.0 (see also third panel of Fig. 4c), suggesting that this encoder copies its input of the next noisy  $o^2$  as the hidden state. This copying approach, however, trades off accuracy in  $h^1$  and reward prediction compared to encoders that do not learn this inconsequential solution for the hidden state. The DVAE-based encoders perform optimally in both noise settings. Notably, the DVAE 1-step Hindsight Encoder achieves the same optimal performance as the theoretically-derived DVAE Full Hindsight Encoder due to absence of cascaded hidden factors in our environment.

Finally, we evaluate our DVAE Full Hindsight Encoder on chain and full structured (lower-triangular adjacency matrix) transition graphs with  $d_S = 5$ , while keeping the rest of the environment setup unchanged. Fig. 5 shows the evolution of the CMI matrix during training in the setting of noisy observations for both transition graph structures. The CMI matrix initially has all elements set to the predefined threshold  $\delta$  and gradually decreases for unconnected factor pairs in the transition while increasing for connected factor pairs. The final binary matrix, obtained by applying the threshold to binarize the CMI matrix, converges to the ground truth adjacency matrix.

## 5. Discussion

We have demonstrated that the proposed DVAE-based hindsight encoder effectively identifies hidden state factors and learns the causal transition graph in a factored-POMDP, outperforming both history-based and typical hindsight-based encoders. This approach shows particular promise in settings with access to full offline trajectories. In biological scenarios, our technique is reminiscent of “trajectory replay” in rodent planning, where neural patterns associated with past experiences are replayed in both forward and reverse directions (Ólafsdóttir et al., 2018). Thus, our method holds value for applications where offline trajectories can be leveraged. In online settings, a causal model initially trained on offline trajectories could support more accurate model rollouts within frameworks like Dyna (Sutton, 1991; Sutton et al., 2012; Peng et al., 2018) or Model Predictive Control

(MPC) (Chua et al., 2018; Wang et al., 2019; Moerland et al., 2023), offering an advantage over models trained solely with history-based or other hindsight-based approaches.

In our formulation, we identified deterministic hidden components of factored state transitions, and, using the Reparametrization Lemma, isolated stochastic effects as unobserved exogenous noise per factor. Future work could refine our framework by also inferring the exogenous noise at each time step through dedicated noise encoders, following the identification of deterministic hidden factors. While our DVAE 1-step Hindsight Encoder was sufficient for a single hidden factor, extending it to scenarios with multiple cascaded hidden factors, with only the last hidden factor influencing an observable factor, may require additional future information for effective latent identification. Moreover, expanding this approach to continuous state-action spaces would link our work to DVAE research on latent dynamics in stochastic driven dynamical systems (Girin et al., 2020). Addressing these areas would support further scaling and generalization of the framework.

## Acknowledgments

C. Han, A. Gilra and E. Vasilaki acknowledge the CHIST-ERA grant for the “Causal Explanations in Reinforcement Learning (CausalXRL)” project (CHIST-ERA-19-XAI-002), by the Engineering and Physical Sciences Research Council (EPSRC), United Kingdom (grant reference EP/V055720/1) for supporting the work. D. Basu acknowledges the CHIST-ERA grant for the CausalXRL project (CHIST-ERA-19-XAI-002) by L’Agence Nationale de la Recherche, France (grant reference ANR21-CHR4-0007), as well as the ANR JCJC for the REPUBLIC project (ANR-22-CE23-0003-01), and the PEPR project FOUNDRY (ANR23-PEIA-0003) for supporting the work. C. Han and E. Vasilaki acknowledge the grant for the “Magnetic Architectures for Reservoir Computing Hardware (MARCH)” project, by the EPSRC, United Kingdom (grant reference EP/V006339/1) for supporting the work. C. Han, M. Mangan and E. Vasilaki acknowledge the grant for the “Active learning and selective attention for robust, transparent and efficient AI (ActiveAI)” project, by the EPSRC, United Kingdom (grant reference EP/S030964/1) for supporting the work.

## References

- Andrea Baisero and Christopher Amato. Learning complementary representations of the past using auxiliary tasks in partially observable reinforcement learning. In *AAMAS*, pages 1762–1764, 2020.
- Lars Buesing, Theophane Weber, Yori Zwols, Nicolas Heess, Sebastien Racaniere, Arthur Guez, and Jean-Baptiste Lespiau. Woulda, Coulda, Shoulda: Counterfactually-Guided Policy Search. September 2018. URL <https://openreview.net/forum?id=BJG0voC9YQ>.
- Kurtland Chua, Roberto Calandra, Rowan McAllister, and Sergey Levine. Deep reinforcement learning in a handful of trials using probabilistic dynamics models. *Advances in neural information processing systems*, 31, 2018.
- Raj Ghugare, Homanga Bharadhwaj, Benjamin Eysenbach, Sergey Levine, and Ruslan Salakhutdinov. Simplifying model-based rl: learning representations, latent-space models, and policies with one objective. *arXiv preprint arXiv:2209.08466*, 2022.

- Laurent Girin, Simon Leglaive, Xiaoyu Bie, Julien Diard, Thomas Hueber, and Xavier Alameda-Pineda. Dynamical variational autoencoders: A comprehensive review. *arXiv preprint arXiv:2008.12595*, 2020.
- Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144*, 2016.
- Daniel Jarrett, Corentin Tallec, Florent Altché, Thomas Mesnard, Rémi Munos, and Michal Valko. Curiosity in hindsight: intrinsic exploration in stochastic environments. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *ICML'23*, pages 14780–14816, Honolulu, Hawaii, USA, July 2023. JMLR.org. URL <https://icml.cc/virtual/2023/poster/24131>.
- Nan Rosemary Ke, Aniket Didolkar, Sarthak Mittal, Anirudh Goyal, Guillaume Lajoie, Stefan Bauer, Danilo Rezende, Yoshua Bengio, Michael Mozer, and Christopher Pal. Systematic evaluation of causal discovery in visual model based reinforcement learning. *arXiv preprint arXiv:2107.00848*, 2021.
- Diederik P. Kingma and Max Welling. An introduction to variational autoencoders. *Foundations and Trends® in Machine Learning*, 12(4):307–392, 2019. ISSN 1935-8245. doi: 10.1561/22000000056. URL <http://dx.doi.org/10.1561/22000000056>.
- Michael Littman and Richard S Sutton. Predictive representations of state. *Advances in neural information processing systems*, 14, 2001.
- Chris J Maddison, Andriy Mnih, and Yee Whye Teh. The concrete distribution: A continuous relaxation of discrete random variables. *arXiv preprint arXiv:1611.00712*, 2016.
- Thomas M Moerland, Joost Broekens, Aske Plaat, Catholijn M Jonker, et al. Model-based reinforcement learning: A survey. *Foundations and Trends® in Machine Learning*, 16(1):1–118, 2023.
- Tianwei Ni, Benjamin Eysenbach, Erfan Seyedsalehi, Michel Ma, Clement Gehring, Aditya Mahajan, and Pierre-Luc Bacon. Bridging State and History Representations: Understanding Self-Predictive RL, April 2024. URL <http://arxiv.org/abs/2401.08898>. arXiv:2401.08898 [cs].
- H Freyja Ólafsdóttir, Daniel Bush, and Caswell Barry. The role of hippocampal replay in memory and planning. *Current Biology*, 28(1):R37–R50, 2018.
- Frans Oliehoek, Stefan Witwicki, and Leslie Kaelbling. A sufficient statistic for influence in structured multiagent environments. *Journal of Artificial Intelligence Research*, 70:789–870, 2021.
- Baolin Peng, Xiujun Li, Jianfeng Gao, Jingjing Liu, Kam-Fai Wong, and Shang-Yu Su. Deep dyna-q: Integrating planning for task-completion dialogue policy learning. *arXiv preprint arXiv:1801.06176*, 2018.
- Jayakumar Subramanian, Amit Sinha, Raihan Seraj, and Aditya Mahajan. Approximate Information State for Approximate Planning and Reinforcement Learning in Partially Observed Systems.

- Journal of Machine Learning Research*, 23:1–83, 2022. URL <https://www.jmlr.org/papers/v23/20-1165.html>.
- Richard S Sutton. Dyna, an integrated architecture for learning, planning, and reacting. *ACM Sigart Bulletin*, 2(4):160–163, 1991.
- Richard S Sutton, Csaba Szepesvári, Alborz Geramifard, and Michael P Bowling. Dyna-style planning with linear function approximation and prioritized sweeping. *arXiv preprint arXiv:1206.3285*, 2012.
- Tingwu Wang, Xuchan Bao, Ignasi Clavera, Jerrick Hoang, Yeming Wen, Eric Langlois, Shunshi Zhang, Guodong Zhang, Pieter Abbeel, and Jimmy Ba. Benchmarking model-based reinforcement learning. *arXiv preprint arXiv:1907.02057*, 2019.
- Zizhao Wang, Xuesu Xiao, Zifan Xu, Yuke Zhu, and Peter Stone. Causal dynamics learning for task-independent state abstraction. *arXiv preprint arXiv:2206.13452*, 2022.
- Amy Zhang, Rowan McAllister, Roberto Calandra, Yarín Gal, and Sergey Levine. Learning invariant representations for reinforcement learning without reconstruction. *arXiv preprint arXiv:2006.10742*, 2020.
- K. J Åström. Optimal control of Markov processes with incomplete state information. *Journal of Mathematical Analysis and Applications*, 10(1):174–205, February 1965. ISSN 0022-247X. doi: 10.1016/0022-247X(65)90154-X. URL <https://www.sciencedirect.com/science/article/pii/0022247X6590154X>.

## Appendix A. DVAE for Factored-POMDP

### A.1. Log-likelihood decomposition

The detailed derivation from Eq. (2) to Eq. (3) is provided as follows:

$$\begin{aligned} & \mathbb{E}_{p(o_{1:T}|a_{1:T})} [\log p_\theta(o_{1:T}|a_{1:T})] \\ &= \mathbb{E}_{p(o_{1:T}|a_{1:T})} \left[ \mathbb{E}_{q_\phi(h_{1:T}|o_{1:T}, a_{1:T})} [\log p_\theta(o_{1:T}|a_{1:T})] \right] \end{aligned} \quad (14)$$

$$= \mathbb{E}_{p(o_{1:T}|a_{1:T})} \left[ \mathbb{E}_{q_\phi(h_{1:T}|o_{1:T}, a_{1:T})} \left[ \log \frac{p_\theta(o_{1:T}, h_{1:T}|a_{1:T})}{p_\theta(h_{1:T}|o_{1:T}, a_{1:T})} \right] \right] \quad (15)$$

$$= \mathbb{E}_{p(o_{1:T}|a_{1:T})} \left[ \mathbb{E}_{q_\phi(h_{1:T}|o_{1:T}, a_{1:T})} \left[ \log \frac{p_\theta(o_{1:T}, h_{1:T}|a_{1:T}) q_\phi(h_{1:T}|o_{1:T}, a_{1:T})}{q_\phi(h_{1:T}|o_{1:T}, a_{1:T}) p_\theta(h_{1:T}|o_{1:T}, a_{1:T})} \right] \right] \quad (16)$$

$$\begin{aligned} &= \mathbb{E}_{p(o_{1:T}|a_{1:T})} \left[ \mathbb{E}_{q_\phi(h_{1:T}|o_{1:T}, a_{1:T})} \left[ \log \frac{p_\theta(o_{1:T}, h_{1:T}|a_{1:T})}{q_\phi(h_{1:T}|o_{1:T}, a_{1:T})} \right] \right. \\ & \quad \left. + \mathbb{E}_{q_\phi(h_{1:T}|o_{1:T}, a_{1:T})} \left[ \log \frac{q_\phi(h_{1:T}|o_{1:T}, a_{1:T})}{p_\theta(h_{1:T}|o_{1:T}, a_{1:T})} \right] \right] \end{aligned} \quad (17)$$

$$\begin{aligned} &= \mathbb{E}_{p(o_{1:T}|a_{1:T})} \left[ \underbrace{\mathbb{E}_{q_\phi(h_{1:T}|o_{1:T}, a_{1:T})} [\log p_\theta(o_{1:T}, h_{1:T}|a_{1:T}) - \log q_\phi(h_{1:T}|o_{1:T}, a_{1:T})]}_{\ell_{\text{VLB}}(\theta, \phi; o_{1:T}, a_{1:T})} \right] \\ & \quad + D_{\text{KL}}(q_\phi(h_{1:T}|o_{1:T}, a_{1:T}) \parallel p_\theta(h_{1:T}|o_{1:T}, a_{1:T})) \end{aligned} \quad (18)$$

### A.2. Variational Lower Bound (VLB) of DVAE-based framework for learning latent hidden and transition dynamics of factored-POMDP

**Generative model (transition model).** The generative model for the entire state sequence in the Eq. (4) can be factorized as:

$$\begin{aligned} p_\theta(o_{1:T}, h_{1:T}|a_{1:T}) &= \prod_{t=0}^{T-1} p_\theta(o_{t+1}, h_{t+1}|o_{1:t}, h_{1:t}, a_{1:T}) \\ &= \prod_{t=0}^{T-1} p_{\theta_o}(o_{t+1}|o_{1:t}, h_{1:t+1}, a_{1:T}) p_{\theta_h}(h_{t+1}|o_{1:t}, h_{1:t}, a_{1:T}) \\ &= \prod_{t=0}^{T-1} p_{\theta_o}(o_{t+1}|o_t, h_t, a_t) p_{\theta_h}(h_{t+1}|o_t, h_t, a_t) \end{aligned} \quad (19)$$

where each term in the product is simplified using d-separation in the unrolled transition graph from  $t = 1$  to  $T$  (see Fig. 1c). Here,  $\theta = \theta_o \cup \theta_h$  represents the parameters of the generative model. Note that the observation likelihood  $p_{\theta_o}(o_{t+1}|o_t, h_t, a_t)$  and the hidden prior  $p_{\theta_h}(h_{t+1}|o_t, h_t, a_t)$  in the generative model corresponds to the transition models of the observed and hidden states, respectively.

**Inference model (hidden encoder).** Similarly, we factorize the posterior distribution of the generative model as follows:

$$p_\theta(h_{1:T}|o_{1:T}, a_{1:T}) = \prod_{t=0}^{T-1} p_\theta(h_{t+1}|h_{1:t}, o_{1:T}, a_{1:T})$$

$$= \prod_{t=0}^{T-1} p_{\theta}(h_{t+1}|h_t, o_{t:T}, a_{t:T}) \quad (20)$$

and consider that the inference model, parameterized by  $\phi$ , captures the exact factorized structure of the posterior distribution in Eq. (20).

$$q_{\phi}(h_{1:T}|o_{1:T}, a_{1:T}) = \prod_{t=0}^{T-1} q_{\phi}(h_{t+1}|h_t, o_{t:T}, a_{t:T}) \quad (21)$$

Specifically, the hidden encoder  $q_{\phi}(h_t|h_{t-1}, o_{t-1:T}, a_{t-1:T})$  combines information from the Markovian past, through  $h_{t-1}$ ,  $o_{t-1}$  and  $a_{t-1}$ , with information from the present and future observations  $o_{t:T}$  and actions  $a_{t:T}$  to encode the current hidden state  $h_t$ . It is important to note that we assume Markovianity for the forward transitions but not for the backward transitions. Consequently, the hidden encoder depends on all future information, rather than just the immediate next-step information used in the hindsight-based encoder by Jarrett et al. (2023).

**Variational Lower Bound.** By substituting the decomposed forms of both the generative model from Eq. (19) and the inference model from Eq. (21) into the general form of VLB defined in Eq. (4), we obtain:

$$\begin{aligned} \ell_{\text{VLB}}(\theta, \phi; o_{1:T}, a_{1:T}) &= \sum_{t=0}^{T-1} \mathbb{E}_{q_{\phi}(h_{1:t}|o_{1:T}, a_{1:T})} [\log p_{\theta_o}(o_{t+1}|o_t, h_t, a_t) \\ &\quad - D_{\text{KL}}(q_{\phi}(h_{t+1}|h_t, o_{t:T}, a_{t:T}) \parallel p_{\theta_h}(h_{t+1}|o_t, h_t, a_t))] \end{aligned} \quad (22)$$

By using the factorization in Eq. (21), the expectation in the above VLB can be expressed as a cascade of expectations over conditional distributions of individual hidden states at different time steps:

$$\begin{aligned} \mathbb{E}_{q_{\phi}(h_{1:t}|o_{1:T}, a_{2:T})} [f(h_t)] &= \mathbb{E}_{q_{\phi}(h_1|o_{1:T}, a_{1:T})} [ \\ &\quad \mathbb{E}_{q_{\phi}(h_2|h_1, o_{1:T}, a_{1:T})} [ \\ &\quad \mathbb{E}_{q_{\phi}(h_3|h_2, o_{2:T}, a_{2:T})} [\dots \\ &\quad \mathbb{E}_{q_{\phi}(h_t|h_{t-1}, o_{t-1:T}, a_{t-1:T})} [f(h_t)] \dots]]] \end{aligned} \quad (23)$$

Here,  $f(h_t)$  represents an arbitrary function of  $h_t$ . Each intractable expectation in this sequence can be approximated using a Monte Carlo estimate. This involves iteratively sampling from  $q_{\phi}(h_{\tau}|h_{\tau-1}, o_{\tau-1:T}, a_{\tau-1:T})$  for  $\tau = 1$  to  $t$ , employing the same reparameterization trick used in standard VAEs (Maddison et al., 2016; Jang et al., 2016; Kingma and Welling, 2019). Additionally, the VLB in Eq. (22), which is defined for a single data sequence, can be extended by averaging over a mini-batch of training data sequences, thereby approximating the expected VLB with respect to the true data distribution.

Furthermore, by expressing  $o_t = (o_t^1, \dots, o_t^{d_o})$ ,  $h_t = (h_t^1, \dots, o_t^{d_h})$  and  $s_t = (o_t, h_t)$  and using the factorized forms of both the transition models and the hidden encoder, we have:

$$p_{\theta_o}(o_{t+1}|o_t, h_t, a_t) = \prod_{j=1}^{d_o} p_{\theta_h}(o_{t+1}^j|s_t, a_t), \quad (24)$$

$$p_{\theta_h}(h_{t+1}|o_t, h_t, a_t) = \prod_{j=1}^{d_H} p_{\theta_h}(h_{t+1}^j | s_t, a_t), \quad (25)$$

$$q_{\phi}(h_{t+1}|h_t, o_{t:T}, a_{t:T}) = \prod_{j=1}^{d_H} q_{\phi}(h_{t+1}^j | h_t, o_{t:T}, a_{t:T}) \quad (26)$$

Eq. (8) is obtained by substituting the above expressions into Eq. (22).

### A.3. Conditional mutual information

Starting from the definition of conditional mutual information, we have:

$$I(s_t^i; s_{t+1}^j | s_t \setminus s_t^i, a_t) = \mathbb{E}_{p(s_t, a_t, s_{t+1}^j)} \left[ \log \frac{p(s_t^i, s_{t+1}^j | s_t \setminus s_t^i, a_t)}{p(s_t^i | s_t \setminus s_t^i, a_t) p(s_{t+1}^j | s_t \setminus s_t^i, a_t)} \right] \quad (27)$$

$$= \mathbb{E}_{p(s_t, a_t, s_{t+1}^j)} \left[ \log \frac{p(s_{t+1}^j | s_t, a_t) p(s_t^i | s_t \setminus s_t^i, a_t)}{p(s_t^i | s_t \setminus s_t^i, a_t) p(s_{t+1}^j | s_t \setminus s_t^i, a_t)} \right] \quad (28)$$

$$= \mathbb{E}_{p(s_t, a_t, s_{t+1}^j)} \left[ \log \frac{p(s_{t+1}^j | s_t, a_t)}{p(s_{t+1}^j | s_t \setminus s_t^i, a_t)} \right] \quad (29)$$

$$= \mathbb{E}_{p(s_t, a_t)} \left[ \mathbb{E}_{p(s_{t+1}^j | s_t, a_t)} \left[ \log \frac{p(s_{t+1}^j | s_t, a_t)}{p(s_{t+1}^j | s_t \setminus s_t^i, a_t)} \right] \right] \quad (30)$$

$$= \mathbb{E}_{p(s_t, a_t)} \left[ D_{\text{KL}}(p(h_{t+1}^j | s_t, a_t) \parallel p(h_{t+1}^j | s_t \setminus s_t^i, a_t)) \right] \quad (31)$$

where Eqs. (29) and (31) correspond to Eqs. (12) and (13), respectively.

### A.4. Neural Network-Based Parameterization

The hidden encoder  $q_{\phi}(h_t | h_{t-1}, o_{t-1:T}, a_{t-1:T})$  is implemented using a backward RNN to capture current and future dependencies, and an MLP to model Markovian past dependencies. A combiner function (CF) is then employed to merge the outputs of the MLP and the RNN (its internal state) to produce parameters (e.g., logits) of the distribution of the current hidden state:

$$\overleftarrow{g}_t = \text{RNN}_{\phi_{\overleftarrow{g}}}(\overleftarrow{g}_{t+1}, [o_t, a_t]), \quad (32)$$

$$e_t = \text{MLP}_{\phi_e}(h_{t-1}, o_{t-1}, a_{t-1}), \quad (33)$$

$$f_t = \text{CF}_{\phi_f}(e_t, \overleftarrow{g}_t), \quad (34)$$

$$q_{\phi}(h_t | h_{t-1}, o_{t-1:T}, a_{t-1:T}) = \text{Dist}(h_t; f_t) \quad (35)$$

where  $\text{CF}_{\phi_f}$  is a feedforward combining network parameterized by  $\phi_f$ . Thus, the parameters of the hidden encoder are  $\phi = \phi_{\overleftarrow{g}} \cup \phi_e \cup \phi_f$ .

The transition model for the observed states  $p_{\theta_o}(o_{t+1} | o_t, h_t, a_t)$  and the hidden states  $p_{\theta_h}(h_{t+1} | o_t, h_t, a_t)$  are implemented using factor-wise masked MLPs (MMLPs) following Wang et al. (2022):

$$m_t = \text{MMLP}_{\theta_o}(o_t, h_t, a_t), \quad (36)$$

$$p_{\theta_o}(o_{t+1} | o_t, h_t, a_t) = \text{Dist}(o_t; m_t), \quad (37)$$

$$n_t = \text{MMLP}_{\theta_h}(o_t, h_t, a_t), \quad (38)$$

$$p_{\theta_h}(h_{t+1}|o_t, h_t, a_t) = \text{Dist}(h_t; n_t) \quad (39)$$

where  $m_t$  and  $n_t$  are the outputs of the masked MLPs parameterized by  $\theta_o$  and  $\theta_h$ , respectively. The distributions  $\text{Dist}(o_t; m_t)$  and  $\text{Dist}(h_t; n_t)$  represent the probability distributions of  $o_{t+1}$  and  $h_{t+1}$  parameterized by  $m_t$  and  $n_t$ .

The architecture of the DVAE model is illustrated in Fig. 6.

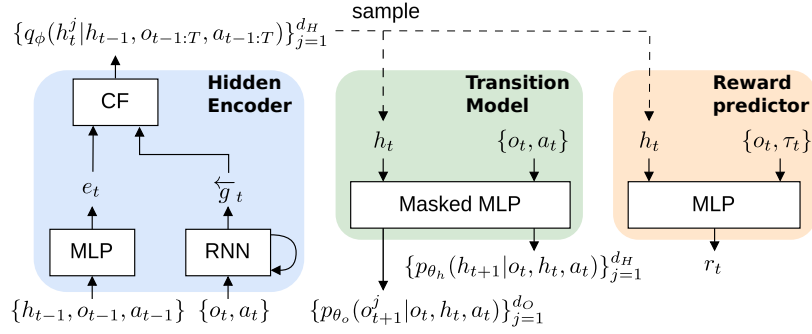


Figure 6: Model architecture illustrating the computational graph for encoding, sampling and prediction processes.