
Towards Principled Representation Learning from Videos for Reinforcement Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We study pre-training representations for decision-making using video data, which
2 is abundantly available for tasks such as game agents and software testing. Even
3 though significant empirical advances have been made on this problem, a theoretical
4 understanding remains absent. We initiate the theoretical investigation into
5 principled approaches for representation learning and focus on learning the latent
6 state representations of the underlying MDP using video data. We study two types
7 of settings: one where there is iid noise in the observation, and a more challenging
8 setting where there is also the presence of exogenous noise, which is non-iid noise
9 that is temporally correlated, such as the motion of people or cars in the background.
10 We study three commonly used approaches: autoencoding, temporal contrastive
11 learning, and forward modeling. We prove upper bounds for temporal contrastive
12 learning and forward modeling in the presence of only iid noise. We show that
13 these approaches can learn the latent state and use it to do efficient downstream RL
14 with polynomial sample complexity. When exogenous noise is also present, we
15 establish a lower bound result showing that the sample complexity of learning from
16 video data can be exponentially worse than learning from action-labeled trajectory
17 data. This partially explains why reinforcement learning with video pre-training is
18 hard. We evaluate these representational learning methods in three visual domains,
19 yielding results that are consistent with our theoretical findings.

20 1 Introduction

21 Representations pre-trained on large amounts of offline data have led to significant advances in
22 machine learning domains such as natural language processing [Liu et al., 2019, Brown et al., 2020]
23 and multi-modal learning [Lin et al., 2021, Radford et al., 2021]. This has naturally prompted a
24 similar undertaking in reinforcement learning (RL) with the goal of training a representation model
25 that can be used in a policy to solve a downstream RL task. The natural choice of data for RL
26 problems is trajectory data, which contains the agent’s observation along with actions taken by
27 the agent and the rewards received by it [Sutton and Barto, 2018]. A line of work has proposed
28 approaches for learning representations with trajectory data in both offline [Uehara et al., 2021, Islam
29 et al., 2022] and online learning settings [Nachum et al., 2018, Bharadhwaj et al., 2022]. However,
30 unlike text and image data, which are abundant on the internet or naturally generated by users,
31 trajectory data is comparatively limited and expensive to collect. In contrast, video data, which
32 only contains a sequence of observations (without any action or reward labeling), is often plentiful,
33 especially for domains such as gaming and software. This motivates a line of work considering
34 learning representations for RL using video data [Zhao et al., 2022]. *But is there a principled
35 foundation underlying these approaches? Are representations learned from video data as useful
36 as representations learned from trajectory data?* We initiate a theoretical understanding of these

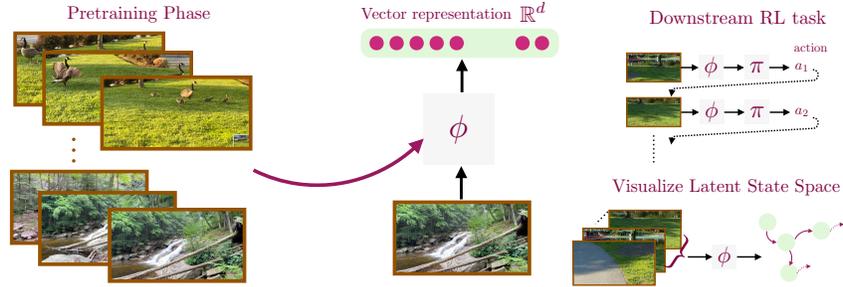


Figure 1: A flowchart of our video pre-training phase. **Left:** We assume access to a large set of videos (or, unlabeled episodes). **Center:** A representation learning method is used to train a model ϕ which maps an observation to a vector representation. **Right:** This representation can be used in a downstream task to do reinforcement learning or visualize the latent world state.

37 approaches to show when and how these approaches yield representations that can be used to solve a
 38 downstream RL task efficiently.

39 Consider a representation learning pipeline shown in Figure 1. We are provided videos, or equivalently
 40 a sequence of observations, from agents navigating in the world. We make no assumption about the
 41 behavior of the agent in the video data. They can be trying to solve one task, many different tasks, or
 42 none at all. This video data is used to learn a model ϕ that maps any given observation to a vector
 43 representation. This representation is subsequently used to perform downstream RL — defining a
 44 policy on top of the learned representation and only training the policy for the downstream task. We
 45 can also use this representation to define a dynamics model or a critique model. The representation
 46 can also help visualize the agent state space or dynamics for the purpose of debugging.

47 A suitable representation for performing RL efficiently is aligned with the underlying dynamics of
 48 the world. Ideally, the representation captures the latent agent state, which contains information about
 49 the world relevant to decision-making while ignoring any noise in the observation. For example,
 50 in Figure 1, ignoring noise such as the motion of geese in the background is desirable if the task
 51 involves walking on the pavement. We distinguish between two types of noise: (1) temporally
 52 independent noise that occurs at each time step independent of the history, (2) temporally dependent
 53 noise, or exogenous noise, that can evolve temporally but in a manner independent of the agent’s
 54 actions (such as the motion of geese in Figure 1).

55 A range of approaches have been developed that provably recover the latent agent state from observa-
 56 tions using trajectory data [Misra et al., 2020, Efroni et al., 2022] which contains actions. However,
 57 for many domains there is relatively little trajectory data that exists naturally, making it expensive
 58 to scale these learning approaches. In contrast, video data is more naturally available but these
 59 prior provable approaches do not work with video data. On the other hand, it is unknown whether
 60 approaches that empirically work with video data provably recover the latent representation and lead
 61 to efficient RL. Motivated by this, we build a theoretical understanding of three such video-based
 62 representation learning approaches: *autoencoder* which trains representations by reconstructing
 63 observations, *forward modeling* which predicts future observations, and *temporal contrastive* learning
 64 which trains a representation to determine if a pair of observations are causally related or not.

65 Our first theoretical result shows that in the absence of exogenous noise, forward modeling and
 66 temporal contrastive learning approaches both provably work. Further, they lead to efficient down-
 67 stream RL that is strictly more sample-efficient than solving these tasks without any pre-training.
 68 Our second theoretical result establishes a lower bound showing that in the presence of exogenous
 69 noise, any compact and frozen representation that is pre-trained using video data cannot be used to do
 70 efficient downstream RL. In contrast, if the trajectory data was available, efficient pre-training would
 71 be possible. This establishes a statistical gap showing that video-based representation pre-training
 72 can be exponentially harder than trajectory-based representation pre-training.

73 We empirically test our theoretical results in three visual domains: GridWorld (a navigation domain),
 74 ViZDoom basic (a first-person 3D shooting game), and ViZDoom Defend The Center (a more
 75 challenging first-person 3D shooting game). We evaluate the aforementioned approaches along with
 76 ACRO [Islam et al., 2022], a representation pre-trained using trajectory data and designed to filter out
 77 exogenous noise. We observe that in accordance with our theory, both forward modeling and temporal

78 contrastive learning succeed at RL when there is no exogenous noise. However, in the presence
 79 of exogenous noise, their performance degrades. Specifically, we find that temporal contrastive
 80 learning is especially prone to fail in the presence of exogenous noise, as it can rely exclusively
 81 on such noise to optimally minimize the contrastive loss. While we find that forward modeling is
 82 somewhat robust to exogenous noise, however, as exogenous noise increases, its performance quickly
 83 degrades as well. While any finite-sample guarantees for the autoencoding method remain an open
 84 question, empirically, we find that the performance of autoencoder-based representation learning is
 85 unpredictable. On the other hand, ACRO continues to perform well, highlighting a disadvantage of
 86 video pre-training. The code for all experiments will be made available at <url-redacted>.

87 2 Representation Learning for RL using Video Dataset

88 We assume access to a dataset \mathcal{D} of n videos $\mathcal{D} = \{(x_1^{(i)}, x_2^{(i)}, \dots, x_H^{(i)})\}_{i=1}^n$ where $x_j^{(i)}$ is the j^{th}
 89 observation (or frame) of the i^{th} video. We are provided a decoder class $\Phi = \{\phi : \mathcal{X} \rightarrow [N]\}$,
 90 and our goal is to learn a decoder $\phi \in \Phi$ that captures task-relevant information in the underlying
 91 state $\phi^*(x)$ while throwing away as much exogenous noise as possible. Instead of proposing a new
 92 algorithm, we analyze the following three classes of well-known video-based representation learning
 93 methods. Our goal is to understand whether these methods provably learn useful representations.

94 **Autoencoder.** This approach first maps a given observation x to an abstract state $\phi(x)$ using a decoder
 95 $\phi \in \Phi$, and then uses it to reconstruct the observation x with the aid of a reconstruction model class
 96 $\mathcal{Z} = \{z : [N] \rightarrow \mathcal{X}\}$. Formally, we optimize the following loss:

$$\ell_{\text{auto}}(z, \phi) = \frac{1}{nH} \sum_{i=1}^n \sum_{h=1}^H \|z(\phi(x_h^{(i)})) - x_h^{(i)}\|_2^2. \quad (1)$$

97 In practice, autoencoders are typically implemented using a Vector Quantized bottleneck trained in a
 98 Variational AutoEncoder manner, which is called the VQ-VAE approach [Oord et al., 2017].

99 **Forward Modeling.** This approach is similar to the autoencoder approach but instead of re-
 100 constructing the input observation, we reconstruct a future observation using a model class
 101 $\mathcal{F} = \{f : [N] \times [K] \rightarrow \Delta(\mathcal{X})\}$ where N is the output size of the decoder class Φ and $K \in \mathbb{N}$
 102 is a hyperparameter representing the forward time steps from the current observation. We collect a
 103 dataset of *multistep transitions* $\mathcal{D}_{\text{for}} = \{(x^{(i)}, k^{(i)}, x'^{(i)})\}_{i=1}^n$ sampled iid using the video dataset \mathcal{D}
 104 where the observation $x^{(i)}$ is sampled randomly from the i^{th} video, $k^{(i)} \in [K]$, and $x'^{(i)}$ is the frame
 105 $k^{(i)}$ -steps ahead of $x^{(i)}$ in the i^{th} video. We distinguish between two types of sampling procedures,
 106 one where $k^{(i)}$ is always a fixed given value $k \in [K]$, and one where $k^{(i)} \sim \text{Unf}([K])$. Given the
 107 dataset \mathcal{D}_{for} , we optimize the following loss:

$$\ell_{\text{for}}(f, \phi) = \frac{1}{n} \sum_{i=1}^n \ln f(x'^{(i)} | \phi(x^{(i)}), k^{(i)}). \quad (2)$$

108 **Temporal Contrastive Learning.** Finally, this approach trains the decoder ϕ to learn to separate a
 109 pair of temporally causal observations from a pair of temporally *acausal* observations. We collect
 110 a dataset of $\mathcal{D}_{\text{temp}} = \{(x^{(i)}, k^{(i)}, x'^{(i)}, z^{(i)})\}_{i=1}^{\lfloor n/2 \rfloor}$ tuples using the multistep transitions dataset
 111 \mathcal{D}_{for} . We use 2 multistep transitions to create a single datapoint for $\mathcal{D}_{\text{temp}}$ to keep the datapoints
 112 independent. To create the i^{th} datapoint for $\mathcal{D}_{\text{temp}}$, we use the multistep transitions $(x^{(2i)}, k^{(2i)}, x'^{(2i)})$
 113 and $(x^{(2i+1)}, k^{(2i+1)}, x'^{(2i+1)})$ and sample $z^{(i)} \sim \text{Unf}(\{0, 1\})$. If $z^{(i)} = 1$, then our i^{th} datapoint
 114 is a causal observation pair $(x^{(2i)}, k^{(2i)}, x'^{(2i)}, z^{(i)})$, otherwise, it is an acausal observation pair
 115 $(x^{(2i)}, k^{(2i)}, x^{(2i+1)}, z^{(i)})$. Depending on how we sample k , we collect a different dataset \mathcal{D}_{for} , and
 116 accordingly a different dataset $\mathcal{D}_{\text{temp}}$. Given the dataset $\mathcal{D}_{\text{temp}}$, we optimize the following loss using a
 117 regression model g belonging to a model class $\mathcal{G} = \{g : \mathcal{X} \times [K] \times \mathcal{X} \rightarrow [0, 1]\}$:

$$\ell_{\text{temp}}(g, \phi) = \frac{1}{\lfloor n/2 \rfloor} \sum_{i=1}^{\lfloor n/2 \rfloor} \left(z^{(i)} - g(\phi(x^{(i)}), k^{(i)}, x'^{(i)}) \right)^2. \quad (3)$$

118 **Practical Implementations.** We use the aforementioned description of methods for theoretical
 119 analysis. However, their practical implementations differ in a few notable ways. Most importantly

120 we either use a continuous vector representation $\phi : \mathcal{X} \rightarrow \mathbb{R}^d$ for modeling Φ , or apply a Vector
 121 Quantized (VQ) bottleneck [Oord et al., 2017] on top of the vector representation to model a discrete-
 122 representation decoder. We also optimize the loss using minibatches and use square loss for training
 123 forward modeling and SimCLR loss [Chen et al., 2020] for contrastive learning. We experimentally
 124 show that our theoretical findings extend to these practical implementations.

125 3 Is Video Based Representation Learning Provably Correct?

126 In this section, we present our main theoretical results. We first prove that both forward modeling and
 127 temporal contrastive methods succeed when there is no exogenous noise. We then establish a lower
 128 bound showing that video-based representation learning is exponentially harder than trajectory-based
 129 representation learning. We defer all proofs to the Appendix and only provide a sketch here.

130 3.1 Upper Bound in Block MDP Setting

131 We start by stating our theoretical setting and our main assumptions.

132 **Theoretical Setting.** We assume a Block MDP setting and access to a dataset $\mathcal{D} =$
 133 $\left\{ (x_1^{(i)}, x_2^{(i)}, \dots, x_H^{(i)}) \right\}_{i=1}^n$ of n independent and identically distributed (iid) videos sampled from
 134 data distribution D . We denote the probability of a video as $D(x_1, x_2, \dots, x_H)$. We assume that D
 135 is generated by a mixture of Markovian policies Π_D , i.e., the generative procedure for D is to sample
 136 a policy $\pi \in \Pi_D$ with some probability and then generate an entire episode using it. We assume
 137 that observations encode time steps. This can be trivially accomplished by simply concatenating the
 138 time step information to the observation. We also assume that the video data has good state space
 139 coverage and that the data is collected by *noise-free policies*.

140 **Assumption 1** (Requirements on Data Collection). *There exists an $\eta_{min} > 0$ such that if s is a state*
 141 *reachable at time step h by some policy in Π , then $D(\phi^*(x_h) = s) \geq \eta_{min}$. Further, we assume that*
 142 *every data collection policy $\pi \in \Pi_D$ is noise-free, i.e., $\pi(a | x_h) = \pi(a | \phi^*(x_h))$ for all (a, x_h) .*

143 **Justification for Assumption 1** In practice, we expect this assumption to hold for tasks such as
 144 gaming, or software debugging, where video data is abundant and, therefore, can be expected
 145 to provide good coverage of the underlying state space. This assumption is far weaker than the
 146 assumption in batch RL which also requires actions and rewards to be labeled, which makes it more
 147 expensive to collect data that has good coverage [Chen and Jiang, 2019]. Further, unlike imitation
 148 learning from observations (ILO) [Torabi et al., 2019], we don't require that these videos provide
 149 demonstrations of the desired behavior. E.g., video streaming of games is extremely common on the
 150 internet, and one can get many hours of video data this way. However, this data wouldn't come with
 151 actions (which will be mouse or keyboard strokes) or reward labeling, and the game levels or tasks
 152 in the data can be different or even unrelated to the downstream tasks we want to solve. As such, the
 153 video data do not provide demonstrations of the desired task. Further, as the video data is typically
 154 generated by humans, we can expect the data collection policies to be noise-free, as these policies are
 155 realized by humans who would not make decisions based on noise. E.g., a human player is unlikely
 156 to turn left due to the background motion of leaves that is unrelated to the game's control or objective.

157 We analyze the temporal contrastive learning and forward modeling approaches and derive upper
 158 bounds for these methods in Block MDPs. While autoencoder-based approaches sometimes do
 159 well in practice, it is an open question whether finite-sample bounds exist for them and we leave
 160 their theoretical analysis to future work and instead evaluate them empirically. In addition to the
 161 decoder class Φ , we assume a function class \mathcal{F} to model f for forward modeling and \mathcal{G} to model g
 162 for temporal contrastive learning. We make a realizability assumption on these function classes.

163 **Assumption 2** (Realizability). *There exists $f^* \in \mathcal{F}$, $g^* \in \mathcal{G}$ and $\phi_{for}, \phi_{temp} \in \Phi$ such that $f^*(X' |$
 164 $\phi_{for}(x), k) = \mathbb{P}_{for}(X' | x, k)$ and $g^*(z | \phi_{temp}(x), k, x') = \mathbb{P}_{temp}(z = 1 | x, k, x')$ on the appropriate
 165 support, and where \mathbb{P}_{for} and \mathbb{P}_{temp} are respectively the Bayes classifier for the forward modeling and
 166 temporal contrastive learning methods.*

167 **Justification for Assumption 2.** Realizability is a typical assumption made in theoretical analysis of
 168 RL algorithms [Agarwal et al., 2020]. Intuitively, the assumption states that the function classes are
 169 expressive enough to represent the Bayes classifier of their problem. In practice, this is usually not a

170 concern as we will use expressive deep neural networks to model these function classes. We will
 171 empirically show the feasibility of this assumption in our experiments.

172 Finally, we assume that our data distribution has the required information to separate the latent states.
 173 We state this assumption formally below and then show settings where this is true.

174 **Assumption 3** (Margin Assumption). *We assume that the margins β_{for} and β_{temp} defined below:*

$$\beta_{\text{for}} = \inf_{s_1, s_2 \in \mathcal{S}, s_1 \neq s_2} \mathbb{E}_k [\|\mathbb{P}_{\text{for}}(X' | s_1, k) - \mathbb{P}_{\text{for}}(X' | s_2, k)\|_{\text{TV}}]$$

$$\beta_{\text{temp}} = \inf_{s_1, s_2 \in \mathcal{S}, s_1 \neq s_2} \frac{1}{2} \mathbb{E}_{k, s'} [\|\mathbb{P}_{\text{temp}}(z = 1 | s_1, k, s') - \mathbb{P}_{\text{temp}}(z = 1 | s_2, k, s')\|],$$

175 *are strictly positive, and where in the definition of β_{temp} , we sample s' from the video data distribution*
 176 *and k is sampled according to our data collection procedure.*

177 **Justification for Assumption 3.** This assumption states that we need margins (β_{for}) for forward
 178 modeling and (β_{temp}) for temporal contrastive learning. A common scenario where these assump-
 179 tions are true is when for any pair of different states s_1, s_2 , there is a third state s' that is reachable
 180 from one but not the other. If the video data distribution D supports all underlying transitions, then
 181 this immediately implies that $\|\mathbb{P}_{\text{for}}(X' | s_1, k) - \mathbb{P}_{\text{for}}(X' | s_2, k)\|_{\text{TV}} > 0$ which implies $\beta_{\text{for}} > 0$.
 182 This scenario occurs in almost all navigation tasks. Specifically, it occurs in the three domains we
 183 experiment with. While it is less clear, under this assumption we also have $\beta_{\text{temp}} > 0$.

184 We now state our main result for forward modeling under Assumption 1-3.

185 **Theorem 1** (Forward Modeling Result). *Fix $\epsilon > 0$ and $\delta \in (0, 1)$ and let \mathcal{A} be any prov-*
 186 *ably efficient RL algorithm for tabular MDPs with sample complexity $n_{\text{samp}}(S, A, H, \epsilon, \delta)$. If*
 187 *n is $\text{poly}\{S, H, 1/\eta_{\text{min}}, 1/\beta_{\text{for}}, 1/\epsilon, \ln(1/\delta), \ln|\mathcal{F}|, \ln|\Phi|\}$ for a suitable polynomial, then forward*
 188 *modeling learns a decoder $\hat{\phi} : \mathcal{X} \rightarrow [|\mathcal{S}|]$. Further, running \mathcal{A} on the tabular MDP with*
 189 *$n_{\text{samp}}(S, A, H, T, \epsilon/2, \delta/4)$ episodes returns a latent policy $\hat{\varphi}$. Then there exists a bijective mapping*
 190 *$\alpha : \mathcal{S} \rightarrow [|\mathcal{S}|]$ such that with probability at least $1 - \delta$ we have:*

$$\forall s \in \mathcal{S}, \quad \mathbb{P}_{x \sim q(\cdot|s)} \left(\hat{\phi}(x) = \alpha(s) \mid \phi^*(x) = s \right) \geq 1 - \frac{4S^3 H^2}{\eta_{\text{min}}^2 \beta_{\text{for}}} \sqrt{\frac{1}{n} \ln \left(\frac{|\mathcal{F}| \cdot |\Phi|}{\delta} \right)},$$

191 *and the learned observation-based policy $\hat{\varphi} \circ \hat{\phi} : x \mapsto \hat{\varphi}(\hat{\phi}(x))$ is ϵ -optimal, i.e.,*

$$V(\pi^*) - V(\hat{\varphi} \circ \hat{\phi}) \leq \epsilon.$$

192 *Finally, the number of online episodes used in the downstream RL task is given by*
 193 *$n_{\text{samp}}(S, A, H, \epsilon_o/2, \delta_o/4)$ and doesn't scale with the complexity of function classes Φ and \mathcal{F} .*

194 The result for temporal contrastive is identical to Theorem 1 but instead of β_{for} we have β_{temp} and
 195 instead of \mathcal{F} we have \mathcal{G} . These upper bounds provide the desired result which shows that not only
 196 can we learn the right representation and near-optimal policy but also do so without online episodes
 197 scaling with $\ln|\Phi|$. Typically, the function class for forward modeling \mathcal{F} is much more complex than
 198 \mathcal{G} , however, as we show in Appendix C.5, the margin for forward modeling β_{for} is larger than for
 199 contrastive learning β_{temp} leading to a trade-off between these two approaches.

200 3.2 Learning from Video is Exponentially Harder than Learning from Trajectory Data

201 When online RL is possible, there exist algorithms Misra et al. [2020], Efroni et al. [2022] that can
 202 learn an accurate latent state decoder $\hat{\phi}$ with high probability and use it to learn near-optimal policies.
 203 These methods train the decoder using online trajectory data. This begs the following question: *Is it*
 204 *possible to learn a latent state decoder that is useful for performing RL using offline video data?* As
 205 the next result shows, this is not always the case.

206 **Theorem 2** (Lower Bound for Video). *Suppose $|\mathcal{S}|, |\mathcal{A}|, H \geq 2$. Then, for any $\epsilon \in (0, 1)$, any*
 207 *algorithm \mathcal{A}_1 that outputs a state decoder ϕ with $\phi_h : \mathcal{X} \rightarrow [L]$, $L \leq 2^{1/4\epsilon-1}$, $\forall h \in [H]$ given a*
 208 *video dataset \mathcal{D} sampled from some MDP and satisfies Assumption 1, and any online RL algorithm*
 209 *\mathcal{A}_2 uses that state decoder ϕ in its interaction with such an MDP (i.e., \mathcal{A}_2 only observes states*

210 through ϕ) and output a policy $\hat{\pi}$, there exists an MDP instance M in a class of MDPs which satisfies
 211 Assumption 3 and is PAC learnable with $\tilde{O}(\text{poly}(|\mathcal{S}|, |\mathcal{A}|, H, 1/\varepsilon))$ complexity, such that

$$V_M(\pi_M^*) - V_M(\hat{\pi}) > \varepsilon,$$

212 regardless of the size of the video dataset \mathcal{D} for algorithm \mathcal{A}_1 and the number of episodes of
 213 interaction for algorithm \mathcal{A}_2 .

214 The basic idea behind that hard instance construction is that, without the action information, it is
 215 impossible for the learning agent to distinguish between endogenous states and exogenous noise. For
 216 example, consider an image consisting of $N \times N$ identical mazes but where the agent controls just
 217 one maze. Other mazes contain other agents which are exogenous for our purpose. In the absence of
 218 actions, we cannot tell which maze is the one we are controlling and must memorize the configuration
 219 of all $N \times N$ mazes which grow exponentially with N . Another implication from that hard instance
 220 is – if the margin condition (Assumption 3) is violated, the exponentially large state decoder is also
 221 required for the regular block MDP without exogenous noise; a detailed discussion can also be found
 222 in Section C.3. We also discuss settings where we may be able to efficient-learning with just video
 223 data with additional assumptions in Appendix C.4.

224 4 Experimental Results and Discussion

225 We empirically evaluate the above video-based representation learning methods on three visual
 226 environments: a gridworld environment and two VizDoom environments. We defer the results on one
 227 of the VizDoom environments along with additional experimental details and results to Appendix D.
 228 Our main goal is to validate our theoretical findings by evaluating these methods in the presence and
 229 absence of exogenous noise and comparing their performance with a trajectory-based method.

230 4.1 Experimental Details

231 **GridWorld.** We consider navigation in a 12×12 Minigrid environment [Chevalier-Boisvert et al.,
 232 2023]. The agent (red triangle) can only observe an area around itself, and the goal is to reach the key
 233 quickly (Fig. 3). The position of the agent and key randomizes each episode.

234 **VizDoom Defend the Center** This is a first-person shooting game [Wydmuch et al., 2018, Kempka
 235 et al., 2016], in which the player needs to kill a variety of monsters to score (Fig. 5). The episode
 236 ends when the monster is killed or after 500 steps.

237 **Exogenous Noise.** For all domains, the observation is an RGB image. We add exogenous noise
 238 to it by superimposing 10 generated diamonds of a particular size. The color and position of
 239 these diamonds are our exogenous state. At the start of each episode, we randomly generate these
 240 diamonds, after which they move in a deterministic path. We also test the setting in which there
 241 is exogenous noise in the reward. We compute a score based on just the exogenous noise and add
 242 it to the reward presented to the agent. However, the agent is still evaluated on the original reward.

243 **Model and Learning.** Our decoder class Φ is a convolutional neural network. We use a deconvolu-
 244 tional neural network to model f and h . We experimented with both using a vector representation for
 245 ϕ and also using a VQ-bottleneck to discretize the embeddings. We use PPO to do downstream RL
 246 and keep ϕ frozen during the RL training. We also visualize the learned representations by training
 247 a decoder on them and fixing ϕ to reconstruct the input observations. We then look at the generated
 248 images to see what information from the observation is preserved by the representation.

249 **ACRO.** We also evaluate the learned representations against ACRO [Islam et al., 2022] which
 250 uses trajectory data. This approach learns representation ϕ by predicting action given a pair of
 251 observations $\mathbb{E}[\ln p(a_h | \phi(x_h), x_{h+k}, k)]$. ACRO is designed to filter out exogenous noise as this
 252 information is not predictive of the action. Our goal is to test if we get much better representations
 253 if we have access to trajectory data instead of video data.

254 4.2 Empirical Results and Discussion

255 We present our main empirical results in Fig. 2 and Fig. 4 and discuss the results below.

256 **Forward modeling and temporal contrastive both work when there is no exogenous noise.** In
 257 accordance with Theorem 1, we observe that in the case of both GridWorld (Figure 2) and ViZDoom

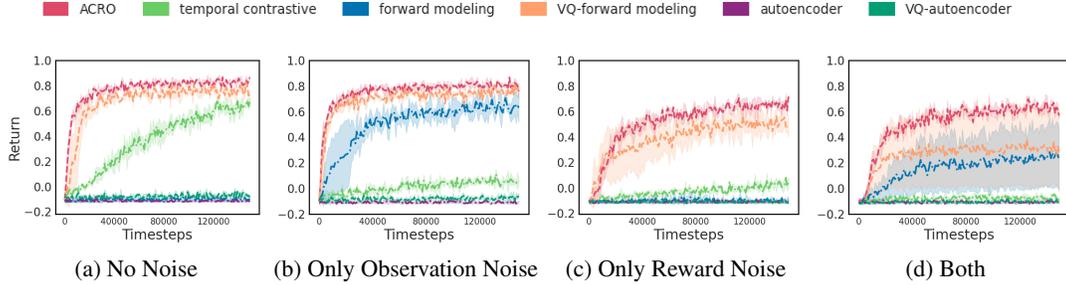


Figure 2: RL experiments in the GridWorld environment.

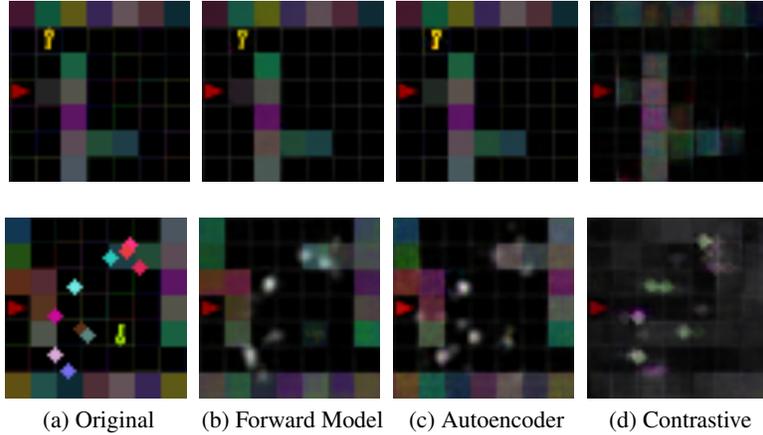


Figure 3: Decoded image reconstructions for different methods in the GridWorld environment. We train a reconstruction model on top of frozen learned representations ϕ trained with a given video-based method. **Top row:** shows an example from the setting where there is no exogenous noise. **Bottom row:** shows an example with exogenous noise (colored diamond shapes).

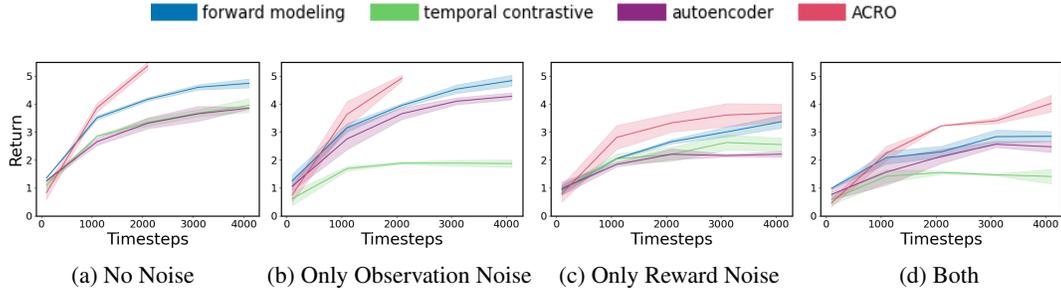


Figure 4: RL experiments using different latent representations for the ViZDoom Defend the Center environment.

258 Defend the Center (Figure 4), these approaches learn a decoder ϕ that lead to success with RL
 259 in the absence of any exogenous noise. For GridWorld, we find support for this result with VQ
 260 bottleneck during representation learning (Fig. 2(a)) whereas for ViZDoom Defend the Center,
 261 we find support for this result even without the use of a VQ bottleneck (Fig. 4(a)). These results
 262 are further supported via qualitative evaluation through image decoding from the learned latent
 263 representations (Fig. 3) which show that these representations can recover critical elements like walls.
 264 We find that autoencoder performs well in ViZDoom Defend the Center but not in gridworld, which
 265 aligns with a lack of any theoretical understanding of autoencoders.

266 **Performance with exogenous noise.** We find that in the presence of exogenous noise (Figure 2, Fig-
 267 ure 4), representations from forward modeling achieve a lower performance specially in gridworld,
 268 whereas temporal contrastive representations completely fail. One hypothesis for the stark failure of

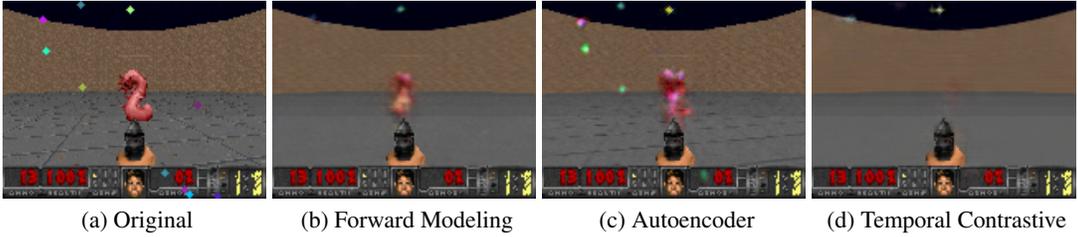


Figure 5: Decoded image reconstructions for different methods in ViZDoom Defend the Center.

temporal contrastive learning is that the agent can tell whether two observations are causal or not, by simply focusing on the noisy diamonds that move in a predictive manner. Therefore, the contrastive learning loss can be reduced by focusing entirely on the exogenous noise. Whereas, forward modeling is more robust as it needs to predict future observations, and the agent’s state is more helpful for doing that than noise. This shows in the reconstructions (Figure 3(b)(d), Figure 5(b)(d)). As expected, the reconstructions for forward modeling continue to capture state-relevant information, whereas for temporal contrastive they focus on noise and miss relevant state information. In Appendix C.6, we formally prove that there exists an instance where forward modeling can recover the latent state for low-levels of exogenous noise, whereas temporal contrastive cannot do so for any level of exogenous noise.

Comparison with ACRO. Finally, we draw a comparison between the performance of video-pretrained representation and ACRO which uses trajectory data. ACRO achieves the strongest performance across all tasks (Figure 2, Figure 4). Additionally, we also observe that as we increase the size of the exogenous noise elements in the observation space (Figure 6), the performance of forward modeling, the overall best video-based approach, degrades more drastically compared to ACRO. This agrees with our theoretical finding (Theorem 2) that learning representations from video-based data is significantly harder than trajectory-based data when exogenous noise is present.

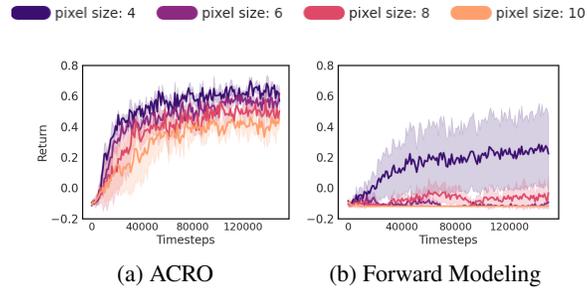


Figure 6: RL performance with varying size for exogenous noise in the GridWorld environment.

5 Conclusion

Videos are a naturally available source of data for training representations for RL. In this work, we study whether existing video-based representation learning methods are provably effective for downstream RL tasks. We provide both upper and lower bounds for these methods in two theoretical settings and provide empirical validation of our findings on three visual domains. Using our theoretical tools to develop better video-based representation learning methods and extending our analysis to other formal settings are natural future work directions.

References

- Alekh Agarwal, Sham Kakade, Akshay Krishnamurthy, and Wen Sun. Flambe: Structural complexity and representation learning of low rank mdps. *Advances in neural information processing systems*, 33:20095–20107, 2020.
- Arthur Aubret, Markus R. Ernst, Céline Teulière, and Jochen Triesch. Time to augment self-supervised visual representation learning. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=o8xdgmwCP81>.

- 310 Bowen Baker, Ilge Akkaya, Peter Zhokov, Joost Huizinga, Jie Tang, Adrien Ecoffet, Brandon
311 Houghton, Raul Sampedro, and Jeff Clune. Video pretraining (vpt): Learning to act by watching
312 unlabeled online videos. *Advances in Neural Information Processing Systems*, 35:24639–24654,
313 2022.
- 314 Homanga Bharadhwaj, Mohammad Babaeizadeh, Dumitru Erhan, and Sergey Levine. Information
315 prioritization through empowerment in visual model-based rl. *arXiv preprint arXiv:2204.08585*,
316 2022.
- 317 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
318 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
319 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 320 Jinglin Chen and Nan Jiang. Information-theoretic considerations in batch reinforcement learning. In
321 *International Conference on Machine Learning*, pages 1042–1051. PMLR, 2019.
- 322 Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for
323 contrastive learning of visual representations. In *International conference on machine learning*,
324 pages 1597–1607. PMLR, 2020.
- 325 Maxime Chevalier-Boisvert, Bolun Dai, Mark Towers, Rodrigo de Lazcano, Lucas Willems, Salem
326 Lahlou, Suman Pal, Pablo Samuel Castro, and Jordan Terry. Minigrid & miniworld: Modular &
327 customizable reinforcement learning environments for goal-oriented tasks. *CoRR*, abs/2306.13831,
328 2023.
- 329 Simon Du, Akshay Krishnamurthy, Nan Jiang, Alekh Agarwal, Miroslav Dudik, and John Langford.
330 Provably efficient rl with rich observations via latent state decoding. In *International Conference*
331 *on Machine Learning*, pages 1665–1674. PMLR, 2019.
- 332 Yonathan Efroni, Dipendra Misra, Akshay Krishnamurthy, Alekh Agarwal, and John Langford.
333 Provably filtering exogenous distractors using multistep inverse dynamics. In *International*
334 *Conference on Learning Representations*, 2022. URL [https://openreview.net/forum?](https://openreview.net/forum?id=RQLLzMCefQu)
335 [id=RQLLzMCefQu](https://openreview.net/forum?id=RQLLzMCefQu).
- 336 Alejandro Escontrela, Ademi Adeniji, Wilson Yan, Ajay Jain, Xue Bin Peng, Ken Goldberg, Young-
337 woon Lee, Danijar Hafner, and Pieter Abbeel. Video prediction models as rewards for reinforcement
338 learning. *arXiv preprint arXiv:2305.14343*, 2023.
- 339 Sara A Geer. *Empirical Processes in M-estimation*, volume 6. Cambridge university press, 2000.
- 340 Wonjoon Goo and Scott Niekum. One-shot learning of multi-step tasks from observation via activity
341 localization in auxiliary video. In *2019 international conference on robotics and automation*
342 *(ICRA)*, pages 7755–7761. IEEE, 2019.
- 343 Zhaohan Guo, Shantanu Thakoor, Miruna Pîslar, Bernardo Avila Pires, Florent Alché, Corentin
344 Tallec, Alaa Saade, Daniele Calandriello, Jean-Bastien Grill, Yunhao Tang, et al. Byol-explore:
345 Exploration by bootstrapped prediction. *Advances in neural information processing systems*, 35:
346 31855–31870, 2022.
- 347 Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning
348 behaviors by latent imagination. *arXiv preprint arXiv:1912.01603*, 2019. URL [https://arxiv.](https://arxiv.org/pdf/1912.01603.pdf)
349 [org/pdf/1912.01603.pdf](https://arxiv.org/pdf/1912.01603.pdf).
- 350 Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains
351 through world models. *arXiv preprint arXiv:2301.04104*, 2023.
- 352 Riashat Islam, Manan Tomar, Alex Lamb, Yonathan Efroni, Hongyu Zang, Aniket Didolkar, Dipendra
353 Misra, Xin Li, Harm van Seijen, Remi Tachet des Combes, et al. Agent-controller representations:
354 Principled offline rl with rich exogenous information. *arXiv preprint arXiv:2211.00164*, 2022.
- 355 Nan Rosemary Ke, Amanpreet Singh, Ahmed Touati, Anirudh Goyal, Yoshua Bengio, Devi Parikh,
356 and Dhruv Batra. Learning dynamics model in reinforcement learning by incorporating the long
357 term future, 2019.

- 358 Michael Kearns and Satinder Singh. Near-optimal reinforcement learning in polynomial time.
359 *Machine learning*, 49:209–232, 2002.
- 360 Michał Kempka, Marek Wydmuch, Grzegorz Runc, Jakub Toczek, and Wojciech Jaśkowski. Vizdoom:
361 A doom-based ai research platform for visual reinforcement learning. In *2016 IEEE conference on*
362 *computational intelligence and games (CIG)*, pages 1–8. IEEE, 2016.
- 363 Alex Lamb, Riashat Islam, Yonathan Efroni, Aniket Didolkar, Dipendra Misra, Dylan Foster, Lekan
364 Molu, Rajan Chari, Akshay Krishnamurthy, and John Langford. Guaranteed discovery of control-
365 lable latent states with multi-step inverse models. *arXiv preprint arXiv:2207.08229*, 2022. URL
366 <https://arxiv.org/pdf/2207.08229.pdf>.
- 367 Xudong Lin, Gedas Bertasius, Jue Wang, Shih-Fu Chang, Devi Parikh, and Lorenzo Torresani.
368 Vx2text: End-to-end learning of video-based text generation from multimodal inputs. In *Proceed-*
369 *ings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7005–7015,
370 2021.
- 371 Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike
372 Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining
373 approach. *arXiv preprint arXiv:1907.11692*, 2019.
- 374 Zakaria Mhammedi, Dylan J Foster, and Alexander Rakhlin. Representation learning with multi-step
375 inverse kinematics: An efficient and optimal approach to rich-observation rl. *arXiv preprint*
376 *arXiv:2304.05889*, 2023.
- 377 Vincent Micheli, Eloi Alonso, and François Fleuret. Transformers are sample-efficient world
378 models. In *The Eleventh International Conference on Learning Representations*, 2023. URL
379 <https://openreview.net/forum?id=vhFulAcb0xb>.
- 380 Dipendra Misra, Mikael Henaff, Akshay Krishnamurthy, and John Langford. Kinematic state abstrac-
381 tion and provably efficient rich-observation reinforcement learning. In *International conference on*
382 *machine learning*, pages 6961–6971. PMLR, 2020.
- 383 Ofir Nachum, Shixiang Gu, Honglak Lee, and Sergey Levine. Near-optimal representation learning
384 for hierarchical reinforcement learning. In *International Conference on Learning Representations*,
385 2018.
- 386 Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learning.
387 *arXiv preprint arXiv:1711.00937*, 2017.
- 388 Nikhil Parthasarathy, SM Eslami, João Carreira, and Olivier J Hénaff. Self-supervised video pretrain-
389 ing yields strong image representations. *arXiv preprint arXiv:2210.06433*, 2022.
- 390 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
391 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
392 models from natural language supervision. In *International conference on machine learning*, pages
393 8748–8763. PMLR, 2021.
- 394 Harshit Sikchi, Akanksha Saran, Wonjoon Goo, and Scott Niekum. A ranking game for imitation
395 learning. *arXiv preprint arXiv:2202.03481*, 2022.
- 396 Vlad Sobal, Jyothir SV, Siddhartha Jalagam, Nicolas Carion, Kyunghyun Cho, and Yann LeCun.
397 Joint embedding predictive architectures focus on slow features. *arXiv preprint arXiv:2211.10831*,
398 2022.
- 399 Nitish Srivastava, Elman Mansimov, and Ruslan Salakhudinov. Unsupervised learning of video
400 representations using lstms. In Francis Bach and David Blei, editors, *Proceedings of the*
401 *32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine*
402 *Learning Research*, pages 843–852, Lille, France, 07–09 Jul 2015. PMLR. URL <https://proceedings.mlr.press/v37/srivastava15.html>.
403
- 404 Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.

- 405 Yunhao Tang, Zhaohan Daniel Guo, Pierre Harvey Richemond, Bernardo Avila Pires, Yash Chandak,
406 Remi Munos, Mark Rowland, Mohammad Gheshlaghi Azar, Charline Le Lan, Clare Lyle, András
407 György, Shantanu Thakoor, Will Dabney, Bilal Piot, Daniele Calandriello, and Michal Valko.
408 Understanding self-predictive learning for reinforcement learning. In *Proceedings of the 40th*
409 *International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning*
410 *Research*, pages 33632–33656. PMLR, 23–29 Jul 2023. URL [https://proceedings.mlr.](https://proceedings.mlr.press/v202/tang23d.html)
411 [press/v202/tang23d.html](https://proceedings.mlr.press/v202/tang23d.html).
- 412 Faraz Torabi, Garrett Warnell, and Peter Stone. Recent advances in imitation learning from ob-
413 servation. *arXiv preprint arXiv:1905.13566*, 2019. URL [https://arxiv.org/pdf/1905.](https://arxiv.org/pdf/1905.13566.pdf)
414 [13566.pdf](https://arxiv.org/pdf/1905.13566.pdf).
- 415 Masatoshi Uehara, Xuezhou Zhang, and Wen Sun. Representation learning for online and offline rl
416 in low-rank mdps. *arXiv preprint arXiv:2110.04652*, 2021.
- 417 Tongzhou Wang, Simon S Du, Antonio Torralba, Phillip Isola, Amy Zhang, and Yuandong Tian. De-
418 noised mdps: Learning world models better than the world itself. *arXiv preprint arXiv:2206.15477*,
419 2022. URL <https://arxiv.org/pdf/2206.15477.pdf>.
- 420 Marek Wydmuch, Michał Kempka, and Wojciech Jaśkowski. Vizdoom competitions: Playing doom
421 from pixels. *IEEE Transactions on Games*, 11(3):248–259, 2018.
- 422 Weirui Ye, Yunsheng Zhang, Pieter Abbeel, and Yang Gao. Become a proficient player with limited
423 data through watching pure videos. In *The Eleventh International Conference on Learning*
424 *Representations*, 2022.
- 425 Amy Zhang, Rowan McAllister, Roberto Calandra, Yarin Gal, and Sergey Levine. Learning
426 invariant representations for reinforcement learning without reconstruction. *arXiv preprint*
427 *arXiv:2006.10742*, 2020.
- 428 Tony Zhao, Siddharth Karamcheti, Thomas Kollar, Chelsea Finn, and Percy Liang. What makes
429 representation learning from videos hard for control? *RSS Workshop on Scaling Robot Learning*,
430 2022. URL <https://api.semanticscholar.org/CorpusID:252635608>.

431 **A Preliminaries and Overview**

432 In this section, we provide a formal overview of our learning setup and problem statement.

433 **Mathematical Notation.** We use $[N]$ for $N \in \mathbb{N}$ to define the set $\{1, 2, \dots, N\}$. We assume all sets
 434 to be countable. For a given set \mathcal{U} , we denote its cardinality by $|\mathcal{U}|$ and define $\Delta(\mathcal{U})$ as the space of
 435 all distributions over \mathcal{U} . We denote the uniform distribution over \mathcal{U} by $\text{Unf}(\mathcal{U})$. Finally, $\text{poly}\{\cdot\}$
 436 denotes a term that scales polynomially in the listed quantities.

437 **Block MDPs.** We study episodic RL in Block Markov Decision Processes (Block MDP) [Du et al.,
 438 2019]. A Block MDP is defined by the tuple $(\mathcal{X}, \mathcal{S}, \mathcal{A}, T, R, q, \mu, H)$ where \mathcal{X} is a set of observations
 439 that can be infinitely large, \mathcal{S} is a finite set of *latent* states, and \mathcal{A} is a set of finite actions. The
 440 transition dynamics $T : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ define transitions in the latent state space. The reward
 441 function $R : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ assigns a reward $R(s, a)$ if action a is taken in the latent state s . When
 442 the agent visits a state s , it receives an observation $x \sim q(\cdot | s)$ sampled from an emission function
 443 $q : \mathcal{S} \rightarrow \Delta(\mathcal{X})$. This emission process contains temporally independent noise but no exogenous
 444 noise. Finally, $\mu \in \Delta(\mathcal{S})$ is the distribution over the initial latent state and H is the horizon denoting
 445 the number of actions per episode. The agent interacts with a block MDP environment by repeatedly
 446 generating an episode $(x_1, a_1, r_1, \dots, x_H, a_H, r_H)$ where $s_1 \sim \mu$ and for all $h \in [H]$ we have
 447 $x_h \sim q(\cdot | s_h)$, $r_h = R(s_h, a_h)$, and $s_{h+1} \sim T(\cdot | s_h, a_h)$, and all actions $\{a_h\}_{h=1}^H$ are taken by the
 448 agent. The agent never directly observes the latent states (s_1, s_2, \dots, s_H) .

449 A key assumption in Block MDPs is that two different latent states cannot generate the same
 450 observation. This is called the *disjoint emission property* and holds in many game and OS settings.
 451 Formally, this property allows us to define a decoder $\phi^* : \mathcal{X} \rightarrow \mathcal{S}$ that maps an observation to the
 452 unique state that can generate it. The agent does not have access to ϕ^* . If the agent had access to ϕ^* ,
 453 one could map each observation from an infinitely large space to the finite latent state space, which
 454 allows the use of classical finite RL methods [Kearns and Singh, 2002].

455 **Exogenous Block MDPs (Ex-Block MDP).** We also consider RL in Exogenous Block MDPs (Ex-
 456 Block MDPs) that extend Block MDPs to include exogenous noise [Efroni et al., 2022]. An Ex-Block
 457 MDP is defined by $(\mathcal{X}, \mathcal{S}, \Xi, \mathcal{A}, T, T_\xi, R, q, H, \mu, \mu_\xi)$ where $\mathcal{X}, \mathcal{S}, \mathcal{A}, T, R, H$ and μ have the same
 458 meaning and type as in Block MDPs. The additional quantities include Ξ which is the space of
 459 exogenous noise and can be infinitely large. We use the notation $\xi \in \Xi$ to denote the exogenous noise.
 460 For the setting in Fig. 1, the exogenous noise variable ξ captures variables such as the position of
 461 geese, the position of leaves on the trees in the background, and lighting conditions. The exogenous
 462 noise ξ changes with time according to the transition function $T_\xi : \Xi \rightarrow \Delta(\Xi)$ and is at start sampled
 463 from μ_ξ . Note that unlike the agent state $s \in \mathcal{S}$, the exogenous noise $\xi \in \Xi$, evolves independently
 464 of the agent’s action and does not influence the evolution of the agent’s state. The emission process
 465 $q : \mathcal{S} \times \Xi \rightarrow \Delta(\mathcal{X})$ in Ex-Block MDP uses both the current agent state and exogenous noise, to
 466 generate the observation at a given time. For example, the image generated by the agent’s camera
 467 contains information based on the agent’s state (e.g., agent’s position and orientation), along with
 468 exogenous noise (e.g., the position of geese). Similar to the Block MDP, we assume there exists
 469 *unknown* decoders $\phi^* : \mathcal{X} \rightarrow \mathcal{S}$ and $\phi_\xi^* : \mathcal{X} \rightarrow \Xi$ that can map an observation to the current agent
 470 state s and exogenous ξ respectively.

471 **Provable RL.** We assume access to a policy class $\Pi = \{\pi : \mathcal{X} \rightarrow \mathcal{A}\}$ where a policy $\pi \in \Pi$ allows the
 472 agent to take actions. For a given policy π , we use $\mathbb{E}_\pi[\cdot]$ to denote expectation taken over an episode
 473 generated by sampling actions from π . We define the value of a policy $V(\pi) = \mathbb{E}_\pi \left[\sum_{h=1}^H r_h \right]$
 474 as the expected total reward or expected return. Our goal is to learn a near-optimal policy $\hat{\pi}$, i.e.,
 475 $\sup_{\pi \in \Pi} V(\pi) - V(\hat{\pi}) \leq \epsilon$ with probability at least $1 - \delta$ for a given tolerance parameter $\epsilon > 0$
 476 and failure probability $\delta \in (0, 1)$, using number of episodes that scale polynomially in $1/\epsilon, 1/\delta$, and
 477 other relevant quantities. We will call such an algorithm as provably efficient. There exist several
 478 provably efficient RL approaches for solving Block MDPs [Mhammedi et al., 2023, Misra et al.,
 479 2020], and Ex-Block MDPs [Efroni et al., 2022]. These approaches typically assume access to a
 480 decoder class $\Phi = \{\phi : \mathcal{X} \rightarrow [N]\}$ and attempt to learn ϕ^* using it. These algorithms don’t use
 481 any pre-training and instead directly interact with the environment and learn a near-optimal policy
 482 by using samples that scale with $\text{poly}(S, A, H, \ln|\Phi|, 1/\epsilon, 1/\delta)$. Crucially, the dependence on $\ln|\Phi|$
 483 cannot be removed. The decoder class Φ and all other function classes in this work are assumed to

484 have bounded statistical complexity measures. For simplicity, we will assume that these function
485 classes are finite and derive guarantees that scale logarithmically in their size (e.g., $\ln|\Pi|$).¹

486 **Representation Pre-training using Videos.** RL algorithms for the above settings require online
487 episodes that scale with $\ln|\Phi|$ which is expensive for real-world problems where Φ is represented by a
488 complex neural network. Offline RL approaches Uehara et al. [2021] offer a substitute for expensive
489 online interactions but require access to labeled episodes (with actions and rewards) that are not
490 naturally available in many settings such as games and software. In contrast, we focus on pre-training
491 the decoder ϕ using video data which is naturally available in these settings.

492 **Problem Statement.** We are given two hyperparameters $\epsilon > 0$ and $\delta \in (0, 1)$ and a sufficiently large
493 dataset of videos. We are also given a decoder class $\Phi = \{\phi : \mathcal{X} \rightarrow [N]\}$ containing decoders that
494 map an observation to one of the N possible *abstract states*. During the pre-training phase, we learn
495 a decoder $\phi \in \Phi$ using the video data. We then freeze ϕ and use it to do RL in a downstream task.
496 Instead of using any particular choice of algorithm for RL, we assume we are given a provably efficient
497 tabular RL algorithm \mathcal{A} . We convert the observation-based RL problem to a tabular MDP problem by
498 converting an observation x to its abstract state representation $\phi(x)$ using the frozen learned decoder
499 ϕ . The algorithm \mathcal{A} uses $\phi(x)$ instead of x and outputs an *abstract policy* $\varphi : [N] \rightarrow \mathcal{A}$. We want
500 that $\sup_{\pi \in \Pi} V(\pi) - V(\varphi \circ \phi) \leq \epsilon$ with probability at least $1 - \delta$, where $\varphi \circ \phi : x \mapsto \varphi(\phi(x))$ is
501 our learned policy. We also require the number of online episodes in the downstream RL phase to
502 not scale with the size of the decoder class Φ . This allows us to minimize expensive online episodes
503 while using naturally available offline video data for pre-training.

504 B Additional Related Work

505 **Representation Learning for Reinforcement Learning** A line of research on recurrent state space
506 models is essentially concerned with the next-frame approach, although typically with conditioning on
507 actions. Moreover, to model uncertainty in the observations, a latent variable with a posterior depend-
508 ing on the current observation (or even a sequence of future observations) is typically introduced. [Ke
509 et al., 2019] considered learning such a sequential prediction model which predicts observations and
510 conditions on actions. They used a latent variable with a posterior depending on future observations
511 to model uncertainty. These representations were used for model-predictive control and improved
512 imitation learning. Dreamer [Hafner et al., 2019, 2023] uses the next-frame objective but also condi-
513 tions on actions. The IRIS algorithm [Micheli et al., 2023] uses the next-frame objective but uses the
514 transformer architecture, again conditioning on actions. The InfoPower approach [Bharadhwaj et al.,
515 2022] combines a one-step inverse model with a temporal contrastive objective. Sobal et al. [2022]
516 explored using semi-supervised objectives for learning representations in RL, yet used action-labeled
517 data. Wang et al. [2022] used a decoupled recurrent neural network approach to learn to extract
518 endogenous states, but relied on action-labeled data to achieve the factorization. Deep Bisimulation
519 for Control [Zhang et al., 2020] introduced an objective to encourage observations with similar value
520 functions to map to similar representations.

521 Self-prediction methods such as BYOL-explore [Guo et al., 2022] proposed learning reward-free
522 representations for exploration, but depended on open-loop prediction of future states conditioned on
523 actions. An analysis paper studied a simplified action-free version of the self-prediction objective
524 [Tang et al., 2023] and showed results in the absence of using actions, although this has not been
525 instantiated empirically to our knowledge.

526 A further line of work from theoretical reinforcement learning has examined provably efficient
527 objectives for discovering representations. Efroni et al. [2022] explored representation learning in
528 the presence of exogenous noise, establishing a sample efficient algorithm. However Efroni et al.
529 [2022] and the closely related work on filtering exogenous noise required actions [Lamb et al., 2022,
530 Islam et al., 2022]. Other theoretical work on learning representations for RL has required access to
531 action-labeled data [Misra et al., 2020].

532 **Representation Learning from Videos** Self-supervised representation learning from videos has
533 a long history. Srivastava et al. [2015] used recurrent neural networks with a pixel prediction
534 objective on future frames. Parthasarathy et al. [2022] explored temporal contrastive objectives for
535 self-supervised learning from videos. They also found that the features learned well aligned with

¹Our theoretical analyses can be extended to other complexity metrics such as Rademacher complexity.

536 human perceptual priors, despite the model not being explicitly trained to achieve such alignment.
 537 Aubret et al. [2023] applied temporal contrastive learning to videos of objects being manipulated in a
 538 3D space, showing that this outperformed standard augmentations used in computer vision.

539 **Using Video Data for Reinforcement Learning** The VIPER method [Escontrela et al., 2023] uses
 540 a pre-trained autoregressive generative model over action-free expert videos as a reward signal for
 541 training an imitation learning agent. The Video Pre-training (VPT) algorithm [Baker et al., 2022]
 542 trained an inverse kinematics model on a small dataset of Minecraft videos and used the model
 543 to label a large set of unlabeled Minecraft videos from the internet. This larger dataset was then
 544 used for imitation learning and reinforcement learning for downstream tasks. Zhao et al. [2022]
 545 explicitly studied the challenges in using videos for representation learning in RL, identifying five
 546 key factors: task mismatch, camera configuration, visual feature shift, sub-optimal behaviors in the
 547 data, and robot morphology. Goo and Niekum [2019] learn reward functions for multi-step tasks
 548 from videos by leveraging a single video segmented with action labels (one-shot learning). Sikchi
 549 et al. [2022] propose a two-player ranking game between a policy and a reward function to satisfy
 550 pairwise performance rankings between behaviors. Their proposed method achieves state-of-the-art
 551 sample efficiency and can solve previously unsolvable tasks in the learning from observation (no
 552 actions) setting.

553 Recently some approaches have also considered recovering *latent actions* from video data using an
 554 encoder-decoder approach [Ye et al., 2022]. In general, the lower bound in Theorem 2 applies to these
 555 methods and they do not provably work in the hard instances with exogenous noise. For example,
 556 the latent actions can capture *exogenous noise* instead of actions, if the former is more predictive of
 557 changes in the observations. However, in simpler cases such as 3D games, where the agent’s action is
 558 typically most predictive of changes in observations, or in settings with no exogenous noise, one can
 559 expect these approaches to do well.

560 C Proofs of Theoretical Statements

561 We state our setting and general assumptions before presenting method specific results. We also
 include a table of notations in Table 1.

Notation	Description
$[N]$	Denotes the set $\{1, 2, \dots, N\}$
$\Delta(\mathcal{U})$	Denotes the set of all distributions over a set \mathcal{U}
$\text{Unf}(\mathcal{U})$	Uniform distribution over \mathcal{U}
$\text{supp}(\mathbb{P})$	Support of a distribution $\mathbb{P} \in \Delta(\mathcal{U})$, i.e., $\text{supp}(\mathbb{P}) = \{x \in \mathcal{U} \mid \mathbb{P}(x) > 0\}$.
\mathcal{X}	Observation space
\mathcal{S}	Latent endogenous state
\mathcal{A}	Action space
$T : \mathcal{S} \rightarrow \mathcal{A} \rightarrow \Delta(\mathcal{S})$	Transition dynamics
$R : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$	Reward function
μ	Start state distribution
H	Horizon indicating the maximum number of actions per episode
$\phi^* : \mathcal{X} \rightarrow \mathcal{S}$	Endogenous state decoder

Table 1: Description for mathematical notations.

562
 563 We are given a dataset $\mathcal{D} = \left\{ (x_1^{(i)}, x_2^{(i)}, \dots, x_H^{(i)}) \right\}_{i=1}^n$ of n independent and identically distributed
 564 (iid) unlabeled episodes. We will use the word video and unlabeled episodes interchangeably. We
 565 assume the underlying data distribution is D . We denote the probability of an unlabeled episode as
 566 $D(x_1, x_2, \dots, x_H)$. We assume that D is generated by a mixture of Markovian policies Π_D , i.e., the
 567 generative procedure for D is to sample a policy $\pi \in \Pi_D$ with probability Θ_π and then generate
 568 an entire episode using it. For this reason, we will denote $D = \Theta \circ \Pi_D$ where Θ is the mixture
 569 distribution. We assume *no direct knowledge* about either Π_D or Θ , other than that the set of policies
 570 in Π_D are Markovian. We define the *underlying* distribution over the action-labeled episode as

571 $D(x_1, a_1, x_2, \dots, x_H, a_H)$, of which the agent *only* gets to observe the (x_1, x_2, \dots, x_H) . We will
 572 use the notation D to refer to any distribution that is derived from the above joint distribution.

573 We assume that observations encode time steps. This can be trivially accomplished by simply
 574 concatenating the time step information to the observation. This also implies that observations from
 575 different time steps are different. Because of this property, we can assume that the Markovian policies
 576 used to realize D were time homogenous, i.e., they only depend on observation and not observation
 577 and timestep pair (this is because we include timesteps in the observation). Therefore, for all $h \in [H]$
 578 and $k \in \mathbb{N}$ we have:

$$D(x_{h+k} = x' \mid x_h = x) = D(x_{k+1} = x' \mid x_1 = x) \quad (4)$$

579 We denote $D(x_h)$ to define the marginal distribution over an observation x_h , and $D(x_h, x_{h+k})$ to
 580 denote the marginal distribution over a pair of observations (x_h, x_{h+k}) in the episode. We similarly
 581 define $D(x_h, a_h)$ as the distribution over observation action pairs (x_h, a_h) .

582 We assume that the video data has good coverage. This is stated formally below:

583 **Assumption 4** (State Coverage by D). *Given our policy class Π , there exists an $\eta_{\min} > 0$ such that if*
 584 $\sup_{\pi \in \Pi} \mathbb{P}_{\pi}(s_h = s) > 0$ *for some $s \in \mathcal{S}$, then we assume $D(\phi^*(x_h) = s) \geq \eta_{\min}$.*

585 In practice, Assumption 4 can be satisfied since videos are more easily available than labeled episodes
 586 and we can hope that a large diverse collection of videos can provide reasonable coverage over the
 587 underlying state action space. E.g., for tasks like gaming, one can use hours of streaming data from
 588 many users.

589 Further, we also assume that the data policy depends only on the endogenous state. Recall that for an
 590 observation $x \in \mathcal{X}$, its endogenous state is given by $\phi^*(x) \in \mathcal{S}$.

591 **Assumption 5** (Noise-Free Video Distribution). *For any h , $\pi \in \Pi_D$, $x_h \in \text{supp } \mathbb{P}_{\pi}$ and $a \in \mathcal{A}$, we*
 592 *have*

$$\pi(a \mid x_h) = \pi(a \mid \phi^*(x_h)).$$

593 **Justification of Noise-Free Policy.** Typically, video data is created by humans. E.g., a human may
 594 be playing a game and the video data is collected by recording the user’s screen. A user is unlikely to
 595 take actions relying on iid or exogenous noise in the observation process. Therefore, the collected
 596 data can be expected to obey the noise-free assumption.

597 **Multi-step transition.** We choose to analyze a multi-step variant of standard temporal contrastive
 598 and forward modeling algorithms that train on a dataset of pairs of observations (x, x') that can be
 599 variable time steps apart. As our proof will show, this gives the algorithms more expressibility and
 600 allows them to learn correct representations for some problems that their single-step variants (i.e., the
 601 observations are adjacent) or fixed time-step variants (i.e., the observations are fixed time steps apart)
 602 cannot solve. We will use the variable k to denote the time steps by which these observations differ.
 603 Formally, we will call (x, k, x') as a multi-step transition where x was observed at some time step h ,
 604 and x' was observed at $h + k$. For the single-step variant of the algorithms, we have $k = 1$. For the
 605 fixed multi-step variant, we have $k > 1$ but k is fixed. Finally, in the general multi-step variant, we
 606 will assume that k is picked from $\text{Unf}([K])$ where K is a fixed upper bound.

607 **Extending episode to $H + K$.** When using $k > 1$, we may want to collect a multi-step transition
 608 (x, k, x') where $x = x_H$ to allow learning state representation for time step H . However, at this point,
 609 we don’t have time steps left to observe x_{H+k} . We alleviate this by assuming that we can allow an
 610 episode to run till $H + K$ if necessary. In practice, this is not a problem where the algorithm sets the
 611 horizon and not the environment. However, if we cannot go past H , then we can instead assume that
 612 all states are reachable by the time step $H - K$ and so their state representation can be learned when
 613 x is selected at x_{H-K} . In our analysis ahead, we make the former setting that the episodes can be
 614 extended to $H + K$, but it can be easily rephrased to work with the other setting.

615 For both the forward model and the temporal contrastive approach, we assume access to a dataset
 616 $\mathcal{D}_{\text{for}} = ((x^{(i)}, k^{(i)}, x'^{(i)}))_{i=1}^n$ of pairs of observations. We define a few different distributions that
 617 can be used to generate this set. For a given $k \in [K]$, we define a distribution D_k over k -step separate

618 observations as:

$$D_k(X = x, X' = x') = \frac{1}{H} \sum_{h=1}^H D(x_h = x, x_{h+k} = x') \quad (5)$$

619 We can sample $(x, k, x') \sim D_k(X, X')$ by sampling an episode $(x_1, x_2, \dots, x_H) \sim D$, and then
620 sampling a $h \sim \text{Unf}([H])$, and choosing $x = x_h$ and $x' = x_{h+k}$.

621 We also define a distribution D_{unf} where we also sample k uniformly over available choices:

$$D_{\text{unf}}(X = x, k, X' = x') = \frac{1}{K} D_k(x_h = x, x_{h+k} = x') \quad (6)$$

622 We can sample $(x, k, x') \sim D_{\text{unf}}(X, X')$ by sampling an episode $(x_1, x_2, \dots, x_H) \sim D$, and then
623 sampling $h \in [H]$, and sampling $k \in [K]$, and choosing (x_h, x_{h+k}) as the selected pair.

624 We define a useful notation $\rho \in \Delta(\mathcal{X})$ as:

$$\rho(X = x) = \frac{1}{H} \sum_{h=1}^H D(x_h = x). \quad (7)$$

625 The distribution $\rho(X)$ is a good distribution to sample from as it covers states across all time steps.
626 Finally, because of Assumption 4, we have the following:

$$\forall s \in \mathcal{S}, \quad \rho(s) \geq \frac{\eta_{\min}}{H} \quad (8)$$

627 This is because we assume every state $s \in \mathcal{S}$, is visited at some time step t , and so we have
628 $D(s_t = s) \geq \eta_{\min}$, and $\rho(s) = \frac{1}{H} \sum_{h=1}^H D(s_h = s) \geq \frac{1}{H} D(s_t = s) \geq \frac{\eta_{\min}}{H}$.

629 It can be easily verified that for both $D_k(X, X')$ and $D_{\text{unf}}(X, X')$, their marginals over X is given
630 by $\rho(X)$. Both D_k and D_{unf} satisfy the noise-free property. We prove this using the next two
631 Lemma.s

632 **Lemma 1** (Property of Noise-Free policy). *Let π be a policy such that for any $x \in \mathcal{X}$, we have
633 $\pi(a | x) = \pi(a | \phi^*(x))$. Then for any $h \in [H]$ and $k \in [K]$ we have $\mathbb{P}_\pi(x_{h+k} = x' | x_h = x)$ only
634 depend on $\phi^*(x)$ and this common value is defined by $\mathbb{P}_\pi(x_{h+k} | s_h = \phi^*(x))$.*

635 *Proof.* The proof is by induction on k . For $k = 1$ we have:

$$\mathbb{P}_\pi(x_{h+1} = x' | x_h = x) = \sum_{a \in \mathcal{A}} T(x' | x, a) \pi(a | x_h = x) = \sum_{a \in \mathcal{A}} T(x' | \phi^*(x), a) \pi(a | x_h = \phi^*(x)),$$

636 and as the right hand side only depends on $\phi^*(x)$, the base case is proven. For the general case, we
637 have:

$$\begin{aligned} \mathbb{P}_\pi(x_{h+k} = x' | x_h = x) &= \sum_{\tilde{x} \in \mathcal{X}} \mathbb{P}_\pi(x_{h+k} = x', x_{h+k-1} = \tilde{x} | x_h = x) \\ &= \sum_{\tilde{x} \in \mathcal{X}} \mathbb{P}_\pi(x_{h+k} = x' | x_{h+k-1} = \tilde{x}) \mathbb{P}_\pi(x_{h+k-1} = \tilde{x} | x_h = x) \\ &= \sum_{\tilde{x} \in \mathcal{X}} \mathbb{P}_\pi(x_{h+k} = x' | x_{h+k-1} = \tilde{x}) \mathbb{P}_\pi(x_{h+k-1} = \tilde{x} | x_h = \phi^*(x)), \end{aligned}$$

638 where the second step uses the fact that π is Markovian and the last step uses the inductive case for
639 $k - 1$. \square

640 **Lemma 2** (Distribution over Pairs). *Let $k \in [K]$, $x \in \text{supp } \rho(X)$, then the distribution $D_k(X' | x)$
641 only depends on $\phi^*(x)$. This allows us to define $D_k(X' | \phi^*(x))$ as this common value. Similarly,
642 the distribution $D_{\text{unf}}(X' | x, k)$ depends only on $\phi^*(x)$ and k . We define this common value as
643 $D_{\text{unf}}(X' | \phi^*(x), k)$.*

644 *Proof.* For any k we have:

$$\begin{aligned}
D_k(X = x, X' = x') &= \frac{1}{H} \sum_{h=1}^H D(x_h = x, x_{h+k} = x') \\
&= \frac{1}{H} \sum_{h=1}^H \sum_{\pi \in \Pi_D} \Theta_\pi \mathbb{P}_\pi(x_h = x, x_{h+k} = x') \\
&= \frac{1}{H} \sum_{h=1}^H \sum_{\pi \in \Pi_D} \Theta_\pi \mathbb{P}_\pi(x_h = x) \mathbb{P}(x_{h+k} = x' \mid x_h = x) \\
&= \frac{1}{H} \sum_{h=1}^H \sum_{\pi \in \Pi_D} \Theta_\pi \mathbb{P}_\pi(x_h = x) \mathbb{P}_\pi(x_{h+k} = x' \mid s_h = \phi^*(x)), \quad (\text{using Lemma 1}) \\
&= \frac{q(x \mid \phi^*(x))}{H} \sum_{h=1}^H \sum_{\pi \in \Pi_D} \Theta_\pi \mathbb{P}_\pi(s_h = \phi^*(x)) \mathbb{P}_\pi(x_{h+k} = x' \mid s_h = \phi^*(x))
\end{aligned}$$

645 The marginal $D_k(X = x)$ is given by:

$$D_k(X = x) = \frac{1}{H} \sum_{h=1}^H \sum_{\pi \in \Pi_D} \Theta_\pi q(x \mid \phi^*(x)) \mathbb{P}_\pi(s_h = \phi^*(x)) = \frac{q(x \mid \phi^*(x))}{H} \sum_{h=1}^H D_k(s_h = \phi^*(x)).$$

646 The conditional $D_k(X' = x' \mid X = x)$ is given by:

$$\begin{aligned}
D_k(X' = x' \mid X = x) &= \frac{D_k(X = x, X' = x')}{D_k(x)} \\
&= \frac{\sum_{h=1}^H \sum_{\pi \in \Pi_D} \Theta_\pi \mathbb{P}_\pi(s_h = \phi^*(x)) \mathbb{P}_\pi(x_{h+k} = x' \mid s_h = \phi^*(x))}{\sum_{h=1}^H D_k(s_h = \phi^*(x))}
\end{aligned}$$

647 Therefore, the conditional $D_k(X' = x' \mid X = x)$ only depends on $\phi^*(x)$, and we define this common
648 value as $D_k(X' = x' \mid s = \phi^*(x))$.

649 The proof for D_{unf} is similar. We can use the property of D_k that we have proven to get:

$$\begin{aligned}
D_{\text{unf}}(X' = x' \mid X = x, k) &= \frac{D_{\text{unf}}(X = x, k, X' = x')}{\sum_{\tilde{x} \in \mathcal{X}} D_{\text{unf}}(X = x, k, X' = \tilde{x})} \\
&= \frac{D_k(X = x, X' = x')}{\sum_{\tilde{x} \in \mathcal{X}} D_k(X = x, X' = \tilde{x})} \\
&= \frac{D_k(X' = x' \mid X = x)}{\sum_{\tilde{x} \in \mathcal{X}} D_k(X' = \tilde{x} \mid X = x)} \\
&= \frac{D_k(X' = x' \mid X = \phi^*(x))}{\sum_{\tilde{x} \in \mathcal{X}} D_k(X' = \tilde{x} \mid X = \phi^*(x))}.
\end{aligned}$$

650 Therefore, $D_{\text{unf}}(X' = x' \mid X = x, k)$ only depends on $\phi^*(x)$. We will define the common values as
651 $D_{\text{unf}}(X' = x' \mid s = \phi^*(x), k)$. □

652 Lemma 2 allows us to define $D_k(x' \mid \phi^*(x))$ and $D_{\text{unf}}(x' \mid \phi^*(x), k)$, as the distribution only
653 depends on the latent state.

654 C.1 Upper Bound for the Forward Model Baseline

655 Let $\mathcal{D}_{\text{for}} = \{(x^{(i)}, k^{(i)}, x'^{(i)})\}_{i=1}^n$ be a pair of iid multi-step observations. We will collect this
656 dataset in one of three ways:

- 657 1. Single step ($k = 1$), in this case we will sample $(x^{(i)}, x'^{(i)}) \sim D_k(X, X')$. As explained
658 before, we can get this sample using the episode data. We save $(x^{(i)}, k, x'^{(i)})$ as our sample.

659 2. Fixed multi-step. We use a fixed $k > 1$, and sample $(x^{(i)}, x'^{(i)}) \sim D_k(X, X')$. We save
 660 $(x^{(i)}, k, x'^{(i)})$ as our sample.

661 3. Variable multi-step. We sample $(x, k, x') \sim D_{\text{unf}}(X, k, X')$ and use it as our sample.

662 We will abstract these three choices using a general notion of $D_{pr} \in \Delta(\mathcal{X} \times [K] \times \mathcal{X})$. In the
 663 first two cases, we assume we have point-mass distribution over k and given this k , we sample
 664 from $D_k(X, X')$. We will assume $(x^{(i)}, k^{(i)}, x'^{(i)}) \sim D_{pr}$. We can create \mathcal{D}_{for} from the dataset
 665 \mathcal{D} of n episodes sampled from D using the sampling procedures explained earlier. Note that as
 666 marginals over both $D_k(X)$ and $D_{\text{unf}}(X)$ is $\rho(X)$, therefore, the marginals over $D_{pr}(X)$ is also
 667 $\rho(X)$. Additionally, we will define $D_{pr}(k)$ as the marginal over k which is either point-mass in the
 668 first two sampling procedures and $\text{Unf}([K])$ in the third procedure.

669 We assume access to two function classes. The first is a decoder class $\Phi_N : \mathcal{X} \rightarrow [N]$ where N is a
 670 given number that satisfies $N \geq |\mathcal{S}|$. The second is a conditional probability class $\mathcal{F} : [N] \times [K] \rightarrow$
 671 $\Delta(\mathcal{X})$.

672 **Assumption 6.** (*Realizability of Φ and \mathcal{F}*) We assume that there exists $\phi^\circ \in \Phi_N$ and $f^\circ \in \mathcal{F}$ such
 673 that $f^\circ(x' | \phi^\circ(x), k) = D_{pr}(x' | x, k) = D_{pr}(x' | \phi^*(x), k)$ for all $(x, k) \sim D_{pr}(\cdot, \cdot)$.

674 This assumption firstly is non-vacuous as $D_{pr}(x' | x) = D_{pr}(x' | \phi^*(x))$, and therefore, we can
 675 apply a bottleneck function ϕ and still assume realizability. For example, we can assume that $\tilde{\phi}$ is the
 676 same as ϕ^* up to the relabeling of its output, and $\tilde{f}(x' | i) = D_{pr}(x' | s)$.

677 Let $\hat{f}, \hat{\phi}$ be the empirical solution to the following maximum likelihood problem.

$$\hat{f}, \hat{\phi} = \arg \max_{f \in \mathcal{F}, \phi \in \Phi_N} \frac{1}{n} \sum_{i=1}^n \ln f(x'^{(i)} | \phi(x^{(i)}), k^{(i)}) \quad (9)$$

678 Note that when k is fixed (we sample from D_k), then information theoretically there is no advantage
 679 of condition on k and it can be dropped from optimization.

680 As we are in a realizable setting (Assumption 6), we can use standard maximum likelihood guarantees
 681 to get the following result.

682 **Proposition 3** (Generalization Bound). *Fix $\delta \in (0, 1)$, then with probability at least $1 - \delta$, we have:*

$$\mathbb{E}_{(x,k) \sim D_{pr}} \left[\left\| D_{pr}(X' | x, k) - \hat{f}(X' | \hat{\phi}(x), k) \right\|_{\text{TV}}^2 \right] \leq \Delta^2(n; \delta),$$

683 where $\Delta^2(n; \delta) = \frac{2}{n} \ln \left(\frac{|\Phi| \cdot |\mathcal{F}|}{\delta} \right)$.

684 For proof see Chapter 7 of Geer [2000].

685 Finally, we assume that the forward modeling objective is expressive to allow the separation of states.
 686 While, this seems like assuming that the objective works, our goal is to establish a formal notion of
 687 the margin so we can verify it later in different settings to see when it holds.

688 **Assumption 7.** (*Forward Modeling Margin*). We assume there exists a $\beta_{\text{for}} \in (0, 1)$ such that:

$$\inf_{s_1, s_2 \in \mathcal{S}, s_1 \neq s_2} \mathbb{E}_{k \sim D_{pr}} [\|D_{pr}(X' | s_1, k) - D_{pr}(X' | s_2, k)\|_{\text{TV}}] \geq \beta_{\text{for}}$$

689 Note that this defines two types of margin depending on D_{pr} . When k is a fixed value, the margin is
 690 given by:

$$\beta_{\text{for}}^{(k)} = \inf_{s_1, s_2 \in \mathcal{S}, s_1 \neq s_2} \|D_{pr}(X' | s_1, k) - D_{pr}(X' | s_2, k)\|_{\text{TV}}.$$

691 When we sample $k \sim \text{Unf}([K])$ then the margin is given by:

$$\beta_{\text{for}}^{(u)} = \inf_{s_1, s_2 \in \mathcal{S}, s_1 \neq s_2} \frac{1}{K} \sum_{k=1}^K \|D_{pr}(X' | s_1, k) - D_{pr}(X' | s_2, k)\|_{\text{TV}}.$$

692 We will use the abstract notion β_{for} for forward margin which will be equal to $\beta_{\text{for}}^{(k)}$ or $\beta_{\text{for}}^{(u)}$ depending
 693 on our sampling procedure. It is easy to see that $\beta_{\text{for}}^{(u)} = \frac{1}{K} \sum_{k=1}^K \beta_{\text{for}}^{(k)}$.

694 We are now ready to state our first main result.

695 **Proposition 4 (Recovering Endogenous State).** Fix $\delta \in (0, 1)$, then with probability at least $1 - \delta$
696 we learn $\hat{\phi}$ that satisfies:

$$\mathbb{P}_{x_1, x_2 \sim \rho} \left(\phi^*(x_1) \neq \phi^*(x_2) \wedge \hat{\phi}(x_1) = \hat{\phi}(x_2) \right) \leq \frac{2\Delta(n, \delta)}{\beta_{\text{for}}}.$$

697 *Proof.* We start with a coupling argument where we sample x_1, x_2 independently from $D_{pr}(X)$
698 which is the same as $\rho(X)$.

$$\begin{aligned} & \mathbb{E}_{x_1, x_2 \sim D_{pr}, k \sim D_{pr}} \left[\mathbf{1} \left\{ \hat{\phi}(x_1) = \hat{\phi}(x_2) \right\} \left\| D_{pr}(X' | x_1, k) - D_{pr}(X' | x_2, k) \right\|_{\text{TV}} \right] \\ & \leq \mathbb{E}_{x_1, x_2 \sim D_{pr}, k \sim D_{pr}} \left[\mathbf{1} \left\{ \hat{\phi}(x_1) = \hat{\phi}(x_2) \right\} \left\| \hat{f}(X' | \hat{\phi}(x_1), k) - D_{pr}(X' | x_1, k) \right\|_{\text{TV}} \right] \\ & \quad + \mathbb{E}_{x_1, x_2 \sim D_{pr}, k \sim D_{pr}} \left[\mathbf{1} \left\{ \hat{\phi}(x_1) = \hat{\phi}(x_2) \right\} \left\| \hat{f}(X' | \hat{\phi}(x_1), k) - D_{pr}(X' | x_2, k) \right\|_{\text{TV}} \right] \end{aligned}$$

699 We bound these two terms separately

$$\begin{aligned} & \mathbb{E}_{x_1, x_2 \sim D_{pr}, k \sim D_{pr}} \left[\mathbf{1} \left\{ \hat{\phi}(x_1) = \hat{\phi}(x_2) \right\} \left\| \hat{f}(X' | \hat{\phi}(x_1), k) - D_{pr}(X' | x_1, k) \right\|_{\text{TV}} \right] \\ & \leq \sqrt{\mathbb{E}_{x_1, x_2 \sim D_{pr}, k \sim D_{pr}} \left[\mathbf{1} \left\{ \hat{\phi}(x_1) = \hat{\phi}(x_2) \right\} \right]} \cdot \sqrt{\mathbb{E}_{x_1, x_2 \sim D_{pr}, k \sim D_{pr}} \left[\left\| \hat{f}(X' | \hat{\phi}(x_1), k) - D_{pr}(X' | x_1, k) \right\|_{\text{TV}}^2 \right]} \\ & = \sqrt{\mathbb{E}_{x_1, x_2 \sim D_{pr}} \left[\mathbf{1} \left\{ \hat{\phi}(x_1) = \hat{\phi}(x_2) \right\} \right]} \cdot \sqrt{\mathbb{E}_{(x, k) \sim D_{pr}} \left[\left\| \hat{f}(X' | \hat{\phi}(x), k) - D_{pr}(X' | x) \right\|_{\text{TV}}^2 \right]} \\ & \leq b \cdot \Delta, \end{aligned}$$

700 where $b = \sqrt{\mathbb{E}_{x_1, x_2 \sim D_{pr}} \left[\mathbf{1} \left\{ \hat{\phi}(x_1) = \hat{\phi}(x_2) \right\} \right]}$ and the second step uses Cauchy-Schwarz inequality.
701 It is straightforward to verify that $b \in [0, 1]$. We bound the second term similarly

$$\begin{aligned} & \mathbb{E}_{x_1, x_2 \sim D_{pr}, k \sim D_{pr}} \left[\mathbf{1} \left\{ \hat{\phi}(x_1) = \hat{\phi}(x_2) \right\} \left\| \hat{f}(X' | \hat{\phi}(x_1), k) - D_{pr}(X' | x_2, k) \right\|_{\text{TV}} \right] \\ & = \mathbb{E}_{x_1, x_2 \sim D_{pr}} \left[\mathbf{1} \left\{ \hat{\phi}(x_1) = \hat{\phi}(x_2) \right\} \left\| \hat{f}(X' | \hat{\phi}(x_2), k) - D_{pr}(X' | x_2, k) \right\|_{\text{TV}} \right] \\ & \leq b \cdot \Delta, \end{aligned}$$

702 where the second step uses the crucial coupling argument that we can replace x_1 with x_2 because of
703 the indicator $\mathbf{1} \left\{ \hat{\phi}(x_1) = \hat{\phi}(x_2) \right\}$, and the last step follows as we reduce it to the first term except
704 we switch the names of x_1 and x_2 . Combining the two upper bounds we get:

$$\mathbb{E}_{x_1, x_2 \sim D_{pr}, k \sim D_{pr}} \left[\mathbf{1} \left\{ \hat{\phi}(x_1) = \hat{\phi}(x_2) \right\} \left\| D_{pr}(X' | x_1, k) - D_{pr}(X' | x_2, k) \right\|_{\text{TV}} \right] \leq 2b \cdot \Delta$$

705 or, equivalently,

$$\mathbb{E}_{x_1, x_2 \sim D_{pr}} \left[\mathbf{1} \left\{ \hat{\phi}(x_1) = \hat{\phi}(x_2) \right\} \underbrace{\mathbb{E}_{k \sim D_{pr}} \left[\left\| D_{pr}(X' | x_1, k) - D_{pr}(X' | x_2, k) \right\|_{\text{TV}} \right]}_{:=\Gamma(x_1, x_2)} \right] \leq 2b \cdot \Delta$$

706 Let $\Gamma(x_1, x_2) = \mathbb{E}_{k \sim D_{pr}} \left[\left\| D_{pr}(X' | x_1, k) - D_{pr}(X' | x_2, k) \right\|_{\text{TV}} \right]$. For any two observations, if
707 $\phi^*(x_1) = \phi^*(x_2)$, then $\left\| D_{pr}(X' | x_1) - D_{pr}(X' | x_2) \right\|_{\text{TV}} = 0$, and therefore, $\Gamma(x_1, x_2) = 0$
708 because of Lemma 2. Otherwise, $\Gamma(x_1, x_2)$ is at least β_{for} , by Assumption 6. Combining these two
709 observations we get:

$$\Gamma(x_1, x_2) \geq \beta_{\text{for}} \mathbf{1} \left\{ \phi^*(x_1) \neq \phi^*(x_2) \right\}$$

710 Combining the previous two inequalities we get:

$$\mathbb{E}_{x_1, x_2 \sim D_{pr}} \left[\mathbf{1} \left\{ \hat{\phi}(x_1) = \hat{\phi}(x_2) \wedge \phi^*(x_1) \neq \phi^*(x_2) \right\} \right] \leq \frac{2b \cdot \Delta}{\beta_{\text{for}}}$$

711 This directly gives

$$\mathbb{P}_{x_1, x_2 \sim D_{pr}} \left(\hat{\phi}(x_1) = \hat{\phi}(x_2) \wedge \phi^*(x_1) \neq \phi^*(x_2) \right) \leq \frac{2b\Delta}{\beta_{\text{for}}} \leq \frac{2\Delta}{\beta_{\text{for}}}.$$

712 The proof is completed by recalling that marginal $D_{pr}(X)$ is the same as $\rho(X)$. \square

713 Proposition 4 shows that the learned $\hat{\phi}$ has one-sided error. If it merges two observations, then with
714 high probability they are not from the same state. As $N = |\mathcal{S}|$, we will show below that the reverse is
715 also true.

716 **Theorem 5.** *If $N = |\mathcal{S}|$, then there exists a bijection $\alpha : [N] \rightarrow \mathcal{S}$ such that for any $s \in \mathcal{S}$ we have:*

$$\mathbb{P}_{x \sim q(\cdot | s)} \left(\hat{\phi}(x) = \alpha(s) \mid \phi^*(x) = s \right) \geq 1 - \frac{4N^3 H^2 \Delta}{\eta_{\text{min}}^2 \beta_{\text{for}}},$$

717 provided $\Delta < \frac{\eta_{\text{min}}^2 \beta_{\text{for}}}{N^2 H^2}$.

718 *Proof.* We define a few shorthand below for any $j \in [N]$ and $\tilde{s} \in \mathcal{S}$

$$\begin{aligned} \mathbb{P}(j, \tilde{s}) &= \mathbb{P}_{x \sim \rho} \left(\hat{\phi}(x) = j \wedge \phi^*(x) = \tilde{s} \right) \\ \rho(j) &= \mathbb{P}_{x \sim \rho} \left(\hat{\phi}(x) = j \right) \\ \rho(\tilde{s}) &= \mathbb{P}_{x \sim \rho} \left(\phi^*(x) = \tilde{s} \right). \end{aligned}$$

719 It is easy to verify that $\mathbb{P}(j, \tilde{s})$ is a joint distribution with $\rho(j)$ and $\rho(\tilde{s})$ as its marginals.

720 Fix $i \in [N]$ and $s \in \mathcal{S}$.

$$\begin{aligned} &\mathbb{P}_{x_1, x_2 \sim \rho} \left(\hat{\phi}(x_1) = \hat{\phi}(x_2) \wedge \phi^*(x_1) \neq \phi^*(x_2) \right) \\ &= \mathbb{P}_{x_1, x_2 \sim \rho} \left(\bigcup_{\tilde{s} \in \mathcal{S}, j \in [N]} \left\{ \hat{\phi}(x_1) = j \wedge \hat{\phi}(x_2) = j \wedge \phi^*(x_1) = \tilde{s} \wedge \phi^*(x_2) \neq \tilde{s} \right\} \right) \\ &\geq \mathbb{P}_{x_1, x_2 \sim \rho} \left(\hat{\phi}(x_1) = i \wedge \hat{\phi}(x_2) = i \wedge \phi^*(x_1) = s \wedge \phi^*(x_2) \neq s \right) \\ &= \mathbb{P}_{x_1 \sim \rho} \left(\hat{\phi}(x_1) = i \wedge \phi^*(x_1) = s \right) \mathbb{P}_{x_2 \sim \rho} \left(\hat{\phi}(x_2) = i \wedge \phi^*(x_2) \neq s \right) \\ &= \mathbb{P}_{x \sim \rho} \left(\hat{\phi}(x) = i \wedge \phi^*(x) = s \right) \left(\sum_{s' \in \mathcal{S}} \mathbb{P}_{x \sim \rho} \left(\hat{\phi}(x) = i \wedge \phi^*(x) = s' \right) - \mathbb{P}_{x \sim \rho} \left(\hat{\phi}(x) = i \wedge \phi^*(x) = s \right) \right) \\ &= \mathbb{P}(i, s) \left(\sum_{s' \in \mathcal{S}} \mathbb{P}(i, s') - \mathbb{P}(i, s) \right) \\ &= \mathbb{P}(i, s) (\rho(i) - \mathbb{P}(i, s)). \end{aligned}$$

721 Combining this with Proposition 4, we get:

$$\forall i \in [N], s \in \mathcal{S}, \quad \mathbb{P}(i, s) (\rho(i) - \mathbb{P}(i, s)) \leq \Delta' := \frac{2\Delta}{\beta_{\text{for}}}$$

722 where we have used a shorthand $\Delta' = 2\Delta/\beta_{\text{for}}$. We define a mapping $\alpha : \mathcal{S} \rightarrow [N]$ where for any
723 $s \in \mathcal{S}$:

$$\alpha(s) = \arg \max_{j \in [N]} \mathbb{P}(j, s) \tag{10}$$

724 We immediately have:

$$\mathbb{P}(\alpha(s), s) = \max_{j \in [N]} \mathbb{P}(j, s) \geq \frac{1}{N} \sum_{j=1}^N \mathbb{P}(j, s) = \frac{1}{N} \rho(s) \geq \frac{\eta_{\text{min}}}{NH}, \tag{11}$$

725 where we use the fact that max is greater than average in the first inequality, and Equation 8. Further,
 726 for every $s \in \mathcal{S}$, we have:

$$\mathbb{P}(\alpha(s), s) (\rho(\alpha(s)) - \mathbb{P}(\alpha(s), s)) \leq \Delta'.$$

727 Plugging the lower bound $\mathbb{P}(\alpha(s), s) \geq \frac{\eta_{\min}}{NH}$, we get:

$$\mathbb{P}(\alpha(s), s) \geq \rho(\alpha(s)) - \frac{NH\Delta'}{\eta_{\min}}. \quad (12)$$

728 We now show that if $\Delta' < \frac{\eta_{\min}^2}{2N^2H^2}$, then $\alpha(s)$ is a bijection. Let s_1 and s_2 be such that $\alpha(s_1) =$
 729 $\alpha(s_2) = i$. Then using the above Equation 12 we get $\mathbb{P}(i, s_1) \geq \rho(i) - \frac{NH\Delta'}{\eta_{\min}}$ and $\mathbb{P}(i, s_2) \geq$
 730 $\rho(i) - \frac{NH\Delta'}{\eta_{\min}}$. We have:

$$\rho(i) = \sum_{\tilde{s} \in \mathcal{S}} \mathbb{P}(i, \tilde{s}) \geq \mathbb{P}(i, s_1) + \mathbb{P}(i, s_2) \geq 2\rho(i) - \frac{2NH\Delta'}{\eta_{\min}}$$

731 This implies $\frac{2NH\Delta'}{\eta_{\min}} \geq \rho(i)$ but as $\rho(i) = \rho(\alpha(s_1)) \geq \mathbb{P}(\alpha(s_1), s_1) \geq \frac{\eta_{\min}}{NH}$ (Equation 11), we get
 732 $\frac{2NH\Delta'}{\eta_{\min}} \geq \frac{\eta_{\min}}{NH}$ or $\Delta' \geq \frac{\eta_{\min}^2}{2N^2H^2}$. However, as we assume that $\Delta' < \frac{\eta_{\min}^2}{2N^2H^2}$, therefore, this is a
 733 contradiction. This implies $\alpha(s_1) \neq \alpha(s_2)$ for any two different states s_1 and s_2 . Since we assume
 734 $|N| = |\mathcal{S}|$, this implies α is a bijection.

735 Fix $s \in \mathcal{S}$ and let $i \neq \alpha(s)$. As α is a bijection, let $\tilde{s} = \alpha^{-1}(i)$, we can show that $\mathbb{P}(i, s)$ is small:

$$\mathbb{P}(i, s) \leq \rho(i) - \mathbb{P}(i, \tilde{s}) = \rho(\alpha(\tilde{s})) - \mathbb{P}(\alpha(\tilde{s}), \tilde{s}) \leq \frac{NH\Delta'}{\eta_{\min}} \quad (13)$$

736 where we use $s \neq \tilde{s}$ and Equation 12.

737 This allows us to show that $\mathbb{P}(\alpha(s) | s)$ is high as follows:

$$\begin{aligned} \mathbb{P}(\alpha(s) | s) &= \frac{\mathbb{P}(\alpha(s), s)}{\rho(s)} = \frac{\mathbb{P}(\alpha(s), s)}{\mathbb{P}(\alpha(s), s) + \sum_{i=1, i \neq \alpha(s)}^N \mathbb{P}(i, s)} \\ &\geq \frac{\mathbb{P}(\alpha(s), s)}{\rho(\alpha(s)) + \frac{N^2H\Delta'}{\eta_{\min}}} \\ &\geq \frac{\rho(\alpha(s)) - \frac{NH\Delta'}{\eta_{\min}}}{\rho(\alpha(s)) + \frac{N^2H\Delta'}{\eta_{\min}}} \\ &= 1 - \frac{\left(\frac{N^2H\Delta'}{\eta_{\min}} + \frac{NH\Delta'}{\eta_{\min}} \right)}{\rho(\alpha(s)) + \frac{N^2H\Delta'}{\eta_{\min}}} \\ &\geq 1 - \frac{2N^2H^2\Delta'}{\eta_{\min}\rho(\alpha(s))} \\ &\geq 1 - \frac{2N^3H^2\Delta'}{\eta_{\min}^2}, \end{aligned}$$

738 where the first inequality uses Equation 13 and $\rho(\alpha(s)) \geq \mathbb{P}(\alpha(s), s)$, second inequality uses
 739 Equation 12, and the last step uses $\rho(\alpha(s)) \geq \mathbb{P}(\alpha(s), s) \geq \frac{\eta_{\min}}{NH}$.

740 The proof is completed by noting that:

$$\mathbb{P}_{x \sim q(\cdot | s)}(\hat{\phi}(x) = \alpha(s)) = \mathbb{P}_{x \sim \rho}(\hat{\phi}(x) = \alpha(s) | \phi^*(x) = s) = \mathbb{P}(\alpha(s) | s).$$

741

□

742 Let \mathcal{A} be a PAC RL algorithm for tabular MDPs. We assume that this algorithm’s sample complexity
 743 is given by $n_{\text{samp}}(S, A, H, \epsilon, \delta)$ where S and A are the size of the state space and action space of
 744 the tabular MDP, H is the horizon, and (ϵ, δ) are the typical PAC RL hyperparameters denoting
 745 tolerance and failure probability. Formally, the algorithm \mathcal{A} interacts with a tabular MDP \mathbb{M} for
 746 $n_{\text{samp}}(S, A, H, \epsilon, \delta)$ episodes and outputs a policy $\hat{\varphi} : \mathcal{S} \times [H] \rightarrow \mathcal{A}$ such that with probability at
 747 least $1 - \delta$ we have:

$$\sup_{\varphi \in \Psi_{\text{all}}} V_{\mathbb{M}}(\varphi) - V_{\mathbb{M}}(\hat{\varphi}) \leq \epsilon,$$

748 where Ψ_{all} is the space of all policies of the type $\mathcal{S} \times [H] \rightarrow \mathcal{A}$.

749 We assume that we are given knowledge of the desired (ϵ, δ) hyperparameters in the downstream RL
 750 task during the representation pre-training phase so we can use the right amount of data.

751 **Induced Finite MDP.** The latent MDP inside a block MDP is a tabular MDP with state space \mathcal{S} ,
 752 action space \mathcal{A} , horizon H , transition dynamics T , reward function R , and a start state distribution
 753 of μ . If we directly had access to this latent MDP, say via the true decoding function ϕ^* , then we
 754 can apply the algorithm \mathcal{A} and learn the optimal latent policy φ^* which we can couple with ϕ^*
 755 and learn the optimal observation-based policy. Formally, we write this observation-based policy as
 756 $\varphi \circ \phi^* : \mathcal{X} \times [H] \rightarrow \mathcal{A}$ given by $\varphi(\phi^*(x), h)$. We don’t have access to ϕ^* , but we have access to
 757 $\hat{\phi}$ that with high probability for a given x outputs a state which is same as $\phi^*(x)$ up to the learned
 758 α -bijection. We, therefore, define the induced MDP $\hat{\mathbb{M}}$ as the finite MDP with state space $\hat{\mathcal{S}}$, action
 759 space \mathcal{A} , transition function \hat{T} , reward function \hat{R} and start state distribution $\hat{\mu}$. These same as the
 760 latent Block MDP but where the true state s is replaced by $\alpha(s)$. It is this induced $\hat{\mathbb{M}}$ that the tabular
 761 MDP algorithm \mathcal{A} will see with high probability.

762 **Proposition 6 (PAC RL Bound).** *Let \mathcal{A} be a PAC RL algorithm for tabular MDPs and n_{samp} is its
 763 sample complexity. Let $\hat{\phi} : \mathcal{X} \rightarrow [N]$ be a decoder pre-trained using video data and $\alpha : \mathcal{S} \rightarrow [N]$ is
 764 a bijection such that:*

$$\forall s \in \mathcal{S}, \quad \mathbb{P}_{x \sim q(\cdot|s)} \left(\hat{\phi}(x) = \alpha(s) \right) \geq 1 - \vartheta,$$

765 then let $\hat{\varphi}$ be the policy returned by \mathcal{A} on the tabular MDP induced by $\hat{\phi}(x)$. Then we have with
 766 probability at least $1 - \delta - n_{\text{samp}}(S, A, H, \epsilon, \delta)H\vartheta$:

$$\sup_{\pi \in \Pi} V(\pi) - V(\varphi \circ \hat{\phi}) \leq \epsilon + 2H^2\vartheta$$

767 *Proof.* The algorithm runs for $n_{\text{samp}}(S, A, H, \epsilon, \delta)$ episodes. This implies the agent visits
 768 $n_{\text{samp}}(S, A, H, \epsilon, \delta)H$ many latent states. If the decoder maps every such state s to the correct
 769 permutation $\alpha(s)$, then the tabular MDP algorithm is running as if it ran on the induced MDP $\hat{\mathbb{M}}$. The
 770 probability of failure is bounded by $n_{\text{samp}}(S, A, H, \epsilon, \delta)H\vartheta$ as all these failures are independent
 771 given the state. Further, the failure probability of the tabular MDP algorithm itself is δ . This leads to
 772 the total failure probability of $\delta + n_{\text{samp}}(S, A, H, \epsilon, \delta)H\vartheta$.

773 Let Π be the set of observation-based policies we are competing with and which includes the optimal
 774 observation-based policy π^* . We can write $\sup_{\pi \in \Pi} V(\pi) = V_{\hat{\mathbb{M}}}(\varphi^*)$ where we use the subscript $\hat{\mathbb{M}}$
 775 to denote that the latent policy is running in the induced MDP $\hat{\mathbb{M}}$. Further, for any latent policy φ we
 776 have $V(\varphi \circ \alpha \circ \phi^*) = V_{\hat{\mathbb{M}}}(\varphi)$ as the decoder $\alpha \circ \phi^* : x \mapsto \alpha(\phi^*(x))$ give me access to the true state
 777 of the induced MDP $\hat{\mathbb{M}}$. Then with probability at least $1 - \delta$, we have:

$$V_{\hat{\mathbb{M}}}(\varphi^*) - V_{\hat{\mathbb{M}}}(\hat{\varphi}) \leq \epsilon$$

778 This allows us to bound the sub-optimality of the learned observation-based policy $\hat{\varphi} \circ \hat{\phi}$ as:

$$\begin{aligned} \sup_{\pi \in \Pi} V(\pi) - V(\hat{\varphi} \circ \hat{\phi}) &= V(\varphi^* \circ \alpha \circ \phi^*) - V(\hat{\varphi} \circ \alpha \circ \phi^*) + V(\hat{\varphi} \circ \alpha \circ \phi^*) - V(\hat{\varphi} \circ \hat{\phi}) \\ &= V_{\hat{\mathbb{M}}}(\varphi^*) - V_{\hat{\mathbb{M}}}(\hat{\varphi}) + V(\hat{\varphi} \circ \alpha \circ \phi^*) - V(\hat{\varphi} \circ \hat{\phi}) \\ &\leq \epsilon + V(\hat{\varphi} \circ \alpha \circ \phi^*) - V(\hat{\varphi} \circ \hat{\phi}) \end{aligned}$$

779 Here we use $\hat{\varphi} \circ \alpha \circ \phi^*$ to denote an observation-based policy that takes action as $\hat{\varphi}(\alpha(\phi^*(x)), h)$.

780 We bound $V(\hat{\varphi} \circ \alpha \circ \phi^*) - V(\hat{\varphi} \circ \hat{\phi})$ below. Let $\mathcal{E}_h = \{\hat{\phi}(x_h) = \alpha(\phi^*(x_h))\}$ and $\mathcal{E} = \bigcap_{h=1}^H \mathcal{E}_h$ be
 781 two events. We have $\mathbb{P}(\mathcal{E}_h) \geq 1 - \vartheta$. Further, using union bound we have $\mathbb{P}(\mathcal{E}^c) = \mathbb{P}(\bigcup_{h=1}^H \mathcal{E}_h^c) \leq$
 782 $\sum_{h=1}^H \mathbb{P}(\mathcal{E}_h^c) \leq H\vartheta$.

783 We first prove an upper bound on $V(\hat{\varphi} \circ \alpha \circ \phi^*)$:

$$\begin{aligned} V(\hat{\varphi} \circ \alpha \circ \phi^*) &= \mathbb{E}_{\hat{\varphi} \circ \alpha \circ \phi^*} \left[\sum_{h=1}^H r_h \right] \\ &= \mathbb{E}_{\hat{\varphi} \circ \alpha \circ \phi^*} \left[\sum_{h=1}^H r_h \mid \mathcal{E} \right] \mathbb{P}_{\hat{\varphi} \circ \alpha \circ \phi^*}(\mathcal{E}) + \mathbb{E}_{\hat{\varphi} \circ \alpha \circ \phi^*} \left[\sum_{h=1}^H r_h \mid \mathcal{E}^c \right] \mathbb{P}_{\hat{\varphi} \circ \alpha \circ \phi^*}(\mathcal{E}^c) \\ &\leq \mathbb{E}_{\hat{\varphi} \circ \alpha \circ \phi^*} \left[\sum_{h=1}^H r_h \mid \mathcal{E} \right] + H^2\vartheta \\ &= \mathbb{E}_{\hat{\varphi} \circ \hat{\phi}} \left[\sum_{h=1}^H r_h \mid \mathcal{E} \right] + H^2\vartheta \end{aligned}$$

784 Here we have used the fact that value of any policy is in $[0, H]$ since the horizon is H and the rewards
 785 are in $[0, 1]$.

786 We next prove a lower bound on $V(\hat{\varphi} \circ \hat{\phi})$:

$$\begin{aligned} V(\hat{\varphi} \circ \hat{\phi}) &= \mathbb{E}_{\hat{\varphi} \circ \hat{\phi}} \left[\sum_{h=1}^H r_h \right] \\ &= \mathbb{E}_{\hat{\varphi} \circ \hat{\phi}} \left[\sum_{h=1}^H r_h \mid \mathcal{E} \right] \mathbb{P}_{\hat{\varphi} \circ \hat{\phi}}(\mathcal{E}) + \mathbb{E}_{\hat{\varphi} \circ \hat{\phi}} \left[\sum_{h=1}^H r_h \mid \mathcal{E}^c \right] \mathbb{P}_{\hat{\varphi} \circ \hat{\phi}}(\mathcal{E}^c) \\ &\geq \mathbb{E}_{\hat{\varphi} \circ \hat{\phi}} \left[\sum_{h=1}^H r_h \mid \mathcal{E} \right] \mathbb{P}_{\hat{\varphi} \circ \hat{\phi}}(\mathcal{E}) \\ &\geq \mathbb{E}_{\hat{\varphi} \circ \hat{\phi}} \left[\sum_{h=1}^H r_h \mid \mathcal{E} \right] - \mathbb{E}_{\hat{\varphi} \circ \hat{\phi}} \left[\sum_{h=1}^H r_h \mid \mathcal{E} \right] H\vartheta \\ &\geq \mathbb{E}_{\hat{\varphi} \circ \hat{\phi}} \left[\sum_{h=1}^H r_h \mid \mathcal{E} \right] - H^2\vartheta \end{aligned}$$

787 Combining the two upper bounds we get:

$$V(\hat{\varphi} \circ \alpha \circ \phi^*) - V(\hat{\varphi} \circ \hat{\phi}) \leq \mathbb{E}_{\hat{\varphi} \circ \hat{\phi}} \left[\sum_{h=1}^H r_h \mid \mathcal{E} \right] + H^2\vartheta - \mathbb{E}_{\hat{\varphi} \circ \hat{\phi}} \left[\sum_{h=1}^H r_h \mid \mathcal{E} \right] + H^2\vartheta \leq 2H^2\vartheta$$

788 Therefore, with probability at least $1 - \delta - n_{\text{samp}}(S, A, H, \epsilon, \delta)H\vartheta$, learn a policy $\hat{\varphi} \circ \hat{\phi}$ such that:

$$\sup_{\pi \in \Pi} V(\pi) - V(\hat{\varphi} \circ \hat{\phi}) \leq \epsilon + 2H^2\vartheta.$$

789

□

790 **Theorem 7** (Wrapping up the proof.). Fix $\epsilon_o > 0$ and $\delta_o \in (0, 1)$ and let \mathcal{A} be any PAC RL algorithm
 791 for tabular MDPs with sample complexity $n_{\text{samp}}(S, A, H, \epsilon, \delta)$. If n satisfies:

$$n = \mathcal{O} \left(\left\{ \frac{N^4 H^4}{\eta^4 \min^* \beta^2} + \frac{N^6 H^8}{\epsilon_o^2 \eta^4 \min^* \beta^2} + \frac{N^6 H^6 n_{\text{samp}}^2(S, A, H, \epsilon_o/2, \delta_o/4)}{\delta_o^2 \eta^4 \min^* \beta^2} \right\} \ln \left(\frac{|\mathcal{F}| |\Phi|}{\delta_o} \right) \right),$$

792 then forward modeling learns a decoder $\hat{\phi} : \mathcal{X} \rightarrow N$. Further, running \mathcal{A} on the tabular MDP with
 793 induced by $\hat{\phi}$ with hyperparameters $\epsilon = \epsilon_o/2$, $\delta = \delta_o/4$, returns a latent policy $\hat{\varphi}$. Then there exists

794 a bijective mapping $\alpha : \mathcal{S} \rightarrow [|\mathcal{S}|]$ such that with probability at least $1 - \delta$ we have:

$$\forall s \in \mathcal{S}, \quad \mathbb{P}_{x \sim q(\cdot|s)} \left(\hat{\phi}(x) = \alpha(s) \mid \phi^*(x) = s \right) \geq 1 - \frac{4N^3 H^2 \Delta}{\eta_{\min}^2 \beta_{\text{for}}},$$

795 and

$$V(\pi^*) - V(\hat{\phi} \circ \hat{\phi}) \leq \epsilon_{\circ}$$

796 Further, the amount of online interactions in the downstream RL is given by
797 $n_{\text{samp}}(S, A, H, \epsilon_{\circ}/2, \delta_{\circ}/4)$ and doesn't scale with $\ln|\Phi|$.

798 *Proof.* We showed in Theorem 5 that we learn a $\hat{\phi}$ such that:

$$\mathbb{P}_{x \sim q(\cdot|s)} \left(\hat{\phi}(x) = \alpha(s) \mid \phi^*(x) = s \right) \geq 1 - \frac{4N^3 H^2 \Delta}{\eta_{\min}^2 \beta_{\text{for}}},$$

799 provided $\Delta < \frac{\eta_{\min}^2 \beta_{\text{for}}}{N^2 H^2}$.

800 Let $\vartheta = \frac{4N^3 H^2 \Delta}{\eta_{\min}^2 \beta_{\text{for}}}$. Then from Proposition 6 we learn a $\hat{\phi}$ such that:

$$V(\pi^*) - V(\hat{\phi} \circ \hat{\phi}) \leq \epsilon + 2H^2 \vartheta,$$

801 with probability at least $1 - \delta - n_{\text{samp}}(S, A, H, \epsilon, \delta)H\vartheta$. The failure probability $\delta -$
802 $n_{\text{samp}}(S, A, H, \epsilon, \delta)H\vartheta$ was when condition in Theorem 5 holds which holds with δ probability.
803 Hence, total failure probability is:

$$2\delta + n_{\text{samp}}(S, A, H, \epsilon, \delta)H\vartheta.$$

804 We set δ both in our representation learning analysis and in PAC RL to $\delta_{\circ}/4$. We also set ϵ in the
805 PAC RL algorithm to $\epsilon_{\circ}/2$. This means the PAC RL algorithm runs for $n_{\text{samp}}(S, A, H, \epsilon_{\circ}/2, \delta_{\circ}/4)$
806 episodes.

807 We enforce $\vartheta \leq \frac{\delta_{\circ}}{2n_{\text{samp}}(S, A, H, \epsilon_{\circ}/2, \delta_{\circ}/4)H}$. Then the total failure probability becomes:

$$2\delta_{\circ}/4 + \delta_{\circ}/4 + \delta_{\circ}/2 \leq \delta_{\circ}$$

808 We also enforce $2H^2\vartheta \leq \epsilon_{\circ}/2$. The sub-optimality of the PAC RL policy is given by:

$$\epsilon_{\circ}/2 + \epsilon_{\circ}/2 \leq \epsilon_{\circ}$$

809 This gives us our derived PAC RL bound.

810 We now accumulate all conditions:

$$\begin{aligned} \Delta &= \sqrt{\frac{2}{n} \ln \left(\frac{4|\mathcal{F}||\Phi|}{\delta_{\circ}} \right)} \\ \vartheta &= \frac{4N^3 H^2 \Delta}{\eta_{\min}^2 \beta_{\text{for}}} \\ \Delta &< \frac{\eta_{\min}^2 \beta_{\text{for}}}{N^2 H^2} \\ \vartheta &\leq \frac{\delta_{\circ}}{2n_{\text{samp}}(S, A, H, \epsilon_{\circ}/2, \delta_{\circ}/4)H} \\ 2H^2\vartheta &\leq \epsilon_{\circ}/2 \end{aligned}$$

811 This simplifies to

$$\begin{aligned} \Delta &\leq \frac{\eta_{\min}^2 \beta_{\text{for}}}{N^2 H^2} \\ \Delta &\leq \frac{\delta_{\circ} \eta_{\min}^2 \beta_{\text{for}}}{8N^3 H^3 n_{\text{samp}}(S, A, H, \epsilon_{\circ}/2, \delta_{\circ}/4)} \\ \Delta &\leq \frac{\epsilon_{\circ} \eta_{\min}^2 \beta_{\text{for}}}{16N^3 H^4} \end{aligned}$$

812 Or,

$$n = \mathcal{O} \left(\left\{ \frac{N^4 H^4}{\eta_{\text{for}}^4 \beta_{\text{for}}^2} + \frac{N^6 H^8}{\epsilon_{\circ}^2 \eta_{\text{for}}^4 \beta_{\text{for}}^2} + \frac{N^6 H^6 n_{\text{samp}}^2 (S, A, H, \epsilon_{\circ}/2, \delta_{\circ}/4)}{\delta_{\circ}^2 \eta_{\text{for}}^4 \beta_{\text{for}}^2} \right\} \ln \left(\frac{|\mathcal{F}| |\Phi|}{\delta_{\circ}} \right) \right)$$

813 This completes the proof. \square

814 C.2 Upper Bound for the Temporal Contrastive Approach

815 We first convert our video dataset \mathcal{D} into a dataset suitable for contrastive learn-
 816 ing. We first split the datasets into $\lfloor n/2 \rfloor$ pairs of videos. For each video pair
 817 $\left\{ \left(x_1^{(2l)}, x_2^{(2l)}, \dots, x_H^{(2k)} \right), \left(x_1^{(2l+1)}, x_2^{(2l+1)}, \dots, x_H^{(2l+1)} \right) \right\}$, we create a tuple (x, x', k, z) where
 818 $z \in \{0, 1\}$ as follows. As in forward modeling, we will either use a fixed value of k , or sample
 819 $k \in \text{Unf}([K])$. We denote this general distribution over k by $\omega \in \Delta([K])$ which is either point mass,
 820 or $\text{Unf}([K])$. We sample $k \sim \omega$ and $z \sim \text{Unf}(\{0, 1\})$ and $h \in \text{Unf}([H])$. We set $x = x_h^{(2l)}$. If $z = 1$,
 821 then we set $x' = x_{h+k}^{(2l)}$, otherwise, we sample $h' \sim \text{Unf}(\{0, 1\})$ and select $x' = x_{h'}^{(2l)}$. This way, we
 822 collect a dataset $\mathcal{D}_{\text{cont}}$ of $\lfloor n/2 \rfloor$ tuples (x, k, x', z) . We view a tuple (x, k, x', z) as a *real observation*
 823 *pair* when $z = 1$, and a *fake observation pair* when $z = 0$. Note that our sampling process leads to
 824 all data points being iid.

825 We define the distribution $D_{\text{cont}}(X, k, X', Z)$ as the distribution over (x, k, x', z) . We can express
 826 this distribution as:

$$\begin{aligned} D_{\text{cont}}(X = x, k, X' = x', Z = 1) &= \frac{\omega(k)}{2H} \sum_{h=1}^H D(x = x_h, x' = x_{h+k}) \\ &= \frac{\omega(k)}{2} \rho(x) D(x_{k+1} = x' \mid x_1 = x) \\ D_{\text{cont}}(X = x, X' = x', Z = 0) &= \frac{\omega(k)}{2H^2} \sum_{h=1}^H D(x = x_h) \sum_{h'=1}^H D(x' = x_{h'}) \\ &= \frac{\omega(k)}{2} \rho(x) \rho(x') \end{aligned}$$

827 where we use the time homogeneity of D and definition of ρ . We will use a shorthand to denote
 828 $D(x_{k+1} = x' \mid x_1 = x)$ as $D(x' \mid x, k)$ in this analysis. It is easy to verify that $D(x' \mid x, k) =$
 829 $D(x' \mid \phi^*(x), k)$. The marginal distribution $D_{\text{cont}}(x, k, x')$ is given by:

$$D_{\text{cont}}(x, k, x') = \frac{\omega(k) \rho(x)}{2} (D(x' \mid x, k) + \rho(x')) \quad (14)$$

830 Note that $D_{\text{cont}}(X)$ is the same as $\rho(X)$.

831 We will use D_{cont} for any marginal and conditional distribution derived from $D_{\text{cont}}(X, k, X', Z)$.
 832 We assume a model class $\mathcal{G} : \mathcal{X} \times [K] \times [N] \rightarrow [0, 1]$ that we use for solving the prediction
 833 problem. We will also reuse the decoder class $\phi : \mathcal{X} \rightarrow [N]$ that we defined earlier, and we will
 834 assume that $N = |\mathcal{S}|$. This can be relaxed by doing clustering or working with a different induced
 835 MDP (e.g., see the clustering algorithm in Misra et al. [2020]). However, this is not the main point of
 836 the analysis.

837 We define the expected risk minimizer of the squared loss problem below:

$$\hat{g}, \hat{\phi} = \arg \min_{g \in \mathcal{G}, \phi \in \Phi} \frac{1}{\lfloor n/2 \rfloor} \sum_{i=1}^{\lfloor n/2 \rfloor} \left(g(\phi(x^{(i)}), k^{(i)}, x'^{(i)}) - z^{(i)} \right)^2 \quad (15)$$

838 We express the Bayes classifier of this problem below:

839 **Lemma 3** (Bayes Classifier). *The Bayes classifier of the problem posed in Equation 15 is given by*
 840 $D_{\text{cont}}(z = 1 \mid x, k, x')$ *which satisfies:*

$$D_{\text{cont}}(z = 1 \mid x, k, x') = \frac{D(\phi^*(x') \mid \phi^*(x), k)}{D(\phi^*(x') \mid \phi^*(x), k) + \rho(\phi^*(x'))}.$$

841 *Proof.* We can express the Bayes classifier as:

$$\begin{aligned}
D_{\text{cont}}(z = 1 | x, k, x') &= \frac{D_{\text{cont}}(x, k, x', z = 1)}{D_{\text{cont}}(x, k, x', z = 1) + D_{\text{cont}}(x, k, x', z = 0)} \\
&= \frac{\omega(k)/2\rho(x)D(x' | x)}{\omega(k)/2\rho(x)D(x' | x) + \omega(k)/2\rho(x)\rho(x')} \\
&= \frac{D(x' | x, k)}{D(x' | x, k) + \rho(x')} \\
&= \frac{D(x' | \phi^*(x), k)}{D(x' | \phi^*(x), k) + \rho(x')} \\
&= \frac{q(x' | \phi^*(x))D(\phi^*(x') | \phi^*(x), k)}{q(x' | \phi^*(x))D(\phi^*(x') | \phi^*(x), k) + q(x' | \phi^*(x))\rho(\phi^*(x'))} \\
&= \frac{D(\phi^*(x') | \phi^*(x), k)}{D(\phi^*(x') | \phi^*(x), k) + \rho(\phi^*(x'))}.
\end{aligned}$$

842

□

843 **Assumption 8 (Realizability).** *There exists $g^* \in \mathcal{G}$ and $\phi^\circ \in \Phi$ such that for all $(x, k, x') \in$*
844 *$\text{supp } D_{\text{cont}}(X, k, X')$, we have $D_{\text{cont}}(z = 1 | x, k, x') = g^*(\phi^\circ(x), k, x')$.*

845 We will use the shorthand to denote $g^*(x, k, x') = g^*(\phi^\circ(x), k, x')$.

846 As before, we start with typical square loss guarantees in the realizable setting.

847 **Theorem 8.** *Fix $\delta \in (0, 1)$. Under realizability (Assumption 8), the ERM solution of $\hat{f}, \hat{\phi}$ in Eq. (15)*
848 *satisfies:*

$$\mathbb{E}_{(x, k, x') \sim D_{\text{cont}}} \left[\left(\hat{g}(\hat{\phi}(x), k, x') - g^*(x, k, x') \right)^2 \right] \leq \Delta_{\text{cont}}^2 = \frac{2}{n} \ln \frac{|\mathcal{G}| \cdot |\Phi|}{\delta}$$

849 For proof see Proposition 12 in Misra et al. [2020].

850 We will prove a coupling result similar to the case for forward modeling. However, to do this, we
851 need to define a coupling distribution:

$$D_{\text{coup}}(X_1 = x_1, X_2 = x_2, k, X' = x') = \omega(k)D_{\text{cont}}(X = x_1)D_{\text{cont}}(X = x_2)D_{\text{cont}}(X' = x')$$

852 We will derive a useful importance ratio bound.

$$\frac{D_{\text{coup}}(x_1, k, x')}{D_{\text{cont}}(x_1, k, x')} = \frac{2\rho(x_1)\rho(x')}{\rho(x_1)D(x' | x_1, k) + \rho(x_1)\rho(x')} \leq 2 \tag{16}$$

853 We now prove an analogous result to Proposition 4.

854 **Theorem 9 (Coupling for Temporal Contrastive Learning).** *With probability at least $1 - \delta$ we have:*

$$\mathbb{E}_{(x_1, x_2, k, x') \sim D_{\text{coup}}} \left[\mathbf{1} \left\{ \hat{\phi}(x_1) = \hat{\phi}(x_2) \right\} \left| g^*(x_1, k, x') - g^*(x_2, k, x') \right| \right] < 4\Delta_{\text{cont}}(n, \delta)$$

855 *Proof.* We start with triangle inequality:

$$\begin{aligned}
&\mathbb{E}_{(x_1, x_2, k, x') \sim D_{\text{coup}}} \left[\mathbf{1} \left\{ \hat{\phi}(x_1) = \hat{\phi}(x_2) \right\} \left| g^*(x_1, k, x') - g^*(x_2, k, x') \right| \right] \\
&\leq \mathbb{E}_{(x_1, x_2, k, x') \sim D_{\text{coup}}} \left[\mathbf{1} \left\{ \hat{\phi}(x_1) = \hat{\phi}(x_2) \right\} \left| g^*(x_1, k, x') - \hat{g}(\hat{\phi}(x_1), k, x') \right| \right] + \\
&\quad \mathbb{E}_{(x_1, x_2, k, x') \sim D_{\text{coup}}} \left[\mathbf{1} \left\{ \hat{\phi}(x_1) = \hat{\phi}(x_2) \right\} \left| \hat{g}(\hat{\phi}(x_1), k, x') - g^*(x_2, k, x') \right| \right]
\end{aligned}$$

856 We bound the first term as:

$$\begin{aligned}
& \mathbb{E}_{(x_1, x_2, k, x') \sim D_{\text{coup}}} \left[\mathbf{1} \left\{ \hat{\phi}(x_1) = \hat{\phi}(x_2) \right\} \left| g^*(x_1, k, x') - \hat{g}(\hat{\phi}(x_1), k, x') \right| \right] \\
& \leq \underbrace{\sqrt{\mathbb{E}_{(x_1, x_2, k, x') \sim D_{\text{coup}}} \left[\mathbf{1} \left\{ \hat{\phi}(x_1) = \hat{\phi}(x_2) \right\} \right]}}_{:=b} \cdot \sqrt{\mathbb{E}_{(x_1, x_2, k, x') \sim D_{\text{coup}}} \left[\left| g^*(x_1, k, x') - \hat{g}(\hat{\phi}(x_1), k, x') \right|^2 \right]} \\
& = b \sqrt{\mathbb{E}_{(x_1, k, x') \sim D_{\text{coup}}} \left[\left(g^*(x_1, k, x') - \hat{g}(\hat{\phi}(x_1), k, x') \right)^2 \right]} \\
& = b \sqrt{\mathbb{E}_{(x_1, k, x') \sim D_{\text{cont}}} \left[\frac{D_{\text{coup}}(x_1, k, x')}{D_{\text{cont}}(x_1, k, x')} \left(g^*(x_1, k, x') - \hat{g}(\hat{\phi}(x_1), k, x') \right)^2 \right]} \\
& \leq b \sqrt{2 \mathbb{E}_{(x_1, k, x') \sim D_{\text{cont}}} \left[\left(g^*(x_1, k, x') - \hat{g}(\hat{\phi}(x_1), k, x') \right)^2 \right]} \\
& \leq \sqrt{2} b \Delta_{\text{cont}}.
\end{aligned}$$

857 where we use Cauchy-Schwartz's inequality in the first step and Equation 16 in the second inequality.
858 The second term is bounded as:

$$\begin{aligned}
& \mathbb{E}_{(x_1, x_2, k, x') \sim D_{\text{coup}}} \left[\mathbf{1} \left\{ \hat{\phi}(x_1) = \hat{\phi}(x_2) \right\} \left| \hat{g}(\hat{\phi}(x_1), k, x') - g^*(x_2, k, x') \right| \right] \\
& = \mathbb{E}_{(x_1, x_2, k, x') \sim D_{\text{coup}}} \left[\mathbf{1} \left\{ \hat{\phi}(x_1) = \hat{\phi}(x_2) \right\} \left| \hat{g}(\hat{\phi}(x_2), k, x') - g^*(x_2, k, x') \right| \right] \\
& = \mathbb{E}_{(x_1, x_2, k, x') \sim D_{\text{coup}}} \left[\mathbf{1} \left\{ \hat{\phi}(x_1) = \hat{\phi}(x_2) \right\} \left| \hat{g}(\hat{\phi}(x_1), k, x') - g^*(x_1, k, x') \right| \right] \\
& \leq \sqrt{2} b \Delta_{\text{cont}}.
\end{aligned}$$

859 where we use the coupling argument in the first step and then reduce it to the first term using
860 symmetric of (x_1, x_2) in D_{coup} . Combining the upper bounds of the two terms and using $b \leq 1$ and
861 $2\sqrt{2} < 4$ completes the proof. \square

862 **Assumption 9** (Temporal Contrastive Margin). *We assume that there exists a $\beta_{\text{temp}} > 0$ such that*
863 *for any two different states s_1 and s_2 :*

$$\frac{1}{2} \mathbb{E}_{k \sim \omega, s' \sim \rho} [|g^*(s_1, k, s') - g^*(s_2, k, s')|] \geq \beta_{\text{temp}}$$

864 The factor of $\frac{1}{2}$ is chosen for comparison with forward modeling as will become clear later at the end
865 of the proof. As before, if k is fixed, the margin is given by

$$\beta_{\text{temp}}^{(k)} := \frac{1}{2} \inf_{s_1 \neq s_2; s_1, s_2 \in \mathcal{S}} \mathbb{E}_{s' \sim \rho} [|g^*(s_1, k, s') - g^*(s_2, k, s')|]$$

866 and when $k \sim \text{Unf}([K])$ the margin is given by

$$\beta_{\text{temp}}^{(u)} := \frac{1}{2} \inf_{s_1 \neq s_2; s_1, s_2 \in \mathcal{S}} \mathbb{E}_{k \sim \text{Unf}([K]), s' \sim \rho} [|g^*(s_1, k, s') - g^*(s_2, k, s')|]$$

867 We directly have $\beta_{\text{temp}}^{(u)} \geq \frac{1}{K} \sum_{k=1}^K \beta_{\text{temp}}^{(k)}$.

Lemma 4.

$$\mathbb{P}_{x_1, x_2 \sim \rho} \left(\hat{\phi}(x_1) = \hat{\phi}(x_2) \wedge \phi^*(x_1) \neq \phi^*(x_2) \right) \leq \frac{2\Delta_{\text{cont}}(n, \delta)}{\beta_{\text{temp}}}$$

868 *Proof.* We start with the left-hand side in Theorem 9.

$$\begin{aligned}
& \mathbb{E}_{(x_1, k, x_2, x') \sim D_{\text{coup}}} \left[\mathbf{1} \left\{ \hat{\phi}(x_1) = \hat{\phi}(x_2) \right\} |g^*(x_1, k, x') - g^*(x_2, k, x')| \right] \\
&= \mathbb{E}_{(x_1, x_2) \sim D_{\text{coup}}} \left[\mathbf{1} \left\{ \hat{\phi}(x_1) = \hat{\phi}(x_2) \right\} \mathbb{E}_{k \sim \omega, x' \sim \rho} [|g^*(x_1, k, x') - g^*(x_2, k, x')|] \right] \\
&= \mathbb{E}_{(x_1, x_2) \sim \rho} \left[\mathbf{1} \left\{ \hat{\phi}(x_1) = \hat{\phi}(x_2) \right\} \mathbb{E}_{k \sim \omega, s' \sim \rho} [|g^*(x_1, k, s') - g^*(x_2, k, s')|] \right] \\
&\geq 2\beta_{\text{temp}} \mathbb{E}_{(x_1, x_2) \sim \rho} \left[\mathbf{1} \left\{ \hat{\phi}(x_1) = \hat{\phi}(x_2) \wedge \phi^*(x_1) \neq \phi^*(x_2) \right\} \right] \\
&= 2\beta_{\text{temp}} \mathbb{P}_{(x_1, x_2) \sim \rho} \left[\hat{\phi}(x_1) = \hat{\phi}(x_2) \wedge \phi^*(x_1) \neq \phi^*(x_2) \right],
\end{aligned}$$

869 where we use the definition of β_{temp} , the fact that marginal over $D_{\text{coup}}(X)$ is ρ , and that $g^*(x, k, x')$
870 only depends on $\phi^*(x')$ and $\phi^*(x)$ (Lemma 3). Combining with the inequality proved in Theorem 9,
871 completes the proof. \square

872 We have now reduced this analysis to an almost identical one to the forward analysis case (Proposition
873 4). We can, therefore, use the same steps and derive identical bounds. All what changes is that
874 β_{for} is replaced by β_{temp} and in Δ we replace $\ln|\mathcal{F}|$ with $\ln|\mathcal{G}|$. At this point, we can clarify that
875 the factor of $\frac{1}{2}$ was chosen in the definition of β_{temp} so that β_{for} can be replaced by β_{temp} rather
876 than $\frac{\beta_{\text{temp}}}{2}$ which will make it harder to compare margins, as we will do later.

877 C.3 Proof of Lower Bound for Exogenous Block MDPs

878 *thm:exo_lower_bound.* We present a hard instance using a family of exogenous block MDPs, with
879 $H = 2$, $\mathcal{A} = \{1, 2\}$, and a single binary endogenous factor and $d - 1$ exogenous binary factors for
880 each level, where each endogenous and exogenous factor. We first fix an absolute constant $p \in [0, 1]$.
881 Each MDP M_i is indexed by $i \in [d]$, and is specified as follows:

- 882 • **State space:** The state is represented by $x_h := [s_h 1, s_h 2, \dots, s_h d]$, where the superscript
883 denotes different factors. For MDP M_i , only the i -th factor $s_h i$ is an endogenous state for
884 all h , and the other factors are exogenous. Each factor has values of $\{0, 1\}$.
- 885 • **Transition:** For the MDP instance M_i : it has
 - 886 1. For the i -th factor (endogenous factor), $\mathbb{P}(s_2 i \mid s_1 i, a) = [s_2 i = (s_1 i = a)]$. That is,
887 the endogenous states have deterministic dynamics. If $s_1 i = a$, then it transitions to
888 $s_2 i = 1$, otherwise it transitions to $s_2 i = 0$.
 - 889 2. For the j -th factor with $j \neq i$ (exogenous factor), $\mathbb{P}(s_2 j \mid s_1 j) = (1 - p)(s_2 j =$
890 $s_1 j) + p(s_2 j \neq s_1 j)$ for any $s_2 j$ and $s_1 j$. That is, the j -th factor has probability of
891 $1 - p$ of transiting to the same state (i.e., $s_1 j = 0 \rightarrow s_2 j = 0$ or $s_1 j = 1 \rightarrow s_2 j = 1$),
892 and probability of p of transiting to the different state (i.e., $s_1 j = 0 \rightarrow s_2 j = 1$ or
893 $s_1 j = 1 \rightarrow s_2 j = 0$).

894 Note that the MDP terminates at $h = 2$.

- 895 • **Initial state distribution and reward:** The marginal distribution of $s_1 j$ is uniformly
896 distributed at random over $\{0, 1\}$ for all $j \in [d]$, and all factors are independent from each
897 other. For MDP M_i , the agent only receive reward signal after taking action at $h = 2$, with
898 $R(s_2 i, a) = s_2 i$. That is, it always reward 1 at $s_2 i = 1$ and reward 0 at $s_2 i = 0$ no matter
899 which action it takes.
- 900 • **Data collection policy for video data:** We assume that the data collection policy always
901 pick action 0 with probability p and action 1 with probability $1 - p$ for all states.

902 Now we use the following two steps to establish the proof.

903 **Uninformative video data for learning the state decoder** Since video data only contains state
 904 information, from the MDP family construction above, we can easily verify that all MDP instances in
 905 such a family will have an identical video data distribution, *regardless of the choice of constant p* .
 906 This implies that the video data is uninformative for the agent to distinguish the MDP instance from
 907 the MDP family. Now, we assume \mathcal{D}_i is the video data from the instance M_i , and ϕ_i is the state
 908 decoder learned from an arbitrary algorithm \mathcal{A}_1 with \mathcal{D}_i . Then, for any arbitrary algorithm \mathcal{A}_2 that
 909 uses the state decoder ϕ_i in its execution, it is equivalent to such an \mathcal{A}_2 that uses the state decoder ϕ_j
 910 in its execution, where j can be selected arbitrarily from $[d]$.

911 **State decoder requiring exponential length** Without loss of generality, we further restrict the
 912 state decoder ϕ used in the execution of \mathcal{A}_2 for all MDP instance to be some $\phi_h : \mathcal{X} \rightarrow [L]$, where
 913 $h \in \{1, 2\}$ and $L \leq 2^d$. Then we will argue that there must exist a $k \in [d]$, such that

$$\sum_{x_1, \tilde{x}_1 \in \mathcal{X}} \mathbb{P}(\phi_1(x_1) = \phi_1(\tilde{x}_1) \vee (s_1 k \neq \tilde{s}_1 k)) > \frac{2^d - L}{d2^d}, \quad (17)$$

914 where $x_1 := [s_1 1, s_1 2, \dots, s_1 d]$ and $\tilde{x}_1 := [\tilde{s}_1 1, \tilde{s}_1 2, \dots, \tilde{s}_1 d]$. Note that, Eq. (17) means there must
 915 be a probability of at least $2^{d-L}/d2^d$ that ϕ_1 will incorrectly group two different $s_1 k$ together.

916 We now prove Eq. (17). Based on the construct above, we know that $|\mathcal{X}| = 2^d$, and each state in
 917 \mathcal{X} has the same occupancy for x_1 based on the defined initial state distribution (this holds for all
 918 instances in the MDP family, as we are now only talking about the initial state x_1). Thus, we have

$$\sum_{x_1, \tilde{x}_1 \in \mathcal{X}} \mathbb{P}[\phi_1(x_1) = \phi_1(\tilde{x}_1) \vee (s_1 1 = \tilde{s}_1 1) \vee (s_1 2 = \tilde{s}_1 2) \vee \dots \vee (s_1 d = \tilde{s}_1 d)] \leq \frac{L}{2^d}, \quad (18)$$

919 because we defined $\phi_1 : \mathcal{X} \rightarrow [L]$, it means that such ϕ_1 is only able to distinguish the number of L
 920 different states from \mathcal{X} . Then, we obtain

$$\sum_{j \in [d]} \sum_{x_1, \tilde{x}_1 \in \mathcal{X}} \mathbb{P}(\phi_1(x_1) = \phi_1(\tilde{x}_1) \vee (s_1 j \neq \tilde{s}_1 j)) \quad (19)$$

$$= \sum_{x_1, \tilde{x}_1 \in \mathcal{X}} \mathbb{P}(\phi_1(x_1) = \phi_1(\tilde{x}_1)) \quad (20)$$

$$- \sum_{x_1, \tilde{x}_1 \in \mathcal{X}} \mathbb{P}[\phi_1(x_1) = \phi_1(\tilde{x}_1) \vee (s_1 1 = \tilde{s}_1 1) \vee (s_1 2 = \tilde{s}_1 2) \vee \dots \vee (s_1 d = \tilde{s}_1 d)] \quad (21)$$

$$= \frac{2^d - L}{2^d}. \quad (\text{by Eq. (18)})$$

$$\implies \max_{j \in [d]} \sum_{x_1, \tilde{x}_1 \in \mathcal{X}} \mathbb{P}(\phi_1(x_1) = \phi_1(\tilde{x}_1) \vee (s_1 j \neq \tilde{s}_1 j)) > \frac{2^d - L}{d2^d}. \quad (22)$$

921 So this proves Eq. (17).

922 From Eq. (17), we know that for the MDP instance M_k , ϕ_1 will have probability at least $2^{d-L}/2 \cdot d2^d$
 923 to mistake the endogenous state, which implies that for any policy that is represented using the state
 924 decoder ϕ , it must have sub-optimality at least $2^{d-L}/2 \cdot d2^d$. Therefore, it is easy to verify that, for any
 925 $\varepsilon > 0$, we can simply pick $d = 1/4\varepsilon$, and obtain

$$\text{sub-optimality} > \frac{2^d - L}{2 \cdot d2^d} \geq \varepsilon, \quad \forall L \leq 2^{1/4\varepsilon - 1}.$$

926 Then, any arbitrary algorithm \mathcal{A}_2 that uses the state decoder ϕ in its execution, where $\phi_h : \mathcal{X} \rightarrow [L]$
 927 can be chosen arbitrarily for $h \in \{1, 2\}$ and $L \leq 2^{1/4\varepsilon - 1}$, must have sub-optimality larger than ε .

928 **Additional characteristics of MDP family and video data** Note that, by combining the arguments
 929 of uninformative video data and a state decoder requiring exponential length, we obtain impossible
 930 results. We now discuss the following:

931 1. The margin condition defined in Assumption 3 regarding the constructed MDPs

932 2. The PAC learnability of the constructed MDPs

933 3. The coverage condition of video data.

934 For the defined margin condition of forward modeling, we have: for the MDP instance M_i with
 935 constant p , we can bound the forward margin as below (\mathbb{P}_{for} denotes the video distribution)

$$\begin{aligned}
 & \|\mathbb{P}_{\text{for}}(X_2 \mid s_1 i = 0) - \mathbb{P}_{\text{for}}(X_2 \mid s_1 i = 1)\|_{\text{TV}} \\
 &= \frac{1}{2} \sum_{X_2} |\mathbb{P}_{\text{for}}(X_2 \mid s_1 i = 0) - \mathbb{P}_{\text{for}}(X_2 \mid s_1 i = 1)| \\
 &= \frac{1}{2} \sum_{X_2} |\mathbb{P}_{\text{for}}(s_2 i = 0 \mid s_1 i = 0) \mathbb{P}(X_2 \mid s_2 i = 0) + \mathbb{P}_{\text{for}}(s_2 i = 1 \mid s_1 i = 0) \mathbb{P}(X_2 \mid s_2 i = 1) \\
 &\quad - \mathbb{P}_{\text{for}}(s_2 i = 0 \mid s_1 i = 1) \mathbb{P}(X_2 \mid s_2 i = 0) + \mathbb{P}_{\text{for}}(s_2 i = 1 \mid s_1 i = 1) \mathbb{P}(X_2 \mid s_2 i = 1)| \\
 &= \frac{1}{2} \sum_{X_2} |(1 - 2p) [\mathbb{P}(X_2 \mid s_2 i = 0) - \mathbb{P}(X_2 \mid s_2 i = 1)]|. \\
 &\stackrel{(a)}{=} \frac{|1 - 2p|}{2} \sum_{X_2} \mathbb{P}(X_2 \mid s_2 i = 0) + \frac{|1 - 2p|}{2} \sum_{X_2} \mathbb{P}(X_2 \mid s_2 i = 1) \\
 &= |1 - 2p|,
 \end{aligned}$$

936 where step (a) is because $s_2 i$ is a part of X_2 , and then we know $\mathbb{P}(X_2 \mid s_2 i = 0)$ and $\mathbb{P}(X_2 \mid s_2 i = 1)$
 937 cannot be nonzero simultaneously. So picking $p \neq 0.5$ implies positive forward margin.

938 For the temporal contrastive learning, it is easy to verify that $|\mathbb{P}_{\text{for}}(z = 1 \mid s_1 i = 1, X_2) - \mathbb{P}_{\text{for}}(z =$
 939 $1 \mid s_1 i = 1, X_2)| = |1 - 2p|$, so picking $p \neq 0.5$ also implies positive margin for temporal contrastive
 940 learning.

941 As for the PAC learnability, since the latent dynamics of our constructed MDPs are deterministic,
 942 they are provably PAC learnable by Efroni et al. [2022].

943 As for the coverage property of the video data, it is easy to verify

$$\max_{\pi \in \Pi, x_1 \in \mathcal{X}} \frac{\mathbb{P}_{\pi}(x_1, a_1)}{\mathbb{P}_{\text{for}}(x_1, a_1)} = \max_{\pi \in \Pi, x_2 \in \mathcal{X}} \frac{\mathbb{P}_{\pi}(x_2)}{\mathbb{P}_{\text{for}}(x_2)} = \max\{1/p, 1/1-p\}.$$

944 Therefore, we can simply pick $p = 1/3$ and obtain the desired MDP and video data properties. This
 945 completes the proof. \square

946 **Addition remark of Theorem 2** In the proof of Theorem 2, if we pick $p = 0.5$ for that hard
 947 instance, the constructed MDP family reduces to a block MDP without exogenous noise, but the
 948 margin becomes 0 for both forward modeling and temporal contrastive learning. Therefore, it implies
 949 that either the exogenous noise or zero forward margin could make the learnability of the problem
 950 impossible.

951 C.4 Can we get efficient learning under additional assumptions?

952 Our lower bound suggests that one can in general not learn efficient and correct representations with
 953 just video data. However, it may be possible in some cases to do so with an additional assumption.
 954 We highlight one example here and defer a proper formal analysis to future work. One path to success
 955 is when the gold decoder results in the best-in-class error. A domain where this can happen is when
 956 the endogenous state is more predictive of x' than any other $\ln|\mathcal{S}|$ bits of information in x . E.g., in
 957 a navigation domain, there can be many sources of noise in the background, but memorizing all of
 958 them can easily overwhelm the decoder's model capacity. Instead focusing solely on modeling the
 959 agent's state can simplify the task of predicting the future.

960 Recently some approaches have also considered recovering *latent actions* from video data using an
 961 encoder-decoder approach [Ye et al., 2022]. In general, the lower bound in Theorem 2 applies to these
 962 methods and they do not provably work in the hard instances with exogenous noise. For example,
 963 the latent actions can capture *exogenous noise* instead of actions, if the former is more predictive of
 964 changes in the observations. However, in simpler cases such as 3D games, where the agent's action is

965 typically most predictive of changes in observations, or in settings with no exogenous noise, one can
 966 expect these approaches to do well.

967 C.5 Relation Between Margins

968 We defined margins β_{for} for forward modeling and β_{temp} for temporal contrastive learning. The
 969 larger the values of these margins, the more easy it is to separate observations from different
 970 endogenous states. This can be directly inferred from the sample complexity bounds which scale
 971 inversely with these margins. In particular, both β_{for} and β_{temp} depend on the way we sample the
 972 multi-step variable k . We consider two special cases: one where $k \in [K]$ is fixed, we instantiate these
 973 margins as $\beta_{\text{for}}^{(k)}$ and $\beta_{\text{temp}}^{(k)}$, and second where k is uniformly sampled from $[K]$ and we instantiate
 974 those margins as $\beta_{\text{for}}^{(u)}$ and $\beta_{\text{temp}}^{(u)}$.

975 A natural question is how these margins are related. The sample complexity bounds of forward
 976 modeling and temporal contrastive are almost identical except for the difference in margins (β_{for} vs
 977 β_{temp}) and the function classes (\mathcal{F} vs \mathcal{G}). If the function classes were of similar complexity, then
 978 having a larger margin will make it easier to learn the right representation.²

979 **Theorem 10** (Margin Relation). *For any Block MDP and $K \in \mathbb{N}$, the margins*
 980 $\beta_{\text{for}}^{(k)}, \beta_{\text{for}}^{(u)}, \beta_{\text{temp}}^{(k)}, \beta_{\text{temp}}^{(u)} > 0$ *are related as:*

$$\begin{aligned} \frac{1}{K} \beta_{\text{for}}^{(k)} &\leq \beta_{\text{for}}^{(u)} \\ \frac{1}{K} \beta_{\text{temp}}^{(k)} &\leq \beta_{\text{temp}}^{(u)} \\ \frac{\eta_{\min}^2}{4H^2} \beta_{\text{for}}^{(k)} &\leq \beta_{\text{temp}}^{(k)} \leq \beta_{\text{for}}^{(k)} \\ \frac{\eta_{\min}^2}{4H^2} \beta_{\text{for}}^{(u)} &\leq \beta_{\text{temp}}^{(u)} \leq \beta_{\text{for}}^{(u)}. \end{aligned}$$

981 *Proof.* We first prove the first two relations. Fix any $k \in [K]$ then,

$$\begin{aligned} \beta_{\text{for}}^{(u)} &= \inf_{s_1 \neq s_2, s_1, s_2 \in \mathcal{S}} \mathbb{E}_{k' \sim \text{Unif}([K])} [\|D_{pr}(X' | s_1, k') - D_{pr}(X' | s_2, k')\|_{\text{TV}}], \\ &\geq \frac{1}{K} \sum_{k'=1}^K \inf_{s_1 \neq s_2, s_1, s_2 \in \mathcal{S}} \|D_{pr}(X' | s_1, k') - D_{pr}(X' | s_2, k')\|_{\text{TV}}, \\ &\geq \frac{1}{K} \inf_{s_1 \neq s_2, s_1, s_2 \in \mathcal{S}} \|D_{pr}(X' | s_1, k) - D_{pr}(X' | s_2, k)\|_{\text{TV}}, \\ &= \frac{1}{K} \beta_{\text{for}}^{(k)}. \end{aligned}$$

982 Similarly,

$$\begin{aligned} \beta_{\text{temp}}^{(u)} &= \frac{1}{2} \inf_{s_1 \neq s_2, s_1, s_2 \in \mathcal{S}} \mathbb{E}_{k' \sim \text{Unif}([K]), s' \sim \rho} [|g^*(s_1, k', s') - g^*(s_2, k', s')|], \\ &\geq \frac{1}{2K} \sum_{k'=1}^K \inf_{s_1 \neq s_2, s_1, s_2 \in \mathcal{S}} \mathbb{E}_{s' \sim \rho} [|g^*(s_1, k', s') - g^*(s_2, k', s')|], \\ &\geq \frac{1}{2K} \inf_{s_1 \neq s_2, s_1, s_2 \in \mathcal{S}} \mathbb{E}_{s' \sim \rho} [|g^*(s_1, k, s') - g^*(s_2, k, s')|], \\ &= \frac{1}{K} \beta_{\text{temp}}^{(k)}. \end{aligned}$$

²This inference has to be made with a caveat that since we are comparing upper bounds, we cannot guarantee this to hold.

983 We now prove the next two relations. We will prove these bounds for a generic distribution $\omega \in$
 984 $\Delta([K])$ over k . Recall that ω is point-mass over k for $\beta_{\text{temp}}^{(k)}$ and $\text{Unf}([K])$ for $\beta_{\text{temp}}^{(u)}$. We denote
 985 our generic margins as β_{for} and β_{temp} for $k \sim \omega$. We use a shorthand notation $W_k(s, s') =$
 986 $\frac{\rho(s')}{D_{pr}(s'|s, k) + \rho(s')}$ for a given pair of states s, s' and integer $k \in [K]$. It is easy to see that $W_k(s, s') \leq 1$
 987 as $D_{pr}(s' | s, k), \rho(s') \in (0, 1]$. Further, we have $W_k(s, s') \geq \frac{\rho(s')}{2} \geq \frac{\eta_{\min}}{2H}$ where we use
 988 $D_{pr}(s' | s, k), \rho(s') \in (0, 1]$, and Equation 8.

989 We have $g^*(s, k, s') = D_{\text{cont}}(z = 1 | s, k, s') = g^*(s, k, s') = \frac{D_{pr}(s'|s, k)}{D_{pr}(s'|s, k) + \rho(s')}$ using the definition
 990 of D_{cont} in Lemma 3 and Assumption 8. We can use the shorthand W_k and the definition of g^* to
 991 show

$$\begin{aligned}
 \beta_{\text{temp}} &= \frac{1}{2} \inf_{s_1 \neq s_2, s_1, s_2 \in \mathcal{S}} \mathbb{E}_{k \sim \omega, s' \sim \rho} [|g^*(s_1, k, s') - g^*(s_2, k, s')|], \\
 &= \frac{1}{2} \inf_{s_1 \neq s_2, s_1, s_2 \in \mathcal{S}} \sum_{k=1}^K \omega(k) \sum_{s' \in \mathcal{S}} \rho(s') |g^*(s_1, k, s') - g^*(s_2, k, s')|, \\
 &= \frac{1}{2} \inf_{s_1 \neq s_2, s_1, s_2 \in \mathcal{S}} \sum_{k=1}^K \omega(k) \sum_{s' \in \mathcal{S}} W_k(s_1, s') W_k(s_2, s') |D_{pr}(s' | s_1, k) - D_{pr}(s' | s_2, k)|.
 \end{aligned} \tag{23}$$

992 As $W_k(s_1, s') \leq 1$ and $W_k(s_2, s') \leq 1$ we have

$$\begin{aligned}
 \beta_{\text{for}} &= \frac{1}{2} \inf_{s_1 \neq s_2, s_1, s_2 \in \mathcal{S}} \sum_{k=1}^K \omega(k) \sum_{s' \in \mathcal{S}} \underbrace{W_k(s_1, s')}_{\leq 1} \underbrace{W_k(s_2, s')}_{\leq 1} |D_{pr}(s' | s_1, k) - D_{pr}(s' | s_2, k)|, \\
 &\leq \frac{1}{2} \inf_{s_1 \neq s_2, s_1, s_2 \in \mathcal{S}} \sum_{k=1}^K \omega(k) \sum_{s' \in \mathcal{S}} |D_{pr}(s' | s_1, k) - D_{pr}(s' | s_2, k)|, \\
 &= \inf_{s_1 \neq s_2, s_1, s_2 \in \mathcal{S}} \mathbb{E}_{k \sim \omega} [\|D_{pr}(s' | s_1, k) - D_{pr}(s' | s_2, k)\|_{\text{TV}}] \\
 &= \beta_{\text{for}}.
 \end{aligned}$$

993 This gives us $\beta_{\text{temp}}^{(k)} \leq \beta_{\text{for}}^{(k)}$ and $\beta_{\text{temp}}^{(u)} \leq \beta_{\text{for}}^{(u)}$. Finally, we prove the lower bounds. Starting
 994 from Equation 23 and using $W_k(s_1, s') \geq \frac{\eta_{\min}}{2H}$ and $W_k(s_2, s') \leq \frac{\eta_{\min}}{2H}$ we get the following:

$$\begin{aligned}
 \beta_{\text{for}} &= \frac{1}{2} \inf_{s_1 \neq s_2, s_1, s_2 \in \mathcal{S}} \sum_{k=1}^K \omega(k) \sum_{s' \in \mathcal{S}} \underbrace{W_k(s_1, s')}_{\geq \eta_{\min}/2H} \underbrace{W_k(s_2, s')}_{\geq \eta_{\min}/2H} |D_{pr}(s' | s_1, k) - D_{pr}(s' | s_2, k)|, \\
 &\geq \frac{\eta_{\min}^2}{4H^2} \cdot \frac{1}{2} \inf_{s_1 \neq s_2, s_1, s_2 \in \mathcal{S}} \sum_{k=1}^K \omega(k) \sum_{s' \in \mathcal{S}} |D_{pr}(s' | s_1, k) - D_{pr}(s' | s_2, k)|, \\
 &= \frac{\eta_{\min}^2}{4H^2} \inf_{s_1 \neq s_2, s_1, s_2 \in \mathcal{S}} \mathbb{E}_{k \sim \omega} [\|D_{pr}(s' | s_1, k) - D_{pr}(s' | s_2, k)\|_{\text{TV}}] \\
 &= \frac{\eta_{\min}^2}{4H^2} \beta_{\text{for}}.
 \end{aligned}$$

995 This gives us $\beta_{\text{temp}}^{(k)} \geq \frac{\eta_{\min}^2}{4H^2} \beta_{\text{for}}^{(k)}$ and $\beta_{\text{temp}}^{(u)} \geq \frac{\eta_{\min}^2}{4H^2} \beta_{\text{for}}^{(u)}$ which completes the proof. \square

996 The main finding of the above theorem is that forward modeling has a higher margin than temporal
 997 contrastive learning. However, typically the function class used for forward modeling has a higher
 998 statistical complexity than those for temporal contrastive learning as the latter is solving a simpler
 999 binary classification problem than generating an observation.

1000 **C.6 Why temporal contrastive learning is more susceptible to exogenous noise than forward**
 1001 **modeling**

1002 Theorem 2 shows that in the presence of exogenous noise, no video-based representation learning
 1003 approach can be efficient in the worst case. However, this result only presents a worst-case analysis.
 1004 In this section, we show an instance-dependent analysis. The main finding is that the temporal
 1005 contrastive approach is very susceptible to even the smallest amount of exogenous noise, while
 1006 forward modeling is more robust to the presence of exogenous noise. However, both approaches fail
 1007 when there is a significant amount of exogenous noise, consistent with Theorem 2.

1008 **Problem Instance.** We consider a Block MDP with exogenous noise with a state space of $\mathcal{S} =$
 1009 $\{0, 1\}$, action space of $\mathcal{A} = \{0, 1\}$ and exogenous noise space of $\xi = \{0, 1\}$. We consider $H = 1$
 1010 with a uniform distribution over s_1 and ξ_1 , i.e., the start state s_1 and the start exogenous noise
 1011 variable ξ_1 are chosen uniformly from $\{0, 1\}$. The transition dynamics are deterministic and given as follows:
 1012 given action $a_1 \in \{0, 1\}$ and state $s_1 \in \{0, 1\}$, we deterministically transition to $s_2 = 1 - s_1$ if
 1013 $s_1 = a_1$, otherwise, we remain in $s_2 = s_1$. The exogenous noise variable deterministically transitions
 1014 from ξ_1 to $\xi_2 = 1 - \xi_1$. The reward function is given by $R(s_2, s_1) = \mathbf{1}\{s_2 = s_1\}$. We use the
 1015 indicator notation $\mathbf{1}\{\mathcal{E}\}$ to denote 1 if the condition \mathcal{E} is true and 0 otherwise. The observation
 1016 space is given by $\mathcal{X} = \{0, 1\}^{m+2}$ where $(m+2)$ is the dimension of observation space. Given the
 1017 endogenous state s and exogenous noise ξ , the environment generates an observation stochastically
 1018 as $x = [\xi, v_1, \dots, v_l, w_1, \dots, w_{m-l}, s]$ where $v_i \sim \text{psamp}(\cdot | \xi)$ and $w_j \sim \text{psamp}(\cdot | s)$ for all
 1019 $i \in [l]$ and $j \in [m-l]$. The distribution $\text{psamp}(u | s)$ generates $u = s$ with a probability 0.8 and
 1020 $u = 1 - s$ with a probability 0.2. The hyperparameter l is a fixed integer controlling what portion of
 1021 the observation is generated by the exogenous noise compared to the endogenous state. If $l = 1$, we
 1022 only have a small amount of exogenous noise, while if $l = m - 1$ we have the maximal amount of
 1023 exogenous noise. The state s and exogenous noise ξ are both decodable from the observation x . The
 1024 optimal policy achieves a return of 1 and takes action $a_1 = 1$ if $s_1 = 0$ and $a_1 = 0$ if $s_1 = 1$. As the
 1025 optimal policy depends on the value of s_1 , we must learn the latent state to realize the optimal policy.

1026 **Learning Setting.** We assume a decoder class $\Phi = \{\phi^*, \phi_\xi^*\}$ consisting of the true decoder ϕ^*
 1027 and the incorrect decoder ϕ_ξ^* which maps observation to the exogenous noise ξ . Both decoders take
 1028 an observation and map it to a value in $\{0, 1\}$. We assume access to an arbitrarily large dataset \mathcal{D}
 1029 consisting of tuples (x_1, x_2) collecting iid using a fixed data policy π_{data} . This policy takes action
 1030 $a_1 = 0$ in $s_1 = 0$ and action $a_1 = 1$ in $s_1 = 1$. Let $D(x_1, x_2)$ be the data distribution induced by
 1031 π_{data} . We will use D to define other distributions induced by $D(x_1, x_2)$, for example $D(x_2)$ or
 1032 $D(s_2)$. We also assume access to two model classes $\mathcal{F} : \{0, 1\} \rightarrow \Delta(\mathcal{X})$ and $\mathcal{G} : \{0, 1\}^2 \rightarrow [0, 1]$.
 1033 We assume these model classes are finite and contain certain constructions that we define later.

1034 **Overview:** As we increase the value of l , the amount of exogenous noise in the environment
 1035 increases. We will prove that irrespective of the value of l , temporal contrastive learning assigns
 1036 the same loss for both the correct decoder ϕ^* and the incorrect decoder ϕ_ξ^* . In contrast, the forward
 1037 modeling approach is able to prefer ϕ^* over ϕ_ξ^* when the noise is limited, specifically, when $l < m/2$.
 1038 This will establish that temporal contrastive is very susceptible to exogenous noise whereas forward
 1039 modeling is more robust. However, both approaches provably fail when there is $l \geq m/2$.

1040 As we have $H = 1$, we will denote x_2, s_2, ξ_2 by x', s', ξ' and x_1, s_1, ξ_1 by x, s, ξ respectively.
 1041 Note that unless specified otherwise, s and ξ are the endogenous state and exogenous noise of the
 1042 observation x . Similarly, s' and ξ' are the endogenous state and exogenous noise of x' . We will also
 1043 use a shorthand $q(x')$ to denote the emission probability $q(x' | \phi_\xi^*(x'), \phi^*(x'))$ given its endogenous
 1044 state and exogenous noise. We first state the conditional data distribution $D(x' | x)$.

$$\begin{aligned} D(x' | x) &= q(x') T_\xi(\xi' | \xi) \sum_{a \in \mathcal{A}} T(s' | s, a) \pi_{\text{data}}(a | s), \\ &= q(x') \mathbf{1}\{\xi' = 1 - \xi\} \mathbf{1}\{s' = 1 - s\}, \end{aligned} \quad (24)$$

1045 where we use $T_\xi(\xi' | \xi) = \mathbf{1}\{\xi' = 1 - \xi\}$ and $\sum_{a \in \mathcal{A}} T(s' | s, a) \pi_{\text{data}}(a | s) = \mathbf{1}\{s' = 1 - s\}$
 1046 which follows from the definition of π_{data} . Note that $D(x' | x)$ only depends on x via s, ξ , therefore,
 1047 we can define $D(x' | x) = D(x' | s, \xi)$.

1048 Let \tilde{x} be an observation variable with endogenous state \tilde{s} and exogenous noise $\tilde{\xi}$, i.e., $\tilde{s} = \phi^*(\tilde{x})$ and
 1049 $\tilde{\xi} = \phi_\xi^*(\tilde{x})$. We use this to derive the marginal data distribution ρ over x' as follows:

$$\begin{aligned} \rho(x') &= \sum_{s, \xi \in \{0,1\}} D(x', s, \xi) = \sum_{s, \xi \in \{0,1\}} D(x' | s, \xi) \mu(s) \mu_\xi(\xi), \\ &= \frac{q(x')}{4} \sum_{s, \xi \in \{0,1\}} \mathbf{1}\{\xi' = 1 - \xi\} \mathbf{1}\{s' = 1 - s\}, \\ &= \frac{q(x')}{4}, \end{aligned} \quad (25)$$

1050 where in the second step uses the fact that μ and μ_ξ are uniform and Eq. (24). We are now ready to
 1051 prove our desired result.

1052 **Temporal contrastive learning cannot distinguish between good and bad decoder for all $l \in$**
 1053 **$[m - 1]$.** We first recall that temporal contrastive learning approach use the given observed data
 1054 (x_1, x_2) to compute a set of real and fake observation tuples. This is collected into a dataset (x, x', z)
 1055 where $z = 1$ indicates that $(x_1 = x, x_2 = x')$ was observed in the dataset, and $z = 0$ indicates that
 1056 $(x_1 = x, x_2 = x')$ was not observed, or is an imposter. We sample z uniformly in $\{0, 1\}$. The fake
 1057 data is constructed by take $x = x_1$ from one tuple and $x' = x_2$ from another observed tuple. We start
 1058 by computing the optimal Bayes classifier for the temporal contrastive learning approach using the
 1059 definition of Bayes classifier in Lemma 3.

$$D_{\text{cont}}(z = 1 | x, x') = \frac{D(x' | x)}{D(x' | x) + \rho(x')} = \frac{\mathbf{1}\{s' = 1 - s\} \mathbf{1}\{\xi' = 1 - \xi\}}{\mathbf{1}\{s' = 1 - s\} \mathbf{1}\{\xi' = 1 - \xi\} + 1/4},$$

1060 where we use Lemma 3 in the first step and Eqs. (24) and (25) in the second step. Recall that $z = 1$
 1061 denotes whether a given observation tuple (x, x') is real rather than an imposter/false. Note that since
 1062 we have $k = 1$, as it is a $H = 1$ problem, we drop the notation k from all terms.

1063 The marginal distribution over (x, x') for the temporal contrastive is given by Eq. (14) which in our
 1064 case instantiates to:

$$\begin{aligned} D_{\text{cont}}(x, x') &= \frac{D(x)}{2} \{D(x' | x) + \rho(x')\}, \\ &= \frac{1}{8} q(x') q(x) \{\mathbf{1}\{s' = 1 - s\} \mathbf{1}\{\xi' = 1 - \xi\} + 1/4\}, \end{aligned} \quad (26)$$

1065 where we use Eqs. (24) and (25), and $D(x) = q(x) \mu(s) \mu_\xi(\xi) = q(x)/4$.

1066 Let $g \in \mathcal{G}$ be any classifier head. Given a decoder ϕ , we define $g \circ \phi : (x, x') \mapsto g(\phi(x), \phi(x'))$ as a
 1067 model for temporal contrastive learning, with an expected contrastive loss of:

$$\begin{aligned} \ell_{\text{cont}}(g, \phi^*) &= \mathbb{E}_{(x, x') \sim D_{\text{cont}}, z \sim D_{\text{cont}}(\cdot | x, x')} \left[(z - g(\phi^*(x), \phi^*(x')))^2 \right] \\ &= \mathbb{E}_{(x, x') \sim D_{\text{cont}}} \left[D_{\text{cont}}(z = 1 | x, x') (1 - 2g(\phi^*(x), \phi^*(x'))) + g(\phi^*(x), \phi^*(x'))^2 \right] \\ &= \frac{1}{8} \sum_{s, \xi, s', \xi'} \left\{ \mathbf{1}\{s' = 1 - s\} \mathbf{1}\{\xi' = 1 - \xi\} + \frac{1}{4} \right\} \left(\frac{\mathbf{1}\{s' = 1 - s\} \mathbf{1}\{\xi' = 1 - \xi\}}{\mathbf{1}\{s' = 1 - s\} \mathbf{1}\{\xi' = 1 - \xi\} + \frac{1}{4}} (1 - g(s, s')) + g(s, s')^2 \right) \end{aligned}$$

1068 Similarly, the expected temporal contrastive loss of the model $g \circ \phi^*$ with the bad decoder ϕ_ξ^* is given
 1069 by:

$$\begin{aligned} \ell_{\text{cont}}(g, \phi_\xi^*) &= \mathbb{E}_{(x, x') \sim D_{\text{cont}}, z \sim D_{\text{cont}}(\cdot | x, x')} \left[(z - g(\phi_\xi^*(x), \phi_\xi^*(x')))^2 \right] \\ &= \frac{1}{8} \sum_{s, \xi, s', \xi'} \left\{ \mathbf{1}\{s' = 1 - s\} \mathbf{1}\{\xi' = 1 - \xi\} + \frac{1}{4} \right\} \left(\frac{\mathbf{1}\{s' = 1 - s\} \mathbf{1}\{\xi' = 1 - \xi\}}{\mathbf{1}\{s' = 1 - s\} \mathbf{1}\{\xi' = 1 - \xi\} + \frac{1}{4}} (1 - g(\xi, \xi')) + g(\xi, \xi')^2 \right) \end{aligned}$$

1070 Note that by interchanging s with ξ and s' with ξ' , we can show $\ell_{\text{cont}}(g, \phi_\xi^*) = \ell_{\text{cont}}(g, \phi^*)$.
1071 Therefore, $\inf_{g \in \mathcal{G}} \ell_{\text{cont}}(g, \phi_\xi^*) = \inf_{g \in \mathcal{G}} \ell_{\text{cont}}(g, \phi^*)$. This implies that for any value of l , the
1072 temporal contrastive loss assigns the same loss to the good decoder ϕ^* and the bad decoder ϕ_ξ^* .
1073 Hence, in practice, temporal contrastive cannot distinguish between the good and bad decoder and
1074 may converge to the latter leading to poor downstream performance. This convergence to the bad
1075 decoder may happen if it is easier to overfit to noise. For example, in our gridworld example, it
1076 is possibly easier for the model to overfit to the predictable motion of noise than understand the
1077 underlying dynamics of the agent. This is observed in Fig. 3 where the representation learned via
1078 temporal contrastive tends to overfit to the noisy exogenous pixels and perform poorly on downstream
1079 RL tasks (Fig. 2).

1080 **Forward modeling learns the good decoder if $l < \lfloor m/2 \rfloor$.** We likewise analyze the expected
1081 forward modeling loss of the good and bad decoder. For any $f \in \mathcal{F}$, we have $f(x' | u)$ as the
1082 generator head that acts on a given decoder's output $u \in \{0, 1\}$ and generates the next observation x' .
1083 If we use the good decoder ϕ^* , then we cannot predict the exogenous noise ξ or ξ' which can be
1084 either 0 or 1 with equal probability. This implies that for the l noisy bits v_1, \dots, v_l in x' , the best
1085 prediction is that each one has an equal probability of taking 0 or 1. To see this, fix $i \in [l]$ and recall
1086 that $\mathbb{P}(v_i = \xi' | \xi') = 0.8$ and $\mathbb{P}(v_i = 1 - \xi' | \xi') = 0.2$. As ξ' has equal probability of taking value
1087 0 or 1, therefore, $\mathbb{P}(v_i = u) = \sum_{\xi' \in \{0,1\}} \mathbb{P}(v_i = u | \xi')^{1/2} = \frac{0.8+0.2}{2} = 0.5$. However, since we can
1088 deterministically predict s' , therefore, we can predict the true distribution over w_j for all $j \in [m-l]$.
1089 Let f_{good} be this generator head. Formally, we have:

$$f_{\text{good}}(x' | \phi^*(x)) = \underbrace{(1/2)}$$

due to $x'_1 = \xi' \cdot \underbrace{(1/2)^l}$

due to $v_{1:l} \cdot \underbrace{\prod_{j=l+2}^{m+1} \text{psamp}(x'_j | 1 - \phi^*(x))}$

due to $w_{1:m-l} \cdot \underbrace{\mathbf{1}\{x'_{m+2} = 1 - \phi^*(x)\}}$

1090 due to $x'_{m+2} = s'$

1091 The Bayes distribution is given by:

$$\begin{aligned} D(x' | x) &= q(x') \cdot \mathbf{1}\{\phi^*(x') = 1 - \phi^*(x)\} \cdot \mathbf{1}\{\phi_\xi^*(x') = 1 - \phi_\xi^*(x)\} \\ &= \mathbf{1}\{x'_1 = 1 - \phi_\xi^*(x)\} \cdot \prod_{i=1}^l \text{psamp}(x'_{i+1} | 1 - \phi_\xi^*(x)) \cdot \prod_{j=l+2}^{m+1} \text{psamp}(x'_j | 1 - \phi^*(x)) \mathbf{1}\{x'_{m+2} = 1 - \phi^*(x)\}. \end{aligned}$$

1092 As we are optimizing the log-loss, we look at the expected KL divergence ℓ_{kl} between the $D(x' | x)$
 1093 and $f_{\text{good}}(x' | \phi^*(x))$ which gives:

$$\begin{aligned}
 & \ell_{kl}(f_{\text{good}}, \phi^*) \\
 &= \mathbb{E}_x \left[\sum_{x'} D(x' | x) \ln \frac{D(x' | x)}{f_{\text{good}}(x' | \phi^*(x))} \right] \\
 &= \mathbb{E}_x \left[\sum_{x'} D(x' | x) \ln \frac{\mathbf{1}\{x'_1 = 1 - \phi_\xi^*(x)\} \cdot \prod_{i=1}^l \text{psamp}(x'_{i+1} | 1 - \phi_\xi^*(x))}{(1/2)^{l+1}} \right] \\
 &= (l+1) \ln(2) + \mathbb{E}_x \left[\sum_{x'} D(x' | x) \ln \left(\mathbf{1}\{x'_1 = 1 - \phi_\xi^*(x)\} \cdot \prod_{i=1}^l \text{psamp}(x'_{i+1} | 1 - \phi_\xi^*(x)) \right) \right] \\
 &= (l+1) \ln(2) + \mathbb{E}_x \left[\sum_{i=1}^l \sum_{x'_{i+1} \in \{0,1\}} \text{psamp}(x'_{i+1} | 1 - \phi^*(x)) \ln \text{psamp}(x'_{i+1} | 1 - \phi^*(x)) \right] \\
 &= (l+1) \ln(2) - lH(\text{psamp}),
 \end{aligned}$$

1094 where $H(\text{psamp})$ denotes the conditional entropy given by $-1/2 \sum_{s \in \{0,1\}} \sum_{v \in \{0,1\}} \text{psamp}(v |$
 1095 $s) \ln \text{psamp}(v | s)$. As $\text{psamp}(u | u) = 0.8$ and $\text{psamp}(1 - u | u) = 0.2$, we have $H(\text{psamp}) =$
 1096 $-0.8 \ln(0.8) - 0.2 \ln(0.2) \approx 0.500$. Plugging this in, we get $\ell_{kl}(f_{\text{good}}, \phi^*) = l \ln(2) - 0.5l + \ln(2) =$
 1097 $\ln(2) + 0.193l$.

1098 Finally, the analysis when we use the ϕ_ξ^* decoder is identical to above. In this case, we can predict
 1099 $\phi_\xi^*(x')$ and correctly predict the psamp distribution over all the l -noisy bits $v_{1:l}$. However, for the
 1100 $w_{1:m-l}$ bits and the $x'[m+2]$, our best bet is to predict a uniform distribution. We capture this by
 1101 the generator f_{bad} which gives:

$$f_{\text{bad}}(x' | \phi^*(x)) = \underbrace{(1/2)}$$

due to $x'_{m+2} = s' \cdot \underbrace{(1/2)^{m-l}}$

due to $w_{1:m-l} \cdot \underbrace{\prod_{i=2}^{l+1} \text{psamp}(x'_i | 1 - \phi_\xi^*(x))}$

due to $v_{1:l} \cdot \underbrace{\mathbf{1}\{x'_1 = 1 - \phi_\xi^*(x)\}}$

1102 due to $x'_1 = \xi'$

1103 The expected KL loss $\ell_{kl}(f_{\text{bad}}, \phi_\xi^*)$ can be computed almost exactly as before and is equal to
 1104 $\ln(2) + 0.193(m-l)$. We can see that for $\ell_{kl}(f_{\text{good}}, \phi^*) < \ell_{kl}(f_{\text{bad}}, \phi_\xi^*)$ we must have $\ln(2) +$
 1105 $0.193l < \ln(2) + 0.193(m-l)$, or equivalently, $l < m/2$. This completes the analysis.

1106 D Additional Experimental Details

1107 D.1 Details of Experimental Setup

1108 All results are reported with mean and standard error computed over 3 seeds. All the code for this
 1109 work was run on A100, V100, P40 GPUs, with a compute time of approx. 12 hours for grid world
 1110 experiments and 6 hours for ViZDoom experiments. Data collection for gridworld was done using
 1111 a mixture of random walks, optimal trajectories, deviation from optimal trajectories, and walks to
 1112 randomly chosen goal positions. Data collection for Vizdoom was done via pretrained PPO policies
 1113 along with random walks for diversity in the observation space.

Hyperparameter	Value
batch size	128
learning rate	0.001
epochs	400
# of exogenous variables	10
exogenous pixel size	4
# of VQ heads	2
VQ codebook size	100
VQ codebook temperature	0
VQ codebook dimension	32
VQ bottleneck dimension	1024

Table 2: Hyperparameters used for experiments with the GridWorld and ViZDoom domains.

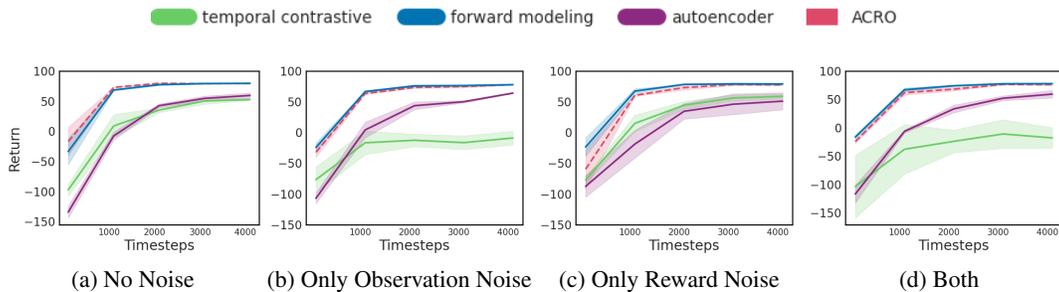


Figure 7: RL experiments using different latent representations for the ViZDoom environment.

1114 **GridWorld Details.** We consider navigation in a 12×12 Minigrid environment [Chevalier-Boisvert
 1115 et al., 2023]. The agent is represented as a red triangle and can take three actions: move forward, turn
 1116 left, and turn right (Figure 3). The agent needs to reach a yellow key. The position of the agent and key
 1117 randomizes each episode. The agent only observes an area around itself (as an agent-centric-view).
 1118 Horizon $H = 12$, and the agent gets a reward of $+1.0$ for reaching the goal and -0.01 in other cases.

1119 **ViZDoom Defend The Center Details.** We test with a ViZDoom environment called Defend the
 1120 Center [Wydmuch et al., 2018, Kempka et al., 2016], which is a first-person shooting game (Figure 5).
 1121 The map is a large circle. A player is spawned in the exact center. 5 monsters are spawned along the
 1122 wall. Monsters are killed after a single shot. After dying, each monster is respawned after some time.
 1123 The episode ends when the player dies. The reward scheme is as follows: $+1$ for killing a monster and
 1124 -1 for death.

1125 **Hyperparameters.** In Table 2, we report the hyperparameter values used for experiments in this
 1126 work with the GridWorld and ViZDoom environments.

1127 D.2 Results on an Additional Domain

1128 **ViZDoom Basic.** We use an additional basic ViZDoom environment [Wydmuch et al., 2018, Kempka
 1129 et al., 2016], which is a first-person shooting game (Figure 8). The player needs to kill a monster
 1130 to win. The map of the environment is a rectangle with gray walls, ceiling, and floor. The player is
 1131 spawned along the longer wall in the center. A red, circular monster is spawned randomly somewhere
 1132 along the opposite wall. The player can take one of three actions at each time step (left, right, shoot).
 1133 One hit is enough to kill the monster. The episode finishes when the monster is killed or on timeout.
 1134 The reward scheme is as follows: $+101$ for shooting the enemy, -1 per time step, and -5 for missed
 1135 shots. Results for this environment are shown in Figure 7 and Figure 8 and further validate our
 1136 findings from theory and experiments.

1137 D.3 Additional Ablations

1138 **Harder Exogenous Noise.** Figure 6 showed the results when we increase the size of the exogenous
 1139 noise variables (diamond shapes overlayed on the image) in the gridworld domain while keeping the

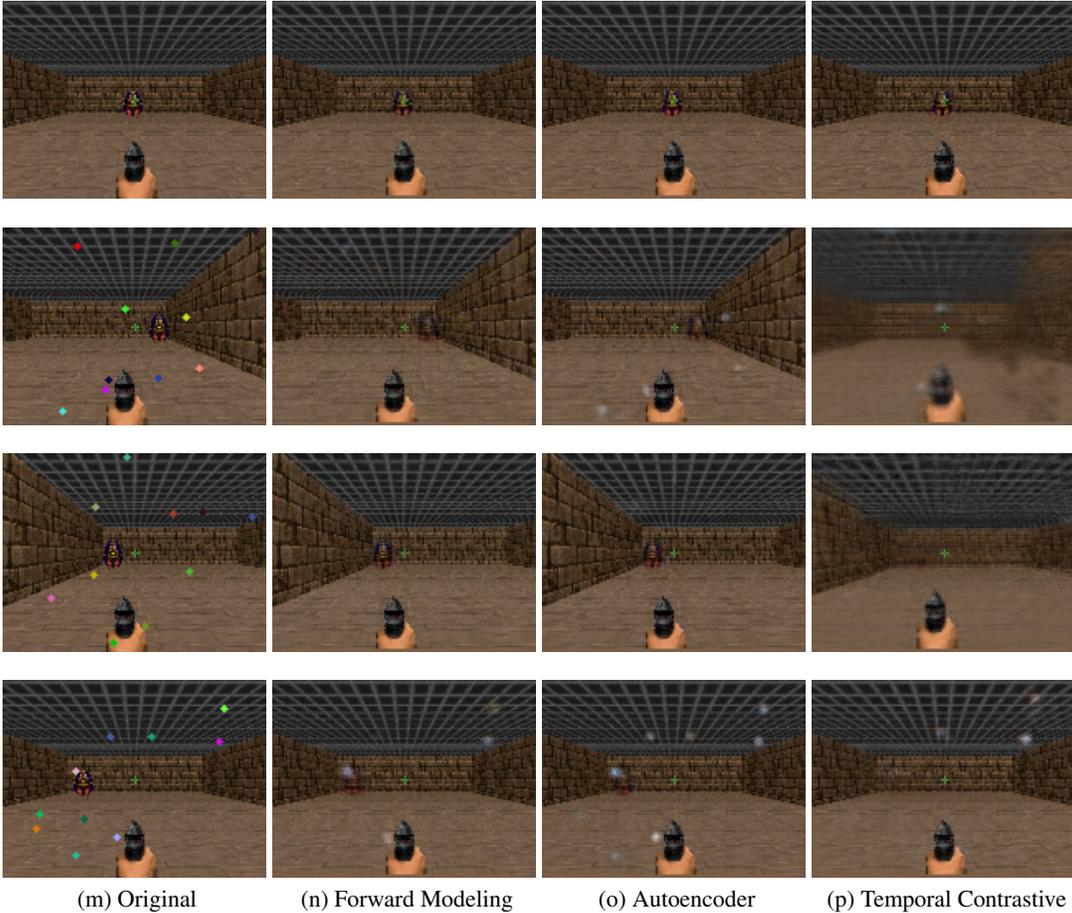


Figure 8: Decoded image reconstructions from different latent representation learning methods in the ViZDoom environment. We train a decoder on top of frozen representations trained with the three video pre-training approaches.

1140 number of exogenous variables fixed at 10. We also increase the number of exogenous noise variables
 1141 in the gridworld domain, while keeping their sizes fixed at 4 pixels and present the results in Figure 9.
 1142 Both results show significant degradation in the performance of video-based representation learning
 1143 methods whereas ACRO which uses trajectory data continues to perform well. This supports one of
 1144 our main theoretical results that exogenous noise poses a challenge for video-based representation
 1145 learning.

1146 **I.I.D. Noise in Gridworld.** We evaluate iid noise in the gridworld domain. We use the diamond-
 1147 shaped exogenous noise that we used in Figure 2, however, at each time step, we randomly sample
 1148 the color and position of each diamond, independent of the agent’s history. Figure 10(a) shows the
 1149 result for forward modeling and Figure 10(b) shows the same for ACRO. We also ablate the number
 1150 of noisy diamonds. As expected, forward modeling and ACRO can learn a good policy while the
 1151 increase in the number of noisy diamonds (num noise var) only slightly decreases their performance.

1152 **I.I.D. Noise in the Basic ViZDoom environment.** We evaluate the representation learning methods
 1153 on the basic ViZDoom domain but with independent and identically distributed (iid) noise. We add iid
 1154 Gaussian noise to each pixel sampled from a 0 mean Gaussian distribution with a standard deviation of
 1155 0.001. Based on theory, we expect temporal contrastive objectives to be substantially better at filtering
 1156 out Gaussian iid noise, which is validated experimentally for the basic ViZDoom Environment
 1157 (Figure 11(a)). Figure 11(b) refers to the basic ViZDoom result for convenient comparison.

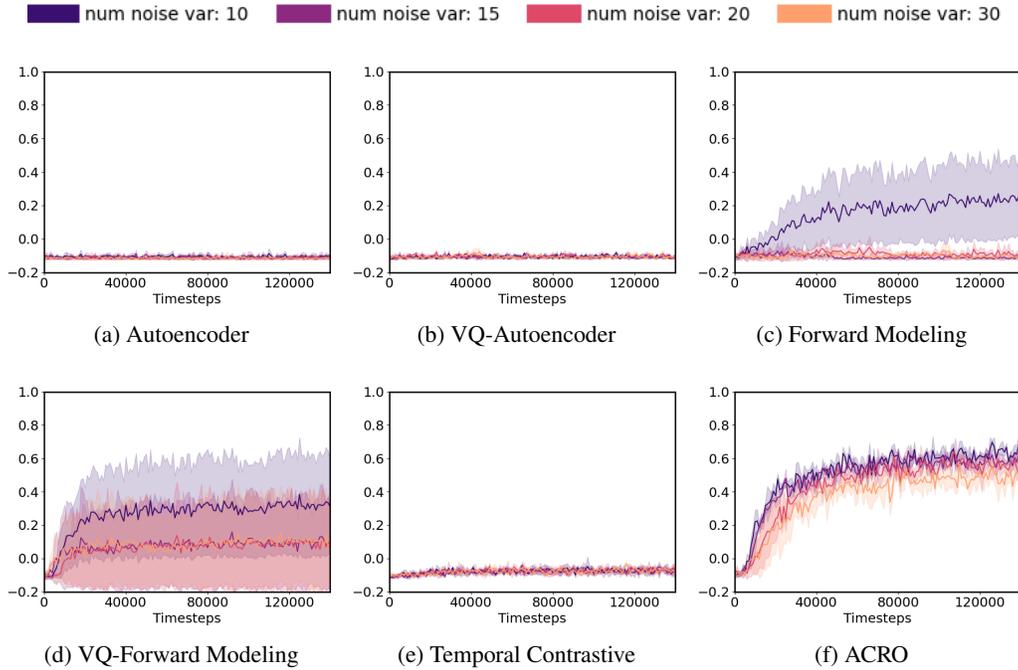


Figure 9: Gridworld experiments with exogenous noise of size 4 and different the number of exogenous noise variables. Several video-based representation learning methods struggle to learn as the number of exogenous noise variables increases, whereas ACRO which uses trajectory data, still performs well.

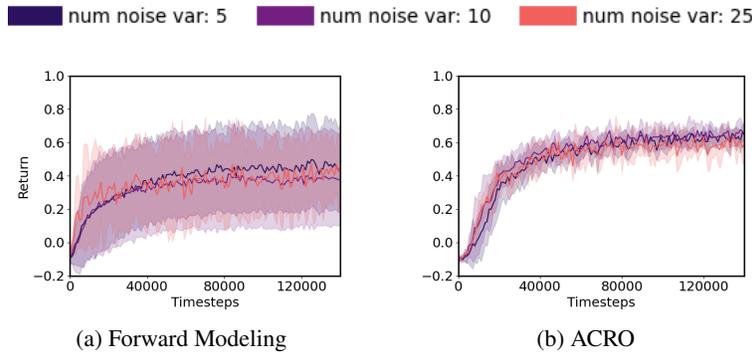


Figure 10: Experiments with iid noise for the Gridworld environment. ‘Num noise var’ denotes the number of noisy diamonds constituting the exogenous noise.

1158 **Additional reconstructions.** We show additional image reconstructions Figure 12 for the Grid-
 1159 World environment and in Figure 13 for the ViZDoom Defend the Center environment. We highlight
 1160 that important parts of the observation space are recovered successfully by the forward modeling
 1161 approach under varying levels of exogenous noise, whereas temporal contrastive learning often fails.

1162

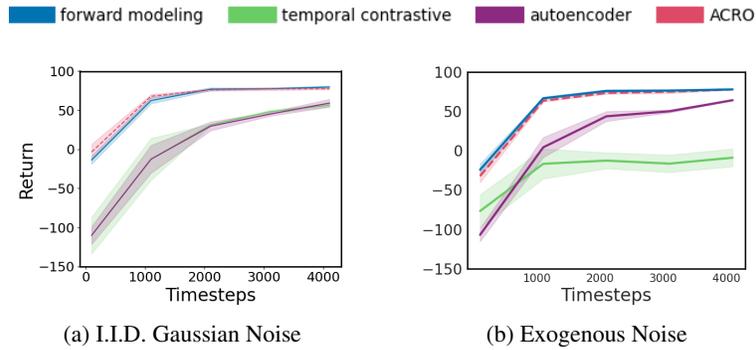


Figure 11: Experiments with (a) Gaussian iid noise for the ViZDoom environment and (b) exogenous noise.

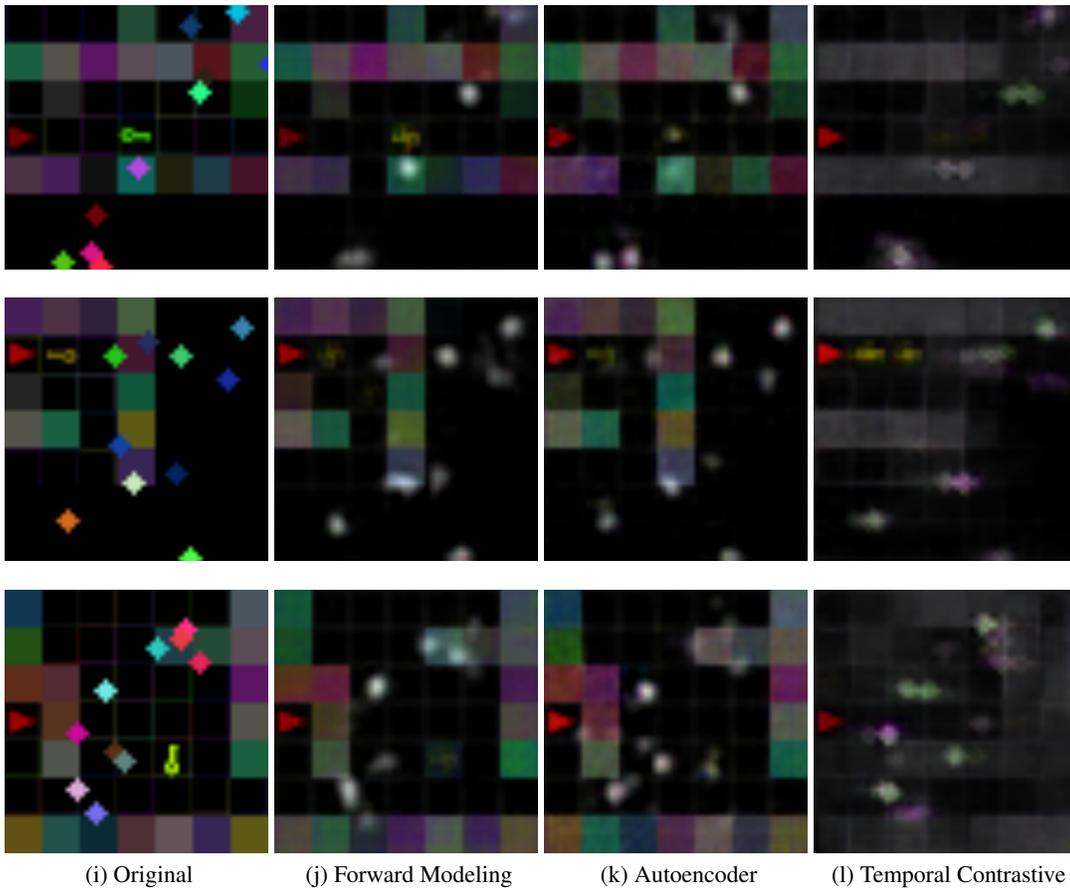


Figure 12: Decoded image reconstructions from different latent representation learning methods in the GridWorld environment. We train a decoder on top of frozen representations trained with the three video pre-training approaches.

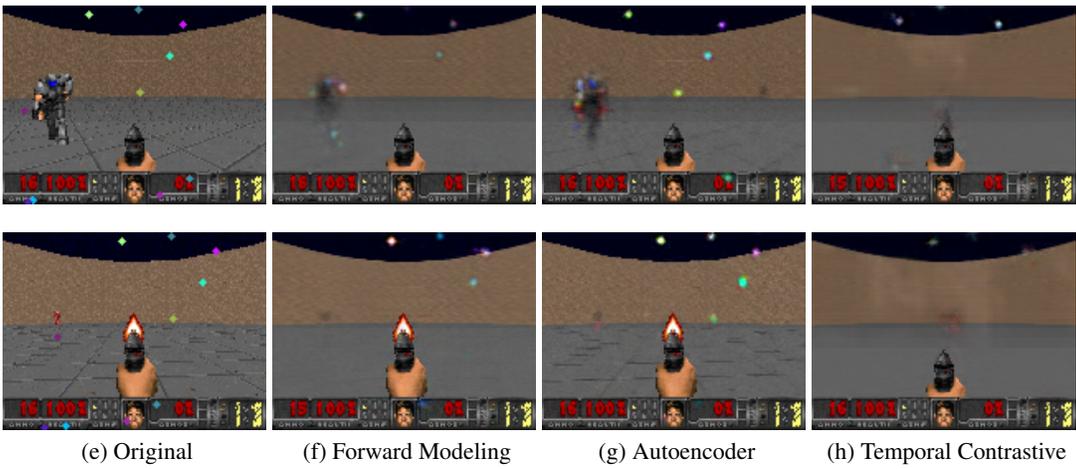


Figure 13: Decoded image reconstructions from different latent representation learning methods in the ViZDoom Defend the Center environment. We train a decoder on top of frozen representations trained with the three video pre-training approaches.