
Appendix for “Learning World Models with Identifiable Factorization”

Anonymous Author(s)

Affiliation

Address

email

1 Appendix organization:

2

3	S1 Theoretical Proofs	2
4	S1.1 Proof of Proposition 1	2
5	S1.2 Proof of Proposition 2	2
6	S1.3 Proof of Proposition 3	3
7	S1.4 Proof of Theorem 1	3
8	S2 Derivation of the objective function	7
9	S2.1 Discussions	9
10	S3 Environment descriptions	11
11	S3.1 Synthetic data	11
12	S3.2 Modified Cartpole	11
13	S3.3 Variant of Robodesk	11
14	S3.4 Variants of DeepMind Control Suite	12
15	S4 Experimental Details	13
16	S4.1 Synthetic Dataset	13
17	S4.1.1 Extra Results.	14
18	S4.2 Modified Cartpole	14
19	S4.3 Variant of Robodesk	14
20	S4.4 Variants of Deep Mind Control Suite	15
21	S4.4.1 Extra Results.	15
22	S4.5 Visualization for DMC	16
23	S5 Comparison between IFactor and Denoised MDP	18

24

25 S1 Theoretical Proofs

26 S1.1 Proof of Proposition 1

27 Proposition 1 shows that s_t^r , which has directed paths to $r_{t+\tau}$ (for $\tau \geq 0$), is minimally sufficient for
 28 policy learning that aims to maximize the future reward and can be characterized by conditional
 29 dependence with the cumulative reward variable R_t .

30 **Proposition 1.** *Under the assumption that the graphical representation, corresponding to the*
 31 *environment model, is Markov and faithful to the measured data, $s_t^r \subseteq \mathbf{s}_t$ is a minimal subset of state*
 32 *dimensions that are sufficient for policy learning, and $s_{i,t} \in s_t^r$ if and only if $s_{i,t} \not\perp R_t | a_{t-1}, s_{t-1}^r$.*

33 We first give the definitions of the Markov condition and the faithfulness assumption, which will be
 34 used in the proof.

35 **Definition 1** (Global Markov Condition [1, 2]). *The distribution p over a set of variables \mathbf{V} satisfies*
 36 *the global Markov property on graph G if for any partition (A, B, C) such that if B d -separates A*
 37 *from C , then $p(A, C|B) = p(A|B)p(C|B)$.*

38 **Definition 2** (Faithfulness Assumption [1, 2]). *There are no independencies between variables that*
 39 *are not entailed by the Markov Condition in the graph.*

40 Below, we give the proof of Proposition 1.

41 *Proof.* The proof contains the following three steps.

- 42 • In step 1, we show that a state dimension $s_{i,t}$ is in s_t^r , that is, it has a directed path to $r_{t+\tau}$, if
 43 and only if $s_{i,t} \not\perp R_t | a_{t-1}, \mathbf{s}_{t-1}$.
- 44 • In step 2, we show that for $s_{i,t}$ with $s_{i,t} \not\perp R_t | a_{t-1}, \mathbf{s}_{t-1}$, if and only if $s_{i,t} \not\perp R_t | a_{t-1}, s_{t-1}^r$.
- 45 • In step 3, we show that s_t^r are minimally sufficient for policy learning.

46 **Step 1:** We first show that if a state dimension $s_{i,t}$ is in s_t^r , then $s_{i,t} \not\perp R_t | a_{t-1}, \mathbf{s}_{t-1}$.

47 We prove it by contradiction. Suppose that $s_{i,t}$ is independent of R_t given a_{t-1} and \mathbf{s}_{t-1} . Then according
 48 to the faithfulness assumption, we can see from the graph that $s_{i,t}$ does not have a directed path to $r_{t+\tau}$,
 49 which contradicts the assumption, because, otherwise, a_{t-1} and \mathbf{s}_{t-1} cannot break the paths between
 50 $s_{i,t}$ and R_t which leads to the dependence.

51 We next show that if $s_{i,t} \not\perp R_t | a_{t-1}, \mathbf{s}_{t-1}$, then $s_{i,t} \in s_t^r$.

52 Similarly, by contradiction suppose that $s_{i,t}$ does not have a directed path to $r_{t+\tau}$. From the graph, it is
 53 easy to see that a_{t-1} and \mathbf{s}_{t-1} must d -separate the path between $s_{i,t}$ and R_t . According to the Markov
 54 assumption, $s_{i,t}$ is independent of R_t given a_{t-1} and \mathbf{s}_{t-1} , which contradicts to the assumption. Since
 55 we have a contradiction, it must be that $s_{i,t}$ has a directed path to $r_{t+\tau}$, i.e. $s_{i,t} \in s_t^r$.

56 **Step 2:** In step 1, we have shown that $s_{i,t} \not\perp R_t | a_{t-1}, \mathbf{s}_{t-1}$, if and only if it has a directed path to $r_{t+\tau}$.
 57 From the graph, it is easy to see that for those state dimensions which have a directed path to $r_{t+\tau}$,
 58 a_{t-1} and \mathbf{s}_{t-1} cannot break the path between $s_{i,t}$ and R_t . Moreover, for those state dimensions which
 59 do not have a directed path to $r_{t+\tau}$, a_{t-1} and s_{t-1}^r are enough to break the path between $s_{i,t}$ and R_t .

60 Therefore, for $s_{i,t}$, $s_{i,t} \not\perp R_t | a_{t-1}, \mathbf{s}_{t-1}$, if and only if $s_{i,t} \not\perp R_{t+1} | a_{t-1}, s_{t-1}^r$.

61 **Step 3:** In the previous steps, it has been shown that if a state dimension $s_{i,t}$ is in s_t^r , then $s_{i,t} \not\perp$
 62 $R_t | a_{t-1}, s_{t-1}^r$, and if a state dimension $s_{i,t}$ is not in s_t^r , then $s_{i,t} \perp R_t | a_{t-1}, s_{t-1}^r$. This implies that s_t^r are
 63 minimally sufficient for policy learning to maximize the future reward. \square

64 S1.2 Proof of Proposition 2

65 Moreover, the proposition below shows that s_t^a , which receives an edge from a_{t-1} , can be directly
 66 controlled by actions and can be characterized by conditional dependence with the action variable.

67 **Proposition 2.** Under the assumption that the graphical representation, corresponding to the
 68 environment model, is Markov and faithful to the measured data, $s_t^a \subseteq \mathbf{s}_t$ is a minimal subset of state
 69 dimensions that are sufficient for direct control, and $s_{i,t} \in s_t^a$ if and only if $s_{i,t} \not\perp a_{t-1} | \mathbf{s}_{t-1}$.

70 Below, we give the proof of Proposition 2.

71 *Proof.* The proof contains the following two steps.

- 72 • In step 1, we show that a state dimension $s_{i,t}$ is in s_t^a , that is, it receives an edge from a_{t-1} , if
 73 and only if $s_{i,t} \not\perp a_{t-1} | \mathbf{s}_{t-1}$.
- 74 • In step 2, we show that s_t^a contains a minimally sufficient subset of state dimensions that can
 75 be directly controlled by actions.

76 **Step 1:** We first show that if a state dimension $s_{i,t}$ is in s_t^a , then $s_{i,t} \not\perp a_{t-1} | \mathbf{s}_{t-1}$.

77 We prove it by contradiction. Suppose that $s_{i,t}$ is independent of a_{t-1} given \mathbf{s}_{t-1} . Then according to
 78 the faithfulness assumption, we can see from the graph that $s_{i,t}$ does not receive an edge from a_{t-1} ,
 79 which contradicts the assumption, because, otherwise, \mathbf{s}_{t-1} cannot break the paths between $s_{i,t}$ and
 80 a_{t-1} which leads to the dependence.

81 We next show that if $s_{i,t} \not\perp a_{t-1} | \mathbf{s}_{t-1}$, then $s_{i,t} \in s_t^a$.

82 Similarly, by contradiction suppose that $s_{i,t}$ does not receive an edge from a_{t-1} . From the graph,
 83 it is easy to see that \mathbf{s}_{t-1} must break the path between $s_{i,t}$ and a_{t-1} . According to the Markov
 84 assumption, $s_{i,t}$ is independent of a_{t-1} given \mathbf{s}_{t-1} , which contradicts to the assumption. Since we have
 85 a contradiction, it must be that $s_{i,t}$ has an edge from a_{t-1} .

86 **Step 2:** In the previous steps, it has been shown that if a state dimension $s_{i,t}$ is in s_t^a , then $s_{i,t} \not\perp$
 87 $a_{t-1} | \mathbf{s}_{t-1}$, and if a state dimension $s_{i,t}$ is not in s_t^a , then $s_{i,t} \perp a_{t-1} | \mathbf{s}_{t-1}$. This implies that s_t^a is
 88 minimally sufficient for one-step direct control. \square

89 S1.3 Proof of Proposition 3

90 Furthermore, based on Proposition 1 and Proposition 2, we can further differentiate s_t^{ar} , $s_t^{\bar{ar}}$, $s_t^{\bar{r}}$ from
 91 s_t^r and s_t^a , which is given in the following proposition.

92 **Proposition 3.** Under the assumption that the graphical representation, corresponding to the
 93 environment model, is Markov and faithful to the measured data, we can build a connection between
 94 the graph structure and statistical independence of causal variables in the RL system, with (1) $s_{i,t} \in s_t^{ar}$
 95 if and only if $s_{i,t} \not\perp R_t | a_{t-1}, s_{t-1}^r$ and $s_{i,t} \not\perp a_{t-1} | \mathbf{s}_{t-1}$, (2) $s_{i,t} \in s_t^{\bar{ar}}$ if and only if $s_{i,t} \not\perp R_t | a_{t-1}, s_{t-1}^r$ and
 96 $s_{i,t} \perp a_{t-1} | \mathbf{s}_{t-1}$, (3) $s_{i,t} \in s_t^{\bar{r}}$ if and only if $s_{i,t} \perp R_t | a_{t-1}, s_{t-1}^r$ and $s_{i,t} \not\perp a_{t-1} | \mathbf{s}_{t-1}$, and (4) $s_{i,t} \in s_t^{\bar{r}}$ if
 97 and only if $s_{i,t} \perp R_t | a_{t-1}, s_{t-1}^r$ and $s_{i,t} \perp a_{t-1} | \mathbf{s}_{t-1}$.

98 *Proof.* This proposition can be easily proved by leveraging the results from Propositions 1 and 2. \square

99 S1.4 Proof of Theorem 1

100 According to the causal process in the RL system (as described in Eq.1 in [3]), we can build the
 101 following mapping from latent state variables \mathbf{s}_t to observed variables o_t and future cumulative reward
 102 R_t :

$$[o_t, R_t] = f(s_t^r, s_t^{\bar{r}}, \eta_t), \quad (1)$$

103 where

$$\begin{aligned} o_t &= f_1(s_t^r, s_t^{\bar{r}}), \\ R_t &= f_2(s_t^r, \eta_t). \end{aligned} \quad (2)$$

104 Here, note that to recover s_t^r , it is essential to take into account all future rewards $r_{t:T}$, because any
 105 state dimension $s_{i,t} \in \mathbf{s}_t$ that has a directed path to the future reward $r_{t+\tau}$, for $\tau > 0$, is involved in
 106 s_t^r . Hence, we consider the mapping from s_t^r to the future cumulative reward R_t , and η_t represents
 107 residuals, except s_t^r , that have an effect to R_t .

The following theorem shows that the different types of states s_t^{ar} , $s_t^{\bar{ar}}$, $s_t^{a\bar{r}}$, and $s_t^{\bar{a}\bar{r}}$ are blockwise identifiable from observed image variable o_t , reward variable r_t , and action variable a_t , under reasonable and weak assumptions.

Theorem 1. Suppose that the causal process in the RL system and the four categories of latent state variables can be described as that in Section 2 and illustrated in Figure 1(c). Under the following assumptions

A1. The mapping f in Eq. 1 is smooth and invertible with smooth inverse.

A2. For all $i \in \{1, \dots, d_o + d_R\}$ and $j \in \mathcal{F}_{i,:}$, there exist $\{\tilde{\mathbf{s}}_t^{(l)}\}_{l=1}^{|\mathcal{F}_{i,:}|}$, so that $\text{span}\{\mathbf{J}_f(\tilde{\mathbf{s}}_t^{(l)})_{i,:}\}_{l=1}^{|\mathcal{F}_{i,:}|} = \mathbb{R}_{\mathcal{F}_{i,:}}^{d_s}$, and there exists a matrix T with its support identical to that of $\mathbf{J}_{\hat{f}}^{-1}(\hat{\mathbf{s}}_t)\mathbf{J}_f(\tilde{\mathbf{s}}_t)$, so that $[\mathbf{J}_f(\tilde{\mathbf{s}}_t^{(l)})T]_{j,:} \in \mathbb{R}_{\mathcal{F}_{i,:}}^{d_s}$.

Then, reward-relevant and controllable states s_t^{ar} , reward-relevant but not controllable states $s_t^{\bar{ar}}$, reward-irrelevant but controllable states $s_t^{a\bar{r}}$, and noise $s_t^{\bar{a}\bar{r}}$, are blockwise identifiable.

In the theorem presented above, Assumption A1 only assumes the invertibility of function f , while functions f_1 and f_2 are considered general and not necessarily invertible. Since the function f is the mapping from all (latent) variables, including noise factors, that influence the observed variables, the invertibility assumption holds reasonably. However, note that it is not reasonable to assume the invertibility of the function f_2 since usually, the reward function is not invertible. Assumption A2, which is also given in [4, 5], aims to establish a more generic condition that rules out certain sets of parameters to prevent ill-posed conditions. Specifically, it ensures that the Jacobian is not partially constant. This condition is typically satisfied asymptotically, and it is necessary to avoid undesirable situations where the problem becomes ill-posed.

Proof. The proof consists of four steps.

1. In step 1, we show that $s_t^a = s_t^{ar} \cup s_t^{a\bar{r}}$ is blockwise identifiable, by using the characterization that the action variable a_t only directly influences s_t^{ar} and $s_t^{a\bar{r}}$.
2. In step 2, we show that $s_t^r = s_t^{ar} \cup s_t^{\bar{ar}}$ is blockwise identifiable, by using the characterization that the future cumulative reward R_t is only influenced by s_t^{ar} and $s_t^{\bar{ar}}$.
3. In step 3, we show that s_t^{ar} is blockwise identifiable, by using the identifiability of $s_t^{ar} \cup s_t^{a\bar{r}}$ and $s_t^{ar} \cup s_t^{\bar{ar}}$.
4. In step 4, we further show the blockwise identifiability of $s_t^{\bar{ar}}$, $s_t^{a\bar{r}}$, and $s_t^{\bar{a}\bar{r}}$.

Step 1: prove the block identifiability of s_t^a .

For simplicity of notation, below, we omit the subscript t .

Let $h := f^{-1} \circ \hat{f}$. We have

$$\hat{\mathbf{s}} = h(\mathbf{s}), \quad (3)$$

where $h = f^{-1} \circ \hat{f}$ is the transformation between the true latent variable and the estimated one, and $\hat{f} : \mathcal{S} \rightarrow \mathcal{X}$ denotes the estimated invertible generating function. Note that as both f^{-1} and \hat{f} are smooth and invertible, h and h^{-1} is smooth and invertible.

Since $h(\cdot)$ is smooth over \mathcal{S} , its Jacobian can be written as follows:

$$J_{h^{-1}} = \begin{bmatrix} A := \frac{\partial s^{\bar{a}}}{\partial \hat{s}^a} & B := \frac{\partial s^{\bar{a}}}{\partial \hat{s}^a} \\ C := \frac{\partial s^a}{\partial \hat{s}^a} & D := \frac{\partial s^a}{\partial \hat{s}^a} \end{bmatrix} \quad (4)$$

The invertibility of h^{-1} implies that $J_{h^{-1}}$ is full rank. Since s^a has changing distributions over the action variable a while $s^{\bar{a}}$ has invariant distributions over different values of a , we can derive that $C = 0$. Furthermore, because $J_{h^{-1}}$ is full rank and C is a zero matrix, D must be of full rank, which implies h_a^{-1} is invertible, where h_a^{-1} denotes the first derivative of h_a^{-1} . Therefore, s^a is blockwise identifiable up to invertible transformations.

150 **Step 2: prove the blockwise identifiability of s_t^r .**

151

152 Recall that we have the following mapping:

$$[o_t, R_t] = f(s_t^r, s_t^{\bar{r}}, \eta_t),$$

153 where

$$\begin{aligned} o_t &= f_1(s_t^r, s_t^{\bar{r}}), \\ R_t &= f_2(s_t^r, \eta_t). \end{aligned}$$

154 Note that here to recover s_t^r , we need to take into account all future rewards, because s_t^r contains
155 all those state dimensions that have a directed path to future rewards $r_{t+1:T}$. η_t represents all other
156 factors, except s_t^r , that influence R_t at time instance t . Further note here we assume the invertibility of
157 f , while f_1 and f_2 are general functions not necessarily invertible.

158 We denote by $\tilde{\mathbf{s}} = (s^r, s^{\bar{r}}, \eta)$. We further denote by the dimension of s^r by d_{s^r} , the dimension of $s^{\bar{r}}$ by
159 $d_{s^{\bar{r}}}$, the dimension of $\tilde{\mathbf{s}}$ by $d_{\tilde{\mathbf{s}}}$, the dimension of o by d_o , and the dimension of R_t by d_R .

160 We denote by \mathcal{F} the support of $\mathbf{J}_f(\mathbf{s})$, by $\hat{\mathcal{F}}$ the support of $\mathbf{J}_{\hat{f}}(\hat{\mathbf{s}})$, and by \mathcal{T} the support of $\mathbf{T}(\mathbf{s})$. We
161 also denote T as a matrix with the same support as \mathcal{T} . The proof technique is similar to that in [4, 5].

162 Since $h := \hat{f}^{-1} \circ f$, we have $\hat{f} = f \circ h^{-1}(\tilde{\mathbf{s}})$. By applying the chain rule repeatedly, we have

$$\mathbf{J}_{\hat{f}}(\hat{\mathbf{s}}) = \mathbf{J}_f(\tilde{\mathbf{s}}) \cdot \mathbf{J}_{h^{-1}}(h(\tilde{\mathbf{s}})). \quad (5)$$

163 With Assumption A2, for any $i \in \{1, \dots, d_o + d_R\}$, there exists $\{\tilde{\mathbf{s}}^{(l)}\}_{l=1}^{|\mathcal{F}_{i,:}|}$, s.t. $\text{span}(\{\mathbf{J}_f(\tilde{\mathbf{s}}^{(l)})_{i,:}\}_{l=1}^{|\mathcal{F}_{i,:}|}) = \mathbf{R}_{\mathcal{F}_{i,:}}^{d_{\tilde{\mathbf{s}}}}$.

164 Since $\{\mathbf{J}_f(\tilde{\mathbf{s}}^{(l)})_{i,:}\}_{l=1}^{|\mathcal{F}_{i,:}|}$ forms a basis of $\mathbf{R}_{\mathcal{F}_{i,:}}^{d_{\tilde{\mathbf{s}}}}$, for any $j_0 \in \mathcal{F}_{i,:}$, we can write canonical basis vector
165 $e_{j_0} \in \mathbf{R}_{\mathcal{F}_{i,:}}^{d_{\tilde{\mathbf{s}}}}$ as:

$$e_{j_0} = \sum_{l \in \mathcal{F}_{i,:}} \alpha_l \cdot \mathbf{J}_g(\tilde{\mathbf{s}}^{(l)})_{i,:}, \quad (6)$$

166 where $\alpha_l \in \mathbb{R}$ is a coefficient.

167 Then, following Assumption A2, there exists a deterministic matrix T such that

$$T_{j_0,:} = e_{j_0}^\top T = \sum_{l \in \mathcal{F}_{i,:}} \alpha_l \cdot \mathbf{J}_g(\tilde{\mathbf{s}}^{(l)})_{i,:} T \in \mathbf{R}_{\hat{\mathcal{F}}_{i,:}}^{d_{\tilde{\mathbf{s}}}}, \quad (7)$$

168 where \in is due to that each element in the summation belongs to $\mathbf{R}_{\hat{\mathcal{F}}_{i,:}}^{d_{\tilde{\mathbf{s}}}}$.

169 Therefore,

$$\forall j \in \mathcal{F}_{i,:}, T_{j,:} \in \mathbf{R}_{\hat{\mathcal{F}}_{i,:}}^{d_{\tilde{\mathbf{s}}}}.$$

170 Equivalently, we have:

$$\forall (i, j) \in \mathcal{F}, \quad \{i\} \times T_{j,:} \subset \hat{\mathcal{F}}. \quad (8)$$

171 We would like to show that \hat{s}^r does not depend on $s^{\bar{r}}$ and η , that is, $T_{i,j} = 0$ for $i \in \{1, \dots, d_{s^r}\}$ and
172 $j \in \{d_{s^r} + 1, \dots, d_{\tilde{\mathbf{s}}}\}$.

173 We prove it by contradiction. Suppose that \hat{s}^r had dependence on $s^{\bar{r}}$, that is, $\exists (j_{s^r}, j_{s^{\bar{r}}}) \in \mathcal{T}$ with
174 $j_{s^r} \in \{1, \dots, d_{s^r}\}$ and $j_{s^{\bar{r}}} \in \{d_{s^r} + 1, \dots, d_{s^r} + d_{s^{\bar{r}}}\}$.

175 Hence, there must exist $i_r \in \{d_o + 1, \dots, d_o + d_R\}$, such that, $(i_r, j_{s^{\bar{r}}}) \in \mathcal{F}$.

176 It follows from Equation 8 that:

$$\{i_r\} \times T_{j_{s^{\bar{r}}},:} \in \hat{\mathcal{F}} \implies (i_r, j_{s^r}) \in \hat{\mathcal{F}}. \quad (9)$$

177 However, due to the structure of \hat{f}_2 , $[\mathbf{J}_{\hat{f}_2}]_{i_r, j_{s^{\bar{r}}}} = 0$, which results in a contradiction. Therefore,
178 such (i_r, j_{s^r}) does not exist and \hat{s}^r does not depend on $s^{\bar{r}}$. The same reasoning implies that \hat{s}^r does
179 not dependent on η . Thus, \hat{s}^r does not depend on $(s^{\bar{r}}, \eta)$. In conclusion, \hat{s}^r does not contain extra
180 information beyond s^r .

181 Similarly, we can show that $(\hat{s}^{\bar{r}}, \hat{\eta})$ does not contain information of s^r .

182 Therefore, there is a one-to-one mapping between s^r and \hat{s}^r .

183 **Step 3: prove the blockwise identifiability of s_l^{ar} .**

184

185 In Step 1 and Step 2, we have shown that both s^a and s^r are blockwise identifiable. That is,

$$\begin{aligned}\hat{s}^r &= h_r(s^r), \\ \hat{s}^a &= h_a(s^a),\end{aligned}\tag{10}$$

186 where h_a and h_r are invertible functions.

187 According to the invariance relation of s^{ar} , We have the following relations:

$$\hat{s}^{ar} = h_r(s^r)_{1:d_{sar}} = h_a(s^a)_{1:d_{sar}}.\tag{11}$$

188 It remains to show that both $\tilde{h}_r := h_r(\cdot)_{1:d_{sar}}$ and $\tilde{h}_a := h_a(\cdot)_{1:d_{sar}}$ do not depend on $s^{\bar{ar}}$ and $s^{a\bar{r}}$ in their arguments.

190 We will prove this by contradiction. Without loss of generality, we suppose $\exists l \in \{1, \dots, d_{\bar{ar}}\}$,
191 $s^{r*} \in \mathcal{S}^r$, s.t., $\frac{\partial \tilde{h}_r}{\partial s_l^{ar}}(s^{r*}) \neq 0$. As h is smooth, it has continuous partial derivatives. Thus, $\frac{\partial \tilde{h}_r}{\partial s_l^{ar}} \neq 0$ holds
192 true in a neighbourhood of s^{r*} , i.e.,

$$\exists \eta > 0, \text{ s.t., } s_l^{\bar{ar}} \rightarrow \tilde{h}_r(s^{ar*}, (s_{-l}^{\bar{ar}*}, s_l^{\bar{ar}*})) \text{ is strictly monotonic on } (s_l^{\bar{ar}*} - \eta, s_l^{\bar{ar}*} + \eta),\tag{12}$$

193 where $s_{-l}^{\bar{ar}}$ denotes variable $s^{\bar{ar}}$ excluding the dimension l .

194 We further define an auxiliary function $\psi : \mathcal{S}^{ar} \times \mathcal{S}^{\bar{ar}} \times \mathcal{S}^{a\bar{r}} \rightarrow \mathbb{R}_{\geq 0}$ as follows:

$$\psi(s^{ar}, s^{\bar{ar}}, s^{a\bar{r}}) := |\tilde{h}_r(s^r) - \tilde{h}_a(s^a)|.\tag{13}$$

195 To obtain the contradiction to the invariance, it remains to show that $\psi > 0$ with a probability greater
196 than zero w.r.t. the true generating process.

197 There are two situations at $(s^{ar*}, s^{\bar{ar}*}, s^{a\bar{r}*})$ where $s^{\bar{ar}*}$ is an arbitrary point in $\mathcal{S}^{\bar{ar}}$:

- 198 • situation 1: $\psi(s^{ar*}, s^{\bar{ar}*}, s^{a\bar{r}*}) > 0$;
- 199 • situation 2: $\psi(s^{ar*}, s^{\bar{ar}*}, s^{a\bar{r}*}) = 0$.

200 In situation 1, we have identified a specific point $\psi(s^{ar*}, s^{\bar{ar}*}, s^{a\bar{r}*})$ that makes $\psi > 0$.

201 In situation 2, Eq. 12 implies that $\forall s_l^{\bar{ar}} \in (s_l^{\bar{ar}*}, s_l^{\bar{ar}*} + \eta)$

$$\psi(s^{ar*}, (s_{-l}^{\bar{ar}*}, s_l^{\bar{ar}}), s^{a\bar{r}*}) > 0.$$

202 Thus, in both situations, we can locate a point $(s^{ar*}, s^{\bar{ar}*'}, s^{a\bar{r}*})$ such that $\psi(s^{ar*}, s^{\bar{ar}*'}, s^{a\bar{r}*}) > 0$, where
203 $s^{\bar{ar}*'} = s^{\bar{ar}*}$ in case 1 and $s_l^{\bar{ar}*'} \in (s_l^{\bar{ar}*}, s_l^{\bar{ar}*} + \eta)$, $s_{-l}^{\bar{ar}*'} = s_{-l}^{\bar{ar}*}$ in situation 2.

204 Since ψ is a composition of continuous functions, it is continuous. As pre-image of open sets are
205 always open for continuous functions, the open set $\mathbb{R}_{>0}$ has an open set $\mathcal{U} \in \mathcal{S}^{ar} \times \mathcal{S}^{\bar{ar}} \times \mathcal{S}^{a\bar{r}}$ as
206 its preimage. Due to $(s^{ar*}, s^{\bar{ar}*'}, s^{a\bar{r}*}) \in \mathcal{U}$, \mathcal{U} is nonempty. As \mathcal{U} is nonempty and open, \mathcal{U} has a
207 Lebesgue measure of greater than zero.

208 As we assume that $p_{s^{ar}, s^{\bar{ar}}, s^{a\bar{r}}}$ is fully supported over the entire domain $\mathcal{S}^{ar} \times \mathcal{S}^{\bar{ar}} \times \mathcal{S}^{a\bar{r}}$, we can deduce
209 that $\mathbb{P}_p[\mathcal{U}] > 0$. That is, $\psi > 0$ with a probability greater than zero, which contradicts the invariance
210 condition. Therefore, it has been shown that $\hat{h}_r(s^r)$ does not depend on $s^{\bar{ar}}$. This proof technique is
211 related to that in [6].

212 Similarly, we can show that $\hat{h}_a(s^a)$ does not depend on $s^{a\bar{r}}$.

213 Finally, the smoothness and invertibility of \hat{h}_r and \hat{h}_a follow from the smoothness and invertibility of
214 h_r and h_a over the entire domain.

215 Therefore, $h_r(h_a)$ is a smooth invertible mapping between s^{ar} and \hat{s}^{ar} . That is, s^{ar} is blockwise
216 invertible.

217 **Step 4: prove the blockwise identifiability of $s_t^{\bar{a}r}$, $s_t^{a\bar{r}}$, and $s_t^{\bar{a}\bar{r}}$.**

218

219 We can use the same technique in Step 3 to show the identifiability of $s_t^{\bar{a}r}$ and $s_t^{a\bar{r}}$. Specifically,
 220 since s_r and s^{ar} are identifiable, we can show that $s_t^{\bar{a}r}$ is identifiable. Similarly, since s^a and s^{ar} are
 221 identifiable, we can show that $s_t^{a\bar{r}}$ is identifiable. Furthermore, since s^{ar} , $s_t^{\bar{a}r}$, and $s_t^{a\bar{r}}$ are identifiable,
 222 we can show that $s_t^{\bar{a}\bar{r}}$ is identifiable \square

223 S2 Derivation of the objective function

224 We start by defining the components of the world mode as follows:

$$\begin{cases} \text{Observation Model:} & p_\theta(o_t | \mathbf{s}_t) \\ \text{Reward Model:} & p_\theta(r_t | s_t^r) \\ \text{Transition Model:} & p_\gamma(\mathbf{s}_t | \mathbf{s}_{t-1}, a_{t-1}) \\ \text{Representation Model:} & q_\phi(\mathbf{s}_t | o_t, \mathbf{s}_{t-1}, a_{t-1}) \end{cases} \quad (14)$$

225 The latent dynamics can be disentangled into four catogories:

$$\begin{array}{ll} \text{Disentangled Transition Model:} & \text{Disentangled Representation Model:} \\ \left\{ \begin{array}{l} p_{\gamma_1}(s_t^{ar} | s_{t-1}^r, a_{t-1}) \\ p_{\gamma_2}(s_t^{\bar{a}r} | s_{t-1}^r) \\ p_{\gamma_3}(s_t^{a\bar{r}} | \mathbf{s}_{t-1}, a_{t-1}) \\ p_{\gamma_4}(s_t^{\bar{a}\bar{r}} | \mathbf{s}_{t-1}) \end{array} \right. & \left\{ \begin{array}{l} q_{\phi_1}(s_t^{ar} | o_t, s_{t-1}^r, a_{t-1}) \\ q_{\phi_2}(s_t^{\bar{a}r} | o_t, s_{t-1}^r) \\ q_{\phi_3}(s_t^{a\bar{r}} | o_t, \mathbf{s}_{t-1}, a_{t-1}) \\ q_{\phi_4}(s_t^{\bar{a}\bar{r}} | o_t, \mathbf{s}_{t-1}) \end{array} \right. \end{array} \quad (15)$$

226 We define the information bottleneck objective for latent dynamics models [7, 8]

$$\max I(\mathbf{s}_{1:T}; (o_{1:T}, r_{1:T}) | a_{1:T}) - \beta \cdot I(\mathbf{s}_{1:T}, i_{1:T} | a_{1:T}), \quad (16)$$

227 where β is scalar and i_t are dataset indices that determine the observations $p(o_t | i_t) = \delta(o_t - \bar{o}_t)$ as in
 228 [9].

229 Maximizing the objective leads to model states that can predict the sequence of observations and
 230 rewards while limiting the amount of information extracted at each time step. We derive the lower
 231 bound of the first term in Equation 16:

$$\begin{aligned} & I(\mathbf{s}_{1:T}; (o_{1:T}, r_{1:T}) | a_{1:T}) \\ &= \mathbb{E}_{q(o_{1:T}, r_{1:T}, \mathbf{s}_{1:T}, a_{1:T})} \left(\sum_t \ln p(o_{1:T}, r_{1:T} | \mathbf{s}_{1:T}, a_{1:T}) - \ln p(o_{1:T}, r_{1:T} | a_{1:T}) \right) \\ & \pm \mathbb{E} \left(\sum_t \ln p(o_{1:T}, r_{1:T} | \mathbf{s}_{1:T}, a_{1:T}) \right) \\ & \geq \mathbb{E} \left(\sum_t \ln p(o_{1:T}, r_{1:T} | \mathbf{s}_{1:T}, a_{1:T}) \right) - \text{KL} \left(p(o_{1:T}, r_{1:T} | \mathbf{s}_{1:T}, a_{1:T}) \parallel \prod_t p_\theta(o_t | s_t) p_\theta(r_t | s_t^r) \right) \\ &= \mathbb{E} \left(\sum_t \ln p_\theta(o_t | \mathbf{s}_t) + \ln p_\theta(r_t | s_t^r) \right). \end{aligned} \quad (17)$$

232 Thus, we obtain the objective function:

$$\mathcal{J}_O^t = \ln p_\theta(o_t | \mathbf{s}_t) \quad \mathcal{J}_R^t = \ln p_\theta(r_t | s_t^r) \quad (18)$$

For the second term in Equation 16, we use the non-negativity of the KL divergence to obtain an upper bound,

$$\begin{aligned}
& I(\mathbf{s}_{1:T}; i_{1:T} \mid a_{1:T}) \\
&= \mathbb{E}_{q(o_{1:T}, r_{1:T}, \mathbf{s}_{1:T}, a_{1:T}, i_{1:T})} \left(\sum_t \ln q(\mathbf{s}_t \mid \mathbf{s}_{t-1}, a_{t-1}, i_t) - \ln p(\mathbf{s}_t \mid \mathbf{s}_{t-1}, a_{t-1}) \right) \\
&= \mathbb{E} \left(\sum_t \ln q_\phi(\mathbf{s}_t \mid \mathbf{s}_{t-1}, a_{t-1}, o_t) - \ln p(\mathbf{s}_t \mid \mathbf{s}_{t-1}, a_{t-1}) \right) \\
&\leq \mathbb{E} \left(\sum_t \ln q_\phi(\mathbf{s}_t \mid \mathbf{s}_{t-1}, a_{t-1}, o_t) - \ln p_\gamma(\mathbf{s}_t \mid \mathbf{s}_{t-1}, a_{t-1}) \right) \\
&= \mathbb{E} \left(\sum_t \text{KL}(q_\phi(\mathbf{s}_t \mid \mathbf{s}_{t-1}, a_{t-1}, o_t) \parallel p_\gamma(\mathbf{s}_t \mid \mathbf{s}_{t-1}, a_{t-1})) \right).
\end{aligned} \tag{19}$$

According to equation 15, we have $p_\gamma = p_{\gamma_1} \cdot p_{\gamma_2} \cdot p_{\gamma_3} \cdot p_{\gamma_4}$ and $q_\phi = q_{\phi_1} \cdot q_{\phi_2} \cdot q_{\phi_3} \cdot q_{\phi_4}$.

$$\begin{aligned}
\text{KL}(q_\phi \parallel p_\gamma) &= \text{KL}(q_{\phi_1} \cdot q_{\phi_2} \cdot q_{\phi_3} \cdot q_{\phi_4} \parallel p_{\gamma_1} \cdot p_{\gamma_2} \cdot p_{\gamma_3} \cdot p_{\gamma_4}) \\
&= \mathbb{E}_{q_\phi} \left(\ln \frac{q_{\phi_1}}{p_{\gamma_1}} + \ln \frac{q_{\phi_2}}{p_{\gamma_2}} + \ln \frac{q_{\phi_3}}{p_{\gamma_3}} + \ln \frac{q_{\phi_4}}{p_{\gamma_4}} \right) \\
&= \text{KL}(q_{\phi_1} \parallel p_{\gamma_1}) + \text{KL}(q_{\phi_2} \parallel p_{\gamma_2}) + \text{KL}(q_{\phi_3} \parallel p_{\gamma_3}) + \text{KL}(q_{\phi_4} \parallel p_{\gamma_4})
\end{aligned} \tag{20}$$

We introduce additional hyperparameters to regulate the amount of information contained within each category of variables:

$$\mathcal{J}_D^t = -\beta_1 \text{KL}(q_{\phi_1} \parallel p_{\gamma_1}) - \beta_2 \text{KL}(q_{\phi_2} \parallel p_{\gamma_2}) - \beta_3 \text{KL}(q_{\phi_3} \parallel p_{\gamma_3}) - \beta_4 \text{KL}(q_{\phi_4} \parallel p_{\gamma_4}). \tag{21}$$

Additionally, we introduce two supplementary objectives to explicitly capture the distinctive characteristics of the four distinct representation categories. Specifically, we characterize the reward-relevant representations by measuring the dependence between s_t^r and R_t , given a_{t-1} and s_{t-1}^r , that is $I(s_t^r, R_t \mid a_{t-1}, s_{t-1}^r)$ (see Figure 7(a)). Note that if there exists a directed edge from s_t^r to a_t in the graphical model, a_t should also be conditioned). To ensure that s_t^r are minimally sufficient for policy training, we maximize $I(s_t^r, R_t \mid a_{t-1}, s_{t-1}^r)$ while minimizing $I(s_t^r, R_t \mid a_{t-1}, s_{t-1}^r)$ to discourage the inclusion of redundant information in s_t^r concerning the rewards:

$$I(s_t^r; R_t \mid a_{t-1}, s_{t-1}^r) - I(s_t^{\bar{r}}; R_t \mid a_{t-1}, s_{t-1}^r). \tag{22}$$

The conditional mutual information can be expressed as the disparity between two mutual information values.

$$\begin{aligned}
I(s_t^r; R_t \mid a_{t-1}, s_{t-1}^r) &= I(R_t; s_t^r, a_{t-1}, s_{t-1}^r) - I(R_t; a_{t-1}, s_{t-1}^r), \\
I(s_t^{\bar{r}}; R_t \mid a_{t-1}, s_{t-1}^r) &= I(R_t; s_t^{\bar{r}}, a_{t-1}, s_{t-1}^r) - I(R_t; a_{t-1}, s_{t-1}^r).
\end{aligned} \tag{23}$$

Combining the above two equations, we eliminated the identical terms, ultimately yielding the following formula

$$I(R_t; s_t^r, a_{t-1}, s_{t-1}^r) - I(R_t; s_t^{\bar{r}}, a_{t-1}, s_{t-1}^r). \tag{24}$$

We use the Donsker-Varadhan representation to express mutual information as a supremum over functions,

$$\begin{aligned}
I(X; Y) &= D_{KL}(p(x, y) \parallel p(x)p(y)) \\
&= \sup_{T \in \mathcal{T}} \mathbb{E}_{p(x, y)}[T(x, y)] - \log \mathbb{E}_{p(x)p(y)}[e^{T(x, y)}].
\end{aligned} \tag{25}$$

We employ mutual information neural estimation [10] to approximate the mutual information value. We represent the function T using a neural network that accepts variables (x, y) as inputs and is parameterized by α . The neural network is optimized through stochastic gradient ascent to find the supremum. Substituting x and y with variables defined in Equation 24, our objective is reformulated as follows:

$$\mathcal{J}_{\text{RS}}^t = \lambda_1 \cdot \{I_{\alpha_1}(R_t; s_t^r, a_{t-1}, \mathbf{sg}(s_{t-1}^r)) - I_{\alpha_2}(R_t; s_t^{\bar{r}}, a_{t-1}, \mathbf{sg}(s_{t-1}^r))\}. \tag{26}$$

256 To incorporate the conditions from the original objective, we apply the stop_gradient operation to
 257 the variable s_{t-1}^r . Similarly, to ensure that the representations s_t^a are directly controllable by actions,
 258 while $s_t^{\bar{a}}$ are not, we maximize the following objective:

$$I(s_t^a; a_{t-1} | \mathbf{s}_{t-1}) - I(s_t^{\bar{a}}, a_{t-1} | \mathbf{s}_{t-1}), \quad (27)$$

259 By splitting the conditional mutual information and eliminating identical terms, we obtain the
 260 following objective function:

$$\mathcal{J}_{AS}^t = \lambda_2 \cdot \{I_{\alpha_3}(a_{t-1}; s_t^a, \mathbf{sg}(\mathbf{s}_{t-1})) - I_{\alpha_4}(a_{t-1}; s_t^{\bar{a}}, \mathbf{sg}(\mathbf{s}_{t-1}))\}. \quad (28)$$

261 where $\alpha_1, \alpha_2, \alpha_3, \alpha_4$ can be obtained by maximizing Equation 25. Intuitively, these two objective
 262 functions ensure that s_t^r is predictive of the reward, while $s_t^{\bar{r}}$ is not; similarly, s_t^a can be predicted by
 263 the action, whereas $s_t^{\bar{a}}$ cannot.

264 Combine the equation 18, equation 21, equation 26 and equation 28, the total objective function is:

$$\begin{aligned} \mathcal{J}_{TOTAL} &= \max_{\phi, \theta, \gamma, \alpha_1, \alpha_3} \min_{\alpha_2, \alpha_4} \mathbb{E}_{q_\phi} \left(\sum_t (\mathcal{J}_O^t + \mathcal{J}_R^t + \mathcal{J}_D^t + \mathcal{J}_{RS}^t + \mathcal{J}_{AS}^t) \right) + \text{const} \\ &= \max_{\phi, \theta, \gamma, \alpha_1, \alpha_3} \min_{\alpha_2, \alpha_4} \mathbb{E}_{q_\phi} \{ \log p_\theta(o_t | \mathbf{s}_t) + \log p_\theta(r_t | s_t^r) \\ &\quad - \sum_{i=1}^4 \beta_i \cdot \text{KL}(q_{\phi_i} \| p_{\gamma_i}) + \lambda_1 \cdot (I_{\alpha_1} - I_{\alpha_2}) + \lambda_2 \cdot (I_{\alpha_3} - I_{\alpha_4}) \} + \text{const}. \end{aligned} \quad (29)$$

265 The expectation is computed over the dataset and the representation model. Throughout the model
 266 learning process, the objectives for estimating mutual information and learning the world model are
 267 alternately optimized.

268 S2.1 Discussions

269 In this subsection, we examine the mutual information constraints in equation 26 and equation 28 and
 270 their relationship with other objectives. Our findings reveal that while other objectives partially fulfill
 271 the desired functionality of the mutual information constraints, incorporating both mutual information
 272 objectives is essential for certain environments.

273 The objective functions can be summarized as follows::

$$\begin{aligned} \mathcal{J}_O^t &= \ln p_\theta(o_t | \mathbf{s}_t), \quad \mathcal{J}_R^t = \ln p_\theta(r_t | s_t^r), \quad \mathcal{J}_D^t = -\text{KL}(q_\phi \| p_\gamma), \\ \mathcal{J}_{RS}^t &= \lambda_1 \cdot \{I_{\alpha_1}(R_t; s_t^r, a_{t-1}, \mathbf{sg}(s_{t-1}^r)) - I_{\alpha_2}(R_t; s_t^{\bar{r}}, a_{t-1}, \mathbf{sg}(s_{t-1}^{\bar{r}}))\}, \\ \mathcal{J}_{AS}^t &= \lambda_2 \cdot \{I_{\alpha_3}(a_{t-1}; s_t^a, \mathbf{sg}(\mathbf{s}_{t-1})) - I_{\alpha_4}(a_{t-1}; s_t^{\bar{a}}, \mathbf{sg}(\mathbf{s}_{t-1}))\}, \end{aligned} \quad (30)$$

274 and the KL divergence term can be further decomposed into 4 components:

$$\begin{aligned} \mathcal{J}_{D_1}^t &= -\beta_1 \cdot \text{KL}(q_{\phi_1}(s_t^{ar} | o_t, s_{t-1}^r, a_{t-1}) \| p_{\gamma_1}(s_t^{ar} | s_{t-1}^r, a_{t-1})) \\ \mathcal{J}_{D_2}^t &= -\beta_2 \cdot \text{KL}(q_{\phi_2}(s_t^{\bar{ar}} | o_t, s_{t-1}^r) \| p_{\gamma_2}(s_t^{\bar{ar}} | s_{t-1}^r)) \\ \mathcal{J}_{D_3}^t &= -\beta_3 \cdot \text{KL}(q_{\phi_3}(s_t^{a\bar{r}} | o_t, \mathbf{s}_{t-1}, a_{t-1}) \| p_{\gamma_3}(s_t^{a\bar{r}} | \mathbf{s}_{t-1}, a_{t-1})) \\ \mathcal{J}_{D_4}^t &= -\beta_4 \cdot \text{KL}(q_{\phi_4}(s_t^{\bar{a}\bar{r}} | o_t, \mathbf{s}_{t-1}) \| p_{\gamma_4}(s_t^{\bar{a}\bar{r}} | \mathbf{s}_{t-1})). \end{aligned} \quad (31)$$

275 Specifically, maximizing I_{α_1} in \mathcal{J}_{RS}^t enhances the predictability of R_t based on the current state s_{t-1}^r
 276 conditioning on (s_{t-1}^r, a_{t-1}) . However, notice that this objective can be partially accomplished by
 277 optimizing \mathcal{J}_R^t . When learning the world model, both the transition function and the reward function
 278 are trained: the reward function predicts the current reward r_t using s_t^r , while the transition model
 279 predicts the next state. These combined predictions contribute to the overall prediction of R_t .

280 Minimizing I_{α_2} in \mathcal{J}_{RS}^t eliminates extraneous reward-related information present in $s_t^{\bar{r}}$. According to
 281 our formulation, $s_t^{\bar{r}}$ can still be predictive of R_t as long as it does not introduce additional predictability

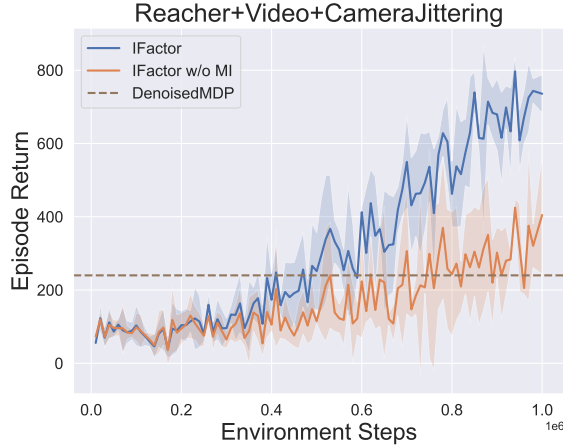


Figure 1: Ablation of the mutual information constraints in the Reacher environment with video background and jittering camera. The dashed brown line illustrates the policy performance of Denoised MDP after 1 million environment steps.

beyond what is already captured by (s_{t-1}^r, a_{t-1}) . This is because we only assume that $s_t^{\bar{r}}$ is **conditionally** independent from R_t when conditioning on s_{t-1}^r . If we don't condition on s_{t-1}^r , it introduces s_{t-1}^r as the confounding factor between $s_t^{\bar{r}}$ and R_t , establishing association between $s_t^{\bar{r}}$ and R_t (refer to Figure 7(a)). Note that the KL divergence constraints govern the information amount within each state category. By amplifying the weight of the KL constraints on $s_t^{\bar{r}}$, the value of I_{α_2} can indirectly be diminished.

By maximizing I_{α_3} and minimizing I_{α_4} in \mathcal{J}_{AS}^t , we ensure that s_t^a can be predicted based on $(a_{t-1}, \mathbf{s}_{t-1})$ while $s_t^{\bar{a}}$ cannot. The KL constraints on s_t^{ar} and $s_t^{a\bar{r}}$ incorporate the action a_{t-1} into the prior and posterior, implicitly requiring that s_t^a should be predictable given a_{t-1} . Conversely, the KL constraints on $s_t^{\bar{a}r}$ and $s_t^{\bar{a}\bar{r}}$ do not include the action a_{t-1} in the prior and posterior, implicitly requiring that s_t^a should not be predictable based on a_{t-1} . However, relying solely on indirect constraints can sometimes be ineffective, as it may lead to entangled representations that negatively impact policy performance (see Figure 5).

Ablation of the mutual information constraints. The inclusion of both \mathcal{J}_{RS}^t and \mathcal{J}_{AS}^t is essential in certain environments to promote disentanglement and enhance policy performance, despite sharing some common objectives. We have observed improved training stability in the variant of Robodesk environment (see Figure 4) and significant performance gains in the Reacher environment with video background and camera jittering (see Figure 1). When two mutual information objectives are removed, we notice that entangled representations emerged in these environments, as depicted in Figure 5. We assign values of 0.1 to λ_1 and λ_2 in the environment of modified Cartpole, variant of Robodesk and Reacher with video background and jittering camera. Empirically, a value of 0.1 has been found to be preferable for both λ_1 and λ_2 . Using a higher value for regularization might negatively impact the learning of representation and transition model. In other DMC environments, the ELBO loss alone has proven effective due to the inherent structure of our disentangled latent dynamics. The choice of hyperparameters $(\beta_1, \beta_2, \beta_3, \beta_4)$ depends on the specific goals of representation learning and the extent of noise interference in the task. If the objective is to accurately recover the true latent variables for understanding the environment, it is often effective to assign equal weights to the four KL divergence terms (for experiments on synthetic data and modified cartpole). When the aim is to enhance policy training stability by mitigating noise, it is recommended to set the values of β_1 and β_2 higher than β_3 and β_4 (for experiments on variants of Robodesk and DMC). Moreover, in environments with higher levels of noise, it is advisable to increase the discrepancy in values between the hyperparameters.

S3 Environment descriptions

S3.1 Synthetic data

For the sake of simplicity, we consider one lag for the latent processes in Section 4. Our identifiability proof can actually be applied for arbitrary lags directly because the identifiability does not rely on the number of previous states. We extend the latent dynamics of the synthetic environment to incorporate a general time-delayed causal effect with $\tau \geq 1$ in the synthetic environment. When $\tau = 1$, it reduces to a common MDP. The ground-truth generative model of the environment is as follows::

$$\left\{ \begin{array}{ll} \text{Observation Model:} & p_{\theta}(o_t | \mathbf{s}_t) \\ \text{Reward Model:} & p_{\theta}(r_t | s_t^r) \\ \text{Transition Model:} & p_{\gamma}(\mathbf{s}_t | \mathbf{s}_{t-\tau:t-1}, a_{t-\tau:t-1}) \end{array} \right. \quad \text{Transition : } \left\{ \begin{array}{l} p_{\gamma_1}(s_t^{ar} | s_{t-\tau:t-1}^r, a_{t-\tau:t-1}) \\ p_{\gamma_2}(s_t^{\bar{ar}} | s_{t-\tau:t-1}^r) \\ p_{\gamma_3}(s_t^{\bar{ar}} | \mathbf{s}_{t-\tau:t-1}, a_{t-\tau:t-1}) \\ p_{\gamma_4}(s_t^{\bar{ar}} | \mathbf{s}_{t-\tau:t-1}) \end{array} \right. \quad (32)$$

Data Generation We generate synthetic datasets with 100, 000 data points according to the generating process in Equation 32, which satisfies the identifiability conditions stated in Theorem 1. The latent variables \mathbf{s}_t have 8 dimensions, where $s_t^{ar} = s_t^{\bar{ar}} = s_t^{\bar{ar}} = s_t^{\bar{ar}} = 2$. At each timestep, a one-hot action of dimension 5, denoted as a_t , is taken. The lag number of the process is set to $\tau = 2$. The observation model $p_{\theta}(o_t | \mathbf{s}_t)$ is implemented using a random three-layer MLPs with LeakyReLU units. The reward model $p_{\theta}(r_t | s_t^r)$ is represented by a random one-layer MLP. It's worth noting that the reward model is not invertible due to the scalar nature of r_t . Four distinct transition functions, namely p_{γ_1} , p_{γ_2} , p_{γ_3} , and p_{γ_4} , are employed and modeled using random one-layer MLP with LeakyReLU units. The process noise is sampled from an i.i.d. Gaussian distribution with a standard deviation of $\sigma = 0.1$. To simulate nonstationary noise for various latent variables in RL, the process noise terms are coupled with the historical information by multiplying them with the average value of all the time-lagged latent variables, as suggested in [11].

S3.2 Modified Cartpole

We have modified the original Cartpole environment by introducing two distractors. The first distractor is an uncontrollable Cartpole located in the upper portion of the image, which does not affect the rewards. The second distractor is a controllable green light positioned below the reward-relevant Cartpole in the lower part of the image, but it is not associated with any rewards. The task-irrelevant cartpole undergoes random actions at each time step and stops moving when its angle exceeds 45 degrees or goes beyond the screen boundaries. The action space consists of three independent degrees of freedom: direction (left or right), force magnitude (10N or 20N), and green light intensity (lighter or darker). This results in an 8-dimensional one-hot vector. The objective of this variant is to maintain balance for the reward-relevant cartpole by applying suitable forces.

S3.3 Variant of Robodesk

The RoboDesk environment with noise distractors [12] is a control task designed to simulate realistic sources of noise, such as flickering lights and shaky cameras. Within the environment, there is a large TV that displays natural RGB videos. On the desk, there is a green button that controls both the hue of the TV and a light. The agent's objective is to manipulate this button in order to change the TV's

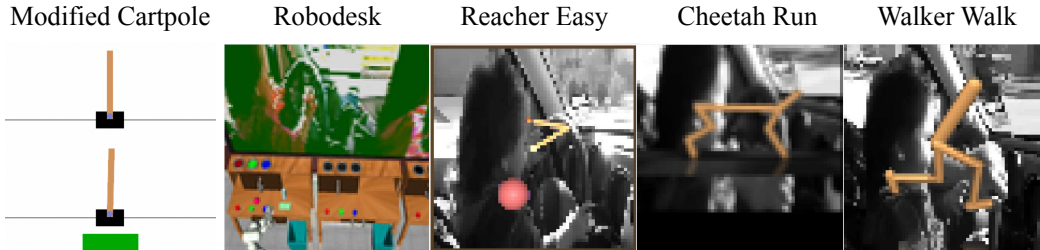


Figure 2: Visualization of the environments used in our experiments.

	Ctrl + Rew	$\overline{\text{Ctrl}} + \text{Rew}$	Ctrl + $\overline{\text{Rew}}$	$\overline{\text{Ctrl}} + \overline{\text{Rew}}$
Modified Cartpole	Agent	Agent	Green Light	Distractor cartpole
Robodesk	Agent, Button, Light on desk	TV content, Button sensor noise	Blocks on desk, Handle on desk, Other movable objects (Green hue of TV)	Jittering and flickering environment lighting, Jittering camera
DMC	Noiseless	Agent	(Agent)	—
	Video Background	Agent	(Agent)	Background
	Video Background + Noisy Sensor	Agent	(Agent) Background	—
	Video Background + Camera Jittering	Agent	(Agent)	Background Jittering camera

Table 1: Categorization of various types of information in the environments we evaluated. We use the red color to emphasize the categorization difference between IFactor and Denoised MDP. Unlike Denoised MDP that assumes independent latent processes, IFactor allows for causally-related processes. Therefore, in this paper, the term "controllable" refers specifically to one-step controllability, while "reward-relevant" is characterized by the conditional dependence between s_t^* and the cumulative reward variable R_t when conditioning on (s_{t-1}, a_{t-1}) . Following this categorization, certain agent information can be classified as (one-step) uncontrollable (including indirectly controllable and uncontrollable factors) but reward-relevant factors, such as some position information determined by the position and velocity in the previous time-step rather than the action. On the other hand, the green hue of TV in Robodesk is classified as controllable but reward-irrelevant factors, as they are independent of the reward given the state of the robot arm and green button, aligning with the definition of $s_t^{a\bar{r}}$.

349 hue to green. The agent’s reward is determined based on the greenness of the TV image. In this
350 environment, all four types of information are present (see Table 1).

351 S3.4 Variants of DeepMind Control Suite

352 Four variants [12] are introduced for each DMC task:

- 353 • **Noiseless:** Original environment without distractors.
- 354 • **Video Background:** Replacing noiseless background with natural videos [13] ($\overline{\text{Ctrl}} + \overline{\text{Rew}}$).
- 355 • **Video Background + Sensor Noise:** Imperfect sensors sensitive to intensity of a background
356 patch ($\text{Ctrl} + \text{Rew}$).
- 357 • **Video Background + Camera Jittering:** Shifting the observation by a smooth random
358 walk ($\overline{\text{Ctrl}} + \overline{\text{Rew}}$).

359 The video background in the environment incorporates grayscale videos from Kinetics-400, where
360 pixels with high blue channel values are replaced. Camera jittering is introduced through a smooth
361 random walk shift using Gaussian-perturbing acceleration, velocity decay, and pulling force. Sensor
362 noise is added by perturbing a specific sensor based on the intensity of a patch in the background video.
363 The perturbation involves adding the average patch value minus 0.5. Different sensors are perturbed
364 for different environments. These sensor values undergo non-linear transformations, primarily piece-
365 wise linear, to compute rewards. While the additive reward noise model may not capture sensor
366 behavior perfectly, it is generally sufficient as long as the values remain within moderate ranges and
367 stay within one linear region. (Note: the variants of Robodesk and DMC are not the contributions
368 of this paper. We kindly refer readers to the paper of Denoised MDP [12] for a more detailed
369 introduction.)

S4 Experimental Details

Computing Hardware We used a machine with the following CPU specifications: Intel(R) Xeon(R) Silver 4110 CPU @ 2.10GHz; 32 CPUs, eight physical cores per CPU, a total of 256 logical CPU units. The machine has two GeForce RTX 2080 Ti GPUs with 11GB GPU memory.

Reproducibility We’ve included the code for the framework and all experiments in the supplement. We plan to release our code under the MIT License after the paper review period.

S4.1 Synthetic Dataset

Hyperparameter Selection and Network Structure We adopt a similar experimental setup to TDRL [11], while extending it by decomposing the dynamics into four causally related latent processes proposed in this paper (refer to Equation 32). For all experiments, we assign $\beta_1 = \beta_2 = \beta_3 = \beta_4 = 0.003$ as the weights for the KL divergence terms. In this particular experiment, we set λ_1 and λ_2 to 0 because the utilization of the ELBO loss alone has effectively maximized J'_{RS} and J'_{AS} , as illustrated in Figure 3. Here, J'_{RS} represents $I_{\alpha_1} - I_{\alpha_2}$, and J'_{AS} represents $I_{\alpha_3} - I_{\alpha_4}$. The network structure employed in this experiment is presented in Table 2.

Training Details The models are implemented in PyTorch 1.13.1. The VAE network is trained using AdamW optimizer for 100 epochs. A learning rate of 0.001 and a mini-batch size of 64 are used. We have used three random seeds in each experiment and reported the mean performance with standard deviation averaged across random seeds.

Table 2: Architecture details. BS: batch size, T: length of time series, o_dim: observation dimension, s_dim: latent dimension, s_t^{ar} _dim: latent dimension for s_t^{ar} , $s_t^{\bar{ar}}$ _dim: latent dimension for $s_t^{\bar{ar}}$, s_t^{ar} _dim: latent dimension for s_t^{ar} , $s_t^{\bar{ar}}$ _dim: latent dimension for $s_t^{\bar{ar}}$ ($s_dim = s_t^{ar_dim} + s_t^{\bar{ar}}_dim + s_t^{ar_dim} + s_t^{\bar{ar}}_dim$), LeakyReLU: Leaky Rectified Linear Unit.

Configuration	Description	Output
1. MLP-Obs-Encoder Observation Encoder for Synthetic Data		
Input: $o_{1:T}$	Observed time series	$BS \times T \times o_dim$
Dense	128 neurons, LeakyReLU	$BS \times T \times 128$
Dense	128 neurons, LeakyReLU	$BS \times T \times 128$
Dense	128 neurons, LeakyReLU	$BS \times T \times 128$
Dense	Temporal embeddings	$BS \times T \times s_dim$
2. MLP-Obs-Decoder Observation Decoder for Synthetic Data		
Input: $\hat{s}_{1:T}$	Sampled latent variables	$BS \times T \times s_dim$
Dense	128 neurons, LeakyReLU	$BS \times T \times 128$
Dense	128 neurons, LeakyReLU	$BS \times T \times 128$
Dense	o_dim neurons, reconstructed $\hat{o}_{1:T}$	$BS \times T \times o_dim$
3. MLP-Reward-Decoder Reward Decoder for Synthetic Data		
Input: $\hat{s}_{1:T}$	Sampled latent variables	$BS \times T \times s_dim$
Dense	1 neurons, LeakyReLU	$BS \times T \times 1$
4. Disentangled Prior for s_t^{ar} Nonlinear Transition Prior Network		
Input	Sampled latents and actions $s_{1:T}^r, a_{1:T}$	$BS \times T \times (s_t^r_dim + a_dim)$
Dense	s_t^{ar} _dim neurons, prior output	$BS \times T \times s_t^{ar_dim}$
5. Disentangled Prior for $s_t^{\bar{ar}}$ Nonlinear Transition Prior Network		
Input	Sampled latent variable sequence $s_{1:T}^r$	$BS \times T \times s_t^r_dim$
Dense	$s_t^{\bar{ar}}$ _dim neurons, prior output	$BS \times T \times s_t^{\bar{ar}}_dim$
6. Disentangled Prior for $s_t^{ar\bar{r}}$ Nonlinear Transition Prior Network		
Input	Sampled latents and actions $s_{1:T}, a_{1:T}$	$BS \times T \times (s_dim + a_dim)$
Dense	$s_t^{ar\bar{r}}$ _dim neurons, prior output	$BS \times T \times s_t^{ar\bar{r}}_dim$
7. Disentangled Prior for $s_t^{\bar{ar}\bar{r}}$ Nonlinear Transition Prior Network		
Input	Sampled latent variable sequence $s_{1:T}$	$BS \times T \times s_dim$
Dense	$s_t^{\bar{ar}\bar{r}}$ _dim neurons, prior output	$BS \times T \times s_t^{\bar{ar}\bar{r}}_dim$

388 S4.1.1 Extra Results.

389 During the training process, we record the estimation value of four mutual information (MI) terms.
 390 The corresponding results are presented in Figure 3. Despite not being explicitly incorporated into
 391 the objective function, the terms $I_{\alpha_1} - I_{\alpha_2}$ and $I_{\alpha_3} - I_{\alpha_4}$ exhibit significant maximization. Furthermore,
 392 the estimation values of I_{α_2} and I_{α_4} are found to be close to 0. These findings indicate that the state
 393 variable s_t^T contains little information about the reward, and the predictability of s_t^T by the action is
 also low.

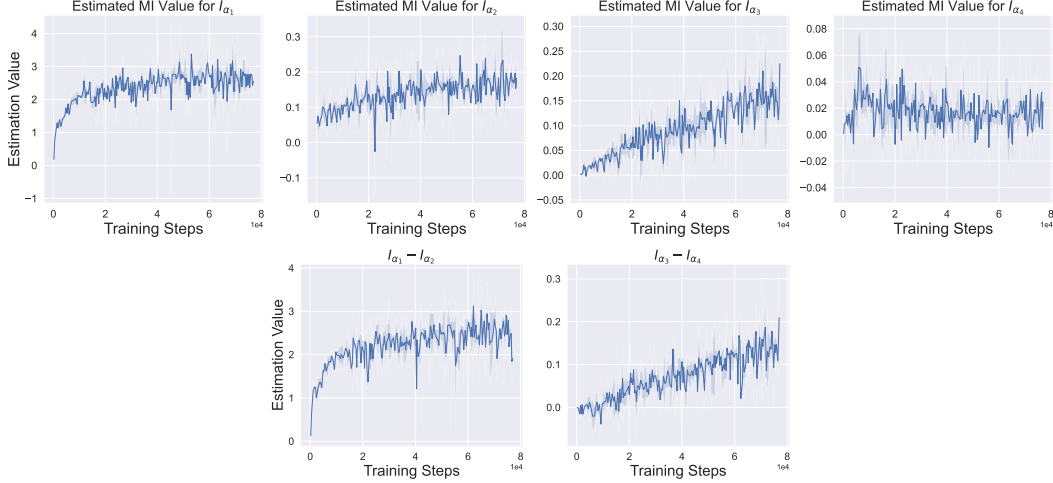


Figure 3: Estimation of the value of four mutual information terms and their differences in experiments on synthetic data.

395 S4.2 Modified Cartpole

396 In the modified Cartpole environment, we configure the values as follows: $\beta_1 = \beta_2 = \beta_3 = \beta_4 = 0.1$
 397 and $\lambda_1 = \lambda_2 = 0.1$. Recurrent State Space Model (RSSM) uses a deterministic part and a stochastic
 398 part to represent latent variables. The deterministic state size for four dynamics are set to be (15, 15,
 399 15, 15), and the stochastic state size are set to be (2, 2, 1, 4). The architecture of the encoder and
 400 decoder for observation is shown in Table 3 and Table 4 (64×64 resolution). Reward model uses
 401 3-layer MLPs with hidden size to be 100 and four mutual information neural estimators are 4-layer
 402 MLPs with hidden size to be 128.

403 S4.3 Variant of Robodesk

404 In the variant of Robodesk, we conduct experiments with the following hyperparameter settings:
 405 $\beta_1 = \beta_2 = 2, \beta_3 = \beta_4 = 0.25$, and $\lambda_1 = \lambda_2 = 0.1$. For the four dynamics, we set the deterministic
 406 state sizes to (120, 40, 40, 40), and the stochastic state sizes to (30, 10, 10, 10). Denoised MDP
 407 utilizes two latent processes with deterministic state sizes [120, 120] and stochastic state sizes [20,
 408 10]. For the mutual information neural estimators, we employ 4-layer MLPs with a hidden size of
 409 128. To ensure a fair comparison, we align the remaining hyperparameters and network structure
 410 with those in the Denoised MDP. We reproduce the results of the Denoised MDP using their released
 411 code, maintaining consistency with their paper by employing the default hyperparameters. In order to
 412 evaluate the impact of the Mutual Information (MI) constraints, we conduct an ablation study. The
 413 results are shown in Figure 4. The constraints \mathcal{J}_{RS}^t and \mathcal{J}_{AS}^t are observed to stabilize the training
 414 process of IFactor. The results of IFactor are averaged over 5 runs, while the results of Denoised
 415 MDP and IFactor without MI are averaged over three runs.

416 **Policy learning based on the learned representations by IFactor** We retrain policies using the
 417 Soft Actor-Critic algorithm [14] with various combinations of the four learned latent categories
 418 as input. We wrap the original environment with visual output using our representation model to
 419 obtain compact features. In this process, both deterministic states and stochastic states are utilized

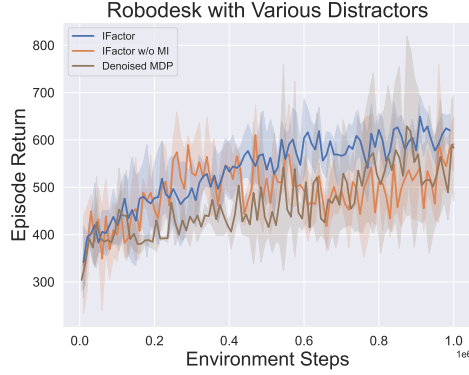


Figure 4: Comparison between IFactor and Denoised MDP in the variant of Robodesk environment.

to form the feature. For instance, when referring to s_t^r , we use both the deterministic states and stochastic states of s_t^r . The implementation of SAC algorithm is based on Stable-Baselines3[15], with a learning rate of 0.0002. Both the policy network and Q network consist of 4-layer MLPs with a hidden size of 256. We use the default hyperparameter settings in Stable-Baselines3 for other parameters.

S4.4 Variants of Deep Mind Control Suite

In the noiseless DMC environments, we set $\beta_1 = \beta_2 = \beta_3 = \beta_4 = 1$. For the DMC environments with video background, we set $\beta_1 = \beta_2 = 1$ and $\beta_3 = \beta_4 = 0.25$. In the DMC environments with video background and noisy sensor, we set $\beta_1 = \beta_2 = 2$ and $\beta_3 = \beta_4 = 0.25$. Lastly, for the DMC environments with video background and jittering camera, we set $\beta_1 = \beta_2 = 1$ and $\beta_3 = \beta_4 = 0.25$. Regarding the Reacher environment with video background and jittering camera, we set $\lambda_1 = \lambda_2 = 0.1$ for our experiments. For the other environments, we set $\lambda_1 = \lambda_2 = 0$. The deterministic state sizes for the four dynamics are set to (120, 120, 60, 60), while the stochastic state sizes are set to (20, 20, 10, 10). The four mutual information neural estimators utilize a 4-layer MLPs with a hidden size of 128. We align the other hyperparameters and network structure with those used in the Denoised MDP for a fair comparison.

Operator	Input Shape	Kernel Size	Stride	Padding
Input	[3, 96, 96]	—	—	—
Conv. + ReLU	[32, 47, 47]	4	2	0
Conv. + ReLU	[64, 22, 22]	4	2	0
Conv. + ReLU	[128, 10, 10]	4	2	0
Conv. + ReLU	[256, 4, 4]	4	2	0
Conv. + ReLU *	[256, 2, 2]	3	1	0
Reshape + FC	[1024]	—	—	—

Table 3: The encoder architecture designed for observation resolution of (96×96) . Its output is then fed into other networks for posterior inference. The default activation function used in the network is RELU. For observations with a resolution of (64×64) , the last convolutional layer(*) is removed.

S4.4.1 Extra Results.

Figure 1 demonstrates the notable improvement in policy performance in the Reacher environment with video background and jittering camera due to the inclusion of the constraints J_{RS}^t and J_{AS}^t . To

Operator	Input Shape	Kernel Size	Stride	Padding
Input	[input_size]	—	—	—
FC + ReLU + Reshape	[1024, 1, 1]	—	—	—
Conv. Transpose + ReLU *	[128, 3, 3]	3	1	0
Conv. Transpose + ReLU	[128, 9, 9]	5	2	0
Conv. Transpose + ReLU	[64, 21, 21]	5	2	0
Conv. Transpose + ReLU	[32, 46, 46]	6	2	0
Conv. Transpose + ReLU	[3, 96, 96]	6	2	0

Table 4: The decoder architecture designed for (96×96) -resolution observation. For (64×64) -resolution observation, the first transpose convolutional layer(*) is removed.

	Environment Steps	Action Repeat	Train Every	Collection Intervals	Batch Size	Sequence Length	Horizon
Modified Cartpole	200,000	1	5	5	20	30	8
Robodesk	1,000,000	2	1000	100	50	50	15
DMC	1,000,000	2	1000	100	25	50	12

Table 5: Some hyperparameters of our method in the environment of Modified Cartpole, Robodesk and DMC. *Environment Steps* represents the number of interactions between the agent and the environment. *Action Repeat* determines how many times an agent repeats an action in a step. *Train Every* specifies the environment step between adjacent training iterations. *Collection Intervals* defines the number of times the model is trained in each training iteration (including world models, policy networks and value networks). *Batch Size* refers to the number of trajectories in each mini-batch. *Sequence Length* denotes the length of the chunk used in training the world models. *Horizon* determines the length of dreaming when training the policy using the world model. Hyperparameters are aligned with those used in the Denoised MDP for fair comparison.

439 further investigate how they affects the model learning, we record the estimation values of four Mutual
440 Information terms throughout the training process, as depicted in Figure 5. The results indicate that
441 both $I_{\alpha_1} - I_{\alpha_2}$ and $I_{\alpha_3} - I_{\alpha_4}$ are maximized for both IFactor and IFactor without MI. However, IFactor
442 exhibits a significantly higher rate of maximizing $I_{\alpha_3} - I_{\alpha_4}$ compared to IFactor without MI. This
443 increased maximization leads to greater predictability of s_t^a by the action, ultimately contributing to
444 the observed performance gain.

445 S4.5 Visualization for DMC

446 In this experiment, we investigate five types of representations, which can be derived from the
447 combination of four original disentangled representation categories. Specifically, s_t^a is the controllable
448 and reward relevant representation. $s_t^r = (s_t^{ar}, s_t^{\bar{ar}})$ is the reward-relevant representation. $s_t^{\bar{ar}}$ is the
449 controllable but reward-irrelevant representation. $s_t^{\bar{ar}}$ is the uncontrollable and reward-irrelevant
450 representation (noise). $s_t^{\bar{r}} = (s_t^{ar}, s_t^{\bar{ar}})$ is the reward-irrelevant representation. Only representations
451 of s_t^r are used for policy optimization. We retrain 5 extra observation decoders to reconstruct the
452 original image, which can precisely characterize what kind of information each type of representation
453 contains, surpassing the limitations of the original decoder that is used in latent traversal. The
454 visualization results are shown in Figure 6. It can be observed that s_t^{ar} captures the movement of
455 the agent partially but not well enough; s_t^r captures the movement of the agent precisely but $s_t^{\bar{r}}$
456 fails (Reacher and Cheetah) or captures extra information of the background (Walker). This finding
457 suggests that s_t^r contains sufficient information within the original noisy observation for effective
458 control, while effectively excluding other sources of noise.

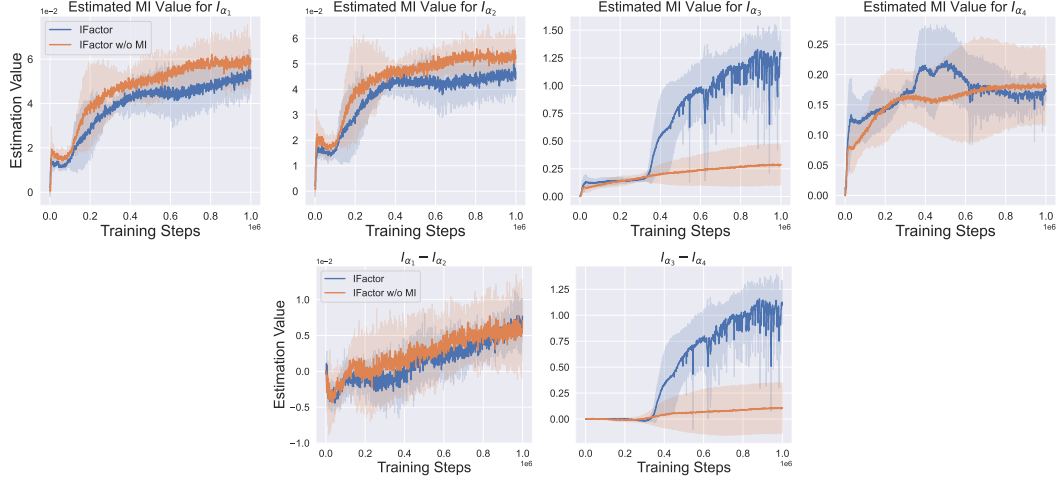


Figure 5: Estimation of the value of four mutual information terms and their differences in the Reacher Easy environment with video background and jittering camera.

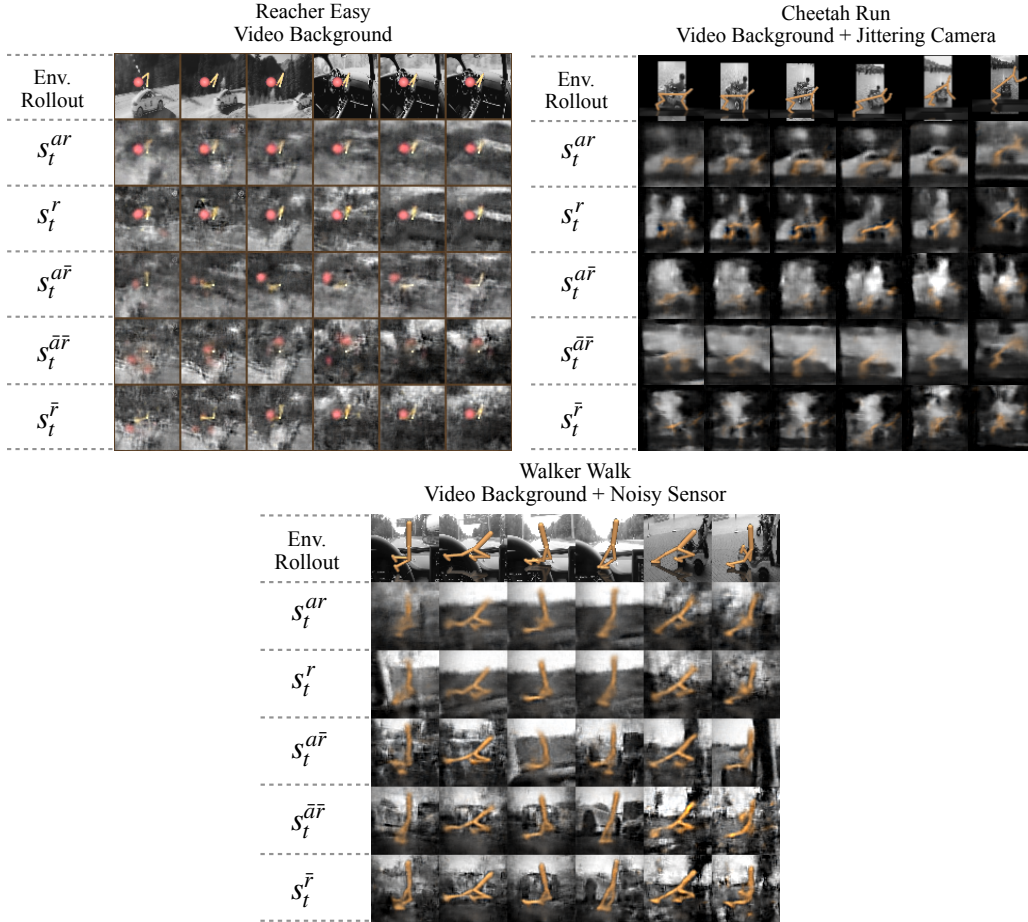


Figure 6: Visualization of the DMC variants and the factorization learned by IFactor.

459 S5 Comparison between IFactor and Denoised MDP

460 While both IFactor and Denoised MDP share the common aspect of factorizing latent variables
 461 based on controllability and reward relevance, it is crucial to recognize the numerous fundamental
 462 distinctions between them.

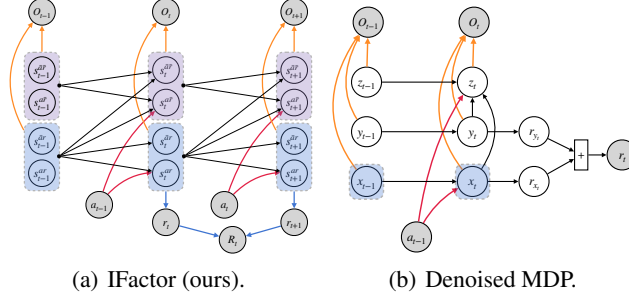


Figure 7: Graphical illustration of our world model and Denoised MDP.

463 First and foremost, Denoised MDP adopts a rather stringent assumption by solely considering three
 464 types of latent variables and assuming independent latent processes for x_t and y_t . However, this strict
 465 assumption may not hold in many scenarios where uncontrollable yet reward-relevant factors exhibit
 466 dependencies on controllable and reward-relevant factors. Take, for instance, the case of car driving:
 467 the agent lacks control over surrounding vehicles, yet their behaviors are indirectly influenced by the
 468 agent’s actions. In contrast, our approach encompasses four types of causally related latent variables
 469 while only assuming conditional independence when conditioning on the state in the previous time
 470 step. This assumption holds true naturally within the MDP framework.

471 Secondly, Denoised MDP is limited to factoring out additive rewards solely for x_t , disregarding the
 472 possibility of non-additive effects in many uncontrollable yet reward-relevant factors. In contrast, our
 473 method embraces the inclusion of non-additive effects of s_t^{ar} on the reward, which is more general.

474 Thirdly, Denoised MDP uses only controllable and reward-relevant latent variables for policy opti-
 475 mization, which we show in the theoretical analysis that it is generally insufficient. In contrast, our
 476 method utilize both controllable and uncontrollable reward-relevant factors for policy training.

477 Finally, Denoised MDP makes the assumption of an intermediate causal effect from x_t to z_t and
 478 from y_t to z_t , which is inherently unidentifiable without further intervention. It is worth noting that
 479 imposing interventions on the latent states is unrealistic in most control tasks, as agents can only
 480 choose actions at specific states and cannot directly intervene on the state itself. In contrast, our
 481 method assumes that there exists no intermediate causal effect for latent variables. In conjunction with
 482 several weak assumptions, we provide a proof of block-wise identifiability for our four categories
 483 of latent variables. This property serves two important purposes: (1) it ensures the removal of
 484 reward-irrelevant factors and the utilization of minimal and sufficient reward-relevant variables for
 485 policy optimization, and (2) it provides a potential means for humans to comprehend the learned
 486 representations within the reinforcement learning (RL) framework. Through latent traversal of the
 487 four types of latent variables, humans can gain insights into the specific kind of information that each
 488 category of representation contains within the image.

489 From the perspective of model structure, it is worth highlighting that the architecture of both the
 490 transition model (prior) and the representation model (posterior) in IFactor differs from that of
 491 Denoised MDP. The structure of prior and posterior of IFactor is shown as follows:

$$\begin{array}{ll}
 \text{Prior:} & \text{Posterior:} \\
 \left\{ \begin{array}{l} p_{\gamma_1}(s_t^{ar} | s_{t-1}^r, a_{t-1}) \\ p_{\gamma_2}(s_t^{\bar{ar}} | s_{t-1}^r) \\ p_{\gamma_3}(s_t^{ar\bar{r}} | \mathbf{s}_{t-1}, a_{t-1}) \\ p_{\gamma_4}(s_t^{\bar{ar}\bar{r}} | \mathbf{s}_{t-1}) \end{array} \right. & \left\{ \begin{array}{l} q_{\phi_1}(s_t^{ar} | o_t, s_{t-1}^r, a_{t-1}) \\ q_{\phi_2}(s_t^{\bar{ar}} | o_t, s_{t-1}^r) \\ q_{\phi_3}(s_t^{ar\bar{r}} | o_t, \mathbf{s}_{t-1}, a_{t-1}) \\ q_{\phi_4}(s_t^{\bar{ar}\bar{r}} | o_t, \mathbf{s}_{t-1}) \end{array} \right. \quad (33)
 \end{array}$$

492 While Denoised MDP has the following prior and posterior:

$$\begin{array}{cc}
\text{Prior:} & \text{Posterior:} \\
\left\{ \begin{array}{l} p_{\gamma_1}(x_t | x_{t-1}, a_{t-1}) \\ p_{\gamma_2}(y_t | y_{t-1}) \\ p_{\gamma_3}(z_t | x_t, y_t, z_{t-1}) \end{array} \right. & \left\{ \begin{array}{l} p_{\phi_1}(x_t | x_{t-1}, y_{t-1}, z_{t-1}, o_t, a_{t-1}) \\ p_{\phi_2}(y_t | x_{t-1}, y_{t-1}, z_{t-1}, o_t, a_{t-1}) \\ p_{\phi_3}(z_t | x_t, y_t, o_t, a_{t-1}) \end{array} \right. \quad (34)
\end{array}$$

A notable distinction can be observed between Denoised MDP and IFactor in terms of the assumptions made for the prior and posterior structures. Denoised MDP assumes independent priors for x_t and y_t , whereas IFactor only incorporates conditional independence, utilizing s_{t-1}^r as input for the transition of both s_t^{ar} and $s_t^{\bar{ar}}$. Moreover, the posterior of y_t receives a_{t-1} as input, potentially implying controllability. Similarly, the posterior of x_t incorporates z_{t-1} as input, which may introduce noise from z_{t-1} into x_t . These implementation details can deviate from the original concept. In contrast, our implementation ensures consistency between the prior and posterior, facilitating a clean disentanglement in our factored model.

From the perspective of the objective function, IFactor incorporates two supplementary mutual information constraints, namely $\mathcal{J}_{\text{RS}}^t$ and $\mathcal{J}_{\text{AS}}^t$, to promote disentanglement and improve policy performance.

References

- [1] P. Spirtes, C. Glymour, and R. Scheines. *Causation, Prediction, and Search*. Springer-Verlag Lectures in Statistics, 1993.
- [2] J. Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge University Press, Cambridge, 2000.
- [3] Biwei Huang, Chaochao Lu, Liu Leqi, José Miguel Hernández-Lobato, Clark Glymour, Bernhard Schölkopf, and Kun Zhang. Action-sufficient state representation learning for control with structural constraints. In *Proceedings of the 39th International Conference on Machine Learning*, volume 162, pages 9260–9279, Baltimore, Maryland, 2022.
- [4] Sébastien Lachapelle, Pau Rodríguez, Yash Sharma, Katie Everett, Rémi Le Priol, Alexandre Lacoste, and Simon Lacoste-Julien. Disentanglement via mechanism sparsity regularization: A new principle for nonlinear ICA. In *Proceedings of the 1st Conference on Causal Learning and Reasoning*, volume 177, pages 428–484, Eureka, CA, 2022.
- [5] Yujia Zheng, Ignavier Ng, and Kun Zhang. On the identifiability of nonlinear ica: Sparsity and beyond. In *Advances in Neural Information Processing Systems 35*, New Orleans, Louisiana, 2022.
- [6] Julius von Kügelgen, Yash Sharma, Luigi Gresele, Wieland Brendel, Bernhard Schölkopf, Michel Besserve, and Francesco Locatello. Self-supervised learning with data augmentations provably isolates content from style. In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems*, pages 16451–16467, Virtual Event, 2021.
- [7] Danijar Hafner, Timothy P. Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning behaviors by latent imagination. In *Proceedings of the 8th International Conference on Learning Representations*, Addis Ababa, Ethiopia, 2020.
- [8] Naftali Tishby, Fernando C. N. Pereira, and William Bialek. The information bottleneck method. *CoRR*, physics/0004057, 2000.
- [9] Alexander A. Alemi, Ian Fischer, Joshua V. Dillon, and Kevin Murphy. Deep variational information bottleneck. In *Proceedings of the 5th International Conference on Learning Representations*, Toulon, France, 2017.
- [10] Mohamed Ishmael Belghazi, Aristide Baratin, Sai Rajeswar, Sherjil Ozair, Yoshua Bengio, R. Devon Hjelm, and Aaron C. Courville. Mutual information neural estimation. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80, pages 530–539, Stockholm, Sweden, 2018.

- 537 [11] Weiran Yao, Guangyi Chen, and Kun Zhang. Temporally disentangled representation learning.
538 In *Advances in Neural Information Processing Systems 35*, New Orleans, Louisiana, 2022.
- 539 [12] Tongzhou Wang, Simon Du, Antonio Torralba, Phillip Isola, Amy Zhang, and Yuandong Tian.
540 Denoised mdps: Learning world models better than the world itself. In *Proceedings of the 39th*
541 *International Conference on Machine Learning*, pages 22591–22612, Baltimore, Maryland,
542 2022.
- 543 [13] Amy Zhang, Rowan Thomas McAllister, Roberto Calandra, Yarin Gal, and Sergey Levine.
544 Learning invariant representations for reinforcement learning without reconstruction. In *Pro-*
545 *ceedings of the 9th International Conference on Learning Representations*, Virtual Event,
546 Austria, 2021.
- 547 [14] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy
548 maximum entropy deep reinforcement learning with a stochastic actor. In *Proceedings of the*
549 *35th International conference on machine learning*, pages 1861–1870, Stockholm, Sweden,
550 2018.
- 551 [15] Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah
552 Dormann. Stable-baselines3: Reliable reinforcement learning implementations. *Journal of*
553 *Machine Learning Research*, 22:268:1–268:8, 2021.