
Temporal Predictive Coding For Model-Based Planning In Latent Space

Tung Nguyen^{*1} Rui Shu^{*2} Tuan Pham^{*1} Hung Bui¹ Stefano Ermon²

Abstract

High-dimensional observations are a major challenge in the application of model-based reinforcement learning (MBRL) to real-world environments. To handle high-dimensional sensory inputs, existing approaches use representation learning to map high-dimensional observations into a lower-dimensional latent space that is more amenable to dynamics estimation and planning. In this work, we present an information-theoretic approach that employs *temporal* predictive coding to encode elements in the environment that can be predicted across time. Since this approach focuses on encoding temporally-predictable information, we implicitly prioritize the encoding of task-relevant components over nuisance information within the environment that are provably task-irrelevant. By learning this representation in conjunction with a recurrent state space model, we can then perform planning in latent space. We evaluate our model on a challenging modification of standard DMControl tasks where the background is replaced with natural videos that contain complex but irrelevant information to the planning task. Our experiments show that our model is superior to existing methods in the challenging complex-background setting while remaining competitive with current state-of-the-art models in the standard setting.

1. Introduction

Learning to control from high dimensional observations has been made possible due to the advancements in reinforcement learning (RL) and deep learning. These advancements have enabled notable successes such as solving video games (Mnih et al., 2015; Lample & Chaplot, 2017) and continuous control problems (Lillicrap et al., 2016) from pixels. How-

^{*}Equal contribution ¹VinAI Research ²Stanford University. Correspondence to: Tung Nguyen <v.tungnd13@vinai.io>, Rui Shu <ruishu@stanford.edu>.

ever, it is well known that performing RL directly in the high-dimensional observation space is sample-inefficient and may require a large amount of training data (Lake et al., 2017). This is a critical problem, especially for real-world applications. Recent model-based RL works (Kaiser et al., 2020; Ha & Schmidhuber, 2018; Hafner et al., 2019; Zhang et al., 2019; Hafner et al., 2020) proposed to tackle this problem by learning a world model in the latent space, and then applying RL algorithms in the latent world model.

The existing MBRL methods that learn a latent world model typically do so via reconstruction-based objectives, which are likely to encode task-irrelevant information, such as of the background. In this work, we take inspiration from the success of contrastive learning and propose a model that employs temporal predictive coding for planning from pixels. The use of temporal predictive coding circumvents the need for reconstruction-based objectives and prioritizes the encoding of temporally-predictable components of the environment, thus making our model robust to environments dominated by nuisance information. Our primary contributions are as follows:

- We propose a temporal predictive coding approach for planning from high-dimensional observations and theoretically analyze its ability to prioritize the encoding of task-relevant information.
- We show experimentally that temporal predicting coding allows our model to outperform existing state-of-the-art models when dealing with complex environments dominated by task-irrelevant information, while remaining competitive on standard DeepMind control (DMC) tasks. Additionally, we conduct detailed ablation analyses to characterize the behavior of our model.

2. Motivation

The motivation and design of our model are largely based on two previous works (Shu et al., 2020; Hafner et al., 2020). In this section, we briefly go over the relevant concepts in each work as well as how they motivate our paper. Shu et al. (2020) proposed PC3, an information-theoretic approach that uses contrastive predictive coding (CPC) to learn a latent space amenable to locally-linear control. Their framework uses a CPC objective between the latent states of two

consecutive time steps. They then use the latent dynamics F as the variational device to define the following lower bound $\ell_{\text{cpc}}(E, F)$:

$$\mathbb{E} \frac{1}{K} \sum_i \ln \frac{F\left(E\left(o_{t+1}^{(i)}\right) \mid E\left(o_t^{(i)}\right), a_t^{(i)}\right)}{\frac{1}{K} \sum_j F\left(E\left(o_{t+1}^{(j)}\right) \mid E\left(o_t^{(j)}\right), a_t^{(j)}\right)}, \quad (1)$$

in which E is the encoder, which encodes a high-dimensional observation into a latent state. We make particular note of their choice to define the CPC objective for latent states *across* time steps instead of between the frame and its corresponding state—as is popular in the existing literature on contrastive learning for RL (Hafner et al., 2020; Sriniwas et al., 2020; Ding et al., 2020). We shall refer to these respective objectives henceforth as *temporal* predictive coding (TPC) and *static* predictive coding (SPC) respectively. Shu et al. (2020) motivated temporal predictive coding via a theory of predictive suboptimality, which shows that the representation learned by temporal predictive coding is no worse at future-observation prediction than the representation learned by a latent variable model when applied to the same task.

In this work, we take a stronger position and argue that temporal predictive coding learns an information-theoretically *superior* representation to conventional latent variable models and static predictive coding models. This is because conventional latent variables models (i.e. temporal variational autoencoders equipped with a Gaussian decoder) and static predictive coding models both seek to encode as much information as possible about the original high-dimensional observations. In contrast, temporal predictive coding *only* encourages the encoding of temporally-predictable information within the environment, which we shall prove is sufficient for optimal decision making (Section 3.2). This property of temporal predictive coding thus makes this objective particularly suitable for handling RL problems dominated by nuisance information, which has recently become an active area of research (Ding et al., 2020; Zhang et al., 2020; Ma et al., 2020).

However, since PC3 only tackles the problem from an optimal control perspective, it is not readily applicable to RL problems. Indeed, PC3 requires a depiction of the goal to perform control, and also the ability to teleport to random locations of the state space to collect data, which are impractical in many problems. On the other hand, Dreamer (Hafner et al., 2020) achieves state-of-the-art performance on many RL tasks, but learns the latent space using a reconstruction-based objective. And while the authors reported a contrastive approach that yielded inferior performance to their reconstruction-based approach, their contrastive approach employed a static predictive coding objective.

The goal of our paper is thus two-fold. First, we propose a

model that leverages the concepts in PC3 and apply them to the Dreamer paradigm and RL setting. In particular, we show how to incorporate temporal predictive coding into the training of the recurrent state space model used in Dreamer. Second, we will demonstrate both theoretically and empirically that our resulting model is robust to RL environments that are dominated by nuisance information.

3. Temporal Predictive Coding for Planning

To plan in an unknown environment, we need to model the environment dynamics from experience. Following the Dreamer paradigm, we do so by iteratively collecting new data and using those data to train the world model. In this section, we focus on presenting the proposed model (its components and objective functions), a theoretical analysis of temporal predictive coding, followed lastly by some practical considerations when implementing the method.

3.1. Model Definition

Since we aim to learn a latent dynamics model for planning, we shall define an encoder E to embed high-dimensional observations into a latent space, a latent dynamics F to model the world in this space, and a reward function,

$$\begin{aligned} \text{Encoder:} & \quad E(o_t) = s_t \\ \text{Latent dynamics:} & \quad F(s_t \mid s_{<t}, a_{<t}) = p(s_t \mid s_{<t}, a_{<t}) \\ \text{Reward function:} & \quad R(r_t \mid s_t) = p(r_t \mid s_t) \end{aligned} \quad (2)$$

in which t is the discrete time step, $\{o_t, a_t, r_t\}_{t=1}^T$ are data sequences with image observations o_t , continuous action vectors a_t , scalar rewards r_t , and s_t denotes the latent state at time t . We model the transition dynamics using a recurrent neural network with a deterministic state $h_t = \text{RNN}(h_{t-1}, s_{t-1}, a_{t-1})$, which summarizes information about the past, followed by the stochastic state model $p(s_t \mid s_{<t}, a_{<t}) = p(s_t \mid h_t)$. Following Hafner et al. (2019), we refer to this RNN as the recurrent state space model (RSSM). In practice, we use a deterministic encoder, and Gaussian distribution for dynamics and reward functions. The graphical model is presented in Figure 1. We now introduce the three key components of our model: *temporal predictive coding*, *consistency*, and *reward prediction*.

Temporal predictive coding with an RSSM Instead of performing pixel prediction to learn E and F as in Hafner et al. (2019; 2020), we maximize the mutual information between the past latent codes and actions against the future latent code $I(E(O_t); E(O_{<t}), A_{<t})$. This objective prioritizes the encoding of predictable components from the environment, which potentially helps avoid encoding nuisance information when dealing with complex image observations. Unlike in Shu et al. (2020), which defined

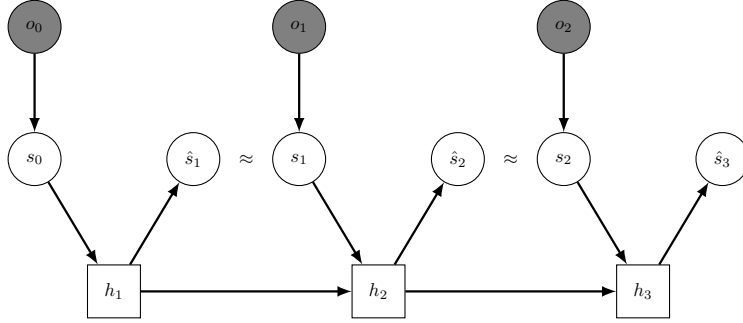


Figure 1. The graphical model of TPC, in which we employ a recurrent neural network to model the dynamics. We omit the action and reward for simplicity. The approximation of \hat{s}_t with s_t is done by using contrastive predictive coding and maximum likelihood.

the TPC objective in a Markovian setting, the use of an RSSM means that we are maximizing the mutual information between the latent code at any time step t and the entire historical trajectory latent codes and actions prior to t . To handle a trajectory over T steps, we define an objective that sums over every possible choice of $t \in \{2, \dots, T\}$,

$$\sum_{t=2}^T I(E(O_t); E(O_{<t}), A_{<t}). \quad (3)$$

To estimate this quantity, we employ contrastive predictive coding (CPC) proposed by Oord et al. (2018). We perform CPC by introducing a critic function f to construct the lower bound at a particular time step t ,

$$I(E(O_t); E(O_{<t}), A_{<t}) \geq \mathbb{E} \frac{1}{K} \sum_i \ln \frac{\exp f(E(o_t^{(i)}), E(o_{<t}^{(i)}), a_{<t}^{(i)})}{\frac{1}{K} \sum_j \exp f(E(o_t^{(i)}), E(o_{<t}^{(j)}), a_{<t}^{(j)})} =: \ell_{\text{tpc}}^{(t)}, \quad (4)$$

where the expectation is over K i.i.d. samples of $(o_t, o_{<t}, a_{<t})$. Note that Eq. (4) uses past $(E(o_{<t}^{(j)}), a_{<t}^{(j)})$ from unrelated trajectories as an efficient source of negative samples for the contrastive prediction of the future latent code $E(o_t^{(i)})$. Following Shu et al. (2020), we choose to tie f to our recurrent latent dynamics model F ,

$$\exp f(s_t, s_{<t}, a_{<t}) = F(s_t | s_{<t}, a_{<t}). \quad (5)$$

There are two favorable properties of this particular design. First, it is parameter-efficient since we can circumvent the instantiation of a separate critic f . Moreover, it takes advantage of the fact that an optimal critic is the true latent dynamics induced by the encoder E (Poole et al., 2019; Ma & Collins, 2018). We denote our objective as $\ell_{\text{tpc}}(E, F) = \sum_{t=2}^T \ell_{\text{tpc}}^{(t)}(E, F)$.

Enforcing latent dynamics consistency Although the true dynamics is an optimal critic for the CPC bound, maximizing this objective only does not ensure the learning of

a latent dynamics model F that is consistent with the true latent dynamics, due to the non-uniqueness of the optimal critic (Shu et al., 2020). Since an accurate dynamics is crucial for planning in the latent space, we additionally introduce a consistency objective, which encourages the latent dynamics model to maintain a good prediction of the future latent code given the past latent codes and actions. Similar to the recurrent CPC objective, we optimize for consistency at every time step in the trajectory,

$$\ell_{\text{cons}}(E, F) = \sum_{t=2}^T \mathbb{E}_{p(o_t, o_{<t}, a_{<t})} \ln F(E(o_t) | E(o_{<t}), a_{<t}). \quad (6)$$

Reward prediction Finally, we train the reward function by maximizing the likelihood of the true reward value conditioned on the latent state $s_t = E(o_t)$,

$$\ell_{\text{reward}}(E, R) = \sum_{t=1}^T \mathbb{E}_{p(o_t)} \ln R(r_t | E(o_t)). \quad (7)$$

3.2. Theoretical Analysis

In contrast to a reconstruction-based objective, which explicitly encourages the encoder to behave injectively on the space of observations, our choice of mutual information objective as specified in Eq. (4) may discard information from the observed scene. In this section, we wish to formally characterize the information discarded by the temporal predictive coding (TPC) objective and argue that any information discarded by an optimal encoder under the TPC objective is provably task-irrelevant.

Lemma 1. Consider an optimal encoder and reward predictor pair (E^*, R^*) where

$$\arg \max_E I(E(O_t); E(O_{<t}), A_{<t}) = E^* \quad (8)$$

$$D_{\text{KL}}(p(r_t | o_t) \| R^*(r_t | E^*(o_t))) = 0. \quad (9)$$

Let $\pi(a_t | E^*(o_{\leq t}), a_{<t})$ denote an E^* -restricted policy whose access to the observations $o_{<t}$ is restricted by E^* .

Let $\pi_{\text{aux}}(a_t \mid E^*(o_{\leq t}), E'(o_{\leq t}), a_{<t})$ denote an (E^*, E') -restricted policy which has access to auxiliary information about $o_{<t}$ via some encoder E' . Let $\eta(\pi)$ denote the expected cumulative reward achieved by a policy π over a finite horizon T . Then there exists no encoder E' where the optimal E^* -restricted policy underperforms the optimal (E^*, E') -restricted policy,

$$\nexists E' \text{ s.t. } \eta(\pi^*) < \eta(\pi_{\text{aux}}^*). \quad (10)$$

Intuitively, since E^* optimizes TPC, any excess information contained in E' about o_t (not already accounted for by E^*) must be *temporally-unpredictive*—it is neither predictable from the past nor predictive of the future. The excess information conveyed by E' is effectively nuisance information and thus cannot be exploited to improve the agent’s performance. It is therefore permissible to dismiss the excess information in E' as being task-irrelevant. We provide a proof formalizing this intuition in Appendix B.

It is worth noting that TPC does not actively penalize the encoding of nuisance information, nor does temporally-predictive information necessarily mean it will be task-relevant. However, if the representation space has limited capacity to encode information about o_t (e.g., due to dimensionality reduction, a stochastic encoder, or an explicit information bottleneck regularizer), TPC will favor temporally-predictive information over temporally-unpredictive information—and in this sense, thus favor potentially task-relevant information over provably task-irrelevant information. This is in sharp contrast to a reconstruction objective, which makes no distinction between these two categories of information contained in the observation o_t .

3.3. Practical Implementation Details

Avoiding map collapse Naively optimizing ℓ_{tpc} and ℓ_{cons} can lead to a trivial solution, where the latent map collapses into a single point to increase the consistency objective arbitrarily. In the previous work, Shu et al. (2020) resolved this by adding Gaussian noise to the future encoding, which balances the latent space retraction encouraged by ℓ_{cons} with the latent space expansion encouraged by ℓ_{tpc} . However, we found this trick insufficient and thus further incorporate a variant of static predictive coding. Our static predictive coding defines a lower bound for $I(E(O_t); O_t)$ where we employ a Gaussian distribution $\mathcal{N}(E(O_t) \mid E(O_t), \sigma I)$ as our choice of variational critic for $f(E(O_t), O_t)$. Crucially, our variant keeps the variance σ fixed so that the static predictive coding objective solely serves the role of encouraging map expansion.

Smoothing the dynamics model Since our encoder is deterministic, the dynamics always receives clean latent codes

as inputs during training. However, in behavior learning, we roll out multiple steps towards the future from a stochastic dynamics. These roll outs are susceptible to a cascading error problem, which hurts the value estimation and policy learning. To resolve this issue, we smooth the dynamics by adding Gaussian noise to the inputs of the recurrent dynamics during training. The noise-adding procedure is as follows: assume the dynamics outputs $\hat{s}_t \sim \mathcal{N}(\mu_t, \sigma_t^2)$ as the prediction at time step t , we then add $\epsilon_t \sim \mathcal{N}(0, \sigma_t^2)$ to $s_t = E(o_t)$ and feed it to the latent dynamics, and repeat for every time step $1 \leq t \leq T$. We call this dynamics-associated noise, which ensures that the latent dynamics can handle the amount of noise that it produces when rolling out. The overall objective of our model is

$$\begin{aligned} \max_{E, F, R} \quad & \lambda_1 \ell_{\text{tpc}}(E, F) + \lambda_2 \ell_{\text{cons}}(E, F) \\ & + \lambda_3 \ell_{\text{spc}}(E) + \lambda_4 \ell_{\text{reward}}(E, R). \end{aligned} \quad (11)$$

4. Behavior Learning

Following Hafner et al. (2020), we use latent imagination to learn a parameterized policy for control. For self-containedness, this section gives a summary of this approach. Given the latent state $s_t = E(o_t)$, we roll out multiple steps into the future using the learned dynamics F , estimate the return and perform backpropagation through the dynamics to maximize this return, which in turn improves the policy.

Action and value models Two components needed for behavior learning are the action model and the value model, which both operate on the latent space. The value model estimates the expected imagined return when following the action model from a particular state s_τ , and the action model implements a policy, conditioned on s_τ , that aims to maximize this return. With imagine horizon H , we have

$$\begin{aligned} \text{Action model: } \quad & a_\tau \sim \pi(a_\tau \mid s_\tau) \\ \text{Value model: } \quad & v(s_\tau) \approx \mathbb{E}_{\pi(\cdot \mid s_\tau)} \sum_{\tau=t}^{t+H} \gamma^{\tau-t} r_\tau \end{aligned} \quad (12)$$

Value estimation To learn the action and value model, we need to estimate the state values of imagined trajectories $\{s_\tau, a_\tau, r_\tau\}_{\tau=t}^{t+H}$, where s_τ and a_τ are sampled according to the dynamics and the policy. In this work, we use value estimation presented in Sutton & Barto (2018),

$$\begin{aligned} V_N^k(s_\tau) & \doteq \mathbb{E}_{q_\theta, q_\phi} \left(\sum_{n=\tau}^{h-1} \gamma^{n-\tau} r_n + \gamma^{h-\tau} v_\psi(s_h) \right) \\ V_\lambda(s_\tau) & \doteq (1-\lambda) \sum_{n=1}^{H-1} \lambda^{n-1} V_N^n(s_\tau) + \lambda^{H-1} V_N^H(s_\tau) \\ & \text{with } h = \min(\tau + k, t + H). \end{aligned} \quad (13)$$

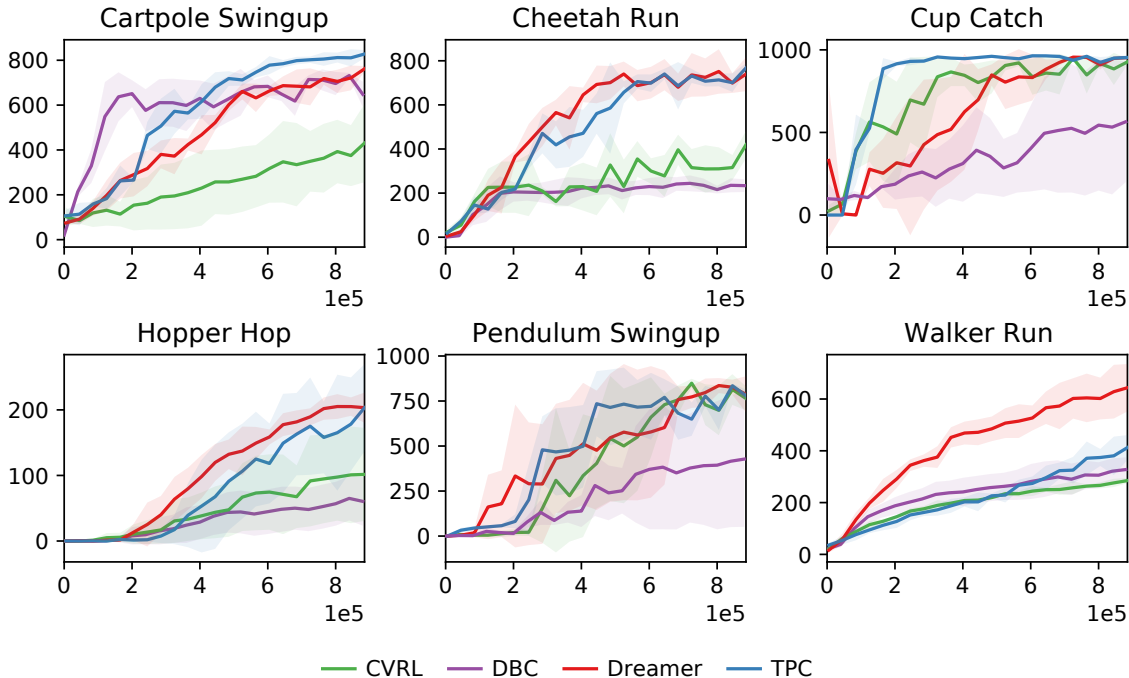


Figure 2. TPC vs the baselines in standard DMC tasks. The return is computed in 1000 environment steps. Each task is run with 3 seeds.

which allows us to estimate the return beyond the imagine horizon. The value model is then optimized to regress this estimation, while the action model is trained to maximize the value estimation of all states s_τ along the imagined trajectories,

$$\begin{aligned} \max_{\pi} \mathbb{E} \left(\sum_{\tau=t}^{t+H} V_{\lambda}(s_{\tau}) \right) \\ \min_v \mathbb{E} \left(\sum_{\tau=t}^{t+H} \frac{1}{2} \|v_{\psi}(s_{\tau}) - V_{\lambda}(s_{\tau})\|^2 \right). \end{aligned} \quad (14)$$

5. Experiments

In this section, we empirically evaluate the proposed model (which we shall simply refer to as TPC for convenience) in various settings. First, we design experiments to compare the relative performance of our model with the current best model-based method in several standard control tasks. Second, we evaluate its ability to handle a more realistic but also more complicated scenario, in which we replace the background of the environment with a natural video. We discuss how our model is superior in the latter case compared to other existing methods, while remaining competitive in the standard setting. Finally, we conduct ablation studies to demonstrate the importance of the components of our model.

Control tasks For the standard setting, we test our model on 6 DeepMind Control (DMC) tasks (Tassa et al., 2018): Cartpole Swingup, Cheetah Run, Walker Run, Pendulum Swingup, Hopper Hop and Cup Catch. In the natural background setting, we replace the background of each data trajectory with a video taken from the kinetics dataset (Kay et al., 2017). We split the original dataset into two separate sets for training and evaluation to also test the generalization of each method.

Baseline methods We compare TPC with Dreamer (Hafner et al., 2020), CVRL (Ma et al., 2020), and DBC (Zhang et al., 2020). We compare against Dreamer since it is the current state of the art model-based method for planning from pixels and directly inspired our design of TPC. CVRL is a closely-related contrastive-based model recently developed that also touted its ability to handle nuisance backgrounds. Finally, we compare with DBC, a leading model-free method for learning representations that are invariant to distractors, to showcase how temporal predictive coding fares against bisimulation for learning task-relevant representations for downstream control. Each method is evaluated by the environment return in 1000 steps. For the baselines, we use the best set of hyperparameters as reported in their paper. We run each task with 3 different seeds for each model.

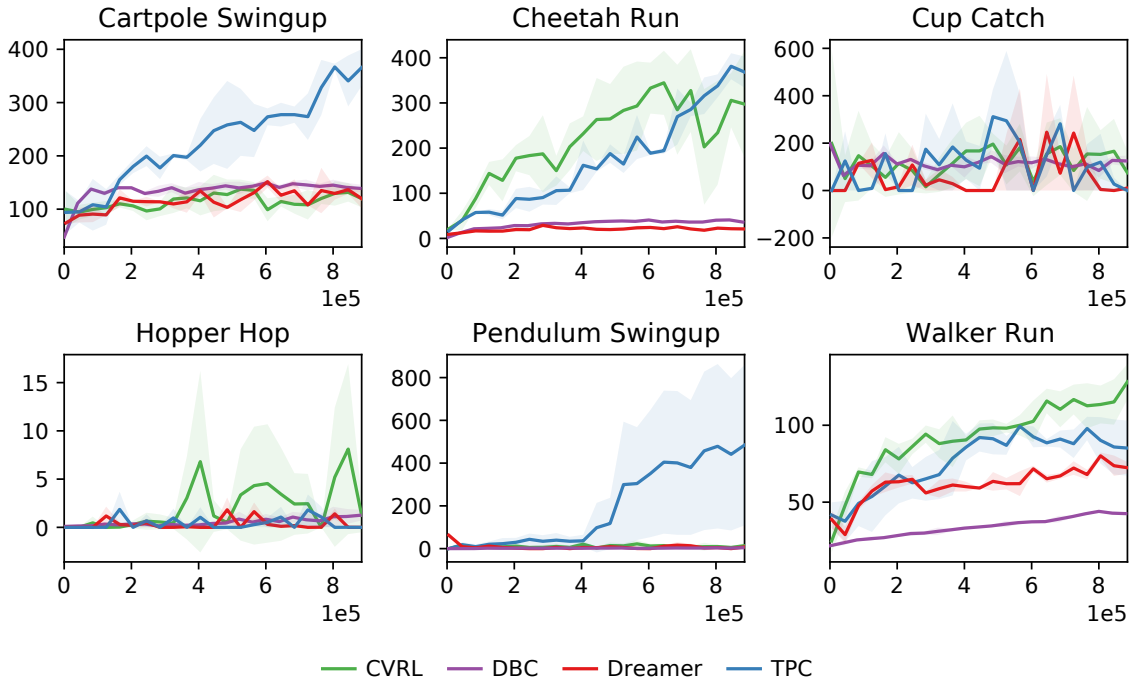


Figure 3. TPC vs the baselines in the natural backgrounds setting. The return is computed in 1000 environment steps. Each task is run with 3 seeds.

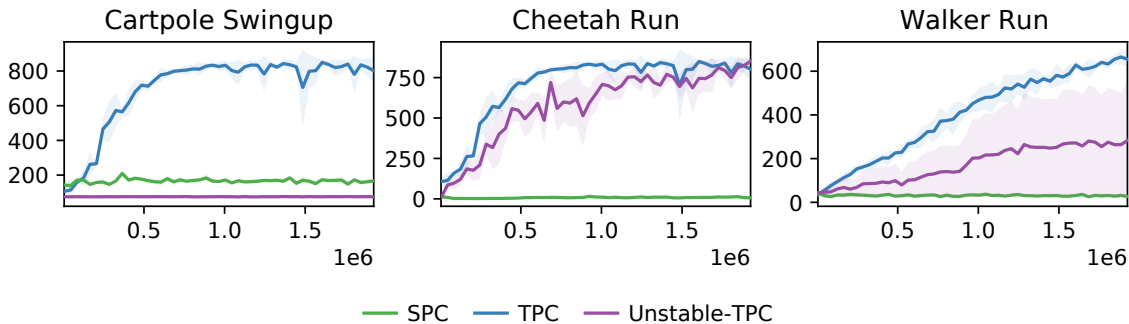


Figure 4. TPC vs two different variants, where we omit the temporal predictive coding in SPC and static predictive coding in Unstable-TPC. The return is computed in 1000 environment steps. Each task is run with 3 seeds.

5.1. Comparisons in Standard Setting

We demonstrate the performance of all methods on standard control tasks in Figure 2. TPC is competitive with Dreamer in all but one task, while CVRL and DBC both underperformed Dreamer on four tasks. TPC’s strong performance is in contrast to what was previously observed in Dreamer (Hafner et al., 2020), where they showed the inferiority of a *static* contrastive learning approach (which applies CPC between an image observation and its corresponding latent code) compared to their reconstruction-based approach. Our results thus show that the use of *temporal* contrastive predictive coding is critical for learning a latent space that is suitable for behavior learning.

5.2. Comparisons in Natural Background Setting

In this setting, we evaluate the robustness of TPC versus the baselines in dealing with complicated, natural backgrounds. The performance of all models in this setting is shown in Figure 3. Dreamer fails to achieve meaningful performance across all six tasks. This is due to the use of reconstruction loss, which forces the world model to reconstruct every single pixel in the observations, thus leading to the encoding of task-irrelevant information.

TPC outperforms all the baselines significantly on two of the six tasks, while being competitive to CVRL on Walker Run and Cheetah Run (see Appendix C for more environment steps, where TPC outperforms CVRL on three tasks).

All methods fail to work on Hopper Hop and Cup Catch¹ after 1×10^6 environment steps. Hopper Hop is a very challenging task even in the standard setting. Furthermore, the agents in these two tasks are tiny compared to the complex backgrounds, thus making it difficult for the model to distill task-relevant information. We note that the discrepancy between our results for CVRL and results reported in their paper (Ma et al., 2020) is due to the fact that while they used Imagenet (Russakovsky et al., 2015) for backgrounds in the original experiments, we use videos from kinetics dataset (Kay et al., 2017) instead.

DBC fails to work on all six tasks, which is in contrast to their reported performance in (Zhang et al., 2020). We note that, in their experiments, the authors used a single video for both training and testing. Their poor performance in our setting suggests that DBC is incapable of dealing with multiple distracting scenes, and also does not generalize well to unseen distractors.

Overall, the empirical results show that our proposed model is superior to all the baselines. In practice, it is desirable to have model which can work well for both standard (easy) setting and natural background (complex, distracting) setting. TPC satisfies this, since it performs on-par with Dreamer and significantly better than CVRL and DBC on standard control tasks, and in natural background setting which Dreamer fails, TPC works reasonably well and outperforms their counterparts CVRL and DBC.

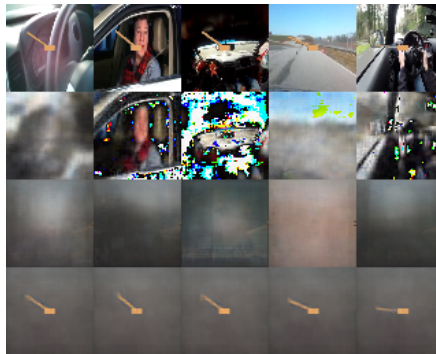
5.3. Ablation Analysis

We conduct an ablation analysis to evaluate the importance of each CPC objective employed in TPC. To do so, we compare the original model with two variants, where we omit the temporal predictive coding or static predictive coding objective respectively (added to prevent map collapse; see Section 3.3). The relative performance of these models on three control tasks is shown in Figure 4. The original TPC achieves the best performance, while the variant with only static predictive coding (denoted SPC) fails across all three tasks. TPC without static predictive coding (Unstable-TPC) is unstable, since it faces the map collapsing problem, which leads to poor performance.

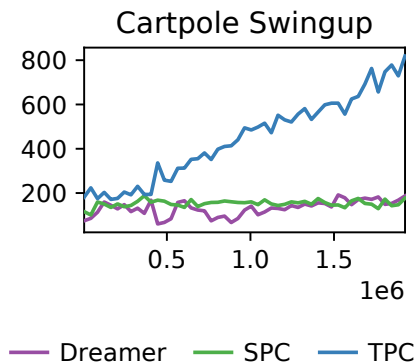
5.4. Analysis in the Random Background Setting

To demonstrate that temporal predictive coding is able to filter out temporally-unpredictive information, we conduct an experiment where the background is randomly chosen for each time step. In this setting, the only source of temporally-predictive information is the agent itself. Using the representations learned by Dreamer, TPC, and SPC, we then

¹TPC and CVRL work on Cup Catch after 2×10^6 steps. The results are shown in Appendix C.



(a) Reconstructions of a trajectory. 1st row: ground-truth data, 2nd row: Dreamer, 3rd row: SPC, 4th row: TPC



(b) Performance of TPC versus Dreamer and SPC in the random background setting.

Figure 5. TPC and the baselines in the random background setting

trained auxiliary decoders that try to reconstruct the original observations. Figure 5 shows that the representation learned by TPC is unable to reconstruct the background but can reconstruct the agent with high fidelity. In contrast, Dreamer’s representation can reconstruct most of the image but not the agent itself, and SPC’s reconstructions are largely uninterpretable. Furthermore, of the three models, only TPC successfully performs in the random background setting.

While the random background setting experiment aligns with our theoretical analysis and clearly demonstrates TPC’s robustness to temporally-unpredictive nuisance information, it also exposes a particular limitation in our theoretical understanding of TPC’s performance in the natural background setting, since natural videos are in fact a source of *temporally-predictive* nuisance information. We provide two possible explanations for TPC’s empirically strong performance in the natural background setting. First, the natural videos may contain considerable temporally-unpredictive information that TPC successfully dismisses. Second, the natural videos largely consist of hard-to-predict informa-

tion that our recurrent neural network deems effectively unpredictable. We note that the latter potential explanation appears intimately related to the concept of “usable information” (Xu et al., 2020) whenever we instantiate a mutual information estimator with finite-computational power—and is worthy of further research scrutiny.

6. Related Work

Learning latent space for model-based RL via reconstruction Latent world models can be learned by jointly training the encoder and the latent dynamics model with observation reconstruction loss. Learning Controllable Embedding (LCE) approach, including E2C (Watter et al., 2015), RCE (Banijamali et al., 2018) and PCC (Levine et al., 2020), uses randomly collected data to pre-train a Markovian latent dynamic model that is specifically designed for locally-linear control, then run offline optimal control on top of the learned latent space. CARL (Cui et al., 2020) extends these works for Soft Actor-Critic (Haarnoja et al., 2018) and also proposes an online version, in which they iteratively learn the model and a parameterized policy. World Models (Ha & Schmidhuber, 2018) learn a recurrent latent dynamic model in a two-stage process to evolve their linear controllers in imagination. PlaNet (Hafner et al., 2019) jointly learns a recurrent state space model (RSSM) and plans in latent space using the cross entropy method, while Dreamer (Hafner et al., 2020) uses RSSM to iteratively learn the model and the policy by backpropagating through the dynamics. SOLAR (Zhang et al., 2019) models the dynamics as time-varying linear-Gaussian with quadratic costs and controls using guided policy search. However, training world models with reconstruction loss has several drawbacks: it requires a decoder as an auxiliary network for predicting images, and by reconstructing every single pixel, those methods are potentially vulnerable to task-irrelevant information such as an irrelevant background.

Learning latent space for model-based RL via contrastive learning An alternative framework for learning latent world models is contrastive learning. Contrastive learning is a self-supervised learning technique that aims to learn representations by contrasting positive samples against negative samples without having to reconstruct images (Oord et al., 2018; Chen et al., 2020a). Recently proposed contrastive learning methods have achieved significant successes in learning representations purely from unlabeled data, which include works by (Chen et al., 2020a;b; Bachman et al., 2019; Hénaff et al., 2019; He et al., 2020; Tian et al., 2019). Poole et al. (2019) has also established a close connection between contrastive learning and mutual information maximization. In the context of RL, recent works have proposed to use this framework to accelerate RL from pixels in two distinct directions: 1) cast contrastive

learning as an auxiliary representation learning task, and use model-free RL methods on top of the learned latent space (Oord et al., 2018; Srinivas et al., 2020); and 2) use contrastive learning in conjunction with learning a latent dynamics for planning in the latent space (Shu et al., 2020; Ding et al., 2020; Hafner et al., 2020).

Learning to control from pixels with distractors To work well in real-world scenarios, a control-from-pixel method needs to be robust against distractors in the observation. Recent works have been proposed to address this problem, including CVRL (Ma et al., 2020) and DBC (Zhang et al., 2020). Similar to TPC, CVRL employs a contrastive-based objective to learn the latent representation and latent dynamics that enable direct planning in latent space. However, they aim to maximize the mutual information between the current observation o_t and its corresponding latent state s_t . This in theory still encourages the model to encode as much information in the observation as possible, including rich but task-irrelevant information such as the background. Our proposed TPC, in contrast, employs a *temporal* predictive coding objective, which helps discourage the encoding of provably task-irrelevant information, as shown in Section 3.2. Taking a different approach, DBC was proposed to learn invariant representations for RL by forcing the latent space to preserve the on-policy bisimulation metric. This is in theory better than our contrastive objective, since it is capable of filtering out task-irrelevant but temporally-predictive information. However, while being theoretically sound, their implemented objective does not adhere to the true bisimulation metric. Specifically, they never solved the fixed point equation to compute the metric, but instead used a heuristic version of it, where they replaced the ground metric d with Euclidean distance while comparing two distributions of the next state, with d being the current bisimulation metric.

7. Conclusion

In this work, we propose a temporal predictive coding approach for planning in latent space. We show theoretically and experimentally that temporal predictive coding prioritizes the encoding of task-relevant components over temporally-unpredictive—and thus provably-irrelevant—information. This is critically different from reconstruction-based objectives as well as static contrastive learning objectives that only maximize the mutual information between the current observation and its latent code, which indiscriminately favor the encoding of both task-relevant and irrelevant information. Our experiments show that temporally predictive coding outperforms state-of-the-art models in environments dominated by task-irrelevant information while remaining competitive on standard DMC tasks.

Acknowledgements

This research was supported by NSF (#1651565, #1522054, #1733686), ONR (N00014-19-1-2145), AFOSR (FA9550-19-1-0024), FLI, and Amazon AWS.

References

- Bachman, P., Hjelm, R. D., and Buchwalter, W. Learning representations by maximizing mutual information across views. In *Advances in Neural Information Processing Systems*, pp. 15535–15545, 2019.
- Banijamali, E., Shu, R., Ghavamzadeh, M., Bui, H., and Ghodsi, A. Robust locally-linear controllable embedding. volume 84 of *Proceedings of Machine Learning Research*, pp. 1751–1759, Playa Blanca, Lanzarote, Canary Islands, 09–11 Apr 2018. PMLR. URL <http://proceedings.mlr.press/v84/banijamali18a.html>.
- Chen, T., Kornblith, S., Norouzi, M., and Hinton, G. A simple framework for contrastive learning of visual representations. *arXiv preprint arXiv:2002.05709*, 2020a.
- Chen, T., Kornblith, S., Swersky, K., Norouzi, M., and Hinton, G. Big self-supervised models are strong semi-supervised learners. *arXiv preprint arXiv:2006.10029*, 2020b.
- Cui, B., Chow, Y., and Ghavamzadeh, M. Control-aware representations for model-based reinforcement learning. *arXiv preprint arXiv:2006.13408*, 2020.
- Ding, Y., Clavera, I., and Abbeel, P. Mutual information maximization for robust plannable representations. *arXiv preprint arXiv:2005.08114*, 2020.
- Ha, D. and Schmidhuber, J. Recurrent world models facilitate policy evolution. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 31*, pp. 2450–2462. Curran Associates, Inc., 2018. URL <http://papers.nips.cc/paper/7512-recurrent-world-models-facilitate-policy-evolution.pdf>.
- Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. volume 80 of *Proceedings of Machine Learning Research*, pp. 1861–1870, Stockholm, Sweden, 10–15 Jul 2018. PMLR. URL <http://proceedings.mlr.press/v80/haarnoja18b.html>.
- Hafner, D., Lillicrap, T., Fischer, I., Villegas, R., Ha, D., Lee, H., and Davidson, J. Learning latent dynamics for planning from pixels. volume 97 of *Proceedings of Machine Learning Research*, pp. 2555–2565, Long Beach, California, USA, 09–15 Jun 2019. PMLR. URL <http://proceedings.mlr.press/v97/hafner19a.html>.
- Hafner, D., Lillicrap, T., Ba, J., and Norouzi, M. Dream to control: Learning behaviors by latent imagination. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=S110TC4tDS>.
- He, K., Fan, H., Wu, Y., Xie, S., and Girshick, R. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9729–9738, 2020.
- Hénaff, O. J., Srinivas, A., De Fauw, J., Razavi, A., Doersch, C., Eslami, S., and Oord, A. v. d. Data-efficient image recognition with contrastive predictive coding. *arXiv preprint arXiv:1905.09272*, 2019.
- Kaiser, L., Babaeizadeh, M., Milos, P., Osinski, B., Campbell, R. H., Czechowski, K., Erhan, D., Finn, C., Koza-kowski, P., Levine, S., Mohiuddin, A., Sepassi, R., Tucker, G., and Michalewski, H. Model based reinforcement learning for atari. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=S1xCPJHtDB>.
- Kay, W., Carreira, J., Simonyan, K., Zhang, B., Hillier, C., Vijayanarasimhan, S., Viola, F., Green, T., Back, T., Natsev, P., et al. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, 2017.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Lake, B. M., Ullman, T. D., Tenenbaum, J. B., and Gershman, S. J. Building machines that learn and think like people. *Behavioral and brain sciences*, 40, 2017.
- Lample, G. and Chaplot, D. S. Playing fps games with deep reinforcement learning. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, AAAI’17*, pp. 2140–2146. AAAI Press, 2017.
- Levine, N., Chow, Y., Shu, R., Li, A., Ghavamzadeh, M., and Bui, H. Prediction, consistency, curvature: Representation learning for locally-linear control. In *International Conference on Learning Representations*, 2020. URL https://openreview.net/forum?id=BJxG_0EtDS.
- Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., and Wierstra, D. Continuous control with deep reinforcement learning. In *ICLR (Poster)*, 2016. URL <http://arxiv.org/abs/1509.02971>.

- Ma, X., Chen, S., Hsu, D., and Lee, W. S. Contrastive variational model-based reinforcement learning for complex observations. *arXiv preprint arXiv:2008.02430*, 2020.
- Ma, Z. and Collins, M. Noise contrastive estimation and negative sampling for conditional models: Consistency and statistical efficiency. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 3698–3707, Brussels, Belgium, October–November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1405. URL <https://www.aclweb.org/anthology/D18-1405>.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G., et al. Human-level control through deep reinforcement learning. *nature*, 518(7540): 529–533, 2015.
- Oord, A. v. d., Li, Y., and Vinyals, O. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- Poole, B., Ozair, S., Van Den Oord, A., Alemi, A., and Tucker, G. On variational bounds of mutual information. volume 97 of *Proceedings of Machine Learning Research*, pp. 5171–5180, Long Beach, California, USA, 09–15 Jun 2019. PMLR. URL <http://proceedings.mlr.press/v97/poole19a.html>.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3): 211–252, 2015.
- Shu, R., Nguyen, T., Chow, Y., Pham, T., Than, K., Ghavamzadeh, M., Ermon, S., and Bui, H. H. Predictive coding for locally-linear control. *arXiv preprint arXiv:2003.01086*, 2020.
- Srinivas, A., Laskin, M., and Abbeel, P. Curl: Contrastive unsupervised representations for reinforcement learning. *arXiv preprint arXiv:2004.04136*, 2020.
- Sutton, R. S. and Barto, A. G. *Reinforcement learning: An introduction*. MIT press, 2018.
- Tassa, Y., Doron, Y., Muldal, A., Erez, T., Li, Y., Casas, D. d. L., Budden, D., Abdolmaleki, A., Merel, J., Lefrancq, A., et al. Deepmind control suite. *arXiv preprint arXiv:1801.00690*, 2018.
- Tian, Y., Krishnan, D., and Isola, P. Contrastive multiview coding. *arXiv preprint arXiv:1906.05849*, 2019.
- Watter, M., Springenberg, J., Boedecker, J., and Riedmiller, M. Embed to control: A locally linear latent dynamics model for control from raw images. In *Advances in neural information processing systems*, pp. 2746–2754, 2015.
- Xu, Y., Zhao, S., Song, J., Stewart, R., and Ermon, S. A theory of usable information under computational constraints. *arXiv preprint arXiv:2002.10689*, 2020.
- Zhang, A., McAllister, R., Calandra, R., Gal, Y., and Levine, S. Learning invariant representations for reinforcement learning without reconstruction. *arXiv preprint arXiv:2006.10742*, 2020.
- Zhang, M., Vikram, S., Smith, L., Abbeel, P., Johnson, M., and Levine, S. SOLAR: Deep structured representations for model-based reinforcement learning. volume 97 of *Proceedings of Machine Learning Research*, pp. 7444–7453, Long Beach, California, USA, 09–15 Jun 2019. PMLR. URL <http://proceedings.mlr.press/v97/zhang19m.html>.

Supplementary Materials to Temporal Predictive Coding For Model-Based Planning In Latent Space

A. Hyper Parameters

A.1. Standard setting

Dreamer, CVRL and TPC share the following hyperparameters:

Model components

- Latent state dimension: 30
- Recurrent state dimension: 200
- Activation function: ELU
- The action model outputs a tanh mean scaled by a factor of 5 and a softplus standard deviation for the Normal distribution that is then transformed using tanh (Haarnoja et al., 2018)

Learning updates

- Batch size: 50 for Dreamer and CVRL, 250 for TPC
- Trajectories length: 50
- Optimizer: Adam (Kingma & Ba, 2014) with learning rates 6×10^{-4} for world model, 8×10^{-5} for value and action.
- Gradient update rate: 100 gradient updates every 1000 environment steps.
- Gradient clipping norm: 100
- Imagination horizon: 15
- $\gamma = 0.99$ and $\lambda = 0.95$ for value estimation

Environment interaction

- The dataset is initialized with $S = 5$ episodes collected using random actions.
- We iterate between 100 training steps and collecting 1 episode by executing the predicted mode action with $\mathcal{N}(0, 0.3)$ exploration noise.
- Action repeat: 2
- Environment steps: 2×10^6

Additionally,

- Dreamer and CVRL clip the KL below 3 nats
- CVRL uses a bi-linear model for the critic function in the contrastive loss: $f_{\theta}(s_t, o_t) = \exp(z_t^T W_{\theta} s_t)$, where z_t is an embedding vector for observation o_t and W_{θ} is a learnable weight matrix parameterized by θ .

- TPC has a fixed set of coefficient in the overall objective for all control tasks: $\lambda_1 = 1, \lambda_2 = 0.1, \lambda_3 = 1, \lambda_4 = 1$. We use a fixed Gaussian noise $\epsilon \sim \mathcal{N}(0, 0.2^2)$ to add to the future latent code when computing temporal CPC, as suggested in (Shu et al., 2020), and also use 0.2 as the fixed variance in static CPC.

In TPC, we also use a target network for the value model and update this network every 100 gradient steps. Note that we also tried to use target value network for Dreamer, but it does not improve the results, as suggested by their original paper (Hafner et al., 2020).

Hyperparameters for DBC We use the same set of hyperparameters as reported in the paper (Zhang et al., 2020)

- Replay buffer capacity: 1000000
- Batch size: 128
- Discount γ : 0.99
- Optimizer: Adam
- Critic learning rate: 10^{-5}
- Critic target update frequency: 2
- Critic Q-function soft-update rate τ_Q : 0.005
- Critic encoder soft-update rate τ_ϕ : 0.005
- Actor learning rate: 10^{-5}
- Actor update frequency: 2
- Actor log stddev bounds: $[-5, 2]$
- Encoder learning rate: 10^{-5}
- Decoder learning rate: 10^{-5}
- Decoder weight decay: 10^{-7}
- Temperature learning rate: 10^{-4}
- Temperature Adam’s β_1 : 0.9
- Init temperature: 0.1

Hyperparameters search for TPC TPC has four hyperparameters that can be tuned: $\lambda_1, \lambda_2, \lambda_3$ and λ_4 , which are coefficients for the TPC objective, consistency objective, SPC objective and reward prediction objective, respectively. Since λ_1, λ_3 and λ_4 do not conflict with each other, we fixed them to 1 and only tuned λ_2 in our experiments. We performed grid search for λ_2 in range $\{0.05, 0.1, 0.2\}$ ² on the *Cartpole Swingup* task and then used the same set of hyperparameters for all the remaining tasks.

A.2. Natural background setting

To further encourage the model to focus on task-relevant information from observations, we additionally tune the weight λ_4 of the reward loss in the training objective for both Dreamer and TPC. In each control task they share the same reward coefficient, which is specified in the table below. CVRL and DBC have the same hyperparameters as in the standard setting.

²Larger values of λ_2 lead to representation collapse.

Table 1. Reward coefficients for different tasks in the natural backgrounds setting

Task	Reward coefficient
Cartpole Swingup, Cup Catch	1000
Cheetah Run, Walker Run, Pendulum Swingup, Hopper Hop	100

B. Proof of Lemma 1

Our goal is to show that, under the conditions in Lemma 1,

$$\eta(\pi^*) \geq \eta(\pi_{\text{aux}}^*) \quad (15)$$

for any choice of auxiliary encoder E' .

We start by denoting $s_t = E^*(o_t)$ and $s'_t = E'(o_t)$. Note that the performance of π_{aux} can be written as

$$\eta(\pi_{\text{aux}}) = \mathbb{E}_{(\pi_{\text{aux}}, p)} r(o_{1:T}) \quad (16)$$

$$= \sum_{\tau_{\text{aux}}} r(o_{1:T}) \prod_t p(o_t | o_{<t}, a_{<t}) p(s_t, s'_t | o_t) \pi_{\text{aux}}(a_t | s_{\leq t}, s'_{\leq t}, a_{<t}), \quad (17)$$

where τ_{aux} denotes the full trajectory of $(o, s, s', a)_{1:T}$ and $r(o_t)$ evaluates the reward at o_t (for simplicity, we shall assume $p(r_t | s_t)$ is deterministic. Since $D_{\text{KL}}(p(r_t | o_t) || R^*(r_t | E^*(o_t))) = 0$, we can rewrite as

$$\eta(\pi_{\text{aux}}) = \sum_{\tau_{\text{aux}}} R^*(s_{1:T}) \prod_t p(o_t | o_{<t}, a_{<t}) p(s_t, s'_t | o_t) \pi_{\text{aux}}(a_t | s_{\leq t}, s'_{\leq t}, a_{<t}), \quad (18)$$

where, with a slight abuse of notation, we note that $R^*(E^*(o_t)) = r(o_t)$. We now further rewrite $\pi_{\text{aux}}(a_t | s_{\leq t}, s'_{\leq t}, a_{<t})$ as

$$p(a_t | s_{\leq t}, s'_{\leq t}, a_{<t}, \pi_{\text{aux}}), \quad (19)$$

and subsequently collapse the expression of the performance as

$$\eta(\pi_{\text{aux}}) = \sum_{(o, s, s', a)_{1:T}} R^*(s_{1:T}) p(o_{1:T}, s_{1:T}, s'_{1:T}, a_{1:T} | \pi_{\text{aux}}) \quad (20)$$

$$= \sum_{(s, s', a)_{1:T}} R^*(s_{1:T}) p(s_{1:T}, s'_{1:T}, a_{1:T} | \pi_{\text{aux}}), \quad (21)$$

where the last step arises from marginalization of $o_{1:T}$. Note by chain rule that $p(s_{1:T}, s'_{1:T}, a_{1:T} | \pi_{\text{aux}})$ becomes

$$\prod_t p(s_t | s_{<t}, s'_{<t}, a_{<t}, \pi_{\text{aux}}) p(s'_t | s_{\leq t}, s'_{<t}, a_{<t}, \pi_{\text{aux}}) p(a_t | s_{\leq t}, s'_{\leq t}, a_{<t}, \pi_{\text{aux}}). \quad (22)$$

By analyzing the Markov blankets in $p(s_{1:T}, s'_{1:T}, a_{1:T} | \pi_{\text{aux}})$, we can simplify the above expression to

$$\prod_t p(s_t | s_{<t}, s'_{<t}, a_{<t}) p(s'_t | s_{\leq t}, s'_{<t}, a_{<t}) p(a_t | s_{\leq t}, s'_{\leq t}, a_{<t}, \pi_{\text{aux}}). \quad (23)$$

Note that we omit the dependency on π_{aux} in the first two terms since, given only the history of past actions and observations, the next observation does not depend on our choice of policy but only on the environment dynamics.

Since E^* is optimal under the MI objective, we note that

$$I(S_{<t}, S'_{<t}, A_{<t}; S_t, S'_t) = I(S_{<t}, A_{<t}; S_t). \quad (24)$$

Eq. (24) implies that $s'_{<t}$ is independent of s_t given $(s_{<t}, a_{<t})$, and that $(s_{<t}, s'_{<t}, a_{<t})$ is independent of s'_t given s_t . This allow us to further simplify Eq. (23) to

$$\prod_t p(s_t | s_{<t}, a_{<t}) p(s'_t | s_t) \pi_{\text{aux}}(a_t | s_{\leq t}, s'_{\leq t}, a_{<t}). \quad (25)$$

Thus, the performance expression equates to

$$\eta(\pi_{\text{aux}}) = \sum_{\tau_{\text{aux}}} R^*(s_{1:T}) \prod_t p(s_t | s_{<t}, a_{<t}) p(s'_t | s_t) \pi_{\text{aux}}(a_t | s_{\leq t}, s'_{\leq t}, a_{<t}). \quad (26)$$

Note by way of similar reasoning (up to and including Eq. (23)) that

$$\eta(\pi) = \sum_{\tau} R^*(s_{1:T}) \prod_t p(s_t | s_{<t}, a_{<t}) \pi(a_t | s_{\leq t}, a_{<t}). \quad (27)$$

By comparing Eq. (26) and Eq. (27), we see that $s'_{1:T}$ effectively serves as a source of noise that makes π_{aux} behave like a stochastic policy depending on the seed choice for $s'_{1:T}$. To take advantage of this, we introduce a reparameterization of s' as ϵ such that

$$\eta(\pi_{\text{aux}}) = \sum_{\tau_{\text{aux}}} R^*(s_{1:T}) \prod_t p(s_t | s_{<t}, a_{<t}) p(\epsilon_t) \pi_{\text{aux}}(a_t | s_{\leq t}, \epsilon_{\leq t}, a_{<t}) \quad (28)$$

$$= \mathbb{E}_{p(\epsilon_{1:T})} \sum_{(s,a)_{1:T}} R^*(s_{1:T}) \prod_t p(s_t | s_{<t}, a_{<t}) \pi_{\text{aux}}(a_t | s_{\leq t}, \epsilon_{\leq t}, a_{<t}) \quad (29)$$

$$\leq \max_{\epsilon_{1:T}} \sum_{(s,a)_{1:T}} R^*(s_{1:T}) \prod_t p(s_t | s_{<t}, a_{<t}) \pi_{\text{aux}}(a_t | s_{\leq t}, \epsilon_{\leq t}, a_{<t}) \quad (30)$$

$$\leq \max_{\pi} \eta(\pi), \quad (31)$$

where the last inequality comes from defining a policy

$$\pi' := \pi_{\text{aux}}(a_t | o_{\leq t}, \epsilon_{\leq t}^*, a_{<t}) \quad (32)$$

and noting that the performance of π' must be bounded by the performance of π^* . \square

C. Additional Results

C.1. Comparison with CVRL after 2×10^6 environment steps

In Figure 6, we compare the performance of TPC and CVRL after 2×10^6 environment steps in both the standard setting and the natural background setting. TPC learns faster and achieves much higher rewards compared to CVRL in the standard setting. In the natural control setting, TPC outperforms in 4 out of 6 tasks and is competitive in Walker Run. Both methods do not work on Hopper Hop.

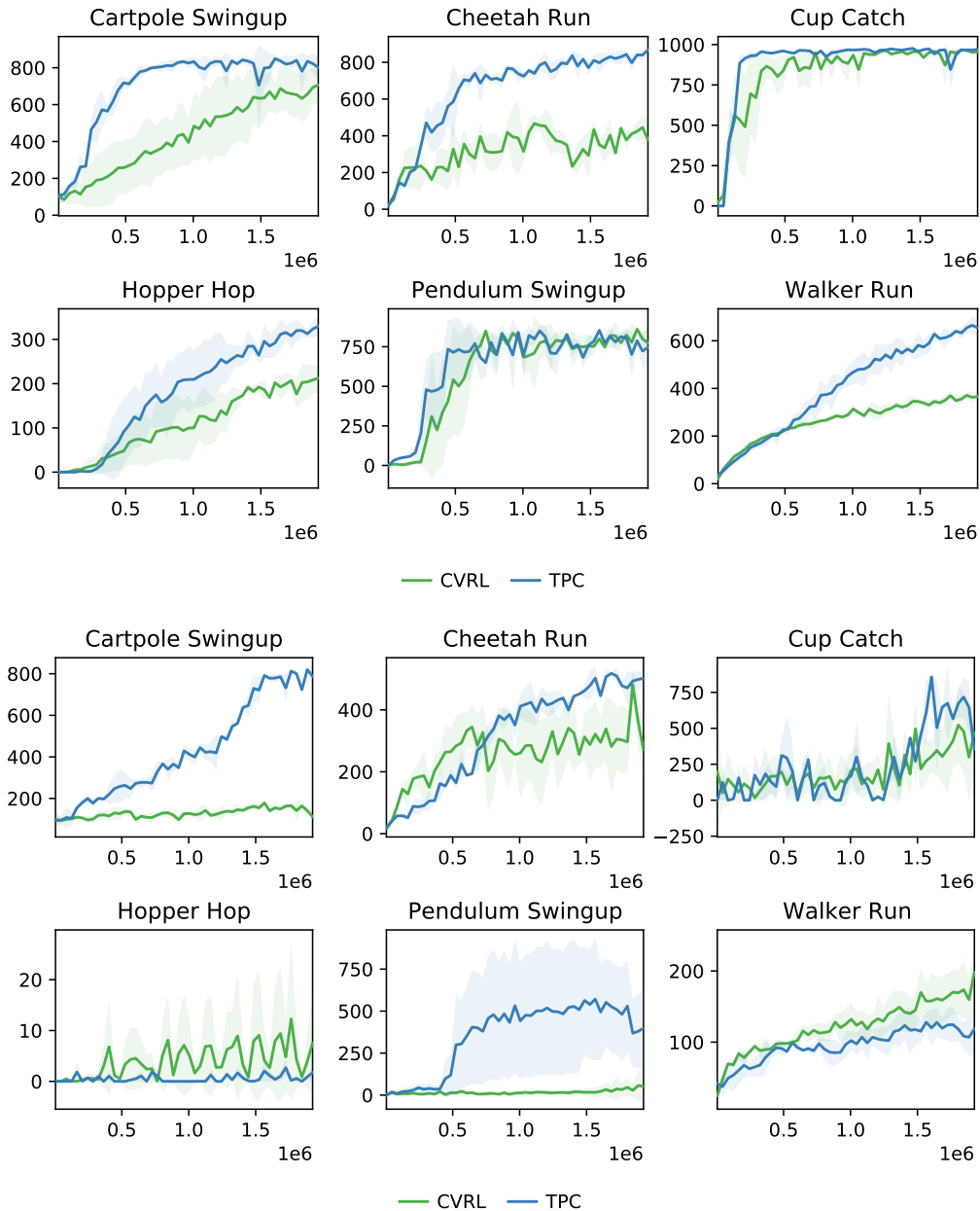


Figure 6. TPC vs CVRL after 2×10^6 in the standard setting (top) and background setting (bottom)

C.2. Comparison with CURL in the natural background setting

As shown in Figure 7, TPC outperforms CURL significantly on 3 of 6 tasks, while CURL performs better on Walker Run. On Hopper Hop and Cup Catch, both methods fail to make progress after 1 million environment steps.

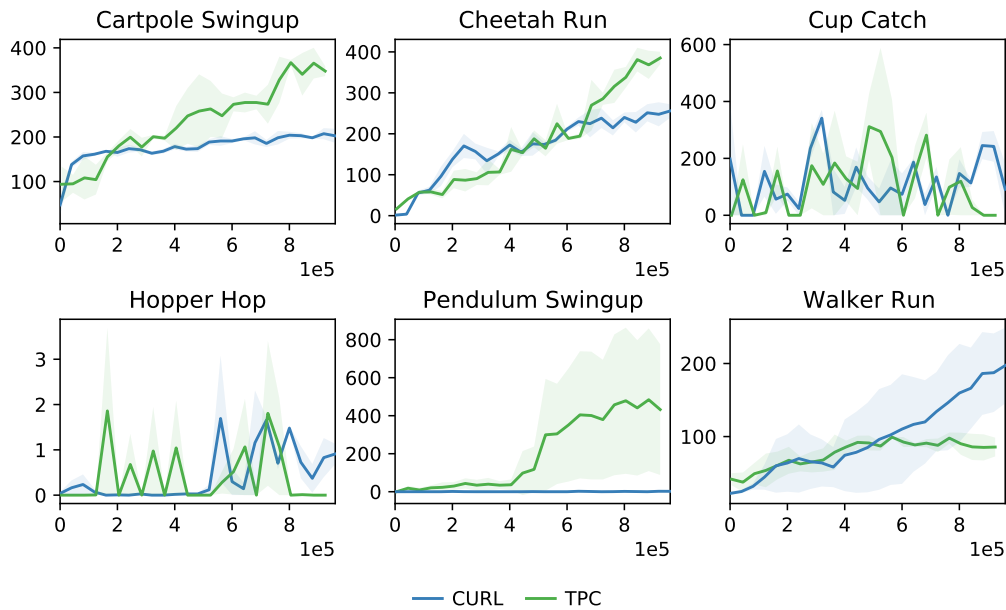


Figure 7. TPC versus CURL in the natural background setting. Each task is run with 3 seeds.

C.3. Importance of dynamics smoothing

We run TPC in the standard setting without dynamics smoothing to investigate the empirical importance of this component. As shown in Figure 8, TPC’s performance degrades significantly without dynamics smoothing. Without smoothing, the dynamics model cannot handle noisy rollouts during test-time planning, leading to poor performance. Dynamics smoothing prevents this by enabling test-time robustness against cascading error.

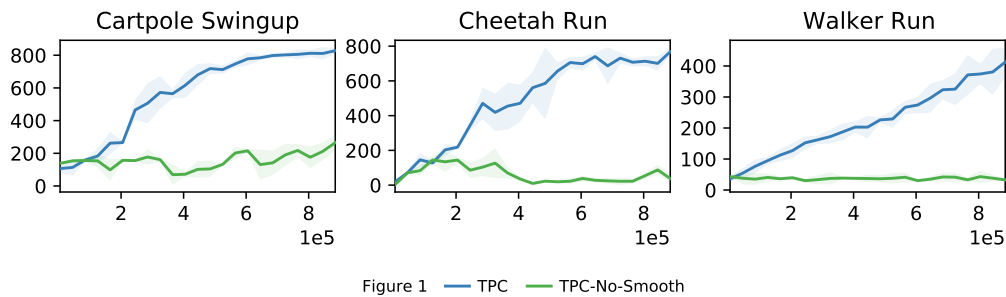


Figure 1 — TPC — TPC-No-Smooth

Figure 8. Performance of TPC in the standard setting without dynamics smoothing.

C.4. Learning a separate reward model

Joint learning of reward is crucial for all models in the natural background setting (see Appendix A.2). Since background information is also temporally-predictive, increasing the weight of reward loss encourages the model to focus more on the components that are important for reward learning. However, in the random background setting, since all temporally-predictive information is task-relevant, TPC *uniquely* can learn the reward separately, as shown in Figure 9.

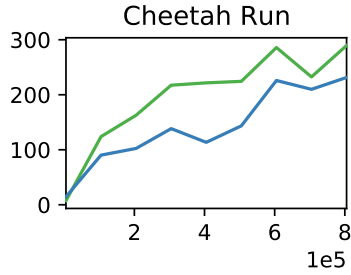


Figure 2 — TPC — TPC-R-Sep

Figure 9. Performance of TPC in the random background setting with a separately trained reward model.

C.5. Reconstructions in the natural background setting

We conduct experiments to investigate what information the encoder in different models learns to encode during training in the natural background setting. To do that, we train auxiliary decoders that try to reconstruct the original observations from the representations learned by Dreamer and TPC. As shown in Figure 10, Dreamer (2nd row) tries to encode both the agent and the background. In contrast, TPC (3rd row) prioritizes encoding the agent, which is task-relevant, over the background.

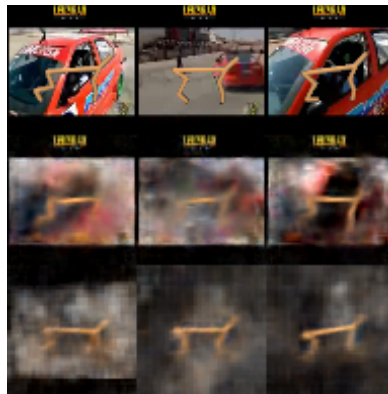


Figure 10. Observation reconstruction of TPC versus Dreamer in the natural background setting

C.6. The simplistic motion background setting



Figure 11. The top row shows a sample sequence of data, and the bottom row shows the reconstruction of TPC.

As discussed in Section 3.2, TPC can capture certain task-irrelevant information. However, TPC can choose to not encode the background *whenever* encoding only the agent is sufficient to maximize the mutual information. In the experiments, we found that forcing the model to predict well the reward helps the encoder focus more on the agent, which can be done by increasing the weight of reward loss. Dreamer, in contrast, must encode as much information about the observation as

possible to achieve a good reconstruction loss. To elaborate on this, we conducted an experiment where we replaced the natural background with a simplistic motion, easily predictable background, which is depicted in Figure 11. Figure 12 shows that TPC works well in this setting, and outperforms Dreamer significantly.

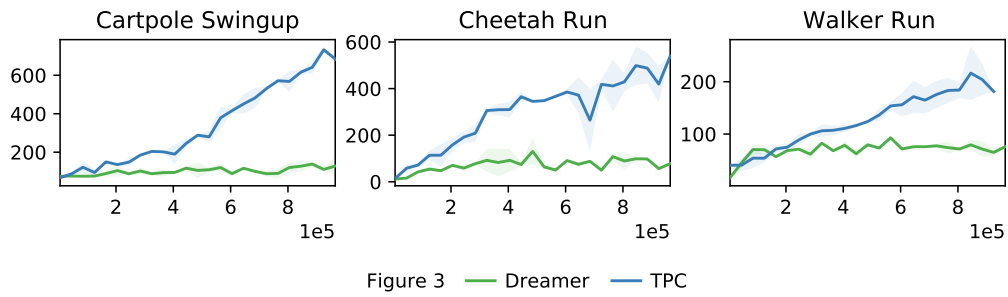


Figure 12. TPC versus Dreamer in the simplistic motion background setting