LEARNING TASK INFORMED ABSTRACTIONS

Xiang Fu, Ge Yang, Pulkit Agrawal, Tommi Jaakkola

Computer Science and Artificial Intelligence Laboratory Massachusetts Institute of Technology {xiangfu,geyang,pulkitag,tommi}@csail.mit.edu

ABSTRACT

Current model-based reinforcement learning methods struggle when operating from complex visual scenes due to their inability to prioritize task-relevant features. To mitigate this problem, we propose learning **Task Informed Abstractions** (TIA) that separates reward-correlated visual features from background distractions. For learning TIA, we introduce the formalism of Task Informed MDP (TiMDP), which is realized by training two models that learn visual features via cooperative reconstruction, but one model is adversarially dissociated from the reward signal. Empirical evaluation shows that TIA leads to significant performance gains over state-of-the-art methods on many visual control tasks where natural and unconstrained visual distractions pose a formidable challenge.

1 INTRODUCTION

Consider results of a simple experiment reported in Figure 1. We train state-of-the-art model-based reinforcement learning agents (Hafner et al., 2020) operating from visual inputs to solve two versions of the Cheetah Run task (Tassa et al., 2018): one with a simple and the other with a visually complex background (Zhang et al., 2021). For each version we train three model variants that only differ in in the capacity of their world model, which contain $0.5 \times$ (small), $1\times$ (medium) and $2\times$ (large) of the parameters of the original implementation. Not surprisingly, the performance with the simple background is only marginally affected by model



Figure 1: Comparison of the performance of a state-of-theart model-based RL algorithm, *Dreamer*, on two versions of the *Cheetah Run* with vs. without visual distraction. Performance is reported for three models of different sizes $(0.5\times, 1\times, 2\times$ of original Dreamer). Results show that even the smallest model has sufficient capacity to capture taskrelevant features when observations are distractor-free (gray), but when the scene is complex (red), task-irrelevant features consume most of the model capacity. Error bars indicate one standard deviation.

capacity, which shows that even the smallest model is sufficient for learning the features necessary for solving the task. Performance is much worse with the complex background, but it increases monotonically with the model size. Given that the task-relevant information (i.e., joint information of the *cheetah*) is the same in both, the performance improvements with model size indicates that the increase in representational capacity is used to encode the complex background. The background conveys no information about the task and therefore interferes with the learning of task-relevant information by consuming model capacity. Here "relevant" refers to those features that are needed to predict the optimal actions, whereas "irrelevant" refers to everything else that makes up the observation.

There are two main components to a model-based learner: (i) a forward dynamics model that predicts future events resulting from executing a sequence of actions from the current state and (ii) a reward predictor for evaluating possible future states. The policy performance critically depends on the prediction accuracy of this model, which is intimately tied to the feature space in which the future

^{*}Equal contribution.

is predicted. Similar to the complex background version of *Cheetah run*, the real world contains observations that are full of irrelevant content. Therefore to learn directly in the real-world, a compact feature space that only captures task-relevant information could make the learning problem much easier. Without this bias towards "task-relevant features", spurious features that confounds with the task would unnecessarily increase the data requirement or lead to training issues. Larger models will be also be needed for complex domains, despite that the dynamics of the core control problem is simple.

A popular choice for feature learning is to reconstruct the raw observations (Kingma & Welling, 2014; Kingma et al., 2014; Watter et al., 2015; Hafner et al., 2020). Often these features are encouraged to be *disentagled* (Bengio, 2013; Higgins et al., 2016; 2017) to identify distinct factors of variation. Since disentanglement simply re-formats the input space, the disentangled feature space would still contain irrelevant information and does not address the core issue of learning task-relevant features. Because the model-based agent also predicts rewards from the feature space, one might expect the combination of disentanglement and reward prediction sufficient to incentivize learning of task-relevant features. However, rewards provide insufficient supervision for feature learning (Yarats et al., 2019). For instance, just knowing the center of mass of a humanoid moving forward is sufficient to predict the reward, whereas it would require the full pose to come up with the optimal action. In a nutshell reconstruction captures too much information, whereas reward-prediction captures too little. Several works attempt to combine these two training signals (Hafner et al., 2020; Oh et al., 2017) but still struggle to learn in complex visual scenarios. Since the goal of the agent is to maximize the expected return, predicting the value function instead of one-step reward may aid in learning all the relevant information (Silver et al., 2017; Oh et al., 2017). However, because the value function is learned simultaneously with the model, it is not stationary and may not provide a stable training signal.

These challenges inspired several works to investigate feature learning methods that neither rely on reconstruction nor solely depend on rewards. One line of work biases the learned features to only capture controllable parts of the environment using an inverse model that predicts actions from a pair of states (Agrawal et al., 2015; Jayaraman & Grauman, 2015; Agrawal et al., 2016; Pathak et al., 2017), or using metrics such as empowerment (Klyubin et al., 2005; Gregor et al., 2016). To understand their shortcoming, consider the scenario of the arm pushing an object. Here both the arm and the object are controllable. While it is easy to capture the part that is directly controllable (e.g., the arm), capturing all controllable features (i.e., arm and the object) without imposing a reconstruction loss is non-trivial. Another idea that has shown promise is the bisimulation metric (Ferns et al., 2011; Zhang et al., 2021). Because supervision in bisimulation comes solely from rewards, it is subject to the same insufficiency mentioned earlier. Another possibility is to use contrastive learning (Chen et al., 2020; Oord et al., 2018), but without additional constraints, these methods do not distinguish between relevant and irrelevant features.

The ongoing discussion illustrates the fundamental challenge in learning task-relevant features: some objectives (e.g., reconstruction) capture too much information, whereas others (e.g., rewards, inverse models, empowerment) capture too little. Using a weighted loss function that combines these objectives has been empirically found not to learn task-relevant features (see Figure 1). In this work, we revisit feature learning by using reconstruction and the reward but propose to explicitly "explain away" irrelevant features by constructing a co-operative two-player game between two models. These models, dubbed as task and distractor, learn task-relevant (s_t^+) and irrelevant features (s_t^-) of the observation (o_t) respectively. Similar to prior work, we force the task model to learn task-relevant features (s_t^-) by predicting the reward. But unlike past work, we also force the distractor model to learn task-irrelevant features (s_t^-) via adversarial dissociation with the reward signal. However, both models cooperate to reconstruct o_t by maximizing $p(o_t | s_t^+, s_t^-)$.

Our method models a Markov decision process (MDP) of a specific factored structure, which we call **Task Informed MDP** (**TiMDP**) (see Figure 2b). It is worth noting that TiMDP is structurally similar to the *relaxed block MDP* (Zhang et al., 2020) formulation in partitioning the state-space into two separate components. However, Zhang et al. 2020 neither proposes a practical method for segregating relevant information nor does it provide any experimental validation in the context of learning from complex visual inputs. We evaluate our method on a custom ManyWorld environment, a suite of control tasks that specifically test the robustness of learning to visual distractions (Zhang et al., 2021) and ATARI games. The results convincingly demonstrate that our method, which we call

<u>**T**</u>ask <u>**I**</u>nformed <u>**A**</u>bstractions (TIA), successfully learns relevant features and outperforms existing state-of-the-art methods.

2 PRELIMINARIES

A Markov Decision Process is represented as the tuple $\langle S, O, A, T, r, \gamma, \rho_0 \rangle$ where O is a highdimensional observation space. A is the space of actions. S is the state space. ρ_0 is the initial state distribution. $r : S \mapsto \mathbb{R}$ is the scalar reward. The goal of RL is to learn a policy $\pi^*(a|s)$ that maximizes cumulative reward $\mathcal{J}_{\pi} = \arg \max_{\pi} \mathbb{E} \sum_t \gamma^{t-1} r_t$ discounted by γ .

Our primary contribution is in learning a feature space for forward dynamics and it is agnostic to the specific choice of the model-based algorithm. We choose to build upon the state-of-the-art method known as Dreamer (Hafner et al., 2020). The main components of this model are:

Representation model:
$$p_{\theta}(s_t \mid o_t, s_{t-1}, a_{t-1})$$
(1)Observation model: $q_{\theta}(o_t \mid s_t)$ $q_{\theta}(s_t \mid s_{t-1}, a_{t-1})$ Transition model: $q_{\theta}(s_t \mid s_{t-1}, a_{t-1})$ $q_{\theta}(s_t \mid s_t)$ Reward model: $q_{\theta}(r_t \mid s_t)$

Model Learning Dreamer (Hafner et al., 2019) makes future predictions in a feature space that is supervised by three signals: (a) image reconstruction $\begin{bmatrix} \mathcal{J}_{O}^{t} \doteq \ln q(o_{t} \mid s_{t}) \end{bmatrix}$, (b) reward prediction $\begin{bmatrix} \mathcal{J}_{R}^{t} \doteq \ln q(r_{t} \mid s_{t}) \end{bmatrix}$ and (c) dynamics regularization $\begin{bmatrix} \mathcal{J}_{D}^{t} \doteq -\beta \text{KL}(p(s_{t} \mid s_{t-1}, a_{t-1}, o_{t}) \parallel q(s_{t} \mid s_{t-1}, a_{t-1}) \end{bmatrix}$. The overall objective is:

$$\mathcal{J}_{\text{Dreamer}} \doteq \mathbb{E}_{\tau} \left[\sum_{t} \mathcal{J}_{\text{O}}^{t} + \mathcal{J}_{\text{R}}^{t} + \mathcal{J}_{\text{D}}^{t} \right]$$
(2)

optimized over the agent's experience τ . To achieve competitive performance on ATARI, a few modifications that are incorporated in this variant *DreamerV2* described in (Hafner et al., 2021).

Policy Learning Dreamer uses the learned forward dynamics model to train a policy using an actor-critic formulation described below:

Action model:
$$a_{\tau} \sim q_{\phi}(a_{\tau} \mid s_{\tau})$$

Value model: $v_{\psi}(s_{\tau}) \approx \mathbb{E}q(\cdot \mid s_{\tau}) \sum_{\tau=t}^{t+H} \gamma^{\tau-t} r_{\tau}$ (3)

The action model is trained to maximize cumulative rewards over a fixed horizon H. Both the action and value models are learned using imagined rollouts from the forward dynamics model. We refer the reader to (Hafner et al., 2020) for more details.

3 LEARNING TASK INFORMED ABSTRACTIONS

Task Informed MDP The state space of an MDP is often not directly observed in the real world, where high-dimensional observations are produced by latent states. We present the graphical model of this common scenario in Figure 2a. As discussed earlier the raw visual observations typically contain both task-relevant and irrelevant features. In order to explicitly segregate these factors, we propose to factor the state space S into two components: a task-relevant component S^+ and a task-irrelevant component S^- . We assume that the reward is fully determined by the task-relevant component $r : S^+ \mapsto \mathbb{R}$, and the task-irrelevant component contains no information about the reward: $MI(r_t; s_t^-) = 0$ at each time step t.

In the most general case, s_{t+1}^- can depend on s_t^+ and s_{t+1}^+ can depend on s_t^- . However, in many realistic scenarios the task vs distractor distinction often follows factored forward dynamics (Guestrin et al., 2003; Pitis et al., 2020) which greatly simplifies the learning model. For this reason we further incorporate this factored structure into our formulation through the assumption: $p(s_{t+1}|s_t, a_t) = p(s_{t+1}^+|s_t^+, a_t)p(s_{t+1}^-|s_t^-, a_t)$.

Our method involves learning two models: one model captures the task-relevant state component s_t^+ , which we call the *task model*. The other model captures the task-irrelevant state component



Figure 2: (a) The graphical model of an MDP. (b) Task-Informed MDP (TiMDP). The state space decomposes into two components: s_t^+ captures the task-relevant features, whereas s_t^- captures the task-irrelevant features. The cross-terms between $s^{+/-}$ are removed by imposing a factored MDP inductive bias. The red arrow indicates an adversarial loss to discourage s^- from picking up reward relevant information.



Figure 3: Components of Task Informed Abstraction Learning. (a) From the dataset of past experience, TIA uses the reward to factor the MDP into a task-relevant world model and a task-irrelevant one. (b) Only the forward dynamics in s_t^+ is used during policy learning. The Policy is trained using back-propagation through time. Note that the images are shown just for demonstration purposes and are not generated during policy learning.

 s_t^- , which we call the *distractor model*. The learning objective for these two models are denoted by \mathcal{J}_P and \mathcal{J}_S (task and distractor), and expanded in Equation (4). A visual illustration is provided in Figure 3a.

$$\begin{aligned}
\mathcal{J}_{\mathrm{P}} &\doteq \mathbb{E}_{p} \left(\sum_{t} \left(\mathcal{J}_{\mathrm{Oj}}^{t} + \mathcal{J}_{\mathrm{R}}^{t} + \mathcal{J}_{\mathrm{D}}^{t} \right) \right) \\
\mathcal{J}_{\mathrm{S}} &\doteq \mathbb{E}_{p} \left(\sum_{t} \left(\mathcal{J}_{\mathrm{Oj}}^{t} + \mathcal{J}_{\mathrm{Os}}^{t} + \mathcal{J}_{\mathrm{Radv}}^{t} + \mathcal{J}_{\mathrm{Ds}}^{t} \right) \right) \\
\mathcal{J}_{\mathrm{Oj}}^{t} &\doteq \ln q(o_{t} \mid s_{t}^{+}, s_{t}^{-}) \qquad \mathcal{J}_{\mathrm{Os}}^{t} &\doteq \lambda_{Os} \ln q(o_{t} \mid s_{t}^{-}) \\
\mathcal{J}_{\mathrm{R}}^{t} &\doteq \ln q(r_{t} \mid s_{t}^{+}) \qquad \mathcal{J}_{\mathrm{Radv}}^{t} &\doteq -\lambda_{\mathrm{Radv}} \max_{q} \ln q(r_{t} \mid s_{t}^{-}) \\
\mathcal{J}_{\mathrm{D}}^{t} &\doteq -\beta \mathrm{KL} \left(p(s_{t}^{+} \mid s_{t-1}^{+}, a_{t-1}, o_{t}) || q(s_{t}^{+} \mid s_{t-1}^{+}, a_{t-1}) \right) \\
\mathcal{J}_{\mathrm{Ds}}^{t} &\doteq -\beta \mathrm{KL} \left(p(s_{t}^{-} \mid s_{t-1}^{-}, a_{t-1}, o_{t}) || q(s_{t}^{-} \mid s_{t-1}^{-}, a_{t-1}) \right)
\end{aligned} \tag{4}$$

We will explain each component in the following section.

Reward Dissociation for the distractor model is accomplished via the adversarial objective \mathcal{J}_{Radv}^t . This is a minimax setup where we interleave optimizing the distractor model's reward prediction head (for multiple iterations/training step) with training the distractor model. While the reward prediction head is trained to minimize the reward prediction loss $-\ln q(r_t|s_t^-)$, the distractor model maximizes this objective, so as to exclude reward-correlated information from its learned features (Ganin & Lempitsky, 2015). The reward prediction loss is computed using $\ln \mathcal{N}(r_t; \hat{r}_t, 1)$, where $\mathcal{N}(\cdot; \mu, 1)$ is the unit Gaussian, and \hat{r}_t is the predicted reward.



Figure 4: Visualizing the information represented by Dreamer Hafner et al. (2020) and the task and distractor models of our method on several environments. (a) In the ManyWorld environments, Dreamer mistakes the distractor (yellow) for the target object (blue). The task model of TIA isolates the target object (blue) and the goal (red). (b, c) Dreamer's capacity is consumed at reconstructing the irrelevant video background, and it fails to capture the agent's outline, which is the task-relevant information. In all domains, Dreamer reconstruction tries to capture every pixel of the raw observation but misses task-relevant information. TIA is able to capture task-relevant information with the task model and task-irrelevant information with the distractor model.

Cooperative Reconstruction By jointly reconstructing the image, the distractor model that's biased towards capturing task-irrelevant information will enable the task model to focus on task-relevant features. We implement joint reconstruction through the objective \mathcal{J}_{Oj}^t . Starting with a sequence of observation and actions $\{o_{[<t]}, a_{[<t]}\}$, we first pass this sequence through the two separate recurrent state space model (RSSM, Hafner et al. 2019) to produce the states s_t^+ and s_t^- , which are then used to decode two images \hat{o}_t^+ and \hat{o}_t^- . Given the observation o_t , the joint reconstruction is achieved through a learned mixing where each model additionally decodes a $64 \times 64 \times 3$ tensor. These two tensors are concatenated channel-wise before being passed through a 1×1 convolution layer followed by sigmoid activation to obtain a $64 \times 64 \times 1$ mask M_t with the value between (0, 1). The final reconstruction is obtained through $\hat{o}_t = \hat{o}_t^+ \odot M_t + \hat{o}_t^- \odot (1 - M_t)$, where \odot denotes element-wise product. The reconstruction objective is computed as $\ln \mathcal{N}(o_t; \hat{o}_t, 1)$.

Distractor-model-only Reconstruction A degenerate case exists where a distractor model that captures no information at all can still satisfy the two objectives above, letting the task model reconstruct the entire observation by itself. To avoid such degeneracy, we add an additional image decoder to encourage the distractor model to capture as much information from the observation as possible. This is implemented via the objective \mathcal{J}_{Os}^t .

Policy Learning is similar to Dreamer, except that we replace the world model with the task model. This way, the forward unroll only involves the S^+ subspace.

Action model:
$$a_{\tau} \sim q_{\phi}(a_{\tau} \mid s_{\tau}^{+})$$

Value model: $v_{\psi}(s_{\tau}^{+}) \approx \mathbb{E}_{q(\cdot \mid s_{\tau}^{+})}(\sum_{\tau=t}^{t+H} \gamma^{\tau-t} r_{\tau})$ (5)

An illustration of the policy learning stage is in Figure 3b.

4 **EXPERIMENTS**

Our empirical evaluation aims to answer if our method outperforms existing methods when learning in environments that contain irrelevant information in the form of distractor objects or visually complex backgrounds. For this purpose we make use of three environments that are described in Section 4.1. Baseline wise, we compare our method against several alternatives described in Section 4.2.

4.1 Environments



Figure 5: Our method consistently outperforms Dreamer and other baseline methods in a variety of visual control tasks with distraction. The curves show mean and standard deviation, over five seeds for TIA, Dreamer, and Dreamer-Inverse. Results for DBC and DeepMDP are adopted from results reported in Zhang et al. (2021) and used ten seeds. Our method is effective for both ManyWorld environments (a,b), which contains confusing distracting objects that look similar to the task-relevant components; and the DMC tasks with natural video backgrounds (c,d,e), where the distracting background contains rich information that would consume significant model capacity to capture.

ManyWorld (Figure 4a) We want a test environment where one can vary the level of distraction in a controlled manner. For this purpose, we introduce ManyWorld, a physics domain where one can control the number of objects and their dynamics. The task is to move a target block (in blue) to a location indicated visually by a translucent red sphere. Other objects act as distractors. We turn off collision between objects so that they do not physically interfere with the target object, but occlusion does occur. The visual similarity between the objects introduces confusion, requiring additional effort to resolve when learning a world model.



Figure 6: (Left) Raw observation of ATARI Robotank. (Right) The task model of TIA emphasizes task-relevant information such as the crosshair and the radar for tracking enemies, while ignoring task-irrelevant information such as textures in the raw observation.

Kinematic Control with Natural Video Distraction

(Figure 4b) We consider the DeepMind Control (DMC) suite with natural video background from the Kinetics dataset (Kay et al., 2017) used in prior work (Zhang et al., 2021), which was introduced specifically to test learning under natural visual distraction. These control tasks involve different types of challenges, such as long planning horizon (Hopper), contact and collision (Walker), and larger state/action space (Cheetah). The natural video backgrounds in this test suite contain a large number of factors of variation, adding additional burden to world-model learning and negatively affecting policy optimization downstream. It is a challenging domain for both model-based and model-free algorithms to master. The video backgrounds are from the class "driving car" and are grayscale as in Zhang et al. (2021).

Atari Learning Environments (ALE) (Figure 6) are a standard benchmark for vision-based control. The visuals of these games naturally contain many distractor objects that are irrelevant to the game objective. Our limited compute resources only allowed us to experiment on five games. Each seed takes a week on a V100 Volta GPU. We present results on games where state-of-the-art model-based algorithms perform significantly worse than model-free algorithms or human performance in the hope of closing this gap.

4.2 **BASELINE METHODS**

We include both model-based and model-free baselines, a few proposed specifically to tackle learning in the presence of task-irrelevant distractions. In particular we compare against **Dreamer** (Hafner et al., 2020) which is a state of the art model-based algorithm on DMC. On ALE, we compare against an improved variant **Dreamer(V2)** (Hafner et al., 2021). We compare against a strong model-free method, **Deep Bisimulation for Control (DBC)** (Zhang et al., 2021), which uses the bisimulation metric and is developed specifically to tackle task-irrelevant distractions. Finally, we also include **DeepMDP** (Gelada et al., 2019) which learns a forward model with the sole purpose of acquiring a



Figure 7: ATARI performance at 50M steps. TIA significantly improves policy learning in Demon Attack, Robotank, and Yars Revenge, in which DreamerV2 fails to learn or has inferior sample complexity. We also add the performance of DreamerV2 at 100M and 200M steps, along with the performance of two model-free algorithms, DQN and Rainbow DQN, at 200M steps.

representation, then uses model-free, distributional Q learning for the policy. The DeepMDP and DBC results are adapted from Zhang et al. (2021).

Representation learning through an inverse model Inverse model take the observations o_t and o_{t+1} as input and predict the intervening action a_t (Agrawal et al., 2015; Jayaraman & Grauman, 2015). To investigate if features learned by inverse model suffice for learning task-relevant features, we constructed the *Dreamer-Inverse* model. In this model, the learning objective becomes the following: $\mathcal{J}_{\text{Inverse}}^t = \mathcal{J}_{\text{inv}}^t + \mathcal{J}_R^t + \beta \mathcal{J}_D^t$ where $\mathcal{J}_{\text{inv}}^t \doteq \ln q(a_t \mid s_t, s_{t+1})$ is the inverse model objective and \mathcal{J}_D^t is the dynamics regularizer described in Section 2.

4.3 CAN TIA DISASSOCIATE TASK-IRRELEVANT INFORMATION?

We first evaluate our method on the ManyWorld domain, where relevant information comprises the agent (blue block) and the goal (red sphere). Figure 4a provides a qualitative comparison of the information represented by TIA and the baseline method of *Dreamer*. In many cases, *Dreamer* mistakes the distractor (yellow) for the agent (blue). On the other hand, the task model of TIA isolates the agent (blue) and the goal (red), while the distractor model of TIA captures the distractors (yellow and green). This demonstrates that indeed s_t^+ captures task-relevant information and successfully ignores the rest. Quantitative evaluation depicted in Figure 5 (a) and (b) demonstrates that TIA outperforms the baselines, and the performance gap increases with the number of distractors.

Next, we considered the DMC domains with natural video distractions. We evaluate on three environments, *Cheetah Run*, *Walker Run* and *Hopper Stand* since they are visually different and cover a variety of learning challenges described in Section 4.1. The image reconstruction results Figure 4b and Figure 4c show that Dreamer performs poorly in capturing the full state of the agents and is distracted by the background. In contrast, the task model of our TIA method accurately recovers the relevant part of the raw observation, which happens to be the agent's body in these examples. Quantitative performance reported in Figure 5 (c,d,e) clearly shows that TIA outperforms strong baseline methods described in Section 4.2. Overall, these results suggest that TIA is the new state-of-the-art in learning from cluttered observations.

4.4 **RESULTS ON ATARI**

In previous test environments, the irrelevant factors were manually injected. Due to human intervention, it is possible these tasks were biased, which our method exploited. To investigate performance in more natural settings, we experimented on ATARI games that innately contain significant visual distractions. E.g., images of *robotank* game contain several visual signatures that change during the game: the number of enemy tanks destroyed, the rotating radar scanner, the green sprites, and so on. We evaluated performance on five ATARI games that are known to be challenging for model-based methods (Hafner et al., 2020) without any hyperparameter tuning (i.e., just a single value chosen based on the intuition described in Section 4.5).

The results reported in Figure 7 demonstrate that we substantially outperform the strong baseline of DreamerV2. Furthermore, for the games of *Chopper Command*, *Demon Attack* and *Yars Revenge* we either match or outperform strong model-free baselines of DQN/Rainbow trained for 200M

steps, while our method is only trained for 50M steps. These results convincingly demonstrate the superiority of our method. Figure 6 shows the image reconstruction of TIA's task model in *Robotank*.

4.5 HYPERPARAMETER SELECTION

The two important hyperparameters are $\lambda_{\rm Rady}$ and $\lambda_{\rm Os}$. One particular mode of failure is when the distractor model takes over the reconstruction. It strips the task model of task-relevant information, thus preventing the policy from learning any meaningful behavior. Our reward dissociation scheme relies on informative reward signals to work. Yet, at the beginning of training, the reward collected by a random policy tend to be sparse and noisy, making λ_{Radv} less effective at preventing a dominant distractor model. This scenario suggests using a large λ_{Rady} at the beginning of training and slowly increasing the weight λ_{Os} for the distractor reconstruction loss.



Figure 8: **Reward Dissociation During Learning** We plot the negative log-likelihood loss of the reward prediction $-\ln p(r)$, (lower bounded by 0.92) of a mean predictor, the reward prediction module of the primary model, and the reward prediction module of the secondary model. The features from the primary model contain sufficient information for reward prediction. The performance of the reward predictor for the secondary model follows the same trend as the mean predictor, indicating that the features learned by the secondary model are reward-independent.

The other extreme is for the distractor model

to collapse into degeneracy, where it fails to capture any information from the observations. TIA degenerates into Dreamer in this case. We can increase λ_{Os} so that the distractor model is encouraged to capture more.

To gain more insight into the reward dissociation process, we want to know how much information of the reward does the distractor model capture during learning. We use errors in predicting the reward, measured with the log-likelihood (as unit Gaussian) in Figure 8. We estimate the upper bound of this prediction error using the trailing average of the reward. This corresponds to an uninformed reward predictor (the "Mean Predictor") that always guesses the average. The reward prediction loss of the distractor model remains small (in blue), slightly above the loss lower bound, which equals to $-\ln \mathcal{N}(0; 0, 1) = -\ln \frac{1}{\sqrt{2\pi}} \approx 0.92$.

5 CONCLUSION AND DISCUSSION

In this work, we have shown that the TiMDP formulation that explains away task-irrelevant information can successfully learn from cluttered visual inputs. Our approach of learning $\underline{\mathbf{T}}$ ask $\underline{\mathbf{I}}$ nformed $\underline{\mathbf{A}}$ betractions (TIA) outperforms previous state-of-the-art model-based RL methods on multiple standard benchmarks.

An issue worth mentioning is that while one set of hyperparameters worked well across ATARI games, in the DMC suite, the choice of hyperparameters λ_{Radv} and λ_{Os} is domain-dependent. We discussed good practices for choosing these hyperparameters in Section 4.5. Based on these practices, it might be possible to automatically tune the hyperparameters by considering the reconstruction and reward-prediction loss of the two models that constitute TIA.

Our goal in this work is to build agents that operate from complex visual imagery. While we outperform previous methods, all of our evaluation is on simulated data. We plan to test our method on real-world data in the future. Another area of potential investigation is to characterize performance as a function of the sparsity of reward signals. We hypothesize that a low-noise reward signal is a key factor for the robustness of TIA, and the performance might drop in scenarios with sparse rewards. Developing methods that can overcome this "potential" challenge is another avenue for future work.

ACKNOWLEDGEMENT

This work is supported by the MIT-IBM grant on adversarial learning of multi-modal and structured data, the DARPA Machine Common Sense Program, and the National Science Foundation under Cooperative Agreement PHY-2019786 (The NSF AI Institute for Artificial Intelligence and Fundamental Interactions, http://iaifi.org/). The authors also acknowledge the MIT SuperCloud and the Lincoln Laboratory Supercomputing Center for providing interactive HPC resources.

REFERENCES

- Rishabh Agarwal, Marlos C. Machado, Pablo Samuel Castro, and Marc G Bellemare. Contrastive behavioral similarity embeddings for generalization in reinforcement learning. In *International Conference on Learning Representations*, 2021.
- Pulkit Agrawal, Joao Carreira, and Jitendra Malik. Learning to see by moving. In *Proceedings of the IEEE international conference on computer vision*, pp. 37–45, 2015.
- Pulkit Agrawal, Ashvin V Nair, Pieter Abbeel, Jitendra Malik, and Sergey Levine. Learning to poke by poking: Experiential learning of intuitive physics. In *Advances in Neural Information Processing Systems*, pp. 5074–5082, 2016.
- Hyojin Bahng, Sanghyuk Chun, Sangdoo Yun, Jaegul Choo, and Seong Joon Oh. Learning de-biased representations with biased representations. In *International Conference on Machine Learning*, pp. 528–539. PMLR, 2020.
- Ershad Banijamali, Rui Shu, Hung Bui, Ali Ghodsi, et al. Robust locally-linear controllable embedding. In *International Conference on Artificial Intelligence and Statistics*, pp. 1751–1759. PMLR, 2018.
- Yoshua Bengio. Deep learning of representations: Looking forward. In Adrian-Horia Dediu, Carlos Martín-Vide, Ruslan Mitkov, and Bianca Truthe (eds.), *Statistical Language and Speech Processing - First International Conference, SLSP 2013, Tarragona, Spain, July 29-31, 2013. Proceedings*, volume 7978 of *Lecture Notes in Computer Science*, pp. 1–37. Springer, 2013. doi: 10.1007/978-3-642-39593-2_1.
- Remi Cadene, Corentin Dancette, Matthieu Cord, Devi Parikh, et al. Rubi: Reducing unimodal biases for visual question answering. In *Advances in neural information processing systems*, pp. 841–852, 2019.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pp. 1597–1607. PMLR, 2020.
- Christopher Clark, Mark Yatskar, and Luke Zettlemoyer. Don't take the easy way out: Ensemble based methods for avoiding known dataset biases. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan (eds.), *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, pp. 4067–4080. Association for Computational Linguistics, 2019. doi: 10.18653/v1/D19-1418.
- Christopher Clark, Mark Yatskar, and Luke Zettlemoyer. Learning to model and ignore dataset bias with mixed capacity ensembles. In Trevor Cohn, Yulan He, and Yang Liu (eds.), *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: Findings, EMNLP 2020, Online Event, 16-20 November 2020*, pp. 3031–3045. Association for Computational Linguistics, 2020. doi: 10.18653/v1/2020.findings-emnlp.272.
- Thomas L. Dean, Robert Givan, and Sonia M. Leach. Model reduction techniques for computing approximately optimal solutions for markov decision processes. In Dan Geiger and Prakash P. Shenoy (eds.), UAI '97: Proceedings of the Thirteenth Conference on Uncertainty in Artificial Intelligence, Brown University, Providence, Rhode Island, USA, August 1-3, 1997, pp. 124–131. Morgan Kaufmann, 1997.

- Simon Du, Akshay Krishnamurthy, Nan Jiang, Alekh Agarwal, Miroslav Dudik, and John Langford. Provably efficient rl with rich observations via latent state decoding. In *International Conference* on Machine Learning, pp. 1665–1674. PMLR, 2019.
- Vladimir Feinberg, Alvin Wan, Ion Stoica, Michael I Jordan, Joseph E Gonzalez, and Sergey Levine. Model-based value estimation for efficient model-free reinforcement learning. arXiv preprint arXiv:1803.00101, 2018.
- Norm Ferns, Prakash Panangaden, and Doina Precup. Metrics for finite markov decision processes. In *UAI*, volume 4, pp. 162–169, 2004.
- Norm Ferns, Prakash Panangaden, and Doina Precup. Bisimulation metrics for continuous markov decision processes. *SIAM Journal on Computing*, 40(6):1662–1714, 2011.
- Norman Ferns and Doina Precup. Bisimulation metrics are optimal value functions. In *UAI*, pp. 210–219. Citeseer, 2014.
- Chelsea Finn and Sergey Levine. Deep visual foresight for planning robot motion. In 2017 IEEE International Conference on Robotics and Automation (ICRA), pp. 2786–2793. IEEE, 2017.
- Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In *International conference on machine learning*, pp. 1180–1189. PMLR, 2015.
- Carles Gelada, Saurabh Kumar, Jacob Buckman, Ofir Nachum, and Marc G. Bellemare. Deepmdp: Learning continuous latent space models for representation learning. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings* of Machine Learning Research, pp. 2170–2179. PMLR, 2019.
- Robert Givan, Thomas L. Dean, and Matthew Greig. Equivalence notions and model minimization in markov decision processes. *Artif. Intell.*, 147(1-2):163–223, 2003. doi: 10.1016/S0004-3702(02) 00376-4.
- Karol Gregor, Danilo Jimenez Rezende, and Daan Wierstra. Variational intrinsic control. *arXiv* preprint arXiv:1611.07507, 2016.
- Carlos Guestrin, Daphne Koller, Ronald Parr, and Shobha Venkataraman. Efficient solution algorithms for factored mdps. *Journal of Artificial Intelligence Research*, 19:399–468, 2003.
- David Ha and Jürgen Schmidhuber. Recurrent world models facilitate policy evolution. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett (eds.), Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada, pp. 2455–2467, 2018.
- Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and James Davidson. Learning latent dynamics for planning from pixels. In *International Conference on Machine Learning*, pp. 2555–2565. PMLR, 2019.
- Danijar Hafner, Timothy P. Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning behaviors by latent imagination. In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net, 2020.
- Danijar Hafner, Timothy P Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering atari with discrete world models. In *International Conference on Learning Representations*, 2021.
- He He, Sheng Zha, and Haohan Wang. Unlearn dataset bias in natural language inference by fitting the residual. In Colin Cherry, Greg Durrett, George F. Foster, Reza Haffari, Shahram Khadivi, Nanyun Peng, Xiang Ren, and Swabha Swayamdipta (eds.), *Proceedings of the 2nd Workshop on Deep Learning Approaches for Low-Resource NLP, DeepLo@EMNLP-IJCNLP 2019, Hong Kong, China, November 3, 2019*, pp. 132–142. Association for Computational Linguistics, 2019. doi: 10.18653/v1/D19-6115.

- Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. 2016.
- Irina Higgins, Arka Pal, Andrei Rusu, Loic Matthey, Christopher Burgess, Alexander Pritzel, Matthew Botvinick, Charles Blundell, and Alexander Lerchner. Darla: Improving zero-shot transfer in reinforcement learning. In *International Conference on Machine Learning*, pp. 1480–1490. PMLR, 2017.
- Dinesh Jayaraman and Kristen Grauman. Learning image representations tied to ego-motion. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1413–1421, 2015.
- Lukasz Kaiser, Mohammad Babaeizadeh, Piotr Milos, Blazej Osinski, Roy H. Campbell, Konrad Czechowski, Dumitru Erhan, Chelsea Finn, Piotr Kozakowski, Sergey Levine, Afroz Mohiuddin, Ryan Sepassi, George Tucker, and Henryk Michalewski. Model based reinforcement learning for atari. In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net, 2020.
- Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, 2017.
- Diederik P Kingma and Max Welling. Stochastic gradient vb and the variational auto-encoder. In *Second International Conference on Learning Representations, ICLR*, volume 19, 2014.
- Durk P Kingma, Shakir Mohamed, Danilo Jimenez Rezende, and Max Welling. Semi-supervised learning with deep generative models. In Z Ghahramani, M Welling, C Cortes, N D Lawrence, and K Q Weinberger (eds.), *Advances in Neural Information Processing Systems* 27, pp. 3581–3589. Curran Associates, Inc., 2014.
- Alexander S Klyubin, Daniel Polani, and Chrystopher L Nehaniv. Empowerment: A universal agent-centric measure of control. In 2005 IEEE Congress on Evolutionary Computation, volume 1, pp. 128–135. IEEE, 2005.
- Michael Laskin, Kimin Lee, Adam Stooke, Lerrel Pinto, Pieter Abbeel, and Aravind Srinivas. Reinforcement learning with augmented data. In Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, 2020a.
- Michael Laskin, Aravind Srinivas, and Pieter Abbeel. CURL: contrastive unsupervised representations for reinforcement learning. In Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event, volume 119 of Proceedings of Machine Learning Research, pp. 5639–5650. PMLR, 2020b.
- Kuang-Huei Lee, Ian Fischer, Anthony Liu, Yijie Guo, Honglak Lee, John Canny, and Sergio Guadarrama. Predictive information accelerates learning in RL. In Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, 2020.
- Bogdan Mazoure, Remi Tachet des Combes, Thang Doan, Philip Bachman, and R. Devon Hjelm. Deep reinforcement and infomax learning. In Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, 2020.
- R McCallum. Reinforcement learning with selective perception and hidden state. 1997.
- Junhyuk Oh, Satinder Singh, and Honglak Lee. Value prediction network. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (eds.), Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA, pp. 6118–6128, 2017.

- Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- Deepak Pathak, Pulkit Agrawal, Alexei A Efros, and Trevor Darrell. Curiosity-driven exploration by self-supervised prediction. In *Proceedings of the 34th International Conference on Machine Learning*, pp. 2778–2787, 2017.
- Karl Pertsch, Oleh Rybkin, Frederik Ebert, Shenghao Zhou, Dinesh Jayaraman, Chelsea Finn, and Sergey Levine. Long-horizon visual planning with goal-conditioned hierarchical predictors. *Advances in Neural Information Processing Systems*, 33, 2020.
- Silviu Pitis, Elliot Creager, and Animesh Garg. Counterfactual data augmentation using locally factored dynamics. In Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, 2020.
- Sébastien Racanière, Théophane Weber, David P Reichert, Lars Buesing, Arthur Guez, Danilo Rezende, Adria Puigdomenech Badia, Oriol Vinyals, Nicolas Heess, Yujia Li, et al. Imaginationaugmented agents for deep reinforcement learning. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pp. 5694–5705, 2017.
- Roberta Raileanu, Max Goldstein, Denis Yarats, Ilya Kostrikov, and Rob Fergus. Automatic data augmentation for generalization in deep reinforcement learning. *arXiv preprint arXiv:2006.12862*, 2020.
- Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Simon Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, et al. Mastering atari, go, chess and shogi by planning with a learned model. *Nature*, 588(7839):604–609, 2020.
- Pierre Sermanet, Corey Lynch, Yevgen Chebotar, Jasmine Hsu, Eric Jang, Stefan Schaal, Sergey Levine, and Google Brain. Time-contrastive networks: Self-supervised learning from video. In 2018 IEEE International Conference on Robotics and Automation (ICRA), pp. 1134–1141. IEEE, 2018.
- David Silver, Hado van Hasselt, Matteo Hessel, Tom Schaul, Arthur Guez, Tim Harley, Gabriel Dulac-Arnold, David P. Reichert, Neil C. Rabinowitz, André Barreto, and Thomas Degris. The predictron: End-to-end learning and planning. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the* 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017, volume 70 of Proceedings of Machine Learning Research, pp. 3191–3199. PMLR, 2017.
- Aravind Srinivas, Allan Jabri, Pieter Abbeel, Sergey Levine, and Chelsea Finn. Universal planning networks: Learning generalizable representations for visuomotor control. In *International Conference on Machine Learning*, pp. 4732–4741. PMLR, 2018.
- Austin Stone, Oscar Ramirez, Kurt Konolige, and Rico Jonschkowski. The distracting control suite–a challenging benchmark for reinforcement learning from pixels. *arXiv preprint arXiv:2101.02722*, 2021.
- Adam Stooke, Kimin Lee, Pieter Abbeel, and Michael Laskin. Decoupling representation learning from reinforcement learning. *arXiv preprint arXiv:2009.08319*, 2020.
- Aviv Tamar, Sergey Levine, Pieter Abbeel, Yi Wu, and Garrett Thomas. Value iteration networks. In Daniel D. Lee, Masashi Sugiyama, Ulrike von Luxburg, Isabelle Guyon, and Roman Garnett (eds.), Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain, pp. 2146–2154, 2016.
- Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David Budden, Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, Timothy Lillicrap, and Martin Riedmiller. Deepmind control suite. January 2018.

- Niklas Wahlström, Thomas B Schön, and Marc Peter Deisenroth. From pixels to torques: Policy learning with deep dynamical models. *arXiv preprint arXiv:1502.02251*, 2015.
- Haohan Wang, Songwei Ge, Zachary Lipton, and Eric P Xing. Learning robust global representations by penalizing local predictive power. In Advances in Neural Information Processing Systems, pp. 10506–10518, 2019.
- Manuel Watter, Jost Springenberg, Joschka Boedecker, and Martin Riedmiller. Embed to control: A locally linear latent dynamics model for control from raw images. *Advances in neural information processing systems*, 28:2746–2754, 2015.
- Wilson Yan, Ashwin Vangipuram, Pieter Abbeel, and Lerrel Pinto. Learning predictive representations for deformable objects using contrastive estimation. *arXiv preprint arXiv:2003.05436*, 2020.
- Denis Yarats, Amy Zhang, Ilya Kostrikov, Brandon Amos, Joelle Pineau, and Rob Fergus. Improving sample efficiency in model-free reinforcement learning from images. October 2019.
- Denis Yarats, Ilya Kostrikov, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. In *International Conference on Learning Representations*, 2021.
- Amy Zhang, Clare Lyle, Shagun Sodhani, Angelos Filos, Marta Kwiatkowska, Joelle Pineau, Yarin Gal, and Doina Precup. Invariant causal prediction for block mdps. In *International Conference* on Machine Learning, pp. 11214–11224. PMLR, 2020.
- Amy Zhang, Rowan Thomas McAllister, Roberto Calandra, Yarin Gal, and Sergey Levine. Learning invariant representations for reinforcement learning without reconstruction. In *International Conference on Learning Representations*, 2021.

A OTHER RELATED WORKS

State-of-the-art deep reinforcement learning algorithms often jointly optimize the the discounted return together with an auxiliary representation learning objective such as image reconstruction (Watter et al., 2015; Wahlström et al., 2015) or contrastive learnpredictive coding (CPC, Sermanet et al. 2018; Oord et al. 2018; Yan et al. 2020; Lee et al. 2020; Mazoure et al. 2020). In model-based reinforcement learning methods, reward and value prediction (Oh et al., 2017; Racanière et al., 2017; Silver et al., 2017; Tamar et al., 2016; Feinberg et al., 2018; Schrittwieser et al., 2020) are also shown to improve performance, along with learning a world model. More recently, data augmentations are found to improve sample complexity (Laskin et al., 2020a; Yarats et al., 2021; Raileanu et al., 2020; Laskin et al., 2020b; Stooke et al., 2020) by taking advantage of domain knowledge of symmetry transformations in the space of data. A recent benchmark (Stone et al., 2021) however, shows that data augmentation helps but still performs poorly in the presence of complex visual distractions. We do find data augmentation to be an orthogonal and complementary approach to our proposal, which focuses on the way to inform representation learning for RL, which feature is more useful, and therefore should be learned first.

A more principled approach exploits additional structure in the real world to learn state-abstractions. The most representative are block MDPs (Du et al., 2019), bisimulation (Givan et al., 2003; Dean et al., 1997) and bisimulation metrics (Ferns et al., 2004; 2011; Ferns & Precup, 2014), including recent impressive empirical gains on natural scene domains (Zhang et al., 2021; Gelada et al., 2019; Agarwal et al., 2021). The work presented here fall under the broad umbrella of learning state and temporal abstractions and is the most closely related to *utile distinction* (McCallum, 1997) which is limited to finite-state machines. Like *utile distinction*, TIA is distinct from DBC and bisimulation embeddings in that we rely on the task specification for feature separation, and we are not concerned about generalization across a set of MDPs related by dynamics.

The spiritual nearest neighbor to our work is the idea of using a primary-bias duo for model debiasing in supervised learning (Clark et al., 2019; Cadene et al., 2019; Wang et al., 2019; He et al., 2019; Bahng et al., 2020; Clark et al., 2020). These methods remove dataset bias by imposing independence constraints between the primary model and a model that's biased by design. Our approach focuses on RL and uses two models for removing task-irrelevant features for policy learning. Our implementation extends recent work in model-based RL from pixels (Watter et al., 2015; Finn & Levine, 2017; Banijamali et al., 2018; Ha & Schmidhuber, 2018; Kaiser et al., 2020; Hafner et al., 2020; 2021; Srinivas et al., 2018; Pertsch et al., 2020).

B PERFORMANCE GAP WHEN LEARNING IN THE PRESENCE OF VISUAL DISTRACTIONS

When complex visual distractors are present in the observations, state of the art model-based agents struggle to maintain their original sample efficiency and asymptotic performance. We produce this gap in details with Figure 9, where we include three additional (six in total) domains from the DeepMind Control suite. We arrange the domains in the order of difficulty levels, from easy to hard for the model-based approach. This gap is smaller on the two easiest domains, *Walker stand* and *Walker walk*; and is much more pronounced in *Walker run, Hopper stand*, and *Cheetah run*. On the most challenging task, *Finger spin*, where the model-based algorithm dreamer Hafner et al. (2020) performs much worse than model-free approaches, the gap almost disappears – making the domain ill-suited for testing our proposal. Note that the maximum episodic return on these tasks is calibrated to 1000 Tassa et al. (2018). We additionally provide the ratio for the number of pixels that are replaced by the background video for each domain, as a proxy for how much visual distraction is introduced: Walker 63%, Hopper 77%, Cheetah: 83%, and Finger 92%.

Our intent with this paper and the proposal, learning *Task Informed Abstractions*, is to close this gap, such that model-based agents can retain its performance even when learning in the presence of complex visual distractions.



Figure 9: The data efficiency and performance gap when learning with video distraction backgrounds.

C BRIDGING THE GAP BY LEARNING TASK INFORMED ABSTRACTIONS

We produce the complete result on DeepMind Control Suite including the 3 additional tasks from above where we only expect marginal improvements. We setup the experiments by matching the total number of parameters, so that each one of TIA's two models are only half in size as the single world model from dreamer. This comparison is disadvantage to our method, because the smaller task model runs the risk of being too small to capture the set of task-relevant features in its entirety.



Figure 10: **Performance of Task Informed Abstractions** on three additional DeepMind control domains.

Despite of this disadvantage in model capacity, task informed abstraction is able to reduce the gap on *Walker run* while making the gap significantly smaller on *Hopper stand* and *Cheetah run* (see Figure 9). On *Walker stand* and *Walker walk* the gap is small to begin with, therefore our method cannot bring much benefit. Finally, despite being a poor choice for testing our method, we include results on *Finger spin* for the sake of completeness.

D UNDERSTANDING FAILURE CASES

In Section 4.5 we provide the principled approach to tuning the hyperapameters λ_{Radv} and λ_{Os} , where we balance these two terms to avoid either one of the two models taking over the entire reconstruction. We label these two extreme cases, where either the distractor model, or the task model takes over, as type I, and type II. We provide detailed renderings of these failure cases below Figure 11, in comparison to a successfully learned world model.

The first column in Figure 11a shows a successfully trained agent whose task model is able to perfectly reconstruct the walker agent. In type I failure mode (Figure 11b) the distractor model would take over the entire reconstruction, causing the task model to lose its grasp on the task-relevant features. This would prevent the policy from learning any useful behavior, which makes the collected reward practically random. With random rewards, the term \mathcal{J}_{Radv} is ineffective at dissociating task-relevant features, causing training to fail. The policy performance under this scenario is usually close to that of a random policy. When this happens, we want to increase λ_{Radv} to dissociate the distractor model from the reward, or decrease λ_{Os} that weighs the distractor reconstruction term.

In the type II failure case (see Figure 11c), the task model takes over the reconstruction. In this case the model degenerates into dreamer without separating out the task-irrelevant features, and the performance is close to a dreamer agent with a smaller model. In this case we would increase the weight λ_{Os} on the distractor reconstruction to encourage it to learn more features.

By tuning these two parameters λ_{Radv} and λ_{Os} , we were able to avoid these two failure modes, and stabilize training on a wide variety of domains. A future direction would be to tune these to parameters automatically using the signals mentioned in Section 4.5.



Figure 11: Detailed rendering from one successfully trained walker agent, and two failure cases of type 1 and type 2. In type 1 failure case, the distractor model takes over the entire reconstruction, rendering the task model ineffective for policy learning (random policy). In the type 2 failure case, the task model attempts to capture all factors of variation by itself, thus failing to perfectly reconstruct the image, also leading to sub-par policy performance in the lower 400.

E MODEL SIZES AND ARCHITECTURE DETAILS

For fair comparison, on DMC and ManyWorld we match the total number of parameters of our two world models combined with that of a single, large Dreamer model. The total number of parameters is 10 million on DMC, and 1 million on ManyWorld. We divide the two models equally in size, with an additional image reconstruction head for the distractor model. *Dreamer-Inverse* has the same size for all model components as the Dreamer model except it replaces the deconvolution head with an action prediction head for learning the inverse dynamics. We scale the model size by changing the width of each layer in the networks, without changing the overall architecture or the depth of the networks. All other hyperparameters such as learning rates are kept the same as Hafner et al. (2020).

On the ATARI Learning Environments we compare against the state-of-the-art on this domain, DreamerV2 Hafner et al. (2021), which uses a world model of 20 million trainable parameters. We made the task model the same size as the original implementation while adding a smaller distractor model which contains 12 million parameters. A key difference between Dreamer and DreamerV2 is that the latter has an additional prediction head for the discount factor γ_t besides the standard reward prediction head. This discount factor head plays an instrumental role in allowing DreamerV2 to solve ATARI games, therefore we additionally dissociate the distractor model from information about the discount factor, by adding an additional adversarial prediction loss. We use the same scale for the discount factor γ as that for reward: $\lambda_{\gamma adv} = \lambda_{Radv}$. All other hyperparameters such as learning rates are kept the same as Hafner et al. (2021).

We use an input size of $32 \times 32 \times 3$ in ManyWorld, $64 \times 64 \times 3$ in DeepMind Control Suite, and $64 \times 64 \times 1$ in ATARI games. We use grayscale for the natural video backgrounds, the same as previous work Zhang et al. (2021).

F HYPERPARAMETERS

For fair comparison, we did not tweak existing hyperparameters from Dreamer and used identical settings as Hafner et al. (2020) and Hafner et al. (2021). Our reward-dissociation scheme introduces two new hyperparameters λ_{Radv} and λ_{Os} . We scale the reward dissociation loss via λ_{Radv} such that the term matches reconstruction losses in magnitude. For this reason, the differences in scale

in Table 1 mostly reflect the differences in input image sizes. We tweaked λ_{Os} to stabilize training. Detailed settings for each domain are in Table 1.

Domain and Task	$\lambda_{ m Radv}$	λ_{Os}
ManyWorld, 1 Distractor	600.0	2.0
ManyWorld, 2 Distractor	150.0	2.0
Hopper Stand	30k	2.0
Cheetah Run	20k	1.5
Walker Run	20k	0.25
Walker Walk	20k	0.25
Walker Stand	20k	0.25
Finger Spin	30k	2.5
All ATARI games	2k	1.0

Table 1: Hyperparameters

G TRANSFER TO NOVEL DISTRACTIONS

The task informed abstraction we introduce in this paper improves learning when distractions are present. To adapt to out-of-distribution scenarios unseen during training, additional architectural changes that reject distracting image features on the fly may be required. To provide a baseline and intuitions for this future direction, we evaluate how well existing agents perform under this type of domain shift. In Table 2 we take agents that are trained (1) without video background, (2) with background videos from the *driving car* class or (3) with white noise backgrounds, and evaluate against background videos from a different class, *walking the dog* (labeled as *transfer*, see Table 2).

 Table 2: DeepMind Control Transfer Performance
 transfer to the video class walking the dog as background

Trainin	g Condition	Drmr, No Bg	Drmr, Video	TIA, Video	Drmr, Noise	TIA, Noise
Hopper Stand	In-domain Transfer	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	$\begin{array}{c} 183.8 \pm 162.1 \\ 186.2 \pm 142.2 \end{array}$	$\begin{array}{c} 596.4 \pm 234.1 \\ 629.4 \pm 231.5 \end{array}$	$\begin{array}{c} 769.7 \pm 205.4 \\ 357.4 \pm 206.7 \end{array}$	$\begin{array}{c} 744.8 \pm 75.8 \\ 354.9 \pm 136.9 \end{array}$
Walker Run	In-domain Transfer	$\begin{vmatrix} 728.2 \pm 37.8 \\ 127.4 \pm 34.0 \end{vmatrix}$	520.2 ± 84.4 530.9 ± 76.9	625.3 ± 64.7 645.3 ± 78.4	$\begin{array}{c} 737.6 \pm 26.7 \\ 341.1 \pm 108.9 \end{array}$	$696.9 \pm 43.9 \\ 345.5 \pm 138.0$
Cheetah Run	In-domain Transfer	$\begin{vmatrix} 876.3 \pm 36.0 \\ 21.1 \pm 7.7 \end{vmatrix}$	325.7 ± 96.6 312.6 ± 115.7	$\begin{array}{c} 556.6 \pm 167.7 \\ 557.4 \pm 194.6 \end{array}$	$\begin{array}{c} 754.9 \pm 67.0 \\ 227.0 \pm 75.1 \end{array}$	$\begin{array}{c} 734.2 \pm 163.4 \\ 309.5 \pm 233.2 \end{array}$

The Dreamer agent trained with no background distraction fails to transfer its performance when background videos are introduced at test time, which is expected. In the second experiment we train both dreamer and TIA with video background, but test using videos from a different category. Dreamer did not learn as well as TIA as indicated by its poor performance in the training environments, but both methods retain their training performance post-transfer, unaffected by the change in the background video. As control, we also train both methods using white noise as the background. The training and transfer performance are both identical between the two methods, and the transfer performance is worse than performance on white noise.

We additionally evaluate ManyWorld agents that are trained with (1) no distractor, (2) one distractor, or (3) two distractors, with an additional distractor (three distractors, see Table 3).

 Table 3: ManyWorld Transfer Performance
 transfer to three distraction blocks.

Trainin	g Condition	Drmr, 0	Drmr, 1	Drmr, 2	TIA, 1	TIA, 2
ManyWorld	In-domain Transfer	$\begin{array}{ } 246.0 \pm 3.5 \\ 192.6 \pm 18.9 \end{array}$	242.9 ± 5.4 198.4 ± 27.4	$\begin{array}{c} 217.4 \pm 29.3 \\ 192.0 \pm 35.1 \end{array}$	246.1 ± 1.7 185.7 ± 21.4	245.8 ± 1.8 233.4 ± 6.5

Both results show that while TIA learns better from cluttered scenes, mechanism to reject unseen backgrounds at decision time is required to transfer successfully. This points to the incorporation of attention as a great avenue for future work.