

Learning Representations for Pixel-based Control: What Matters and Why?

Anonymous Authors¹

Abstract

Learning representations for pixel-based control has garnered significant attention recently in reinforcement learning. A wide range of methods have been proposed to enable efficient learning, leading to sample complexities similar to those in the full state setting. However, moving beyond carefully curated pixel data sets (centered crop, appropriate lighting, clear background, etc.) remains challenging. In this paper, we adopt a more difficult setting, incorporating background distractors, as a first step towards addressing this challenge. We start by exploring a simple baseline approach that does not use metric-based learning, data augmentations, world-model learning, or contrastive learning. We then analyze when and why previously proposed methods are likely to fail or reduce to the same performance as the baseline in this harder setting and why we should think carefully about extending such methods beyond the well curated environments. Our results show that finer categorization of benchmarks on the basis of characteristics like density of reward, planning horizon of the problem, presence of task-irrelevant components, etc., is crucial in evaluating algorithms. Based on these observations, we propose different metrics to consider when evaluating an algorithm on benchmark tasks. We hope such a data-centric view can motivate researchers to rethink representation learning when investigating how to best apply RL to real-world tasks.

1. Introduction

Learning useful representations for downstream tasks is a key component for success in rich observation environments [14, 39, 47, 54, 55]. Consequently, a significant amount of work proposes various representation learning objectives

that can be tied to the original reinforcement learning (RL) problem. Such auxiliary objectives include the likes of contrastive learning losses [42, 36, 11], state similarity metrics like bisimulation or policy similarity [62, 61, 1], and pixel reconstruction losses [29, 20, 25, 24]. On a separate axis, data augmentations have been shown to provide huge performance boosts when learning to control from pixels [35, 33]. Each of these methods has been shown to work well for particular settings and hence displayed promise to be part of a general purpose representation learning toolkit. Unfortunately, these methods were proposed with different motivations and tested on different tasks, making the following question hard to answer:

What really matters when learning representations for downstream control tasks?

Learning directly from pixels offers much richer applicability than when learning from carefully constructed states. Consider the example of a self-driving car, where it is nearly impossible to provide a complete state description of the position and velocity of all objects of interest, such as road edges, highway markers, other vehicles, etc. In such real world applications, learning from pixels offers a much more feasible option. However, this requires algorithms that can discern between task-relevant and task-irrelevant components in the pixel input, i.e., learn good representations. Focusing on task-irrelevant components can lead to brittle or non-robust behavior when put in slightly different environments. For instance, billboard signs over buildings in the background have no dependence on the task in hand while a self-driving car tries to change lanes. However, if such task-irrelevant components are not discarded, they can lead to sudden failure when the car drives through a different environment, say a forest where there are no buildings or billboards. Avoiding brittle behavior is therefore key to efficient deployment of artificial agents in the real world.

There has been a lot of work recently that tries to learn efficiently from pixels. A dominant idea throughout prior work has been that of attaching an auxiliary loss to the standard RL objective, with the exact mechanics of the loss varying for each method [29, 62, 36]. A related line of work learns representations by constructing world models directly from pixels [45, 41, 22, 25]. We show that these work well when the world model is simple. However, as the world

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

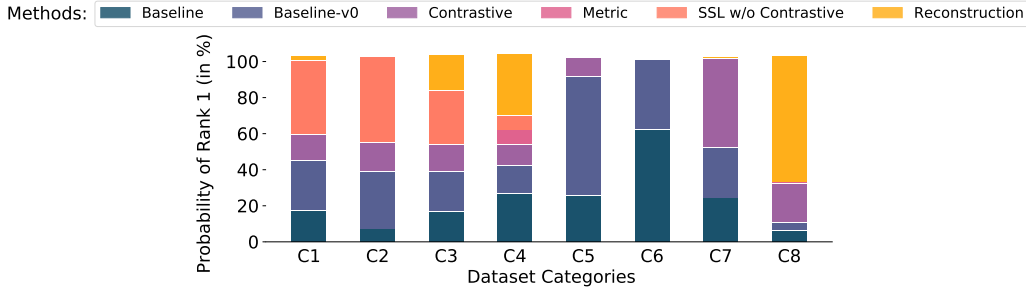


Figure 1: **Comparing pixel-based RL methods across finer categorizations** of evaluation benchmarks. Each category ‘Cx’ denotes different data-centric properties of the evaluation benchmark (e.g., C1 refers to discrete action, dense reward, without distractors, and with data randomly cropped [33, 35]). Exact descriptions of each category and the algorithms are provided in Table 4 and Table 5. Baseline-v0 refers to applying the standard deep RL agent (e.g., Rainbow DQN [53] and SAC [23]); Baseline refers to adding reward and transition prediction to baseline-v0, as described in Section 3; Contrastive includes algorithms such as PI-SAC [38] and CURL [36]; Metric denotes the state metric losses such as DBC [62]; SSL w/o Contrastive includes algorithms such as SPR [46]; Reconstruction includes DREAMER [25] and TIA [19]. For a given method, we always consider the best performing algorithm. Every method leads to varied performance across data categories, making a comparison which is an *average across all categories* highly uninformative.

model gets even slightly more complicated, which is true of the real world and imitated in simulation with the use of video distractors [60, 31, 48], such approaches can fail. For other methods, it is not entirely clear what component/s in auxiliary objectives can lead to failure, thus making robust behavior hard to achieve. Another distinct idea is of using data augmentations [35, 33] over the original observation samples, which seem to be quite robust across different environments. However, as we will show, a lot of the success of data augmentations is an artifact of how the benchmark environments save data, which is not replicable in the real world [48], thus resulting in failure¹. It is important to note that some of these methods are not designed for robustness but instead for enhanced performance on particular benchmarks. For instance, the ALE [7] benchmark involves simple, easy to model objects, and it becomes hard to discern if methods that perform well are actually good candidates for answering ‘what really matters for robust learning in the real world.’

Contributions. In this paper, we explore the major components responsible for the successful application of various representation learning algorithms. Based on recent work in RL theory for learning with rich observations [17, 4, 9], we hypothesize certain key components to be responsible for sample efficient learning. We test the role these play in previously proposed representation learning objectives and then consider an exceedingly simple *baseline* (see Figure 2) which takes away the extra “knobs” and instead combines two simple but key ideas, that of reward and transition prediction. We conduct experiments across multiple settings, including the MuJoCo domains from DMC Suite [51] with natural distractors [60, 31, 48], and Atari100K [30] from ALE [7]. Following this, we identify the failure modes of previously proposed objectives and highlight why they result in comparable or worse performance than the consid-

ered baseline. Our observations suggest that relying on a particular method across multiple evaluation settings does not work, as the efficacy varies with the exact details of the task, even within the same benchmark (see Figure 1). We note that a finer categorization of available benchmarks based on metrics like density of reward, presence of task-irrelevant components, inherent horizon of tasks, etc., play a crucial role in determining the efficacy of a method. We list such categorizations as suggestions for more informative future evaluations. The findings of this paper advocate for a more data-centric view of evaluating RL algorithms [13], largely missing in current practice. We hope the findings and insights presented in this paper can lead to better representation learning objectives for real-world applications.

2. Related Work

Prior work on **auxiliary objectives** includes the Horde architecture [50], UVFA [44] and the UNREAL agent [29]. These involve making predictions about features or pseudo-rewards, however only the UNREAL agent used these predictions for learning representations. Even so, the benchmark environments considered there always included only task-relevant pixel information, thus not pertaining to the hard setting we consider in this work. Representations can also be fit so as to obey certain state similarities. If these **state metrics** preserve the optimal policies and are easy to learn/given a priori, such a technique can be very useful. Recent works have shown that we can learn efficient representations either by learning the metrics like that in bisimulation [18, 62, 61, 8], by recursively sampling states [10] or by exploiting sparsity in dynamics [52]. **Data augmentations** modify the input image into different distinct views, each corresponding to a certain type of modulation in the pixel data. These include cropping, color jitter, flip, rotate, random convolution, etc. Latest works [35, 58] have shown that augmenting the states samples in the replay buffer with

¹It is also hard to pick exactly which data augmentation will work for a particular environment or task [43].

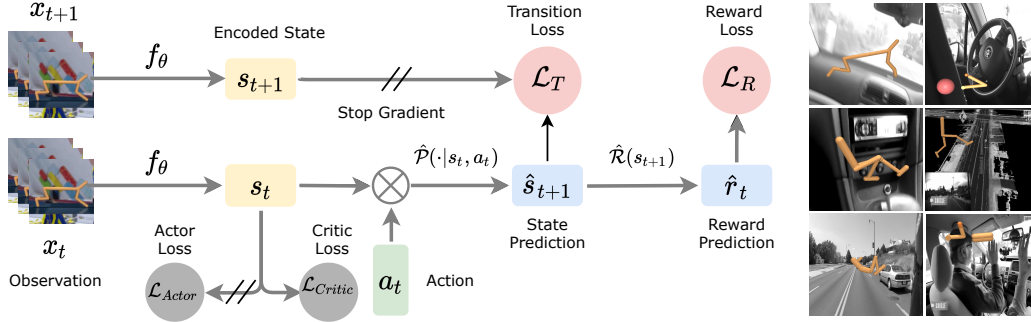


Figure 2: **(Left) Baseline for control over pixels.** We employ two losses besides the standard actor and critic losses, one being a reward prediction loss and the other a latent transition prediction loss. The encoded state s_t is the learnt representation. Gradients from both the transition/reward prediction and the critic are used to learn the representation, whereas the actor gradients are stopped. In the ALE setting, the actor and critic losses are replaced by a Rainbow DQN loss [53]. **(Right) Natural Distractor** in the background for standard DMC setting (left column) and custom off-center setting (right column). More details about the distractors can be found in Appendix 2.

such techniques alone can lead to impressive gains when learning directly from pixels. Recently, Stone et al. [48] illustrated the advantages and failure cases of augmentations. **Contrastive learning** involves optimizing for representations such that positive pairs (those coming from the same sample) are pulled closer while negative pairs (those coming from different samples) are pushed away [42, 11]. The most widely used method to generate positive/negative pairs is through various data augmentations [36, 46, 38, 49, 57]. However, temporal structure can induce positive/negative pairs as well. In such a case, the positive pair comes from the current state and the actual next state while the negative pair comes from the current state and any other next state in the current batch [42]. Other ways of generating positive/negative pairs can be through learnt state metrics [1] or encoding instances [25]. Another popular idea for learning representations is learning world models [22] in the pixel space. This involves learning prediction models of the world in the pixel space using a **pixel reconstruction** loss [20, 24, 25]. Other methods that do not explicitly learn a world model involve learning representations using reconstruction based approaches like autoencoders [56].

Quite a few papers in the past have analysed different sub-topics in RL through large scale studies. Engstrom et al. [15] and Andrychowicz et al. [3] have focused on analysing different policy optimization methods with varying hyperparameters. Our focus is specifically on representation learning methods that improve sample efficiency in pixel-based environments. Henderson et al. [27] showed how RL methods in general can be susceptible to lucky seedings. Recently, Agarwal et al. [2] proposed statistical metrics for reliable evaluation. Despite having similar structure, our work is largely complimentary to these past investigations. Babaeizadeh et al. [5] analysed reward and transition but only focused on the Atari 200M benchmark and pixel reconstruction methods. In comparison, our work is spread across multiple evaluation benchmarks, and our results show that

reconstruction can be a fine technique only in a particular benchmark category.

3. Method

We model the RL problem using the framework of contextual decision processes (CDPs), a term introduced in Krishnamurthy et al. [34] to broadly refer to any sequential decision making task where an agent must act on the basis of rich observations (context) x_t to optimize long-term reward. The true state of the environment s_t is not available and the agent must construct it on its own, which is required for acting optimally on the downstream task. Furthermore, the emission function which dictates what contexts are observed for a given state is assumed to only inject noise that is uncorrelated to the task in hand, i.e. it only changes parts of the context that are irrelevant to the task [62, 48]. Consider again the example of people walking on the sides of a road while a self-driving car changes lanes. Invariance to parts of the context that have no dependence on the task, e.g. people in the background, is an important property for any representation learning algorithm since we cannot expect all situations to remain exactly the same when learning in the real world. A more detailed description of the setup and all the prior methods used is provided in Appendix 1.

We start by exploring the utility of two fundamental components in RL, that of reward and transition prediction, in learning representations. A lot of prior work has incorporated these objectives either individually or in the presence of more nuanced architectures. Here, our aim is to start with the most basic components and establish their importance one by one. Particularly, we use a simple soft actor-critic setup (taking inspiration from SAC-AE [56]) as the base architecture, and attach the reward and transition prediction modules to it (See Figure 2). Note that the transition network is over the encoded state s_t and not over the observations [37]. Moreover, the transition model is fit between

the encoded state and the reward model. Unless noted otherwise, we call this architecture as the *baseline* for all our experiments. Details about the implementation, network sizes and all hyperparameters is provided in Appendix 3 and Appendix 4 (Table 3) respectively.

4. Empirical Study

In this section, we analyze the baseline architecture across six DMC tasks: Cartpole Swingup, Cheetah Run, Finger Spin, Hopper Hop, Reacher Easy, and Walker Walk. A common observation in our experiments is that the baseline is able to reduce the gap to more sophisticated methods significantly, sometimes even outperforming them in certain cases. This highlights that the baseline might serve as a stepping stone for other methods to build over. We test the importance of having both the reward and transition modules individually, by removing each of them one by one.

4.1. Reward Prediction

Figure 3 (left) shows a comparison of ‘with vs without reward prediction’. All other settings are kept unchanged and the only difference is the reward prediction. When the reward model is removed, there remains no grounding objective for the transition model. This results in a representation collapse as the transition model loss is minimized by the trivial representation which maps all observations to the same encoded state leading to degraded performance. This hints at the fact that without a valid grounding objective (in this case from predicting rewards), learning good representations can be very hard. Note that it is not the case that there is no reward information available to the agent, since learning the critic does provide enough signal to learn efficiently when there are no distractions present. However, in the presence of distractions the signal from the critic can be extremely noisy since it is based on the current value functions, which are not well developed in the initial stages of training. One potential fix for such a collapse is to not use the standard maximum likelihood based approaches for the transition model loss and instead use a contrastive version of the loss, which has been shown to learn general representations in the self-supervised learning setting. We test this later in the paper and observe that although it does help prevent collapse, the performance is still heavily inferior to when we include the reward model. Complete performances for individual tasks are shown in Appendix 8.1.

Linear Reward Predictor. We also compare to the case when the reward decoder is a linear network instead of the standard 1 layer MLP. We see that performance decreases significantly in this case as shown in Figure 3 (middle), but still does not collapse like in the absence of reward prediction. We hypothesize that the reward model is potentially

removing useful information for predicting the optimal actions. Therefore, when it is attached directly to the encoded state, i.e., in the linear reward predictor case, it might force the representation to only preserve information required to predict the reward well, which might not always be enough to predict the optimal actions well. For instance, consider a robot locomotion task. The reward in this case only depends on one variable, the center of mass, and thus the representation module would only need to preserve that in order to predict the reward well. However, to predict optimal actions, information about all the joint angular positions and velocities is required, which might be discarded if the reward model is directly attached to the encoded state. This idea is similar to why contrastive learning objectives in the self-supervised learning setting always enforce consistency between two positive/negative pairs *after* projecting the representation to another space. It has been shown that enforcing consistency in the representation space can remove excess information, which hampers final performance [11]. We indeed see a similar trend in the RL case as well.

4.2. Transition Prediction

Similarly, Figure 3 (right) shows a comparison of ‘with vs without transition prediction’. The transition model loss enforces temporal consistencies among the encoded states. When this module is removed, we observe a slight dip in performance across most tasks, with the most prominent drop in cartpole as shown in Appendix 8.1 (Figure 16). This suggests that enforcing such temporal consistencies in the representation space is indeed an important component for robust learning, but not a sufficient one. To examine if the marginal gain is an artifact of the exact architecture used, we explored other architectures in Appendix 8.2 but did not observe any difference in performance.

4.3. Connections to Value-Aware Learning

The baseline introduced above also resembles a prominent idea in theory, that of learning value aware models [17, 4]. Value-aware learning advocates for learning a model by fitting it to the value function of the task in hand, instead of fitting it to the true model of the world. The above baseline can be looked at as doing value aware learning in the following sense: the grounding to the representation is provided by the reward function, thus defining the components responsible for the task in hand and then the transition dynamics are learnt only for these components and not for all components in the observation space. There remains one crucial difference though. Value aware methods learn the dynamics based on the value function (multi-step) and not the reward function (1-step), since the value function captures the long term nature of the task in hand. To that end, we also test a more exact variant of the value-aware setup where we use the critic function as the target for optimizing the transi-

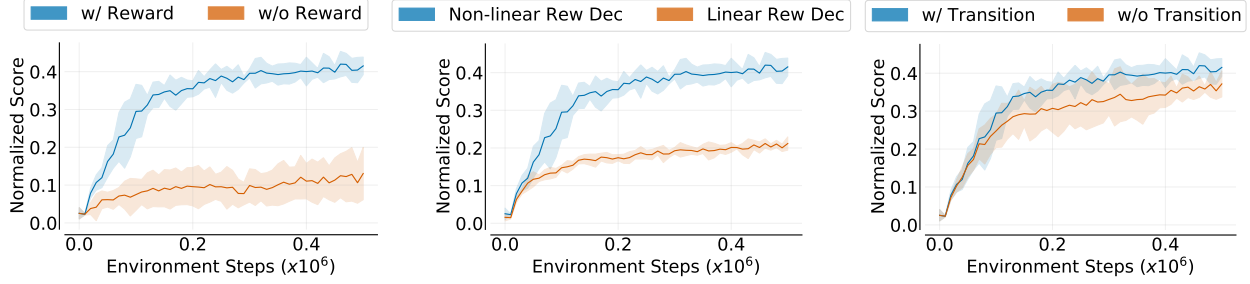


Figure 3: **Baseline Ablations.** Average normalized performance across six standard domains from DMC. Mean and std err. for 5 runs. **Left plot:** Baseline with vs without reward prediction **Middle plot:** Baseline with non-linear vs linear reward predictor/decoder. **Right plot:** Baseline with vs without transition prediction.

tion prediction, both with and without a reward prediction module (Table 1). Complete performances are provided in Appendix 8.8. We see that the value aware losses perform worse than the baseline. A potential reason for this could be that since the value estimates are noisy when using distractors, directly using these as targets inhibits learning a stable latent state representation. Indeed, more sophisticated value aware methods such as in Temporal Predictive Coding [40] lead to similar scores as the baseline.

Table 1: **Truly value-aware objectives.** We report average final score after 500K steps across six standard domains from DMC.

	Average Scores
Baseline	0.42 ± 0.02
Value-aware (w/ reward)	0.36 ± 0.03
Value-aware (w/o reward)	0.23 ± 0.03

5. Comparison

So far, we have discussed why the two modules we identify as being vital for minimal and robust learning are actually necessary. Now we ask what other components could be added to this architecture which might improve performance, as has been done in prior methods. We then ask when do these added components actually improve performance, and when do they fail. More implementation details are provided in Appendix 3.

Metric Losses. Two recent works that are similar to the baseline above are DBC [62] and MiCO [10], both of which learn representations by obeying a distance metric. DBC learns the metric by estimating the reward and transition models while MiCO uses transition samples to directly compute the metric distance. We compare baseline’s performance with DBC as shown in Figure 4 (left). Note that without the metric loss, DBC is similar to the baseline barring architectural differences such as the use of probabilistic transition models in DBC compared to deterministic models in the baseline. Surprisingly, we observe that the performance of the baseline exceeds that of DBC. To double check, we ran a version of DBC without the metric loss. Again, the “without metric” version lead to superior performance than

the “with metric” one (DBC).

Data Augmentations. A separate line of work has shown strong results when using data augmentations over the observation samples. These include the RAD [35] and DRQ [33] algorithms, both of which differ very minimally in their implementations. We run experiments for three different augmentations— ‘crop’, ‘flip’, and ‘rotate’. The ‘crop’ augmentation always crops the image by some shifted margin from the center. Interestingly, the image of the robot is also always centered, thus allowing ‘crop’ to always only remove background or task-irrelevant information and never remove the robot or task-relevant information. This essentially amounts to not having background distractors and thus we see that this technique performs quite well as shown in Figure 4 (middle). However, augmentations that do not explicitly remove the distractors, such as rotate and flip, lead to similar performance as the baseline. This suggests that augmentations might not be helpful when distractor information cannot be removed, or when we do not know where the objects of interest lie in the image, something true of the real world. We test this by shifting the robot to the side, thus making the task-relevant components off-center and by zooming out i.e. increasing the amount of irrelevant information even after cropping. We see that performance of ‘crop’ drops drastically in this case, showcasing that most of the performance gains from augmentations can be attributed to how the data is collected and not to the algorithm itself. Additional ablations are provided in Appendix 8.3.

Table 2: **RAD additional ablations.** We report average final score after 500K steps across Cheetah Run and Walker Walk domains from DMC. This illustrates that the performance of augmentations is susceptible to quality of data. Also, for the “Zoomed Out” setting, it is worth noting that both *crop* and *flip* settle to the same score.

	Standard	Off-center	Zoomed Out
RAD Crop	0.34 ± 0.14	0.30 ± 0.08	0.23 ± 0.10
RAD Flip	0.27 ± 0.08	0.29 ± 0.07	0.23 ± 0.07

Contrastive and SSL w/o Contrastive Losses. A lot of recent methods also deploy contrastive losses (for example, CPC [42]) to learn representations, which essentially refers to computing positive/negative pairs and pushing to-

²DBC [62] performance data is taken from their publication.

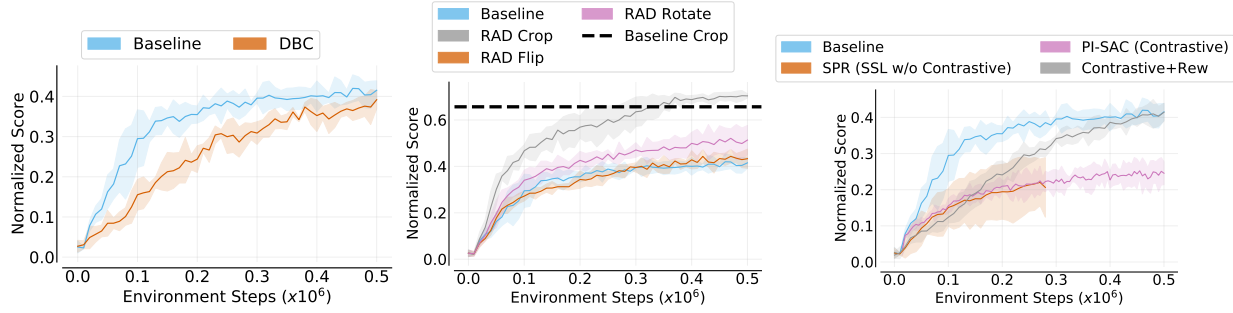


Figure 4: **Baseline Ablations.** Average normalized performance across six standard domains from DMC. Mean and std err. for 5 runs. **Left plot:** Baseline vs state metric losses (DBC [62]). The performance of baseline is compared with bisimulation metrics employed by DBC². **Middle plot:** Data Augmentations. Cropping removes irrelevant segments while flip and rotate do not, performing similar to the baseline. Baseline with random crop performs equally as good as RAD. **Right plot:** Contrastive and SSL w/o Contrastive. We replace the transition loss of the baseline with a contrastive version (Contrastive + Rew). Further, we consider simple contrastive (PI-SAC [38]) and SSL w/o contrastive (variant of SPR [46] for DMC) losses as well.

gether/pulling apart representations respectively. In practice, this can be done for any kind of loss function, such as the encoding function f_θ [25], or using random augmentations [36, 38], so on and so forth. Therefore, we test a simple modification to the baseline, that of using the contrastive variant of the transition prediction loss than the maximum likelihood version. We see, in Figure 4 (right), that the contrastive version leads to inferior results than the baseline, potentially suggesting that contrastive learning might not add a lot of performance improvement, particularly when there is grounding available from supervised losses. A similar trend has been witnessed in the self-supervised literature with methods like SIMSIAM [12], BARLOW TWINS [59], and BYOL [21] getting similar or better performance than contrastive methods like SIMCLR [11]. Complete performances are provided in Appendix 8.5.

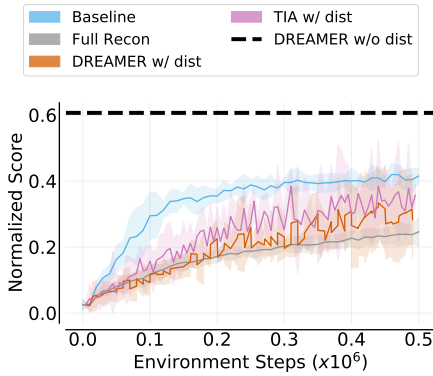


Figure 5: **Pixel reconstruction.** Average normalized performance across six DMC domains with distractors. Baseline achieves better performance than SOTA methods like DREAMER and TIA³.

SPR [46] is known to be a prominent algorithm in the ALE domain, leading to the best results overall. SPR deploys a specific similarity loss for transition prediction motivated by BYOL [21]. We follow the same setup and test a variant of the baseline which uses the cosine similarity loss from

³TIA [19] performance data is taken from their publication.

SPR and test its performance on DMC based tasks. We again show in Figure 4 (right) that there is very little or no improvement in performance as compared to the baseline performance. This again suggests that the same algorithmic idea can have an entirely different performance just by changing the evaluation setting⁴ (ALE to DMC).

Learning World Models. We test DREAMER [25], a state of the art model-based method that learns world models through pixel reconstruction on two settings, with and without distractors. Although the performance in the “without distractors” case is strong, we see that with distractors, DREAMER fails on some tasks, while performing inferior to the baseline in most tasks (see Figure 5). This suggests that learning world models through reconstruction might only be a good idea when the world models are fairly simple to learn. If world models are hard to learn, as is the case with distractors, reconstruction based learning can lead to severe divergence that results in no learning at all. We also compare against the more recently introduced method from Fu et al. [19]. Their method, called TIA[19] incorporates several other modules in addition to DREAMER and learns a decoupling between the distractor background and task relevant components. We illustrate the performance of each of the above algorithms in Figure 5 along with a version where we add full reconstruction loss to the baseline. Interestingly, TIA still fails to be superior to the baseline, particularly for simpler domains like Cartpole. Complete performances are provided in Appendix 8.6.

Relevant Reconstruction and Sparse Rewards. Since thus far we only considered dense reward based tasks, using the reward model for grounding is sufficient to learn good representations. More sophisticated auxiliary tasks consid-

⁴The SPR version without augmentations actually uses two separate ideas for improvement in performance, a cosine similarity transition prediction loss and a separate convolution encoder for the transition network, making it hard to attribute gains over the base DER [53] to just transition loss.

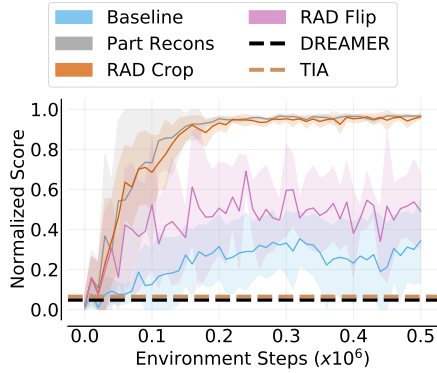


Figure 6: **Reconstruction and augmentations for sparse settings.** Normalized performance for ball-in-cup catch domain from DMC. ⁵Part Recons. in Figure 6 amounts to reconstructing the DMC robot over a solid black background. ⁶As also evident by TIA’s [19] performance for DMC ball-in-cup catch experiments. ⁷We use the DRQ(ϵ) version from Agarwal et al. [2] for fair

Atari 100K. We study the effect of techniques discussed thus far for the Atari 100K benchmark, which involves 26 Atari games and compares performance relative to human-achieved scores at 100K steps or 400K frames. We consider the categorization proposed by Bellemare et al. [6] based on the nature of reward (dense, human optimal, score exploit and sparse) and implement two versions of the baseline algorithm, one with both the transition and reward prediction modules and the other with only reward prediction. Our average results over all games show that the baseline performs comparably to CURL [36], SimPLe [30], DER [53], and OTR [32] while being quite inferior to DRQ⁷ [33, 2] and

⁵Part Recons. in Figure 6 amounts to reconstructing the DMC robot over a solid black background.

⁶As also evident by TIA’s [19] performance for DMC ball-in-cup catch experiments.

⁷We use the DRQ(ϵ) version from Agarwal et al. [2] for fair

SPR [46]. Since our implementation of the baseline is over the DER code, similar performance to DER might suggest that the reward and transition prediction do not help much in this benchmark. Note that ALE does not involve the use of distractors and so learning directly from the RL head (DQN in this case) should be enough to encode information about the reward and the transition dynamics in the representation. This comes as a stark contrast to the without distractors case in DMC Suite, where transition and reward prediction still lead to better performance. Such differences can also be attributed to the continuous *vs* discrete nature of DMC and ALE benchmarks. More interestingly, we find that when plotting the average performance for only the dense reward environments, the gap in performance between DER and SPR/DRQ decreases drastically. Note that SPR builds over DER but DRQ builds over OTR.

We further delve into understanding the superior performance of SPR and DRQ. In particular, SPR combines a cosine similarity transition prediction loss with data augmentations. To understand the effect of each of these individually, we run SPR without data augmentations, referring to this version by SPR^{**}⁸. We see that SPR^{**} leads to performance similar to the baseline and the DER agent, suggesting that such a self-supervised loss may not lead to gains when run without data augmentations. Finally, we take the DER agent and add data augmentations to it (from DRQ). This is shown as DER + AUG in Figure 7. We see that this leads to collapse, with the worst performance across all algorithms. Note that DRQ builds over OTR and performs quite well whereas when the same augmentations are used with DER, which includes a distributional agent in it, we observe a collapse. This again indicates that augmentations can change data in a fragile manner, sometimes leading to enhanced performance with certain algorithms, while failing with other algorithms. Segregating evaluation of algorithms based on these differences is therefore of utmost importance. We show the individual performance on all 25 games in Appendix 8.5 (Table 7).

6. Discussion

The above description of results on DMC Suite and Atari 100K point to a very interesting observation, that evaluation of different algorithms is very much correlated with a finer categorization of the evaluation benchmark, and not the whole benchmark itself. Specifically, focusing on finer categorizations such as density of reward, inherent horizon of the problem, presence of irrelevant and relevant task components, discreteness *vs* continuity of actions etc. is vital

evaluation and denote it as DRQ.

⁸Note that this is different from the SPR without augmentations version reported in Schwarzer et al. [46] since that version uses dropout as well which is not a fair comparison.

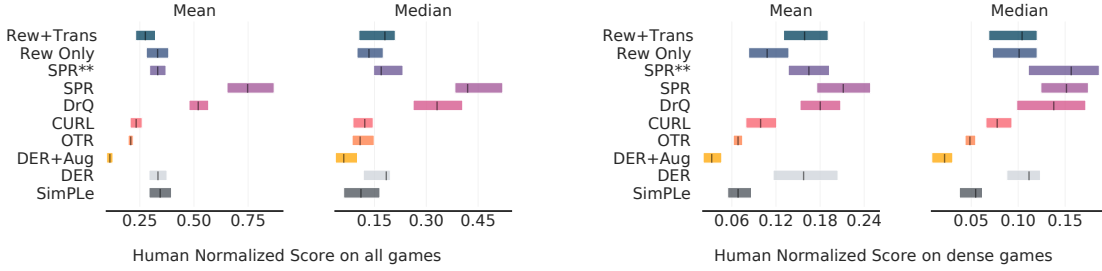


Figure 7: **Atari 100K**. Human normalized performance (mean/median) across 25 games from the Atari 100K benchmark. Mean and 95% confidence interval for 5 runs. **Left plot**: Comparison for all 25 games. **Right plot**: Comparison for only dense reward games (7 games from Table 7).

in recognizing if certain algorithms are indeed *better* than others. Figure 1 stands as a clear example of such discrepancies. These observations pave the way for a better evaluation protocol for algorithms, one where we rank algorithms for different categories, each governed by a specific data-centric property of the evaluation benchmark. Instead of saying that algorithm X is better than algorithm Y in benchmark Z, our results advocate for an evaluation methodology which claims algorithm X to be better than algorithm Y in dense reward, short horizon problems (considered from benchmark Z), i.e. enforcing less emphasis on the benchmark itself and more on certain properties of a subset of the benchmark. Having started with the question of what matters when learning representations over pixels, our experiments and discussion clearly show that largely it is the data-centric properties of the evaluation problems that matter the most.

7. Conclusion

In this paper we explore what components in representation learning methods matter the most for robust performance. As a starting point, we focused on the DMC Suite with distractors and the Atari 100k benchmark. Our results show that a simple baseline, one involving a reward and transition prediction modules can be attributed to a lot of performance benefits in DMC Suite with distractors. We then analysed why and when existing methods fail to perform as good or better than the baseline, also touching on similar observations on the ALE simulator. Some of our most interesting findings are as follows:

- Pixel reconstruction is a sound technique in the absence of clutter in the pixels, but suffers massively when distractors are added. In particular, DREAMER and adding a simple pixel reconstruction loss leads to worse performance than the baseline in DMC Suite (Figure 5).
- Contrastive losses in and of itself do not seem to provide gains when there is a supervised loss available in place of it. We observe that replacing the supervised state prediction loss of the baseline by the InfoNCE contrastive loss does not lead to performance improve-

ments over the baseline in DMC Suite (Figure 4 right plot). On the other hand, using contrastive losses with data augmentations can lead to more robust improvements [38, 16].

- Certain augmentations (‘crop’) do well when data is centered while dropping in performance when data is off-center or when cropping fails to remove considerable amounts of task-irrelevant information. Other augmentations (‘flip’ and ‘rotate’) show the opposite behavior (RAD ablations on DMC Suite in Table 2).
- SSL w/o contrastive losses does not provide much gains when used alone. With data augmentations, they lead to more significant gains. For Atari100k, Figure 7 shows that SPR, a state of the art non contrastive method leads to similar performance as the base DER agent when used without data augmentations (denoted by SPR**). Using the SPR inspired loss in DMC Suite also did not lead to gains over the baseline (in Figure 4 right plot).
- Augmentations are susceptible to collapse in the presence of distributional Q networks. Figure 7 shows that ‘crop’ and ‘intensity’ augmentations added to the DER agent lead to a complete failure in performance in Atari100k.

These results elicit the observation that claiming dominance over other methods for an entire benchmark may not be an informative evaluation methodology. Instead, focusing the discussion to a more data-centric view, one where specific properties of the environment are considered, forms the basis of a much more informative evaluation methodology. We argue that as datasets become larger and more diverse, the need for such an evaluation protocol would become more critical. We hope this work can provide valuable insights in developing better representation learning algorithms and spur further discussion in categorizing evaluation domains in more complex scenarios, such as with real world datasets and over a wider class of algorithmic approaches.

References

- [1] Agarwal, R., Machado, M. C., Castro, P. S., and Bellemare, M. G. Contrastive behavioral similarity embeddings for generalization in reinforcement learning. In *International Conference on Learning Representations*, 2021.
- [2] Agarwal, R., Schwarzer, M., Castro, P. S., Courville, A., and Bellemare, M. G. Deep reinforcement learning at the edge of the statistical precipice. *arXiv preprint arXiv:2108.13264*, 2021.
- [3] Andrychowicz, M., Raichuk, A., Stańczyk, P., Orsini, M., Girgin, S., Marinier, R., Hussenot, L., Geist, M., Pietquin, O., Michalski, M., et al. What matters in on-policy reinforcement learning? a large-scale empirical study. *arXiv preprint arXiv:2006.05990*, 2020.
- [4] Ayoub, A., Jia, Z., Szepesvari, C., Wang, M., and Yang, L. Model-based reinforcement learning with value-targeted regression. In *International Conference on Machine Learning*, pp. 463–474. PMLR, 2020.
- [5] Babaeizadeh, M., Saffar, M. T., Hafner, D., Erhan, D., Kannan, H., Finn, C., and Levine, S. On trade-offs of image prediction in visual model-based reinforcement learning. 2020.
- [6] Bellemare, M., Srinivasan, S., Ostrovski, G., Schaul, T., Saxton, D., and Munos, R. Unifying count-based exploration and intrinsic motivation. *Advances in neural information processing systems*, 29:1471–1479, 2016.
- [7] Bellemare, M. G., Naddaf, Y., Veness, J., and Bowling, M. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279, 2013.
- [8] Biza, O., Platt, R., van de Meent, J.-W., and Wong, L. L. Learning discrete state abstractions with deep variational inference. *arXiv preprint arXiv:2003.04300*, 2020.
- [9] Castro, P. S. Scalable methods for computing state similarity in deterministic markov decision processes. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 10069–10076, 2020.
- [10] Castro, P. S., Kastner, T., Panangaden, P., and Rowland, M. Mico: Learning improved representations via sampling-based state similarity for markov decision processes. *CoRR*, abs/2106.08229, 2021. URL <https://arxiv.org/abs/2106.08229>.
- [11] Chen, T., Kornblith, S., Norouzi, M., and Hinton, G. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pp. 1597–1607. PMLR, 2020.
- [12] Chen, X. and He, K. Exploring simple siamese representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15750–15758, 2021.
- [13] Co-Reyes, J. D., Sanjeev, S., Berseth, G., Gupta, A., and Levine, S. Ecological reinforcement learning. *arXiv preprint arXiv:2006.12478*, 2020.
- [14] Du, S., Krishnamurthy, A., Jiang, N., Agarwal, A., Dudik, M., and Langford, J. Provably efficient rl with rich observations via latent state decoding. In *International Conference on Machine Learning*, pp. 1665–1674. PMLR, 2019.
- [15] Engstrom, L., Ilyas, A., Santurkar, S., Tsipras, D., Janoos, F., Rudolph, L., and Madry, A. Implementation matters in deep rl: A case study on ppo and trpo. In *International conference on learning representations*, 2019.
- [16] Fan, J. and Li, W. Robust deep reinforcement learning via multi-view information bottleneck. *arXiv preprint arXiv:2102.13268*, 2021.
- [17] Farahmand, A.-m., Barreto, A., and Nikovski, D. Value-aware loss function for model-based reinforcement learning. In *Artificial Intelligence and Statistics*, pp. 1486–1494. PMLR, 2017.
- [18] Ferns, N., Panangaden, P., and Precup, D. Bisimulation metrics for continuous markov decision processes. *SIAM J. Comput.*, 40(6):1662–1714, December 2011. ISSN 0097-5397. doi: 10.1137/10080484X. URL <https://doi.org/10.1137/10080484X>.
- [19] Fu, X., Yang, G., Agrawal, P., and Jaakkola, T. Learning task informed abstractions. In Meila, M. and Zhang, T. (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 3480–3491. PMLR, 18–24 Jul 2021. URL <http://proceedings.mlr.press/v139/fu21b.html>.
- [20] Gelada, C., Kumar, S., Buckman, J., Nachum, O., and Bellemare, M. G. DeepMDP: Learning continuous latent space models for representation learning. In Chaudhuri, K. and Salakhutdinov, R. (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 2170–2179. PMLR, 09–15 Jun 2019. URL <https://proceedings.mlr.press/v97/gelada19a.html>.

- [21] Grill, J.-B., Strub, F., Althé, F., Tallec, C., Richemond, P. H., Buchatskaya, E., Doersch, C., Pires, B. A., Guo, Z. D., Azar, M. G., et al. Bootstrap your own latent: A new approach to self-supervised learning. *arXiv preprint arXiv:2006.07733*, 2020.
- [22] Ha, D. and Schmidhuber, J. World models. *arXiv preprint arXiv:1803.10122*, 2018.
- [23] Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In Dy, J. and Krause, A. (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 1861–1870. PMLR, 10–15 Jul 2018. URL <https://proceedings.mlr.press/v80/haarnoja18b.html>.
- [24] Hafner, D., Lillicrap, T., Fischer, I., Villegas, R., Ha, D., Lee, H., and Davidson, J. Learning latent dynamics for planning from pixels. In *International Conference on Machine Learning*, pp. 2555–2565. PMLR, 2019.
- [25] 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=S1lOTC4tDS>.
- [26] Hafner, D., Lillicrap, T. P., Norouzi, M., and Ba, J. Mastering atari with discrete world models. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=0oabwyZbOu>.
- [27] Henderson, P., Islam, R., Bachman, P., Pineau, J., Precup, D., and Meger, D. Deep reinforcement learning that matters. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- [28] Hessel, M., Modayil, J., Van Hasselt, H., Schaul, T., Ostrovski, G., Dabney, W., Horgan, D., Piot, B., Azar, M., and Silver, D. Rainbow: Combining improvements in deep reinforcement learning. In *Thirty-second AAAI conference on artificial intelligence*, 2018.
- [29] Jaderberg, M., Mnih, V., Czarnecki, W. M., Schaul, T., Leibo, J. Z., Silver, D., and Kavukcuoglu, K. Reinforcement learning with unsupervised auxiliary tasks. *arXiv preprint arXiv:1611.05397*, 2016.
- [30] Kaiser, L., Babaeizadeh, M., Milos, P., Osinski, B., Campbell, R. H., Czechowski, K., Erhan, D., Finn, C., Kozakowski, P., Levine, S., et al. Model-based reinforcement learning for atari. *arXiv preprint arXiv:1903.00374*, 2019.
- [31] 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.
- [32] Kielak, K. Importance of using appropriate baselines for evaluation of data-efficiency in deep reinforcement learning for atari. *arXiv preprint arXiv:2003.10181*, 2020.
- [33] Kostrikov, I., Yarats, D., and Fergus, R. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. *CoRR*, abs/2004.13649, 2020. URL <https://arxiv.org/abs/2004.13649>.
- [34] Krishnamurthy, A., Agarwal, A., and Langford, J. Pac reinforcement learning with rich observations. *arXiv preprint arXiv:1602.02722*, 2016.
- [35] Laskin, M., Lee, K., Stooke, A., Pinto, L., Abbeel, P., and Srinivas, A. Reinforcement learning with augmented data. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M. F., and Lin, H. (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 19884–19895. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/e615c82aba461681ade82da2da38004a-Paper.pdf>.
- [36] Laskin, M., Srinivas, A., and Abbeel, P. Curl: Contrastive unsupervised representations for reinforcement learning. *Proceedings of the 37th International Conference on Machine Learning, Vienna, Austria, PMLR 119*, 2020. arXiv:2004.04136.
- [37] Lee, A. X., Nagabandi, A., Abbeel, P., and Levine, S. Stochastic latent actor-critic: Deep reinforcement learning with a latent variable model. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M. F., and Lin, H. (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 741–752. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/08058bf500242562c0d031ff830ad094-Paper.pdf>.
- [38] Lee, K.-H., Fischer, I., Liu, A., Guo, Y., Lee, H., Canny, J., and Guadarrama, S. Predictive information accelerates learning in rl. *arXiv preprint arXiv:2007.12401*, 2020.
- [39] 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.

- [40] Nguyen, T., Shu, R., Pham, T., Bui, H., and Ermon, S. Temporal predictive coding for model-based planning in latent space. *arXiv preprint arXiv:2106.07156*, 2021.
- [41] Oh, J., Guo, X., Lee, H., Lewis, R., and Singh, S. Action-conditional video prediction using deep networks in atari games. *arXiv preprint arXiv:1507.08750*, 2015.
- [42] Oord, A. v. d., Li, Y., and Vinyals, O. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- [43] Raileanu, R., Goldstein, M., Yarats, D., Kostrikov, I., and Fergus, R. Automatic data augmentation for generalization in deep reinforcement learning. *arXiv preprint arXiv:2006.12862*, 2020.
- [44] Schaul, T., Horgan, D., Gregor, K., and Silver, D. Universal value function approximators. In *International conference on machine learning*, pp. 1312–1320. PMLR, 2015.
- [45] Schmidhuber, J. Formal theory of creativity, fun, and intrinsic motivation (1990–2010). *IEEE Transactions on Autonomous Mental Development*, 2(3):230–247, 2010.
- [46] Schwarzer, M., Anand, A., Goel, R., Hjelm, R. D., Courville, A., and Bachman, P. Data-efficient reinforcement learning with self-predictive representations. *arXiv preprint arXiv:2007.05929*, 2020.
- [47] Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., Hubert, T., Baker, L., Lai, M., Bolton, A., et al. Mastering the game of go without human knowledge. *nature*, 550(7676):354–359, 2017.
- [48] Stone, A., Ramirez, O., Konolige, K., and Jonckhowski, R. The distracting control suite—a challenging benchmark for reinforcement learning from pixels. *arXiv preprint arXiv:2101.02722*, 2021.
- [49] Stooke, A., Lee, K., Abbeel, P., and Laskin, M. Decoupling representation learning from reinforcement learning. In *International Conference on Machine Learning*, pp. 9870–9879. PMLR, 2021.
- [50] Sutton, R. S., Modayil, J., Delp, M., Degris, T., Pilarski, P. M., White, A., and Precup, D. Horde: A scalable real-time architecture for learning knowledge from unsupervised sensorimotor interaction. In *The 10th International Conference on Autonomous Agents and Multiagent Systems-Volume 2*, pp. 761–768, 2011.
- [51] 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.
- [52] Tomar, M., Zhang, A., Calandra, R., Taylor, M. E., and Pineau, J. Model-invariant state abstractions for model-based reinforcement learning. *arXiv preprint arXiv:2102.09850*, 2021.
- [53] van Hasselt, H. P., Hessel, M., and Aslanides, J. When to use parametric models in reinforcement learning? In Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL <https://proceedings.neurips.cc/paper/2019/file/1b742ae215adff18b75449c6e272fd92d-Paper.pdf>.
- [54] Wahlström, N., Schön, T. B., and Deisenroth, M. P. From pixels to torques: Policy learning with deep dynamical models. *arXiv preprint arXiv:1502.02251*, 2015.
- [55] Watter, M., Springenberg, J. T., Boedecker, J., and Riedmiller, M. Embed to control: A locally linear latent dynamics model for control from raw images. *arXiv preprint arXiv:1506.07365*, 2015.
- [56] Yarats, D., Zhang, A., Kostrikov, I., Amos, B., Pineau, J., and Fergus, R. Improving sample efficiency in model-free reinforcement learning from images. *arXiv preprint arXiv:1910.01741*, 2019.
- [57] Yarats, D., Fergus, R., Lazaric, A., and Pinto, L. Reinforcement learning with prototypical representations. *arXiv preprint arXiv:2102.11271*, 2021.
- [58] Yarats, D., Kostrikov, I., and Fergus, R. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=GY6-6sTvGaf>.
- [59] Zbontar, J., Jing, L., Misra, I., LeCun, Y., and Deny, S. Barlow twins: Self-supervised learning via redundancy reduction. *arXiv preprint arXiv:2103.03230*, 2021.
- [60] Zhang, A., Wu, Y., and Pineau, J. Natural environment benchmarks for reinforcement learning. *arXiv preprint arXiv:1811.06032*, 2018.
- [61] Zhang, A., Lyle, C., Sodhani, S., Filos, A., Kwiatkowska, M., Pineau, J., Gal, Y., and Precup, D. Invariant causal prediction for block MDPs.

In III, H. D. and Singh, A. (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 11214–11224. PMLR, 13–18 Jul 2020. URL <https://proceedings.mlr.press/v119/zhang20t.html>.

[62] Zhang, A., McAllister, R., Calandra, R., Gal, Y., and Levine, S. Learning Invariant Representations for Reinforcement Learning without Reconstruction. *arXiv*, jun 2020. ISSN 23318422. URL <http://arxiv.org/abs/2006.10742>.

[63] Zhang, A., Sodhani, S., Khetarpal, K., and Pineau, J. Learning Robust State Abstractions for Hidden-Parameter Block MDPs. pp. 1–22, 2020. URL <http://arxiv.org/abs/2007.07206>.