

Natural Environment Benchmarks for Reinforcement Learning

Amy Zhang

McGill University
Facebook AI Research
amy.x.zhang@mail.mcgill.ca

Yuxin Wu

Facebook AI Research
yuxinwu@fb.com

Joelle Pineau

McGill University
Facebook AI Research
jpineau@cs.mcgill.ca

Abstract

While current benchmark reinforcement learning (RL) tasks have been useful to drive progress in the field, they are in many ways poor substitutes for learning with real-world data. By testing increasingly complex RL algorithms on low-complexity simulation environments, we often end up with brittle RL policies that generalize poorly beyond the very specific domain. To combat this, we propose three new families of benchmark RL domains that contain some of the complexity of the natural world, while still supporting fast and extensive data acquisition. The proposed domains also permit a characterization of generalization through fair train/test separation, and easy comparison and replication of results. Through this work, we challenge the RL research community to develop more robust algorithms that meet high standards of evaluation.

Introduction

The field of Reinforcement Learning (RL) has exploded in recent years, with over 10K research papers published per year for the last six years (Henderson et al., 2018). The availability of challenging and widely-used benchmarking domains, such as the Atari Learning Environment (Bellemare et al., 2013) and MuJoCo (Todorov, Erez, and Tassa, 2012), has been a major driver of progress. By allowing the community to rally around a class of domains, these benchmarks enable fair and easy-to-share comparison of methods, which is useful to properly evaluate progress and ideas. The widespread use of benchmark datasets has had similar effect, in terms of driving progress, on several other subfields of AI (LeCun and Cortes, 2010; Krizhevsky, Nair, and Hinton, b; Deng et al., 2009). In other areas of science, from physics to biology, the use of simulators and models is also common practice.

More recently however, over-reliance on our current RL benchmarks has been called into question (Henderson et al., 2018). Results showing serious brittleness of methods suggest that either our algorithms are not sufficiently robust, or that our simulators are not sufficiently diverse to induce interesting learned behaviors. While there is a wealth of work on the former, very few research groups are paying attention to the latter, with the result that we devise increasingly rich algorithms, but continue to test them on synthetic domains of limited complexity which are a poor indicator of

real-world performance.

Most benchmarks and datasets used to evaluate machine learning algorithms (excluding RL) consist of data acquired from the real-world, including images, sound, human-written text. There are cases where synthetic data is considered in early phases of research, but most of the work is done on real-world data. In contrast, almost all of RL is done with low-complexity synthetic benchmarks. Of course some work uses robots and other physical systems, but the cost and complexity of data acquisition and platform sharing is prohibitive, and therefore such work can rarely be independently replicated.

The aim of this paper is to explore a new class of RL simulators that incorporate signals acquired from the natural (real) world as part of the state space. The use of natural signal is motivated by several observations. First, in comparison to just injecting random noise into the simulator, linking the state to a real-world signal ensures we have more meaningful task characteristics. Second, by sourcing a component of the state space from the real-world we can achieve fair train/test separation, which is a long-standing challenge for RL¹. Yet the tasks we propose remain fast and simple to use; in contrast to other work that might require a common robot infrastructure (Kober, Bagnell, and Peters, 2013) or animal model (Guez et al., 2008) or actual plant eco-system (Hall et al.), our set of tasks requires only a computer and Internet connection. The domains are easy to install, large quantities of data can be rapidly acquired, and the domains lend themselves to fair evaluations and comparisons.

In this paper we describe three families of natural environment RL domains, and we provide benchmark performance for several common RL algorithms on these domains. The three families of domains include two visual reasoning tasks, where an RL agent is trained to navigate inside natural images to classify images and localize objects, and a variant of the Atari Learning Environment that incorporates natural video in the background. In the process, we also uncover weaknesses of existing benchmarks that may not be well-recognized in the community. The primary goal of this work is to encourage the community to tackle RL domains beyond

¹In simulation domains, RL agents effectively typically train & test with the same simulator; if the simulator parameters are altered between training and evaluation then it is assumed to be an instance of transfer learning.

current short-description-length simulators, and to develop methods that are effective and robust in domains with natural conditions. Some of these new tasks also require RL to achieve higher-order cognition, for example combining the problems of image understanding and task solving.

Motivation

Consider one of the most widely used simulators for RL benchmarking: the Atari Learning Environment (Bellemare et al., 2013). In the words of the authors: *ALE is a simple object-oriented framework that allows researchers and hobbyists to develop AI agents for Atari 2600 games. It is built on top of the Atari 2600 emulator Stella.* The original Atari source code for some of these games is less than 100KB², the game state evolves in a fully deterministic manner, and there is no further injection of noise to add complexity. Even the core physics engine code for MuJoCo is around 1MB³, which simulates basic physical dynamics of the real world. Thus we argue that the inherent complexity of most ALE games and current physics engines, as defined by the description length of the domain, is trivially small.

Now compare this to a robot that has to operate in the real-world. The space of perceptual inputs depends on the robot’s sensors & their resolution. A standard Bumblebee stereo vision camera⁴ will generate over 10MB per second. Now consider that this robot is deployed in a world with zettabytes (=10²¹ bytes) of human-made information⁵, and where each human body may contain upwards of 150 zettabytes⁶. Clearly, RL algorithms have a long way to go before they can tackle the real-world in all its beautiful complexity.

While we strongly support the deployment and evaluation of RL algorithms in real-world domains, there are good reasons to explore protocols that allow replicable evaluation of RL algorithms in a fair and standardized way. This is the primary goal of this work. We aim to propose a set of benchmark RL domains that (a) contain some of the complexity of the natural world, (b) support fast and plentiful data acquisition, (c) allow fair train/test separation, and (d) enable easy replication and comparison.

Technical Setting

In reinforcement learning, an agent interacts with an environment modeled as a Markov Decision Process (MDP) (Bellman, 1957), which can be represented by a 6-tuple $(\mathcal{S}, \mathcal{A}, p_0(\mathcal{S}), T, R, \gamma)$, where:

- \mathcal{S} is the set of states,
- \mathcal{A} is the set of actions,

²<http://www.atariage.com/2600/programming/index.html>

³Supplied by Emo Todorov.

⁴<https://www.ptgrey.com/bumblebee2-firewire-stereo-vision-camera-systems>

⁵<https://blogs.cisco.com/sp/the-zettabyte-era-officially-begins-how-much-is-that>

⁶<https://bitesizebio.com/8378/how-much-information-is-stored-in-the-human-genome/>

- $p_0(\mathcal{S})$ is the initial state distribution,
- $T(S_{t+1}|S_t, A_t)$ is the probability of transitioning from state S_t to S_{t+1} , $S_t, S_{t+1} \in \mathcal{S}$ after action $A_t \in \mathcal{A}$,
- $R(r_{t+1}|S_t, A_t)$ is the probability of receiving reward $r_{t+1} \in \mathbb{R}$ after executing action A_t while in state S_t ,
- $\gamma \in [0, 1)$ is the discount factor.

Value-based methods aim to learn the value function of each state or state-action pair of the optimal policy π . We denote the state value function for a particular policy π as $V_\pi(s), \forall s \in \mathcal{S}$. The state-action value function is denoted $Q_\pi(s, a), \forall s, a \in (\mathcal{S}, \mathcal{A})$. In order to find the value functions corresponding to the optimal policy π^* , we have the update functions:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t)], \quad (1)$$

$$V(s_t) = \max_a Q(s_t, a), \quad (2)$$

which will converge to optimal $Q^*(S_t, A_t)$ and $V^*(S_t)$.

When learning an estimate $\hat{Q}(\cdot, \cdot | \omega)$ parameterized by ω of the optimal value function Q^* with temporal-difference methods we use the gradient update:

$$\omega_{t+1} \leftarrow \omega_t + \alpha[r_{t+1} + \gamma \max_a \hat{Q}(s_{t+1}, a; \omega) - \hat{Q}(s_t, a_t; \omega)] \nabla_\omega \hat{Q}(s_t, a_t; \omega).$$

The optimal policy is found by acting greedily over the optimal value function at each state

$$\pi^*(s) = \arg \max_a Q^*(s, a). \quad (3)$$

Learning the state-action value function with this bootstrapping method is called *Q-learning* (Watkins and Dayan, 1992). Value-based methods are *off-policy* in that they can be trained with samples not taken from the policy being learned.

Policy-based methods are methods that directly learn the policy as a parameterized function π_θ rather than learn the value function explicitly, where the parameters of the function are θ . Policy gradients use REINFORCE (Williams, 1992) with the update function

$$\theta_{t+1} \leftarrow \theta_t + \alpha G_t \frac{\nabla_\theta \pi(A_t | S_t, \theta_t)}{\pi(A_t | S_t, \theta_t)}, \quad (4)$$

where α is the step size, and $G_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots$ the return. A more general version of REINFORCE uses a baseline $b(S_t)$ to minimize the variance of the update:

$$\theta_{t+1} = \theta_t + \alpha (G_t - b(S_t)) \frac{\nabla_\theta \pi(A_t | S_t, \theta_t)}{\pi(A_t | S_t, \theta_t)}. \quad (5)$$

This baseline can be an estimate of the state value, learned separately in tabular form or as a parameterized function with weights ω . If the state value function is updated with bootstrapping like in value-based methods, then it is an actor-critic method.

Actor-Critic methods are hybrid value-based and policy-based methods that directly learn both the policy (actor) and

the value function (critic) (Konda and Tsitsiklis, 2000). The new update for actor-critic is:

$$\theta_{t+1} \leftarrow \theta_t + \alpha(G_t - \hat{V}(S_t)) \frac{\nabla_{\theta} \pi(A_t | S_t, \theta_t)}{\pi(A_t | S_t, \theta_t)} \quad (6)$$

$$= \theta_t + \alpha(R_{t+1} + \gamma \hat{V}(S_{t+1}) - \hat{V}(S_t)) \frac{\nabla_{\theta} \pi(A_t | S_t, \theta_t)}{\pi(A_t | S_t, \theta_t)}, \quad (7)$$

where $\hat{V}(\cdot)$ is a parameterized estimate of the optimal value function. The corresponding update for $\hat{V}(\cdot)$ is very similar to that in Q-learning (Watkins and Dayan, 1992):

$$\hat{V}(S_t) \leftarrow \hat{V}(S_t) + \alpha[r_{t+1} + \gamma \hat{V}(S_{t+1}) - \hat{V}(S_t)] \quad (8)$$

When learning an estimate $\hat{V}(\cdot; \omega)$ parameterized by ω of the optimal value function V^* with temporal-difference methods, we use the gradient update:

$$\omega_{t+1} \leftarrow \omega_t + \alpha[r_{t+1} + \gamma \hat{V}(s_{t+1}; \omega) - \hat{V}(s_t; \omega)] \nabla_{\omega} \hat{V}(s_t; \omega). \quad (9)$$

Popular RL Algorithms

Advantage Actor Critic (A2C). Mnih et al. (2016) propose an on-policy method based on actor-critic with several parallel actors which replaces the value estimate with the *advantage* $A_{\pi}(a, s) = Q_{\pi}(a, s) - V_{\pi}(s)$.

Actor Critic using Kronecker-Factored Trust Region (ACKTR). Wu et al. (2017) uses *trust region optimization* with a *Kronecker-factored approximation (K-FAC)* (Martens and Grosse, 2015) with actor-critic. *Trust region optimization* (Schulman et al., 2015) is an approach where the update is clamped at a maximum learning rate η_{\max} . *K-FAC* is an invertible approximation of the Fisher information matrix of the neural network representing the policy by block partitioning the matrix according to the layers of the neural network, then approximating these blocks as Kronecker products of smaller matrices. Martens and Grosse (2015) show that this approximation is efficient to invert and preserves gradient information. ACKTR is a constrained optimization problem with a constraint that the policy does not move too far in the update, measured with KL-divergence. It also computes steps using the natural gradient direction as opposed to gradient direction. However, computing the exact second derivative is expensive, so Wu et al. (2017) instead use K-FAC as an approximation.

Proximal Policy Optimization (PPO). Schulman et al. (2017) propose a family of policy gradient methods that also use trust region optimization to clip the size of the gradient and multiple epochs of stochastic gradient ascent for each policy update. PPO uses a penalty to constrain the update to be close to the previous policy.

Deep Q-Network (DQN) Mnih et al. (2013) modify Q-learning to use a deep neural network with parameters ω_t to model the state-action value function. The authors introduce a few tricks to stabilize training, mainly using a separate network Q' to compute the target values, which is implemented as an identical neural network but with different parameters ω'_t copied over from ω_t at fixed intervals. The second trick

is *experience replay*, or keeping a buffer of prior experience for batch training. The new gradient update to ω_t is:

$$\omega_{t+1} \leftarrow \omega_t + \alpha[r_{t+1} + \gamma \max_a Q(s_{t+1}, a; \omega') - Q(s_t, a_t; \omega)] \nabla_{\omega} Q(s_t, a_t; \omega) \quad (10)$$

Related Work

Simulation Environments and Benchmarks for RL

There has been many recent proposed simulation engines that try to bridge the gap between simulation and reality by creating more and more realistic but still rendered pixel-level observation spaces (Brodeur et al., 2017; Kolve et al., 2017).

The current set of benchmark tasks for RL such as the Atari Learning Environment (Bellemare et al., 2013) and OpenAI gym (Plappert et al., 2018) are primarily composed of deterministic settings. A larger issue is that even in the tasks with larger state spaces, such as the pixel state in Atari, the description of the rules can be modeled in a small instruction set (lines of code or natural language rules). The real world is not deterministic, in part because the stochasticity comes from unobserved variables, but is also not directly model-able in a few lines of rules or code.

Duan et al. (2016) released a set of benchmark tasks for continuous RL, pointing out that existing algorithms that work well on discrete tasks (Bellemare et al., 2013) wouldn't necessarily transfer well to continuous, high dimensional action spaces in the control domain.

Benchmarks are necessary to evaluate various proposed algorithms and compare them against each other. However, the current suite of available tasks conflate the difficulty of visual comprehension with that of finding an optimal policy, and are a black box for determining how algorithms are actually solving the task. Our results show that visual comprehension is still a difficult task even though we can achieve record scores in Atari on pixel observation. We must take a step back and focus on tasks that can partition the various dimensions along which RL tasks are difficult.

Henderson et al. (2018) point out issues of reproducibility in deep RL, and we also find that implementations on top of different frameworks (Abadi et al., 2016; Paszke et al., 2017) as built by Dhariwal et al. (2017) and Kostrikov (2018) have very different results.

Rajeswaran et al. (2017) show that the recent improvements in deep RL are not necessarily due to deep neural networks, and that similar improvements can be seen with linear models. They also propose widening the initial state distribution to generate more robust policies. We take this a step further by proposing to widen the state space of the MDP through the introduction of natural signal.

RL for Classical Computer Vision Tasks

There has been much recent work bridging RL and computer vision (CV), including object detection (Caicedo and Lazebnik, 2015), object tracking (Yun et al., 2017; Zhang et al., 2017), and object classification (Zhang, Ballas, and Pineau, 2018). They show it is possible to use RL techniques to perform object localization, and that using RL to localize is a

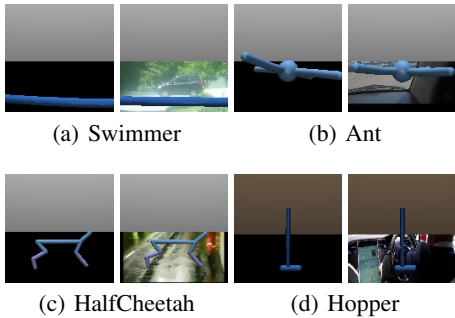


Figure 1: Mujoco frames with original black ground plane (left) and natural video embedded as background in replacement of the ground plane (right).

promising direction of research (LeCun, Bengio, and Hinton, 2015). These works show that RL has been successfully applied to visual comprehension tasks, but often using many domain-specific tricks that are not carried over in RL applications. Our work is the first to evaluate on the visual comprehension tasks with state-of-the-art algorithms designed for RL.

There has also been some work applying CV techniques to solve RL problems in robotics (Rusu et al., 2016), games (Mnih et al., 2015), and navigation of maps Mirowski et al. (2018) via pixel-level observation spaces. These typically consist of applying CNNs to RL tasks to process the low-level pixel state, but only with medium-sized convolutional neural networks or fully-connected networks composed of 2-3 layers.

New Benchmark RL Tasks with Natural Signal

We aim to develop RL benchmarks that capture more of the complexity of the real world, without prohibitive resource and time costs. We consider three families of tasks, the first two are based on visual reasoning tasks and the third is a variant of existing RL benchmarks.

Visual Reasoning using RL

The first set of proposed tasks consist of gridworld environments overlaid on a natural image. These environments show how we can transform traditionally supervised learning tasks to basic RL navigation tasks with natural signal that requires visual comprehension. We illustrate this with a few examples (MNIST, CIFAR10 and CIFAR100 for classification; Cityscapes for localization), but the main win here is that we can leverage any existing image dataset. Each of these datasets has a pre-defined train/test split which we respect (train RL agents on training set images; evaluate on test set images) to extract fair generalization measures. These new domains contain several real-world images with natural complexity, and are easily downloadable for easy replication, thus meeting the desiderata outlined in our motivation above.

Agent navigation for image classification. We propose an image classification task starting with a masked image where the agent starts at a random location on the image. It



Figure 2: Atari frames, original (left), Gaussian noise (center), and with natural video embedded as background (right).

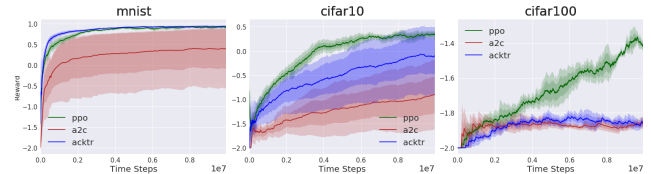


Figure 3: Agent navigation for image classification results. Variance computed across 5 seeds. Note the difference in scale on y-axis.

can unmask windows of the image by moving in one of 4 directions: {UP, DOWN, LEFT, RIGHT}. At each timestep it also outputs a probability distribution over possible classes \mathcal{C} . The episode ends when the agent correctly classifies the image or a maximum of 20 steps is reached. The agent receives a -0.1 reward at each timestep that it misclassifies the image. The state received at each time step is the full image with unobserved parts masked out.

We evaluate on MNIST (LeCun and Cortes, 2010), CIFAR10 (Krizhevsky, Nair, and Hinton, a), and CIFAR100 (Krizhevsky, Nair, and Hinton, b), all of which consist of 60k images, 28x28 and 32x32 respectively, and with 10 classes apiece. MNIST is grayscale (single channel), and CIFAR10 and CIFAR100 are 3 channel RGB. To scale the difficulty of the problem, we can change the window size w of the agent and maximum number of steps per episode M .

Agent navigation for object localization. Given the segmentation mask of an object in an image, the agent has to move to sit on top of the object. There are again 4 possible actions at each timestep, with a time limit of 200 steps. We can further complicate the task with several objects and an additional input of which object class the goal is.

We use the Cityscapes (Cordts et al., 2016) dataset for object detection with a window size $w = 10$ which controls the difficulty of the task. The Cityscapes dataset consists of 50k 256x256 images and 30 classes. The window size dictates the footprint of the agent. The episode ends if the footprint overlaps with the desired object in the image. The agent is dropped in the center of the image and is given a class label representing a goal object to find and navigate to. The episode ends when the agent is on top of the desired object, for which the environment gives a reward of 1, or if the maximum of 200 steps is reached. There is no penalty for each step in this task – reward is 0 at each step the agent is not on the desired object.

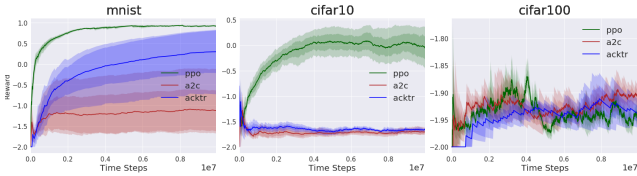


Figure 4: Agent navigation for image classification results with ResNet-18 trunk. Variance computed across 5 seeds.

Natural Video RL Benchmarks

We also propose a modification to existing RL benchmark tasks to incorporate natural signal. Effectively, we take Atari (Bellemare et al., 2013) tasks from OpenAI gym (Plappert et al., 2018) and add natural videos as the background of the observed frames.

We used videos of driving cars from the Kinetics dataset (Kay et al., 2017) and created a mask of the Atari frames by filtering for black pixels (0, 0, 0), substituting the video frame for the black background. To maintain optical flow we used consecutive frames from randomly chosen videos for the background and randomly sampled from the same set of 840 videos for train and test.

We do the same for MuJoCo tasks in OpenAI gym (Plappert et al., 2018). The default MuJoCo uses a low-dimensional state space consisting of position and velocity of each joint. Instead, we consider PixelMuJoCo, where the observation space consists of a camera tracking the agent. Lillicrap et al. (2015) also use a pixel version of MuJoCo and demonstrate similar performance to the low-dimensional version. In our new benchmark, we substitute the floor plane of the PixelMuJoCo tasks (Ant, Swimmer, Hopper, and HalfCheetah) with the same video frames as in the Atari domains. We have included results for PixelMuJoCo but do not include them in our proposed set of benchmarks because we have discovered policies learned for MuJoCo are *open-loop*, and completely ignore the observation input.

After applying these changes, the state space for these environments drastically increases, and the problem becomes one of visually comprehending the scene in order to attend to the objects corresponding to the game, and ignoring the objects in the video. Example frames for Atari and PixelMuJoCo with natural signal can be seen in Figures 2, 1.

Results

In this section we provide benchmark performance of existing popular RL algorithms on the new proposed domains.

Visual Reasoning

For the proposed visual reasoning tasks, we run both a small convolutional neural network (CNN) as commonly used for pixel RL tasks and Resnet-18 (He et al., 2015) on MNIST (LeCun and Cortes, 2010), CIFAR10 (Krizhevsky, Nair, and Hinton, a), and Cityscapes (Cordts et al., 2016).

The CNN consists of 3 convolutional layers and a fully connected layer, of varying filter sizes and strides to contend

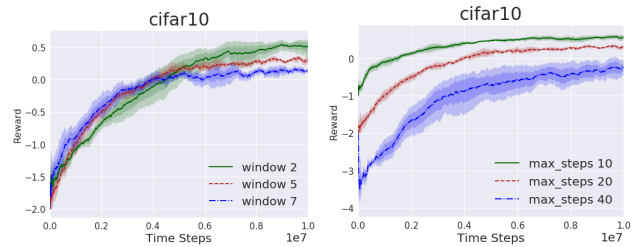


Figure 5: CIFAR10 with PPO, varying window size (left) with fixed maximum number of steps $M = 20$ and maximum number of steps per episode (right) with fixed window size $w = 5$.

with the images from different datasets being different sizes. These layers are interpolated with ReLUs. More detail about model architecture can be found in the Appendix.

Agent navigation for image classification. Results for the image classification task on MNIST, CIFAR10, and CIFAR100 are found in Figures 3, 4 for the 3-layer CNN and ResNet-18, respectively. We see that PPO and ACKTR are able to achieve average reward close to 1 on MNIST, which means the agent is able to accurately classify the digit without needing to take many steps. Performance is worse on CIFAR10 and CIFAR100, as expected, because the datasets consist of more difficult visual concepts. We see the same performance drop across datasets in supervised learning (He et al., 2015). A2C consistently performs worst across all datasets and trunk models.

More interestingly, we also see performance drop when moving from a 3-layer CNN to ResNet-18 across all 3 datasets. PPO is still able to achieve the same performance on MNIST and CIFAR10, which both have 10 classes, but ACKTR and A2C suffer dramatically. None of the methods work well with ResNet-18 and across 100 classes. This conflates two more difficult problems – the action space is now $10\times$ larger and there are $10\times$ more concepts to learn.

We can alter the difficulty of this task along two dimensions – varying the window size w of the agent, or the maximum number of steps per episode M . In experiments, we try values of $w \in [2, 5, 7]$ and $M \in [10, 20, 40]$. Results are in Figure 5. Performance increases with fewer number of steps, which corresponds to more immediate rewards and therefore an easier RL task. Initially, a larger window size performs better, as expected, but as training continues the smaller window $w = 2$ dominates.

We see that the large models that have been successful in SL for object classification suffer in RL tasks. However, accurate object classification and localization are necessary components of how humans solve many tasks, and how we expect RL algorithms to learn to solve tasks. Here we show by isolating the visual comprehension component of common RL tasks how well current state-of-the-art RL algorithms actually perform, and that simple plug and play of successful SL vision models does not give us the same gains when applied in an RL framework. We thus expect this new family of tasks to spur new innovation in RL algorithms.

Agent navigation for object localization. Results for the

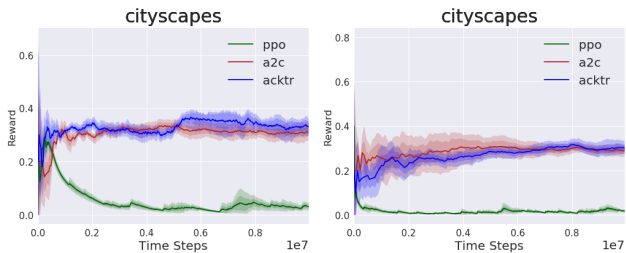


Figure 6: Agent navigation for image detection results. 3-layer CNN (left), ResNet-18 (right).

object detection task on Cityscapes is found in Figure 6. Object detection is a much more difficult task with again a small drop in performance when moving from the 3-layer CNN trunk to ResNet-18.

Here we see that PPO completely fails to learn, whereas it beat out both A2C and ACKTR in the classification task. But both A2C and ACKTR are not able to navigate to the desired object in the image more than 40% of the time.

Both of these vision tasks demonstrate what others have also found (Irpan, 2018; Rajeswaran et al., 2017) – that deep models do not perform well in the RL framework in the same way they do in the SL framework. Clearly, this opens up many interesting directions for future research.

Natural Signal in RL Tasks

For both Atari (Bellemare et al., 2013) and PixelMuJoCo (Plappert et al., 2018) we follow the preprocessing done in Mnih et al. (2015) and resize the frames to 84×84 , convert to grayscale, and perform frame skipping and sticky actions for 4 consecutive steps. We also stack four consecutive frames for each observation. For algorithm implementations we use OpenAI Baselines (Dhariwal et al., 2017) and Ilya Kostrikov’s implementation (Kostrikov, 2018).

As baseline, we compare with the original Atari and PixelMuJoCo tasks with static black background⁷.

PixelMuJoCo. For PixelMuJoCo, we evaluate on Hopper, Swimmer, Ant, and HalfCheetah. There are results reported by Lillicrap et al. (2015) for PixelMuJoCo with DDPG, but the rewards are not directly comparable, and the pixel environment they use is different from the renderer provided in OpenAI gym based on the visualizations in the paper, and has not been open sourced. PixelMuJoCo results are in Figure 7. We see similar performance across baseline and natural, with small performance gaps apparently especially in HalfCheetah and Hopper. We suspect this is actually caused by the policy falling into a local optima where it ignores the observed state entirely – the fact that Lillicrap et al. (2015) also report similar results for the low-dimensional state space and pixel space for MuJoCo tasks also points to this conclusion.

To test our hypothesis, we replace the observation with Gaussian noise. Results are shown in Figure 7 as Pure Noise. Even in the case where there is *no information* returned in the observation—it is pure iid Gaussian noise—PPO and

⁷Refer to Fig. 2 for visualizations of what the original and modified Atari games look like.

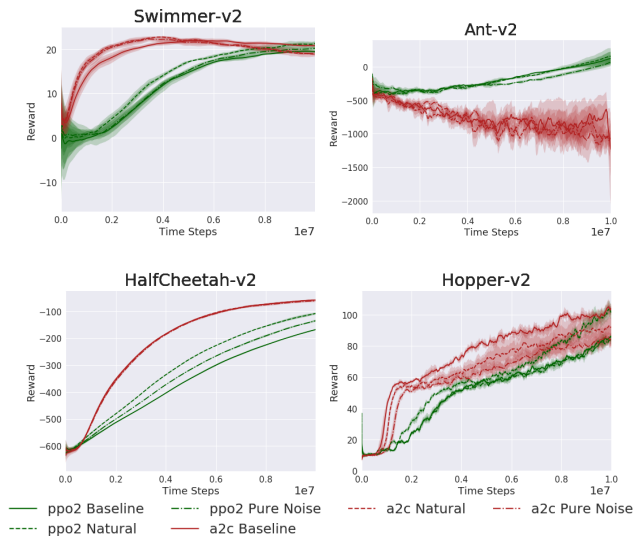


Figure 7: Natural signal in PixelMuJoCo results, using Dhariwal et al. (2017) code implementation. Variance computed across 5 seeds.

A2C are able to learn as good a policy as when actual state information is provided in the observation. Our results show that current RL algorithms are solving MuJoCo tasks as an *open-loop control system*, where it completely ignores the output when deciding the next action. These results suggest that MuJoCo is perhaps not a strong benchmark for RL algorithms, and a good test for open-control policies is substituting the observation with pure noise.

Atari. For Atari, we selected 16 environments (mainly ones with black backgrounds for ease of filtering) and evaluated PPO, ACKTR, A2C, and DQN on both the default environment and with injected video frames. Full results can be seen in Fig. 11 in the Appendix, we have only included 4 environments in the main paper (Fig. 8) because of space constraints. We see much larger gaps in performance for many games, showing that visual comprehension is more crucial in these environments. The addition of natural noise with visual flow causes the policy to completely fail in some games, while causing only a small drop in performance in others. In these cases it is again possible that the policy is treating the task as an open-loop control problem.

To see how much the performance difference is caused by the addition of *natural signal* as opposed to just changing these environments to no longer be deterministic, we also evaluate on a few Atari environments where the background is replaced with random i.i.d. Gaussian noise changing from frame to frame (Figure 8). We see with all four games that we have best performance with the original static black background (Baseline), similar to reduced performance with Gaussian noise (Noise), and worst performance with video (Natural). However, the performance difference varies. In Amidar, performance with random noise is very close to baseline for ACKTR and PPO, but suffers a massive drop with natural signal. In Beamrider, the algorithms all fail to obtain good policies with any addition of random noise or natural signal.

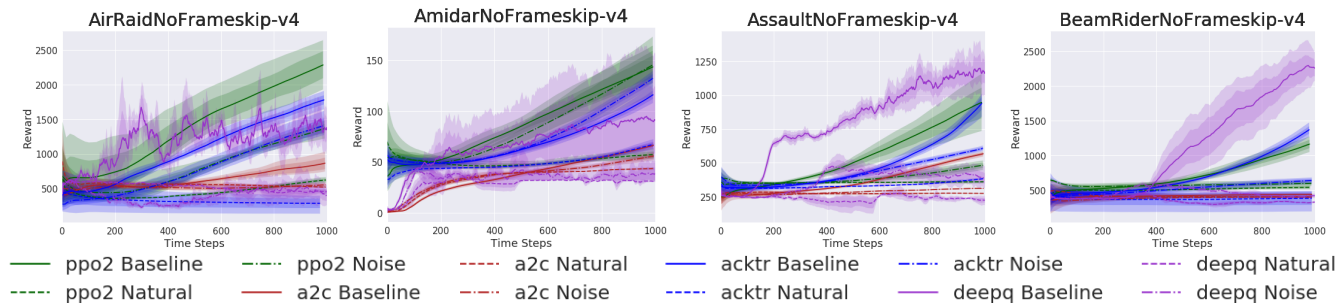


Figure 8: Atari frames with baseline, Gaussian noise, and natural video embedded as background, using Dhariwal et al. (2017) code implementation. Variance computed across 5 seeds.

We see that the Atari tasks are complex enough, or require different enough behavior in varying states that the policy cannot just ignore the observation state, and instead learns to parse the observation to try to obtain a good policy. In most games it is able to do this with static background, suffers with random noise background, and fails completely with natural signal background.

Discussion

We have proposed three new families of benchmark tasks to dissect performance of RL algorithms. The first two are domains that test visual comprehension by bringing traditional supervised learning tasks into the RL framework. In the process, we have shown that naive plug and play of successful vision models fails in the RL setting. This suggests that the end-to-end frameworks espoused for RL currently are not successful at implicitly learning visual comprehension.

The third family of tasks call for evaluating RL algorithms via incorporating signal from the natural world, by injecting frames from natural video into current RL benchmarks. We have shown that performance deteriorates drastically in this setting across several state-of-the-art RL optimization algorithms and trunk models. With this new set of tasks, we call for new algorithm development to be more robust to natural noise. We also observe that state-of-the-art performance on the PixelMuJoCo domain can be achieved with an open-loop policy, making it an odd choice for an RL benchmark. Based on these results, we have also proposed replacing the observation with pure noise as a test for open-loop policies.

As a side note, we note that we were able to achieve the same results reported by Dhariwal et al. (2017) for some of the games but not all with the default hyperparameters provided and their code, and also saw large differences in performance when comparing results using Kostrikov (2018) vs. Dhariwal et al. (2017) implementations. Results across 16 Atari games are shown in Figures 11 and 12 in the Appendix.

The first set of tasks in object recognition and localization are tests of how well models can learn visual comprehension in an RL setting. Only once we have achieved good results in those tasks can we move onto tasks with more difficult dynamics and be sure that the performance is caused by visual comprehension as opposed to memorizing trajectories.

Beyond piping natural signal into the state space through

the observations, another type of noise in the real world is noise in the action effects and dynamics. Transitions from one state to the next exhibit noise from imperfect actuators and sensors. It is still an open question how we can inject natural dynamics signal into simulated environments.

References

- Abadi, M.; Barham, P.; Chen, J.; Chen, Z.; Davis, A.; Dean, J.; Devin, M.; Ghemawat, S.; Irving, G.; Isard, M.; Kudlur, M.; Levenberg, J.; Monga, R.; Moore, S.; Murray, D. G.; Steiner, B.; Tucker, P.; Vasudevan, V.; Warden, P.; Wicke, M.; Yu, Y.; and Zheng, X. 2016. Tensorflow: A system for large-scale machine learning. In *12th USENIX Symposium on Operating Systems Design and Implementation (OSDI 16)*, 265–283.
- Bellemare, M. G.; Naddaf, Y.; Veness, J.; and Bowling, M. 2013. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research* 47:253–279.
- Bellman, R. 1957. A markovian decision process. *Journal of Mathematics and Mechanics* 6(5):679–684.
- Brodeur, S.; Perez, E.; Anand, A.; Golemo, F.; Celotti, L.; Strub, F.; Rouat, J.; Larochelle, H.; and Courville, A. C. 2017. Home: a household multimodal environment. *CoRR* abs/1711.11017.
- Caicedo, J. C., and Lazebnik, S. 2015. Active object localization with deep reinforcement learning. *CoRR* abs/1511.06015.
- Cordts, M.; Omran, M.; Ramos, S.; Rehfeld, T.; Enzweiler, M.; Benenson, R.; Franke, U.; Roth, S.; and Schiele, B. 2016. The cityscapes dataset for semantic urban scene understanding. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Deng, J.; Dong, W.; Socher, R.; jia Li, L.; Li, K.; and Feifei, L. 2009. Imagenet: A large-scale hierarchical image database. In *In CVPR*.
- Dhariwal, P.; Hesse, C.; Klimov, O.; Nichol, A.; Plappert, M.; Radford, A.; Schulman, J.; Sidor, S.; and Wu, Y. 2017. Openai baselines. <https://github.com/openai/baselines>.

- Duan, Y.; Chen, X.; Houthoofd, R.; Schulman, J.; and Abbeel, P. 2016. Benchmarking deep reinforcement learning for continuous control. *CoRR* abs/1604.06778.
- Guez, A.; Vincent, R.; Avoli, M.; and Pineau, J. 2008. Adaptive treatment of epilepsy via batch-mode reinforcement learning. In *Innovative Applications of Artificial Intelligence*.
- Hall, K.; Albers, H.; Taleghan, M.; and Dietterich, T. Optimal spatial-dynamic management of stochastic species invasions.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2015. Deep residual learning for image recognition. *CoRR* abs/1512.03385.
- Henderson, P.; Islam, R.; Bachman, P.; Pineau, J.; Precup, D.; and Meger, D. 2018. Deep reinforcement learning that matters. In *AAAI*. AAAI Press.
- Irpan, A. 2018. Deep reinforcement learning doesn't work yet. <https://www.alexirpan.com/2018/02/14/rl-hard.html>.
- Kay, W.; Carreira, J.; Simonyan, K.; Zhang, B.; Hillier, C.; Vijayanarasimhan, S.; Viola, F.; Green, T.; Back, T.; Natsev, P.; Suleyman, M.; and Zisserman, A. 2017. The kinetics human action video dataset. *CoRR* abs/1705.06950.
- Kober, J.; Bagnell, J. A.; and Peters, J. 2013. Reinforcement learning in robotics: A survey. *International Journal of Robotics REsearch*.
- Kolve, E.; Mottaghi, R.; Gordon, D.; Zhu, Y.; Gupta, A.; and Farhadi, A. 2017. AI2-THOR: an interactive 3d environment for visual AI. *CoRR* abs/1712.05474.
- Konda, V. R., and Tsitsiklis, J. N. 2000. Actor-critic algorithms. In Solla, S. A.; Leen, T. K.; and Müller, K., eds., *Advances in Neural Information Processing Systems 12*. MIT Press. 1008–1014.
- Kostrikov, I. 2018. Pytorch implementations of reinforcement learning algorithms. <https://github.com/ikostrikov/pytorch-a2c-ppo-acktr>.
- Krizhevsky, A.; Nair, V.; and Hinton, G. Cifar-10 (canadian institute for advanced research).
- Krizhevsky, A.; Nair, V.; and Hinton, G. Cifar-100 (canadian institute for advanced research).
- LeCun, Y., and Cortes, C. 2010. MNIST handwritten digit database.
- LeCun, Y.; Bengio, Y.; and Hinton, G. E. 2015. Deep learning. *Nature* 521(7553):436–444.
- Lillicrap, T. P.; Hunt, J. J.; Pritzel, A.; Heess, N.; Erez, T.; Tassa, Y.; Silver, D.; and Wierstra, D. 2015. Continuous control with deep reinforcement learning. *CoRR* abs/1509.02971.
- Martens, J., and Grosse, R. B. 2015. Optimizing neural networks with kronecker-factored approximate curvature. *CoRR* abs/1503.05671.
- Mirowski, P.; Grimes, M. K.; Malinowski, M.; Hermann, K. M.; Anderson, K.; Teplyashin, D.; Simonyan, K.; Kavukcuoglu, K.; Zisserman, A.; and Hadsell, R. 2018. Learning to navigate in cities without a map. *CoRR* abs/1804.00168.
- Mnih, V.; Kavukcuoglu, K.; Silver, D.; Graves, A.; Antonoglou, I.; Wierstra, D.; and Riedmiller, M. 2013. Playing atari with deep reinforcement learning. [cite arxiv:1312.5602](https://arxiv.org/abs/1312.5602)Comment: NIPS Deep Learning Workshop 2013.
- Mnih, V.; Kavukcuoglu, K.; Silver, D.; Rusu, A. A.; Veness, J.; Bellemare, M. G.; Graves, A.; Riedmiller, M.; Fidjeland, A. K.; Ostrovski, G.; Petersen, S.; Beattie, C.; Sadik, A.; Antonoglou, I.; King, H.; Kumaran, D.; Wierstra, D.; Legg, S.; and Hassabis, D. 2015. Human-level control through deep reinforcement learning. *Nature* 518(7540):529–533.
- Mnih, V.; Puigdomènech Badia, A.; Mirza, M.; Graves, A.; Lillicrap, T. P.; Harley, T.; Silver, D.; and Kavukcuoglu, K. 2016. Asynchronous Methods for Deep Reinforcement Learning. *ArXiv e-prints*.
- Paszke, A.; Gross, S.; Chintala, S.; Chanan, G.; Yang, E.; DeVito, Z.; Lin, Z.; Desmaison, A.; Antiga, L.; and Lerer, A. 2017. Automatic differentiation in pytorch.
- Plappert, M.; Andrychowicz, M.; Ray, A.; McGrew, B.; Baker, B.; Powell, G.; Schneider, J.; Tobin, J.; Chociej, M.; Welinder, P.; Kumar, V.; and Zaremba, W. 2018. Multi-goal reinforcement learning: Challenging robotics environments and request for research.
- Rajeswaran, A.; Lowrey, K.; Todorov, E. V.; and Kakade, S. M. 2017. Towards generalization and simplicity in continuous control. In Guyon, I.; Luxburg, U. V.; Bengio, S.; Wallach, H.; Fergus, R.; Vishwanathan, S.; and Garnett, R., eds., *Advances in Neural Information Processing Systems 30*. Curran Associates, Inc. 6550–6561.
- Rusu, A. A.; Vecerik, M.; Rothörl, T.; Heess, N.; Pascanu, R.; and Hadsell, R. 2016. Sim-to-real robot learning from pixels with progressive nets. *CoRR* abs/1610.04286.
- Schulman, J.; Levine, S.; Moritz, P.; Jordan, M. I.; and Abbeel, P. 2015. Trust region policy optimization. *CoRR* abs/1502.05477.
- Schulman, J.; Wolski, F.; Dhariwal, P.; Radford, A.; and Klimov, O. 2017. Proximal Policy Optimization Algorithms. *ArXiv e-prints*.
- Todorov, E.; Erez, T.; and Tassa, Y. 2012. Mujoco: A physics engine for model-based control. In *IROS*, 5026–5033. IEEE.
- Watkins, C. J. C. H., and Dayan, P. 1992. Q-learning. In *Machine Learning*, 279–292.
- Williams, R. J. 1992. Simple statistical gradient-following algorithms for connectionist reinforcement learning. In *Machine Learning*, 229–256.
- Wu, Y.; Mansimov, E.; Liao, S.; Grosse, R. B.; and Ba, J. 2017. Scalable trust-region method for deep reinforcement learning using kronecker-factored approximation. *CoRR* abs/1708.05144.

Yun, S.; Choi, J.; Yoo, Y.; Yun, K.; and Choi, J. Y. 2017. Action-decision networks for visual tracking with deep reinforcement learning. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 1349–1358.

Zhang, A.; Ballas, N.; and Pineau, J. 2018. A dissection of overfitting and generalization in continuous reinforcement learning. *CoRR* abs/1806.07937.

Zhang, D.; Maei, H.; Wang, X.; and Wang, Y. 2017. Deep reinforcement learning for visual object tracking in videos. *CoRR* abs/1701.08936.

Implementation

Model Architectures

```
if dataset == 'mnist':
    self.main = nn.Sequential(
        nn.Conv2d(num_inputs, 10, 5, stride=2),
        nn.ReLU(),
        nn.Conv2d(10, 20, 5),
        nn.ReLU(),
        nn.Conv2d(20, 10, 5),
        nn.ReLU(),
        Flatten(),
        nn.Linear(360, 512),
        nn.ReLU()
    )
elif dataset == 'cifar10':
    self.main = nn.Sequential(
        nn.Conv2d(num_inputs, 6, 5, stride=2),
        nn.ReLU(),
        nn.Conv2d(6, 16, 5),
        nn.ReLU(),
        nn.Conv2d(16, 6, 5),
        nn.ReLU(),
        Flatten(),
        nn.Linear(216, 512),
        nn.ReLU()
    )
elif dataset == 'cityscapes':
    self.main = nn.Sequential(
        nn.Conv2d(num_inputs, 32, 8, stride=7),
        nn.ReLU(),
        nn.Conv2d(32, 64, 4, stride=4),
        nn.ReLU(),
        nn.Conv2d(64, 32, 3, stride=1),
        nn.ReLU(),
        Flatten(),
        nn.Linear(32 * 7 * 7, 512),
        nn.ReLU()
    )
else: # MuJoCo and Atari
    self.main = nn.Sequential(
        nn.Conv2d(num_inputs, 32, 8, stride=4),
        nn.ReLU(),
        nn.Conv2d(32, 64, 4, stride=2),
        nn.ReLU(),
        nn.Conv2d(64, 32, 3, stride=1),
        nn.ReLU(),
        Flatten(),
        nn.Linear(32 * 7 * 7, 512),
        nn.ReLU()
    )
```

Hyperparameters

Hyperparameters were kept constant for all MuJoCo and Atari tasks.

Sample Environment Frames

More Results

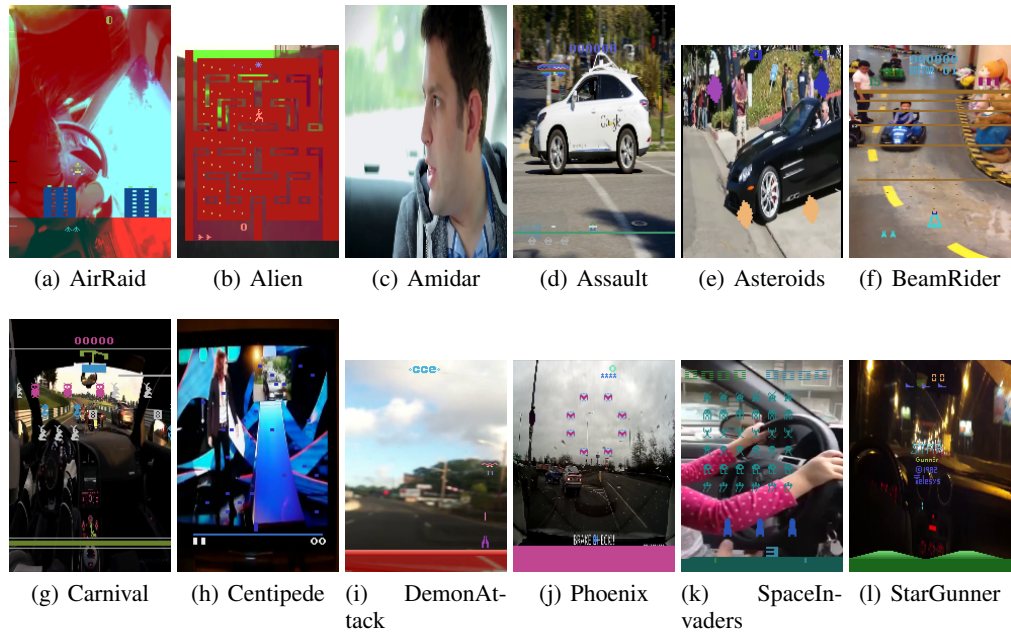


Figure 9: More Atari frames with natural video embedded as background.

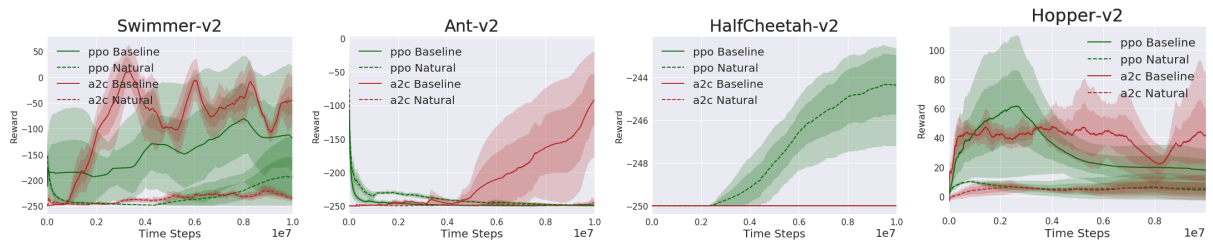


Figure 10: Natural signal in PixelMuJoCo results, using Kostrikov (2018) code implementation. Variance computed across 5 seeds.

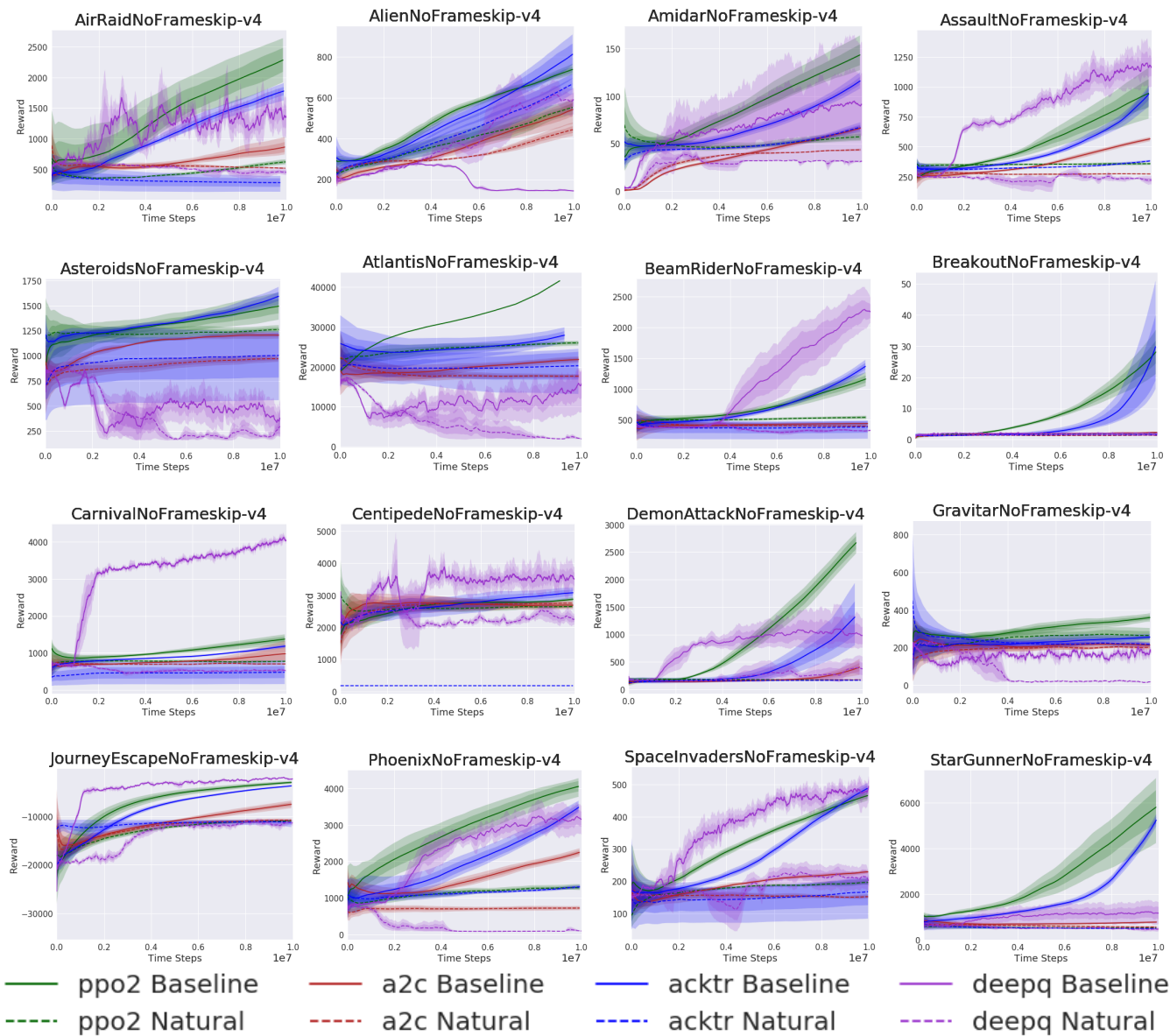


Figure 11: Atari frames with natural video embedded as background, using Dhariwal et al. (2017) code implementation. Variance computed across 5 seeds.

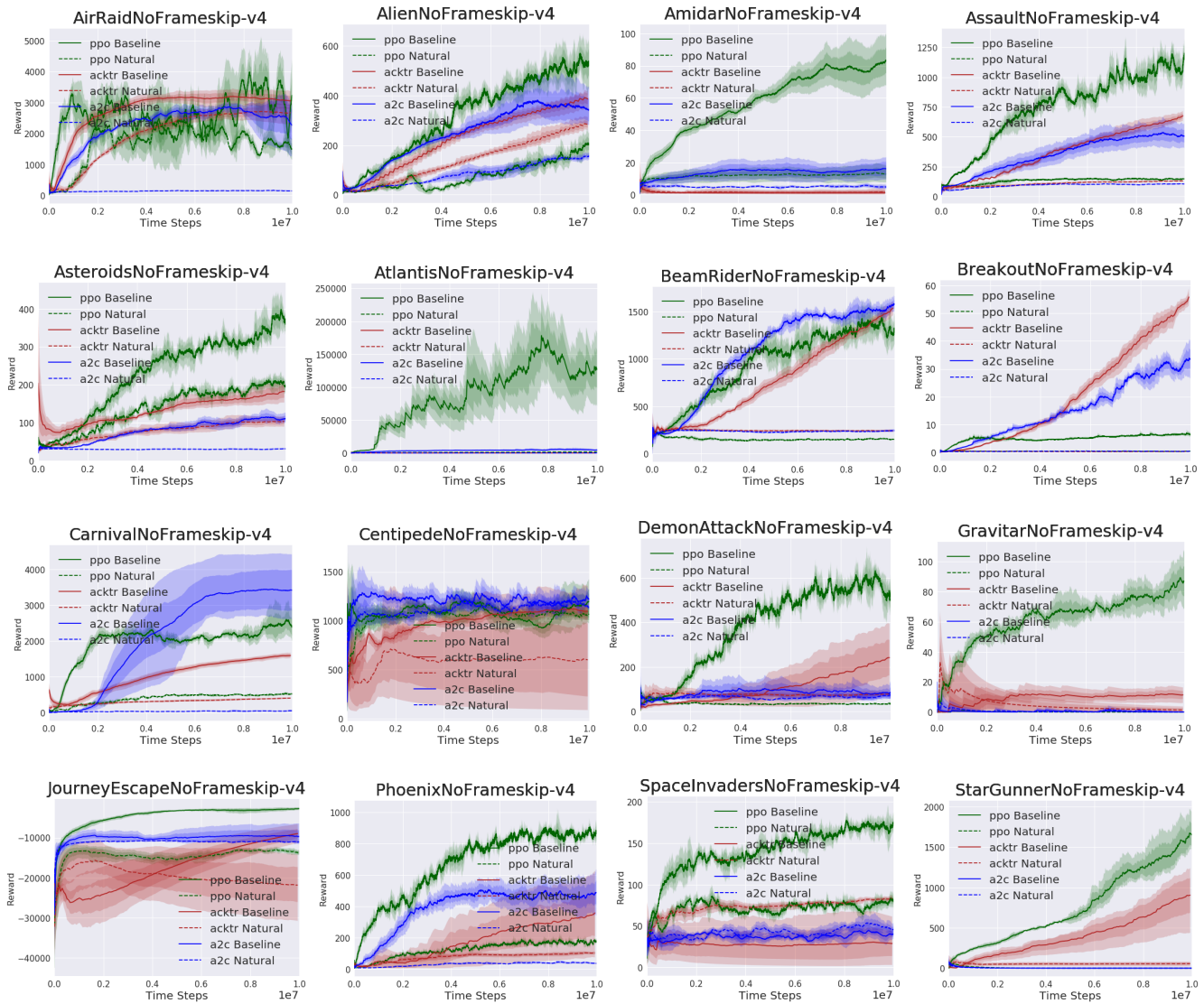


Figure 12: Atari frames with natural video embedded as background, using Kostrikov (2018) code implementation. Variance computed across 5 seeds.