

Asymmetric Actor Critic for Image-Based Robot Learning

Lerrel Pinto^{1,2} Marcin Andrychowicz¹ Peter Welinder¹ Wojciech Zaremba^{1,†} Pieter Abbeel^{1,†}

Abstract—Deep reinforcement learning (RL) has proven a powerful technique in many sequential decision making domains. However, Robotics poses many challenges for RL, most notably training on a physical system can be expensive and dangerous, which has sparked significant interest in learning control policies using a physics simulator. While several recent works have shown promising results in transferring policies trained in simulation to the real world, they often do not fully utilize the advantage of working with a simulator. In this work, we exploit the full state observability in the simulator to train better policies which take as input only partial observations (RGBD images). We do this by employing an *actor-critic* training algorithm in which the critic is trained on full states while the actor (or policy) gets rendered images as input. We show experimentally on a range of simulated tasks that using these asymmetric inputs significantly improves performance. Finally, we combine this method with domain randomization and show real robot experiments for several tasks like picking, pushing, and moving a block. We achieve this simulation to real world transfer without training on any real world data. Videos of these experiments can be found at www.youtube.com/watch?v=b57WTs.

I. INTRODUCTION

Reinforcement learning (RL) coupled with deep neural networks has recently led to successes on a wide range of control problems, including achieving superhuman performance on Atari games [1] and beating the world champion in the classic game of Go [2]. In physics simulators, complex behaviours like walking, running, hopping and jumping have also been shown to emerge [3], [4].

In the context of robotics however, learning complex behaviours faces two unique challenges: scalability and safety. Robots are slow and expensive which makes existing data intensive learning algorithms hard to scale. These physical robots could also damage themselves and their environment while exploring these behaviours. A recent approach to circumvent these challenges is to train on a simulated version of the robot and then transfer to the real robot [5]–[12]. However, this brings about a new challenge: observability.

Simulators have access to the full state of the robot and its surroundings, while in the real world obtaining this full state observability is often infeasible. One option is to infer the full state by visual detectors [13], [14] or state prediction filters [15]. Explicit full state prediction from partial observations is often impossible and this challenge is further exacerbated by the compounding error problem [16]. Another option is to train entirely on rendered partial observations (camera images) of the robot [11], [17]. However, these techniques

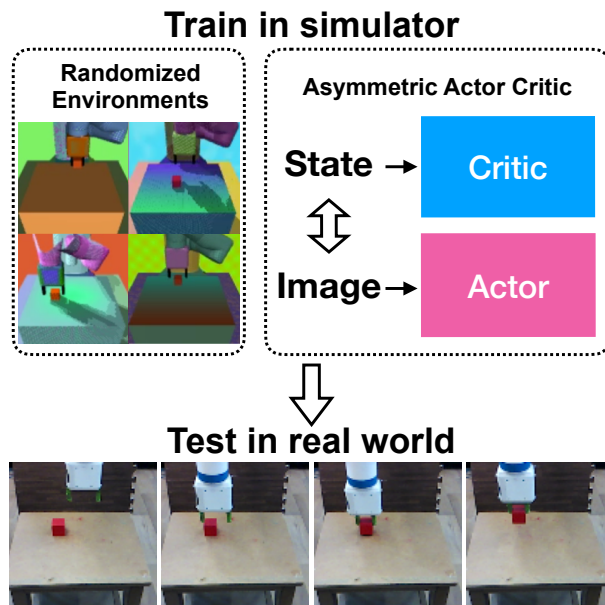


Fig. 1. By training policies with asymmetric inputs for actor-critic along with domain randomization, we learn complex visual policies that can operate in the real without having seen any real world data in training.

are not powerful enough to learn complex behaviours due to the large input dimensionality and partial observability. This leads us to a conundrum, i.e., training on full states is hard since it depends on good state predictors while training on images is hard because of their partial observability and dimensionality. We solve this by learning a policy that relies only on partial observations (RGBD images) but during training we exploit access to the full state.

Physics simulators give us access to both the full state of the system as well as rendered images of the scene. But how can we combine these observations to train complex behaviours faster? In this work, we exploit this access and train an *actor-critic* algorithm [4], [18] that uses asymmetric inputs, i.e. the actor takes visual partial observations as input while the critic takes the underlying full state as input. Since the critic works on full state, it learns the *state-action* value function much faster, which also allows for better updates for the actor. During testing, the actor is employed on the partial observations and does not depend on the full state (the full state is only used during training). This allows us to train an actor/policy network on visual observations while exploiting the availability of full states to train the critic. We experimentally show significant improvements on 2D

¹OpenAI {marcin, pw, woj, pieter}@openai.com

²Carnegie Mellon University lerrelp@cs.cmu.edu

[†]Equal advising

environments like *Particle* and *Reacher* and 3D environments like *Fetch Pick*. To further speed up training, we also demonstrate the utility of using bottlenecks [10].

Another key aspect of this work is to show that these policies learned in a simulator can be transferred to the real robot without additional real world data. Simulators are not perfect representations of the real world. The domain of observations (real camera images) significantly differs from rendered images from a simulator. This makes directly transferring policies from the simulator to the real world hard. However, [11], [17], [19] show how randomizing textures and lighting allows for effective transfer. By combining our asymmetric actor critic training with domain randomization [17], we show that these policies can be transferred to a real robot without any training on the physical system (see Figure 1).

II. RELATED WORK

A. Reinforcement Learning

Recent works in deep reinforcement learning (RL) have shown impressive results in the domain of games [1], [2] and simulated control tasks [3], [4]. The class of RL algorithms our method employs are called *actor-critic* algorithms [18]. Deep Deterministic Policy Gradients (DDPG) [4] is a popular *actor-critic* algorithm that has shown impressive results in continuous control tasks. Although we use DDPG for our base optimizer, our method is applicable to arbitrary *actor-critic* algorithms.

Learning policies in an environment that provides only sparse rewards is a challenging problem due to very limited feedback signal. However it has been shown that sparse rewards often allow for better policies when trained appropriately [14]. Moreover having sparse rewards allows us to circumvent manual shaping of the reward function.

The idea of using different inputs for the actor and critic has been explored previously in the domain of multi-agent learning [20], [21]. However using this in the domain of robot learning and dealing with partial observability hasn't been explored. Exploiting the access to full state in training the critic also draws similarities to Guided Policy Search [22].

B. Transfer from simulation to the real world

Bridging the reality gap in transferring policies trained in a simulator to the real world is an active area of research in the robot learning community. One approach is to make the simulator as close to the real world as possible [23]–[25]. But these methods have had limited success due to the hard system identification problem.

Another approach is domain adaptation from the simulator [5]–[10], since it may be easier to finetune from a simulator policy than training in the real world. However if the simulator differs from the real world by a large factor, the policy trained in simulation can perform very poorly in the real world and finetuning may not be any easier than training from scratch. This limits most of these works to learning simple behaviours. Making policies robust for

physics adaptation [26]–[28] is also receiving interest, but these methods haven't been shown to be powerful enough to work on real robots. Using bottlenecks [10] has been shown to help domain adaptation for simple tasks like *reaching*. In this work, we show how bottlenecks can be exploited for more complex fine manipulation tasks.

A promising approach is domain adaptation by domain randomization [11], [17]. Here the key idea is to train on randomized renderings of the scene, which allows to learn robust policies for transfer. However these works do not show transfer to precise manipulation behaviours. We show that this idea can be extended to complex behaviours when coupled with our asymmetric actor critic.

C. Robotic tasks

We perform real robot experiments on tasks like picking, pushing, and moving a block. The *Picking* task is similar to grasping objects [29], [30], however in this work we learn an end-to-end policy that moves to the object, grasps it and moves the grasped object to its desired position. The focus is hence on the fine manipulation behaviour. The tasks of *Forward Pushing* and *Block Move* are similar to pushing objects [31], [32], however as in the case of *Picking*, this paper focuses on the learning of the fine pushing behaviour.

III. BACKGROUND

Before we discuss our method, we briefly introduce some background and formalism for the RL algorithms used. A more comprehensive introduction can be found in [33].

A. Reinforcement Learning

In this paper we deal with continuous space Markov Decision Processes that can be represented as the tuple $(\mathcal{S}, \mathcal{O}, \mathcal{A}, \mathcal{P}, r, \gamma, \mathbb{S})$, where \mathcal{S} is a set of continuous states and \mathcal{A} is a set of continuous actions, $\mathcal{P} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ is the transition probability function, $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the reward function, γ is the discount factor, and \mathbb{S} is the initial state distribution. \mathcal{O} is a set of continuous partial observations corresponding to states in \mathcal{S} .

An episode for the agent begins with sampling s_0 from the initial state distribution \mathbb{S} . At every timestep t , the agent takes an action $a_t = \pi(s_t)$ according to a deterministic policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$. At every timestep t , the agent gets a reward $r_t = r(s_t, a_t)$, and the state transitions to s_{t+1} , which is sampled accordingly to probabilities $\mathcal{P}(s_t, a_t, \cdot)$. The goal of the agent is to maximize the expected return $E_{\mathbb{S}}[R_0 | \mathbb{S}]$, where the return is the discounted sum of the future rewards $R_t = \sum_{i=t}^{\infty} \gamma^{i-t} r_i$. The Q -function is defined as $Q^{\pi}(s_t, a_t) = E[R_t | s_t, a_t]$. In the partial observability case, the agent takes actions based on the partial observation, $a_t = \pi(o_t)$, where o_t is the observation corresponding to the full state s_t .

B. Deep Deterministic Policy Gradients (DDPG)

Deep Deterministic Policy Gradients (DDPG) [4] is an *actor-critic* RL algorithm that learns a deterministic continuous action policy. The algorithm maintains two neural

networks: the policy (also called the actor) $\pi_\theta : \mathcal{S} \rightarrow \mathcal{A}$ (with neural network parameters θ) and a Q -function approximator (also called the critic) $Q_\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ (with neural network parameters ϕ).

During training, episodes are generated using a noisy version of the policy (called behavioural policy), e.g. $\pi_b(s) = \pi(s) + \mathcal{N}(0, 1)$, where \mathcal{N} is Normal noise. The transition tuples (s_t, a_t, r_t, s_{t+1}) encountered during training are stored in a replay buffer [1]. Training examples sampled from the replay buffer are used to optimize the critic. By minimizing the Bellman error loss $\mathcal{L}_c = (Q(s_t, a_t) - y_t)^2$, where $y_t = r_t + \gamma Q(s_{t+1}, \pi(s_{t+1}))$, the critic is optimized to approximate the Q -function. The actor is optimized by minimizing the loss $\mathcal{L}_a = -E_s[Q(s, \pi(s))]$. The gradient of \mathcal{L}_a with respect to the actor parameters is called the deterministic policy gradient [34] and can be computed by backpropagating through the combined critic and actor networks.

To stabilize the training, the targets for the actor and the critic y_t are computed on separate versions of the actor and critic networks, which change at a slower rate than the main networks. A common practice is to use a Polyak averaged [35] version of the main network.

C. Multigoal RL

We are interested in learning policies that can achieve multiple goals (a universal policy). One way of doing this is by training policies and Q -functions that take as an additional input a goal $g \in \mathcal{G}$ [14], [36], e.g. $a_t = \pi(s_t, g)$. A universal policy can hence be trained by using arbitrary RL algorithms.

Following UVFA [36], the sparse reward formulation $r(s_t, a, g) = [d(s_t, g) < \epsilon]$ will be used in this work, where the agent gets a positive reward when the distance $d(\cdot, \cdot)$ between the current state and the goal is less than ϵ . In the context of a robot performing the task of picking and placing an object, this means that the robot gets a reward only if the object is within ϵ euclidean distance of the desired goal location of the object. Having a sparse reward overcomes the limitation of hand engineering the reward function, which often requires extensive domain knowledge. However, sparse rewards are not very informative and makes it hard to optimize. In order to overcome the difficulties with sparse rewards, we employ a recent method: Hindsight Experience Replay (HER) [14].

D. Hindsight Experience Replay (HER)

HER [14] is a simple method of manipulating the replay buffer used in off-policy RL algorithms that allows it to learn universal policies more efficiently with sparse rewards. After experiencing some episode s_0, s_1, \dots, s_T , every transition $s_t \rightarrow s_{t+1}$ along with the goal for this episode is usually stored in the replay buffer. However with HER, the experienced transitions are also stored in the replay buffer with different goals. These additional goals are states that were achieved later in the episode. Since the goal being pursued does not influence the environment dynamics, we can replay each trajectory using arbitrary goals assuming we use an off-policy RL algorithm to optimize [37].

IV. METHOD

We now describe our method along with the technique of bottlenecks to speed up training. Following this, we also describe domain randomization for transferring simulator learned policies to the real robot.

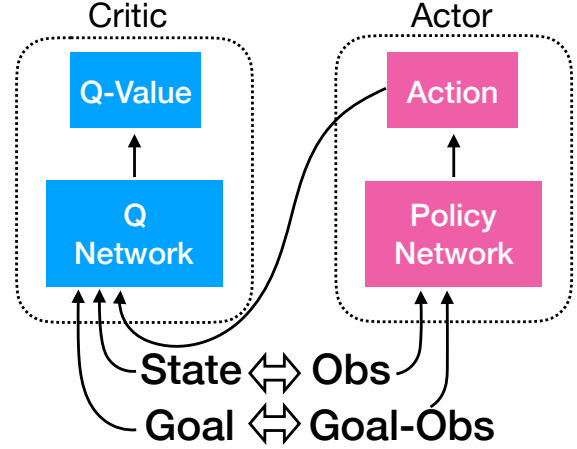


Fig. 2. Having asymmetric inputs, i.e. full states for the critic and partial observations for the actor improves training. In the multi goal setting, the critic additionally requires full goal states while the actor additionally requires partial observations for the goal.

A. Asymmetric Actor Critic

In its essence our method builds on *actor-critic* algorithms [18] by using the full state $s_t \in \mathcal{S}$ to train the critic, while using partial observation $o_t \in \mathcal{O}$ to train the actor (see Figure 2). Note that s_t is the underlying full state for the observation o_t . In our experiments, observations o_t are images taken by an external camera.

Algorithm 1 Asymmetric Actor Critic

```

Initialize actor-critic algorithm  $\mathcal{A}$ 
Initialize replay buffer  $R$ 
for episode= 1,  $M$  do
  Sample a goal  $g$  and an initial state  $s_0$ 
  Render goal observation  $g^o$ 
   $g^o \leftarrow \text{render}(g)$ 
  for  $t = 0, T - 1$  do
    Render image observation  $o_t$ 
     $o_t \leftarrow \text{render}(s_t)$ 
    Obtain action  $a_t$  using behavioural policy:
     $a_t \leftarrow \pi_b(o_t, g^o)$ 
    Execute action  $a_t$ , receive reward  $r_t$  and transition
    to  $s_{t+1}$ 
    Store  $(s_t, o_t, a_t, r_t, s_{t+1}, o_{t+1}, g, g^o)$  in  $R$ 
  end for
for  $n=1, N$  do
  Sample minibatch  $\{s, o, a, r, s', o', g, g^o\}_0^B$  from  $R$ 
  Optimize critic using  $\{s, a, r, s', g\}_0^B$  with  $\mathcal{A}$ 
  Optimize actor using  $\{o, a, r, o', g^o\}_0^B$  with  $\mathcal{A}$ 
end for
end for

```

The algorithm (described in Algorithm 1), begins with initializing the networks for an off-policy actor-critic algorithm \mathcal{A} [37]. In this paper, we use DDPG [4] as the actor-critic algorithm. The replay buffer R used by this algorithm is initialized with no data. For each episode, a goal g and an initial state s_0 are sampled before the rollout begins. g^o is the rendered goal observation. At every timestep t of the episode, a partially observable image of the scene o_t is rendered from the simulator at the full state s_t . The behavioural policy from \mathcal{A} , which is usually a noisy version of the actor is used to generate the action a_t for the agent/robot to take. After taking this action, the environment transitions to the next state s_{t+1} , with its corresponding rendered image o_{t+1} .

Since DDPG relies on a replay buffer to sample training data, we build the replay buffer from the episodic experience $(s_t, o_t, a_t, r_t, s_{t+1}, o_{t+1}, g, g^o)$ previously generated. To improve performance for the sparse reward case, we augment the standard replay buffer by adding hindsight experiences [14].

After the episodic experience has been added to the replay buffer R , we can now train our actor-critic algorithm \mathcal{A} from sampled minibatch of size B from R . This minibatch can be represented as $\{s, o, a, r, s', o', g, g^o\}_0^B$, where s' and o' are the next step full state and next step observation respectively. Since the critic takes full states as input, it is trained on $\{s, a, r, s', g\}_0^B$. Since the actor takes observations as input, it is trained on $\{o, a, r, o', g^o\}_0^B$. We experimentally show that asymmetric inputs for the critic and actor significantly improves performance and allows to transfer more complex manipulation behaviours to real robots.

B. Improvements with bottlenecks

One way of improving the efficiency of training is to use bottlenecks [10]. The key idea is to constrain one of the actor network’s intermediary layers to predict the full state. Since the full state is often of a smaller dimension than the other layers of the network, this state predictive layer is called the bottleneck layer.

C. Randomization for transfer

A powerful technique for domain transfer of policies from rendered images to real world images is domain randomization [11], [17]. The key idea is to randomize visual elements in the scene during the rendering. Learning policies with this randomization allows the policy to generalize to sources of error in the real world and latch on to the important aspects of the observation.

For the purposes of this paper, we randomize the following aspects: texture, lighting, camera location and depth. For textures, random textures are chosen among random RGB values, gradient textures and checker patterns. These random textures are applied on the different physical objects in the scene, like the robot and the table. For lighting randomization, we randomly switch on lighting sources in the scene and also randomize the position, orientation and the specular characteristics of the light. For camera location, we randomize the location of the monocular camera in a

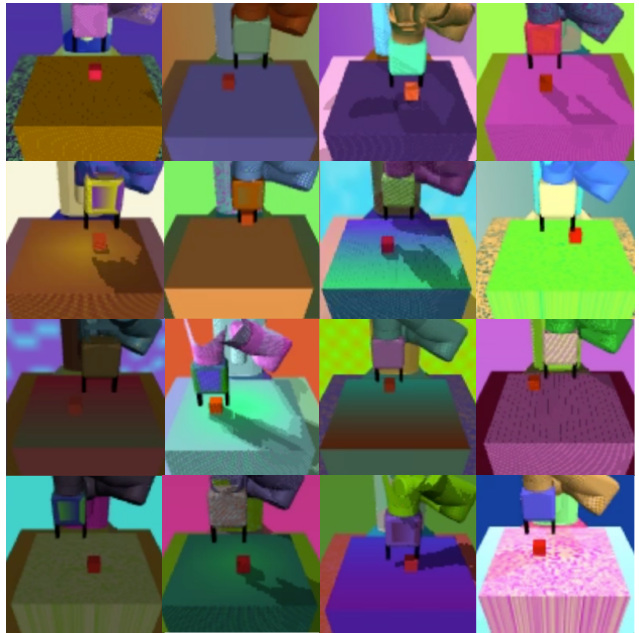


Fig. 3. To enable transfer of policies from the simulator to the real world, we randomize various aspects of the renderer during training. These aspects include textures, lighting and the position of the camera.

box around the expected location of the real world camera. Furthermore, we randomize the orientation and focal length of the camera and add uniform noise to the depth. RGB samples of randomization on the *Fetch Pick* environment can be seen in Figure 3.

V. RESULTS

To show the effectiveness of our method, we experiment on a range of simulated and real robot environments. In this section we first describe the environments. Following this, we discuss comparisons of our methods to baselines and show the utility of our method on improving training. Finally, we discuss real robot experiments.

A. Environments

Since there are no standard environments for multi-goal RL, we create three of our own simulated environments to test our method. The first two environments, *Particle* and *Reacher* are in a 2D workspace. The third environment *Fetch Pick* is in a 3D workspace with a simulated version of the Fetch robot needing to pick up and place a block. All these environments are simulated in the *MuJoCo* [38] physics simulator.

(a) *Particle*: In this 2D environment the goal for the agent is to move the 2D particle to a given location. The state space is 4D and consists of the particle’s location and velocity. The observation space is RGB images ($100 \times 100 \times 3$) from a camera placed above the scene. The action space is the 2D velocity of the particle. This action space allows for control on single RGB observations without requiring memory for velocity (since velocity cannot be inferred from a single RGB

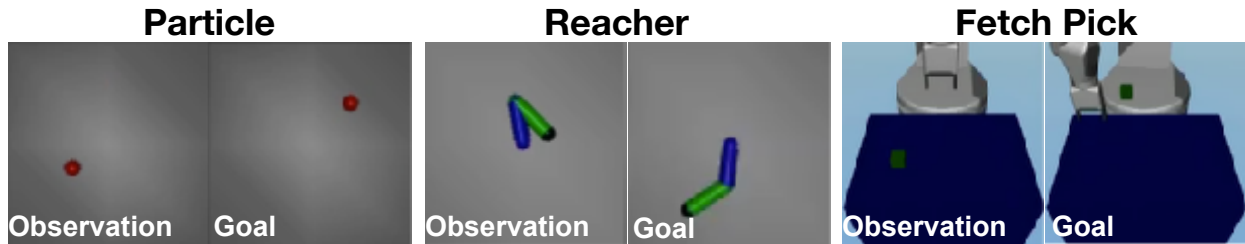


Fig. 4. To evaluate our method, we test on three different environments: *Particle*, *Reacher* and *Fetch Pick*. Since we learn multi goal policies, the policy takes in both the observation at timestep t and the desired goal for the episode.

frame). The agent gets a sparse reward (+1) if the particle is within ϵ of the desired goal position and no reward (0) otherwise. The observation for the goal is an image of the particle in its desired goal position.

(b) *Reacher*: In this 2D environment the goal for the agent is to move the end-effector of a two-link robot arm to a target location. The state space is 4D and consists of the joint positions and velocities. The observation space is RGB images ($100 \times 100 \times 3$) from a camera placed above the scene. The action space is the 2D velocities for the joints. The agent gets a sparse reward (+1) if the end-effector is within ϵ of the desired goal position and no reward (0) otherwise. The observation for the goal is an image of the reacher in its desired goal end-effector position.

(c) *Fetch Pick*: In this 3D environment with the simulated Fetch robot, the goal for the agent is to pick up the block on the table and move it to a given location in the air. The state space consists of the joint positions and velocities of the robot and the block on the table. The observation space is RGBD images ($100 \times 100 \times 4$) from a camera placed in front of the robot. The action space is 4D. Since this problem does not require gripper rotation, we keep it fixed. Three of the four dimensions of the action space specify the desired relative¹ position for the gripper. The last dimension specifies the desired distance between the fingers of the gripper. The agent gets a sparse reward (+1) if the block is within ϵ of the desired goal block position. The observation for the goal is an image of the block in its desired goal block position and the Fetch arm in a random position. To make exploration in this task easier following [14], we record a *single* state in which the box is grasped and start half of the training episodes from this state.

B. Robot evaluation

For our real world experiments we use a 7-DOF Fetch robotic arm², which is equipped with a two fingered parallel gripper. The camera observations for the real world experiments is an off the shelf Intel RealSense R200 camera that can provide aligned RGBD images. Since real depth often contains holes [39], we employ nearest neighbour hole filling to get better depth images [40]. To further improve the depth,

¹The desired gripper position is relative to the current gripper position.

²fetchrobotics.com/platforms-research-development/

we cover/recolor parts of the robot that are black like parts of the torso and parts of the gripper.

We experiment on three tasks for the real robot. The first task is *Pick* which is similar to the simulated task of *Fetch Pick* described in Section V-A(c). The second task is *Forward Push*, where the robot needs to push the block forward³. The third task is *Block Move*, where the robot needs to move the block to the target position on the table. In all tasks the goal is specified by an image of the box in the target location. A video of these experiments can be found in www.goo.gl/b57WTs and sample successes from our method in Figure 7.

The observations for the real robot tasks is an RGBD image from the physical camera placed in front of the robot. The goal observation for the actor is a simulated image describing the desired goal. We note that giving real world observations for the goal observation also works, however for consistency in evaluation, we use a simulated goal observation.

C. Does asymmetric inputs to actor critic help?

To study the effect of asymmetric inputs, we compare to the baseline of using symmetric inputs (images for both the actor and the critic networks). Figure 5 shows a summary of the final episodic rewards, with the x-axis being the number of episodes the agent experiences. As evident from the *Particle* results, asymmetric input versions of both DDPG and HER perform much better than their symmetric counterparts. The simplicity of the *Particle* may explain the similar performance between asymmetric DDPG and HER. *Fetch Pick* is a much harder sparse reward task, which shows the importance of using HER over DDPG. In this case as well, the asymmetric version of HER performs significantly better than the symmetric version.

D. Would imitation learning from an expert policy succeed?

Imitation learning is a powerful technique in robotics [41]. Hence a much stronger baseline is to behaviour clone from an expert policy. To do this we first train an expert policy [41] on full states that performs the task perfectly. Now given this expert policy, we behaviour clone to a policy that takes the partial visual observations as input. We use DAGger [42] for better imitation/cloning.

³The fingers are blocked for this task to avoid grasping.

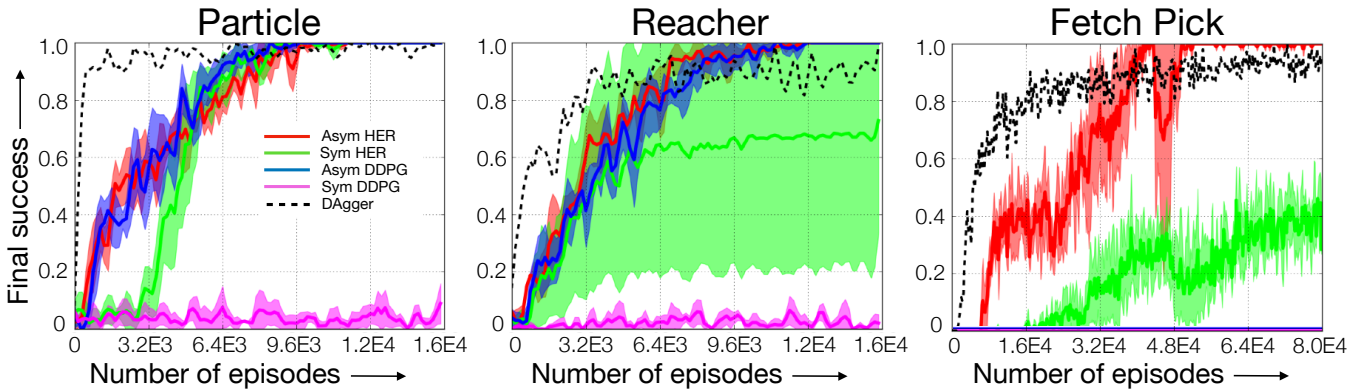


Fig. 5. We show that asymmetric inputs for training outperforms symmetric inputs by significant margins. The shaded region corresponds to ± 1 standard deviation across 5 random seeds. Although the behaviour cloning (BC) by expert imitation baseline (dashed lines) learn faster initially, it saturates to a sub optimal value compared to asymmetric HER. Also note that the BC baseline doesn't include the iterations the expert policy was trained on.

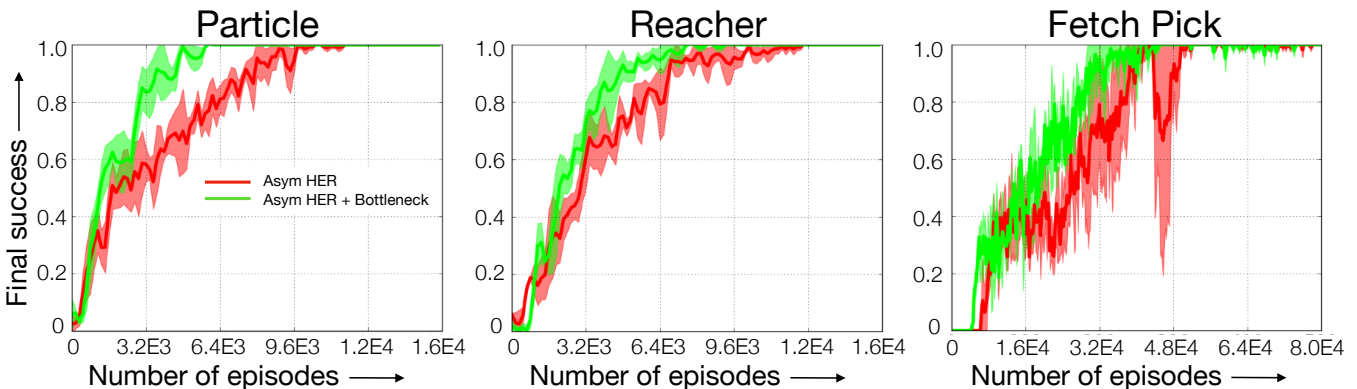


Fig. 6. We show that bottlenecks can be used to further improve training of our method. *Particle* and *Reacher*, the improvements are quite significant. On *Fetch Pick*, we observe more stable training (lower variance denoted by shaded regions).

Figure 5 shows the final episodic rewards of the behaviour cloned policy (in dashed lines) with the x-axis being the number of demonstrations. As expected, the DAgger policy learns much faster initially (since it receives supervision from a much stronger expert policy). However in all the environments, it saturates in performance and is lower than our method (asymmetric HER) for a large number of rollouts. One reason for this might be that behaviour cloning would fail if the expert policy depends on information contained in the full state but not in the partial observation.

E. Can we speed up training?

Another way of incorporating the full state from the simulator is by adding an auxiliary task of predicting the full state from partial observations. By adding a bottleneck layer [10] in the actor and adding an additional L_2 loss between the bottleneck output and the full state, we further speed up training. On our simulated tasks, these bottlenecks in the policy network improve the stability and speed of training (as seen in Figure 6).

F. How well do these policies transfer?

By combining asymmetric HER with domain randomization [17], we show significant performance gains (see Table

I) compared to baselines previously mentioned. Our method succeeds on all the three tasks for all the 5 times the policy was run with different block initializations and goals. We also note that behaviours like *push-grasping* [43] and *re-grasping* [44] emerge from these trained policies which can be seen in the video. Among the baselines we evaluate against, we note that behaviour cloning with DAgger is the only one that performs reasonably (as seen in the video and Table I).

G. How important is domain randomization?

To show the importance of randomization, we perform ablations (last two columns of Table I) by training policies without any randomization and testing them in the real world. We notice that without any randomization, the policies fail to perform on the real robot while performing perfectly in the simulator. Another randomization ablation is by removing the viewpoint randomization while keeping the texture and lighting randomization during training. We notice that apart from the *Block Move* task, removing viewpoint randomization severely affects performance.

Randomizing the observations in training also gives us an added benefit: robustness to distractors. Figure 8 shows the performance of our *Pick* policy, which was trained on

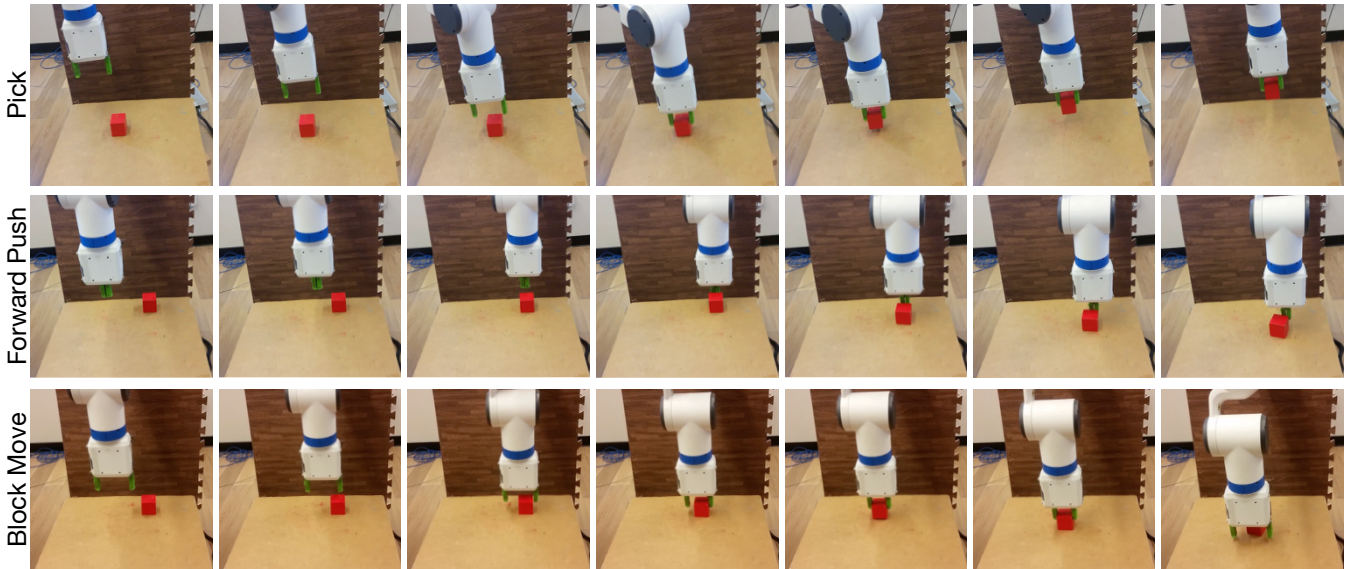


Fig. 7. Successive frames of our asymmetric HER policies on three real robot tasks show how our method can be successfully used for simulation to real transfer of complex policies. Full length experiments can be found in the videos on www.google.com/b57WTs.

TABLE I
COMPARISON OF ASYMMETRIC HER WITH BASELINES AND ABLATIONS

	Asym HER	Baselines				Randomization ablations	
		Asym DDPG	Sym HER	Vanilla BC	DAGger BC	Without any	Without viewpoints
Pick	5/5	0/5	0/5	0/5	3/5	0/5	0/5
Forward Push	5/5	0/5	0/5	1/5	0/5	0/5	0/5
Block Move	5/5	0/5	0/5	0/5	0/5	0/5	4/5

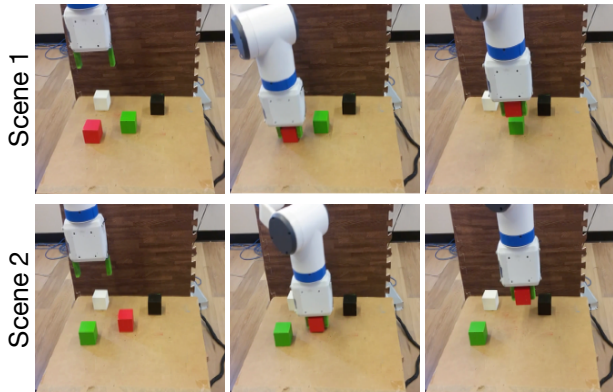


Fig. 8. Domain randomization during training allows the learned policies to be robust to distractors. Here we see how a policy trained to *Pick* the red block is robust to distractor blocks. The difference in the two scenes shown here is that in spite of changing the initial location of the red block, the arm still picks the red block.

a single red block, work even in the presence of distractor blocks.

H. Implementation Details

In this section we provide more details on our training setup. The critic is a fully connected neural network with 3 hidden layers; each with 512 hidden units. The hidden layers use ReLU [45] for the non linear activation. The input to

the critic is the concatenation of the current state s_t , the desired goal state g and the current action a_t . The actor is a convolutional neural network (CNN) with 4 convolutional layers with 64 filters each and kernel size of 2×2 . This network is applied on both the current observation o_t and the goal observation g^o . The outputs of both CNNs are then concatenated and passed through a fully connected neural network with 3 hidden layers. Similar to the critic, the hidden layers have 512 hidden units each with ReLU activation. The output of this actor network is normalized by a tanh activation and rescaled to match the limits of the environment's action space. In order to prevent tanh saturation, we penalize the preactivations in the actor's cost.

During each iteration of DDPG, we sample 16 parallel rollouts of the actor. Following this we perform 40 optimization steps on minibatches of size 128 from the replay buffer of size 10^5 transitions. The target actor and critic networks are updated every iteration with a polyak averaging of 0.98. We use Adam [46] optimization with a learning rate of 0.001 and the default Tensorflow [47] values for the other hyperparameters. We use a discount factor of $\gamma = 0.98$ and use a fixed horizon of $T = 50$ steps. For efficient learning, we also normalize the input states by running averages of the means and standard deviations of encountered states.

The behavioural policy chosen for exploration chooses a uniform random action from the space of valid actions with probability 20%. For the rest 80% probability, the output of

the actor is added with coordinate independent Normal noise with standard deviation equal to 5% of the action range.

VI. CONCLUSION

In this work we introduce asymmetric actor-critic, a powerful way of utilizing the full state observability in a simulator. By training the critic on full states while training its actor on rendered images, we learn vision-based policies for complex manipulation tasks. Our extensive evaluation both in the simulator and on the real world robot shows significant improvements over standard actor-critic baselines. This method's performance is also superior to the much stronger imitation learning with DAGger baseline, even though it was trained without an expert. Coupled with domain randomization, our method is able to learn visual policies that works in the real world while being trained solely in a simulator.

REFERENCES

- [1] V. Mnih *et al.*, "Human-level control through deep reinforcement learning," *Nature*, 2015.
- [2] D. Silver *et al.*, "Mastering the game of go with deep neural networks and tree search," *Nature*, 2016.
- [3] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz, "Trust region policy optimization," in *ICML 2015*.
- [4] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, "Continuous control with deep reinforcement learning," *arXiv preprint arXiv:1509.02971*, 2015.
- [5] M. Cutler and J. P. How, "Efficient reinforcement learning for robots using informative simulated priors," in *ICRA 2015*.
- [6] A. Ghadirzadeh, A. Maki, D. Kragic, and M. Björkman, "Deep predictive policy training using reinforcement learning," *arXiv preprint arXiv:1703.00727*, 2017.
- [7] J. Z. Kolter and A. Y. Ng, "Learning omnidirectional path following using dimensionality reduction," in *RSS*, 2007.
- [8] A. A. Rusu, M. Vecerik, T. Rothörl, N. Heess, R. Pascanu, and R. Hadsell, "Sim-to-real robot learning from pixels with progressive nets," *arXiv preprint arXiv:1610.04286*, 2016.
- [9] P. Christiano, Z. Shah, I. Mordatch, J. Schneider, T. Blackwell, J. Tobin, P. Abbeel, and W. Zaremba, "Transfer from simulation to real world through learning deep inverse dynamics model," *arXiv preprint arXiv:1610.03518*, 2016.
- [10] F. Zhang, J. Leitner, B. Upcroft, and P. Corke, "Vision-based reaching using modular deep networks: from simulation to the real world," *arXiv preprint arXiv:1610.06781*, 2016.
- [11] F. Sadeghi and S. Levine, "(cad)2rl: Real single-image flight without a single real image," *arXiv preprint arXiv:1611.04201*, 2016.
- [12] D. Held, Z. McCarthy, M. Zhang, F. Shentu, and P. Abbeel, "Probabilistically safe policy transfer," *ICRA 2017*.
- [13] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *CVPR 2014*.
- [14] M. Andrychowicz, F. Wolski, A. Ray, J. Schneider, R. Fong, P. Welinder, B. McGrew, J. Tobin, P. Abbeel, and W. Zaremba, "Hindsight experience replay," *arXiv preprint arXiv:1707.01495*, 2017.
- [15] R. E. Kalman *et al.*, "A new approach to linear filtering and prediction problems," *Journal of basic Engineering*, 1960.
- [16] A. Venkatraman, R. Capobianco, L. Pinto, M. Hebert, D. Nardi, and J. A. Bagnell, "Improved learning of dynamics models for control," in *International Symposium on Experimental Robotics*. Springer, 2016.
- [17] J. Tobin, R. Fong, A. Ray, J. Schneider, W. Zaremba, and P. Abbeel, "Domain randomization for transferring deep neural networks from simulation to the real world," *arXiv preprint arXiv:1703.06907*, 2017.
- [18] V. R. Konda and J. N. Tsitsiklis, "Actor-critic algorithms," in *Advances in neural information processing systems*, 2000, pp. 1008–1014.
- [19] U. Viereck, A. t. Pas, K. Saenko, and R. Platt, "Learning a visuomotor controller for real world robotic grasping using easily simulated depth images," *arXiv preprint arXiv:1706.04652*, 2017.
- [20] J. Foerster, G. Farquhar, T. Afouras, N. Nardelli, and S. Whiteson, "Counterfactual multi-agent policy gradients," *arXiv preprint arXiv:1705.08926*, 2017.
- [21] R. Lowe, Y. Wu, A. Tamar, J. Harb, P. Abbeel, and I. Mordatch, "Multi-agent actor-critic for mixed cooperative-competitive environments," *arXiv preprint arXiv:1706.02275*, 2017.
- [22] S. Levine, C. Finn, T. Darrell, and P. Abbeel, "End-to-end training of deep visuomotor policies," *JMLR 2016*.
- [23] S. James and E. Johns, "3d simulation for robot arm control with deep q-learning," *arXiv preprint arXiv:1609.03759*, 2016.
- [24] B. Planche, Z. Wu, K. Ma, S. Sun, S. Kluckner, T. Chen, A. Hutter, S. Zakharov, H. Kosch, and J. Ernst, "Depthsynth: Real-time realistic synthetic data generation from cad models for 2.5 d recognition," *arXiv preprint arXiv:1702.08558*, 2017.
- [25] S. R. Richter, V. Vineet, S. Roth, and V. Koltun, "Playing for data: Ground truth from computer games," in *ECCV 2016*.
- [26] A. Rajeswaran, S. Ghotra, S. Levine, and B. Ravindran, "Epopt: Learning robust neural network policies using model ensembles," *arXiv preprint arXiv:1610.01283*, 2016.
- [27] L. Pinto, J. Davidson, R. Sukthankar, and A. Gupta, "Robust adversarial reinforcement learning," *ICML*, 2017.
- [28] W. Yu, C. K. Liu, and G. Turk, "Preparing for the unknown: Learning a universal policy with online system identification," *arXiv preprint arXiv:1702.02453*, 2017.
- [29] A. Bicchi and V. Kumar, "Robotic grasping and contact: a review," in *ICRA 2000*.
- [30] L. Pinto and A. Gupta, "Supersizing self-supervision: Learning to grasp from 50k tries and 700 robot hours," *ICRA 2016*.
- [31] Z. Balorda, "Reducing uncertainty of objects by robot pushing," in *ICRA 1990*.
- [32] L. Pinto and A. Gupta, "Learning to push by grasping: Using multiple tasks for effective learning," in *ICRA 2017*.
- [33] L. P. Kaelbling, M. L. Littman, and A. W. Moore, "Reinforcement learning: A survey," *Journal of artificial intelligence research*, 1996.
- [34] D. Silver, G. Lever, N. Heess, T. Degris, D. Wierstra, and M. Riedmiller, "Deterministic policy gradient algorithms," in *ICML 2014*.
- [35] B. T. Polyak and A. B. Juditsky, "Acceleration of stochastic approximation by averaging," *SIAM Journal on Control and Optimization*, 1992.
- [36] T. Schaul, D. Horgan, K. Gregor, and D. Silver, "Universal value function approximators," in *ICML 2015*.
- [37] D. Precup, R. S. Sutton, and S. Dasgupta, "Off-policy temporal-difference learning with function approximation," in *ICML*, 2001.
- [38] E. Todorov, T. Erez, and Y. Tassa, "Mujoco: A physics engine for model-based control," in *IJROS*, 2012.
- [39] K.-J. Oh, S. Yea, and Y.-S. Ho, "Hole filling method using depth based in-painting for view synthesis in free viewpoint television and 3-d video," in *PCS 2009*.
- [40] N.-E. Yang, Y.-G. Kim, and R.-H. Park, "Depth hole filling using the depth distribution of neighboring regions of depth holes in the kinect sensor," in *ICSPCC 2012*.
- [41] B. D. Argall, S. Chernova, M. Veloso, and B. Browning, "A survey of robot learning from demonstration," *Robotics and autonomous systems*, 2009.
- [42] S. Ross, G. J. Gordon, and D. Bagnell, "A reduction of imitation learning and structured prediction to no-regret online learning," in *International Conference on Artificial Intelligence and Statistics*, 2011.
- [43] M. Dogar and S. Srinivasa, "A framework for push-grasping in clutter," *Robotics: Science and Systems (RSS)*, 2011.
- [44] T. Schlegl, M. Buss, T. Omata, and G. Schmidt, "Fast dextrous re-grasping with optimal contact forces and contact sensor-based impedance control," in *ICRA 2001*.
- [45] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *NIPS 2012*.
- [46] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [47] M. Abadi *et al.*, "Tensorflow: Large-scale machine learning on heterogeneous distributed systems," *arXiv preprint arXiv:1603.04467*, 2016.