

# OHPL: One-shot Hand-eye Policy Learner

Changjae Oh, Yik Lung Pang, and Andrea Cavallaro

**Abstract**—The control of a robot for manipulation tasks generally relies on object detection and pose estimation. An attractive alternative is to learn control policies directly from raw input data. However, this approach is time-consuming and expensive since learning the policy requires many trials with robot actions in the physical environment. To reduce the training cost, the policy can be learned in simulation with a large set of synthetic images. The limit of this approach is the domain gap between the simulation and the robot workspace. In this paper, we propose to learn a policy for robot reaching movements from a single image captured directly in the robot workspace from a camera placed on the end-effector (a hand-eye camera). The idea behind the proposed policy learner is that view changes seen from the hand-eye camera produced by actions in the robot workspace are analogous to locating a region-of-interest in a single image by performing sequential object localisation. This similar view change enables training of object reaching policies using reinforcement-learning-based sequential object localisation. To facilitate the adaptation of the policy to view changes in the robot workspace, we further present a dynamic filter that learns to bias an input state to remove irrelevant information for an action decision. The proposed policy learner can be used as a powerful representation for robotic tasks, and we validate it on static and moving object reaching tasks.

## I. INTRODUCTION

Data-driven robot learning trains an agent to acquire skills for tasks, such as locomotion [1], grasping [2], and manipulation [1], [3]. However, learning the control policies requires a large amount of training data and expert knowledge. Reinforcement Learning (RL) aims at learning policies by interacting with the environment [2], [4], [5]. The agent acquires generalised policies while exploring the environment and maximising a reward signal. RL with real-world training requires expensive and potentially unsafe active data collection by robots to gather a large number of experiences from scratch.

To address the limitation of active data acquisition in traditional learning-based approaches, simulation may be employed to learn control policies that are then transferred to the real world. The domain gap between the simulation and real world is addressed by domain adaptation [6], [7], domain randomisation [8], meta-learning [9], [10] or learning from depth images [11], [12], [13]. When the policy trained in simulation is adapted to the real-world robot workspace with sim-to-real techniques, precise camera calibration is needed to ensure the coherency between the geometric information from the simulation and real-world [13], [14].

Changjae Oh and Yik Lung Pang equally contributed. The authors are with the Centre for Intelligent Sensing, Queen Mary University of London, E1 4NS, London, United Kingdom. {c.oh, y.l.pang, a.cavallaro}@qmul.ac.uk.

In this paper, we propose a one-shot hand-eye robot policy learner (OHPL) for the object reaching task. OHPL does not require calibration information nor sim-to-real adaptation. Given the first-person view provided by a camera placed on the end-effector, we observe that the hand-eye robot view changes for the reaching task are similar to locating the object into a region-of-interest (RoI) from the first-person-view image. In OHPL, we first capture a single image from the robot workspace and then use RL-based sequential object localisation as a proxy task to learn the control policy for the object reaching task. OHPL learns the policy to generate sequential actions to change the position and size of the RoI (from the image) until it localises the object. To cope with the view changes in the robot workspace, OHPL includes a dynamic filter that generates image-conditioned parameters naturally learned to bias an input state to remove irrelevant information for an action decision. The learned policy can then be deployed for both the static and moving object reaching task in the robot workspace<sup>1</sup>.

## II. RELATED WORK

In this section, we review learning-based robot control, the domain gap problem, and proxy tasks to learn representations without explicit supervision.

Control policies can be learned with RL or in a supervised fashion (e.g., *imitation learning*, IL). With IL, agents acquire skills from demonstrations, which include kinesthetic demonstrations [15], teleoperation [16], and demonstrations from a human [17], [18], [19]. Expert demonstrations can be used directly to train the agent behaviour [10], [20] or be used to estimate a reward function for RL [21]. Although IL can achieve good performance using expert demonstrations to guide policy learning [22], it requires a large amount of training data and policies can be biased to the specific demonstration.

RL supports learning skill policies through trial and error, interactively within an environment by maximising cumulative reward received by the agent. RL can be directly applied in the real world to learn end-to-end policies for motor control [2], [4], [5], but a large number of trials are needed to generalise a policy and manual resetting may be needed between episodes, increasing the cost for data collection [4], [5], [23]. Policies are thus often trained first in simulated environments, then transfer to a real-world task [6], [7], [8]. If there is no significant domain shift between the simulation and test environments, the learned policy can be directly used

<sup>1</sup>Code and video results are available at: <http://cor-smal.eecs.qmul.ac.uk/OHPL.html>

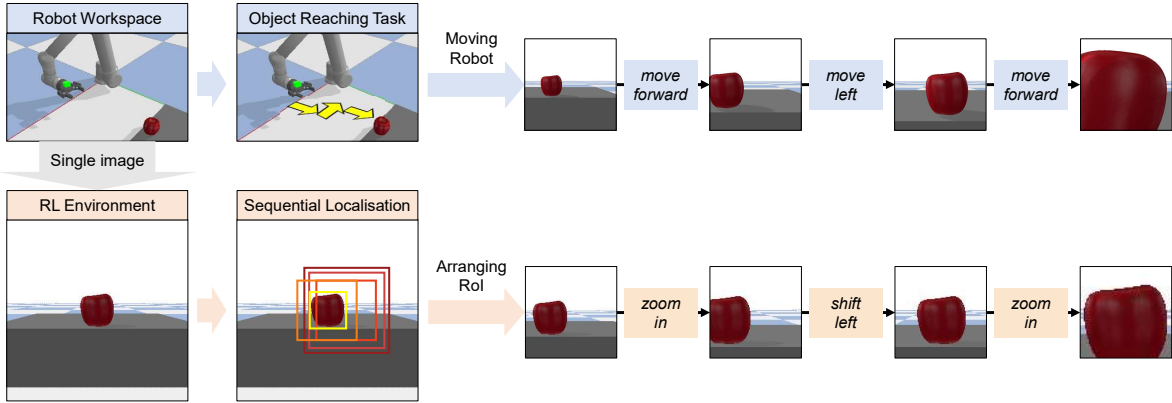


Fig. 1. Sequential input states and corresponding generated actions from (top) hand-eye robot manipulation and (bottom) RL-based object localisation. In (top-left) the robot workspace with hand-eye UR5 robot, our method takes (bottom-left) a single image in first-person view from the camera (green cube) on the gripper as an RL environment for sequential object localisation. In the RL environment, the input state for the RL agent is produced by cropping and resizing the RoI based on the bounding box. That is, the RL agent sequentially localises the object by arranging the position (shift) and size (zoom) of the RoI with discrete actions. The input states for the sequential localisation task show similar view changes to the hand-eye robot reaching task in the robot workspace. The agent thus can be directly deployed to the robot workspace since the proxy task is performed in the nearly same domain.

in the test environment [24]. However, transferring skills is challenging due to both visual and dynamics differences between the two domains [6], [7], [8].

To reduce the domain gap, a manipulation model can be learned with synthetic objects rendered with real object textures [25]. Another approach is to learn a policy from a modality with a relatively small domain gap, e.g. depth images [11], [12], [13]. Visual domain adaptation has also been used for robot manipulation, transferring knowledge from simulation to the real world. The methods include transfer learning [6] and learning policy with generated real-like images [7] or simplified images using generative adversarial networks [26].

Domain randomisation varies the appearance (e.g. textures, lighting and camera position) of sensory inputs during training, thus helping the agent learn relevant features that are invariant across the domain gap [8]. This technique can be applied to various neural network-based robot manipulation approaches with learning feed-forward networks [8], [27], imitation learning [28], [29], and reinforcement learning [26], [30], [31]. While domain randomisation provides good performance in addressing the domain gap, the trained model still has to address an unseen environment during testing. Domain randomisation in reinforcement learning setting may also lead to stability issues [30]. Since these real-to-simulation approaches do not learn the policy directly from the robot workspace, precise calibration is also a core issue to address robot manipulation [14], [13].

When expert supervision is unavailable or insufficient, a proxy task is an alternative to pre-train representations to be used, after finetuning, for specific tasks. For instance, visual representations can be learned from images, without annotations, by solving a proxy task, such as reconstructing input under corruptions such as noise, holes, colour, by performing denoising [32], inpainting [33], and colourisation [34], respectively. The representations can also be learned by classifying input with pseudo-label such as patch

orders [35], segmenting object motions [36], and clustering images [37]. An example proxy task for the RL agent is predicting a future state from the current state and action, which helps the agent to encode the representations that are effective for achieving the goal [38], [39].

### III. LEARNING OBJECT LOCALISATION AS A PROXY TASK

#### A. Problem definition

We cast the reaching task in the robot workspace as a sequential object localisation task in a first-person-view image [40], [41]. The view is changed by a set of discrete *robot actions* (Fig. 1, top). These view changes are similar to sequential object localisation with discrete *image-processing actions* (Fig. 1, bottom). We hypothesise that the model trained with such a task can, if appropriately formulated, naturally acquire skills for the agent to perform the reaching task in the robot workspace.

We formulate the sequential object localisation task by locating the RoI as a Markov Decision Process (MDP). We take an initial observation from the robot, a single image  $\mathcal{I}$ , to form an environment for RL. In the MDP, the RL agent learns to change the position and size of the bounding box over some discrete time steps until the RoI image captures the object in  $\mathcal{I}$ . At each time step  $t$ , the agent receives a state,  $s_t$ , the RoI image resized to the network input size, and selects an action,  $a_t$ , among  $N$  possible actions, according to the policy that maps  $s_t$  to  $a_t$ . After executing the action, a state transition function generates the next state,  $s_{t+1}$  and the agent receives a reward  $r_t$ .

#### B. State and action

The overview of our approach is shown in Fig. 2. At time step  $t$ , a square bounding box is represented by a three-dimensional vector  $b_t = [x_t, y_t, w_t]$  where  $[x_t, y_t]$  and  $w_t$  denote the top-left pixel and the window size of the bounding box, respectively. The RoI image, based on the bounding

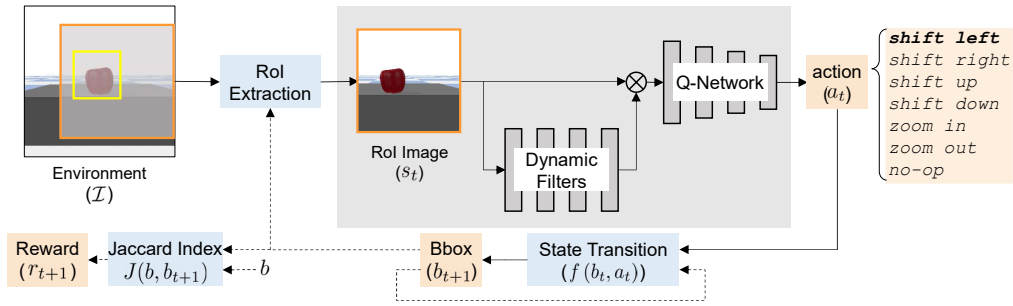


Fig. 2. Overview of the one-shot hand-eye policy Learner (OHPL). Given a single image as the RL environment,  $\mathcal{I}$ , and the bounding box (Bbox) at time step  $t$ ,  $b_t$ , the agent takes an extracted and resized RoI as an input state,  $s_t$ , passes it through the dynamic filter and then performs element-wise multiplication ( $\otimes$ ) with  $s_t$ . The combined input is fed into the Q-Network to determine the action,  $a_t$ . Based on  $a_t$ , the size or position of the next RoI,  $b_{t+1}$ , is changed with the state transition,  $f(b_t, a_t)$  to locate the target object. For instance, the agent generates  $a_t$  as *shift left* to move  $b_t$  (outlined orange) towards the groundtruth bounding box (outlined yellow),  $b$ , in the next step. The process is repeated until the Jaccard index between  $b$  and the current bounding box becomes higher than 0.8 or the agent moves the maximum actions in an episode. The dotted line denotes the process for  $t+1$  step. The process in the grey box can be directly deployed to the robot workspace where the hand-eye robot performs sequential discrete actions.

box, is extracted from the RL environment,  $\mathcal{I} \in \mathbb{R}^{360 \times 360 \times 3}$ , and then resized to the network input size, which becomes the input state  $s_t \in \mathbb{R}^{84 \times 84 \times 3}$  to the agent. The agent consists of convolutional neural networks that reduce dimension of the high-dimensional input state [38], [42]. We set the initial state,  $s_1$ , to an RoI image that includes the object but with random position and size of the bounding box. This random spawning strategy enables the agent to learn the generalised localisation policy in the single image environment  $\mathcal{I}$ .

We consider seven discrete robot actions ( $N = 7$ ): *move left*, *move right*, *move up*, *move down*, *move forward*, *move backward*, and *stop*. These actions and corresponding hand-eye view changes can be imitated in the proxy task that changes the position and size of the RoI in  $\mathcal{I}$  with the following discrete actions: *shift left*, *shift right*, *shift up*, *shift down*, *zoom in*, *zoom out*, and *no-operation (no-op)*. For instance, moving forward and left of the robot corresponds to *zoom in* and *shift left* actions by changing the window size,  $w_t$ , of the RoI and its position.

After the action  $a_t$  is decided from  $s_t$ , the next state  $s_{t+1}$  is obtained by the state transition function  $f(b_t, a_t)$  that arranges the position and size of the current bounding box,  $b_t$ . By arranging the position  $(x_t, y_t)$  with the displacement  $(\Delta x_t, \Delta y_t)$ , and  $w_t$  with the RoI size change  $\Delta w_t$ , we obtain the next bounding box,  $b_{t+1}$  as follows:

$$b_{t+1} = [x_t + \Delta x_t, y_t + \Delta y_t, w_t + \Delta w_t]. \quad (1)$$

The amount of displacement depends on the size of the RoI, which corresponds to the view changes in the robot workspace. For example, the displacement between the current and next state is small when the object is close, while the displacement is large when the object is far. Therefore, we adapt the amount of displacement based on the RoI size as follows:

$$\Delta x_t = \sigma_t w_t, \Delta y_t = \sigma_t w_t, \quad (2)$$

where  $\sigma_t$  controls the ratio between the amount of displacement and the bounding box size. Considering the PID control error that is likely to occur in real robot control we randomly generate  $\sigma_t \in [0.05, 0.15]$  in every step. The agent thus

can learn policies with more generalised cases to reach the goal position. The next state  $s_{t+1}$  can be finally obtained by extracting the region based on the computed  $b_{t+1}$  and resizing. For exception handling, we set  $s_{t+1} = s_t$  when the location of  $b_{t+1}$  is outside of  $\mathcal{I}$  or  $w_t$  is smaller or larger than 20 and 360, respectively.

### C. Reward function

We devise a reward function that encourages an RL agent to localise the object by producing dense rewards. The reward function integrates intrinsic motivations to encourage the agent to naturally move towards the goal to achieve higher Jaccard index, alleviating the problem of sparse extrinsic rewards. For each time step  $t$ , the reward,  $r_t$ , is defined as follows:

$$r_t = \begin{cases} 1 & J(b, b_t) > 0.8 \\ \alpha(J(b, b_t) - 1) & 0.8 \geq J(b, b_t) > 0 \\ -1 & J(b, b_t) = 0, \end{cases} \quad (3)$$

where  $J(b, b_t)$  measures the overlap ratio between the groundtruth bounding box,  $b$ , and the bounding box at  $t$ ,  $b_t$ . We consider that the agent has reached the goal when the overlap ratio between  $b$  and  $b_t$  is higher than 0.8, where the agent receives a positive reward of +1. Instead of giving a reward of -1 to all the states which fail to achieve the goal, we introduce a small negative reward that encourages the agent to move in order to localise the object with higher Jaccard index score. When  $b_t$  partially overlaps with  $b$ , the small negative reward can be obtained by  $\alpha(J(b, b_t) - 1)$ , with  $\alpha = 0.5$ . The higher the Jaccard index, the smaller the negative reward. Lastly, a negative reward of -1 is given when  $b_t$  does not overlap with  $b$ . This reward setting encourages the RoI to not move away from the object during training. The episode is terminated when the agent reaches the goal or has reached the maximum number of steps  $t = T_{max}$ .

### D. Learning policies with Deep Q-Network

We address the learning of the proxy task using the DQN framework [42]. We employ a state encoding network to



Fig. 3. Sequential object localisation examples in (a) the RL environment  $\mathcal{I}$ , where the bounding box is colour-coded from dark blue to yellow to indicate the temporal evolution of the actions, and (b) corresponding sequence of input states for the OHPL agent.

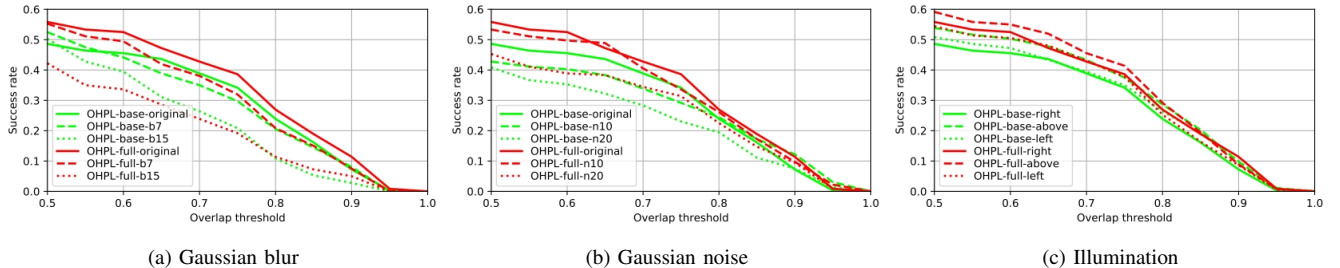


Fig. 4. Sequential object localisation results when varying the overlap thresholds on the test data. To evaluate the model robustness, we tested the sequential object localisation by corrupting the original image by applying (a) Gaussian blur, (b) Gaussian noise (c) and different location of the light source.

approximate the underlying action-value function:

$$Q(s_t, a_t) = r_t + \gamma \max_{a_{t+1}} Q(s_{t+1}, a_{t+1}), \quad (4)$$

where  $\gamma$  is a discount factor that indicates the value of future rewards in the current state. The Q-value  $Q(s_{t+1}, a_{t+1})$  presents how good it is to select  $a_{t+1}$  at  $s_{t+1}$  to receive a high reward.  $Q(s_{t+1}, a_{t+1}) = 0$  when  $s_{t+1}$  is the terminal state. Deep reinforcement learning methods commonly require a large discount factor as the agent may explore thousands of steps before reaching the terminal target. The goal in our environment is around 7-15 steps far away, and we set a low discount factor,  $\gamma = 0.85$  to ensure that the agent follows the intrinsic motivations while optimising for long term success.

Note that discrepancies exist between inputs generated from the 2D single image,  $\mathcal{I}$ , and the 3D robot workspace, due to sensor noise, blur, illumination, and object shape changes. We propose to use a dynamic filter [43] as soft-attention to provide effective information as an input for localisation. The network takes  $s_t$  as an input and generates the feature map which has the same size as  $s_t$ . We then fuse the output feature map and  $s_t$  by multiplying them in an element-wise manner, where the output becomes an input for the Q-Network. Noting that the feature map generated from the dynamic filter is conditioned on the input, the dynamic filter naturally learns to bias the input during the training without any loss function. This process can be considered as estimating the soft-attention based on the input image.

Fig. 3 shows an example of sequential object localisation. Initially, a bounding box is set with a large RoI, and then the agent arranges the RoI position and size to locate the object.

#### IV. VALIDATION

##### A. Models under analysis

We evaluate the OHPL agent on the reaching task when deployed in the robot simulation workspace, PyBullet [45],

but first we analyse the proposed proxy task, RL-based sequential object localisation in the single image environment using two models, namely OHPL-base (vanilla DQN model [42]), and OHPL-full (DQN model with dynamic filter). A single image for each object illuminated from the left is used as the RL environment (see Fig.5). We use 8 synthetic objects, i.e. *baseball*, *cup-H*, *cup-J*, *rubiks cube*, *apple*, *mug*, *can*, *cracker box*, that are adopted from the Yale-Carnegie Mellon University-Berkeley (YCB) benchmark [44], including objects with different shapes, textures, and colours. As each image is employed as a separate RL environment, we generate 8 models each on OHPL-base and OHPL-full, respectively. In each episode, we randomly set a position and size for the initial RoI, while including the object in the RoI. We set  $T_{max} = 50$  for a single episode, where the localisation can be achieved in 7-15 steps. We set the position of the ground-truth bounding box as shown in the rightmost images in Fig. 3. We compare OHPL-base and OHPL-full on the test environment with nine object locations, *middle (M)*, *mid-left (ML)*, *mid-right (MR)*, *top (T)*, *top-left (TL)*, *top-right (TR)*, *bottom (B)*, *bottom-left (BL)*, and *bottom-right (BR)*, and three different illuminations by setting the light source location to left, right, and above. We perform 10 trials for each object location and illumination.

We further test the models by corrupting the test images with Gaussian noise, Gaussian blur, and varying illumination condition, which emulates camera sensor noise, blur, illumination changes in the robot workspace. The Gaussian blur is applied with two kernel sizes, 7 and 15, and the Gaussian noise with two variances, 10 and 20. Noting that the light in the training environment is located to the left-hand side, we also change the light directions to right and above.

We then validate the OHPL agent, trained from the proxy task, in the robot simulation workspace for the object reaching task. We use a 6 DoF UR5 robotic arm [46] with the Robotiq 2F-85 gripper [47]. The arm is controlled by



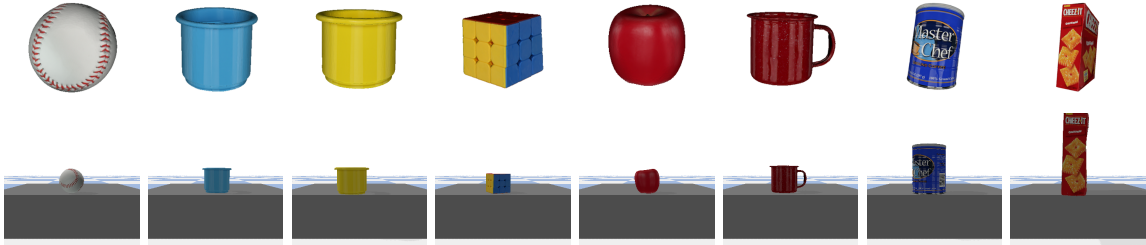


Fig. 5. (Top) Test objects from YCB Benchmark [44] and (bottom) their first-person view as an RL environment,  $\mathcal{I}$ , for training the OHPL agent.

specifying a pose goal in the Cartesian space. Actions issued to the arm have a fixed displacement value of 0.02m. The object is placed on a platform in front of the gripper at about 30 steps away.  $360 \times 360 \times 3$  images are captured at each step. We consider the task completed when the gripper reaches a predefined area in front of the object.

We compare the OHPL agent, *without robot training*, and a model *with robot training*. For the latter, we *train an agent from scratch* (Scratch) following the DQN framework [42]. The Scratch agent first explores randomly for 1000 steps and then follows the  $\epsilon$ -greedy strategy. We pre-defined the explorable region to prevent the agent from exploring states irrelevant to the task and to prevent the robot from colliding with the environment e.g., a grey box to place the target object. We set the maximum steps in an episode as  $T_{max} = 50$  and give a positive extrinsic reward when successful reaching is performed. In addition to the extrinsic reward, at each step, an intrinsic reward is given based on the distance between the gripper and the object to encourage the robot to reach the object and converge faster.

Finally, we validate the OHPL agent by deploying the trained agents to the moving object reaching task where the environment is unseen to both the OHPL agent and the agent with robot training (Scratch). We initialise the object floating without the grey box in the environment. During each step, we randomise the velocity of the object in the x, y and z directions. The same success criteria from the static object reaching task is used. Experiments are performed on all 8 objects in 3 lighting conditions and 9 starting positions. 3 trials are performed for each combination so a total of 648 experiments are performed for each model.

### B. Network setting

In this section, we introduce the network configuration for OHPL agent. For sequential object localisation, an input state is cropped from the RL environment  $\mathcal{I}$  based on the RoI. The hand-eye view is used as the input state for object reaching task in the robot workspace. For both tasks, the input state is then resized to  $84 \times 84 \times 3$ . We adopt the context aggregation network [48] architecture with the rectified linear unit (ReLU) as nonlinear activation for the dynamic filter using 4 convolution layers with 3 intermediate feature maps. We set the filter kernels size to  $3 \times 3$ . The last convolution layer applies a  $1 \times 1$  convolution followed by a sigmoid to predict the final output, which is then combined with the RoI.

TABLE I

SUCCESS RATIO, WITH STANDARD ERROR IN BRACKETS, FOR THE STATIC OBJECT REACHING TASK IN THE ROBOT WORKSPACE WITH AND WITHOUT TRAINING IN THE ROBOT WORKSPACE (ROBOT.), AND WITH (ARCH. FULL) AND WITHOUT (ARCH. BASE) THE DYNAMIC FILTER

Model		Initial position								
Method	Robot. Arch.	M	ML	MR	T	TL	TR	B	BL	BR
Scratch	✓ base	100.0 (0.0)	100.0 (0.0)	87.5 (6.8)	100.0 (0.0)	91.7 (5.6)	100.0 (0.0)	100.0 (0.0)	100.0 (0.0)	95.8 (4.1)
	✓ full	100.0 (0.0)	100.0 (0.0)	100.0 (0.0)	100.0 (0.0)	100.0 (0.0)	87.5 (6.8)	87.5 (6.8)	100.0 (0.0)	87.5 (6.8)
OHPL	base	95.8 (4.1)	95.8 (4.1)	75.0 (8.8)	91.7 (5.6)	95.8 (4.1)	70.8 (9.3)	20.8 (8.3)	8.3 (5.6)	4.2 (4.1)
	full	95.8 (4.1)	79.2 (8.3)	79.2 (8.3)	79.2 (8.3)	87.5 (6.8)	83.3 (7.6)	37.5 (9.9)	54.2 (10.2)	4.2 (4.1)

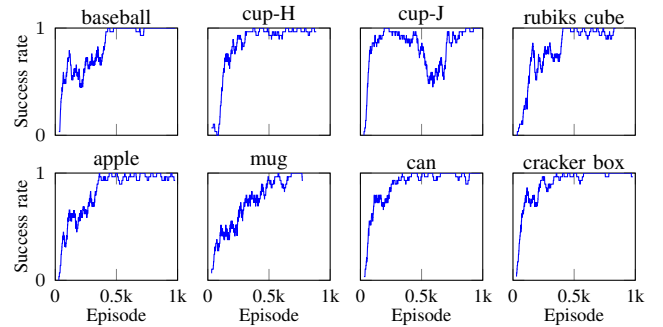


Fig. 6. Training the agent from scratch (Scratch-base) in the robot workspace. We report the success rates (vertical axis) with a running average over past 30 episodes for each episode (horizontal axis) on 8 YCB objects.

The dilation factors of the convolution layers are 1, 2, 4, and 1, respectively. As Q-Network we use the vanilla DQN [42] with minor modification to the final fully connected layer with a single output for each valid action. The experience replay has a  $5k$  memory volume. Note that the dynamic filter makes 0.01% parameter increase to the total number of parameters of the network. The network is optimised by RMSprop [49], where the learning rate is set to 0.001 and decays exponentially. We adopt the  $\epsilon$ -greedy strategy during training and exponentially anneal  $\epsilon$  from 0.95 to 0.05 [42].

### C. Discussion

Fig. 4 shows the results for *sequential object localisation* with varying the image conditions. OHPL-full generally outperforms OHPL-base in localising the object, and the model

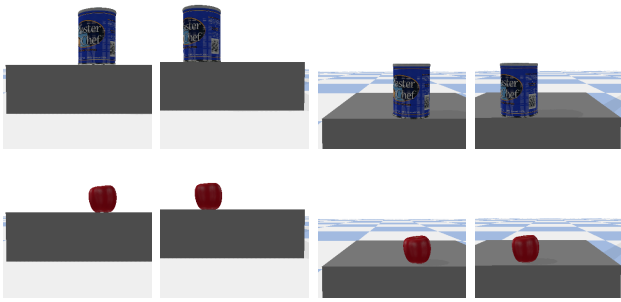


Fig. 7. Examples of hand-eye view from different initial positions, (from left to right) bottom-left (BL), bottom-right (BR), top-left (TL), and top-right (TR). The view from BL and BR is extremely different from the RL environment  $\mathcal{I}$  (Fig. 5), which results in low success of OHP agents.

shows consistent performance in noisy and blurry images. But the performance drops when Gaussian blurring is applied with a  $15 \times 15$  kernel. Also, the performance decreases when the light source is located to right, demonstrating that the performance can be affected when significant light change happens.

Table I shows the evaluation of the *static object reaching task* in the robot workspace. Note that Scratch is the upper bound of OHPL as the agent of Scratch models are trained directly in the robot workspace with robot actions. Scratch models are hence overfitted to the environment as the models have experiences with the real view changes and rewards in the robot workspace. However, to learn the control policies for the object-reaching task, each Scratch model requires more than 500 episodes of training in the robot workspace as shown in Fig. 6. The OHPL agents, without directly learning from a robot workspace and without robot actions, show good results when the initial position is top,  $T$ ,  $TL$  and  $TR$ , or middle,  $M$ ,  $ML$  and  $MR$ , while showing low success rates in bottom cases,  $B$ ,  $BL$ , and  $BR$ . As shown in Fig. 7, the hand-eye view from the bottom position is extremely different from the middle view which is used as the RL environment  $\mathcal{I}$  (see Fig. 5). OHPL-full, with the additional dynamic filter, records better success rates than OHPL-base on average, by achieving better performance in bottom cases.

For a fair comparison, we further evaluate the agents when they are deployed to the *moving object reaching task*. Note that this task can be considered as an unseen scenario to both OHPL and Scratch models. Table II shows the success ratio of the models. Compared to the success ratio of the static object reaching task, Scratch models show significant performance decrease when they are deployed to the moving object reaching task. While OHPL agents trained their control policies within a single view captured from the static object reaching environment, they generally show good performance and are competitive to the Scratch agents.

## V. CONCLUSION

We presented a new approach to learn robot control policies with a proxy task: RL-based object localisation with a single image captured from the robot workspace. The key idea is that view changes of the RoI in a 2D image are similar to the changes of the hand-eye camera view in the

TABLE II  
SUCCESS RATIO, WITH STANDARD ERROR IN BRACKETS, OF THE MODELS DEPLOYED TO THE MOVING OBJECT REACHING TASK IN THE ROBOT WORKSPACE WITH (ARCH. FULL) AND WITHOUT (ARCH. BASE) THE DYNAMIC FILTER.

Model		Object							
Method	Arch.	ball	cup-H	cup-J	rubiks	apple	mug	can	cracker
Scratch	base	28.4 (5.0)	87.7 (3.7)	77.8 (4.6)	92.6 (2.9)	58.0 (5.5)	64.2 (5.3)	85.2 (3.9)	63.0 (5.4)
	full	29.6 (5.1)	75.3 (4.8)	58.0 (5.5)	42.0 (5.5)	71.6 (5.0)	55.6 (5.5)	84.0 (4.1)	98.8 (1.2)
OHPL	base	67.9 (5.2)	82.7 (4.2)	54.3 (5.5)	90.1 (3.3)	69.1 (5.1)	66.7 (5.2)	56.8 (5.5)	66.7 (5.2)
	full	79.0 (4.5)	66.7 (5.2)	84.0 (4.1)	75.3 (4.8)	92.6 (2.9)	91.4 (3.1)	90.1 (3.3)	76.5 (4.7)

robot workspace. We showed that the policy learned from OHPL can be deployed to the robot workspace *directly* as the control policy for the reaching task, without finetuning. We also presented a dynamic filter that works as soft attention of an input image. With only 0.01% parameter increase, the network has better generalisation when the agent is deployed to perform static and moving object reaching task.

Future work includes validating the proposed framework with a physical robot, applying to other robotic tasks, and learning policies from different viewpoints.

## ACKNOWLEDGEMENT

This work is supported in part by the CHIST-ERA program through the project CORSMAL, under the UK EPSRC grant EP/S031715/1, and the UK EPSRC grant NCNR EP/R02572X/1.

## REFERENCES

- [1] N. Ratliff, J. A. Bagnell, and S. S. Srinivasa, "Imitation learning for locomotion and manipulation," in *Proceedings of the IEEE-RAS International Conference on Humanoid Robots (ICHR)*, 2007.
- [2] J. Peters and S. Schaal, "Reinforcement learning of motor skills with policy gradients," *Neural Networks*, vol. 21, no. 4, pp. 682–697, 2008.
- [3] P. Sharma, D. Pathak, and A. Gupta, "Third-person visual imitation learning via decoupled hierarchical controller," in *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- [4] S. Levine, C. Finn, T. Darrell, and P. Abbeel, "End-to-end training of deep visuomotor policies," *The Journal of Machine Learning Research (JMLR)*, vol. 17, no. 1, pp. 1334–1373, 2016.
- [5] S. Levine, P. Pastor, A. Krizhevsky, J. Ibarz, and D. Quillen, "Learning hand-eye coordination for robotic grasping with deep learning and large-scale data collection," *The International Journal of Robotics Research (IJRR)*, vol. 37, no. 4-5, pp. 421–436, 2018.
- [6] 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.
- [7] K. Bousmalis, A. Irpan, P. Wohlhart, Y. Bai, M. Kelcey, M. Kalakrishnan, L. Downs, J. Ibarz, P. Pastor, K. Konolige *et al.*, "Using simulation and domain adaptation to improve efficiency of deep robotic grasping," in *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, 2018.
- [8] 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," in *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2017.
- [9] C. Finn, P. Abbeel, and S. Levine, "Model-agnostic meta-learning for fast adaptation of deep networks," in *Proceedings of the International Conference on Machine Learning (ICML)*, 2017.

- [10] Y. Duan, M. Andrychowicz, B. Stadie, O. J. Ho, J. Schneider, I. Sutskever, P. Abbeel, and W. Zaremba, "One-shot imitation learning," in *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- [11] M. Gualtieri, A. ten Pas, K. Saenko, and R. Platt, "High precision grasp pose detection in dense clutter," in *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2016.
- [12] U. Viereck, A. ten Pas, K. Saenko, and R. Platt, "Learning a visuomotor controller for real world robotic grasping using simulated depth images," *arXiv preprint arXiv:1706.04652*, 2017.
- [13] M. Gualtieri and R. Platt, "Learning 6-dof grasping and pick-place using attention focus," *arXiv preprint arXiv:1806.06134*, 2018.
- [14] M. Yan, I. Frosio, S. Tyree, and J. Kautz, "Sim-to-real transfer of accurate grasping with eye-in-hand observations and continuous control," *arXiv preprint arXiv:1712.03303*, 2017.
- [15] S. Calinon, F. Guenter, and A. Billard, "On learning, representing, and generalizing a task in a humanoid robot," *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, vol. 37, no. 2, pp. 286–298, 2007.
- [16] P. Pastor, H. Hoffmann, T. Asfour, and S. Schaal, "Learning and generalization of motor skills by learning from demonstration," in *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, 2009.
- [17] B. C. Stadie, P. Abbeel, and I. Sutskever, "Third-person imitation learning," *arXiv preprint arXiv:1703.01703*, 2017.
- [18] P. Sermanet, C. Lynch, Y. Chebotar, J. Hsu, E. Jang, S. Schaal, and S. Levine, "Time-contrastive networks: Self-supervised learning from video," in *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, 2018.
- [19] S. Song, A. Zeng, J. Lee, and T. Funkhouser, "Grasping in the wild: Learning 6dof closed-loop grasping from low-cost demonstrations," *IEEE Robotics and Automation Letters (RA-L)*, vol. 5, no. 3, pp. 4978–4985, 2020.
- [20] D. A. Pomerleau, "Efficient training of artificial neural networks for autonomous navigation," *Neural Computation*, vol. 3, no. 1, pp. 88–97, 1991.
- [21] P. Abbeel and A. Y. Ng, "Apprenticeship learning via inverse reinforcement learning," in *Proceedings of the International Conference on Machine Learning (ICML)*, 2004.
- [22] L. Tai, J. Zhang, M. Liu, J. Boedecker, and W. Burgard, "A survey of deep network solutions for learning control in robotics: From reinforcement to imitation," *arXiv preprint arXiv:1612.07139*, 2016.
- [23] D. Kalashnikov, A. Irpan, P. Pastor, J. Ibarz, A. Herzog, E. Jang, D. Quillen, E. Holly, M. Kalakrishnan, V. Vanhoucke, and S. Levine, "Qt-opt: Scalable deep reinforcement learning for vision-based robotic manipulation," *arXiv preprint arXiv:1806.10293*, 2018.
- [24] E. Tzeng, C. Devin, J. Hoffman, C. Finn, P. Abbeel, S. Levine, K. Saenko, and T. Darrell, "Adapting deep visuomotor representations with weak pairwise constraints," *Algorithmic Foundations of Robotics XII*, vol. 13, pp. 688–703, 2016.
- [25] A. Saxena, J. Driemeyer, and A. Y. Ng, "Robotic grasping of novel objects using vision," *The International Journal of Robotics Research (IJRR)*, vol. 27, no. 2, pp. 157–173, 2008.
- [26] S. James, P. Wohlhart, M. Kalakrishnan, D. Kalashnikov, A. Irpan, J. Ibarz, S. Levine, R. Hadsell, and K. Bousmalis, "Sim-to-real via sim-to-sim: Data-efficient robotic grasping via randomized-to-canonical adaptation networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [27] S. James, A. J. Davison, and E. Johns, "Transferring end-to-end visuomotor control from simulation to real world for a multi-stage task," *arXiv preprint arXiv:1707.02267*, 2017.
- [28] S. James, M. Bloesch, and A. J. Davison, "Task-embedded control networks for few-shot imitation learning," *arXiv preprint arXiv:1810.03237*, 2018.
- [29] A. Bonardi, S. James, and A. J. Davison, "Learning one-shot imitation from humans without humans," *IEEE Robotics and Automation Letters (RA-L)*, vol. 5, pp. 3533–3539, 2020.
- [30] J. Matas, S. James, and A. J. Davison, "Sim-to-real reinforcement learning for deformable object manipulation," *arXiv preprint arXiv:1806.07851*, 2018.
- [31] S. Iqbal, J. Tremblay, A. Campbell, K. Leung, T. To, J. Cheng, E. Leitch, D. McKay, and S. Birchfield, "Toward sim-to-real directional semantic grasping," in *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, 2020.
- [32] P. Vincent, H. Larochelle, Y. Bengio, and P.-A. Manzagol, "Extracting and composing robust features with denoising autoencoders," in *Proceedings of the International Conference on Machine Learning (ICML)*, 2008.
- [33] D. Pathak, P. Krahenbuhl, J. Donahue, T. Darrell, and A. A. Efros, "Context encoders: Feature learning by inpainting," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [34] R. Zhang, P. Isola, and A. A. Efros, "Colorful image colorization," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016.
- [35] C. Doersch, A. Gupta, and A. A. Efros, "Unsupervised visual representation learning by context prediction," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2015.
- [36] D. Pathak, R. Girshick, P. Dollár, T. Darrell, and B. Hariharan, "Learning features by watching objects move," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [37] M. Caron, P. Bojanowski, A. Joulin, and M. Douze, "Deep clustering for unsupervised learning of visual features," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018.
- [38] D. Pathak, P. Agrawal, A. A. Efros, and T. Darrell, "Curiosity-driven exploration by self-supervised prediction," in *Proceedings of the International Conference on Machine Learning (ICML)*, 2017.
- [39] T. Lesort, N. Díaz-Rodríguez, J.-F. Goudou, and D. Filliat, "State representation learning for control: An overview," *Neural Networks*, vol. 108, pp. 379–392, 2018.
- [40] S. Yun, J. Choi, Y. Yoo, K. Yun, and J. Y. Choi, "Action-decision networks for visual tracking with deep reinforcement learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [41] Z. Jie, X. Liang, J. Feng, X. Jin, W. Lu, and S. Yan, "Tree-structured reinforcement learning for sequential object localization," in *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, 2016.
- [42] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, "Playing atari with deep reinforcement learning," *arXiv preprint arXiv:1312.5602*, 2013.
- [43] X. Jia, B. De Brabandere, T. Tuytelaars, and L. V. Gool, "Dynamic filter networks," in *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, 2016.
- [44] B. Calli, A. Walsman, A. Singh, S. Srinivasa, P. Abbeel, and A. M. Dollar, "Benchmarking in manipulation research: Using the yale-cmu-berkeley object and model set," *IEEE Robotics Automation Magazine*, vol. 22, no. 3, pp. 36–52, 2015.
- [45] E. Coumans and Y. Bai, "Pybullet, a python module for physics simulation for games, robotics and machine learning," <http://pybullet.org>, 2016–2019.
- [46] "Universal robots." <http://www.universal-robots.com/>, accessed: 2020-05-23.
- [47] "Robotiq 2f-85 gripper." <http://robotiq.com/products/adaptive-robot-gripper/>, accessed: 2020-05-23.
- [48] F. Yu and V. Koltun, "Multi-scale context aggregation by dilated convolutions," *arXiv preprint arXiv:1511.07122*, 2015.
- [49] T. Tieleman and G. Hinton, "Lecture 6.5-rmsprop, coursera: Neural networks for machine learning," *University of Toronto, Technical Report*, 2012.