

# Towards Generalization and Data Efficient Learning of Deep Robotic Grasping

Zhixin Chen, Mengxiang Lin\*, Zhixin Jia, Shibo Jian  
State Key Laboratory of Software Development Environment  
School of Mechanical Engineering and Automation, Beihang University  
Beijing, China  
{zhixinc, linmx, jiazx, jianshibo}@buaa.edu.cn

**Abstract**—Deep reinforcement learning (DRL) has been proven to be a powerful paradigm for learning complex control policy autonomously. Numerous recent applications of DRL in robotic grasping have successfully trained DRL robotic agents end-to-end, mapping visual inputs into control instructions directly, but the amount of training data required may hinder these applications in practice. In this paper, we propose a DRL based robotic visual grasping framework, in which visual perception and control policy are trained separately rather than end-to-end. The visual perception produces physical descriptions of grasped objects and the policy takes use of them to decide optimal actions based on DRL. Benefiting from the explicit representation of objects, the policy is expected to be endowed with more generalization power over new objects and environments. In addition, the policy can be trained in simulation and transferred in real robotic system without any further training. We evaluate our framework in a real world robotic system on a number of robotic grasping tasks, such as semantic grasping, clustered object grasping, moving object grasping. The results show impressive robustness and generalization of our system.

**Index Terms**—Deep reinforcement learning, Visual grasp, Simulation and real-world, Generalization

## I. INTRODUCTION

As one of the most important manipulation skills, vision-based grasping has been studied intensely in the field of robotics for decades. Although robots performing pick and place tasks have been applied successfully in industry, creating autonomous robots for grasping objects in unstructured real-world scenes remains an open question. In recent years, deep reinforcement learning (DRL) has attracted increasing research attention in robotics since its success in video game. Combining with CNNs, DRL maps the features of visual observations directly into control policies by trial-and-error. This provides a general way for robots to learn manipulation skills by using information acquired from cameras [1], [2].

The typical way to train a visual-based DRL agent is in an end-to-end fashion, in which the reward signal of reinforcement learning is used to train both CNNs and policy networks synchronously. However, in order to achieve satisfying performance, large amounts of interaction data are required for training of a DRL agent. For example, to collect enough training data, [3] executed 800,000 robotic grasping attempts in several months with 14 robots, and [4] collected 700 hours

robot grasping data with 50,000 grasping attempts. Moreover, DRL methods try to map raw visual observation into a lower dimensional latent space of control that preserves various types of information about manipulated objects. However, tangled and uninterpretable latent representation restricts the generalization across object and environment and further leads to poor control policy. Most works have to evaluate the trained DRL agents with the similar objects and environments as that of training [5], [2], since the networks trained before have to be fine tuned to adapt to change and hundreds of thousands of robotic manipulation experiences may be needed once more when transferring to a new environment. These limitations will definitely prohibit the use of DRL method in real-world robotic application.

A feasible way to alleviate data requirements is to train DRL agents in simulation, in which the interaction could be speed up easily by using programming techniques, such as multithreading. With this approach, large volumes of experiences can be captured efficiently than that of real world interactions and meanwhile many variants of environments could be constructed for generalization concern. However, there is a huge gap between the simulation and the real-world, which causes the agents trained in simulation to be hardly applied in real world conditions, especially in the context of vision-based robotic manipulation, where illumination changes and varying textures can have significant effects on the quality of the results. Therefore, some DRL based robotic control approaches are only verified in simulation due to difficulties in transferring from simulation to real robots [6], [7]. To alleviate the problem some techniques are proposed to allow for automatic adaptation to the real world environments [8], [9].

In this paper, we propose a DRL-based visual grasping system aiming at improving generalization performance with the least cost of the acquisition of real world experiences. Following the typical visual based DRL paradigm, our framework consists of two major components: a CNN based visual *perception* and a DRL based control *policy*. The perception module extracts features from visual observation (i.e. raw images) and then the features are mapped into the action space by the policy module.

We train the perception and the policy separately instead of

\* Corresponding author

end-to-end. The perception is trained in a supervised setting to produce the semantic and spatial information of the grasped objects. In the meantime, the control policy is trained in a simulation environment where the class and pose of the object to be grasped can be read automatically. Training the policy with the quantitative description of manipulated objects can be beneficial to both generalization and transferability since the information irrelevant for control decision is discarded. In our work, after roughly **30 minutes** of training in simulation, the policy is directly transferred to a real robotic system **without any additional training**. The performance of our system is evaluated on challenging tasks including semantic grasping, clustered object and moving object grasping. The experimental results demonstrate the robustness and generalization of our approach.

## II. RELATED WORK

Simulation environments can provide experiences data much more effective since the simulation could be accelerated by programming. And many robotic manipulation DRL algorithms are verified in simulation environments [1], [10]. Unfortunately, the gap between simulations and real-world makes the the agent trained in simulation can hardly use in physical robot. Many works had tried to bridge reality gap [9], [8]. In [9], the images came from simulations were rendered in randomization, and while the visual perception had seen enough variability over simulations, the real-world images may appear to the model just as another variation. Such randomization made a successful visual grasping model in real-world. [8] unified the visual inputs from simulation and real-world using an adaptation network. The adaptation network was trained to generalize canonical images from randomized images from simulation. And because of the randomization, the trained network could also generalize canonical images from real-world. The generalized canonical images which had been mapped into the same space were used to train the visual grasping DRL agents. Benefited from the adaptation network, DRL agents could be trained in simulation and used in real-world.

Many researches have tried to relieve data inefficiency by improving the efficiency of DRL training process and experiences data generation. Guided policy search (GPS) algorithm [11] converts reinforcement learning into supervised learning, where a trained a local linear controller provided with full state observation (i.g., object poses) served as supervisor. And a global policy parameterized by neural networks derived from supervision. This allows a CNNs policy with 92,000 parameters to be trained in tens of minutes of real-world robotic manipulation time and in test stage the full state is no longer available that the policy trained in a supervised setting could handle several novel, unknown configurations. Another direction to improve sample efficiency is to accelerate model-free reinforcement learning with a learned dynamics models [10]. The learned models can generate synthetic sample data to enrich the agent experiences efficiently which has no need to execute the physical robot, though it needs additional efficient

model learning algorithms [12], [13]. However, the learned model would quite differ from the true dynamics and the induced error would weaken performance of learned policy.

These methods tried to train an optimal policy and visual perception simultaneously in an end-to-end style. However, it can hardly be generalized to different manipulated objects and different execution environments. Since the generalization is relied on the distribution of training data, it requires a huge experience data to achieve usable generalization ability [14]. It is impractical in a robotic grasping task to acquire enough data. An intuitive alternative is to train image representation and reinforcement agent separately [15], [5]. With an auto encoder pretrained by an auxiliary reconstruct loss, the high dimension of image input is embedded into a low dimension, latent space and aggregate useful features before interacting with environment. This way the training of the reinforcement agent networks would be more easily with much less interaction experiences for there is no need to learn the state representation and the training would significantly speed up. However, the latent feature representation has no exact physical meanings and would be lack of interpretability as well the trained policy. From this perspective, meaningful feature representation would significantly improve generalization ability.

## III. FRAMEWORK

We propose a robotic grasping framework based on deep reinforcement learning. Reinforcement learning enables agents (e.g., robots) to learn an optimal policy through interaction with environments by trial-and-error. In doing so, we formulate a robotic grasping problem as a Markov decision process: at time  $t$ , the robot receives the state of target objects and constructs environment state  $s_t$  accordingly. After that, the robot chooses an action  $a_t$  to move itself based on the current policy  $\pi(a_t|s_t)$ . Then the environment transits to a new state  $s_{t+1}$  reacting to  $a_t$  and an immediate reward  $R_t(s_t, a_t, s_{t+1})$  is offered by the environment. The goal of the robot is to find an optimal policy  $\pi^*(a_t|s_t)$  that maximizes the discounted future rewards

$$G_t = \sum_{i=t}^{\infty} \gamma^{i-t} R_i(s_i, a_i, s_{i+1})$$

where  $0 < \gamma < 1$  is the discounted factor which reduces the weight of future rewards.

Similar to recent works [5], [11], our framework is composed of two stages, as shown in Fig.1. Raw RGB images from camera are input into the perception, where object semantic segmentation and pose estimation are made by Mask R-CNN and PCA respectively. The policy is a PPO [16] agent which receives the control quantities of desired objects and decides which action will be taken to execute grasping. The pseudo code of grasping a single object with our framework is presented in Algorithm 1. The details of each component are discussed in the following subsection.

The perception and policy are trained separately. In particular, Mask R-CNN is trained in a supervised way, in which

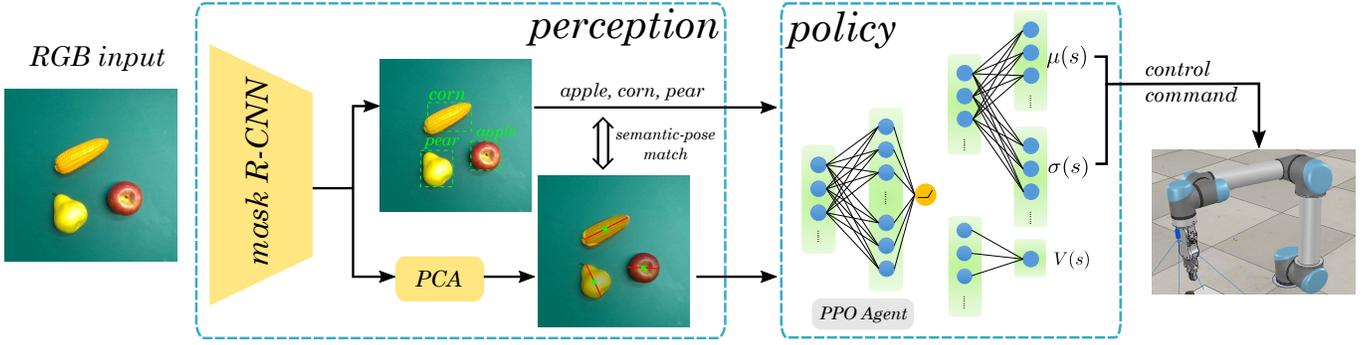


Fig. 1. Our framework for robotic prehensile control.

the labels are constructed manually using a tool LabelMe [17]. There is no training stage for PCA as it is an unsupervised method. The PPO is trained in simulation for fast experience data acquisition. Since the semantic and position information can be read by an interface provided by the simulation environment, the training of PPO could proceed in parallel with that of Mask R-CNN.

---

**Algorithm 1** Separate Perception and Policy

---

**Input:**  $image$

**Perception**

- 1: ▷ detect object class and mask from raw images
- 2:  $mask, class \leftarrow \mathbf{mask-rcnn}(image)$
- 3: ▷ get physical quantities from mask
- 4:  $center, direction \leftarrow \mathbf{PCA}(mask)$

**Policy**

- 5: **repeat**
  - 6:   ▷ decide action via PPO
  - 7:    $action \leftarrow \mathbf{PPO}(center, direction)$
  - 8: **until**  $|action| < \epsilon$
  - 9: ▷ grasp specific object and success check
  - 10:  $success \leftarrow \mathbf{grasp}(class)$
- 

A. Perception

The perception plays a sensor-like role that transforms raw image inputs into physical quantities binding with object semantic information (i.e., object class and its corresponding pose). We should note that this work focuses on a 3DOF grasp [18] given that the workspace of a robot is constrained to a tabletop.

1) *Semantic Segmentation:* To grasp a target object, the robot must know where the target is. To achieve this, we leverage a popular semantic segmentation method Mask R-CNN [19] as the front part of the perception to detect and segment objects from raw images.

Based on Faster R-CNN [20], Mask R-CNN introduces a mask branch at the end of the original network for segmentation tasks. It proceeds in two stages: first, the region proposal

network (RPN) [21] is applied to scan the image to find the area where the target exists; secondly, the object class and its bounding box coordinates are predicted simultaneously. Finally, a pixel-wise binary mask for each bounding box is generated by a fully convolutional network (FCN) indicating whether each pixel in bounding box is the point of the detected object. As a result, masks on the original image exactly cover the areas where the objects exist. The class and mask of an object provide us with a good starting point for pose estimation.

2) *Object Pose:* Since an object mask produced carries information about object pose, learning a DRL policy from masks is in principle possible, as what most DRL based approaches do. In realistic robotic applications however, we can not afford to collect such huge interaction data required by a policy learning algorithm. To avoid this difficulty, we further infer pose for an object instance based on the object mask obtained. In a 3DOF grasp setting, object pose can be represented by a 2-dimensional position coordinates of the object center and the direction of the object. Here, we develop a Principal Component Analysis (PCA) based method to estimate 3D object pose from a pixel-wise binary mask output by Mask R-CNN. In general, PCA is an unsupervised method that could identify the main components of data with largest variances from a big dataset. For our purpose, the center and main direction of a set of pixel points are inferred by using PCA.

The output of Mask R-CNN is an object with its covered mask which contains all pixel points formalized as

$$mask = \{(x_0, y_0), (x_1, y_1), \dots, (x_n, y_n)\}$$

where  $n$  is the number of pixel points in the mask, i.e the number of samples in PCA. Firstly, we calculate the mean point of  $mask$  as the center point of the mask:

$$c = (x, y) = \frac{1}{n} \sum_{i=1}^n (x_i, y_i)$$

Note that, the mean point  $c = (x, y)$  is the geometric center of a mask. After that, all the points in the  $mask$  are subtracted by the  $c$  resulting residual coordinates

$$Res = mask - c$$

Thus, the covariance matrix of  $Res$  and its corresponding eigenvalues and eigenvectors are calculated:

$$\lambda_1, \lambda_2 = \text{eigenvalues}(\text{covMat}(Res))$$

$$\alpha_1, \alpha_2 = \text{eigenvectors}(\text{covMat}(Res))$$

Since the pixel points are in two dimensions, there are totally two eigenvalues and eigenvectors of  $Res$  matrix. The  $\lambda$  and  $\alpha \in \mathbb{R}^{2 \times 1}$  are sorted by the magnitude of eigenvalues in descending order. Finally, the main component with largest variance is calculated:

$$M = Res \cdot \alpha_1 \cdot \alpha_1^\top + c$$

$M \in \mathbb{R}^{n \times 2}$  contains  $n$  points on a straight line. We take two points from  $M$  randomly to construct a straight line.  $\theta$  represents the angle of the straight line respect to the horizontal axis.

Fig.2 shows the results of PCA on a number of objects with various shapes. With the help of a calibrated camera, the position and orientation in pixel coordinates can be mapped into that of a physical coordinate system.

## B. Policy

The policy is a deep reinforcement agent that receives the physical quantities from the perception and decides an optimal action to move the robot. For our framework, we adopt a policy gradient method called Proximal Policy Optimization (PPO) [16] which is favorable for high dimension continuous robotic control problem. PPO significantly improves data efficiency over other policy gradient methods by updating multiple gradient steps with the same trajectory. Moreover, to avoid an increase in data variance, PPO introduces a clipping coefficient which stops the gradient when the difference between the updated and original policy is too large.

1) *State Representation*: We concatenate the output from the perception  $x, y, \theta$  (the pose of a target object) and the robotic configuration  $x_r, y_r, \theta_r$  (the pose of the end effector of a robot) to form the environment state  $s_t$  as the input of the policy network PPO:

$$s_t = (x, y, \theta, x_r, y_r, \theta_r)$$

Through several fully connected layers, PPO finally outputs an action distribution over current state  $s_t$ . Thereafter, the optimal action in the current state  $s_t$  at time step  $t$  is sampled from the action distribution:

$$a_t \sim \pi(a|s_t)$$

The action is a three-dimensional continuous variable instructing the robot's next moving direction and magnitude. The execution of the action would lead to a new environment state  $s_{t+1}$  and an immediate reward  $R_t$  offered by environment.

2) *Reward*: The reward function for learning the policy  $R_t$  is defined as:

$$R_t = \begin{cases} -d_t - 0.1 & \text{away} \\ -d_t + 0.1 & \text{approaching} \\ 1 & \text{grasp success} \\ -1 & \text{grasp failed} \end{cases}$$

where  $d_t$  is the distance between the end effector's current and target position in time step  $t$ .

If  $d_t$  is decreasing compared to the previous  $d_{t-1}$ , the end effector is *approaching* the target and will receive a slightly positive reward addition, and otherwise, it is *away* with a negative reward. In this way, we encourage the end effector to approach and trace the target object as soon as possible. When the policy decides actions bounded in a very small magnitude for several time steps, the policy will decide to execute the grasping, i.e., the end effector moves down on  $z$  coordinate and closes the gripper. A grasping is counted as a success if the gripper is not fully closed. With a larger reward for a successful grasping, the policy could learn the tracing target policy and grasping policy simultaneously.

3) *Training Loss*: PPO is an Actor-Critic style [22] algorithm and typically contains a value function and a policy function. The value function  $V_{s_t}$  which estimate the expected reward from a state  $s_t$  is trained to minimize TD error [23], whose loss function is defined as:

$$L_V = (V(s_t) - (R_t + \gamma V(s_{t+1})))^2$$

where  $\gamma$  is the discounted factor.

The policy function  $\pi(a_t|s_t)$  which decides an optimal action over a state  $s_t$  is trained to maximize a novel surrogate objective [16] according to the value function:

$$L_a = \mathbf{min}(r_t A_t, \text{clip}(r_t, 1 - \epsilon, 1 + \epsilon) A_t)$$

where  $r_t = \frac{\pi(a_t|s_t)}{\pi_{old}(a_t|s_t)}$  is the importance sampling coefficient in which  $\pi_{old}(a_t|s_t)$  is the behavior policy whose parameters are frozen during one update epoch.  $A_t$  is the advantage function [22] which indicates if the reward of current action is above average. And it could be estimated easily by  $A_t = R_t + \gamma V(s_{t+1}) - V(s_t)$  or GAE method [24] according to the value function  $V(s_t)$ . And a *clip* is a function that limits the importance sampling value between  $1 - \epsilon$  and  $1 + \epsilon$  in order to avoid a large step update where  $\epsilon$  is the clipping coefficient which usually equals to 0.2 [16].

Therefore, the final loss function  $L$  becomes

$$L = L_V - L_a$$

and the parameters of network are updated through gradient descent method according to  $L$ .

## IV. EXPERIMENTAL EVALUATION

### A. Implementation

To evaluate our approach, we implemented a visual-based grasping system based on the framework shown in Fig.1. The perception consists of a Mask R-CNN and a PCA procedure, which are executed in pipeline. The input images resized into



Fig. 2. The example results of PCA pose estimation. The black shadows represent the masks produced by Mask R-CNN, and the green points are the center points of the masks, i.e. the position coordinate of the object, and the red lines indicate the main component of the masks, i.e. the orientation of the object in the plane. The range of  $\theta$  is  $[-\frac{\pi}{2}, \frac{\pi}{2}]$ , where the green arrow indicates a positive  $\theta$ , while the red arrow is negative.

600×600 are fed into the Mask R-CNN and a number of object instances covered with their masks are produced. For each mask produced, PCA is invoked to compute its position and orientation as the output of the perception. The implementation of the Mask R-CNN is based on [25]. Instead of using a pre-trained Mask R-CNN model on general objects datasets such as MSCOCO [26], we train the Mask R-CNN on our own dataset considering detection accuracy. To this end, 1000 images of 21 classes of objects are collected and labeled with their mask ground truth manually by the label tool LabelMe [17].

For the policy, three fully-connected layers are stacked together to form a PPO agent. The first layer takes as input a six dimension vector concatenating the object position and the robot position and transforms the input into a 512 dimension latent vector. And the second layer transforms the 512 dimension vector into two streams: a 512 dimension action vector and a 512 dimension value vector. Then one stream is transformed into two 3-dimension vectors, representing the parameters of action distribution  $\mu$  and  $\sigma$ , while another stream is transformed into a scalar representing the value of current environment state. For efficient training of PPO, we

setup a simulation environment in V-REP [27], as shown in Fig.3. Seven classes of objects from a robotic manipulation benchmarks YCB [28] are used for PPO training in simulation, including a detergent bottle, an orange, a round can, a rectangular can, a cup, a pudding box and an electric drill. Since the class and pose of an object in the simulation environment can be obtained directly through software interfaces, PPO could be trained separately, without the help of the perception. The parameters of PPO are learned in a learning rate of  $1e - 5$  using Adam optimization method [29].

Both the training of PPO and Mask R-CNN is done on a PC with a RTX 2080Ti GPU. The average rewards of PPO training in simulation over 5 runs are shown in Fig.4. Very impressive results are obtained after about **30 minutes** training of PPO, indicating by the red arrow in Fig.4. By contrast, the interactions with the same number of episodes would take tens of hours for a real physical robot. The training of Mask R-CNN on our own dataset takes about 10 hours. This training may be not necessary since a pre-trained model on general dataset usually works well in many cases. It is worth noting that all the training above does not require a real-world robot and the trained networks will be transferred into a real-world

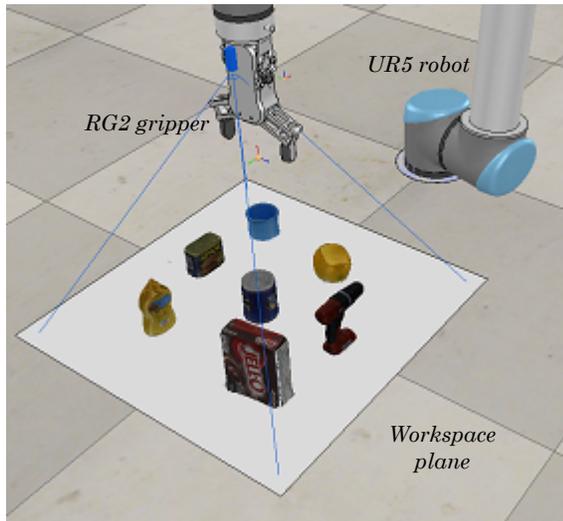


Fig. 3. Simulation environment set up in V-REP [27].

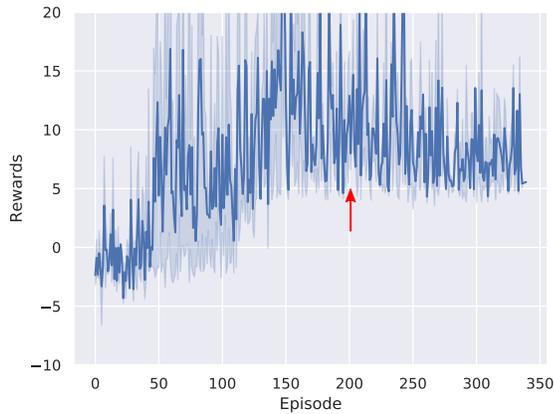


Fig. 4. The process of policy training in simulation. The rewards converged very quickly and the model trained has achieved a good performance after about 30 minutes. The policy model used in all experiments is trained for 200 episodes, as indicated by the red arrow.

robotic grasping system directly.

### B. Real-world Evaluation

The overall goal of our evaluation is to determine whether the trained networks can enable a real world robot to perform various grasping tasks without any further training. To this end, a number of grasp tasks commonly used in our daily life are designed to evaluate the ability to perform grasping skills and generalization over objects and situations. We use an industrial UR5 robot arm with an RG2 gripper to achieve two-finger rigid grip. A RealSense camera [30] is located 100 cm above the work surface, producing RGB images for input. A laptop with a RTX 2080 GPU acceleration is used for real-time robotic control and communication with UR5 via TCP/IP protocol. The experimental hardware platform is shown in Fig.5. It is

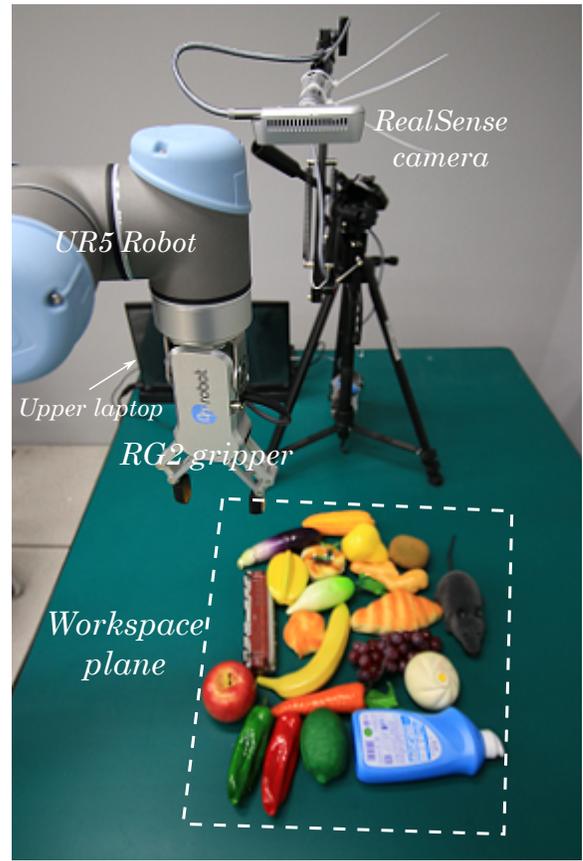


Fig. 5. The hardware setup of our system.

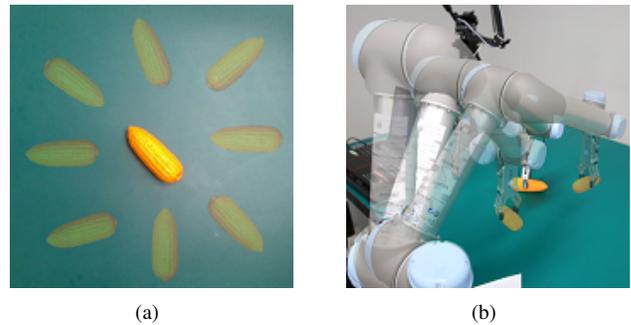


Fig. 6. Sim-to-Real transferring test of the trained policy. (a) Picking up a corn in various position and orientation. (b) Picking up a corn with different robotic configurations.

worth note that the objects used in the experiments are totally different from that of PPO training in simulation.

1) *Sim-to-Real Transfer*: As mentioned before, the trained networks including Mask R-CNN and PPO are transferred into our robotic grasping system without any further training. We first examine the behavior of the system in a controlled manner. As shown in Fig.6, a target object (a corn) is placed on the work surface in various positions and orientations. The robot grasps the target successfully for 20 randomly chosen object locations.

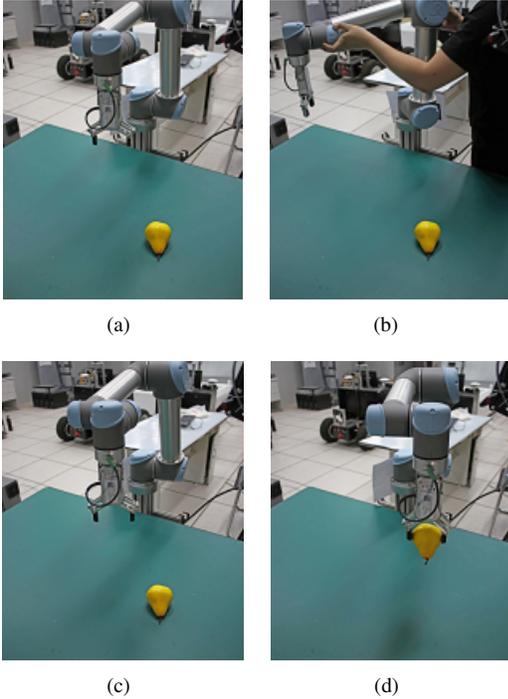


Fig. 7. Robustness test of the trained policy. (a) The robot starts to perform a task. (b) We forced it manually to an unseen configuration in training. (c) The robot finds a proper path back to the right way. (d) The robot grasps the target successfully.

TABLE I  
THE RESULTS OF MULTIPLE OBJECTS REMOVE UP.

scenarios	grasp success	remove completion
brown	100%(112/112)	100%(10/10)
green	100%(120/120)	100%(10/10)
dense	93.7%(104/111)	80%(8/10)

Furthermore, in order to test the robustness of the control policy, we manually introduce external disturbances. As shown in Fig.7, the control policy could find its correct trajectory again and grasp the target successfully after a sudden change on the robot configuration during the robot’s execution, exhibiting good stability and robustness.

2) *Multi-object Grasping*: Multi-object grasping is a common task used to measure the performance for a vision based robotic grasping system. In our test setting, 10-13 objects are placed randomly on the table and the UR5 robot is requested to pick up all objects sequentially and then put them out of the workspace. In addition, the background color of the work surface is shifted from white into brown or green. Two example test settings in different backgrounds are shown in Fig.8. A grasp is successful if an object is grasped and threw aside, while a remove completion means no objects are left on the table. We perform 10 tests for each background and grasp success rate and remove completion rate are presented in Table I.

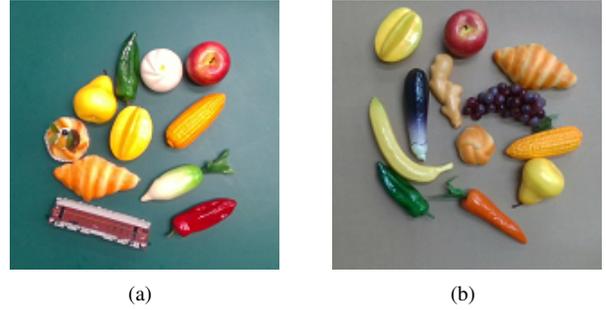


Fig. 8. Multiple object removing up. (a) 11 objects in a green background. (b) 12 objects in a brown background.

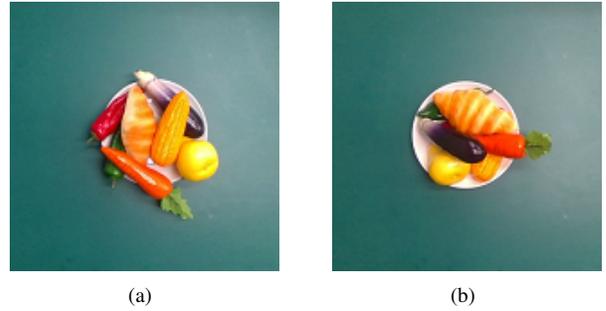


Fig. 9. Dense clustered object removing up. (a) Objects blocked each other but every object is visible. (b) The red pepper is completely invisible.

3) *Clustered object Grasping*: A challenge task in robotic manipulation is to grasp objects clustered closely together. As shown in Fig.9, in a cluster scenario, the objects would block each other and some objects may be completely invisible. In such a task, the order of manipulations really matters if we want to remove all the objects in sequence. To decide the ordering of picking up, we define a mask ratio  $r$  for each object recognized as follows:

$$r = \frac{m}{M}$$

where  $m$  is the recognized mask with possible occlusion and  $M$  is the full mask of the object which is pre-determined. The larger the ratio, the more likely the object is to be picked up firstly as it is less occluded by others. For dense cluster scenarios, we perform 10 tests and grasp success rate and remove completion rate are presented in Table I. The failure cases occur due to misidentifications by Mask R-CNN because of partially visible objects.

4) *Semantic Grasping*: In a semantic grasping task, a robot is instructed to grasp a specified object among a set of candidates. The capability of semantic grasping is essential to allow autonomous robots to perform manipulations in an unstructured environment. Benefiting from the power of Mask R-CNN to detect objects, our system first identifies the class of an object before deciding how to pick up. Similar to the experiment setting in multi-object grasping, 10-13 objects are randomly placed on the table for each trial. For each run, we randomly specify one object and simply count the number of

successful grasps. We perform five trials and the success rate achieves **100% (60/60)**.

5) *Moving Object Grasping*: Grasping a moving object is still a challenging task in visual-based robotic manipulation [31]. Our learning based approach provides a promising way to approach this challenge. To demonstrate the effectiveness of our approach, a case study is conducted in a scenario where a small fake mouse is moving and the robot is ordered to pick up the mouse in motion.

To pick up the moving mouse, our robot needs to be able to track the target continuously and decide to execute a grasping once the action outputted by the control policy is smaller than the preset threshold. In order to further reduce the time between making a grasp decision and closing the gripper, we add a fixed movement in  $z$  direction simultaneously with the movement in the  $x - y$  plane, instead of moving down in  $z$  direction after the decision making. This minor modification significantly improves the successful rate of picking up the mouse in our experiments. The example pictures of robot's execution on picking up a moving mouse are shown in Fig.10. However, due to the computation cost of the system and communication delay between the laptop and the robot, the delay time in our current implementation is about 200ms, which limits the speed of moving objects in our experiments.

## V. CONCLUSION AND FUTURE WORK

We presented a robotic grasping approach that combines visual perception and a DRL based control policy. Comparing with other alternatives, training on a real robot is avoided by decoupling the control from visual perception with the help of a physical representation of objects, which makes them easier to be trained. Moreover, the policy trained in simulation could be transferred to a real system without any further training. Real world experiments on UR5 demonstrate the robustness and generalization over a wide variation in challenging grasping tasks. However, in this work, we only consider 3DOF grasping in which objects are placed on a table and the grasping height is fixed. In future work, we would like to extend this work to a 6DOF grasping. To do so, it will be important to investigate the pose of gripper in 3D shape perception.

## REFERENCES

- [1] S. Gu, E. Holly, T. P. Lillicrap, and S. Levine, "Deep reinforcement learning for robotic manipulation with asynchronous off-policy updates," *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 3389–3396, 2016.
- [2] H. Zhu, A. Gupta, A. Rajeswaran, S. Levine, and V. Kumar, "Dexterous manipulation with deep reinforcement learning: Efficient, general, and low-cost," *2019 International Conference on Robotics and Automation (ICRA)*, pp. 3651–3657, 2018.
- [3] 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," *I. J. Robotics Res.*, vol. 37, pp. 421–436, 2016.
- [4] L. Pinto and A. Gupta, "Supersizing self-supervision: Learning to grasp from 50k tries and 700 robot hours," *2016 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 3406–3413, 2015.
- [5] C. Finn, X. Y. Tan, Y. Duan, T. Darrell, S. Levine, and P. Abbeel, "Deep spatial autoencoders for visuomotor learning," *2016 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 512–519, 2015.
- [6] C. Finn, T. Yu, J. Fu, P. Abbeel, and S. Levine, "Generalizing skills with semi-supervised reinforcement learning," *ArXiv*, vol. abs/1612.00429, 2016.
- [7] D. Ghosh, A. Singh, A. Rajeswaran, V. Kumar, and S. Levine, "Divide-and-conquer reinforcement learning," *ArXiv*, vol. abs/1711.09874, 2017.
- [8] 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 *CVPR*, 2018.
- [9] J. Tobin, R. H. Fong, A. Ray, J. Schneider, W. Zaremba, and P. Abbeel, "Domain randomization for transferring deep neural networks from simulation to the real world," *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 23–30, 2017.
- [10] S. Gu, T. P. Lillicrap, I. Sutskever, and S. Levine, "Continuous deep q-learning with model-based acceleration," *ArXiv*, vol. abs/1603.00748, 2016.
- [11] S. Levine, C. Finn, T. Darrell, and P. Abbeel, "End-to-end training of deep visuomotor policies," *J. Mach. Learn. Res.*, vol. 17, pp. 39:1–39:40, 2015.
- [12] S. Levine and P. Abbeel, "Learning neural network policies with guided policy search under unknown dynamics," in *NIPS*, 2014.
- [13] M. P. Deisenroth and C. E. Rasmussen, "Pilco: A model-based and data-efficient approach to policy search," in *ICML*, 2011.
- [14] C. Devin, P. Abbeel, T. Darrell, and S. Levine, "Deep object-centric representations for generalizable robot learning," *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 7111–7118, 2017.
- [15] S. Lange, M. A. Riedmiller, and A. Voigtländer, "Autonomous reinforcement learning on raw visual input data in a real world application," *The 2012 International Joint Conference on Neural Networks (IJCNN)*, pp. 1–8, 2012.
- [16] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *ArXiv*, vol. abs/1707.06347, 2017.
- [17] B. C. Russell, A. Torralba, K. P. Murphy, and W. T. Freeman, "Labelme: A database and web-based tool for image annotation," *International Journal of Computer Vision*, vol. 77, pp. 157–173, 2005.
- [18] A. ten Pas, M. Gualtieri, K. Saenko, and R. Platt, "Grasp pose detection in point clouds," *The International Journal of Robotics Research*, vol. 36, pp. 1455 – 1473, 2017.
- [19] K. He, G. Gkioxari, P. Dollár, and R. B. Girshick, "Mask r-cnn," *2017 IEEE International Conference on Computer Vision (ICCV)*, pp. 2980–2988, 2017.
- [20] R. B. Girshick, "Fast r-cnn," *2015 IEEE International Conference on Computer Vision (ICCV)*, pp. 1440–1448, 2015.
- [21] S. Ren, K. He, R. B. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, pp. 1137–1149, 2015.
- [22] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. P. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu, "Asynchronous methods for deep reinforcement learning," in *ICML*, 2016.
- [23] R. S. Sutton and A. G. Barto, "Reinforcement learning: An introduction," *IEEE Transactions on Neural Networks*, vol. 16, pp. 285–286, 1988.
- [24] J. Schulman, P. Moritz, S. Levine, M. I. Jordan, and P. Abbeel, "High-dimensional continuous control using generalized advantage estimation," *CoRR*, vol. abs/1506.02438, 2015.
- [25] W. Abdulla, "Mask r-cnn for object detection and instance segmentation on keras and tensorflow," [https://github.com/matterport/Mask\\_RCNN](https://github.com/matterport/Mask_RCNN), 2017.
- [26] T.-Y. Lin, M. Maire, S. J. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *ECCV*, 2014.
- [27] E. Rohmer, S. P. N. Singh, and M. Freese, "V-rep: A versatile and scalable robot simulation framework," *2013 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pp. 1321–1326, 2013.
- [28] B. alli, A. Singh, A. Walsman, S. S. Srinivasa, P. Abbeel, and A. M. Dollar, "The ycb object and model set: Towards common benchmarks for manipulation research," *2015 International Conference on Advanced Robotics (ICAR)*, pp. 510–517, 2015.
- [29] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *CoRR*, vol. abs/1412.6980, 2014.
- [30] L. Keselman, J. I. Woodfill, A. Grunnet-Jepsen, and A. Bhowmik, "Intel realsense stereoscopic depth cameras," *ArXiv*, vol. abs/1705.05548, 2017.

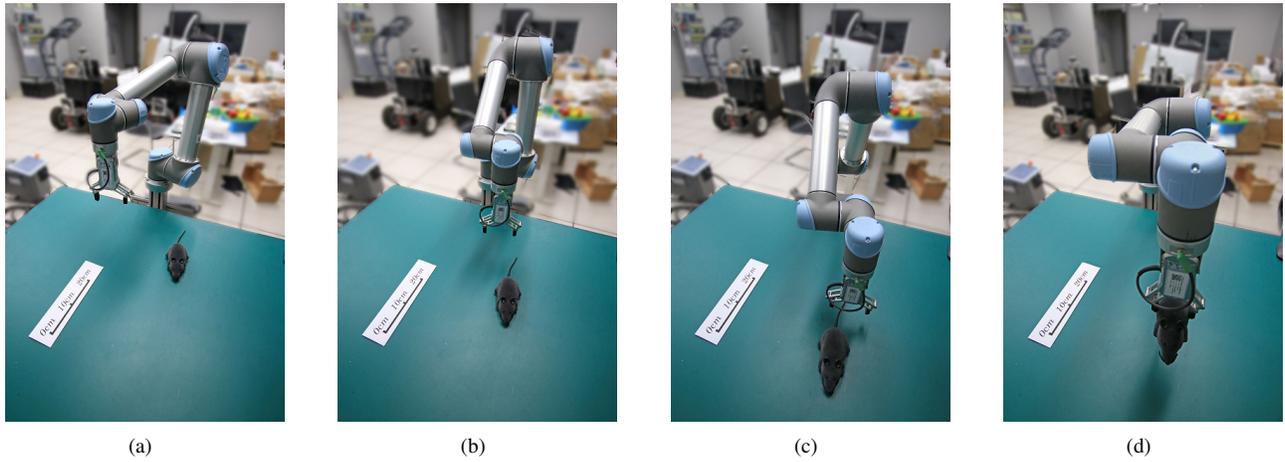


Fig. 10. Moving object grasping. (a) The fake mouse is in the initial position. (b) The mouse is moving forward and the robot is changing its control strategy. (c) The mouse keeps moving and the robot keeps tracing it while gripper moving down simultaneously. (d) The robot grasps the mouse successfully.

- [31] A. Menon, B. J. Cohen, and M. Likhachev, "Motion planning for smooth pickup of moving objects," *2014 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 453–460, 2014.