

# Affordance Learning from Play for Sample-Efficient Policy Learning

Jessica Borja-Diaz\*, Oier Mees\*, Gabriel Kalweit, Lukas Hermann, Joschka Boedecker, Wolfram Burgard

<http://vapo.cs.uni-freiburg.de>

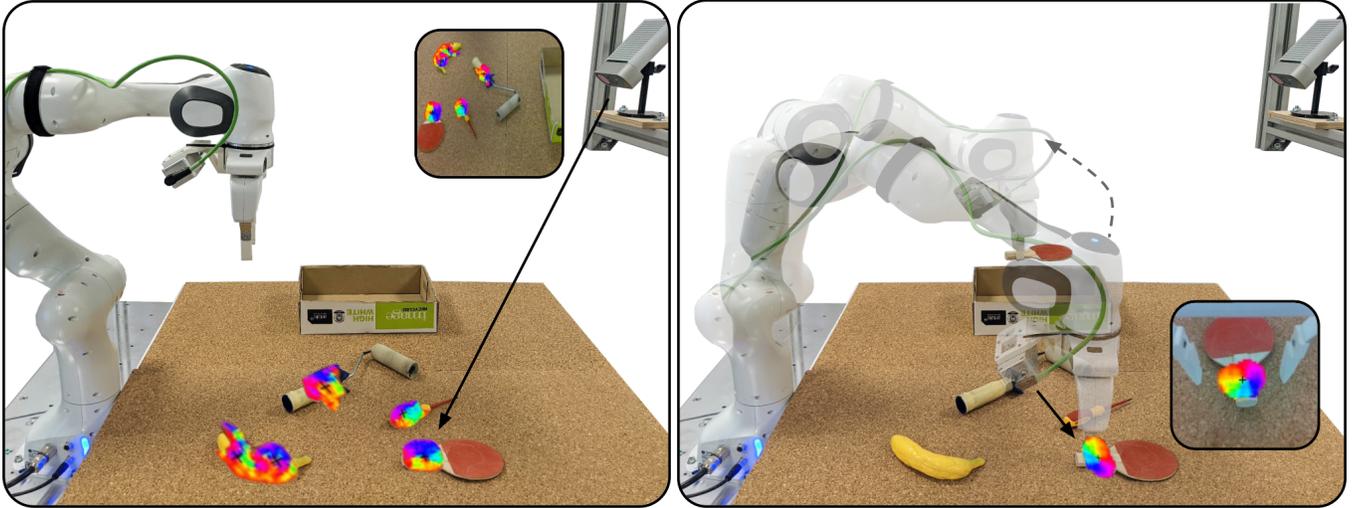


Fig. 1: Real world setup for a tidy up task: our self-supervised visual affordance model guides the robot to the vicinity of actionable regions in the environment with a model-based policy. Once inside this area, we switch to a local reinforcement learning policy, in which we embed our affordance model to favor the same object regions favored by people and to boost sample-efficiency.

**Abstract**—Robots operating in human-centered environments should have the ability to understand how objects function: *what* can be done with each object, *where* this interaction may occur, and *how* the object is used to achieve a goal. To this end, we propose a novel approach that extracts a self-supervised visual affordance model from human teleoperated play data and leverages it to enable efficient policy learning and motion planning. We combine model-based planning with model-free deep reinforcement learning (RL) to learn policies that favor the same object regions favored by people, while requiring minimal robot interactions with the environment. We evaluate our algorithm, Visual Affordance-guided Policy Optimization (VAPO), with both diverse simulation manipulation tasks and real world robot tidy-up experiments to demonstrate the effectiveness of our affordance-guided policies. We find that our policies train  $4\times$  faster than the baselines and generalize better to novel objects because our visual affordance model can anticipate their affordance regions.

## I. INTRODUCTION

Humans have the ability to effortlessly recognize and infer functionalities of objects despite their large variation in appearance and shape. For example, we understand that we need to pull the handle of a drawer to open it or grasp a knife by the handle to use it. This capacity to focus on the most relevant behaviors in a given situation enables efficient decision making by limiting the choices of action that are even considered. Gibson’s theory of affordances [1] provides a way to reason about object function. It suggests that objects have action possibilities, e.g., a mug is “graspable” and a

door is “openable” and has been extensively studied in both the robotics and the computer vision communities [2].

However, the abstract notion of “what actions are possible?” addressed by current affordance learning methods is limited. A robot needs to know *where* are actionable regions in an environment, the specific points on the object that need to be manipulated for a successful interaction, *what* it can achieve with it and *how* the object is used to achieve a goal. Current affordance learning methods have two major problems. First, they are limited by requiring heavy supervision in the form of manually annotated segmentation masks [3]–[6] or expensive interactive exploration [7], [8], restricting their scalability and applicability in practical robotics scenarios. Second, current affordance-augmented robotic systems are limited in the complexity of the actions they model by relying often on predefined action templates [7]–[10]. Together, these limitations naturally restrict the scope of affordance learning systems to a narrow set of objects and robotics applications.

In light of these issues, we propose a method for sample-efficient policy learning of complex manipulation tasks that is guided by a self-supervised visual affordance model. Therefore, we call our algorithm Visual Affordance-guided Policy Optimization (VAPO). Towards overcoming the issues of expensive manual supervision and exploration, we propose to learn affordances that are *grounded* in real human behavior from teleoperated play data [11]. Play data is not random, but rather structured by human knowledge of object affordances (e.g., if people see a drawer in a scene, they tend to open it). Moreover, affordances discovered from unlabeled play are *functional affordances*, priming a robot to approach an object

\*Equal contribution. All authors are with the University of Freiburg, Germany. This work has been supported partly by the German Federal Ministry of Education and Research under contract 01IS18040B-OML

the way a human would. On the other hand, teleoperated play data does not bear the risk of the correspondence problem as opposed to recordings directly from human demonstrations. We hence leverage this visual affordance model to guide a robot to perform complex manipulation tasks. Aside from accelerating learning, a critical advantage of imbuing robots with an object-centric visual affordance prior is generalization: the learned policy generalizes to unseen object instances because our visual affordance model can anticipate their affordance regions.

Our approach decomposes object manipulation into a sample-efficient combination of model-based planning and model-free reinforcement learning, inspired by a recent line of work that aims to combine classical motion planning with machine learning [12]–[14]. Concretely, we first predict object affordances and drive the end-effector from free-space to the vicinity of the afforded region with a model-based method. Once inside this area, the model cannot be trusted and we switch to a RL policy in which the agent is rewarded for interacting with the afforded regions. This way, the local policy has a “human prior” for how to approach an object, but is free to discover its exact grasping strategy. Our self-supervised visual affordance model is leveraged twice to boost sample-efficiency: 1) driving the model-based planner to the vicinity of afforded regions, 2) guiding a local grasping RL policy to favor the same object regions favored by people. Standard model-free RL faces a number of challenges, since the policy must solve two problems: representation learning and task learning from high-dimensional raw observations in a single end-to-end training procedure. As in practice solving *both* problems together is difficult, embedding our visual affordance model within a reinforcement learning loop alleviates the representation learning challenge. The interplay between model-based and model-free policies allows for a sample-efficient division of the robot control learning, without assuming a predefined set of manipulation primitives, 3D object shapes or a tracking system.

## II. RELATED WORK

**Predicting Semantic Representations** To successfully interact with a 3D object, a robot must be able to “understand” it given some perceptual observation. There exists a large body of work in the computer vision community targeting such an understanding in the form of different semantic labels. For example, predicting category labels [15], or more fine-grained output such as semantic keypoints [16], part segmentations [17] or afforded spatial relations [18] can arguably yield more actionable representations e.g. allowing one to infer where “handles” are. However, merely obtaining such semantic labels is clearly not sufficient on its own, a robot must also understand *what* needs to be done and *how* the object is used to achieve a goal.

**Acting with Model-based Planning** Towards obtaining useful information for how to act, some methods aim for representations that can be leveraged by classical robotics techniques. In particular, traditional analytical approaches use knowledge of the 3D object pose [19], [20], shape [15],

[21], gripper configuration, friction coefficients, etc. to determine optimal action trajectories. However, model-based methods rely on an accurate model of the environment and they normally do not handle perception errors and physical interactions naturally [22], limiting their reliability. Our approach uses model-based planning to guide the robot to the vicinity of detected affordance regions and switches then to a local RL policy.

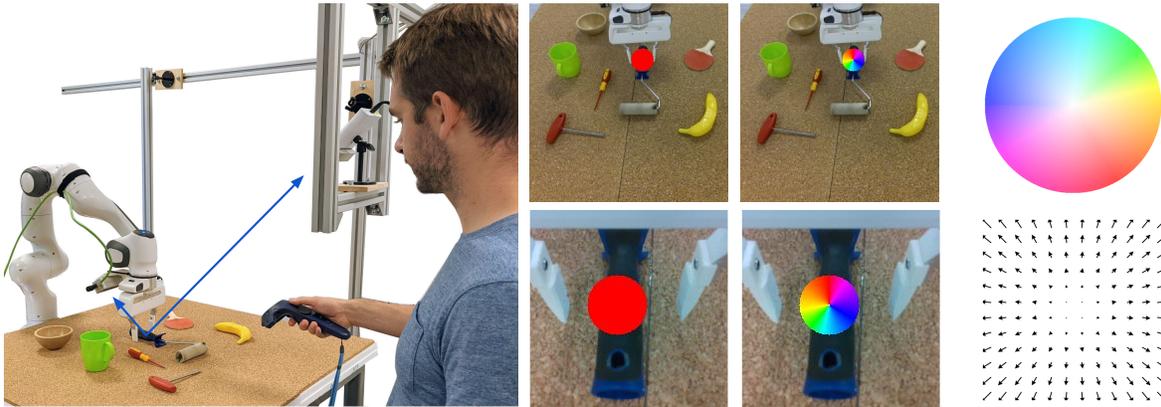
**Reinforcement Learning Grasping** RL models offer a counterpoint to the planning paradigm. Instead of breaking the task into two steps, static grasp synthesis followed by motion planning, it can operate directly from raw sensory inputs in closed-loop feedback control, which are not subject to estimation errors [23], [24]. Unlike model-based methods, RL methods do not require a detailed description of the environment and the task, but rather require access to interaction with the environment and to a reward function. Such binary rewards are easy to describe, but unfortunately they render RL methods extremely sample-inefficient and brittle. Although there have been promising advances in learning data-driven reward functions [25]–[27], for most complex problems of interest, learning RL policies from scratch remains intractable. In contrast, we inject an object-centric visual affordance prior extracted from human teleoperated play data to boost sample efficiency.

**Visual Affordances** Most closely related to our approach is the line of work where visual affordances are learned for object manipulation [5], [6], [28]–[30]. Traditionally, visual affordance learning methods are limited by their requirement of manually drawn segmentation masks or keypoints [3]–[6] and some leverage additional sensing, such as force gloves [31]. Recently, there has been a shift to explore other forms of supervision such as videos [32], a robot’s gripper grasp success/failure [10], [28] or thermal image contact data [30]. In contrast, we leverage a self-supervised signal of a robot’s gripper opening and closing during human teleoperation to learn image-based functional affordances.

## III. APPROACH

The main incentive of our method is to learn sample-efficient policies of complex manipulation tasks that are guided by a self-supervised visual affordance model. Our approach consists of three steps. First, we train a network to discover and learn object affordances in unlabeled play data (Sec. III-A). Second, we divide the space into regions where a model-based policy is reliable and regions where it may have limitations handling perception errors or physical interactions. We leverage the learned affordance model to drive the end-effector from free-space to the vicinity of the afforded region with a model-based policy  $\pi_{mod}$  (Sec. III-B). Third, once inside this area we switch to a local reinforcement learning policy  $\pi_{rl}$ , in which we embed our affordance model to favor the same object regions favored by people and to boost sample-efficiency (Sec. III-C). Thus, our final policy is defined as a mixture:

$$\pi(a|s) = (1 - \alpha(s)) \cdot \pi_{mod}(a|s) + \alpha(s) \cdot \pi_{rl}(a|s), \quad (1)$$



**Fig. 2:** Visualization of our self-supervised object affordance labelling. We leverage a self-supervised signal of a robot’s gripper opening and closing during human teleoperation to project the 3D tool-center-point into the static and gripper cameras. We label the neighboring pixels within a radius around the afforded region with a binary segmentation mask and direction vectors from each pixel towards the affordance region center. On the right we show the color code used to interpret the direction vectors.

where  $\alpha(s) \in [0, 1]$ . We use an estimate of the normalized distance between the robot’s gripper and the affordance region  $\alpha(s)$  to switch between the policies. An overview of the system is given in Figure 1.

#### A. Learning Visual Affordances from Play

Our key insight is to learn about object interactions from play data by leveraging a self-supervised signal of a robot’s gripper opening and closing during human teleoperation, as shown in Figure 2. In this way, without explicit manual segmentation labels, we learn to anticipate not only *where* are regions that afford human-object interactions, but also a powerful prior on *how* humans approach those objects. The only assumption our method makes is an existing robot-camera calibration. We decouple the affordance prediction task into different components.

First, the affordance model  $\mathcal{F}_a$  learns to transform an image  $I$  into a binary segmentation map  $A \in \mathbb{R}^{H \times W}$ , indicating regions that afford an interaction. Second, it estimates 2D pixel coordinates of the affordance region centers by predicting a vector from each affordance pixel towards the center. Estimating the center points of the afforded regions is a key component in order to disambiguate affordances from multiple objects in a scene. Clearly, play data showing people naturally interacting with objects partially reveals the afforded regions in an environment. Thus, in order to discover affordances in unlabeled data the gripper action is used as a heuristic to detect human-object interactions.

Intuitively, if the gripper closes during play, it is indicative of a possible interaction that will start at that position. Thus, we can project the gripper’s 3D point  $p_{gripper}^t$  to a camera image pixel  $u_{gripper}^t$  and label the pixels within a radius  $r$  for the past  $n$  frames as an afforded region. Similarly, if the gripper transitions from being closed to open, it means that an interaction with an object ended at the 3D position  $p_i$ . This allows us to discover a set of interaction points throughout time  $P^k = (p_1, p_2, \dots, p_k)$ , which represent the world coordinates of where interactions have occurred until

timestep  $k$ . To get the full set of interaction locations for a timestep  $t$  we consider the 3D positions from where a grasp will occur and where an interaction has been previously occurred until  $t$ . Finally, each 3D point is projected to a camera image pixel to create the affordance mask label by marking neighboring pixels. The pixel coordinates of the projected points are used as the affordance region centers.

In order to disambiguate affordances from multiple objects in a scene, we let the network estimate 2D pixel coordinates of the affordance region centers by predicting a vector from each affordance pixel towards the center  $V \in \mathbb{R}^{H \times W \times 2}$ . We construct these labels by calculating the displacement from each pixel of the affordance mask to the corresponding projected center. The background pixels are pointed towards a fixed position to avoid false positives.

One limitation of the proposed heuristic is that it assumes users interacting with the environment during play will only close the gripper to perform meaningful interactions. To avoid erroneous labeling due to closing/opening the gripper in free-space without object-interaction, we introduce an additional check by requiring the gripper-width to stay for  $\Delta t$  timesteps in a range between opened and closed.

To train the full affordance model  $F_a$  we apply two different loss functions. For the affordance segmentation  $A$  loss, we use a weighted sum between a cross entropy  $\ell_{ce}$  and a dice loss  $\ell_{dice}$  to account for class imbalance. Similar to Xie *et al.* [33], for the direction prediction we optimize a weighted cosine similarity loss given by:

$$\ell_{dir} = \sum_{i \in \mathcal{O}} \alpha_i (1 - V_i^T \bar{V}_i) + \frac{\lambda_b}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \left( 1 - V_i^T \begin{bmatrix} 0 \\ 1 \end{bmatrix} \right)$$

Where  $V_i, \bar{V}_i$  are the predicted and ground truth unit directions of pixel  $i$  respectively.  $\mathcal{B}, \mathcal{O}$  are the sets of pixels belonging to the background and affordance region classes. The total loss for the affordance model is given by  $w_{ce} \ell_{ce} + w_{dice} \ell_{dice} + w_{dir} \ell_{dir}$ .

### B. From Model-Based to Reinforcement Learning Workspace

Classical motion planning algorithms have difficulty in the presence of stochastic dynamics and high-dimensional systems. RL methods on the other hand offers a promising solution for its ability to learn general policies that can handle complex interactions and high-dimensional observations. However, for most complex problems of interest, learning RL policies from scratch remains intractable. Inspired by recent works that aim to combine both type of controllers [12]–[14], we divide the space into regions where a model-based policy is reliable and regions where it may have limitations handling perception errors or physical interactions.

Concretely, we predict affordances and the corresponding region centers using a static camera image. Given this information of *where* are regions that afford human-object interactions, we localize a chosen pixel region center in 3D and drive the end-effector from free-space to the vicinity of the afforded region with a model-based policy  $\pi_{mod}$  and hand control over to the model-free policy  $\pi_{rl}$ . We use an estimate of the distance between the robot’s gripper and the predicted affordance region center to switch between the policies. Restricting the area where the RL policy is active to the vicinity of regions that afford human-object interactions has the advantage that it makes it more sample-efficient. Besides, this division of labour allows to learn local RL policies by switching to a gripper camera, improving generalization across different locations.

### C. Affordance-guided Reinforcement Learning Grasping

Once the model-based policy  $\pi_{mod}$  has brought the end-effector to the vicinity of a region that affords human-object interactions, we switch to a local gripper-camera based RL policy which we augment with an object-centric visual affordance prior to boost sample efficiency.

**Problem Formulation:** We consider the standard Markov decision process (MDP)  $M = (\mathcal{S}, \mathcal{A}, \mathcal{T}, r, \mu_0, \gamma)$ , where  $\mathcal{S}$  and  $\mathcal{A}$  denote the state space and action space respectively.  $\mathcal{T}(s'|s, a)$  is the probability of transitioning from state  $s$  to state  $s'$  when applying action  $a$ . The actions are drawn from a probability distribution over actions  $\pi(a|s)$  referred to as the agent’s policy.  $r(s, a)$  is the reward received by an agent for executing action  $a$  in state  $s$ ,  $\mu_0$  the initial state distribution, and  $\gamma \in (0, 1)$  the discount factor prioritizing long- versus short-term reward. The goal in RL is to optimize a policy  $\pi(a|s)$  that maximizes the expected discounted return  $\mathbb{E}_{\pi, \mu_0, \mathcal{T}} [\sum_{t=0}^{\infty} \gamma^t r(s, a)]$ .

**Observation Space:** The observation space is composed of two parts: 1) the proprioceptive state including the 3D world coordinates of the end effector, the orientation euler angles and the gripper width. 2) The visual inputs consisting of the current RGB-D image observed by the gripper camera and the binary affordance mask predicted from the corresponding affordance model.

**Action Space:** We use a 7-DOF Franka Emika Panda robot with a parallel gripper both in simulation and in the real world. The action space consists of delta XYZ position, delta euler angles and the binary gripper action.

**Reward:** The reward function should not only signal a successful object interaction, but also guide the exploration process to focus on actionable object regions. To realize this, we leverage the visual affordance model to guide the agent to get close to the affordance centers. This way, the local policy has a “human prior” for how to approach an object, but is free to discover its exact grasping strategy. Given the detected affordance center and the fact that the RL policy only acts locally within a neighborhood, we normalize the euclidean distance between the end effector and the affordance center to create a positive reward  $R_{aff}$  which increases as we get closer to the detected center. Additionally if the agent goes outside the neighborhood, it receives a negative reward  $R_{out}$  and if it successfully manipulates an object it receives a positive reward of  $R_{succ}$ . Our total reward function is:

$$r(s, a) = \lambda_1 R_{succ} + \lambda_2 R_{aff} + \lambda_3 R_{out} \quad (2)$$

## IV. IMPLEMENTATION DETAILS

**Teleoperated play data** During the unscripted teleoperated interactions we record images from two cameras: a static camera that captures the global scene, and a camera mounted on the robot’s gripper. The static camera image has a resolution of  $200 \times 200$  and the gripper camera uses a resolution of  $64 \times 64$ . We label the images of the static camera with a radius of  $r = 10$  pixels around the projected center and the the gripper camera images with  $r = 25$ .

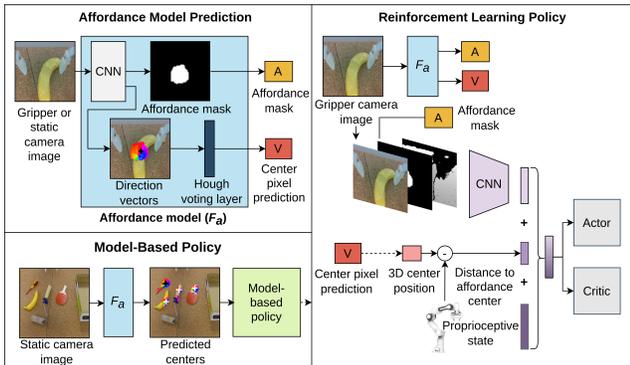
**Affordance model** We use a U-net [34] architecture followed by two parallel branches of convolutional layers that produce the affordance mask and center directions. Similar to Xiang *et al.* [35], we use a Hough voting layer to predict the 2D object centers during inference. The Hough voting layer takes the affordance mask and the direction vectors as input to compute a score for each pixel, indicating its likelihood of being an affordance region center. The location with the maximum score is selected as the object center.

We define a two-stage affordance detection by training separate models for the two cameras as shown in Figure 3. One model is trained with images from a static camera and predicts a spatial interaction hotspots map, indicating actionable regions. Similarly, we train an affordance model with images from a gripper camera, which gives a finer-grained spatial interaction map about where humans tend to interact with each object.

The affordance model should give insight into which parts of an object are relevant for its use. As this is dependent on the shape of the objects rather than the color, we would like the affordance model to be invariant to different colors. For this reason, the images are converted to grayscale before being fed to the networks. Both affordance models are trained with stochastic gradient descent with a learning rate of  $1e-5$  and a batch size of 256. The loss weights are set to  $w_{ce} = 1$ ,  $w_{dice} = 5$ ,  $w_{dir} = 2.5$ .

### Affordance-guided Reinforcement Learning

We train the policy using Soft Actor-Critic [36]. We concatenate the RGB-D images with the inferred affordance mask and pass it through a convolutional neural network



**Fig. 3:** Overview of the full approach. The affordance model takes an image from either camera as input to predict object affordance masks and center pixel predictions (top left). The static camera affordances are used to select a position that the model-based policy will move towards (bottom left). We then switch to a RL policy which takes as input the the predictions of the gripper camera affordance, the robot’s proprioception, the distance to the predicted center, and the current RGB-D image (right).

(CNN) as depicted in Figure 3. The CNN is composed by three convolutional layers with kernel size [8,4,3] respectively and one linear layer to obtain a feature representation of size 16. Then we concatenate the obtained representation to the robot state and the distance to the affordance center. Finally this is passed through four fully connected layers. The critic and actor are implemented following the same architecture without weight-sharing. For the simulation experiments, we train a single policy for all the objects with an episode length of 100 steps during 400K episode steps. This amounts to 30hrs of learning experience. We train for 3 seeds initializations. In the reward function we set  $\lambda_1 = \lambda_2 = \lambda_3 = 1$  and the rewards  $R_{succ} = 200$ ,  $R_{out} = -1$ .

## V. EXPERIMENTAL RESULTS

In this section we seek to answer the following questions: how does our method compare to the baseline policies in terms of sample efficiency and task completion? And, is the proposed approach applicable to a real world tidy-up task?

### A. Experimental Setup

We evaluate our method with both diverse simulation manipulation tasks and real world robot tidy-up experiments.

1) *Simulation:* We evaluate two tasks in simulation: a grasping task and a drawer opening task. The grasping task consists on lifting different objects in a PyBullet simulated environment. The policy is trained over 15 different objects with varying degrees of complexity, such as hammers, knives and power drills, as shown in Figure 4. After the policy executes a close-gripper action, the gripper attempts to lift the object and waits in the air for two seconds. If the object is still in the gripper at the end of this time, we define the grasp as being successful.

VAPO is not exclusive to a grasping task. To show this, we train a policy to open a drawer as shown in Figure 5. Every episode consists of the drawer on a closed position



**(a)** Seen objects during affordance model training. **(b)** Unseen objects during affordance model training.

**Fig. 4:** Objects used in simulation grasping experiments. The objects propose different challenges as some are small or must be grasped in a specific manner, e.g. grasping the frying pan requires to use the handle to be successfully lifted.

and the robot in a neutral position. The episode is deemed successful if the robot opens the drawer at least 15cm.

To train the affordance models we teleoperate the robot using a virtual reality (VR) controller to collect unscripted play data. We gather two hours of human interaction which amounts to  $\sim 100K$  images for each environment to train the static camera and gripper camera affordance models.

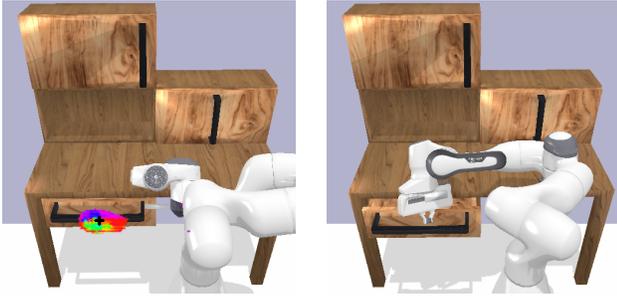
2) *Real world:* For the real world experiment, we setup the environment using a 7-DOF Franka Emika Panda robot. The full setup can be seen in Figure 1. Similar as in simulation, we collect play data by teleoperating the robot using a VR controller as shown in Figure 2. We accumulate 1.5 hours of human interaction, which results in  $\sim 70K$  images and use this to train both the gripper camera and static camera affordance models. The labels for both simulation and real world experiments are obtained as described in Section III-A. We only use the data to train the affordance models and do not need human annotation.

### B. Evaluation Protocol

To test the sample efficiency of the affordance-guided RL policy, we compare against a sparse-reward SAC agent, *local-SAC*. For this baseline, we remove  $R_{aff}$  from the reward function and we modify the observation by removing the affordance mask and distance to the center. This policy still uses  $\pi_{mod}$  to move through free space, but does not use the affordances for interaction. In essence, it is a sparse-reward SAC agent operating with the RGB-D images of the gripper camera in the vicinity of the objects. For all the experiments we show the success rate as the average success over a given number of attempts to complete a task.

### C. Simulation Experiments

We start of by training policies to lift a diverse set of 15 objects on which the affordance model was also trained on. We observe that our approach outperforms the baseline and lifts significantly more objects as it has a strong prior on how objects should be interacted with. We observe that VAPO successfully can grasp objects at the anticipated afforded regions (handle of a pan, power-drill, knife), while the baseline fails to grasp objects of complex (frying pan) or ambiguous (bowl) geometries. This shows the effectiveness of the affordance-guided policy in learning stable functional grasps. Not only does our method learn better, but it is

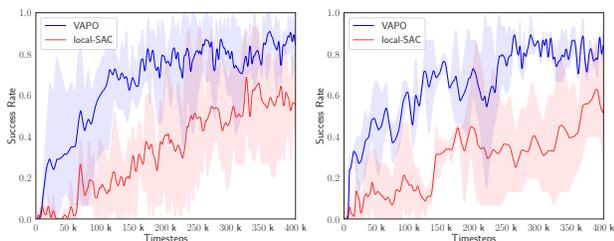


**Fig. 5:** Drawer opening task. On the left, the detected affordance and the corresponding center are shown. On the right we show a rollout of the RL agent opening the drawer.

critically more sample-efficient. After  $\sim 30$  hours (400k timesteps) of robot interaction the baseline reaches a success rate of 0.6, while VAPO matches this performance at 100k steps. This indicates that our method learns up to  $4\times$  faster than the baseline. After training for 400k timesteps, VAPO remains stable at an overall success rate of 0.90.

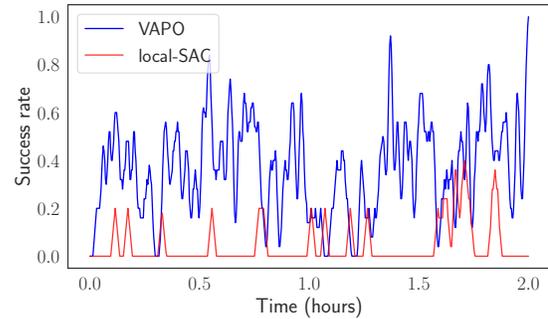
Next we push our affordance model to generalize to unseen objects in two sets of experiments. In the first setting, we train and test the policies on 15 objects which were not seen by the affordance model during training. We observe in Figure 6 that VAPO outperforms the baseline by a large margin in terms of both number of objects lifted and sample-efficiency. In the second setting, we evaluate the trained policies zero-shot on lifting 15 unseen objects. This form of zero-shot evaluation is very challenging, as the objects are unseen for both the affordance model and the RL agent. We report a lifting success of 13/15 for VAPO and 8/15 for the baseline. This demonstrates the effectiveness of imbuing robots with an object-centric visual affordance. Aside from accelerating learning, the visual affordance model generalizes sufficiently to new object shapes and can anticipate their affordance regions, providing a useful object-centric prior.

To analyze if our approach is applicable for more tasks, we conduct experiments on a drawer opening task (Figure 5). We report a success rate over 100 episodes of 0.84 for VAPO and of 0.52 for the baseline. The results are consistent with the previous experiments showing that our method outperforms the baselines, while being more sample-efficient.



**(a)** Policy trained on seen objects by the affordance model. **(b)** Policy trained on **unseen** objects by the affordance model.

**Fig. 6:** VAPO vs. local-SAC for the pick up tasks. In both experiments, our method learns  $4\times$  faster as compared to the baseline and successfully lifts most of the objects.



**Fig. 7:** VAPO vs. local-SAC in real world tidy-up experiments. The success rate over the last ten episodes is shown. After two hours of real world robot interaction, the baseline rarely lifts any objects, while our approach consistently “functionally” grasps all the objects.

#### D. Real World Experiments

We finally evaluate our approach on a real world tidy-up experiment. We show the learning curves for this experiment in Figure 7. We use a 7-DOF Franka Emika Panda robot and run our policy at 20 Hz. We train all methods to pickup four objects: a plastic banana, a screwdriver, a table tennis racket and a paint roller. After two hours of training VAPO is able to consistently “functionally” grasps all the objects, e.g., grasping the objects by the handles, while the SAC baseline rarely achieves to lift any object, despite the agent starting at the same robot pose as our method. This is due to the low number of samples that sparse-reward SAC is trained on, since most success stories of RL in the real world require several orders of magnitude more data [24]. Overall, our results demonstrate the effectiveness of our approach to learn sample-efficient policies by leveraging self-supervised visual affordances.

## VI. CONCLUSION

In this paper, we introduced the novel approach VAPO (Visual Affordance-guided Policy Optimization) as a method for sample-efficient policy learning of manipulation tasks that is guided by a self-supervised visual affordance model. The key advantage of our formulation is the extraction of visual affordances from unlabeled human teleoperated play data to learn a strong prior about *where* actionable regions in an environment are. We distill this knowledge into an interplay between model-based and model-free policies that allows for a sample-efficient division of the robot control learning, without assuming a predefined set of manipulation primitives, 3D object shapes or a tracking system. Our results show that aside from accelerating learning, a critical advantage of imbuing robots with an object-centric visual affordance prior is the ability of policies to generalize to unseen, functionally similar, objects. To the best of our knowledge, this work is the first one to demonstrate the effectiveness of visual affordances to guide model-based policies and closed-loop RL policies to learn robot manipulation tasks in the real world.

## REFERENCES

- [1] J. J. Gibson, "The ecological approach to visual perception," *Boston: Houghton Mifflin*, 1979.
- [2] M. Hassani, S. Khan, and M. Tahtali, "Visual affordance and function understanding: A survey," *arXiv preprint arXiv:1807.06775*, 2018.
- [3] A. Nguyen, D. Kanoulas, D. G. Caldwell, and N. G. Tsagarakis, "Detecting object affordances with convolutional neural networks," in *IROS*, 2016.
- [4] A. Myers, C. L. Teo, C. Fermüller, and Y. Aloimonos, "Affordance detection of tool parts from geometric features," in *ICRA*, 2015.
- [5] A. Nguyen, D. Kanoulas, D. G. Caldwell, and N. G. Tsagarakis, "Object-based affordances detection with convolutional neural networks and dense conditional random fields," in *IROS*, 2017.
- [6] T.-T. Do, A. Nguyen, and I. Reid, "Affordancenet: An end-to-end deep learning approach for object affordance detection," in *ICRA*, 2018.
- [7] K. Mo, L. Guibas, M. Mukadam, A. Gupta, and S. Tulsiani, "Where2act: From pixels to actions for articulated 3d objects," in *ICCV*, 2021.
- [8] T. Nagarajan and K. Grauman, "Learning affordance landscapes for interaction exploration in 3d environments," in *NeurIPS*, 2020.
- [9] A. Zeng, S. Song, K.-T. Yu, E. Donlon, F. R. Hogan, M. Bauza, D. Ma, O. Taylor, M. Liu, E. Romo *et al.*, "Robotic pick-and-place of novel objects in clutter with multi-affordance grasping and cross-domain image matching," in *ICRA*, 2018.
- [10] L. Yen-Chen, A. Zeng, S. Song, P. Isola, and T.-Y. Lin, "Learning to see before learning to act: Visual pre-training for manipulation," in *ICRA*, 2020.
- [11] C. Lynch, M. Khansari, T. Xiao, V. Kumar, J. Tompson, S. Levine, and P. Sermanet, "Learning latent plans from play," in *CoRL*, 2020.
- [12] T. Silver, K. Allen, J. Tenenbaum, and L. Kaelbling, "Residual policy learning," *arXiv preprint arXiv:1812.06298*, 2018.
- [13] M. A. Lee, C. Florensa, J. Tremblay, N. Ratliff, A. Garg, F. Ramos, and D. Fox, "Guided uncertainty-aware policy optimization: Combining learning and model-based strategies for sample-efficient policy learning," in *ICRA*, 2020.
- [14] B. Ichter, P. Sermanet, and C. Lynch, "Broadly-exploring, local-policy trees for long-horizon task planning," in *CoRL*, 2021.
- [15] A. X. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song, H. Su *et al.*, "Shapenet: An information-rich 3d model repository," *arXiv preprint arXiv:1512.03012*, 2015.
- [16] Y. You, Y. Lou, C. Li, Z. Cheng, L. Li, L. Ma, C. Lu, and W. Wang, "Keypointnet: A large-scale 3d keypoint dataset aggregated from numerous human annotations," in *CVPR*, 2020.
- [17] K. Mo, S. Zhu, A. X. Chang, L. Yi, S. Tripathi, L. J. Guibas, and H. Su, "Partnet: A large-scale benchmark for fine-grained and hierarchical part-level 3d object understanding," in *CVPR*, 2019.
- [18] O. Mees, A. Emek, J. Vertens, and W. Burgard, "Learning object placements for relational instructions by hallucinating scene representations," in *ICRA*, 2020.
- [19] Y. Xiang, T. Schmidt, V. Narayanan, and D. Fox, "Posecnn: A convolutional neural network for 6d object pose estimation in cluttered scenes," in *RSS*, 2017.
- [20] J. Tremblay, T. To, B. Sundaralingam, Y. Xiang, D. Fox, and S. Birchfield, "Deep object pose estimation for semantic robotic grasping of household objects," in *CoRL*, 2018.
- [21] O. Mees, M. Tatarchenko, T. Brox, and W. Burgard, "Self-supervised 3d shape and viewpoint estimation from single images for robotics," in *IROS*, 2019.
- [22] O. Mees and W. Burgard, "Composing pick-and-place tasks by grounding language," in *ISER*, 2021.
- [23] S. Levine, C. Finn, T. Darrell, and P. Abbeel, "End-to-end training of deep visuomotor policies," *JMLR*, vol. 17, pp. 1334–1373, 2016.
- [24] D. Kalashnikov, A. Irpan, P. Pastor, J. Ibarz, A. Herzog, E. Jang, D. Quillen, E. Holly, M. Kalakrishnan, V. Vanhoucke *et al.*, "Qt-opt: Scalable deep reinforcement learning for vision-based robotic manipulation," in *CoRL*, 2018.
- [25] P. Sermanet, C. Lynch, Y. Chebotar, J. Hsu, E. Jang, S. Schaal, and S. Levine, "Time-contrastive networks: Self-supervised learning from video," in *ICRA*, 2018.
- [26] A. Singh, L. Yang, K. Hartikainen, C. Finn, and S. Levine, "End-to-end robotic reinforcement learning without reward engineering," in *RSS*, 2019.
- [27] O. Mees, M. Merklinger, G. Kalweit, and W. Burgard, "Adversarial skill networks: Unsupervised robot skill learning from videos," in *ICRA*, 2020.
- [28] I. Lenz, H. Lee, and A. Saxena, "Deep learning for detecting robotic grasps," *IJRR*, vol. 34, no. 4-5, pp. 705–724, 2015.
- [29] 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," *IJRR*, vol. 37, 2018.
- [30] P. Mandikal and K. Grauman, "Dexterous robotic grasping with object-centric visual affordances," in *ICRA*, 2021.
- [31] C. Castellini, T. Tommasi, N. Noceti, F. Odone, and B. Caputo, "Using object affordances to improve object recognition," *IEEE transactions on autonomous mental development*, vol. 3, no. 3, pp. 207–215, 2011.
- [32] T. Nagarajan, C. Feichtenhofer, and K. Grauman, "Grounded human-object interaction hotspots from video," in *ICCV*, 2019.
- [33] C. Xie, Y. Xiang, A. Mousavian, and D. Fox, "Unseen object instance segmentation for robotic environments," *IEEE Transactions on Robotics (T-RO)*, 2021.
- [34] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," *arXiv preprint arXiv:1505.04597*, 2015.
- [35] Y. Xiang, T. Schmidt, V. Narayanan, and D. Fox, "Posecnn: A convolutional neural network for 6d object pose estimation in cluttered scenes," in *RSS*, 2018.
- [36] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, "Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor," in *International conference on machine learning*. PMLR, 2018, pp. 1861–1870.