

Affordance Learning for End-to-End Visuomotor Robot Control

Aleksi Hämmäläinen¹, Karol Arndt¹, Ali Ghadirzadeh^{1,2} and Ville Kyrki¹

Abstract—Training end-to-end deep robot policies requires a lot of domain-, task-, and hardware-specific data, which is often costly to provide. In this work, we propose to tackle this issue by employing a deep neural network with a modular architecture, consisting of separate perception, policy, and trajectory parts. Each part of the system is trained fully on synthetic data or in simulation. The data is exchanged between parts of the system as low-dimensional latent representations of affordances and trajectories. The performance is then evaluated in a zero-shot transfer scenario using Franka Panda robot arm. Results demonstrate that a low-dimensional representation of scene affordances extracted from an RGB image is sufficient to successfully train manipulator policies. We also introduce a method for affordance dataset generation, which is easily generalizable to new tasks, objects and environments, and requires no manual pixel labeling.

I. INTRODUCTION

Recent years have seen widespread research and adoption of deep learning in robotics. Convolutional neural networks became the *de facto* standard way of processing visual input, and generative models are starting to see wider adoption as trajectory generators [1]. Data-driven approaches like deep learning make it possible to learn diverse behaviours with little feature engineering [2] and can be scaled to complex problems if enough data is available.

Despite extensive use of deep learning, a few issues limit its further adoption. Specifically, (i) end-to-end training requires vast amounts of task-specific training data and months of real world experience [3], which are simply unfeasible to provide; (ii) using task- and domain-specific training data means that changing the task, such as transferring it to another robot or changing the task objective, requires often entirely new training data, which is costly to obtain [2]; (iii) the behaviour of deep models can be difficult to predict and interpret, as the intermediate representations learned by the neural network rarely provide meaningful information to humans — despite recent successes in interpreting neural networks, it is still difficult to gain understanding of why the model behaves the way it does.

In this paper, we propose to address those issues with a modular neural network, which consists of three parts: perception, policy and trajectory generation. Each part of the system is, on its own, a neural network trained entirely on synthetic data or in simulation, with no real-world adaptation. Information is exchanged between parts using low-dimensional representations of affordances and trajectories. By using visual affordances, the intermediate representation

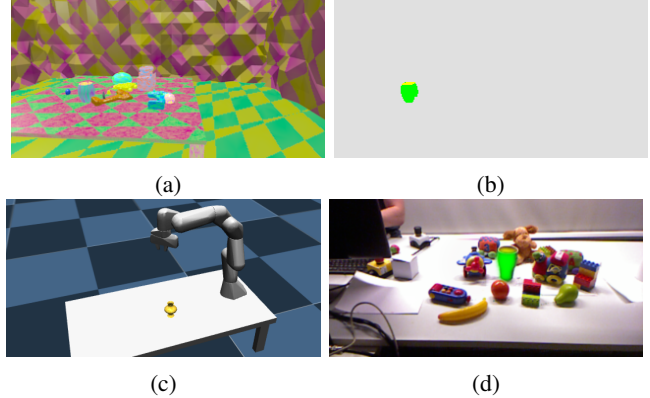


Fig. 1: In our experiments, we train an affordance detection model using synthetic scene images (a) with their respective affordance images (b). The intermediate representations of this model are used to train manipulator policies in MuJoCo (c). The performance of the system is evaluated in a real world setup (d). The cup’s affordances, *wrap-grasp* and *contain*, are respectively marked in green and yellow.

captures semantic information of what can be done with each part of each object [4], which provides the system with generality across different tasks.

The perception part is a variational encoder-decoder structure [5] trained for affordance detection from synthetic RGB images. Those images are generated using randomized textures, object shapes, distractor objects and camera viewpoints to maximize the model’s capability to generalize to new environments [6], [7]. The encoder part outputs a low-dimensional representation of scene affordances. The trajectory part is a variational autoencoder trained on a set of trajectories. The policy part, trained in a physics simulator, maps latent affordance representations to latent representations of trajectories, which are then passed to the trajectory decoder. Using variational autoencoders ensures smoothness of the latent space [8], which, together with the low dimensionality of the latent representation, makes it quick to retrain the policy layers using a small amount of training data.

We demonstrate that the system achieves good performance in the task of inserting a ball into a container both in a simulated environment and on real hardware (the Franka Panda robotic arm), without the need for any real world adaptation of the parts (a zero-shot transfer scenario). We also evaluate its susceptibility to clutter, demonstrating that it can complete the task even in heavily cluttered environments, such as shown in Figure 1d, without having any explicit model of the cup object.

We also propose a new method of generating a domain

*This work was supported by Academy of Finland grant 313966.

¹Aalto University, Espoo, Finland first.last@aalto.fi

²KTH Royal Institute of Technology, Stockholm, Sweden

randomized dataset for affordance detection, which allows to automatically generate vast amounts of data for affordance training without the need of hand-labeling pixels. Perception models trained using this affordance data can then be used both in real world and in simulated environments. The code for generating the dataset and training our models is available online¹.

II. RELATED WORK

In this section, we introduce studies related to the problem of visual perception training for control. The three commonly used approaches to address the problem of perception training are: (1) end-to-end training of the perception and control layers in a single neural network architecture, e.g., [2], [9], [10], (2) perception training using priors and auxiliary tasks, e.g., [1], [11], and (3) perception training with labeled affordances, e.g., [12]–[15].

End-to-end training results in task specific robust features when large enough training datasets are provided [2]. However, even with large datasets, the learned features may not generalize well to similar tasks. One approach to solve this issue is to introduce extra loss functions to extract general multitask features. A popular solution is to *learn low-dimensional state representations* of the visual data [1], [11], [16], [17]. Encoder-decoder architectures (also known as auto-encoders when the input data is restored at the output) are widely used models for this purpose. The hourglass shape of the network abstracts the input data at the bottleneck (the latent space) based on which the output data can be constructed. Spatial autoencoders, introduced by Finn et al., [11], use CNNs to extract a low-dimensional feature vector consisting of spatial image points to localize task-related objects in the input image. Spatial features are suitable for manipulation tasks as they contain information to move the robot end-effector to desired positions. However, the features may not be stable and robust given limited training data [1].

In order to improve the robustness and the generalization power of the extracted features, prior knowledge and auxiliary tasks [9], [10], [18]–[20] are introduced as extra loss functions for the training. Prior knowledge ranges from physical laws [18] to high-level semantic features [9], [10]. Overview of different priors for the perception training problem can be found in [21] and [22]. We train a low-dimensional representation of visual inputs using variational encoder-decoder structures, which, similar to [16], [17], assigns a normal distribution over the latent space conditioned on the observations. This prevents similar states from being scattered in the latent space, which is a suitable property for control purposes. Furthermore, the features can be generally used for different tasks, since they are trained to extract high-level action affordances of the objects irrespective of the task. Similar to [1], we also trained a low-dimensional representation of motor trajectories based on variational autoencoders. This further improves the data efficiency of the proposed method.

Learning visual affordances is an approach to identify a set of functions available to an object given visual observations. Several recent studies worked on understanding affordances at the image pixel level, i.e., to segment pixels with the same affordances [23]–[26]. However, there are few studies that deal with learning affordance representations suitable for control.

Myers et al. [27] demonstrated that local geometric primitives can be used to detect pixel-level affordances. Song et al. [12] proposed to combine a local affordance extraction method based on the SIFT descriptor, and a global object grasp localization to learn grasp affordances for robotic applications. Nguyen et al. [13] proposed a CNN encoder-decoder architecture to learn different affordances from RGB-D images. In a more recent work, they proposed to use an object detector network to narrow down the image regions at which object affordances are extracted. Do et al. [15] proposed to jointly optimize the object detection and the affordance learning losses in a single feed-forward neural network to speed-up the process in the deployment phase. However, these methods do not extract a suitable representation of the action affordances that can be directly used for control purposes. All the methods are based on a post-processing step in which affordance images are further processed to obtain suitable motor actions. In our method, we extract pixel-level affordances together with a low-dimensional representation which can be efficiently used for end-to-end visuomotor policy training.

A major issue with pixel-wise affordance learning is the intensive image labeling work required to collect training datasets. To alleviate this issue, Srikantha and Gall [28] introduced an expectation-maximization method to learn pixel-level affordances given weakly labeled image data. In this work, we propose a different approach to overcome the difficulties of manual pixel-level labeling. We leverage sim-to-real transfer learning based on domain randomization to collect a diverse set of training data which can be used with the real image inputs.

Domain randomization in simulation is a technique for zero-shot sim-to-real transfer learning by randomizing non-essential aspects of a simulated process. This helps to improve the diversity of the training dataset, which, in turn, helps the extraction of a set of task-relevant features that can be transferred to real world problems. Our method resembles the work of [6], [7], [29] in that we render simulated scenes by randomizing textures and appearances of the objects in the scene. An important contribution of our work is that we get large-scale pixel-level affordance labels automatically without extra manual labeling work.

III. METHOD

The system consists of three separately trained parts—affordance perception, policy, and trajectory generation, similar to [1]. The perception part encodes an RGB image o to a low-dimensional affordance representation s . The policy part maps the state s to a desired low-dimensional action vector

¹ <https://github.com/gamleksi/affordancegym>

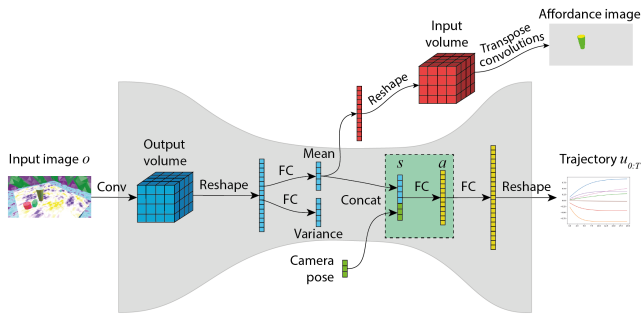


Fig. 2: Structural overview of the system. An input image o is first processed by the affordance encoder (blue). The policy part (green rectangle) produces action a , which is then decoded to a robot trajectory $u_{0:T}$. The affordance image is also generated by the affordance decoder (red).

a . Finally, based on the action a , the trajectory part generates a trajectory of robot joint positions $u_{0:T}$.

This structure is shown in Figure 2. The affordance encoder consists of convolutional and fully connected layers (blue), while the decoder has the inverse structure (fully connected layers followed by transposed convolutions; marked in red). The yellow blocks represent the policy and trajectory generation parts. Information about camera pose is added to the latent space before passing it to the policy to allow it to estimate the 3D pose of the target object w.r.t. the camera. Gray background emphasizes the hourglass shape of the network with a low-dimensional information bottleneck in the middle. Both perception and trajectory blocks are built using variational encoder-decoder structures introduced in the next section.

A. Variational encoder-decoder structures

To facilitate quick training and to improve data efficiency of the learning method, we use low-dimensional representations of affordances and trajectories as the intermediate steps between each part of the system. To obtain these low-dimensional representations in an unsupervised way without using task-specific hand-engineered features, we use variational encoders-decoders — a more general case of a variational autoencoder (VAEs).

Autoencoders are hourglass shaped neural networks with a low-dimensional information bottleneck in the middle, which encodes the *latent representation* of the input vector. Variational autoencoders additionally incorporate a KL-divergence term in their loss functions, which encourages the neighbourhood preservation and disentanglement properties in the latent space [8], [30]. The trade-off between these properties and reconstruction quality can be balanced by the β hyperparameter, introduced by [8]. In our method, these properties make it easier for the policy to find a smooth mapping between the affordances and the actions.

In contrast to a typical VAE application, our affordance perception model uses the latent representation to create an affordance map, instead of reconstructing the input. To emphasize on this, we refer to it as a *variational encoder-decoder structure*.

B. Affordance perception

The perception layer extracts affordance information from an observation o and encodes it as a latent space vector s . Some properties, such as textures and the surrounding environment, have no impact on object affordances and can be ignored by the perception model.

To obtain those low-dimensional affordance representations, we use a variational affordance encoder-decoder (VAED) structure, which encodes an RGB image and generates the corresponding pixel-level affordance labeling. The encoder of VAED consists of a series of convolutional layers followed by one fully connected layer. The mean and covariance are produced by separate fully connected layers. The latent vector, drawn from the posterior distribution, is then passed as an input to the decoder. The decoder has a reversed structure—a fully connected layer followed by a series of transposed convolution layers. The decoder produces a multi-channel probability map, with each channel corresponding to a specific affordance. A pixel value in an affordance channel describes how likely that affordance is to occur in that pixel. The probability output is achieved by using the sigmoid activation function. We use binary cross entropy as the reconstruction loss.

Training the affordance detection network requires corresponding labeled affordance images. The target labels are generated in Blender together with the RGB inputs. We use domain randomization of textures, clutter and lighting to increase the robustness of the model. Even though artificial textures make the input images look unrealistic, high variation allows the model to generalize to real world environments [6].

C. Trajectory generation

To build a generative model for trajectories, we first train a variational autoencoder on a set of task-specific trajectories, similar to [1]. The decoder block of this autoencoder converts low-dimensional latent trajectory representations of actions a to their corresponding trajectories in joint space $u_{0:T}$. This model can therefore be used as a task-specific trajectory generator—it converts each latent action a to a trajectory that is useful for accomplishing the given task. To achieve this, we rely on the observation that trajectories useful for a specific task often exhibit structural similarities, e.g., trajectories useful for pouring liquids into containers consist of motions to a point above the target location followed by a wrist rotation. We generate such trajectories using MoveIt!, an open source motion planning software. The encoder and the decoder of the trajectory VAE consist of three fully connected layers each. We use mean squared error (MSE) as the reconstruction loss.

D. Policy layer

A policy part maps the affordance information s to a desired action vector a . The policy part includes three fully connected layers. In addition to the affordance information s , camera parameters are included as input to the policy network, which allows the policy to account for the relative position of the camera w.r.t the detected affordances. The

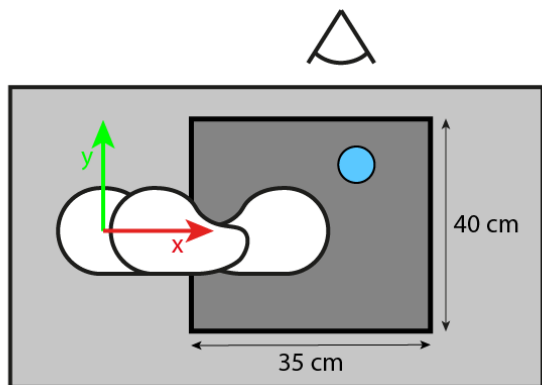


Fig. 3: Experimental setup, as viewed from top. The robot is marked in white. The dark gray area shows the workspace and the blue circle an example location of a target container. The camera is located on the left side (from the robot’s perspective).

TABLE I: Results of the UMD dataset evaluation

Affordance	VAED		Nguyen et al.	
	RGB	RGB-D	RGB	RGB-D
grasp	0.667	0.643	0.719	0.714
cut	0.720	0.654	0.737	0.723
scoop	0.760	0.758	0.744	0.757
contain	0.859	0.858	0.817	0.819
pound	0.757	0.748	0.794	0.806
support	0.792	0.787	0.780	0.803
wrap-grasp	0.774	0.777	0.769	0.767
Average	0.761	0.747	0.766	0.770

policy part can be trained either in a supervised manner or using reinforcement learning. We utilize supervised learning to train the policy.

IV. EXPERIMENTS

The goal of the experiments was to study if a low-dimensional affordance representation can be used to successfully train manipulator policies. In particular, we first evaluated if affordances can be represented in a latent space using the UMD dataset [27], a standard benchmark for the affordance detection task.

To study the approach as a whole, we developed an experimental set-up presented in Figure 3. In this setup, the task of the Franka Panda robot arm was to insert a ball into a container that is located on a table, with the scene observed by an RGB-camera. We used the set-up to evaluate the performance quantified as position accuracy and task success both in simulation and on physical hardware.

A. UMD dataset benchmark

We first evaluated the capabilities of variational affordance encoder-decoder to detect object affordances on the UMD dataset [27], randomly splitting it into training (70%) and validation (30%) data. For this evaluation, we used 20 dimensional latent space and KL-divergence penalty $\beta = 4$. We used F_{β}^w as the evaluation metric [31]. The results are presented in Table I and are compared to the CNN-RGB and CNN-RGBD networks introduced in [13].

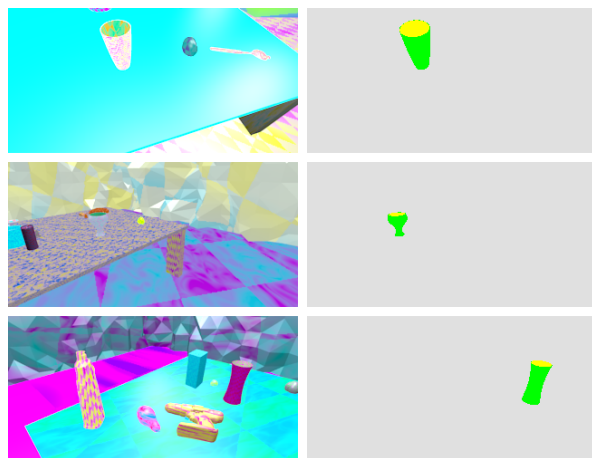


Fig. 4: Samples from the affordance dataset. Green indicates the *wrap-grasp* affordance and yellow indicates the *contain* affordance. Gray pixels mean no affordance at all.

We can notice that for some affordances — *contain* and *wrap-grasp* — our model slightly outperformed the baseline, while for the other affordances the results were the same or slightly worse. Comparing the RGB and RGB-D results, we can also see that discarding the depth information does not have much impact on the performance.

This evaluation shows that the variational encoder-decoder structure using RGB images can be used for affordance detection.

B. Training the perception model

We used Blender’s Python API to generate a domain randomized dataset of 1 million images. Figure 4 shows examples from the generated dataset. The Blender environment was constructed similar to the experimental setup: it included a table with a cup and clutter objects, and the table was surrounded by walls. Each sample includes a rendered RGB image and a corresponding affordance image, which was used as the target output during training.

A simple cylinder model was built in Blender to represent a cup-like object. Its shape was randomly generated by smoothly changing its diameter along the height. Textures of the inner and outer parts of the object were separately randomized. In our setting, two different affordances occur: the outer part has a *wrap-grasp* affordance and the inner part has a *contain* affordance. Clutter objects were located on the table and their affordances were ignored. In total, 66 clutter objects from the Yale-CMU-Berkeley (YCB) Object and Model set [32] were used.

The following features were randomized uniformly within reasonable limits: (1) positions, scales, and textures of the clutter objects and the cup, (2) shape of the cup, (3) texture and scale of the table, (4) camera pose, (5) the number of clutter objects on the table, and (6) the number and positions of lights in the scene.

The affordance detection network had 4 convolutional and corresponding transposed layers. We used KL-divergence penalty $\beta = 4$ and a 10 dimensional latent space throughout the experiments.

C. Training the trajectory model

Training data included evenly distributed samples of trajectories that moved the end-effector on top of the workspace region. The initial pose of the robot was the same as in the test environment. Samples were generated with *RRT** using MoveIt!. The value of the β coefficient was increased while training such that it was initialized with $\beta = 10^{-8}$, and then β was iteratively increased in 400 epoch intervals towards $\beta = 10^{-5}$. By iteratively increasing β , the model learns first to reconstruct trajectories accurately without considering the smoothness and neighborhood preservation properties of the low-dimensional space. Without this, the model got stuck to local minima, where its trajectory reconstruction accuracy was poor.

Trajectories were encoded to 5 dimensional action space, and the number of steps in each trajectory was 24.

D. Policy training

The policy training setup was constructed in simulation to mimic the physical setup (Figure 1c). The relative poses of the camera and a cup to the robot pose corresponded to the experimental setup. In order to train the policy, the affordance perception, policy and trajectory generation parts were combined into a single network, with the weights of the affordance perception and trajectory generation frozen.

The target of the policy was to determine an action that generates a trajectory whose final point lies above the cup. The training loss was calculated as the mean squared error between the target end-effector position and the actual one.

Since the trajectory is expressed in joint space, the forward kinematics solution for the final point u_T was computed. In order to backpropagate the error through this computation, the base-to-hand transformation matrix—which depends on the joint configuration—was expressed as a PyTorch computation graph.

As the latent representation is different for various cup shapes, positions, and camera angles, we used 15 different cup shapes and randomized the camera and cup positions. In total, our training set contained 3 million images with labeled cup positions. The average distance error was 0.9cm with validation data.

E. End-to-end experiments

We evaluated the system with and without clutter in terms of detection accuracy (position error) and task success rate (placing a ball inside the object).

1) *No clutter*: Without clutter, 13 different cup like objects shown in Figure 5 were used. Each of the objects was placed in 10 to 12 positions, with a total number of 143 trials. The diameter of the ball was 4 cm.

The position errors and success rates are reported in Table II. The average error over all objects was 2.2 cm with three objects having over 3 cm average error. The success rate was 100% for 8/13 objects and 73–91% for the rest. The accuracy can be considered good, taking into account that the system has been trained to detect any general cylindrical containers.

The variation between objects seems to be explained by objects' appearance: The objects with the highest average

TABLE II: Experiment results for different positions in clutterless environment. All distances are given in centimeters.

Cup	Radius	Height	x error	y error	Error	Success
blue	3.45	16.60	0.73	2.48	2.68	82%
can	3.32	10.15	0.76	1.23	1.57	100%
can2	4.22	13.76	0.52	1.48	1.63	100%
green	3.92	9.33	0.74	1.00	1.35	100%
white	4.06	10.44	0.67	2.41	2.53	91%
jar	4.75	16.30	1.52	3.20	3.84	82%
red	4.30	8.08	0.59	1.60	1.76	100%
rocket	4.06	13.06	0.54	0.99	1.17	100%
stack1	3.70	10.52	0.52	1.99	2.07	100%
stack2	4.46	8.361	0.56	3.07	3.14	75%
stack3	2.72	8.76	0.51	1.15	1.31	100%
pastel	3.83	10.67	0.56	1.68	1.84	100%
yellow	4.06	5.69	0.95	3.45	3.69	73%
Average	3.91	10.90	0.71	1.97	2.19	92.5%



Fig. 5: Various cups and the ball used throughout the experiments. Left to right, back row: white, pastel, rocket, can2, blue, jar. Front row: yellow, red, green, can, stack1, stack3. stack2 is composed from two bottom layers of stack 1, i.e. the large yellow and green parts.

distance error—yellow and jar—differ significantly from those in the training set. The yellow cup, for example, is smaller than ones in training set, with a height of 5.7 cm. The system therefore expects the cup to be taller, and hence there is systematic error along the y axis. The object labeled as jar is a large aluminum container with a mirror texture on the outside. Even though the training data did not contain any reflective objects, the system still managed to detect the opening and successfully place the ball inside of it 82% of the times.

To study the position dependency of error, all cup positions used for the experiments, together with their average error ellipses are shown in Figure 6. There are no consistent differences in accuracy across positions. However, we can see that the system performs more accurately along the x -axis than the y -axis. The higher error along y can be explained by the lack of depth information in the system and camera placement. While the position of the object along the x -axis can easily be deduced from the input image, calculating the y -position without knowledge of depth or object size is difficult.

2) *Cluttered scene*: To evaluate the system's susceptibility to clutter, we chose three objects (rocket, red, can) from the ones used for previous experiments, providing a variety of different shapes, sizes and texture. The clutter objects were

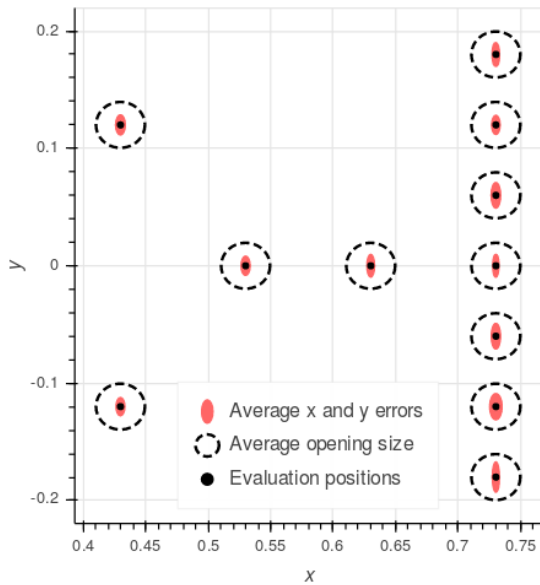


Fig. 6: Evaluation positions with their respective average error ellipses marked.

TABLE III: Results of experiments in the cluttered environment. Columns d , $\%$ and n indicate, respectively, the average position error in cm, percentage of successful trials and the total number of trials.

Clutter objects	Rocket			Can			Red		
	d	$\%$	n	d	$\%$	n	d	$\%$	n
0-6	1.25	97	45	1.64	92	28	2.02	94	69
7-12	1.86	96	30	2.93	76	26	2.4	82	64
13-18	2.84	100	6	4.47	73	30	4.76	63	36
19+	2.27	80	15	6.04	35	14	4.84	60	10

toys from a toy dataset², which were not seen during training. They were randomly shuffled on the table between each trial and their number was steadily increased from 0 to 25. The cup object position was fixed during these experiments.

Example scene images captured during these experiments are shown in Figure 7. The *wrap-grasp* affordance of the cup is visualized in green and the *contain* affordance in yellow. We can see that despite quite heavy clutter the cup was successfully identified in both situations.

The position errors and success rates are presented in Table III. Not surprisingly, the errors increase and success rates decrease with increasing clutter, but the performance deterioration is gradual. Despite the training data having

²<http://irobotics.aalto.fi/software-and-data/toy-dataset>

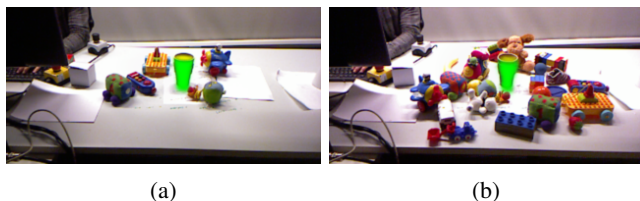


Fig. 7: Successful affordance detection with different amounts of clutter.

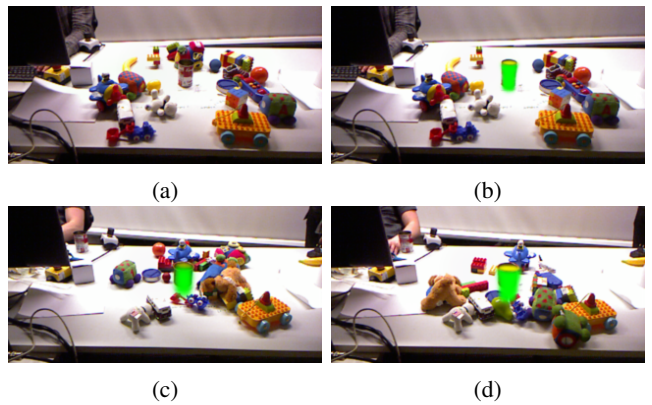


Fig. 8: Examples of failure cases in different clutter configurations: the detection fails with clutter behind the object (a) and succeeds when the clutter is removed (b). With heavy clutter in front of the object its height is incorrectly estimated (c, d)

only 0 to 10 clutter object on the table, our system proved to be resilient to a larger amount of clutter. This is most likely caused by applying heavy texture randomization when the training data was generated. The perception model thus learned to perceive clutter as just a different texture of the table - both of these look identical if there is no depth information available to the model.

Noticing that the task can succeed even in heavily cluttered scenes, such as the one shown in Figure 7b, we further analyzed the data to determine failure cases. We found that the amount of clutter was not as important as its positioning—when clutter objects were located directly in front of the cup, covering its bottom part, the system quite often failed to correctly position the arm on the y axis. This is expected to happen—when the bottom part of the cup is not visible, it is impossible for the model to deduce its height, and thus estimate the y position. The situation is illustrated in Figures 8c and 8d.

Putting object directly behind the cup opening also had a noticeable impact on performance, although it was less significant. In that case, the problems are likely to be caused by the perception model having problems estimating where the opening of the cup ends. It is, however, a difficult task to perform without depth data or stereovision. This is illustrated in Figure 8a, where a distractor object is placed directly behind the opening of a can. As the distractor and the inside of the can have a similar color, it becomes difficult to identify the opening in the image, and so the affordance part fails. When the distractors are removed from the back (Figure 8b), the object is successfully identified again.

V. CONCLUSIONS AND FUTURE WORK

We studied the use of latent space affordance encodings in training manipulator policies. We demonstrated that the affordance encodings can be used to train policies successfully. We also showed that those latent space representations, if trained in a randomized domain, can be made invariant to, e.g., distractor objects and textures, making them applicable

in zero-shot sim-to-real transfer.

We also developed a method for generating a domain-randomized affordance dataset, which works across various domains (heavily randomized images from Blender, simplistic frames from MuJoCo and real world images with different amounts of clutter). The method is easily scalable to new types of objects and affordances, and the software to generate the dataset has been released to the public.

The proposed methods operate on bare RGB images without depth information. Many recent affordance detection methods use depth as an additional cue. However, domain randomization is an essential component for the simulation-to-real transfer. Thus, domain randomization for depth data seems to be an appealing open research problem.

Another recent idea applied in affordance detection is to include semantic information from e.g. object recognition. These methods do not typically have encoder-decoder structure so that the modular approach proposed here could not be directly applied to find corresponding latent representations. Finding such representations for more complex affordance detectors would be useful and remains an open problem.

Finally, the proposed approach is based on feedforward trajectory generation. However, the latent affordance representations could be used also as the state space of feedback policies, for example, to learn closed-loop visual alignment tasks, or even active closed-loop search strategies for objects with particular affordances.

REFERENCES

- [1] Ali Ghadirzadeh, Atsuto Maki, Danica Kragic, and Mårten Björkman. Deep predictive policy training using reinforcement learning. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 2351–2358. IEEE, 2017.
- [2] Sergey Levine, Chelsea Finn, Trevor Darrell, and Pieter Abbeel. End-to-end training of deep visuomotor policies. *The Journal of Machine Learning Research*, 17(1):1334–1373, 2016.
- [3] Sergey Levine, Peter Pastor Sampedro, Alex Krizhevsky, Julian Ibarz, and Deirdre Quillen. Learning hand-eye coordination for robotic grasping with deep learning and large-scale data collection. 2017.
- [4] James J Gibson. The theory of affordances.
- [5] Diederik P Kingma and Max Welling. Auto-encoding variational bayes, 2013. cite arxiv:1312.6114.
- [6] Josh Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and Pieter Abbeel. Domain randomization for transferring deep neural networks from simulation to the real world. In *Intelligent Robots and Systems (IROS), 2017 IEEE/RSJ International Conference on*, pages 23–30. IEEE, 2017.
- [7] Stephen James, Andrew J Davison, and Edward Johns. Transferring end-to-end visuomotor control from simulation to real world for a multi-stage task. *arXiv preprint arXiv:1707.02267*, 2017.
- [8] Irina Higgins, Loïc Matthey, Xavier Glorot, Arka Pal, Benigno Uribe, Charles Blundell, Shakir Mohamed, and Alexander Lerchner. Early visual concept learning with unsupervised deep learning. *CoRR*, abs/1606.05579, 2016.
- [9] Coline Devin, Pieter Abbeel, Trevor Darrell, and Sergey Levine. Deep object-centric representations for generalizable robot learning. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 7111–7118. IEEE, 2018.
- [10] Avi Singh, Larry Yang, and Sergey Levine. Gplac: Generalizing vision-based robotic skills using weakly labeled images. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5851–5860, 2017.
- [11] Chelsea Finn, Xin Yu Tan, Yan Duan, Trevor Darrell, Sergey Levine, and Pieter Abbeel. Deep spatial autoencoders for visuomotor learning. In *2016 IEEE International Conference on Robotics and Automation (ICRA)*, pages 512–519. IEEE, 2016.
- [12] Hyun Oh Song, Mario Fritz, Daniel Goehring, and Trevor Darrell. Learning to detect visual grasp affordance. *IEEE Transactions on Automation Science and Engineering*, 13(2):798–809, 2016.
- [13] Anh Nguyen, Dimitrios Kanoulas, Darwin G Caldwell, and Nikos G Tsagarakis. Detecting object affordances with convolutional neural networks. In *Intelligent Robots and Systems (IROS), 2016 IEEE/RSJ International Conference on*, pages 2765–2770. IEEE, 2016.
- [14] Anh Nguyen, Dimitrios Kanoulas, Darwin G Caldwell, and Nikos G Tsagarakis. Object-based affordances detection with convolutional neural networks and dense conditional random fields. In *Intelligent Robots and Systems (IROS), 2017 IEEE/RSJ International Conference on*, pages 5908–5915. IEEE, 2017.
- [15] Thanh-Toan Do, Anh Nguyen, and Ian Reid. Affordancenet: An end-to-end deep learning approach for object affordance detection. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1–5. IEEE, 2018.
- [16] Manuel Watter, Jost Springenberg, Joshka Boedecker, and Martin Riedmiller. Embed to control: A locally linear latent dynamics model for control from raw images. In *Advances in neural information processing systems*, pages 2746–2754, 2015.
- [17] Herke van Hoof, Nutan Chen, Maximilian Karl, Patrick van der Smagt, and Jan Peters. Stable reinforcement learning with autoencoders for tactile and visual data. In *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 3928–3934. IEEE, 2016.
- [18] Rico Jonschkowski and Oliver Brock. Learning state representations with robotic priors. *Autonomous Robots*, 39(3):407–428, 2015.
- [19] Max Jaderberg, Volodymyr Mnih, Wojciech Marian Czarnecki, Tom Schaul, Joel Z Leibo, David Silver, and Koray Kavukcuoglu. Reinforcement learning with unsupervised auxiliary tasks. *arXiv preprint arXiv:1611.05397*, 2016.
- [20] Piotr Mirowski, Razvan Pascanu, Fabio Viola, Hubert Soyer, Andrew J Ballard, Andrea Banino, Misha Denil, Ross Goroshin, Laurent Sifre, Koray Kavukcuoglu, et al. Learning to navigate in complex environments. *arXiv preprint arXiv:1611.03673*, 2016.
- [21] Ali Ghadirzadeh. *Sensorimotor Robot Policy Training using Reinforcement Learning*. PhD thesis, KTH Royal Institute of Technology, 2018.
- [22] Timothée Lesort, Natalia Díaz-Rodríguez, Jean-François Goudou, and David Filliat. State representation learning for control: An overview. *Neural Networks*, 2018.
- [23] Mohammed Hassanin, Salman Khan, and Murat Tahtali. Visual affordance and function understanding: A survey. *arXiv preprint arXiv:1807.06775*, 2018.
- [24] Ching-Yao Chuang, Jiaman Li, Antonio Torralba, and Sanja Fidler. Learning to act properly: Predicting and explaining affordances from images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 975–983, 2018.
- [25] Timo Luddecke and Florentin Worgotter. Learning to segment affordances. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 769–776, 2017.
- [26] Anirban Roy and Sinisa Todorovic. A multi-scale cnn for affordance segmentation in rgb images. In *European conference on computer vision*, pages 186–201. Springer, 2016.
- [27] Austin Myers, Ching Lik Teo, Cornelia Fermüller, and Yiannis Aloimonos. Affordance detection of tool parts from geometric features. In *ICRA*, pages 1374–1381, 2015.
- [28] Abhilash Srikantha and Juergen Gall. Weakly supervised learning of affordances. *arXiv preprint arXiv:1605.02964*, 2016.
- [29] Stephen James, Paul Wohlhart, Mrinal Kalakrishnan, Dmitry Kalashnikov, Alex Irpan, Julian Ibarz, Sergey Levine, Raia Hadsell, and Konstantinos Bousmalis. Sim-to-real via sim-to-sim: Data-efficient robotic grasping via randomized-to-canonical adaptation networks. *arXiv preprint arXiv:1812.07252*, 2018.
- [30] Christopher P. Burgess, Irina Higgins, Arka Pal, Loïc Matthey, Nick Watters, Guillaume Desjardins, and Alexander Lerchner. Understanding disentangling in -vae. *CoRR*, abs/1804.03599, 2018.
- [31] Ran Margolin, Lihi Zelnik-Manor, and Ayellet Tal. How to evaluate foreground maps. In *2014 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2014, Columbus, OH, USA, June 23-28, 2014*, pages 248–255, 2014.
- [32] Berk Calli, Aaron Walsman, Arjun Singh, Siddhartha Srinivasa, Pieter Abbeel, and Aaron M Dollar. Benchmarking in manipulation research: The ycb object and model set and benchmarking protocols. *arXiv preprint arXiv:1502.03143*, 2015.