

Learning Goal-Oriented Non-Prehensile Pushing in Cluttered Scenes

Nils Dengler

David Großklaus

Maren Bennewitz

Abstract—Pushing objects through cluttered scenes is a challenging task, especially when the objects to be pushed have initially unknown dynamics and touching other entities has to be avoided to reduce the risk of damage. In this paper, we approach this problem by applying deep reinforcement learning to generate pushing actions for a robotic manipulator acting on a planar surface where objects have to be pushed to goal locations while avoiding other items in the same workspace. With the latent space learned from a depth image of the scene and other observations of the environment, such as contact information between the end effector and the object as well as distance to the goal, our framework is able to learn contact-rich pushing actions that avoid collisions with other objects. As the experimental results with a six degrees of freedom robotic arm show, our system is able to successfully push objects from start to end positions while avoiding nearby objects. Furthermore, we evaluate our learned policy in comparison to a state-of-the-art pushing controller for mobile robots and show that our agent performs better in terms of success rate, collisions with other objects, and continuous object contact in various scenarios.

I. INTRODUCTION

Pushing is often used for re-positioning and re-orientating objects since it simplifies the object manipulation in comparison to pick-and-place approaches. Furthermore, pushing allows for moving large, heavy, and irregularly shaped, as well as small and fragile objects to target positions and can be used for reducing uncertainty in the position of objects [1]. Hereby, the term pushing is separated in non-prehensile pushing [2] and prehensile pushing (push-grasp) [3], [4]. For example, in limited space [5], [6] and when dealing with fragile objects, non-prehensile pushing is the preferred manipulation action, since grasping increases the risk of damage. In the past, pushing has been used to separate objects for better grasping [7], [8] or to sort objects from a table into a bin [9] and is assumed to be more time-efficient than grasping to overcome short distances [10]. The range of pushing actions vary between a few centimeters for corrective actions to larger distances, e.g. to place an object in the last row of a shelf. Several approaches aim at predicting the physical properties of objects and use short pushing actions of predefined length [11], [12].

In general, pushing actions should be contact-rich with smooth arm motions. Furthermore contact to other objects in the workspace should be avoided to prevent any damages and

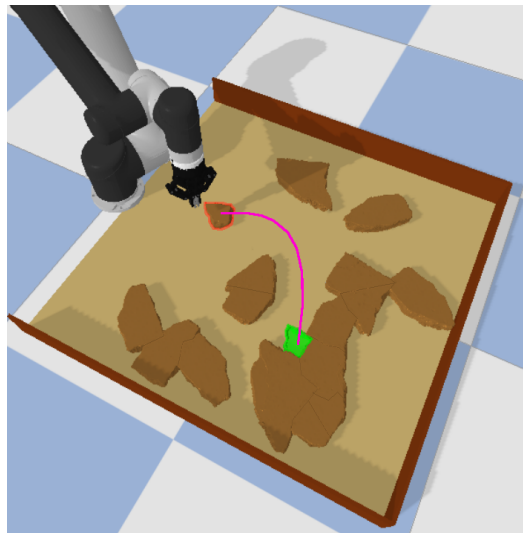


Fig. 1: Targeted application scenario of our system within the RePAIR-project¹. The goal is to push the small fragment to the desired goal pose (green). Shown in magenta is the best pushing path, which maintains a safety distance to the other objects.

changes the configuration of the scene. While for a long time, pushing behaviors were created using expert knowledge in an analytical way, more and more work is focusing on reinforcement learning (RL) to solve this task. Especially the ability to learn from environment interactions and own experiences makes RL a useful way to learn challenging new skills. Start-to-goal pushing with an RL-agent has been tackled before [13] and serves as a benchmark for RL [14], however, pushing in cluttered environments where collisions with other objects have to be avoided is a less researched area. While there are already approaches for mobile bases [15], [16], they have not been transferred to robotic manipulators so far.

In this paper, we present a framework to train an RL-agent that is able to realize obstacle-aware pushing in a contact-rich manner to guide objects with initially unknown dynamics on a planar surface to desired target configurations with a robotic manipulator. As representation of the workspace, we use a depth image taken from a bird's eye view. To reduce the size of the observation space and therefore the complexity, we use the latent space of a variational autoencoder. To accelerate learning, we calculate the optimal 2D path in a grid representation of the environment generated from the depth image. From this path we sample subgoals, which we use as observations to our agent. In addition we use further observations, such as contact information between

All authors are with the Humanoid Robots Lab, University of Bonn, Germany. This work has partially been funded by the European Commission under grant agreement number 964854 –RePAIR – H2020-FETOPEN-2018-2020 and by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany's Excellence Strategy, EXC-2070 – 390732324 – Phenorob.

the end effector of the manipulator and the object as well as the distance to the goal. The output of our system is an incremental motion of the current (x, y, θ) -position of the robot’s end effector. Fig. 1 illustrates a targeted application scenario from the RePAIR-project¹. The goal is to push the small fresco fragment to the desired position in a gentle manner while not damaging it or any other fragment on the assembly table.

The key contributions of our work are the following

- A model-free RL system that learns to generate smooth pushing paths, with contact-rich pushing actions to reach the object’s target positions in cluttered environments, thereby avoiding contact to other, nearby objects.
- A qualitative and quantitative evaluation in simulation in comparison to a state-of-the-art pushing controller [15], which we adapted to our scenario.

As the experiments with a six degrees of freedom robotic arm show, our system leads to reliable pushing, while achieving better performance compared to [15] with respect to success rate, collisions with other objects, and continuous object contact in various scenarios.

II. RELATED WORK

Without an exact model of the object dynamics, it is hard to predict the moving behavior of objects. Recent work has shown that these dynamics can be learned, e.g., Paus *et al.* [11] proposed a system to predict outcomes of pushing actions and used a graph net representing object relations as in- and output and trained their network on two million synthetic samples. To address the problem that typically a large amount of data is needed to train a network, Nematollahi *et al.* [12] proposed an unsupervised learning strategy and used a combination of an inverse and a forward dynamics model. A restriction of both approaches is that they rely on discrete actions. The resulting pushing actions are typically not contact-rich and need a high number of interactions to guide the object towards the goal.

an RL-learning approach that aims at performing (x, y, z) pushing actions in a continuous way was presented by Xu *et al.* [17] who developed an agent that pushes objects by incorporating their physical properties and ensuring that the objects will not fall over, e.g., in the case of an empty bottle. While the motion is stuttering, due to the physical properties, the approach achieves a good success rate in a clutter-free environment. A promising subdomain of general pushing is planar pushing, where the push action is considered only in 2D. Bauza *et al.* [18] proposed a control model that is learned from only a few data points. The authors evaluated their system against an analytical model on the task of following a given object trajectory. As an extension to this system, Hogan *et al.* [19] trained a classifier that evaluates the current pushing behavior to guarantee a smooth trajectory that is close to the given one. Doshi *et al.* [20] proposed a controller that uses differential dynamic programming to generate the motion model. In all three approaches, the end

effector and the object are both tracked with a motion capture tracking system and no further knowledge about the environment is included. For pushing, the authors assume that each object has four possible contact points. The goal was to reach a desired goal location in a predefined orientation with as few as possible contact switches. The assumption of only four contact points constrains the pushing actions and might not lead to the optimal solution.

In terms of goal-oriented pushing, Bejjani *et al.* [4] proposed an approach to push an object in cluttered environments towards a goal configuration where clutter in the scene is intentionally pushed away to clear the path. In contrast to that, we try to avoid the nearby objects as much as possible to avoid any damages. Furthermore, Migimatsu *et al.* [21] designed a task and motion planning system that do not use global but only relative coordinates to determine the position of the end effector to the target, which leads to higher robustness against changes. We also use only relative coordinates in our observation space. Additionally, to improve the learning behavior of our agent, we apply techniques proposed by Lee *et al.* [22] and Lin *et al.* [13] who included force and touch sensor measurements into the observation space to encourage safe pushing actions and increase the convergence of the agent. In particular, our framework also considers information about the contact between the end effector and the object in the observation space.

Further work towards goal-oriented pushing was proposed by Lloyd *et al.* [2] who used data from a depth camera as well as a tactile sensor. While the authors considered no other objects in the scene, they achieved good results with their control-based approach which is also transferable to curved non-flat surfaces. Krivic *et al.* [15] developed a motion controller for a mobile robot that enables reliable pushing of objects of different shapes on the floor through cluttered environments. In this paper, we use [15] as baseline approach and implemented a modified version for object pushing on a planar workspace with a robotic arm.

III. PROBLEM DESCRIPTION

In this work we consider the following problem. In a tabletop environment, a robotic arm is supposed to move an object from its current position to a 2D goal configuration. To achieve this, we consider the end effector (EE) of the arm moving in planar space (x, y, θ) . The robotic arm can be of any degree of freedom (DOF). In addition to the pushing object, there are other objects which need be considered as obstacles and which might obstruct the direct path to the end configuration. The obstacles have to be avoided by the EE and the object at all time. The goal of the RL-agent is to determine the best incremental movement $(\Delta x, \Delta y, \Delta \theta)$ of its EE position at each time step, to move the object with the EE as fast, but also as safe as possible to the goal position while avoiding obstacles on the way. An RGB-D camera is mounted centered above the scene in bird’s eye view to obtain observations of the objects in the workspace.

¹<https://www.repairproject.eu>

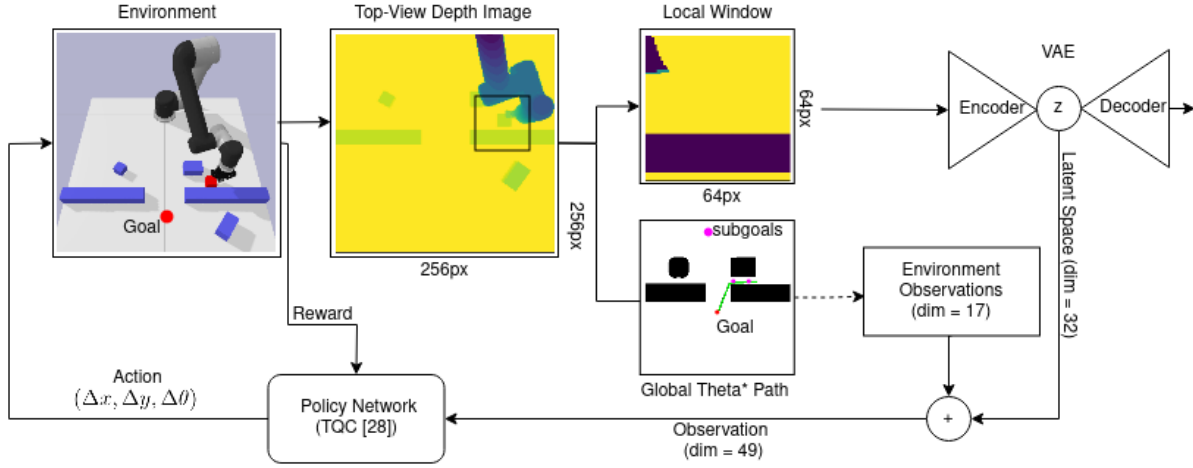


Fig. 2: Overview of our deep reinforcement learning pushing framework. Our system receives a depth image of the environment taken from an RGB-D camera. We calculate an object centered egocentric local window and feed it into the variational auto encoder to get the latent space. Furthermore, a global path from the current object position to the goal position, including subgoals, is calculated. The latent space, the subgoals, and further observations from the environments are used as the concatenated observation for the policy network of the deep RL. The policy network calculates the best 3D incremental motion of the gripper from the observation and the reward it gets from the last environment interaction.

IV. OUR APPROACH

We apply deep reinforcement learning to solve the task described above. This is motivated by the fact that we expect to obtain smoother trajectories as we would get with a pure control-based approach. Especially for traversing narrow passages the lack of parameter tuning can be beneficial. We use a variational auto encoder (VAE) to decouple the feature extraction of the given depth image from the policy learning process [23]. Fig. 2 shows an overview of our proposed system. In the following, we describe the VAE as well as our RL framework in detail.

A. Variational Autoencoder

First, we describe the preprocessing of the input data as well as the network architecture. The networks are implemented and optimized using tensorflow [24].

1) *Preprocessing*: To sense the current world state, i.e., the position of the object and each obstacle, we use a depth image that is gathered from bird’s eye view. To focus on relevant information and ignore distant obstacles that do not influence the next best motion, we use a object centered local window, that is oriented towards the object’s orientation, see Fig. 2. We use a 64x64 pixel window from the original 256x256 image. Since the object’s and the arm’s position are given as individual components in the observation space, we set the corresponding pixels to the background value. We use a convolutional VAE to encode the current normalized local window of the scene into the latent space. We gathered 700k training images in simulation via a random RL-policy and train the network for 10 episodes with a batch size of 256. The dataset contains 10% of blank images to also recognize if no obstacle is around the object.

2) *Network Architecture*: As network architecture we use four convolution layers for the encoder and six deconvolution layers for the decoder. All convolution layers are followed by a batch normalization layer and use the rectified linear unit

function as the layer’s activation function. For the encoder we use max pooling after the first and average pooling after the third layer. The outputs are distributions to directly compute the loss function of our VAE without using any further metric. For the encoder, we use an independent normal distribution and for the decoder an independent Bernoulli distribution. We only use the encoder as part of the observation space, whereas the decoder is ignored. In our experiments, we found that a latent space of 32 works best in terms of training time and feature representation.

B. Reinforcement Learning

RL can be considered as a control problem that can be modeled as a partially observable Markov decision process (POMDP). This means, that the agent cannot determine its exact state s_t at time step t , but has to rely on the current observation o_t to get $s_t \approx o_t(s_{t-1}, a_t)$ for the last state s_{t-1} and the current action a_t . In the end, the goal is to find a stochastic policy $\pi(a_t|o_t)$ that maximizes the expected reward R for each episode, where T is the number of time steps and γ a discount factor.

$$\max \mathbb{E} \left(\sum_{t=0}^T \gamma^t R(s_t, a_t) \right) \quad (1)$$

For the implementation, we followed some ideas proposed by Regier *et al.* [25], which proposed a RL-framework to navigate in cluttered environments with a mobile robot. In the following we define the action and observation space, the reward function, the used RL-algorithm, the experience replay buffer strategy, as well as the learning strategy.

1) *Action Space*: We steer the robot with point control. Therefore, the action space consists of the three values, $(\Delta x, \Delta y, \Delta \theta)$, which are the increment to the current x and y position, as well as the yaw angle θ of the gripper. We

observation	size
Local window latent space	32
EE position at t	5
6D joint angle poses	6
Sub-goal at t-1	2
Sub-goal at t-5	2
Contact with obstacle	1
Object to goal distance	1
Overall:	49

TABLE I: Overview of the observation space.

set the maximum value of $(\Delta x, \Delta y, \Delta \theta)$ to the maximum distance change possible in one predefined time window.

2) *Observation Space*: The observation space of our RL-agent consists of 49 values, as shown in Table I. The EE position is defined as $(x, y, yaw, pitch, roll)$ and given in relative coordinates towards the objects frame. To give the agent an indication of the best path, we include two subgoals, also in relative coordinates, into the observation that we calculate from the current shortest path. The shortest global path is calculated on a binary map, gathered from the depth image, where all obstacles are inflated according to the half of the object’s diameter. Note that the agent never receives the complete shortest global path in its observation and that the shortest path as well as the subgoals are re-calculated at each time step. Therefore, our agent is not constructed as a path following agent but learns the best pushing behavior during training. For the calculation we chose a point after 20% of the global path length of time step t as first subgoal and the subgoal of time step t-1 as second one. We use subgoals from two different time steps to give the agent an indication of the progress it made in any direction. The Boolean value ”contact with obstacle” indicates if the EE touches the object at each time step. During our experiments, we tested different sets of observations and found that the one in Table I leads to the best behavior.

3) *Reward*: Our reward function consists of following three components:

$$r_{dist} = \begin{cases} 50, & \text{if goal reached} \\ -r_{g_dist} - r_{o_dist}, & \text{otherwise} \end{cases} \quad (2)$$

$$r_{collision} = \begin{cases} -10, & \text{if object out of bounds} \\ -5 & \text{if collision occurred} \end{cases} \quad (3)$$

$$r_{touch} = \begin{cases} r_{o_dist}, & \text{contact to object} \\ 0 & \text{otherwise} \end{cases} \quad (4)$$

The first equation encourages the agent toward a faster learning behavior. Therefore, it rewards the agent with a high positive value for accomplishing the task and penalizes higher distances between object and goal as well as object and EE. We use the global path length for r_{g_dist} and the Euclidean distance for r_{o_dist} . To ensure equal rewards through different start goal configurations, we normalize both distances by their initial distance and scale it between 0 and 1. When neither the arm nor the object have moved during a time step and both are at their starting position, the outcome

would be $r_{dist} = -2$. $r_{collision}$ penalizes each collision of the object with clutter in the scene or if the object gets pushed out of the boundaries of the predefined workspace. The last part of the reward r_{touch} considers the suggestion of Lin *et al.* [13] and indicates if the agent has contact with the object. Since we calculate the distance between the EE and the center of the object, a small distance value remains, even if the EE has contact to the object. Therefore, we negate the r_{o_dist} penalty of r_{dist} each time the EE has contact to the object, to encourage a contact-rich behavior.

Together all three parts form the reward function r_{total} of our agent:

$$r_{total} = r_{dist} + r_{collision} + r_{touch} \quad (5)$$

4) *RL-Algorithm*: In this work we use the off-policy algorithm Truncated Quantile Critics (TQC) [28]. During our experiment, TQC led to the best and the most reproducible results. The idea behind TQC is to control the overestimation bias in the critic’s value estimation by using distributional critics [28]. With multiple critics, the points of each of the distributions are used to create a mixture model to increase performance and stability. By truncating the last n points from the mixture of distributions, the overestimation is alleviated. In our work, we use the stable-baselines3 implementation of TQC [29].

5) *Attentive Experience Replay*: The experience replay strategy enables agents to learn from previous experiences they made while interacting with the environment. That means the agent stores action-state transition pairs in a buffer B of size N to reuse previous transitions to update the current policy. In case of the TQC algorithm with a sampling batch size bs the agent uniformly samples bs entries from B and reuses them. While the policy is evolving, some states are more frequently visited than others. For this reason, it is not useful to sample uniformly, because transitions that are rarely visited have a smaller benefit to the update process of the policy than frequently visited ones. Therefore, Sun *et al.* [30] proposed a new strategy to sample an entry from B . With Attentive Experience Replay (AER) they suggest to sample entries according to the similarities between the entry’s state and the current state of the agent. This means that according to the AER strategy, we uniformly sample $k \cdot bs$ entries from B . We then calculate the similarity of each entry to the current state of the agent and use the bs entries with the highest similarity score to update the policy. We use cosine similarity as the similarity measurement, a size of $1e6$ for B , as well as $bs = 512$ and $k = 4$.

6) *Learning Strategy*: As the agent’s learning strategy, we chose curriculum learning, which divides the task into sub-tasks and learns the subtasks one after another in increasing difficulty. We began the training with a maximum start-goal Euclidean distance of 0.06 m and increase it during training to up to 0.6 m. As training environments, we used the scenes shown in Fig. 3. The agent was trained for $7e6$ iterations. Without curriculum learning the agent was not able to learn the task.

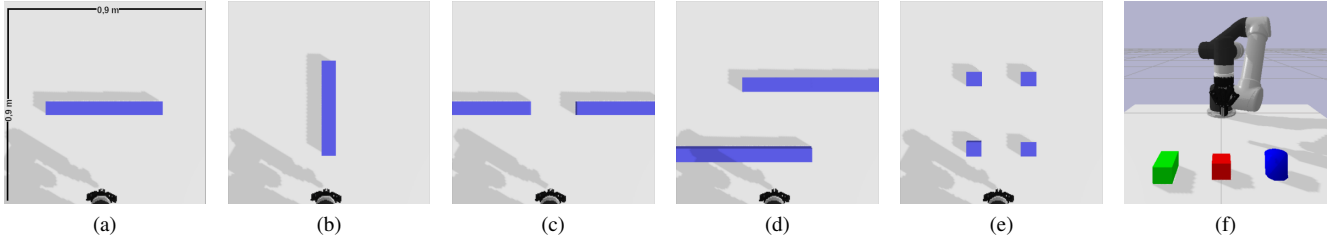


Fig. 3: Figures (a) to (e) depict the different environments used for training and the quantitative evaluation. Figure (f) shows the objects to be pushed. All objects have the same weight but differ in their geometrical shape. As pushing object during training we used the red cube. In a curriculum learning manner, we rotated the obstacle in (a) and (b) and vary its size during the training. Furthermore, the distance between the obstacles in (c) to (e) decreased from 20 cm to 10 cm, making the task more difficult.

Small Cube	Success Rate	Object Contact Rate*	SPL	Path Length
Ours	1.000	0.943 \pm 0.13	0.928	0.429 \pm 0.11
Krivic <i>et al.</i> [15]	0.998	0.870 \pm 0.15	0.918	0.430 \pm 0.10

TABLE II: Quantitative evaluation of straight-line pushing in free space wrt. success rate, object contact, normalized inverse path length (SPL), and path length in meters. The values are the average over 500 runs. The results are in comparison to the approach by Krivic *et al.* [15] where the metrics marked with a * are significant according to the paired t-test with a chosen p-value of 0.05. As shown, our approach performs better in terms of object contact and SPL and equally in terms of success rate.

V. EXPERIMENTS

The goal of our experiments is to demonstrate the performance of our system qualitatively and quantitatively in free space as well as in obstacle-laden environments in terms of success rate, object contact, number of collisions, and shortest path deviation, i.e., normalized inverse path length (SPL) [31]. Furthermore, we provide a comparative evaluation against a state-of-the-art pushing control approach by Krivic *et al.* [15]. We performed the evaluation in pybullet [32] with a 6 DOF UR5² with a Robotiq 2f85 two-finger gripper³. We trained and evaluated our approach on a computer with an i7-6800K six-core CPU at 3.40 GHz and an Nvidia 2070 GPU with 8 GB of memory used for the VAE. For actor and critic, we used a small network with three dense hidden layers of size [512, 256, 128]. For generalization, we used a Gaussian action noise with a standard deviation of 0.4. As global 2D path planner, to sample the sub-goals for the observation space, we used Lazy Theta* [33]. The implementation of our learning framework with all hyperparameters as well as the reimplementation of the baseline approach is available at GitHub⁴.

A. Baseline Approach

To compare our approach to the state of the art for pushing in cluttered environments, we reimplemented the controller proposed by Krivic *et al.* [15]. We implemented the approach as suggested in the paper and used the proposed parameters. Since the original approach was designed for a mobile base and with the assumption that the robot is

always facing the pushing direction, we adapted our implementation to a robotic arm that starts at a position sampled around the object. Due to the sampling, the arm initially dragged the object with it when trying to reposition itself, causing unwanted collisions with obstacles or the object itself. We adjusted this behavior by inverting the pushing direction, once the relocation activation of [15] surpasses a certain threshold $\Psi_{relocate}$. This adjustment enabled the arm to reposition itself more efficiently. The threshold for activating the inversion of the pushing direction was set to $\Psi_{relocate} \geq 0.6$.

B. Quantitative Evaluation

The quantitative evaluation consists of three parts, i.e., pushing in free space, in scenes with obstacles, and in previously unseen, highly cluttered scenes. All metrics except the success rate and the SPL are evaluated only on episodes that both methods could solve successfully. The object contact rate is evaluated for each episode, once the EE first touched the object. Both, object contact rate and collision rate are the average of each episode, averaged over all episodes. For all experiments, we randomly sampled the distance between start and goal within 0.2 to 0.6 m. As pushing object during training we used the red object shown in Fig. 3.

1) *Straight-Line Pushing in Free Space:* We first evaluate straight-line pushing in scenes without other objects to demonstrate the general pushing ability of the two approaches. Therefore, we generated 500 start-goal configurations and compared the results to the shortest path found by Lazy Theta* [33]. As shown in Tab. II, our approach performs slightly better in each metric. Especially, the higher object contact rate shows the impact of our reward function guiding the agent towards the desired contact-rich pushing behavior.

2) *Pushing in Scenes With Obstacles:* Furthermore, we generated five environments which differ in their complexity, as shown in Fig 3 (a) to (e). We used three different types of objects, which were also used in [15], together with the completely unknown complex fragment object shown in Fig. 1, to demonstrate the generalization capabilities. We sampled the orientation and size of the obstacles in (a) and (b) as well as the distance between the obstacles in (c) to (e). In the following we refer to the objects as "small cube" (red), "large cube" (green), and "small cylinder" (blue). For each object,

²<https://www.universal-robots.com/products/ur5-robot/>

³<https://robotiq.com/products/2f85-140-adaptive-robot-gripper>

⁴<https://github.com/NilsDengler/cluttered-pushing>

Small Cube	Success Rate	Object Contact Rate *	Collision Rate *	SPL	Path Length *
Ours	0.980	0.995 \pm 0.02	0.008 \pm 0.04	0.910	0.523 \pm 0.18
Krivic <i>et al.</i> [15]	0.955	0.850 \pm 0.10	0.011 \pm 0.05	0.952	0.513 \pm 0.16

Large Cube	Success Rate	Object Contact Rate *	Collision Rate *	SPL	Path Length *
Ours	0.977	0.995 \pm 0.02	0.007 \pm 0.04	0.910	0.520 \pm 0.18
Krivic <i>et al.</i> [15]	0.957	0.851 \pm 0.10	0.012 \pm 0.05	0.954	0.510 \pm 0.16

Small Cylinder	Success Rate	Object Contact Rate *	Collision Rate *	SPL	Path Length *
Ours	0.967	0.981 \pm 0.05	0.021 \pm 0.06	0.839	0.553 \pm 0.19
Krivic <i>et al.</i> [15]	0.945	0.889 \pm 0.11	0.014 \pm 0.04	0.940	0.512 \pm 0.17

Fragment	Success Rate	Object Contact Rate *	Collision Rate *	SPL	Path Length *
Ours	0.867	0.980 \pm 0.05	0.05 \pm 0.11	0.71	0.630 \pm 0.29
Ours re-trained	0.959	0.984 \pm 0.05	0.01 \pm 0.05	0.83	0.58 \pm 0.23
Krivic <i>et al.</i> [15]	0.953	0.868 \pm 0.11	0.024 \pm 0.07	0.951	0.501 \pm 0.16

TABLE III: Quantitative evaluation wrt. success rate, object contact, collisions, normalized inverse path length (SPL), and path length in meters. The values are the average over 1,000 runs. The results are in comparison to the approach by Krivic *et al.* [15] where the metrics marked with a * are significant according to the paired t-test with a chosen p-value of 0.05. As shown, our approach achieves overall better results in terms of success rate, object contact rate, and collision rate. Please refer to the text for more details.

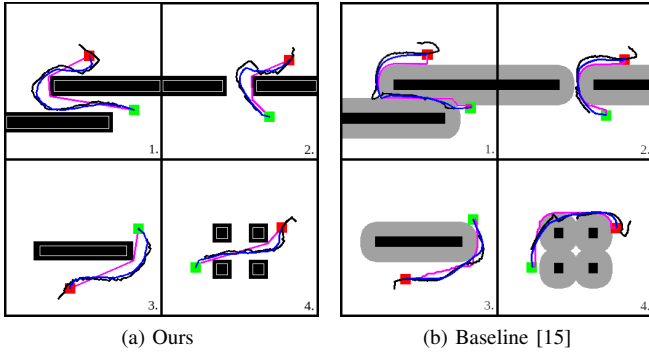


Fig. 4: Qualitative results from the quantitative evaluation of our approach (a) in comparison to the baseline [15] (b). Red indicates the start, green the goal position and magenta is the initial shortest path calculated by Lazy Theta* [33]. The path taken by the end effector is shown in black and the path of the object in blue. The grey area in (b) shows the increased traversal costs around obstacles, used for the baseline approach, while the obstacles in our approach (a) are inflated only by a small amount according to the half of the object's diameter. As can be seen, our agent learned to navigate around objects in a safe distance without strictly following the initial shortest path. Example 4 shows a result where our agent pushes a more efficient path, since it does not rely on any cost map.

we randomly generated 1,000 start-goal configurations within the five environments. As shown in Tab. III, our approach again achieves a significantly higher object contact rate in comparison to the baseline, which shows the benefit of our approach in terms of gentle pushing through contact-rich behavior. Especially for scenarios as in the RePAIR-project¹ gentle, non-abrupt motions are crucial for not damaging any highly fragile objects in the scene. Note that we used different objects for pushing, which the agent never experienced during training. Still, the success rate is consistently high, except for the fragment where our agent still achieves a high success rate, without knowing any dynamics beforehand. In

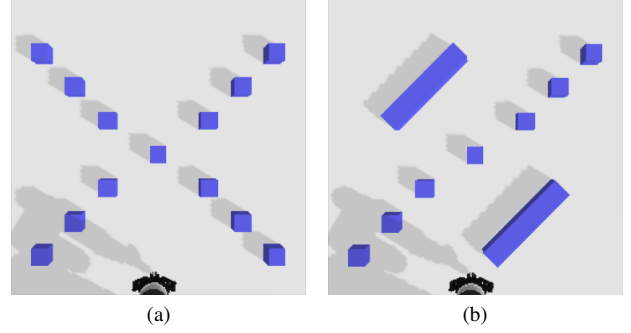


Fig. 5: Unseen, complex environments to further evaluate the performance of our system.

Small Cube	Success Rate	Object Contact Rate *	Collision Rate *	SPL	Path Length
Ours	0.88	0.977 \pm 0.06	0.065 \pm 0.13	0.779	0.492 \pm 0.13
Krivic <i>et al.</i> [15]	0.72	0.566 \pm 0.11	0.01 \pm 0.05	0.720	0.550 \pm 0.18

TABLE IV: Quantitative evaluation in unseen environments with a high density of clutter (Fig. 5) wrt. the success rate, object contact rate, collision rate, the normalized inverse path length (SPL) and the path length in meters in comparison to Krivic *et al.* [15]. The values are the average over 50 runs. The results of the metrics marked with a * are significant according to the paired t-test with a chosen p-value of 0.05. As can be seen, our approach performs better in each metric except the collision rate.

terms of the SPL, the baseline achieves better results while there is no significantly increased path length. This behavior can be explained with the higher obstacle inflation necessary for the baseline approach and is illustrated in Fig. 4 that depicts example trajectories of the experiments. As can be seen, our agent has learned to safely navigate around objects, without strictly following the initial shortest path. This is a key advantage in comparison to the baseline approach, which follows the shortest path as tight as possible due to the properties of the controller method and is crucial if the parameters are not fine-tuned. This explains the lower SPL of our approach. Note that with an increased inflation, similar to the baseline, our SPL results will also increase. Example 4 of Fig. 4 shows a scenario where our agent pushes a more efficient path, since it does not rely on any cost map and therefore on no parameter tuning. As the fragment was never seen during training, we retrained the agent and achieved better overall results as without. This underlines, that our system can be used for serving a general purpose but also retrained to specify on given scenarios.

For all experiments our policy network took on average 0.791 ms for an action prediction of the network and the simulation took 21.91 ms for realizing the action. While our framework has a constant runtime, the runtime of the baseline varies depending on the number of obstacles in the environment and the corresponding greater computational effort.

C. Pushing in Unseen, Complex Environments

Finally, we designed more complex tasks with the goal to evaluate the capabilities of our trained agent in unseen environments with a higher density of clutter. We randomly sam-

pled 50 start-goal configurations of the two scenarios (Fig. 5), which contain many narrow passages. The results in Tab. IV show the good performance in complex and completely unseen environments. Our agent achieved better results than Krivic *et al.* [15] in each metric except the collision rate. Especially the contact rate is significantly increased. As already mentioned, our agent has not been trained on such scenarios, accordingly, the success rate is a bit lower in comparison to the other evaluations with the small cube. Regier *et al.* [25] showed that the success rate will highly increase while the collision rate will decrease, when the agent continues training in the unknown environment for a short time period.

VI. CONCLUSION

In this paper, we presented a novel deep reinforcement learning approach for object pushing in cluttered tabletop environments. We demonstrated the efficacy of our approach in multiple simulated experiments where the results show the increased performance in comparison to an existing control-based method with respect to various metrics. Our agent is able to perform pushing in free space and complex cluttered environments. We showed that the pushing behavior highly benefits from our learning approach in terms of constant object contact and smooth trajectories avoiding obstacles while maintaining equal path length in comparison to the baseline method [15]. The evaluation of the runtime highlights that our system is capable of online pushing. The code of our system can be found on Github⁴ and a video on our web page⁵.

REFERENCES

- [1] M. T. Mason, "Mechanics and planning of manipulator pushing operations," *The International Journal of Robotics Research*, 1986.
- [2] J. Lloyd and N. F. Lepora, "Goal-driven robotic pushing using tactile and proprioceptive feedback," *IEEE Transactions on Robotics*, 2021.
- [3] M. R. Dogar and S. S. Srinivasa, "Push-grasping with dexterous hands: Mechanics and a method," in *Proc. of the IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems (IROS)*. IEEE, 2010.
- [4] W. Bejjani, M. R. Dogar, and M. Leonetti, "Learning physics-based manipulation in clutter: Combining image-based generalization and look-ahead planning," in *Proc. of the IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems (IROS)*. IEEE, 2021.
- [5] A. Cosgun, T. Hermans, V. Emeli, and M. Stilman, "Push planning for object placement on cluttered table surfaces," in *Proc. of the IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems (IROS)*. IEEE, 2011.
- [6] W. Bejjani, "Learning deep policies for physics-based robotic manipulation in cluttered real-world environments," Ph.D. dissertation, University of Leeds, 2021.
- [7] A. Eitel, N. Hauff, and W. Burgard, "Learning to singulate objects using a push proposal network," in *Robotics research*. Springer, 2020, pp. 405–419.
- [8] A. Zeng, S. Song, S. Welker, J. Lee, A. Rodriguez, and T. Funkhouser, "Learning synergies between pushing and grasping with self-supervised deep reinforcement learning," in *Proc. of the IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems (IROS)*. IEEE, 2018.
- [9] M. Ewerton, A. Martínez-González, and J.-M. Odobez, "An efficient image-to-image translation hourglass-based architecture for object pushing policy learning," in *Proc. of the IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems (IROS)*. IEEE, 2021.
- [10] J. Li and D. Hsu, "Push-net: Deep planar pushing for objects with unknown physical properties," in *Proc. of Robotics: Science and Systems (RSS)*, 2018.
- [11] F. Paus, T. Huang, and T. Asfour, "Predicting pushing action effects on spatial object relations by learning internal prediction models," in *Proc. of the IEEE Intl. Conf. on Robotics & Automation (ICRA)*. IEEE, 2020.
- [12] I. Nematollahi, O. Mees, L. Hermann, and W. Burgard, "Hindsight for foresight: Unsupervised structured dynamics models from physical interaction," in *Proc. of the IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems (IROS)*. IEEE, 2020.
- [13] N. Lin, L. Zhang, Y. Chen, Y. Zhu, R. Chen, P. Wu, and X. Chen, "Reinforcement learning for robotic safe control with force sensing," in *WRC Symposium on Advanced Robotics and Automation (WRC SARA)*. IEEE, 2019.
- [14] A. Raffin, "RL baselines3 zoo," <https://github.com/DLR-RM/rl-baselines3-zoo>, 2020.
- [15] S. Krivic and J. Piater, "Pushing corridors for delivering unknown objects with a mobile robot," *Autonomous Robots*, 2019.
- [16] J. Stüber, M. Kopicki, and C. Zito, "Feature-based transfer learning for robotic push manipulation," in *Proc. of the IEEE Intl. Conf. on Robotics & Automation (ICRA)*. IEEE, 2018.
- [17] Z. Xu, W. Yu, A. Herzog, W. Lu, C. Fu, M. Tomizuka, Y. Bai, C. K. Liu, and D. Ho, "Cocoi: Contact-aware online context inference for generalizable non-planar pushing," 2021.
- [18] M. Bauza, F. R. Hogan, and A. Rodriguez, "A data-efficient approach to precise and controlled pushing," in *Conference on Robot Learning*. PMLR, 2018.
- [19] F. R. Hogan, E. R. Grau, and A. Rodriguez, "Reactive planar manipulation with convex hybrid mpc," in *Proc. of the IEEE Intl. Conf. on Robotics & Automation (ICRA)*. IEEE, 2018.
- [20] N. Doshi, F. R. Hogan, and A. Rodriguez, "Hybrid differential dynamic programming for planar manipulation primitives," in *Proc. of the IEEE Intl. Conf. on Robotics & Automation (ICRA)*. IEEE, 2020.
- [21] T. Migimatsu and J. Bohg, "Object-centric task and motion planning in dynamic environments," *IEEE Robotics and Automation Letters (RA-L)*, 2020.
- [22] M. A. Lee, Y. Zhu, P. Zachares, M. Tan, K. Srinivasan, S. Savarese, L. Fei-Fei, A. Garg, and J. Bohg, "Making sense of vision and touch: Learning multimodal representations for contact-rich tasks," *IEEE Transactions on Robotics*, 2020.
- [23] A. Raffin, A. Hill, R. Traoré, T. Lesort, N. Díaz-Rodríguez, and D. Filliat, "Decoupling feature extraction from policy learning: assessing benefits of state representation learning in goal based robotics," *arXiv preprint arXiv:1901.08651*, 2019.
- [24] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, and Z. Chen, "TensorFlow: Large-scale machine learning on heterogeneous systems," 2015, software available from tensorflow.org. [Online]. Available: <https://www.tensorflow.org/>
- [25] P. Regier, L. Gesing, and M. Bennewitz, "Deep reinforcement learning for navigation in cluttered environments," in *Proc. of the Intl. Conf. on Machine Learning and Applications (CMLA)*, 2020.
- [26] 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*, 2018.
- [27] S. Fujimoto, H. Hoof, and D. Meger, "Addressing function approximation error in actor-critic methods," in *International Conference on Machine Learning*, 2018.
- [28] A. Kuznetsov, P. Shvechikov, A. Grishin, and D. Vetrov, "Controlling overestimation bias with truncated mixture of continuous distributional quantile critics," in *International Conference on Machine Learning*, 2020.
- [29] A. Raffin, A. Hill, A. Gleave, A. Kanervisto, M. Ernestus, and N. Dormann, "Stable-baselines3: Reliable reinforcement learning implementations," *Journal of Machine Learning Research*, 2021.
- [30] P. Sun, W. Zhou, and H. Li, "Attentive experience replay," in *Proc. of the Conference on Advancements of Artificial Intelligence (AAAI)*, 2020.
- [31] P. Anderson, A. Chang, D. S. Chaplot, A. Dosovitskiy, S. Gupta, V. Koltun, J. Kosecka, J. Malik, R. Mottaghi, M. Savva, *et al.*, "On evaluation of embodied navigation agents," *arXiv preprint arXiv:1807.06757*, 2018.
- [32] E. Coumans and Y. Bai, "Pybullet, a python module for physics simulation for games, robotics and machine learning," <http://pybullet.org>, 2016–2021.
- [33] A. Nash, S. Koenig, and C. Tovey, "Lazy theta*: Any-angle path planning and path length analysis in 3d," in *Proc. of the Conference on Advancements of Artificial Intelligence (AAAI)*, 2010.

⁵<https://www.hrl.uni-bonn.de/publications/dengler22iros.mp4>