

# Learning a Generative Transition Model for Uncertainty-Aware Robotic Manipulation

Lars Berscheid<sup>1</sup>, Pascal Meißner<sup>2</sup>, and Torsten Kröger<sup>1</sup>

**Abstract**—Robot learning of real-world manipulation tasks remains challenging and time consuming, even though actions are often simplified by single-step manipulation primitives. In order to compensate the removed time dependency, we additionally learn an image-to-image transition model that is able to predict a next state including its uncertainty. We apply this approach to bin picking, the task of emptying a bin using grasping as well as pre-grasping manipulation as fast as possible. The transition model is trained with up to 42 000 pairs of real-world images before and after a manipulation action. Our approach enables two important skills: First, for applications with flange-mounted cameras, picks per hours (PPH) can be increased by around 15% by skipping image measurements. Second, we use the model to plan action sequences ahead of time and optimize time-dependent rewards, e.g. to minimize the number of actions required to empty the bin. We evaluate both improvements with real-robot experiments and achieve over 700 PPH in the YCB Box and Blocks Test.

## I. INTRODUCTION

For real-world manipulation tasks, a robot needs to deal with unknown, stochastic, and contact-rich environments as well as occluded and noisy sensor data. To approach these challenging constraints, a common simplification is to omit time dependency by decomposing actions into single-step, open-loop manipulation primitives. While this is done for classical analytical approaches, it is even more important for learned robotic manipulation. Recent innovations allow robots to learn tasks by interacting with its environment and maximizing the obtained reward. Since data consumption is the fundamental limitation in most learning tasks, the task itself needs to be as simple as possible. Single-step primitives allow to reduce costly training, in particular in the real world.

However, the discarded time dependency is often useful for practical applications. Amongst others, it first allows to plan longer manipulation sequences or second, optimize for time-dependent multiple-step criteria. In this work, we propose to keep learning the manipulation primitives in a single-step and data-efficient manner, while introducing a transition model on top. This visual transition model is learned from pairs of images *before* and *after* an executed manipulation action, and is then able to *predict* the resulting state of an action. As a transition model adds errors into the system, we further predict the uncertainty of the transition model. By estimating the final uncertainty of the manipulation action, the robot is able to manipulate in a *risk-aware* manner.

In this work, we address the task of *bin picking* using a flange-mounted depth camera. It highlights several chal-

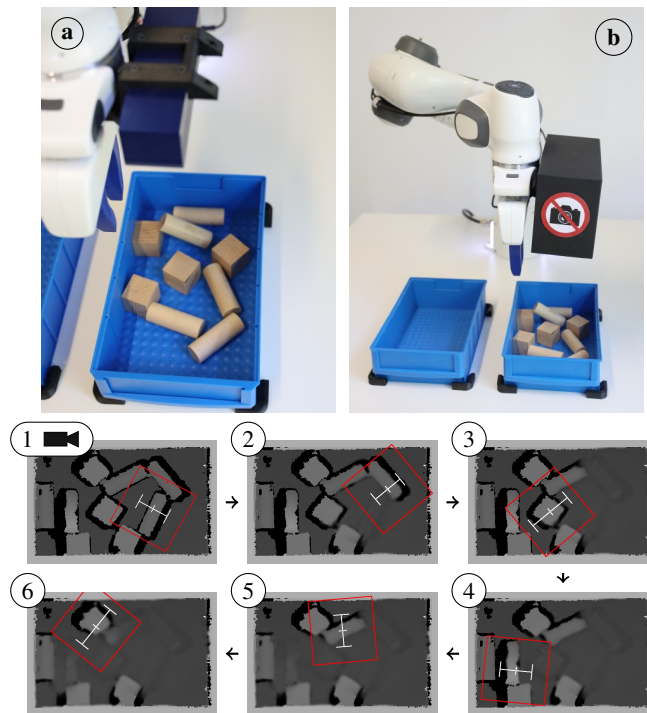


Fig. 1: After a single measurement (a), the robot predicts the next visual state using a learned generative transition model (1-6) and is able to grasp objects blindly considering the predicted uncertainty (b). Image (1) was taken by the flange-mounted depth camera and shows the chosen grasps (white) and the relevant input window for the algorithm (red).<sup>2</sup>

lenges of robotic grasping, e.g. unknown and partially hidden objects, as well as an obstacle-rich environment. On the basis of our prior work, we learn both the task of grasping as well as simple pre-grasping manipulation like shifting or pushing in real-world experiments [1].

We see our contributions as follows: First, we introduce a transition model for predicting both the visual state *and* its pixel-wise uncertainty after a manipulation action. Second, we propagate this uncertainty through the manipulation model. For the latter, we use a common fully-convolutional reward estimation for planar manipulation. The robot is then able to plan actions on the predicted state regarding their uncertainty. Third, we evaluate the novel abilities in various real-world robotic experiments, e.g. by increasing the picks per hour (PPH) in bin picking tasks.

<sup>1</sup>Karlsruhe Institute of Technology (KIT) {lars.berscheid, torsten}@kit.edu

<sup>2</sup>University of Aberdeen pascal.meissner@abdn.ac.uk

<sup>2</sup>Supplementary material is published at <https://pantor.github.io/learning-transition-for-manipulation>

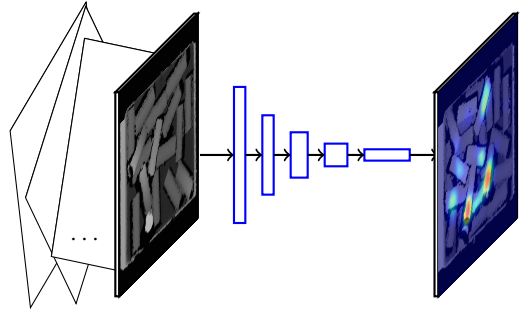
## II. RELATED WORK

Grasping as a foundation for robotic manipulation has been of great interest since the beginning of robotics research. Bohg et al. [2] differentiate between analytical and data-driven approaches. Historically, grasps were synthesized commonly based on analytical constructions of force-closure. For known objects, sampling and ranking based on model-based grasp metrics is still widespread [3]. However, as modeling grasps itself is challenging, it is even harder to do so for pre-grasping manipulation. Moreover, we will focus on the case of manipulation without object model.

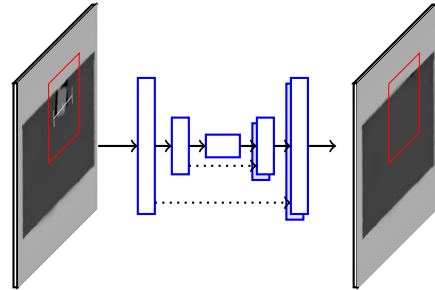
**Learning for manipulation:** In recent years, learning as a purely data-driven method has achieved great progress for manipulation. Kroemer et al. [4] gave a survey about the variety of existing approaches. From our perspective, two major trends have emerged: First, an *end-to-end* approach using a step-wise, velocity-like control of the end effector [5], [6], [7], [8]. Second, the usage of predefined single-step *manipulation primitives* promises a less powerful but more data-efficient learning. This is often combined with planar manipulation and a fully-convolutional neural network (FCNN) as a grasp quality [9] or (more general) an action-value estimator [1], [10]. Similar to *Dex-Net* [9] for bin picking, these approaches map a single image to a single grasp point. Based on the data source, learning for manipulation can be grouped using simulation [8], analytical metrics [9], or real-world interaction [1], [10]. In particular for real-world learning without any transfer gap, data efficiency often limits manipulation to single-step primitives in practice [11].

**Transition model:** In this work, we focus on the *transition model* and its applications for robotic manipulation. Here, model-based reinforcement learning (RL) promises data-efficiency by incorporating a learned transition model into the policy training. While model-based RL have successfully been applied to either non-robotic tasks [12], [13], [14] or low-dimensional non-visual robotic tasks [15], the high-dimensional visual state space hinders adaption for robotic manipulation. Nevertheless, some ideas were explored in this domain: Byravan et al. [16] learned a model predicting rigid body motions from depth images, enabling to generate further rendered images from an updated state model. Boots et al. [17] learned an image-to-image transition model of a moving robot arm. More recently, Finn et al. [18] learned a model estimating pixel flow transformations on images. Combined with model-predictive control, the robot performed nonprehensile manipulation with unknown objects. Similarly, Ebert et al. [19] learned a visual transition model to infer actions to reach an image-specified goal for manipulation tasks.

**Image-to-image translation:** The core task of a visual transition model is image-to-image translation. In recent years, two major approaches emerged in the field of machine learning: First, generative adversarial network (GAN) based methods allow to generate data samples similar to a training set [20]. On top, the *Pix2Pix* architecture by Isola et al. [21] was used in a variety of applications. It is based on a



(a) Manipulation Model  $M$ : The fully-convolutional manipulation neural network (NN) efficiently estimates the reward for a large number of possible actions. For a single rotation and single primitive type, the result can be interpreted as a reward heatmap.



(b) Transition Model  $T$ : The *U-Net* architecture predicts the next image state  $s_{t+1}$  depending on the prior state  $s_t$ , the action  $a_t$ , and its reward  $r_t$ . The result of the predicted window around the action (red) is then patched into the original image.

Fig. 2: We make use of two NNs, one for the *manipulation* policy (a) and one for the *transition* model (b).

conditional GAN combined with a U-Net architecture for the generating neural network (NN). However, it fails to capture stochastic output distributions, as it is a deterministic mapping and very sensitive to mode collapse [21]. Second, variational autoencoder (VAE) based methods try to capture the complete data distribution via a probabilistic latent space. Combining GANs and VAEs, Zhu et al. [22] introduced the *BicycleGAN* that enables multi-modal output distributions and uncertainty estimation. Furthermore, we'll give a brief introduction, but focus on the integration of the BicycleGAN architecture within our robotics application.

## III. LEARNING A TRANSITION MODEL

From Reinforcement learning (RL), we adopt the concept of a Markov decision process (MDP)  $(\mathcal{S}, \mathcal{A}, T^*, r, p_0)$  with the state space  $\mathcal{S}$ , the action space  $\mathcal{A}$ , the transition distribution  $T^*$ , the reward function  $r$  and the initial configuration  $p_0$ . Similar to other data-driven approaches, RL is limited by its data consumption, and even more so for time-dependent tasks resulting in *sparse rewards*. For this reason, we simplify manipulation to a single time step problem. Then, a solution to this MDP is a policy  $\pi : \mathcal{S} \mapsto \mathcal{A}$  mapping the current state  $s \in \mathcal{S}$  to an action  $a \in \mathcal{A}$ . In this work, we explicitly learn a transition model  $T : \mathcal{S} \times \mathcal{A} \mapsto \mathcal{S}$  approximating the true transition distribution  $T^*$ .

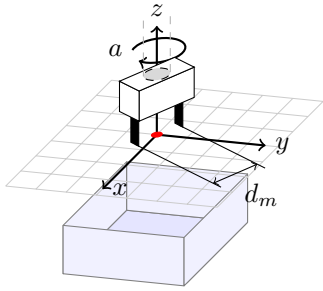


Fig. 3: The manipulation neural network (NN) maps an orthographic depth image to an action type (*grasp* or *shift*) and a planar pose  $(x, y, z, a)$  with gripper width  $d_m$ .

### A. Learning for Robotic Manipulation

Similar to our prior work, we make use of six so-called *manipulation primitives*:

- 1) **Grasps** with four different pre-shaped gripper widths  $d_m$ . Given a grasp point, the robot approaches on a trajectory parallel to its gripper. It then closes the gripper, (possibly) lifts an object, and measures the grasp success by its internal force sensors. We define a binary reward  $r$  for these primitives equal to the grasp success.
- 2) **Shifts** in two different, horizontal directions (relative to the gripper). The robot approaches a given manipulation point, and moves its closed gripper for a few cm, leading to a shifting or pushing motion. The reward is defined by the difference of the maximal grasp reward before and after the action within a pre-defined neighborhood of the action.

We simplify the action space to planar manipulation. Then, each action is described by five parameters: The spatial coordinates  $(x, y, z, a)$  (as shown by Fig. 3) and the manipulation primitive type  $p_t$ .

Let  $s_t \in \mathcal{S}$  be an orthographic image of the scene at time step  $t$ . We align the task-space coordinate system with the image frame by using a top-down view of the scene. This way, orthographic images have a key advantage over other projections: Affine transformations of the image  $s$  correspond to planar poses  $(x, y, a)$  of the robot.

We assume that manipulation (often) depends only on *local* information. We introduce the cropped window  $\hat{s} \in \hat{\mathcal{S}} \subset \mathcal{S}$  at the given affine transformation  $(x, y, a)$ . Its side length is slightly larger than the largest object size. We use a resolution of  $(32 \times 32)$  pixels for the image window. Similar to our prior work [1], [11], the system is learned in a self-supervised manner to maximize the corresponding reward  $r$ . Therefore, a manipulation NN called  $M$  is trained to predict the reward  $r$  of the image windows  $\hat{s}_t$  with given primitive  $p_t$ . We make use of a fully-convolutional NN: During inference, the same manipulation model  $M(s_t)$  estimates rewards  $\psi_t$  for a large number of poses efficiently in a sliding-window approach (Fig. 2a). Rewards for different primitives  $p_t$  are implemented as multiple channels in the output layer. The

height  $z$  is calculated trivially by a model-based controller. The final action  $a_t$  is selected greedily by the policy  $\pi$ .

### B. Simplifications for the Transition Model

Similar to the local simplifications for manipulation, the transition model  $T$  learns only to predict a window  $\hat{s}$  around the action pose. While the window size could in principle match to the overall image  $s$ , we limit the window size to twice the maximal object size. This way, we presume a principle of locality for the state transitions, but it allows to remove all spatial dependencies from the transition model. To estimate the distribution over  $s_{t+1}$ , the state image  $s_t$  is patched with  $\hat{s}_{t+1}$  at the affine transformation  $(x, y, a)$

$$p(s_{t+1}) = T(s_t, a_t) \approx T(\hat{s}_t, p_t, r_t) \quad (1)$$

with the manipulation primitive type  $p_t$ . Additionally, we condition the transition model on the reward  $r_t$  of the action. This is motivated by following arguments:

- By using the estimated reward  $\psi_t$  instead, the calculations of this difficult and high-level feature can be outsourced from the transition model to the already learned manipulation NN. This is in particular important for grasp actions: Here, the next states depends strongly on the binary grasp success. Therefore, we train the manipulation model to learn the grasp success explicitly (as it is measurable) and reuse this information in the transition model.
- Real reward measurements  $r_t$  can be used in the transition model. For example, the grasp reward can be quickly measured by the robot's gripper.
- Note that there is no loss of generality by condition the transition model on the reward, as we can substitute the manipulation model as a reward predictor.

If not stated otherwise, we simplify the transition model to

$$p(s_{t+1}) \approx T(\hat{s}_t, p_t, M(\hat{s}_t)) \quad (2)$$

Let  $\hat{s}'_{t+1}$  denote the prediction with highest probability.

### C. Generative Transition Model Architecture

We employ the BicycleGAN architecture for our transition model [22]. In a nutshell, the architecture is based on a generator, discriminator, and an encoder NN. The generator tries to mimic images from the underlying distribution, taking a noise vector  $z$  from a prior latent distribution as input. While playing a min-max-game, the discriminator aims to distinguish the generated data from the real data. The BicycleGAN introduces an encoder to calculate a representation from an image within the latent space  $z$ . By constraining the latent distribution between the real and generated images, as well as a similarity to a Gaussian, the architecture is able to output a multi-modal distribution. This way, the BicycleGAN architecture is in principle able to capture the stochastic physics as a transition distribution.

We use the generator NN as the transition model  $T$ . It uses a *U-Net* architecture with  $(64 \times 64)$  pixel input and output size (Fig. 2b). We extend the BicycleGAN architecture

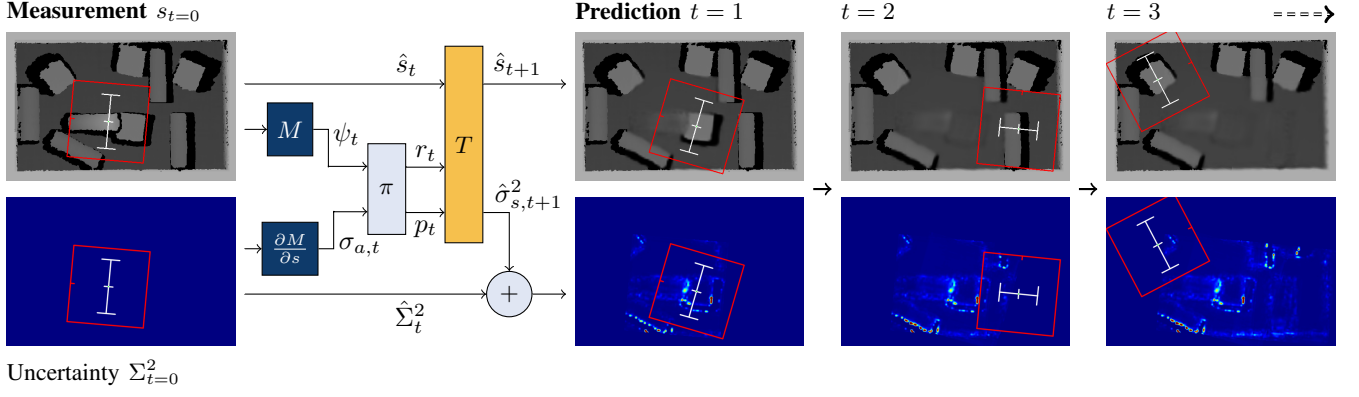


Fig. 4: Iterative multiple-step predictions of the state  $s_t$  with corresponding uncertainty  $\Sigma_t^2$  given an initial measurement  $s_{t=0}$ . The computational flow during a prediction step is as follows: The manipulation neural network (NN)  $M$  estimates rewards  $\psi_t$  for a discrete set of action primitives  $p_t$ , its gradients are used to propagate the cumulative input uncertainty  $\Sigma_t^2$  towards the estimated reward  $\sigma_{a,t}$ . The policy  $\pi$  chooses the final action  $a$  with corresponding reward  $r$  using the lower bound of the estimated reward. Then, the transition model  $T$  predicts the new state  $\hat{s}_{t+1}$  around the action  $a$  with pixel-wise uncertainties  $\hat{\sigma}_{s,t+1}^2$ . The uncertainties are summed up ranging from low (blue) to high (red).

by additionally condition the generator on the manipulation primitive type  $p_t$  (as a one-hot encoded vector) and the reward  $r_t$ .

#### D. Uncertainty Prediction

Given the image distribution  $p(s_{t+1})$ , we are able to estimate a pixel-wise uncertainty of the prediction by sampling from its latent space  $z$ . Let the uncertainty  $\hat{\sigma}_s^2$  be the variance of the prediction  $\hat{s}'_{t+1}$

$$\hat{\sigma}_s^2 = \mathbb{E} \left[ (\hat{s} - \hat{s}'_{t+1})^2 \right] = \sum_i^N p(\hat{s}_i) (\hat{s}_{t+1,i} - \hat{s}'_{t+1})^2 \quad (3)$$

given the number of samples  $N$ . The state uncertainty  $\sigma_{s,t+1}^2$  of the overall image is estimated by zero-padding the uncertainty window  $\hat{\sigma}_s^2$  at  $(x, y, a)$ . The state uncertainty  $\sigma_s^2$  can be propagated towards the estimated reward  $\psi_t$  of the manipulation action. A first-order Taylor series of the manipulation NN  $M$  yields

$$\sigma_a^2 \approx \sum_{i,j} \left| \frac{\partial M(\hat{s}'_{t+1})}{\partial s_{i,j}} \right|^2 \sigma_{s;i,j}^2 \quad (4)$$

summing over all pixels  $(i, j)$ . The approximation presumes a zero covariance between individual pixels and only small uncertainties  $\sigma_s^2$ . A key contribution of our work is to implement (4) fully-convolutionally: The gradient of the fully-convolutional manipulation NN regarding the pixel-wise image input is fully-convolutional itself. Then, the gradient is convoluted over the mean prediction and multiplied with the pixel-wise uncertainty. This results in an efficient computation of both the estimated rewards and its uncertainties  $\psi_t \pm \sigma_a$ . Note that we only consider the input uncertainty; other uncertainties, in particular from  $M$  itself, are neglected in this work.

#### E. From Predicting to Planning

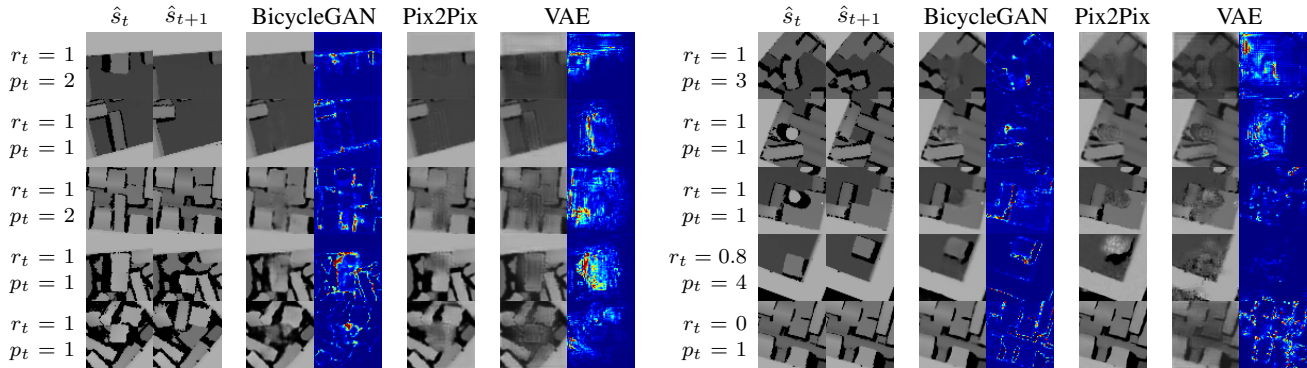
Following, we assume the independence of each state uncertainties  $\hat{\sigma}_{s,t}$  over time  $t$ . Due to the linearity of the variance, we define the cumulative uncertainty

$$\Sigma_t^2 = \sum_{i=0}^t \sigma_{s;i}^2 \quad (5)$$

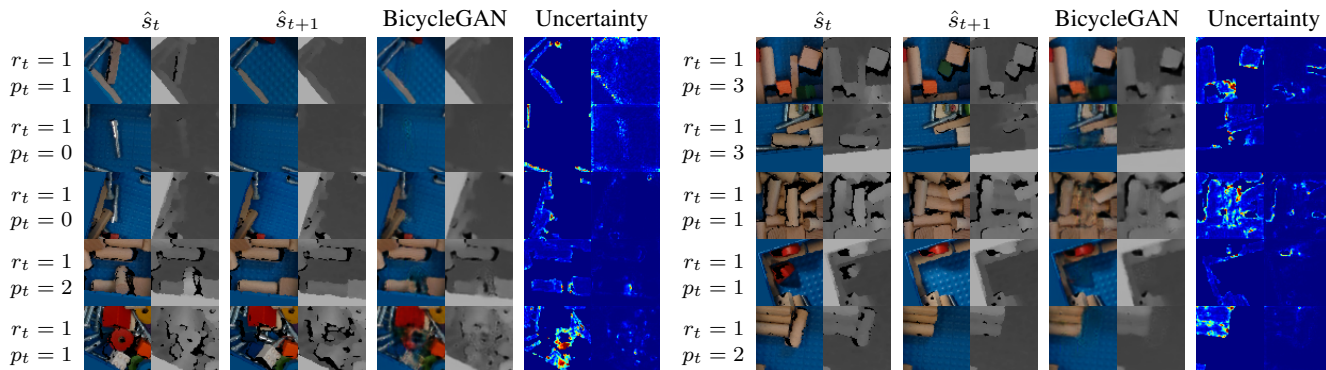
as the accumulated uncertainty of the state. Furthermore, we denote  $t$  as the number of action steps after an image measurement at  $t = 0$ . At  $t = 0$ , we consider the sensor to be perfect and reset the cumulative uncertainty to zero. Moreover, we adapt the policy  $\pi$  to select the action  $a_t$  maximizing the lower confidence bound of the estimated reward  $\psi_t$  given by

$$a_t^* = \arg \max_a \left( M(s_t) - \alpha \left| \frac{\partial M(s_t)}{\partial s} \right| \Sigma_t \right) \quad (6)$$

using the cumulative uncertainty  $\Sigma_t$  in contrast to (4). Let  $\alpha > 0$  denote a parameter for weighing the estimated reward of actions against their confidence. If the uncertainty rises above a threshold  $\sigma_{\text{image}}^2$ , a new image is taken. If the lower bound of the estimated reward falls below a threshold, the bin is assumed to be empty. Furthermore, given both a policy  $\pi$  and a transition model  $T$ , states and actions can be predicted multiple steps ahead by iteratively applying  $\pi$  and  $T$  alternately. Fig. 4 shows the computational flow of a single prediction step, allowing to plan ahead and optimize for multiple-step criteria. As we don't focus on planning algorithms itself, we implemented a simple breadth-first tree search. To allow diverse tree paths, we adapt (6) to sample from the  $N \approx 5$  highest actions uniform randomly. This way, we can predict the overall tree ahead of time and return the action path maximizing the criteria during inference.



(a) Comparison of predictions by the BicycleGAN, Pix2Pix (without uncertainty) and VAE model.



(b) Predictions by the BicycleGAN, for both RGB (not used for downstream manipulation) and depth images.

Fig. 5: Example predictions given an image window  $\hat{s}_t$  before a manipulation action of type  $p_t$  and reward  $r_t$ .  $\hat{s}_{t+1}$  corresponds to the ground truth. The normalized pixel-wise uncertainties are shown from low (blue) to high (red). Action types  $p_t \in \{0, 1, 2, 3\}$  are grasps of different pre-shaped gripper width;  $p_t \in \{4, 5\}$  are shifting motions.

#### IV. EXPERIMENTAL RESULTS

Our experiments are performed on a Franka Panda robot arm using the default gripper and a flange-mounted Ensenso N10 or Framos D435e stereo camera (Fig. 1). We designed and printed custom gripper jaws and attached household anti-slip silicone roll to the robot’s fingertips. Otherwise, the friction of the jaws would not be sufficient for shifting objects. The system uses an Intel Core i7-8700K processor and a NVIDIA GeForce GTX 1070 Ti for computing. We use the *Frankx* library for robot control [23]. Supplementary videos, results, as well as the source code are published at <https://pantor.github.io/learning-transition-for-manipulation>.

##### A. Image Prediction

During around 180h of real-world data collection, the robot attempted 40 000 grasps and 2400 shifts. The robot used two bins with a variety of objects for training. After a successful grasp, the robot moves the objects to the other bin, however returns after filing to take an image *after* the manipulation. After 4000 random actions, the robot used the  $\varepsilon$ -greedy exploration strategy for learning manipulation.  $\varepsilon$  was reduced continuously and reached zero at 30 000 actions, leading to an average reward of  $\bar{r} = 0.55$ . The manipulation

NN was trained around every 500 actions. In the end, we train the BicycleGAN architecture on the complete dataset of state-action-state-tuples  $(s_{t=0}, p_{t=0}, r_{t=0}, s_{t=1})$ . To stabilize the training of the GANs, we smooth the labels and add random noise to the discriminator input. In comparison to the original implementation, we use a training batch size of 32. Otherwise, we keep the training process of the BicycleGAN architecture. Fig. 5 shows example predictions given the manipulation primitive type  $p_t$  and the reward  $r_t$ . We include the ground truth as well as results of the Pix2Pix architecture (without uncertainty estimation) and a VAE for visual comparison. Predicting both a state of  $64 \times 64$  pixels with corresponding uncertainties using  $N = 20$  samples needs around 60 ms.

While our transition model is able to generate realistic images, we find two common failure modes: First, low-level failures like large-scale deviations in the background or blurred edges. Second, wrong high-level features like cropped, blurred, or visual fused objects. The latter case occurs in particular for either shifting actions or occluded layers beneath the grasped object. Regions of high predicted uncertainty can be found primarily near edges, in particular of regions without depth information (black), and within the manipulated objects itself. Examples of iterative predictions

are shown in Fig. 4 and 1.

### B. Grasp Rate

We define the grasp rate as the percentage of successful grasps for grasping 15 objects out of a bin with 25 objects without replacement. For a bin picking scenario with multiple object types, Fig. 6 shows the dependency between the grasp rate and the prediction step  $t$ , equivalent to the number of actions since the last image measurement. The grasp rate

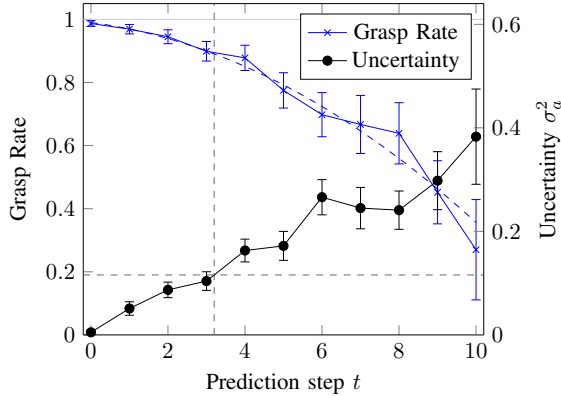


Fig. 6: Grasp rate and uncertainty of the estimated reward  $\sigma_a^2$  depending on the number of prediction steps after an initial image measurement at  $t = 0$ . Based on 1080 grasp attempts in the defined bin picking scenario.

decreases from 98.7% at  $t = 0$  to around 90% for  $t = 3$ . Then, the average uncertainty of the estimated reward  $\sigma_a^2$  increases to around 0.1, resulting in a significant term in (6).

In order to minimize the sum of lost time due to both grasping mistakes and image acquisitions, we estimate an optimal uncertainty threshold  $\sigma_{\text{image}}^2$  given a rough estimate of the execution times of the robot. We identified the durations  $t_G = 6.5\text{ s}$  of the grasp and  $t_I = 2.5\text{ s}$  of the image measurement and related motions. The PPH are then given by

$$\text{PPH} = \frac{\text{Grasp Rate}}{t_G + \frac{1}{t}t_I} \approx \frac{-0.00354t^2 - 0.0228t + 0.987}{6.5 + 2.5t^{-1}} \quad (7)$$

with the optimal prediction step  $t^* \approx 3.2$ . The quadratic regression is plotted in Fig. 6. To optimize the PPH in further experiments, we set the threshold  $\sigma_{\text{image}}^2$  for taking a new image to the corresponding uncertainty  $\sigma_a^2 \approx 0.12$ .

### C. Benchmarking Picks Per Hour

We evaluate our robotic system against the Box and Blocks Test from the YCB benchmark suite [24]. Within 2 min, the robot should grasp and place as many cubes (2.5 cm side length) as possible out of one bin into another. Additionally, we introduce a setting where the robot aims to grasp multiple objects at once. Based on the work of semantic grasping [6], we extend the manipulation NN to predict whether the closed gripper distance will be larger than 4 cm. This corresponds to grasping at least 2 objects; the policy  $\pi$  then maximizes the lower bound of the expected number of

TABLE I: Results for the Box and Blocks Test, including the grasp success rate, the object count after 2 min and the corresponding picks per hour (PPH).

	Grasp Rate	Count	PPH
Heuristic [24]	80%	10.8	324
Random	$(13 \pm 4)\%$	$2.2 \pm 0.7$	$66 \pm 20$
Single	$(97 \pm 2)\%$	$12.8 \pm 0.3$	$384 \pm 10$
Single, Prediction	$(94 \pm 2)\%$	$16.4 \pm 0.5$	$492 \pm 14$
Multiple	$(96 \pm 2)\%$	$20.4 \pm 1.0$	$612 \pm 29$
Multiple, Prediction	$(94 \pm 2)\%$	$23.4 \pm 0.5$	$702 \pm 14$

grasped objects. Table I shows the grasp rate, object count and final PPH for each combination of grasps with/without prediction and of single/multiple objects, in comparison to the heuristic of Calli et al. [24]. Random grasps are evaluated by sampling randomly from the entire action space. We repeated every experiment 5 times. Our robot achieves a peak performance of  $(23.5 \pm 0.5)$  grasps with  $(3.3 \pm 0.2)$  images, corresponding to  $(702 \pm 14)$  PPH. On average, using an approach with prediction improves the PPH by 15%.

### D. Planning

We demonstrate the ability to plan ahead for the task of emptying a bin with as fewest actions as possible. Given an object configuration of three cubes in a row at the side of the bin the robot is able to grasp all objects in 4 steps by shifting the middle cube first. If an outer cube is shifted otherwise, the robot needs at least 5 steps. In this configuration, planning can reduce the average number of action steps from  $4.77 \pm 0.20$  to  $4.14 \pm 0.13$ , measured by 25 successful episodes with a success rate of 78%. Illustrative material can be found on the project website.

### E. Generalization

We qualitatively evaluate the system for objects that were never seen during training (Fig. 7). We find that our system is able to generalize well to compact objects, naturally

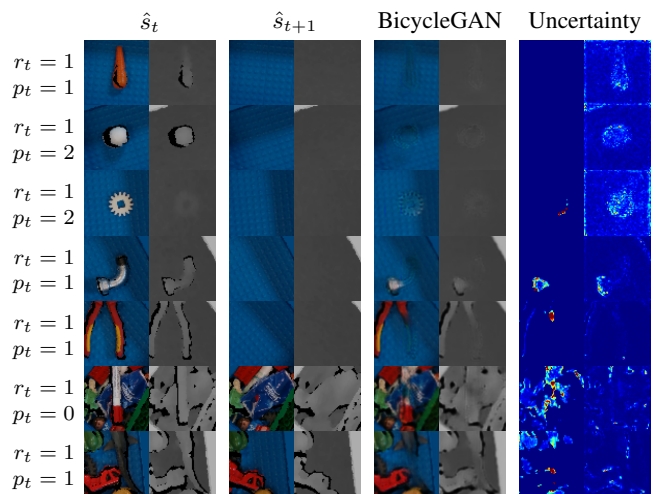


Fig. 7: Example predictions for novel (unknown) objects. The uncertainties are shown from low (blue) to high (red).

depending on the similarity to trained objects. From the perspective of instance segmentation, this problem is very challenging for unknown objects, as the system needs to estimate the connection of objects. For example, the pincers (5th image) are only connected outside of the image. For long and irregular objects, the model often splits single objects into multiple parts. However, it often captures the boundaries for other objects well, e.g. like the red object in the bottom row.

## V. DISCUSSION AND OUTLOOK

We presented an approach to learn a generative transition model for robotic manipulation, for both the next visual state and its pixel-wise uncertainty. This way, the robot is able to plan sequences of single-step manipulation primitives, which were learned data-efficiently without sparse rewards. Furthermore, we demonstrated two more novel skills of our robot: By skipping image measurements, it was able to increase the PPH by around 15% in the Box and Blocks Test [24]. Second, the robot was able to optimize regarding multiple-step criteria in a flexible way, e.g. to minimize the number of actions for emptying a bin. This becomes possible by three major contributions: First, we simplify the transition model by conditioning it on the already learned reward estimation as well as using the spatial invariance of the state space. Second, we propose a model architecture generating samples from a multi-modal distribution, allowing to estimate a prediction uncertainty. Third, the latter was propagated efficiently towards the estimated reward.

Our transition model works in a direct image-to-image setting similar to [17], [18], [19]. In comparison, our model additionally estimates a pixel-wise uncertainty measure and is applied to more difficult, obstacle-rich environments. From a broader perspective, we see several advantages in our approach: In comparison to model-based RL, our transition model does not interact with the single-step policy during training and extends the robotic system *optionally*. This allows to condition the transition model on the estimated reward and reuse the manipulation model learned with model-free RL. Oftentimes, a transition model is used for model-based policy training to reduce sample complexity [14]. However, we find that this is not useful for grasping: The following states depend strongly on the grasp success, which can be measured easily by the gripper and learned explicitly by a manipulation model.

Still, our approach allows for flexible optimization criteria and - even more important for real-world applications - off-policy RL of time-dependent manipulation sequences. In this regard, we learned simple manipulation primitives similar to [10], but were able to combine multiple primitives to more complex action sequences. Although model errors are inevitable for robotic manipulation [18], [19], our uncertainty-aware policy does *actively* avoid these model limitations. However, while our transition model is able to

generate realistic images, we find that it fails to capture high-level features of the stochastic physics. Instead, it usually learns the simplest solution and provides low-level pixel uncertainties.

Eventually, our work leaves room for several improvements. So far, we have ignored the distribution shift between real and predicted images, and in particular its effect on the manipulation model. Despite the principle ability to generalize to unknown objects, we focused our bin picking evaluation on a few object types. We plan to further investigate the limits of generalization in manipulation experiments. Regarding uncertainty estimation, eliminating our sample-based approach could accelerate the algorithm significantly. Moreover, we plan to integrate additional uncertainties, particularly of the manipulation NN, into the policy.

## REFERENCES

- [1] L. Berscheid, P. Meißner, and T. Kröger, "Robot Learning of Shifting Objects for Grasping in Cluttered Environments," in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS 2019)*. IEEE/RSJ, 2019.
- [2] J. Bohg, A. Morales, T. Asfour, and D. Kragic, "Data-Driven Grasp Synthesis - A Survey," *IEEE Transactions on Robotics*, vol. 30, no. 2, pp. 289–309, Apr. 2014.
- [3] A. T. Miller and P. K. Allen, "Grasping! A Versatile Simulator for Robotic Grasping," *IEEE Robotics & Automation Magazine*, vol. 11, no. 4, pp. 110–122, 2004.
- [4] O. Kroemer, S. Niekum, and G. Konidaris, "A review of robot learning for manipulation: Challenges, representations, and algorithms," *arXiv:1907.03146*, 2019.
- [5] S. Levine, P. Pastor, A. Krizhevsky, and D. Quillen, "Learning Hand-Eye Coordination for Robotic Grasping with Large-Scale Data Collection," in *International Symposium on Experimental Robotics*. Springer, 2016, pp. 173–184.
- [6] E. Jang, S. Vijayanarasimhan, P. Pastor, J. Ibarz, and S. Levine, "End-to-end learning of semantic grasping," in *Conference on Robot Learning*, 2017, pp. 119–132.
- [7] 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 *Conference on Robot Learning*, 2018.
- [8] D. Quillen, E. Jang, O. Nachum, C. Finn, J. Ibarz, and S. Levine, "Deep reinforcement learning for vision-based robotic grasping: A simulated comparative evaluation of off-policy methods," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2018, pp. 6284–6291.
- [9] J. Mahler, J. Liang, S. Niyaz, M. Laskey, R. Doan, X. Liu, J. A. Ojea, and K. Goldberg, "Dex-net 2.0: Deep learning to plan robust grasps with synthetic point clouds and analytic grasp metrics," in *Robotics: Science and Systems (RSS)*, 2017.
- [10] 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 *Proceedings of the IEEE International Conference on Intelligent Robots and Systems (IROS)*, 2018.
- [11] L. Berscheid, T. Rühr, and T. Kröger, "Improving Data Efficiency of Self-supervised Learning for Robotic Grasping," in *2019 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2019.
- [12] C. G. Atkeson and J. C. Santamaria, "A comparison of direct and model-based reinforcement learning," in *Proceedings of international conference on robotics and automation*, vol. 4. IEEE, 1997, pp. 3557–3564.
- [13] Ł. Kaiser, M. Babaeizadeh, P. Miłoś, B. Osinski, R. H. Campbell, K. Czechowski, D. Erhan, C. Finn, P. Kozakowski, S. Levine *et al.*, "Model based reinforcement learning for atari," in *International Conference on Learning Representations*, 2019.

- [14] D. Hafner, T. Lillicrap, I. Fischer, R. Villegas, D. Ha, H. Lee, and J. Davidson, "Learning latent dynamics for planning from pixels," in *International Conference on Machine Learning*. PMLR, 2019, pp. 2555–2565.
- [15] P. Abbeel, A. Coates, M. Quigley, and A. Y. Ng, "An application of reinforcement learning to aerobatic helicopter flight," in *Advances in neural information processing systems*, 2007, pp. 1–8.
- [16] A. Byravan and D. Fox, "Se3-nets: Learning rigid body motion using deep neural networks," in *2017 IEEE International Conference on Robotics and Automation (ICRA)*, 2017, pp. 173–180.
- [17] B. Boots, A. Byravan, and D. Fox, "Learning predictive models of a depth camera & manipulator from raw execution traces," in *2014 IEEE International Conference on Robotics and Automation (ICRA)*, 2014, pp. 4021–4028.
- [18] C. Finn and S. Levine, "Deep visual foresight for planning robot motion," in *2017 IEEE International Conference on Robotics and Automation (ICRA)*, May 2017, pp. 2786–2793.
- [19] F. Ebert, C. Finn, S. Dasari, A. Xie, A. Lee, and S. Levine, "Visual foresight: Model-based deep reinforcement learning for vision-based robotic control," *arXiv preprint arXiv:1812.00568*, 2018.
- [20] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in neural information processing systems*, 2014, pp. 2672–2680.
- [21] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-image translation with conditional adversarial networks," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2017, pp. 5967–5976.
- [22] J.-Y. Zhu, R. Zhang, D. Pathak, T. Darrell, A. A. Efros, O. Wang, and E. Shechtman, "Toward multimodal image-to-image translation," in *Advances in neural information processing systems (NIPS)*, 2017, pp. 465–476.
- [23] L. Berscheid and T. Kröger, "Jerk-limited real-time trajectory generation with arbitrary target states," *Robotics: Science and Systems XVII*, 2021.
- [24] B. Calli, A. Walsman, A. Singh, S. Srinivasa, P. Abbeel, and A. M. Dollar, "Benchmarking in manipulation research: Using the yale-cmu-berkeley object and model set," *IEEE Robotics Automation Magazine*, vol. 22, no. 3, pp. 36–52, Sep. 2015.