# Model-free Grasping with Multi-Suction Cup Grippers for Robotic Bin Picking

Philipp Schillinger*, Miroslav Gabriel*, Alexander Kuss*, Hanna Ziesche*, and Ngo Anh Vien*

*Bosch Center for Artificial Intelligence, Renningen, Germany.
Email: firstname.lastname@de.bosch.com

*Abstract*—This paper presents a novel method for model-free prediction of grasp poses for suction grippers with multiple suction cups. Our approach is agnostic to the design of the gripper and does not require gripper-specific training data. In particular, we propose a two-step approach, where first, a neural network predicts pixel-wise grasp quality for an input image to indicate areas that are generally graspable. Second, an optimization step determines the optimal gripper selection and corresponding grasp poses based on configured gripper layouts and activation schemes. In addition, we introduce a method for automated labeling for supervised training of the grasp quality network. Experimental evaluations on a real-world industrial application with bin picking scenes of varying difficulty demonstrate the effectiveness of our method.

## I. INTRODUCTION

Model-free grasping with multi-suction grippers is a key challenge for fully automating many pick-and-place tasks in industry and logistics. Fig. 1 shows a typical robotic bin picking system for warehouse order fullfilment including an industrial robot equipped with a multi-suction gripper, an overhead RGB-D camera and bins containing diverse objects positioned on a conveyor belt. Recent research proposes machine learning methods that enable model-free grasp prediction for a wide variety of unseen objects in unstructured environments [18, 11].

Most approaches focus on grasp prediction for parallel-jaw grippers or single-suction grippers. However, suction grippers with multiple suction cups are rarely studied so far, although they are capable of balancing torques which is favorable for dynamic movements of objects with large dimensions [16]. They also enable lifting of heavier objects without damaging their surfaces by distributing the required grasp force over multiple suction cups. Some multi-suction grippers allow the activation of different suction cup groups to offer flexibility in diverse object portfolios. However, multi-suction grippers are more challenging for grasp prediction approaches to deal with due to their complex geometry and alternative activation patterns.

In this paper, we propose a method for model-free prediction of grasp qualities for suction grippers that is independent of the actual gripper design, e.g., the number and size of suction cups. More specifically, we present a gripper-agnostic grasp quality prediction including a procedure for automatic labeling and supervised training, allowing for generalization to new object shapes. Furthermore, we present



Fig. 1. Robotic bin picking system with six-axis industrial robot, overhead RGB-D camera, tool changer, single-suction gripper, multi-suction gripper, bins containing diverse objects and conveyor belt.

a method for optimal gripper selection and rotation based on the inferred, gripper-agnostic pixel-wise grasp quality prediction combined with footprint images to represent specific gripper design configurations.

Our contributions are as follows: (1) Method for model-free prediction of grasp qualities agnostic to the suction gripper design, including automated labeling for supervised training. (2) Method for optimal gripper selection and orientation by matching of arbitrary suction gripper footprints.

We evaluate contribution (1) by a comparison of our proposed pixel-wise grasp quality prediction with existing methods based on various bin picking scenes of different levels of difficulty. In addition, we demonstrate contribution (2) by real-world experiments of performing multi-suction cup grasp prediction on an industrial bin picking application.

## II. RELATED WORK

Modern robot grasping methods often use deep learning techniques trained on large datasets to predict grasps. Most grasp prediction approaches rely on some form of predicting a pixel-wise grasp quality map that represents the grasp success probability at each pixel. Mahler et al. [15] and Zeng et al. [28] propose to learn a grasp map for suction and parallel-jaw grasps from a supervised dataset using RGB-D or depth as input. Following this approach, GG-CNN [17] and FC-GQ-CNN [22] propose to predict a grasp quality map and a 4 DoF parallel-jaw grasp configuration for each
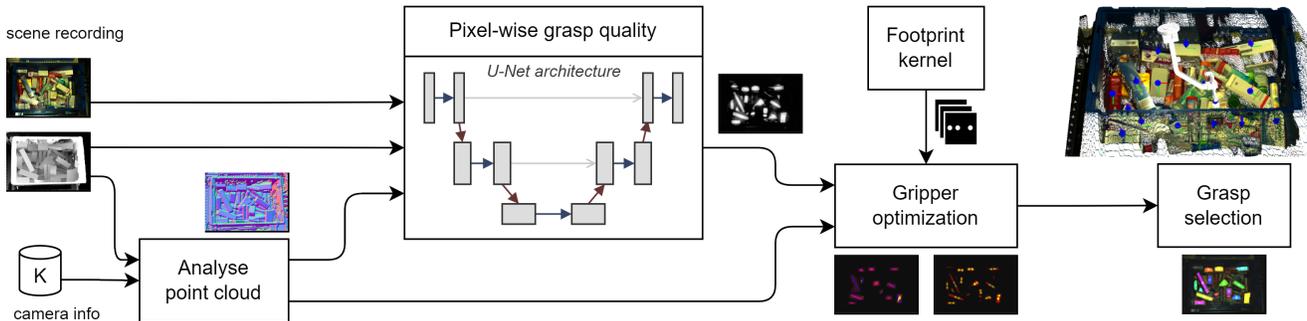
Fig. 2. Summary of our proposed approach. Our method receives an RGB-D scene recording and camera information to derive a pixel-wise grasp quality. Afterwards, by providing gripper geometry information as footprints, a gripper optimization finds the best poses and corresponding grippers.

pixel. Cao et al. [2] propose a pixel-wise grasp map and and a grasp configuration prediction or single-suction grippers. Breyer et al. [1] propose to generate voxel-wise grasps using a truncated signed distance function for a parallel-jaw gripper. A follow-up work optimizes grasp prediction jointly with object shape reconstruction [9]. Other approaches use point clouds as input and predict point-wise grasp qualities and gripper configurations [29, 20, 26, 7, 19, 4, 12].

Various gripper designs are available for robotic grasping. Parallel-jaw and suction grippers are commonly used and effective for regular object sets [3, 6]. In bin picking scenarios, suction grippers have an advantage over parallel-jaw grippers for object reachability. Dex-Net 3.0 [14] improves success rates of Dex-Net 2.0 [15] by training a grasp quality neural network specifically for suction grippers. Shao et al. [24] propose a self-supervised learning method for simulated suction-based picking. Suctionnet-1billion predicts grasps for single-suction grippers through end-to-end training of a pixel-wise prediction network [2]. Jiang et al. [8] jointly learn pixel-wise grasp quality and robot reachability maps for suction vacuum cups. Zeng et al. [27], winners of the Amazon picking challenge, propose learning a pixel-wise grasp map for a hybrid gripper combining parallel-jaw and suction cup functions.

Although there are advanced multi-suction gripper designs in industrial products [13, 16], there is little research on explicitly modelling and using them for robotic grasping. Recent efforts focus on optimizing a single network for the prediction of grasps for different gripper types [10, 23, 25]. However, no prior work has explored learning a pixel-wise grasp map that can be used for both single- and multi-suction grippers without altering the network architecture or the training process.

## III. PROBLEM STATEMENT

Given an image of the scene and a set of multi-suction grippers, our goal is to predict a set of feasible grasp poses which can be used to transfer arbitrary objects from a source to a target bin. In particular, our method receives an RGB-D image as input and infers grasp poses from a multi-channel grasp map where each channel per gripper type encodes pixel-wise grasp quality and rotations.

We denote by $S$ an RGB-D scene image of size $\mathbb{R}^{4 \times h \times w}$, and by $g = (x, y, z, \alpha, \beta, \gamma, t)$ a grasp configuration predicted at position $(x, y, z)$ with orientation $(\alpha, \beta, \gamma)$ and gripper $t \in T$. As grippers, we assume that there is a set $T$ of different types of multi-suction grippers that can be used for a grasp. The problem of predicting multi-suction grasps is to find a mapping $f : S \mapsto g$ for every input image $S$.

We propose a two-step approach for solving this problem as illustrated in Fig. 2. First, a neural network predicts a pixel-wise "graspability" property for the image, denoting how well each individual pixel can be grasped and in the following referred to as grasp quality. Second, an optimization step determines the best gripper and corresponding grasp pose based on the predicted grasp quality.

In contrast to methods that directly approximate $f$, this two-step approach has the benefit that no gripper-specific training data is needed. To the best of our knowledge, a grasp dataset of multi-suction cups does not exist in literature and, in particular, not for our specific multi-suction cup grippers.

## IV. MODEL-FREE GRASP QUALITY PREDICTION

The first part of our proposed method consists of a generic, pixel-wise prediction of graspable surfaces. This prediction can be obtained for a wide range of unknown objects and does not require gripper geometry-specific information. As input, we use high-resolution RGB-D images, e.g. from industry-grade cameras like *Zivid Two* or *Photoneo PhoXi*, and assume that the intrinsic camera parameters are known. The output is a pixel-wise grasp quality prediction $Q$ where each pixel ranges from 0 (not graspable) to 1 (perfectly graspable) and indicates how suitable the respective pixel is for attaching a suction cup.

### A. Grasp Quality Inference

We use a U-Net [21] architecture with a ResNet-34 [5] encoder and a single-channel output $Q$ to infer the pixel-wise grasp quality. The network input is a three-channel image, $I := (S_{\text{Gray}}, S_{\text{Depth}}, S_{\text{Std}})$, which consists of grayscale $S_{\text{Gray}}$, depth $S_{\text{Depth}}$, and the standard deviation of the surface normal vectors $S_{\text{Std}}$. Fig. 3 shows an example of each channel and the resulting network output $Q$.

Using a single grayscale channel instead of three RGB channels largely retains texture information, but reduces the
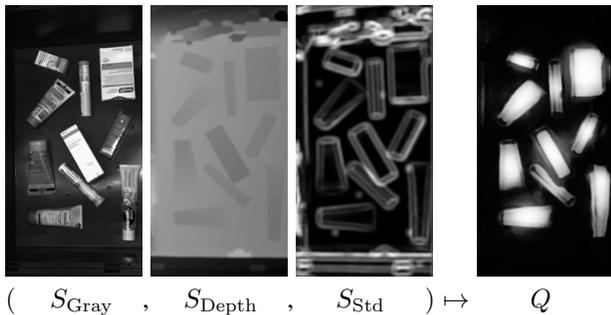
$$( \quad S_{\text{Gray}} \quad , \quad S_{\text{Depth}} \quad , \quad S_{\text{Std}} \quad ) \mapsto \quad Q$$

Fig. 3. Input to the grasp quality network is given by a grayscale channel $S_{\text{Gray}}$, a depth channel $S_{\text{Depth}}$, and the standard deviations of surface normals $S_{\text{Std}}$. Output is a single-channel rating $Q$ how well a suction gripper can be attached to each pixel.

number of channels that provide this information and, more application-specific, prevents the network from overfitting to colored bins as background. Using $S_{\text{Std}}$ as a third channel is motivated by the fact that graspability for suction grippers highly depends on the local surface structure. If the surface is very irregular, a suction cup is less likely to form a sealed vacuum at that point.

To obtain $S_{\text{Std}}$, we first calculate the ordered point cloud from $S_{\text{Depth}}$ and the configured camera matrix $K$ for each pixel $(u, v) \in h \times w$ as

$$S_{\text{Pts}}(u, v) = K^{-1}\big(S_{\text{Depth}}(u, v)[u, v, 1]^{\text{T}}\big)$$

and derive the pixel-wise surface normals, $S_{\text{Normals}}$. Then, $S_{\text{Std}}$ is computed for a small neighborhood and normalized to a value range $S_{\text{Std}}(u, v) \in [0, 1], \forall u, v$.

One practical challenge in calculating $S_{\text{Normals}}$ is that this calculation is susceptible to errors and inaccuracies in the depth image. To address this, we employ two pre-processing steps where the first one fills missing pixels with depth approximations and the second one reduces noise resulting from outlier pixels.

### B. Labeling and Supervised Training

For supervised training of the grasp quality network, we require a dataset that contains RGB-D input data $S$, as well as pixel-wise ground-truth for grasp quality $Q^*$. Annotation of $Q^*$ requires a high effort if done manually for cluttered scenes with many objects or complex geometry. We therefore choose an alternative approach for obtaining approximate labels, $L$, such that $L(u, v) \approx Q^*(u, v)$.

This automatic labeling approach is based on the insight that for singulated objects with simple geometries, the suitability for a pixel to be grasped inversely correlates with the $S_{\text{Std}}$, similar to our initial motivation for using $S_{\text{Std}}$ as an input to the network. Consequently, we calculate one component of the labels by $L_{\text{Std}} := 1 - S_{\text{Std}}$.

Furthermore, we expect grasps closer to the center of mass of an object to be more stable and thus, these pixels should receive a higher grasp quality. To include this property in the training labels, we cluster pixels of graspable surfaces as given by $L_{\text{Std}}$ and calculate a second component of the labels $L_{\text{Dist}}$ as the distance to the respective cluster center.

One limitation of this labeling method is that also the bin and other non-object geometries of the scenes may be considered as graspable areas, solely depending on their surface. For training, we can circumvent this issue by recording a background image of an empty bin, i.e., a scene recording without any objects before placing the training objects in the scene. We then use background subtraction based on depth to mask all non-object pixels $M_{\text{bg}}$ and only consider non-zero labels for object pixels.

The grasp quality labels $L$ for training are thus given by

$$L = \begin{cases} w_{\text{Std}}L_{\text{Std}} + w_{\text{Dist}}L_{\text{Dist}} & \text{where } M_{\text{bg}} > 0 \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

where the weights $w_{\text{Std}}$ and $w_{\text{Dist}}$ can balance the influence of the different label components, but are chosen to be equal in our experiments. The quality metric defined in Eq. (1) distinguishes itself from the approach presented in [8] by using an unnormalized $L_{\text{Dist}}$ score and excluding the residual error of local plane fitting at each pixel, making Eq. (1) more computationally efficient.

Fig. 4 shows three examples from our training dataset. While we observe that this labeling approach works sufficiently well for objects with simple geometries in scenes with only a few instances, the approximation $L(u, v) \approx Q^*(u, v)$ becomes significantly worse for complex geometries and scenes with object instances that are close together. Consequently, we limit training data to such simpler scenes that allow for a good approximation.

The approach can be extended to scenes with a larger number of objects or overlapping objects if instance mask annotations are available, which are much easier to obtain by manual annotation than grasp quality labels. In that case, a background image and object clustering is not required since the background is given by all pixels which are not included in any of the object masks and all other labeling steps can be performed in the same way, including the computation of $S_{\text{Std}}$ and the final labels according to Eq. (1).

During training of the grasp quality prediction, the loss $\mathcal{L}$ for some predicted grasp quality $Q$ and target labels $L$ is given by a pixel-wise mean-squared error

$$\mathcal{L} = w_{\text{bg}}MSE\big(M_{\text{bg}} \circ E\big) + w_{\text{fg}}MSE\big((1 - M_{\text{bg}}) \circ E\big) \quad (2)$$
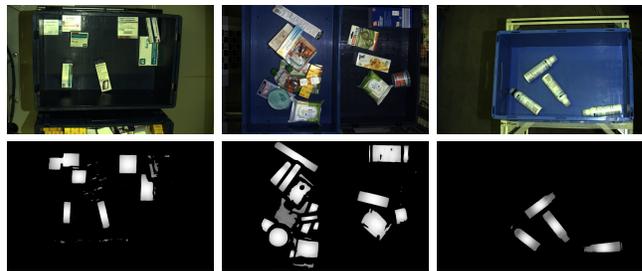


Fig. 4. Three training examples, RGB input shown top and generated labels $L$ shown in the bottom row. Labels become worse for more complex geometries. Missing labels mainly result from invalid depth information.
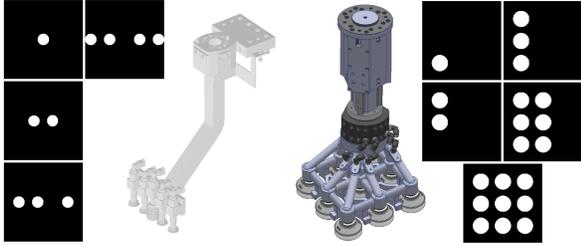
Fig. 5. Example footprints of two different multi-suction grippers that both provide multiple activation patterns to choose from. White areas denote full surface contact, black areas mean no contact.

for some prediction error $E := L - Q$ and $\circ$ denoting pixel-wise multiplication. The background mask $M_{bg}$ balances the loss received for on-object pixels with the one for background pixels with weights $w_{\text{bg}}$ and $w_{\text{fg}}$ calculated per image to reflect the ratio of background and foreground.

For the experiments in this paper, we generated a dataset of around 2,000 proprietary recordings of bin picking scenes collected across various robotic cells with a mixed object portfolio similar to those shown in Fig. 4. Training was then performed for 30 epochs on a *Nvidia V100* GPU with a batch size of 16 and images being down-scaled to a resolution of $1280 \times 800$ pixels. We used stochastic gradient descent with an initial learning rate of $10^{-4}$ and a cosine annealing schedule implemented in *PyTorch*.

## V. GRASP POSE DETECTION

The second part of our proposed method is deriving the full grasp pose from the pixel-wise grasp quality prediction described in Sec. IV. For this we assume minor application knowledge about the type of available grippers which is manually specified as gripper footprints. This does not include further scene understanding or context knowledge such as robot kinematics, bin dimensions, or object models.

### A. Gripper Footprint Matching

We propose a gripper selection and matching based on specified gripper footprints, such as the ones shown in Fig. 5. In this work, we assume that a footprint is always centered at the end-effector pose and that the footprint size is scaled to match the correct pixel-per-mm resolution, for example in our experiment application around two pixels per mm.

To identify grasp poses, including a selection of the best gripper and its orientation, we perform a convolution over the inferred grasp quality $Q$. For this, $n_{\text{r}}$ different discrete rotation steps of $n_{\text{f}}$ different gripper footprints are encoded as separate channels in one combined convolution kernel $F \in \mathbb{R}^{n_{\text{r}} \cdot n_{\text{f}} \times h_F \times w_F}$ of size $h_F \times w_F$. The result of a convolution of $Q$ with $F$ is thus a multi-channel pixel-wise prediction how well each gripper type in each rotation can perform a successful grasp at the respective pixel.

However, accumulating grasp quality like this leaves one issue which we denote by the term "edge wrapping": Consider a box which can be grasped well on two sides with different orientations, but not at the edge between these sides. Using a gripper footprint with two suction

cups that have enough space between the cups might now match one suction cup to one side and the other cup to the other side. This would result in a grasp with a high theoretical graspability but which will fail in practice. A similar example for edge wrapping is shown in Fig. 6.

Edge wrapping again motivates the use of normal vectors for avoiding infeasible grasp proposals. We therefore perform another convolution with the same kernel $F$ over the three-channel image of normal vectors $S_{\text{Normals}}$ to compute the standard deviation of normal vectors for each of the respective gripper footprint areas as

$$F_{\text{Std}} := \left| F * S_{\text{Normals}}^2 - (F * S_{\text{Normals}})^2 \right|^{\frac{1}{2}}. \quad (3)$$

The product of the accumulated grasp quality $F * Q$ with the inverse of the above standard deviation $F_{\text{Std}}$ then denotes the grasp feasibility.

Finally, we compute three single-channel pixel-wise results of this operation. The first result $O_{\text{Type}}$ gives the pixel-wise gripper type and the second result $O_{\text{Rot}}$ denotes the respective gripper rotation, both given by the pixel-wise argmax over all $n_{\text{r}} \cdot n_{\text{f}}$ channels of the convolution result. The last result is given by a pixel-wise grasp quality $O_{\text{Q}}$, calculated as the max over all channels of the convolution result. This is similar to the previously computed grasp quality $Q$, but $O_{\text{Q}}$ now denotes for each pixel how feasible the best possible grasp configuration would be at that pixel in contrast to how well a pixel is suitable for being grasped.

### B. Pixel-to-Pose Transformation

To obtain a list of grasps from the previous pixel-wise results, grasp quality values $Q$ are clustered such that clusters correspond to graspable areas of objects and pixels near the cluster center with highest $O_{\text{Q}}$ are selected for grasping. This has shown to improve robustness compared to directly using the pixels with highest values and ensures that grasps are detected for all objects that receive high (but not necessarily the highest) grasp quality values in a scene.

For each selected pixel $(u, v)$, the grasp pose position is then given by the respective point in the ordered point cloud $x, y, z := S_{\text{Pts}}(u, v)$. The roll and pitch rotations (around the $x$- and $y$-axes of the gripper) are determined by the surface of the object and are thus given by the negative normal vector at the respective pixel $\alpha, \beta := -S_{\text{Normals}}(u, v)$. The yaw rotation around the gripper axis is given by



Fig. 6. Example for the edge wrapping issue. For the purple grasp in the left image, a large three-suction cup gripper is incorrectly selected and ignores scene geometry. Instead, a smaller footprint aligned with the object surface would be correct and is the result when applying the proposed method, as shown in the right image.
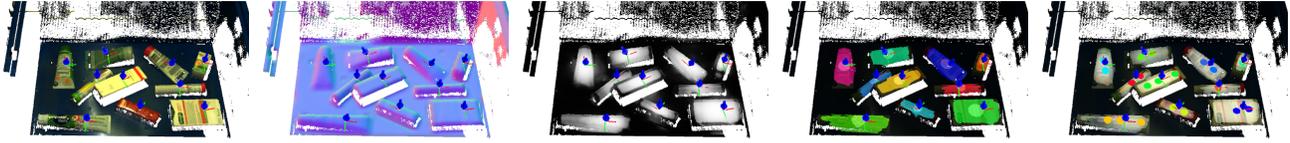
Fig. 7. Summary of intermediate steps of grasp detection on an example scene from an application (outside of the training distribution). From left to right: RGB input, calculated surface normals, inferred grasp quality, clustered graspable objects, and matched footprints.

the gripper rotation determined during footprint matching $\gamma := O_{\mathrm{Rot}}(u, v)$. Finally, the gripper type of the grasp is given by the determined footprint $t = O_{\mathrm{Type}}(u, v)$.

## VI. EXPERIMENTS

We perform two different types of experiments to evaluate the performance of our approach. First, we evaluate grasp quality prediction performance by detecting single-suction grasps for an evaluation dataset and compare it with related methods. Second, we demonstrate our multi-suction grasp detection on an industrial bin picking robot cell.

### A. Single-Suction Grasp Comparison

Due to the lack of a comparable multi-suction gripper work, we determine single-suction grasps on an annotated reference dataset from the target bin picking cell and quantify the performance from the given pixel-wise ground-truth grasp success. We compare our work to the following two methods for single-suction grasp detection:

**Dex-Net.** Satish et al. [22] propose a fully convolutional grasp quality CNN (FC-GQ-CNN) for grasp prediction. We use the pretrained model FCGQCNN-4.0-SUCTION provided by the authors, which was trained on synthetic depth images of objects in clutter with parameters for a Photoneo PhoXi S camera. Published values are based on a very specific combination and placement of RGB-D camera and gripper. Therefore, we apply static cropping close to the bounding box of the bin and a depth value shift to make the depth images similar to the provided example images of the authors. Input images are resized to $640 \times 480$.

**Zeng et al.** [28] introduce a multi-modal grasping framework. They use a separate FCN with residual connections that generates a suctionability map for each pixel. The framework combines RGB and depth data of size $640 \times 480$ using two pre-trained ResNet-101 towers, then concatenates the output features to predict suction affordances. In this paper, we evaluate the framework using two variations: First, we replicated and trained the network with the original dataset of the authors. Second, we retrained on a combination of the authors' dataset and our own dataset with labels and loss as in Sec. IV-B.

The dataset for evaluation is created from eleven types of common objects, see Fig. 8 for examples. For detailed results, we split the dataset into three levels of difficulty. **Simple:** Bins filled with a small number of objects such as boxes or cylinders. **Typical:** Bins containing a heap of a larger number of the objects. **Complex:** Bins filled with a large number of challenging objects like partially transparent objects, stacked and texture-less boxes, or blister packs.

Each category consists of ten RGB-D images and corresponding ground-truth pixel-wise grasp sucess. The grasp success is based on manually annotated object instance masks, combined with a similar method as described in Sec. IV-B for consideration of surface and weight, but with higher emphasis on the distance to the center. To achieve meaningful results, we manually reviewed and corrected the grasp quality to ensure reasonable ground-truth labels.

For each scene, each method is queried for the top 20 grasps based on their grasp quality prediction. Results are listed in Tab. I. We state for each category as *Quality* the average ground-truth grasp quality of feasible grasps (between $0$ and $1$), as *Success* the percentage of feasible grasps among predicted grasps (non-zero grasp quality), as *Objects* the percentage of objects for which grasps are detected (counting at most 20 objects per scene due to top 20 grasps), as *Multi* the average number of grasps predicted for the same object (closer to one is better), and as *None* the percentage of scenes without a single feasible grasp.

It can be seen that our method performs well for the considered application, which motivates its use for multi-suction grasp detection. In particular, it can be observed that our method spreads grasps across objects compared to Dex-Net and Zeng et al., which often predict multiple alternatives for the same object. Especially Zeng et al. achieve a presumably high success rate, but do so by detecting multiple close grasps per object which has limited usefulness in practice. One contribution to this difference likely is the clustering of surface graspability for deriving poses (instead of directly predicting grasp success for pixels), as important for robustly placing multi-suction grippers.

We also achieve a high average grasp quality and while this does not directly indicate the feasibility for multi-suction grasps, it suggests that considering a larger footprint
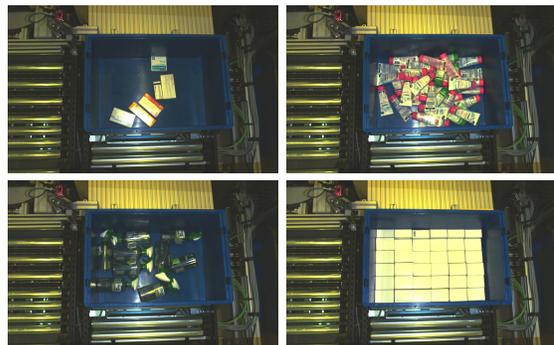


Fig. 8. Examples for categories Simple (*top left*), Typical (*top right*) and Complex (*bottom row*) of the evaluation dataset.

| Simple | Quality | Success | Objects | Multi | None |
|---|---|---|---|---|---|
| Dex-Net | 0.810 | 78.0% | 59.3% | 4.86 | **0%** |
| Zeng et al. | 0.602 | 68.5% | 21.9% | 14.33 | 10% |
| — finetuned | 0.610 | **86.0%** | 22.8% | 15.78 | 10% |
| Ours | **0.894** | 85.8% | **82.2%** | **1.07** | **0%** |

| Typical | Quality | Success | Objects | Multi | None |
|---|---|---|---|---|---|
| Dex-Net | 0.821 | 84.5% | 32.0% | 3.18 | **0%** |
| Zeng et al. | 0.571 | 51.5% | 6.5% | 12.72 | 30% |
| — finetuned | 0.547 | 70.5% | 6.0% | 14.80 | 20% |
| Ours | **0.909** | **90.1%** | **84.5%** | **1.01** | **0%** |

| Complex | Quality | Success | Objects | Multi | None |
|---|---|---|---|---|---|
| Dex-Net | 0.584 | 36.5% | 15.8% | 4.55 | 10% |
| Zeng et al. | 0.299 | 23.5% | 3.9% | 14.17 | 50% |
| — finetuned | 0.601 | **81.5%** | 10.4% | 11.54 | 10% |
| Ours | **0.790** | 56.4% | **28.2%** | **1.36** | **0%** |

geometry also improves the overall quality of single-suction grasps. This might be because the larger gripper geometry forces resulting grasp poses to be more consistently located on highly graspable areas and close towards object centers.

Finally, we would like to emphasize that we do not claim to generally outperform any of the other methods. We verified that our method works sufficiently well on the intended use case and object portfolio compared to related work. In general, applying the proposed gripper footprint optimization can also be done on grasp quality maps predicted by other methods. However, we observed that it works most robustly for the procedure presented in this paper. Still, especially Dex-Net shows a remarkable performance in the presented evaluation, considering that we were able to directly apply pretrained weights with minimal finetuning or preprocessing.

### B. Real-World Multi-Suction Grasp Experiments

We finally demonstrate our approach by detecting poses for multi-suction grasps in the target application of industrial bin picking. The experiment is performed on an industry-grade robotic bin picking cell as shown in Fig. 1. The cell is equipped with two *Zivid Two* overhead RGB-D cameras with a resolution of $1944 \times 1200$ pixels, located around $1.2\,\mathrm{m}$ above the bins. It includes a *Kuka KR10 R1420* robot with a custom-made linear gripper with multiple suction cups and activation schemes, the same as shown in Fig. 5 (left).

The cell is operated by an industrial software stack based on the *Nexeed Automation* framework with grasp poses provided by a *PyTorch*-based implementation of our method that runs on a dedicated IPC of the cell for perception. The IPC runs Ubuntu 20.04 and is equipped with an *Intel Xeon W-1290T* CPU and a *Nvidia Quadro RTX 4000* GPU.

A short video of our experiment is provided online[1]. Various object types similar to those from the grasp quality evaluation are provided in six different bins by the conveyor

[1]See experiment video: https://youtu.be/UZikmSjQy3M

Fig. 9. Visualization of the grasps and corresponding footprints predicted by our approach for an exemplary scene. The executed grasp with selected three-cup footprint is colored in black in the visualization.

belt system. We set an arbitrary sequence of picking orders in which the system composes deliveries from the available objects, simulating typical operations in a logistics center.

In the experiment, detecting multi-suction grasps requires on average $628\,\mathrm{ms}$, of which grasp quality inference takes $53\,\mathrm{ms}$ and gripper footprint optimization takes $372\,\mathrm{ms}$ for the configured four footprints, a maximum rotation of $180°$ (due to symmetry) and a rotation resolution of $5°$. As in the comparison, each detection results in a list of up to 20 grasps from which subsequently, the path planner can select one to execute and considers the respective gripper activation.

Fig. 9 shows the grasp predictions and the performed grasp for one of the scenes from the video. The selected grasp (black) places the footprint centrally on the object and fits three suction cups to increase the robustness. The performed grasp matches the predicted footprint within the error margins of the overall system.

For the complete experiment, the system executed 38 grasps with three failures, i.e., a success rate of $92\%$. In case of a failed grasp attempt, the system automatically executes the next feasible grasp pose. Note that there are no pixel-wise ground-truth labels available in such real-world runs, thus we cannot determine all metrics as in Tab. I for the evaluation dataset. Still, this qualitative experiment verifies practical applicability of the method in an industrial setting, and metrics such as the success rate match the expectations from our previous evaluation.

Overall, it can be seen from the video that choosing multi-suction grasps for larger objects indeed leads to an improved robustness for the grasps. Still, we also observe that the overall system often selects grasps with a single suction cup. This can be attributed to the fact that the path planner is allowed to freely rotate poses for the single-suction grasps due to being rotation-symmetric, which significantly increases the likelihood to find a feasible trajectory.

For smaller objects, the identified graspable areas are sometimes too small for fitting a footprint, a failure case that can be observed for the narrow cylindrical objects where once in the video, only a single grasp pose is detected but deemed infeasible by the path planner. Finally, one concern was that the simple projection of a 2D footprint onto the scene surface might create projection issues on strongly tilted or bent surfaces, but we did not observe practical issues resulting from it.

## VII. Conclusions

In this paper, we proposed a method for detecting and optimizing multi-suction grasp poses for bin picking tasks based on a model-free, gripper-agnostic prediction of pixel-wise graspability values. In addition, we presented an automated procedure for labeling of images for supervised training of the grasp quality network, allowing for a trade-off between annotation quality and labeling effort. For optimizing the selection of an activation pattern and the orientation of a multi-suction grasp, we described a procedure based on a convolution of grasp quality and surface normals with gripper-specific footprints.

In our evaluation and real-world experiments, we observed that the approach reliably predicts poses for one or more suction cups in a realistic setting, leading to feasible and robust grasps performed by the system. To address a broader portfolio of objects and surface properties, future work can include multi-channel grasp quality predictions to denote different surface requirements for gripper selection. On the system side, future work may allow for a closer integration of grasp optimization and motion planning.

## References

[1] M. Breyer, J. J. Chung, L. Ott, R. Siegwart, and J. I. Nieto. Volumetric grasping network: Real-time 6 DOF grasp detection in clutter. In *CoRL*, volume 155, pages 1602–1611. PMLR, 2020.

[2] H. Cao, H.-S. Fang, W. Liu, and C. Lu. Suctionnet-1billion: A large-scale benchmark for suction grasping. *IEEE RA-L*, 6 (4):8718–8725, 2021.

[3] N. Correll, K. E. Bekris, D. Berenson, O. Brock, A. Causo, K. Hauser, K. Okada, A. Rodriguez, J. M. Romano, and P. R. Wurman. Analysis and observations from the first amazon picking challenge. *IEEE T-ASE*, 15(1):172–188, 2016.

[4] H.-S. Fang, C. Wang, M. Gou, and C. Lu. Graspnet-1billion: A large-scale benchmark for general object grasping. In *IEEE CVPR*, pages 11444–11453, 2020.

[5] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *IEEE CVPR*, pages 770–778, 2016.

[6] C. Hernandez, M. Bharatheesha, W. Ko, H. Gaiser, J. Tan, K. van Deurzen, M. de Vries, B. van Mil, J. van Egmond, R. Burger, et al. Team delft's robot winner of the amazon picking challenge 2016. In *Robot World Cup*, pages 613–624. Springer, 2016.

[7] K.-Y. Jeng, Y.-C. Liu, Z. Y. Liu, J.-W. Wang, Y.-L. Chang, H.-T. Su, and W. H. Hsu. GDN: A coarse-to-fine (c2f) representation for end-to-end 6-dof grasp detection. *CoRL*, pages 220–231, 2021.

[8] P. Jiang, J. Oaki, Y. Ishihara, J. Ooga, H. Han, A. Sugahara, S. Tokura, H. Eto, K. Komoda, and A. Ogawa. Learning suction graspability considering grasp quality and robot reachability for bin-picking. *Frontiers in Neurorobotics*, 16, 2022.

[9] Z. Jiang, Y. Zhu, M. Svetlik, K. Fang, and Y. Zhu. Synergies between affordance and geometry: 6-dof grasp detection via implicit representations. In *RSS XVII, Virtual*, 2021.

[10] N. Khargonkar, N. Song, Z. Xu, B. Prabhakaran, and Y. Xiang. Neuralgrasps: Learning implicit representations for grasps of multiple robotic hands. *arXiv preprint arXiv:2207.02959*, 2022.

[11] K. Kleeberger, R. Bormann, W. Kraus, and M. F. Huber. A survey on learning-based robotic grasping. *Current Robotics Reports*, 1(4):239–249, 2020.

[12] Y. Li, L. Schomaker, and S. H. Kasaei. Learning to grasp 3d objects using deep residual u-nets. In *RO-MAN*, pages 781–787. IEEE, 2020.

[13] M. Maggi, G. Mantriota, and G. Reina. Introducing polypus: A novel adaptive vacuum gripper. *Mechanism and Machine Theory*, 167:104483, 2022.

[14] J. Mahler, M. Matl, X. Liu, A. Li, D. Gealy, and K. Goldberg. Dex-net 3.0: Computing robust robot suction grasp targets in point clouds using a new analytic model and deep learning. In *ICRA*, pages 5620–5627. IEEE, 2018.

[15] Jeffrey Mahler, Jacky Liang, Sherdil Niyaz, Michael Laskey, Richard Doan, Xinyu Liu, Juan Aparicio Ojea, and Ken Goldberg. Dex-net 2.0: Deep learning to plan robust grasps with synthetic point clouds and analytic grasp metrics. 2017.

[16] G. Mantriota. Optimal grasp of vacuum grippers with multiple suction cups. *Mechanism and machine theory*, 42 (1):18–33, 2007.

[17] D. Morrison, J. Leitner, and P. Corke. Closing the loop for robotic grasping: A real-time, generative grasp synthesis approach. In *RSS XIV, Pittsburgh, USA*, 2018.

[18] R. Newbury, M. Gu, L. Chumbley, A. Mousavian, C. Eppner, J. Leitner, J. Bohg, A. Morales, T. Asfour, D. Kragic, et al. Deep learning approaches to grasp synthesis: A review. *arXiv preprint arXiv:2207.02556*, 2022.

[19] P. Ni, W. Zhang, X. Zhu, and Q. Cao. Pointnet++ grasping: learning an end-to-end spatial grasp generation algorithm from sparse point clouds. In *ICRA*, pages 3619–3625. IEEE, 2020.

[20] Y. Qin, R. Chen, H. Zhu, M. Song, J. Xu, and H. Su. S4g: Amodal single-view single-shot se (3) grasp detection in cluttered scenes. In *CoRL*, pages 53–65. PMLR, 2020.

[21] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, pages 234–241. Springer, 2015.

[22] V. Satish, J. Mahler, and K. Goldberg. On-policy dataset synthesis for learning robot grasping policies using fully convolutional deep networks. *RA-L*, 4(2):1357–1364, 2019.

[23] L. Shao, F. Ferreira, M. Jorda, V. Nambiar, J. Luo, E. Solowjow, J. A. Ojea, O. Khatib, and J. Bohg. Unigrasp: Learning a unified model to grasp with multifingered robotic hands. *IEEE RA-L*, 5(2):2286–2293, 2020.

[24] Q. Shao, J. Hu, W. Wang, Y. Fang, W. Liu, J. Qi, and J. Ma. Suction grasp region prediction using self-supervised learning for object picking in dense clutter. In *ICMSR*, pages 7–12. IEEE, 2019.

[25] Z. Xu, B. Qi, S. Agrawal, and S. Song. Adagrasp: Learning an adaptive gripper-aware grasping policy. In *ICRA*, pages 4620–4626. IEEE, 2021.

[26] D. Yang, T. Tosun, B. Eisner, V. Isler, and D. Lee. Robotic grasping through combined image-based grasp proposal and 3d reconstruction. In *ICRA*, pages 6350–6356. IEEE, 2021.

[27] A. Zeng, S. Song, K.-T. Yu, E. Donlon, F. R. Hogan, M. Bauza, D. Ma, O. Taylor, M. Liu, E. Romo, et al. Robotic pick-and-place of novel objects in clutter with multi-affordance grasping and cross-domain image matching. In *ICRA*, pages 3750–3757. IEEE, 2018.

[28] A. Zeng, S. Song, K.-T. Yu, E. Donlon, F. R. Hogan, M. Bauza, D. Ma, O. Taylor, M. Liu, E. Romo, et al. Robotic pick-and-place of novel objects in clutter with multi-affordance grasping and cross-domain image matching. *IJRR*, 41(7):690–705, 2022.

[29] B. Zhao, H. Zhang, X. Lan, H. Wang, Z. Tian, and N. Zheng. REGNet: Region-based grasp network for end-to-end grasp detection in point clouds. In *ICRA*, pages 13474–13480. IEEE, 2021.