

# Gradient based Grasp Pose Optimization on a NeRF that Approximates Grasp Success

Gergely Soti<sup>1,2</sup>, Björn Hein<sup>1,2</sup>, and Christian Wurr<sup>1</sup>

<sup>1</sup> Hochschule Karlsruhe – University of Applied Sciences, 76133 Karlsruhe, Germany

<sup>2</sup> Karlsruhe Institute of Technology, 76131 Karlsruhe, Germany  
 gergely.soti@h-ka.de\*

**Abstract.** Current robotic grasping methods often rely on estimating the pose of the target object, explicitly predicting grasp poses, or implicitly estimating grasp success probabilities. In this work, we propose a novel approach that directly maps gripper poses to their corresponding grasp success values, without considering objectness. Specifically, we leverage a Neural Radiance Field (NeRF) architecture to learn a scene representation and use it to train a grasp success estimator that maps each pose in the robot’s task space to a grasp success value. We employ this learned estimator to tune its inputs, i.e., grasp poses, by gradient-based optimization to obtain successful grasp poses. Contrary to other NeRF-based methods which enhance existing grasp pose estimation approaches by relying on NeRF’s rendering capabilities or directly estimate grasp poses in a discretized space using NeRF’s scene representation capabilities, our approach uniquely sidesteps both the need for rendering and the limitation of discretization. We demonstrate the effectiveness of our approach on four simulated 3DoF (Degree of Freedom) robotic grasping tasks and show that it can generalize to novel objects. Our best model achieves an average translation error of 3mm from valid grasp poses. This work opens the door for future research to apply our approach to higher DoF grasps and real-world scenarios.

**Keywords:** robotic grasping, neural scene representation, transfer learning

## 1 Introduction

Research in robotic grasping has explored various approaches such as analytic and data-driven, model-based and model-free, supervised, self supervised and reinforcement learning methods. These methods can be based on different types of sensor data, such as RGB or depth images, and can be designed for different types of grippers [10].

---

\* This research is being conducted as part of the KI5GRob project funded by the German Federal Ministry of Education and Research (BMBF) under project number 13FH579KX9.

Most of these methods are based on object pose estimation, directly estimate a grasp pose or implicitly map grasp poses to their probability of success. However, if we observe ourselves while grasping an object, we might notice, that we intuitively adjust our hand position to increase the chances of a successful grasp and to achieve a good grasp position ultimately. This suggests that the process of grasping can be modeled as an optimization problem that optimizes the pose of our hands to maximize the probability of a successful grasp.

In this work, we introduce a novel approach to robotic grasping. Leveraging VisionNeRF [12], a learned neural network model capable of capturing a 3D scene representation, we create a model that estimates the success of a grasp given a candidate pose. Unlike other – including NeRF-based – grasping methods which directly estimate grasp poses, our approach stands out by formulating grasp pose estimation as a continuous optimization problem. The goal is to maximize the likelihood of successful grasping through gradient-based optimization. We show the efficiency of our proposed approach on four simulated 3DoF robotic grasping tasks. We summarize our contributions as follows:

- We propose a method to explicitly map grasp candidates to their corresponding grasp success value.
- We show the efficacy of applying transfer learning to a trained VisionNeRF to obtain this explicit mapping.
- We propose a novel approach to find valid grasp poses by applying gradient based optimization on the learned grasp success estimator.

## 2 Related Work

### 2.1 Data-driven Robotic Grasping

In recent years, data-driven methods have become the state of the art in the context of robotic object handling. Keypoint detection or dense descriptor-based methods are effective at learning successful grasp poses and can even generalize to object categories, but they often require a large amount of object-specific labeled data to achieve good performance [5, 11, 13, 15, 18]. End-to-end learning models that directly learn to map the robot’s raw sensor input to a desired output offer a promising alternative in unstructured environments [1, 3, 6, 9, 19, 23, 28, 29]. Most of these models directly propose suitable grasp candidates, or estimate the success probability of grasp poses and rank them. These latter models implicitly map grasp poses to success probability, limiting their ability to optimize grasp poses to iterative methods [20] that sample, evaluate, and re-sample grasp candidates to find a better solutions. In contrast, our proposed method explicitly maps grasp poses to grasp success using a neural network, making it differentiable and enabling gradient-based optimization to refine the grasp pose.

### 2.2 NeRFs and NeRF-based Robotic Grasping

Recently, differentiable scene representations, such as Neural Radiance Fields [17], have been increasingly used in the field of robotics also for grasping among

other applications. A NeRF maps a 5-degree-of-freedom (5DoF) pose to an RGB color vector and a so-called density. Color and density are then combined along camera rays via volumetric rendering in order to render novel views for scenes. Various extensions of the NeRFs have been developed for different applications, such as NeuS [21] for surface reconstruction or NeRF-W [16] on unconstrained photo collections of famous landmarks. Plenotrees [26] have been proposed for fast rendering with NeRFs. PixelNeRF [27] and VisionNeRF [12] overcome the need for training a NeRF for each scene, by generalizing over multiple scenes given sparse observations.

Inverse Neural Radiance Fields [25] perform camera pose estimation by inverting a trained NeRF. Starting from an initial camera pose estimate, it uses gradient based optimization to minimize the residual between pixels rendered from an already-trained NeRF and pixels in an observed image. To estimate the 6DoF camera pose, iNeRF casts rays from the camera’s perspective and samples points along them, to finally apply volumetric rendering to get the pixel values and thus the residual. This requires querying the NeRF with different 5DoF poses multiple times. In our method, we use a similar approach, but since we are only interested in estimating 3DoF poses (5DoF with a fixed direction), we can simply use the NeRF’s output at 5DoF poses as an objective.

There are several successful methods that utilize variants NeRFs for robotic grasping. Dex-NeRF [7] uses a NeRF-based model to render high-quality depth images of a scene, which are then fed to DexNet [14] to compute robust grasp poses. Evo-NeRF [8], is similar method, but instead of focusing on improving the depth rendering, the grasp planner network is trained to perform well on the NeRF-rendered depth maps and utilizes a different NeRF implementation to significantly improve training times. GraspNeRF [2] utilizes a multiview NeRF-based approach to estimate a truncated-signed-distance-field in voxels to predict successful grasps. An other approach [24] uses a NeRF to learn dense object descriptors from visual observations, which are then used to track keypoints on objects and calculate grasp poses.

### 2.3 Our Contribution

While these methods demonstrate the effectiveness of utilizing NeRFs in robotic grasping, they typically enhance existing grasp pose estimation techniques with NeRFs’ rendering capabilities or directly estimate grasp poses in a discretized space using NeRFs. However, these methods do not fully exploit the potential of NeRFs for continuous optimization of the grasp pose.

Our method uniquely employs NeRFs to explicitly represent the mapping of grasp poses to grasp success probability. This approach enables a gradient-based optimization method to find optimal grasp poses, providing more fine-grained control over the optimization process. Furthermore, our explicit mapping of grasp poses to grasp success offers a natural representation of the problem, where the gradient directly depicts rigid transformations leading to more successful poses. We believe our method addresses the gaps in the current state of the art and introduces a fresh perspective to the field of robotic grasping.

### 3 Grasp Success Approximation and Optimization

Given an RGB observation of a tabletop scene, the goal of the proposed method is to detect 5-DoF grasps (e.g. with a suction cup) consisting of a position and a direction vector. We assume the camera intrinsics and extrinsics are known for the image. We formulate the 5-DoF grasp detection as an optimization problem, that maximizes grasp success probability over gripper poses. We approximate the function that maps 5-DoF grasp poses  $g = (x, d) \in \mathbf{G}$ , with  $x$  position and  $d$  direction, to their probability of success by the neural network  $\Theta$ . Since neural networks are differentiable, we can solve the problem

$$\max_{g \in \mathbf{G}} \Theta(g, o) \quad (1)$$

by gradient based optimization methods, where  $o = (c, K, RT)$  is an observation containing a camera image with known intrinsics and extrinsics.

In this section we first describe the architecture of  $\Theta$  and how we train it, then we describe the gradient based optimization. Note, that although this formulation is valid for 5DoF grasps, we constrain ourselves to 3DoF grasps (position only with fixed direction) in the evaluation.

#### 3.1 Grasp Success Approximation

NeRFs excel in novel view synthesis and are increasingly being applied in various other tasks that require some sort of scene representation. By using volumetric rendering to compute the loss during training, NeRFs are forced to learn how to consistently represent 3D scenes. In this paper, we demonstrate the potential of this representation for grasp success estimation.

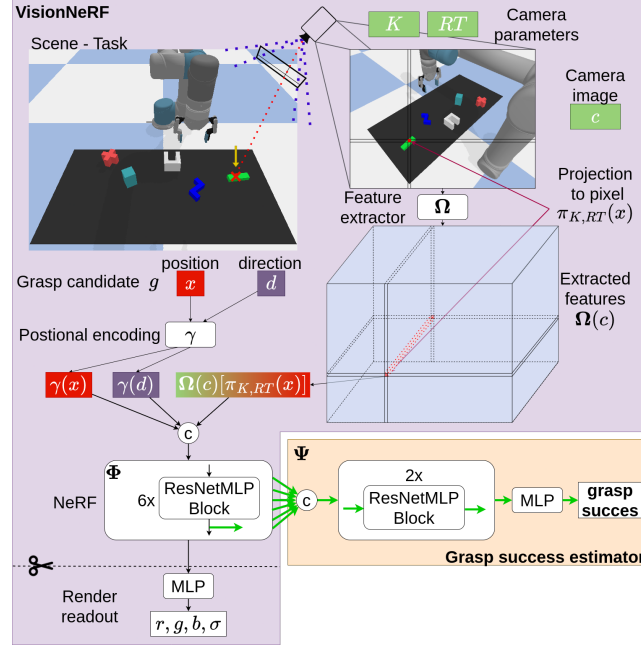
**Architecture** In our approach we use a VisionNeRF [12], a generalized implementation of NeRFs capable of representing multiple scenes by conditioning on observed inputs. To achieve this, a Vision Transformer (ViT) [4] and a Convolutional Neural Network (CNN) are combined to extract global and local features from the input observation, the source image, which are then used to inform the color and density estimator. We denote this combination as  $\Omega$ . While standard NeRFs map a 5DoF pose  $(x, d)$  (corresponding to a 3D point in the scene and the camera’s perspective) to an RGB color vector and density, VisionNeRFs require an additional input: a single camera image  $c$  of the scene with known intrinsics  $K$  and extrinsics  $RT$ .

The NeRF architecture consists of a sequence of ResNetMLP blocks (see [12], for more details) denoted as  $\Phi$ , and a final fully connected layer that outputs the color and density of the 3D point  $x$  given an observation direction  $d$ . To inform this color and density estimator, the feature vector from  $\Omega(c)$  at the projected position of the 3D point onto the camera image  $\pi_{K,RT}(x)$  is concatenated with encoded  $x$  and  $d$  vectors. We use a positional encoding typical for NeRFs:

$$\gamma(p) = (\sin(2^0 \pi p), \cos(2^0 \pi p), \dots, \sin(2^{M-1} \pi p), \cos(2^{M-1} \pi p)) \quad (2)$$

with  $M$  the number of frequency phases. Note, VisionNeRF only applies position encoding to the position vector and not the direction vector, but we also use it on the direction vector, just like the original NeRF implementation. This concatenated vector of  $\gamma(x)$ ,  $\gamma(d)$  and  $\Omega(c)[\pi_{K,RT}(x)]$  is then fed into  $\Phi$  and passed to the final fully connected layer. The output colors and densities of multiple points along a camera ray are then integrated using volumetric rendering to obtain pixel values, thus rendering the target image, as shown in Fig. 1.

To leverage the learned representation, we propose an extension to the VisionNeRF architecture. One potential issue is that the output of  $\Phi$ , which is primarily trained to approximate color and density, can be biased towards these features. To address this, we introduce skip connections after each ResNetMLP block in  $\Phi$ . By concatenating these skip connections with the output of  $\Phi$ , we create the input for our grasp success estimator module, denoted as  $\Psi$ .  $\Psi$  consists of two ResNetMLP blocks and a final fully connected layer, which outputs a grasp success score. Fig. 1 shows the architecture of our model, but for visualization purposes we only depict the position and the direction of our grasp candidate  $g = (x, d)$  in the image. In reality we propagate four 5DoF poses through the network simultaneously, all along  $d$  and centered around  $x$  with 2.5mm spacing. We sum up the output of the grasp success estimator for these poses to obtain the final grasp success of the  $g$ .



**Fig. 1.** The structure of our proposed architecture: a VisionNeRF that estimates color and density for volumetric rendering and a grasp success estimator. Both process 5DoF poses  $(x, d)$  and are informed by the extracted features from the camera image  $c$  that correspond to  $x$

With these we can define the objective function of the optimization problem:

$$\Theta(g, o) = \Psi(\Phi(\gamma(x), \gamma(d), \Omega(c)[\pi_{K, RT}(x)])) \quad (3)$$

with grasp candidate  $g = (x, d)$  and observation  $o = (c, K, RT)$ .

**Training** To get the model to learn to represent the scene, we initially train  $\Omega$  and  $\Phi$  for novel view synthesis via volumetric rendering. The ViT in  $\Omega$  is initialized with pretrained weights from [22]. For training we use the Adam optimizer with a warmup learning rate schedule. The learning rate is increased from 0 to 1e-4 in 10000 steps for  $\Omega$  and similarly for  $\Phi$ , the learning rate is increased from 0 to 1e-5 in 10000 steps.

After training we apply transfer learning to the VisionNeRF by freezing the weights of  $\Omega$  and  $\Phi$  and training only  $\Psi$  to obtain the complete grasp success estimator  $\Theta$ . Categorical cross-entropy loss is used with one successful grasp pose  $g$  for an observation  $o$  as a positive example labeled as 1 and multiple randomly sampled poses as negative examples from the workspace labeled as 0. To obtain a valid grasp pose, we sample a position  $x$  on the top surface of the (prismatic) object - the optimal site for a suction gripper. We set the direction  $d$  perpendicular to this surface. We use the Adam optimizer with learning rate 1e-4 and sample 2047 negative samples.

As baseline, we use the same architecture, but instead of pretraining  $\Omega$  and  $\Phi$ , we only load the ViT pretrained weights. We then train  $\Omega$ ,  $\Phi$  and  $\Psi$  jointly, with the same configurations as described above.

### 3.2 Gradient based Optimization

To solve the optimization problem, we used a gradient base optimization method. We apply the Adam optimizer with a decaying learning rate starting at 0.05 with decay rate 0.8 after each step to minimize the objective function  $-\Theta(g)$  over  $g \in \mathbf{G}$  grasp poses, thus maximizing the estimated grasp success. We initialize the optimization process with  $2^{13}$  random poses as grasp candidates, and the optimization is allowed to run for a maximum of 16 iterations.

Since we constrain ourselves to 3DoF poses only, we fix the direction  $d$  and only optimize  $x$ . This constraint is also applied while training the grasp optimizer by only sampling negative examples with the same direction.

In the context of grasping, the gradient used for optimization corresponds to rigid transformations of the gripper that lead to more successful grasp poses.

## 4 Experiments

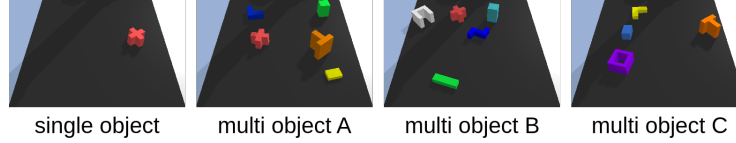
We use a simulated tabletop environment to evaluate the performance of the proposed approach on 3DoF robotic grasping tasks. There are three fixed-pose cameras in the environment that provide camera images as observations. To measure the accuracy of the grasping predictions, we computed the translation

error, which represents the distance between a predicted grasp position and the nearest valid grasp position. Our approach enables the simultaneous optimization of multiple grasp candidates by maximizing their predicted success rate  $\Theta(g, o)$ . We evaluated its performance using two different metrics:

- **best-success**: the translation error of the grasp with the highest predicted success rate
- **lowest-from-5**: the lowest translation error among the five grasp candidates with the highest predicted success rates, which can be roughly understood as if a grasp fails, we can try the next best candidate

For a task, we spawn objects from one of the following sets of objects (see Figure 2):

- **single object**: red cross (0.05)
- **multi object A**: red cross (0.05), green square (0.07), yellow rectangle (0.015), dark blue L-shape (0.03), orange T-shape (0.09)
- **multi object B**: red cross (0.05), turquoise square (0.08), green long rectangle (0.02), white U-shape (0.06), dark blue double-L-shape (0.03)
- **multi object C**: blue rectangle (0.04), yellow L-shape (0.02), orange T-shape (0.07), purple block-ring (0.05)



**Fig. 2.** The different object sets used during training and evaluation.

All objects are prismatic and are characterized by their bases, heights (in meters) and colors. The set multi object A contains objects similar to some of the other multi object sets, while multi object B and multi object C contain mainly different objects.

We define two tasks in which objects need to be grasped:

- **single object grasp**: the object of the single object set is spawned in a random position in the workspace
- **multi object grasp**: five objects are sampled from a multi object set and spawned in random non-overlapping positions in the workspace; the objects need to be removed successively one by one resulting in 5 different scenes for a complete episode

In our experiment we use three different backbones:

- **no-NeRF**: the baseline without pretraining the VisionNeRF
- **single-NeRF**: a VisionNeRF trained on 100 scenes of the single object grasp task, corresponding to 100 complete episodes of the task

- **multi-NeRF**: a VisionNeRF trained on 500 scenes of the multi object grasp task with objects from the multi object A set, corresponding to 100 complete episodes of the task

Both single-NeRF and multi-NeRF are trained for 8000 epochs with batch-size 1. Source and target camera images are sampled from the three fixed-pose camera observations. In contrast, NeRFs are generally trained using many different views which is not realistic in real world setups.

With all three backbones, we train two grasp success estimator modules:

- **single-grasp**: trained on 100 scenes of the single object grasp task, corresponding to 100 complete demonstrations of the task
- **multi-grasp**: trained on 100 scenes of the multi object grasp task with objects from the multi object B set, corresponding to 20 complete demonstrations of the task,

resulting in six models overall. All grasp estimator modules are trained for 250 epochs. The models are evaluated on four tasks:

- **single-object-task**: 50 scenes of the single object grasp task using the single object object set, corresponding to 50 complete episodes
- **multi-object-A-task**: 50 scenes of the multi object task using objects from the multi object A object set corresponding to 10 complete episodes; note, that these objects were also used for training the multi-NeRF module
- **multi-object-B-task**: 50 scenes of the multi object task using objects from the multi object B object set corresponding to 10 complete episodes; note, that these objects were also used for training the multi-grasp module
- **multi-object-C-task**: 50 scenes of the multi object task using objects from the multi object C object set corresponding to 10 complete episodes

For each task we obtain three observations  $o_1, o_2$  and  $o_3$ , one for each camera. For  $\Theta(g, o)$  however, we only need one observation, thus we define two optimization objectives:

- **1-view**:  $\Theta(g, o_1)$
- **3-views**:  $\sum_{i \in [1, 2, 3]} \Theta(g, o_i)$

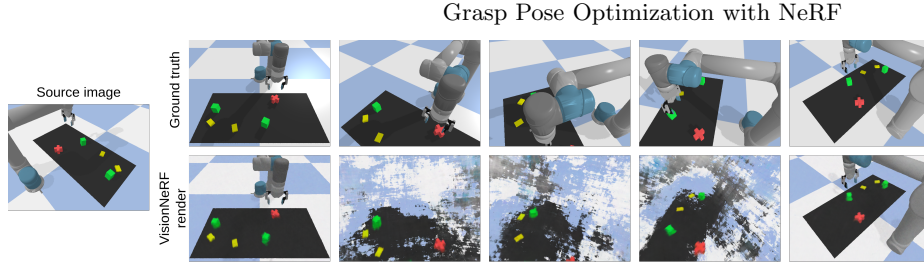
We record the best five grasp candidates with the highest estimated grasp success score after 8, 12 and 16 optimization steps for evaluation.

## 5 Results and Discussion

### 5.1 Qualitative Analysis of Architecture Modules

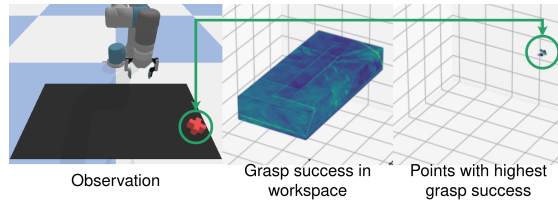
**VisionNeRF** As described above, we trained our VisionNeRF for novel view synthesis only using three perspectives. This leads to a strong bias towards these perspectives during rendering new perspectives given a camera image from a known perspective. When rendering for perspectives that were used during training, the quality of the image is far superior than for other perspectives, however the representation of the objects in the workspace remains mostly consistent, as shown in Fig. 3.





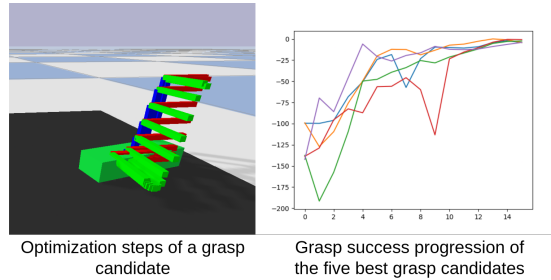
**Fig. 3.** VisionNeRF rendering of new perspectives given a source image with known perspective. The left- and right-most renderings belong to perspectives that were used at training and produce better quality images. For the other perspectives, the static objects in the scene (ground and robot) seems to fall apart, but the objects are rendered consistently even if they were occluded in the ground truth images.

**Grasp success estimator and grasp pose optimization** To ensure that our grasp pose estimation method, which involves gradient-based optimization, is accurate, the learned grasp success estimation function must correctly map 3D (or 5D) space to grasp success. Ideally, the function should assign higher success estimates to points closer to valid grasp positions. Although our learned grasp success estimator has some local maxima that do not correspond to valid grasp poses, the global maxima do, as shown in Fig. 4.



**Fig. 4.** A visualization of the grasp success estimation: the discretized workspace of an instance of the single-grasp-task (left) is mapped to its 3-views optimization objective  $\sum_{i \in [1,2,3]} \Theta(g, o_i)$  (middle). On the right, only the points with the highest success estimation values are shown, also corresponding to the object’s position in the workspace.

Of course, gradient based optimization methods are prone to get stuck in local limits. We overcome this problem by initializing the optimization method with many initial grasp candidates as described in 3.2 and evaluating the grasp candidates that have the highest grasp success estimation at the end of the optimization. Fig. 5 shows the successively improved poses of a grasp candidate during optimization and the estimated grasp success progression of the grasp candidates with the highest estimated grasp success at the end of the optimization. This also suggests, that the gradient does indeed correspond to a rigid transformation that moves the gripper towards better grasp poses.



**Fig. 5.** Grasp pose improvement during optimization (left) and the estimated grasp success progression of the five grasp candidates that have the highest estimated grasp success at the end of the optimization (right).

## 5.2 Robotic Grasping Performance Evaluation

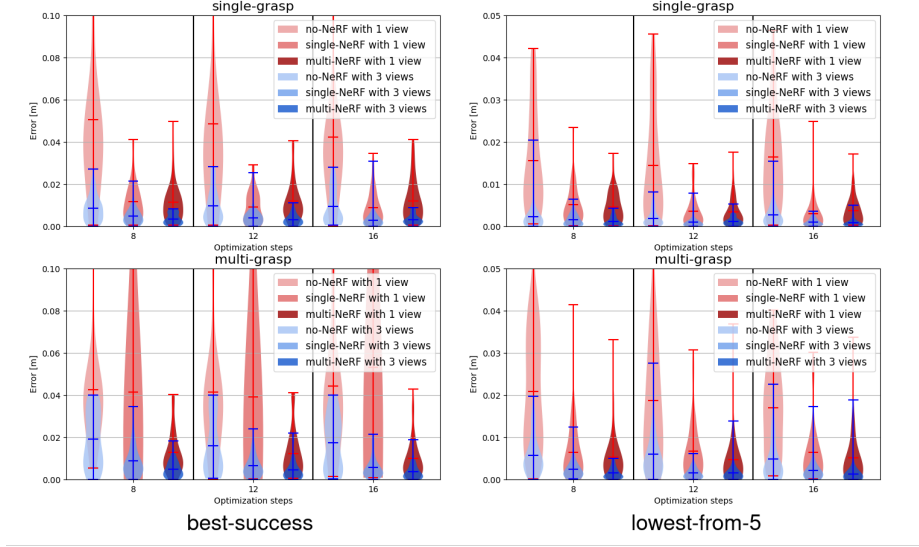
Using 3-views instead of 1-view for the optimization objective reduces the errors of our approach by over an average of 60% for all backbone and grasp success estimator combinations as shown in Fig. 6 using both best-success and lowest-from-5 metrics. Furthermore, for both optimization objectives, the architectures with a pretrained NeRF outperform the models that did not make use of transfer learning, while models with multi-NeRF mostly even outperform models with single-NeRF. A significant exception is observable when models using multi-grasp are evaluated with the best-success metric, where single-NeRF architectures do not outperform models that do not use a pretrained NeRF backbone. This suggests that using a single view in the objective does not depict the reality as reliably as using three views.

In the single-object-task (Fig. 6), architectures that combine a single-NeRF with a single-grasp models can slightly outperform their multi-NeRF counterparts, which is however reasonable, as both single-NeRF and single-grasp models were trained on the same object set, that is used in this task.

We can observe similar behaviour if we observe the results of multi-grasp models with different backbones on the different multi-object-tasks (Fig. 7, right): that models using a NeRF backbone mostly outperform models without a pretrained NeRF backbone and that multi-NeRF outperforms single-NeRF in most cases. Additionally, all models perform best on the multi-object-B-task, which is again reasonable, as the multi-grasp models were trained on the same object set. In case of the multi-object-A-task, the model using multi-NeRF clearly outperforms the other models, which is most likely due to the fact, that the multi-NeRF was also trained on multi object A object set, thus it likely extracts the most descriptive features from scenes with these objects. In the multi-object-C-task the no-NeRF end-to-end model and the single-NeRF model perform similarly and are still outperformed by the multi-NeRF architecture.

When we examine single-grasp models on the same tasks (Fig. 7, left), only multi-object-A-task shows the same pattern regarding backbone configuration. For the multi-object-B task, the architecture without a pretrained NeRF outper-

## single-object-task



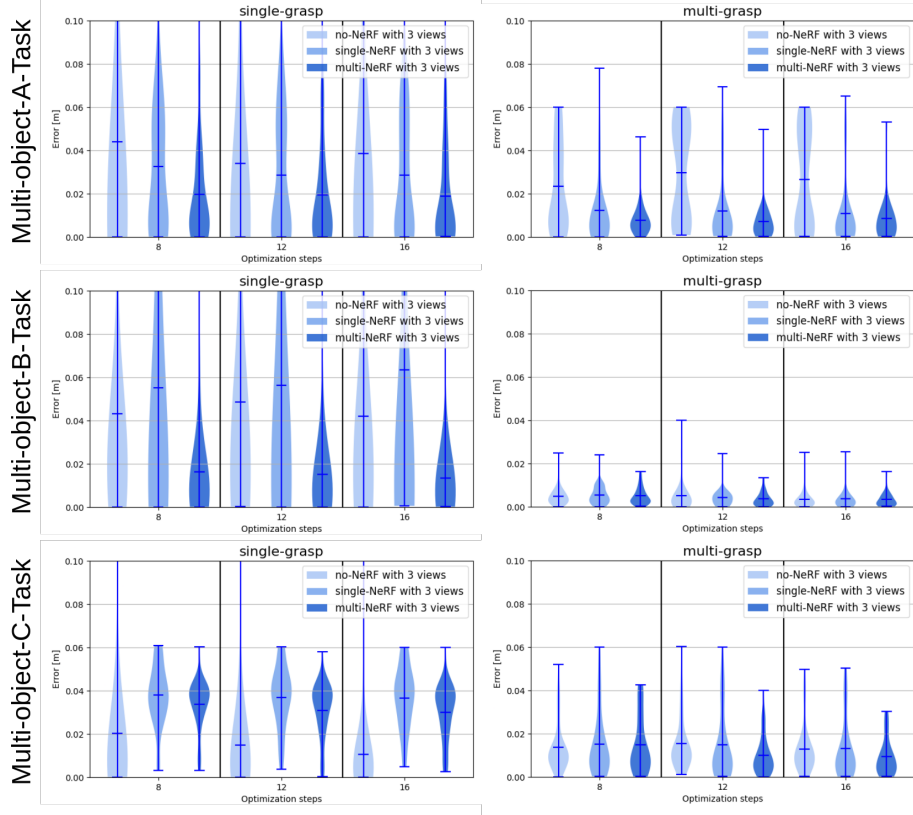
**Fig. 6.** Error distribution of different model configurations on the single-object-task after 8, 12, and 16 optimization steps using 1-view and 3-views as optimization objective for the metrics best-success and lowest-from-5.

forms the single-NeRF architecture, while the model with a multi-NeRF backbone still outperforms both. In case of multi-object-C-task however, the end-to-end architecture delivered the best results.

On average, our models achieved best performance after 16 optimization steps, but only slightly better than after 12 optimization steps. Table 1 summarizes the average errors (in mm) for all models after 16 optimization steps. Considering the best-success metric, for the single-object-task, the single-NeRF with a single-grasp performs best, although it is worth noting, that both multi-NeRF based models have less than 0.7mm larger average error. In case of all multi-object-tasks, the multi-NeRF and the multi-grasp model show better performance, albeit in the multi-object-B-task only slightly better than the other models with a multi-grasp module.

The lowest-from-5 metric models the case when we also consider retrying a failed grasp. As Table 1 demonstrates, the results show similar trends, though not as consistent as for the best-success metric. There is one major outlier: in case of the multi-object-C-task the model combining a single-NeRF with a single-grasp model outperforms all other combinations.

Overall, our results show that it is beneficial to apply transfer learning to a pretrained VisionNeRF model to obtain a model that explicitly maps grasp poses to grasp success. Furthermore, the results suggest that if a VisionNeRF was trained on multiple objects instead of one, thus learning a more descriptive representation of the scene, the obtained grasp success estimator is also better.



**Fig. 7.** Best-success error distribution of different model configurations on all multi-object-tasks after 8, 12, and 16 optimization steps using 3-views as optimization objective.

While single-grasp models are not able to generalize very well to novel objects, the multi-grasp model, which was trained on the multi-object-B object set performed reasonably well on objects from the multi-object-A set, containing partly similar objects, and also on objects from the multi-object-C set containing objects of different shapes and colors. Considering all tasks, the best model is the combination of the multi-NeRF and the multi-grasp model achieving an average error of 3mm.

## 6 Limitations and Future Work

The VisionNeRFs we train only uses three camera perspectives, leading to distorted rendering for other perspectives. This shows, that the learned scene representation is far from perfect. Our method would most likely benefit from a NeRF that is trained with many perspectives. This however, also leads to a major limitation, as such a model would take a huge effort to obtain in a real world

**Table 1.** Average errors of all models using the 3-views optimization objective in mm according to the best-success and lowest-from-5 errors. The single-object-task is denoted as so and the different multi-object-tasks are denoted with mo-X with X referring to the object set they were defined on.

best-success						
	single-grasp			multi-grasp		
	no-NeRF	single-NeRF	multi-NeRF	no-NeRF	single-NeRF	multi-NeRF
so	9.39	<b>2.94</b>	3.17	17.50	5.68	3.61
mo-A	38.46	28.43	18.81	26.47	10.73	<b>8.46</b>
mo-B	41.98	63.30	13.34	3.43	3.59	<b>3.41</b>
mo-C	10.67	36.50	29.98	12.86	13.09	<b>9.33</b>

lowest-from-5						
	single-grasp			multi-grasp		
	no-NeRF	single-NeRF	multi-NeRF	no-NeRF	single-NeRF	multi-NeRF
so	2.70	1.05	<b>0.91</b>	4.87	2.16	1.31
mo-A	22.22	22.87	13.67	13.63	5.75	<b>3.40</b>
mo-B	28.45	28.21	9.38	1.16	<b>1.09</b>	1.29
mo-C	<b>3.53</b>	24.87	25.10	4.66	7.44	5.93

scenario. A possible future work is investigating the possibilities of creating such a model applying sim-2-real transfer learning methods, thus reducing the real world data required.

An other strong limitation is, that our experiments only consider 3DoF grasps and only in simulation. Exploring 5DoF and 6DoF grasps, especially in a real-world experiment, is crucial for a successful model architecture, and is thus also in the scope of possible future works.

For training the grasp estimation models, we used 100 demonstrations from each task. On one hand this does not seem to be that many, if we consider that deep learning architectures usually require an exceptionally large body of data to train on, however for a real world application it would be beneficial if one could reduce the amount of demonstrations to the minimum while retaining the robustness of the method.

## 7 Conclusion

In this work, we propose a unique approach to robotic grasping, employing transfer learning on a trained VisionNeRF to explicitly map grasp poses to their corresponding grasp success. We further applied a gradient-based optimization method on this learned mapping to refine the poses of grasp candidates and thereby attain successful grasp poses. We demonstrated the efficacy of our method on four simulated 3DoF robotic grasping tasks, and showed its ability to generalize to novel objects.

A clear direction for future work is the extension of our method to 5DoF and 6DoF grasps, and to apply it on real world tasks. The methodology we propose here is not limited to robotic grasping, but can be extended to estimate

the success of other types of robotic manipulation or interaction. An intriguing prospect for future development of our work could involve integrating additional criteria, tailored to specific tasks, into the optimization objective. One such criterion could be language conditioning, this could provide a foundation for robots to handle tasks of greater complexity.

## References

1. Berscheid, L., Meißner, P., Kr3ger, T.: Self-supervised learning for precise pick-and-place without object model. *IEEE Robotics and Automation Letters* **5**(3), 4828–4835 (2020)
2. Dai, Q., Zhu, Y., Geng, Y., Ruan, C., Zhang, J., Wang, H.: Graspnerf: Multiview-based 6-dof grasp detection for transparent and specular objects using generalizable nerf. *arXiv preprint arXiv:2210.06575* (2022)
3. Devin, C., Rowghanian, P., Vigorito, C., Richards, W., Rohanimanesh, K.: Self-supervised goal-conditioned pick and place. *arXiv preprint arXiv:2008.11466* (2020)
4. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020)
5. Florence, P.R., Manuelli, L., Tedrake, R.: Dense object nets: Learning dense visual object descriptors by and for robotic manipulation. *arXiv preprint arXiv:1806.08756* (2018)
6. Hundt, A., Killeen, B., Greene, N., Wu, H., Kwon, H., Paxton, C., Hager, G.D.: “good robot!”: Efficient reinforcement learning for multi-step visual tasks with sim to real transfer. *IEEE Robotics and Automation Letters* **5**(4), 6724–6731 (2020)
7. Ichnowski, J., Avigal, Y., Kerr, J., Goldberg, K.: Dex-nerf: Using a neural radiance field to grasp transparent objects. *arXiv preprint arXiv:2110.14217* (2021)
8. Kerr, J., Fu, L., Huang, H., Avigal, Y., Tancik, M., Ichnowski, J., Kanazawa, A., Goldberg, K.: Evo-nerf: Evolving nerf for sequential robot grasping of transparent objects. In: *6th Annual Conference on Robot Learning*
9. Khansari, M., Kappler, D., Luo, J., Bingham, J., Kalakrishnan, M.: Action image representation: Learning scalable deep grasping policies with zero real world data. In: *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 3597–3603. *IEEE* (2020)
10. Kleeberger, K., Bormann, R., Kraus, W., Huber, M.F.: A survey on learning-based robotic grasping. *Current Robotics Reports* **1**, 239–249 (2020)
11. Kulkarni, T.D., Gupta, A., Ionescu, C., Borgeaud, S., Reynolds, M., Zisserman, A., Mnih, V.: Unsupervised learning of object keypoints for perception and control. *Advances in neural information processing systems* **32** (2019)
12. Lin, K.E., Lin, Y.C., Lai, W.S., Lin, T.Y., Shih, Y.C., Ramamoorthi, R.: Vision transformer for nerf-based view synthesis from a single input image. In: *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 806–815 (2023)
13. Liu, X., Jonschkowski, R., Angelova, A., Konolige, K.: Keypose: Multi-view 3d labeling and keypoint estimation for transparent objects. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 11,602–11,610 (2020)

14. Mahler, J., Liang, J., Niyaz, S., Laskey, M., Doan, R., Liu, X., Ojea, J.A., Goldberg, K.: Dex-net 2.0: Deep learning to plan robust grasps with synthetic point clouds and analytic grasp metrics. arXiv preprint arXiv:1703.09312 (2017)
15. Manuelli, L., Gao, W., Florence, P., Tedrake, R.: kpm: Keypoint affordances for category-level robotic manipulation. In: The International Symposium of Robotics Research, pp. 132–157. Springer (2019)
16. Martin-Brualla, R., Radwan, N., Sajjadi, M.S., Barron, J.T., Dosovitskiy, A., Duckworth, D.: Nerf in the wild: Neural radiance fields for unconstrained photo collections. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 7210–7219 (2021)
17. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. In: ECCV (2020)
18. Nagabandi, A., Konolige, K., Levine, S., Kumar, V.: Deep dynamics models for learning dexterous manipulation. In: Conference on Robot Learning, pp. 1101–1112. PMLR (2020)
19. Song, S., Zeng, A., Lee, J., Funkhouser, T.: Grasping in the wild: Learning 6dof closed-loop grasping from low-cost demonstrations. IEEE Robotics and Automation Letters **5**(3), 4978–4985 (2020)
20. Soti, G., Huang, X., Wurrll, C., Hein, B.: Train what you know — precise pick-and-place with transporter networks (2023)
21. Wang, P., Liu, L., Liu, Y., Theobalt, C., Komura, T., Wang, W.: Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. arXiv preprint arXiv:2106.10689 (2021)
22. Wightman, R.: Pytorch image models. <https://github.com/rwightman/pytorch-image-models> (2019). DOI 10.5281/zenodo.4414861
23. Wu, Y., Yan, W., Kurutach, T., Pinto, L., Abbeel, P.: Learning to manipulate deformable objects without demonstrations. arXiv preprint arXiv:1910.13439 (2019)
24. Yen-Chen, L., Florence, P., Barron, J.T., Lin, T.Y., Rodriguez, A., Isola, P.: Nerf-supervision: Learning dense object descriptors from neural radiance fields. In: 2022 International Conference on Robotics and Automation (ICRA), pp. 6496–6503. IEEE (2022)
25. Yen-Chen, L., Florence, P., Barron, J.T., Rodriguez, A., Isola, P., Lin, T.Y.: inerf: Inverting neural radiance fields for pose estimation. In: 2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pp. 1323–1330. IEEE (2021)
26. Yu, A., Li, R., Tancik, M., Li, H., Ng, R., Kanazawa, A.: Plenotrees for real-time rendering of neural radiance fields. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 5752–5761 (2021)
27. Yu, A., Ye, V., Tancik, M., Kanazawa, A.: pixelnerf: Neural radiance fields from one or few images. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 4578–4587 (2021)
28. Zakka, K., Zeng, A., Lee, J., Song, S.: Form2fit: Learning shape priors for generalizable assembly from disassembly. In: 2020 IEEE International Conference on Robotics and Automation (ICRA), pp. 9404–9410. IEEE (2020)
29. Zeng, A., Florence, P., Tompson, J., Welker, S., Chien, J., Attarian, M., Armstrong, T., Krasin, I., Duong, D., Sindhwani, V., et al.: Transporter networks: Rearranging the visual world for robotic manipulation. arXiv preprint arXiv:2010.14406 (2020)