

1 Appendix

2 A Pseudo Code of the Proposed ACE-BNV Policy

3 We summarize the overall NBV planning policy for grasping target object in the following algo-
4 rithm 1.

Algorithm 1 Grasp Affordance Prediction and Next-Best-View Planning

Input: A cluttered scene, a occluded target object given its 3D bounding box \mathbf{T}_{bbox}

Output: A feasible grasp of target object \mathbf{G}

```

for  $t \leq T_{\text{max}}$  do
   $\mathbf{M}_t \leftarrow \mathbf{D}_t$                                  $\triangleright$  Intergrate depth image into TSDF
   $\mathbf{C}_t \leftarrow \text{3D CNN}(\mathbf{M}_t)$                      $\triangleright$  Encode feature
  if  $\mathbf{G}_q^*(\mathbf{T}_{\text{bbox}}, \mathbf{C}_t, \mathbf{O}_{\mathbf{v},t}) \leq q_{\text{max}}$  then
     $\mathbf{O}_{\mathbf{v},t+1} \leftarrow \arg \max_{\mathbf{G}_{\mathbf{v}} \in \mathbf{G}_{\mathbf{v}}} \mathbf{G}_q^*(\mathbf{T}_{\text{bbox}}, \mathbf{G}_{\mathbf{v}})$   $\triangleright$  Evaluate candidate next views
    Move camera to  $\mathbf{O}_{\mathbf{v},t+1}$                          $\triangleright$  Go to the next-best-view
  else
    Execute grasp  $\mathbf{G}_q^*(\mathbf{T}_{\text{bbox}}, \mathbf{C}_t, \mathbf{O}_{\mathbf{v},t})$ 
    Break                                            $\triangleright$  Execute grasp
  end if
end for

```

5 B Network Architecture and Implementation Details

6 We adopt the same encoder with GIGA that takes TSDF $\mathbf{M}_t \in \mathbb{R}^{40 \times 40 \times 40}$ as input and output a fea-
7 ture embedding for each voxel with a 3D CNN layer. Then, the tri-plane feature grids is constructed
8 by projecting each input voxel on a canonical feature plane via orthographic projection. Then,
9 three feature planes are processed with a 2D U-Net that consist of a series of down-sampling and
10 up-sampling 2D convolution layers with skip connections. The output is formulated as the shared
11 tri-plane feature volume $\mathbf{C} \in \mathbb{R}^{3 \times 40 \times 40 \times 32}$, where 32 is the dimension of the feature embedding.

12 Based on the shared tri-plane feature volume, the local feature $\mathbf{C}_{\mathbf{p}}$ of a 3D point $\mathbf{p} = (x, y, z)$
13 is obtained by projecting it to each feature plane and query three features $\mathbf{C}_{\mathbf{p}_x}, \mathbf{C}_{\mathbf{p}_y}, \mathbf{C}_{\mathbf{p}_z}$ at the
14 projected locations using bilinear interpolation, and the local feature $\mathbf{C}_{\mathbf{p}}$ is the concatenated feature
15 of these queried features, i.e., $\mathbf{C}_{\mathbf{p}} = \text{concat}(\mathbf{C}_{\mathbf{p}_x}, \mathbf{C}_{\mathbf{p}_y}, \mathbf{C}_{\mathbf{p}_z})$. We implement our grasp affordance
16 prediction network with a five layer fully-connected network with residual connections. The input
17 dimension of this MLP network is $3 + 96 + 9 \times 96 = 963$ which is composed of view direction
18 unit vector $\mathbf{v} \in \mathbb{R}^3$, the ray feature $\mathbf{C}_{\text{ray}} \in \mathbb{R}^{96}$, and the local geometry feature $\mathbf{C}_{\text{geo}} \in \mathbb{R}^{9 \times 96}$. The
19 output dimension for grasp affordance prediction is $1 + 1 + 1 = 3$ that consist of grasp quality \mathbf{G}_q ,
20 in-plane rotation \mathbf{G}_r , and gripper width \mathbf{G}_w .

21 As for the novel view depth synthesis, we employ a MLP network that takes the 3D point feature
22 $\mathbf{C}_{\mathbf{p}} \in \mathbb{R}^{96}$ as input and output the SDF value of this point, and adpot the same rendering technique
23 with NeuS to synthesize depth images($\eta = 12, \gamma = 5$). We sample 128 rays in depth image in each
24 batch, 64 points on a ray, 4 times importance sampling with 32 points each time. We set the near
25 and far range close to the ground truth depth at the beginning of training, and then gradually relax
26 the range to the maximum range of the implicit feature volume.

27 For experiments in simulation and real word, the size of cubic workspace $L = 30\text{cm}$. The size
28 of the cubic for \mathbf{C}_{vert} is 0.25, which is 7.5cm corresponding to the real world. The points $\mathbf{S} =$
29 $\{\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_n\}$ for \mathbf{C}_{ray} is uniformly sampled with a step of 0.1. The sizes of the three datasets
30 front-observe-front-grasp, front-observe-side-grasp and multi-observe-front-grasp are 1M, 1M and
31 2M grasps, respectively. Each scene contains 240 grasps and $\eta = 12$ ground-truth depth images
32 with the resolution of 480×640 . After data cleaning and balancing, there are about 40% data left.
33 We separate the datasets randomly into 90% training and 10% validation. We train the models with
34 the Adam optimizer and a learning rate of 2×10^{-4} and batch sizes of 128.

35 C Qualitative Results of Real Robot Experiments

36 We present qualitative results in Fig. 1 and 2 and recommend readers watch the supplementary video
 37 for more comprehensive real robot experimental results. Note that our model can select reasonable
 38 next-best-view to observe the occluded target object. We show a representative failure case in Fig. 2,
 39 where small errors in grasp affordance prediction leads to an unsuccessful grasp. This small predic-
 40 tion inaccuracy occurs in most failure experiments. Therefore, in the future, we plan to exploit a
 41 better grasp affordance prediction module to improve the success rate of our method.

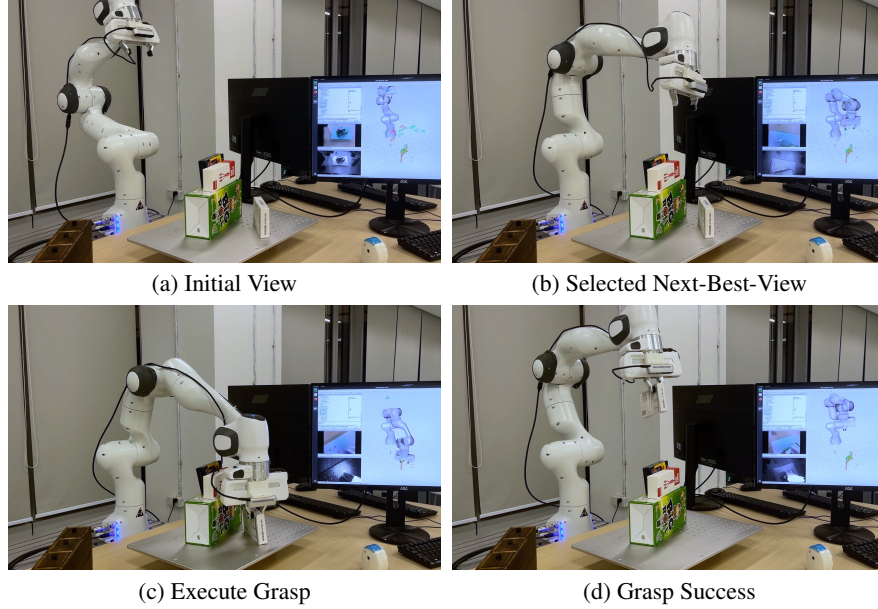


Figure 1: Success Case. The robot planned one new view to observe the target box and successfully grasped it.

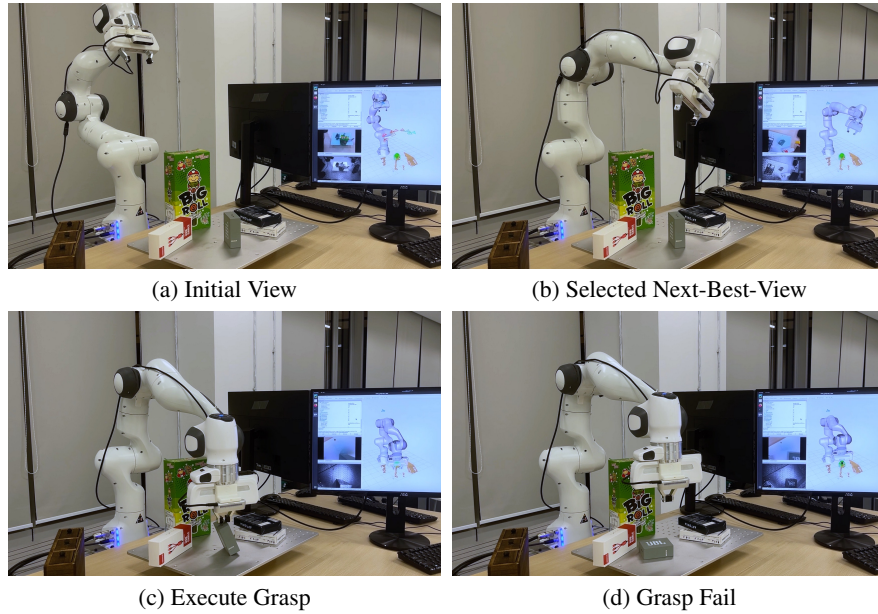


Figure 2: Failure Case. The robot failed to predict accurate grasp affordances of the target object after obtain a new observation. As a result, the grasping is failed.