

---

# Supplementary Material for OBJECT 3DIT: Language-guided 3D-aware Image Editing

---

## 1 Training and Inference Details

We closely follow the training procedure established by [1], with a few modifications. Our approach uses an effective batch size of 1024, which is smaller than the batch size of 1536 used by Zero-1-to-3. This adjustment was necessary because of the additional memory requirements caused by the reintroduction of the CLIP text encoder. This batch size is achieved by using a local batch size of 64 across 40GB NVIDIA RTX A6000 GPUs, along with two gradient accumulation steps. Similar to Zero123, we train on images with a resolution of  $256 \times 256$ , resulting in a latent spatial dimension of  $32 \times 32$ . Following their protocol, we utilize the AdamW optimizer, with a learning rate of  $1e-4$  for all parameters of the model except for those of the concatenation MLP, which uses a learning rate of  $1e-3$ . Our training process runs for a total of 20,000 steps. We then select the best checkpoint based on our metrics computed from an unseen object validation set. As was the case in StableDiffusion, we freeze the CLIP text encoder during training. For inference, we generate images with the DDIM [2] sampler using 200 steps. We do not use classifier-free guidance, i.e. the  $\text{cfg}$  term is set to 1.0.

## 2 Robustness to Severity of Transformation

We analyze the robustness of our method by measuring the performance of the single task rotation model as the complexity of the scene and severity of transformation changes. In Figure 1, we show the average of our Mask PSNR metric as the number of objects in the scene varies from 1 to 4, where a slight drop in performance occurs as the number of objects increases. In Figure 2, we show average Mask PSNR for rotations in a given angle range on a pie chart, where it can be seen that the model does better with smaller angle deviations.

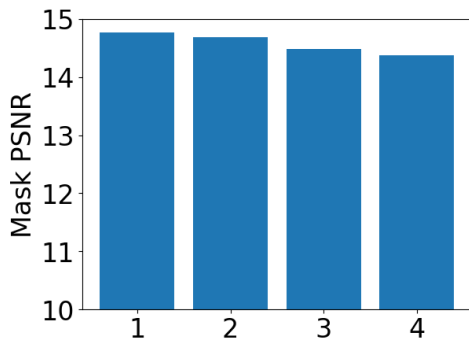


Figure 1: Average Mask PSNR of the single-task rotation model as the number of objects in the scene varies.

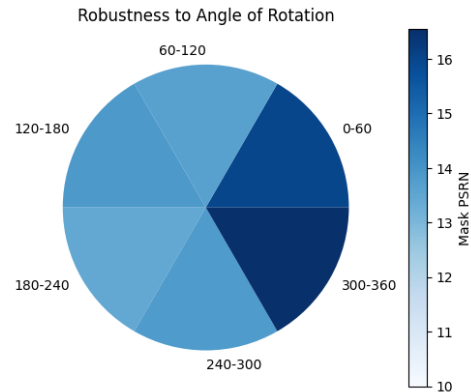


Figure 2: Average Mask PSNR of the single-task rotation model for rotation angles falling within a slice.

Table 1: Summary of key statistics of the OBJECT dataset.

Total objects	62950
Total categories	1613
Object per category median	6.0
Object per category mean	39.03
Object per category std	138.00

### 21 3 Additional Results

22 In Figure 3, we show qualitative results from our multitask model on each of the four editing tasks.  
 23 We also attach three gifs showing our model manipulating a single CLEVR scene. In the rotation  
 24 gif, we prompt the single task model to rotate the "blue box" by increments of 30 degrees. In the  
 25 removal gif, we prompt the model to remove each object with the following descriptions along with a  
 26 visually estimated location: "blue cylinder", "purple sphere", "blue cube", "green cylinder", "brown  
 27 cylinder". For translation, we prompt the model to translate the "blue cube" along the line ranging  
 28 from (0.2, 0.8) to (0.8, 0.2) in increments of (0.05, 0.05).

### 29 4 Dataset Analysis

30 In this section, we provide some details about the composition and statistical makeup of our dataset.  
 31 In Table 1, we show a statistical overview of the dataset, including total number of objects and  
 32 categories, as well as the mean, median, and standard deviation of objects per category. We also  
 33 visualize the distribution of objects across categories, as can be shown in the histogram in Figure 4.  
 34 Finally, we visualize the frequency of category names in the wordcloud in Figure 5.

### 35 References

- 36 [1] R. Liu, R. Wu, B. Van Hoorick, P. Tokmakov, S. Zakharov, and C. Vondrick. Zero-1-to-3:  
 37 Zero-shot one image to 3d object. *arXiv preprint arXiv:2303.11328*, 2023.
- 38 [2] J. Song, C. Meng, and S. Ermon. Denoising diffusion implicit models. *arXiv preprint*  
 39 *arXiv:2010.02502*, 2020.

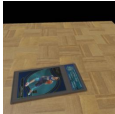
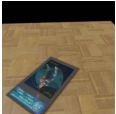
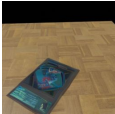





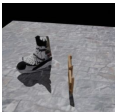
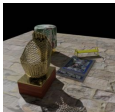
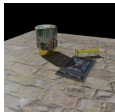
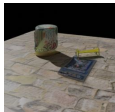



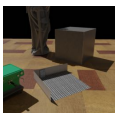




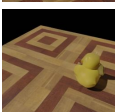
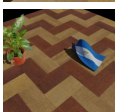
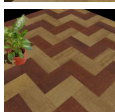
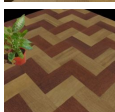









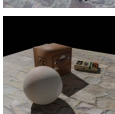
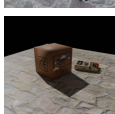
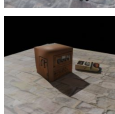
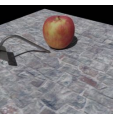


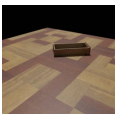
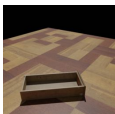
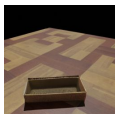
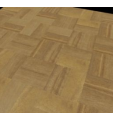





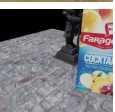
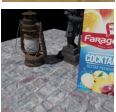
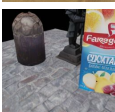



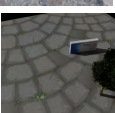
















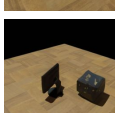
Task: Rotate				Task: Remove			
rotate card by 224.64 degrees				remove robot at (0.52,0.50)			
rotate skate by 219.78 degrees				remove sculpture at (0.31,0.48)			
rotate person by 78.03 degrees				remove box at (0.63,0.66)			
rotate rubber duck by 105.47 degrees				remove flag at (0.73,0.60)			
rotate helmet by 89.57 degrees				remove tire at (0.61,0.81)			
rotate person by 307.27 degrees				remove eyeball at (0.35,0.38)			
Task: Insert				Task: Translate			
insert vase at (0.46,0.32)				move the dollhouse to (0.44,0.30)			
insert cabinet at (0.62,0.45)				move the rock to (0.35,0.33)			
insert a lantern at (0.26,0.73)				move the barrel to (0.58,0.73)			
insert a shield at (0.31,0.63)				move the truck to (0.77,0.36)			
insert a statue at (0.41,0.64)				move the table to (0.46,0.45)			
insert a chair at (0.54,0.80)				move the die to (0.78,0.43)			
Prompt	Input Image	Ground Truth	Ours	Prompt	Input Image	Ground Truth	Ours

Figure 3: Generated examples by the 3DIT multitask model.

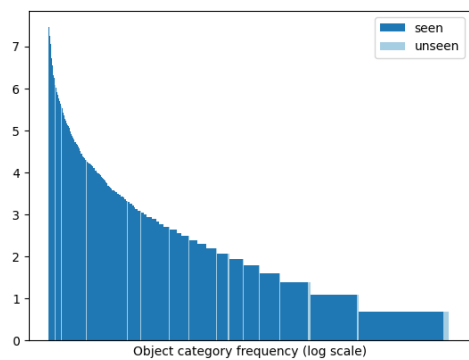


Figure 4: Object categories from seen and unseen splits sorted by frequency, in log scale.



Figure 5: A wordcloud visualizing the frequency of various object category names.