1 We thank the reviewers for their time and insightful comments. We were able to address most of them, which helped
2 improve the paper. We focus on the major comments below, and minor ones will be addressed directly in the paper.

3 **Direct supervision between the two source images [R1, 2, 4].** This is a viable approach; however, it does not
4 generalize as well as our network, as our method reconstructs the source images from unseen target viewpoints, ensuring
5 visual coherence for these views. To support this claim, we designed the proposed approach and retrained the network.
6 Our approach outperforms direct supervision by over 20% (Table/Fig. 1), confirming our approach's superiority.

| Methods | Car | | | Chair | | |
|---|---|---|---|---|---|---|
| | $L_1(\downarrow)$ | SSIM ($\uparrow$) | LPIPS ($\downarrow$) | $L_1(\downarrow)$ | SSIM ($\uparrow$) | LPIPS ($\downarrow$) |
| Supervised | 0.0368 | 0.8481 | 0.1283 | 0.0786 | 0.72 | 0.2501 |
| Ours | **0.0338** | **0.8632** | **0.1029** | **0.0663** | **0.7568** | **0.1766** |

Table 1: Direct supervision vs. Ours



Figure 1: Supervised (left), ours (middle), ground truth (right)

7 **Claim of using no 2D supervision is misleading [R1, 2].** Unsupervised in the paper refers to no target view supervision.
8 Our method does not use any supervision for target viewpoints, only self-supervision on the two source views. We have
9 clarified this point in the paper.

10 **Perceptual evaluation metric [R2].** This is a great point, and we have added the suggested metric to our comparison
11 (Table 1). Full results will be added to the paper.

12 **Model has access to 108 training images during training [R1].** We select two images per object *before* training and
13 only use these for training (LN 187-188). We rewrote the section to avoid confusion.

14 **Evaluation on other datasets than ShapeNet [R1, 3].** This is a great suggestion and we show prelimi-
15 nary results on the real capture dataset (Fig. 2). Additional details and results will be added to the paper.
16

17 **Comparison to SynSin [R4].** SynSin allows rotations of limited range ($\pm 20°$) since
18 it uses a depth point cloud for the 3D features. In our method we generate novel views
19 on the entire viewing hemisphere and therefore did not include this comparison.



Figure 2: Source (left), pre-diction (middle), target (right)

20 **Neural renderer description [R3].** We agree with the reviewer and added additional
21 details of the rendering function in the supplementary material.

22 **Choice of neural rendering function [R3].** This is a very active research area, without
23 a clear consensus on which rendering functions are preferable (see the recent summary
24 report "State of the art in neural rendering"). We chose a rasterization based rendering
25 function as it is fast, differentiable, and the splatting of features to the image plane allows for good gradient flow.

26 **Refinement network and adaptive sampling [R3].** The refinement network fills in gaps and improves rendering
27 quality of the rasterization based neural renderer. Adaptive sampling is an interesting suggestion, and our network
28 would likely benefit from such a step. We will keep this as future work.

29 **Target poses $T_G$ specification [R2].** Poses are represented with spherical coordinates represented by azimuth and
30 elevation. The target pose is selected randomly and given to the network as an input (Fig. 2 in the paper).

31 **More realistic scenes [R3, 4].** This is a research direction that we would like to investigate in future work (LN 254).
32 We plan to investigate model generalization to natural images with background and scenes containing multiple objects.

33 **Bilinear sampling vs. Attentive Neural Processes [R1].** Using "Attentive Neural Processes" looks like a viable
34 option, but we have not investigated this as an alternative and do not know if this would improve performance.

35 **Single image cyclic consistency [R4].** This is technically correct. However, the network collapses to a trivial solution
36 (the identity function) when using a single image. We have added this to the paper.

37 **Supervising the occupancy using binary cross-entropy [R4].** For a single instance, the visual hull can be outside of
38 the object's 3D footprint. However, we find that the multi-view consistency enforces the predicted occupancy features
39 to be coherent with the object's actual 3D structure. We also have preliminary results for extracting the 3D model using
40 the occupancy features.

41 **Performance saturation [R4].** We think that using multiple views will improve performance, but we have not
42 investigated this claim nor the issue of performance saturation. Our approach shows that novel view synthesis is feasible,
43 even if only two images per object are available.

44 **Segmentation network [R1].** We do not use any (pre-trained or not) segmentation networks in our method.