

1 We thank the reviewers for the valuable feedback. While we address the major concerns below, we will include and discuss the
2 missing references (R1, R3), improve and clarify the captions of Figure 2 and Table 2 (R1), and move Figure 3 before Figure 2 (R1)
3 in the revised version.

4 **(R1) Q: Requirement of perfect segmentation and occlusion handling. A:** While obtaining accurate 3D annotations remains
5 a challenging task even for humans, accurate 2D segmentation has become more accessible by either state-of-the-art instance
6 segmentation methods or semi-automatic image processing software. Hence, we believe that learning 3D shapes from 2D supervision
7 can be a much more practical solution for scalable 3D reconstruction. We agree that learning shape representations in the presence
8 of occlusions or imperfect masks remains challenging and it would be a very interesting direction. We will add the above discussion
9 to the limitations and the future work in the revision.

10 **(R1) Q: Marginal improvement of geometric regularization (0.502 vs 0.503). A:** As described in the implementation details
11 (line 231-232), we use $p = 0.8$ as our final model. The ablation study demonstrates that using the geometric regularization with
12 $p = 0.8$ significantly improves 3D IOU from **0.503 to 0.554**. We will highlight the ablation choices in Table 2 to clarify the
13 improvement. In addition, we will provide results with more categories and input viewpoints in the next revision to provide additional
14 evaluations.

15 **(R1) Q: Assumption and background of Equation 1. A:** By viewing the occupancy probability as a thickness of a translucent
16 medium [36], we compute the occupancy probability on the image plane by integrating the multiplication of the transmittance (i.e.,
17 visibility) and occupancy probability at any point along the ray as denoted in Equation 1. Also the visibility of a point along the
18 ray is exponentially decayed with the accumulated density. The same formulation has been used in computer vision literature for
19 differentiable approximation of 3D scene [Rhodin et al. 2015]. We are happy to elaborate more on the formulation in the revision.

20 **(R1) Q: Why implicit value in [0, 1] instead of zero-level set? A:** Without 3D ground truth, it is not possible to obtain the signed
21 distance field in 3D. Thus, we use an occupancy field from $[0, 1]$ which is also used in state-of-the-art 3D reconstruction methods,
22 e.g. Occupancy Network [8], where 0.5 is set as decision boundary. In addition, our approach uses 0.5 with sigmoid function, which
23 is equivalent to the common zero-mean prediction with tanh activation.

24 **(R1) Q: Difference with the Differentiable Ray Consistency (DRC) paper. A:** In DRC, the occupancy fields are stored in explicit
25 voxel grids (32^3), where uniform voxel-wise sampling along ray suffices to propagate gradients to all the voxel grids. In contrast, an
26 implicit shape representation encodes geometry with unlimited resolution at the cost of a dense evaluation of the implicit function,
27 which is not tractable with the uniform sampling and aggregation method proposed in DRC. Hence, we propose an efficient sampling
28 approach for differentiable ray marching with implicit surfaces based on importance, which is one of our main contributions.

29 **(R1) Q: Speed improvement of field probing approach. A:** A naive solution of densely evaluating all points in resolution R_v^3 and
30 aggregating them into pixels in resolution R_p^2 leads to a computational cost of $O(R_v^3)$ and $O(R_p^2 R_v^3)$ for network forwarding and
31 intersection computation, respectively. Such cost is typically intractable for network training with commodity GPUs. Our field
32 probing approach largely reduces the computational complexity to $O(M_v M_p)$ where M_v and M_p stand for the number of 3D anchor
33 points and 2D image pixels respectively with importance sampling. Our experiments show that our method is two magnitudes faster
34 than the naive approach, making it feasible to train a shape inference network with substantially higher resolution.

35 **(R2) Q: Training and inference time; evaluation method and resolution. A:** For each category, our model is trained for one day
36 with a single RTX 2080 Ti GPU. It takes 2s to infer each object at a resolution of 256^3 . For evaluation, we reconstruct the geometry
37 by densely evaluating the implicit function at a regular 3D grid with resolution 256^3 . To compute the 3D IoU, we downsample the
38 voxels to 32^3 using max pooling and compare it with the ground truth. We will include these details in the revision.

39 **(R2) Q: How does the support region radius affect the prediction? A:** A smaller radius of support region helps to reconstruct
40 finer details at the cost of sampling denser points and rays. We found that setting the radius τ as the average distance between any
41 two nearest adjacent anchor points ($\tau = 0.03$) strikes the best trade-off between the prediction accuracy and computational cost with
42 3D IOU of 0.554. Using the same training parameters except the radius, we obtain 3D IOU of 0.515 and 0.502 with radius of 0.01
43 and 0.05 respectively on a chair category. We observed that training with a larger radius tends to ignore details and the one with a
44 smaller radius suffers from insufficient ray assignment. We will add an ablation study of the support region radius in the revision.

45 **(R3) Q: Using anisotropic kernels for shape modeling? A:** We found it non-trivial to use anisotropic kernels for implicit shape
46 learning as it introduces additional parameters to either set or learn (i.e., radius in each axis and axis rotation). We are happy to
47 discuss this as part of future work.

48 **(R3) Q: About hierarchical implicit representation. A:** We believe that our method can be extended to learn hierarchical implicit
49 representations using a spatial data structure such as Octrees and will mention this in the future work section.

50 **(R3) Q: Effect of changing input viewpoints A:** We found that our prediction is reasonably consistent when changing view points.
51 Note that our model is trained with multiview inputs and the numerical evaluations use all 24 views, hence they are not biased w.r.t
52 views. We will provide more results using different viewpoints to better evaluate our approach.

53 **(R3) Q: What does Equation 6 approximate in the limit? A:** We believe R3 meant Equation 5. Our intention of Equation 5 is to
54 encourage local smoothness and the effective smoothed region is controlled by the perturbation Δd . Hence, at the limit, the equation
55 would lead to smoothness for an infinitesimally small volume.

56 **(R3) Q: Additional graphics papers to be discussed. A:** We will include and discuss the suggested papers in the revision.