



Figure 1: Ablation studies on regularization terms. We sequentially add regularizations from left to right. We show the corresponding Chamfer Distance,  $\downarrow$  (CD) & AMIPS,  $\downarrow$  on the bottom.



Figure 2: Scene.

1 We would like to thank reviewers for their detailed comments. We address questions below.

2 **[R1,R2,R3,R4] Ablation Studies:** We provide a qualitative ablation study for different terms by optimizing DEFTET  
 3 to ground truth shape in Fig. 1.  $L_{lap}$  helps maintain the regularity.  $L_{del}$  encourages local deformation.  $L_{vol}$  helps avoid  
 4 flips,  $L_{amips}$  helps avoid needle-like predictions, and  $L_{sm}$  helps smooth the prediction. We will add these in a revision.

5 **R1 Impact of Resolution:** We ablated different resolutions of DEFTET compared to a voxel-based representation in  
 6 Fig. 4. We report the performance of DEFTET at 70 resolution in Table 1 compared to SOTA SDF-based approach [40].

7 **R1 More convincing application:** We partly disagree. Several algorithms have been proposed to differentially render  
 8 an implicit function, yet these methods still use the non-differentiable marching cube algorithm to convert SDF into a 3D  
 9 mesh, which impedes inference speed. DEFTET predicts meshed shape without the time-consuming post-processing.

10 **R2 Unfair comparison to Pixel2Mesh [46] and MeshRCNN [14]:**  $L_{vol}$  and  $L_{amips}$  are volume-based and not suitable  
 11 for [46]&[14].  $L_{lap}$  is also used in [46]&[14].  $L_{sm}$  and  $L_{del}$  can also be used in [46]&[14]. We perform an additional  
 12 experiment on point cloud 3D reconstruction, [46]/[14] achieved 66.06/72.92 IOU, even worse than original results.

13 **R2 No comparison with SOTA SDF approaches:** We wish to clarify that we did compare with SOTA SDF-based  
 14 approaches. For point cloud 3D reconstruction (Tbl. 1, Fig. 3&4), we compared with [40] (ECCV2020), the best  
 15 SDF-based approach at the time of our submission. For single image 3D reconstruction with **2D supervision** (Tbl. 2,  
 16 Fig. 5), where **no** 3D ground truth is utilized, we compared with DVR [37], which is the SOTA SDF-based approach on  
 17 this task. In both tasks, we outperformed the baselines. DISN is not applicable in these two tasks as it is specifically  
 18 designed for single image 3D reconstruction using **3D supervision**, i.e., 3D ground truth is required during training.

19 **R2 Worse result on single image 3D reconstruction:** We first clarify that single image 3d reconstruction can be  
 20 done using 3D or 2D ground truth as supervision, and we apply DEFTET in both regimes. We show results of  
 21 using 2D supervision in the main paper, which has significant improvement, and we report **both quantitative and**  
 22 **qualitative** results of using **3D supervision** in the supplement (Tbl. 2, Fig. 3). Due to the deadline rush, the results  
 23 in supplement were reported with a non converged network – we apologize. We report the true performance as  
 24 2.006/19.463 Chamfer/Hausdorff Mean, which outperforms [46]&[14]. We thank the reviewer for pointing out DISN.  
 25 Since we use the network backbone from OccNet, it is not fair to directly compare with DISN which uses a more  
 26 powerful network structure. We will include a fair comparison in the revision.

27 **R2 Concerns on Eq. 9~11:** We tried both and found that considering all faces converges faster than only considering  
 28 the last face. If  $f_1$  is partially occluded by  $f_2$ , the occluded part would have a much lower visibility due to Eq.10  
 29 ( $m_1 < m_2$ ). Note that we use  $m$  instead of  $D$  to accumulate colors.

30 **R3 Experiments on other categories:** We showed cars in the supplement, and we will add more classes in a revision.

31 **R3 #Vertices:** This indeed is a limitation, as is for the voxel-based approaches, yet less so for DEFTET since we can  
 32 deform the grid geometry. In the future, we plan to adaptively subdivide to get details, akin to oct-tree voxel grids.

33 **R3 Scale to large scene:** We partially agree with the reviewer, but we can increase the resolution of DEFTET for the  
 34 complex scene. We show two optimization results in Fig. 2 by reconstructing multiple objects with DEFTET.

35 **R3 Smoothness term:**  $\mathcal{F}$  here denotes the set of all triangular faces whose  $P_s(v)$  equals to one (L185 in the paper),  
 36 instead of all the faces of the tetrahedralization. We want the neighbouring boundary faces to be approximately parallel.  
 37 We apologize for the typo – there should be a negative sign before cosine. We will revise.

38 **R3 Avoid flips:** We thank the reviewer for raising this question. There is a tradeoff between having no flips and  
 39 performance. We can achieve having no flips by increasing the weight of  $L_{vol}$  but sacrifice the performance. Empirically,  
 40 when doing optimization, we achieve 3.783 CD with 0 flip, and 3.330 CD with 1 flip over all 729162 tetrahedrons.

41 **R3 #Vertices:** The #vertices for  $30^3$  is 4198, and  $70^3$  is 47416, the maximum resolution we can support is  $150^3$  with  
 42 batch size 4 on a 32G V100 GPU. **R3 Related literature:** Thank you for pointing this out, we will add it.

43 **R4 Comparison with DMC [26]:** We agree that we share similarities with DMC – we discussed (L90) and compared  
 44 with (Tbl. 1 & Fig. 3) DMC. We differ in the following aspects: **1. Basic element.** DMC deforms vertices in a cube  
 45 and constraints the deformation to be along the edges within [0,1], while we freely deform tetrahedron’s vertices,  
 46 whose typical displacements are  $\sim 1.18x$  of edge length. Therefore, DMC has a fixed set of possible topologies and  
 47 needs to deal with many corner cases (Fig. 2 in [26]), while we have more flexibility in representing geometric details  
 48 (e.g. the curvy chair back in Fig. 3) and no corner cases. **2. Formulation.** DMC focuses on the surface and uses the  
 49 “occupancy” to derive the probability over possible topologies. While our DEFTET explicitly defines the occupancy  
 50 for one tetrahedron and deforms it which makes training easier than in DMC. We compared with DMC in point cloud  
 51 3D reconstruction task and single image 3D reconstruction using 3D supervision, and achieved significantly better  
 52 performance. DMC is not applicable to single image 3D reconstruction using 2D supervision. The reason why DMC  
 53 runs slower is that the official source code of DMC runs an iterative algorithm to get the mesh after obtaining the  
 54 displacement and occupancy. We refer to the source code for details.

55 **R4 Drawback of laplacian smoothing layer:** We do not notice drawbacks. The computation is negligible.

56 **R4 Comparison with Tet Meshing:** We show that our DEFTET reaches comparable performance to traditional  
 57 tetrahedral meshing methods – a sanity check experiment. We are more accurate than QuarTet in Tbl. 3. The robustness  
 58 (even 1% failing rate) are very important for applications in computer graphics [17]. Note that all of traditional methods  
 59 can not be plugged into deep learning, while DEFTET can easily support it. We apologize for the typo, Tbl. 3 is correct.