

1 We thank the reviewers for their insightful comments. In this rebuttal, we respond to remarks from reviews. If accepted,  
2 we will extend the submission with discussions from below.

### 3 **REVIEWER #1**

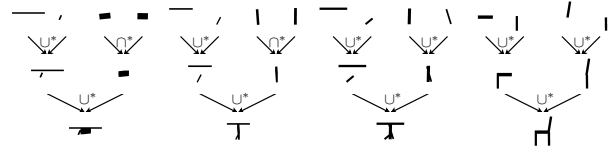
4 **Remark 1** The work lacks discussion about the comparison of interpretability with BSP-Net.

5 **Response** BSP-Net generates multiple convex parts that are difficult to interpret due to an unbounded number of  
6 possible vertices to be used by each convex. These convexes are also problematic to modify from the perspective of a  
7 3D designer. Moreover, their CSG structure is fixed by definition. It uses the intersection of hyperplanes first, and then  
8 the union of convexes.

### 9 **REVIEWER #2**

10 **Remark 1** While the UCSG-NET can predict different CSG trees for different object instances, an empirical evaluation  
11 does not show that.

12 **Response** Our approach is capable of generating diverse  
13 CSG trees for different instances (see Figure on the right).  
14 We found empirically, that the diversity of CSG trees  
15 can be increased with a layer normalization applied for  
16  $\{l \in L | z^{(l)}, h^{(l)}(\hat{V}^{(l)})\}$ . However, this operation slightly



17 degrades quantitative results and we omitted it in our submission. We sacrificed the diversity of CSG trees for accuracy  
18 and reported the qualitative and quantitative results for the best model.

19 **Remark 2** Only a single instance of CSG visualization for each class is shown.

20 **Response** We will provide more various CSG visualizations for each class in the supplementary (including Remark 1).

21 **Remark 3** The proposed approach is inferior to the seemingly more straightforward superquadrics method. The  
22 structure of the ShapeNet dataset may cause it, and further experiments on the ABC dataset would be needed.

23 **Response** Superquadrics [36] method follows the evaluation methodology of Visual Primitives (VP) [37]. These  
24 two approaches use only a union operation. While our approach is inferior to the superquadrics, we focused on  
25 the applicability of our method in 3D design processes where the CSG is commonly used. Therefore, we argue  
26 that referenced approaches are not as versatile as ours. Each work was evaluated on ShapeNet as it is a standard  
27 benchmark. Presenting results on this dataset allows the reader to quickly grasp how the method performs in terms of  
28 the reconstruction. Referring to ABC, we will put effort into analyzing the dataset to use it in our approach.

### 29 **REVIEWER #3**

30 **Remark 1** It is unclear, what low-level representation decisions allowed UCSG-NET, in contrast to CSG-NET, to be  
31 trained in an unsupervised fashion.

32 **Response** The CSG-NET predicts a 3D program. Since the program's instructions are discrete, the model is unable  
33 to be smoothly optimized towards a particular set of instructions. Therefore, the direct training in an unsupervised  
34 manner would be difficult. To solve this limitation, we proposed the following advancements over CSG-NET. Firstly,  
35 we applied smooth boolean operations that are pushed towards discrete forms during the training. Secondly, we used  
36 the SDF representation of shapes. It allowed us to convert shapes into occupancy values and apply differentiable CSG  
37 operations (Fig. 2). These operations are selected dynamically through the attention mechanism (Fig. 3.). Hence,  
38 UCSG-NET does not need any supervision to guide which CSG operations should be selected.

### 39 **REVIEWER #4**

40 **Remark 1** Predicted CSG trees are often redundant and qualitatively dissimilar to human-created ones.

41 **Response** The network predicts such shapes that would not worsen the final reconstruction. Hence, it can predict highly  
42 overlapping boxes. At the same time, we do not force any particular structure of the tree to be learned. Since the task is  
43 ill-posed and there exist infinite CSG trees that reconstruct the same shape, we argue that the space solutions can be  
44 constrained by weak supervision that would enhance the fidelity of predicted CSG trees.

45 **Remark 2** Could the authors explain the motivation of a bottom-up process that groups primitive to form final outputs?

46 **Response** To our knowledge, there are not many machine learning approaches that apply top-down decomposition of  
47 meshes that have proven high reconstruction quality. Moreover, our method was designed in such a way that it can  
48 be used as an extension for other bottom-up approaches in the literature. Therefore, we believe that our work will  
49 encourage future research further to investigate applicability of introduced CSG layers.

50 **Remark 3** Can weak supervision be used in the introduced method?

51 **Response** Applying weak supervision is an interesting future direction. Possibly, it would require introducing a new  
52 term in  $\mathcal{L}_{total}$ . We would also limit the number of possible primitives to be predicted, match primitive predictions with  
53 the ground truth using the Hungarian method, and use cross-entropy loss to optimize predicted masks  $\hat{V}^{(l)}$  directly.

54 **Remark 4** The sentence "contribute to learning disentangled representation of parts" from seems to be unsupported.

55 **Response** To clarify, we found empirically, that our method often separates particular semantic elements of the object  
56 during reconstruction (hence disentanglement) and merges them in the final layer, ex. wings and the hull of an airplane  
57 or legs and the counter of a desk.