# Sketch and Text Guided Diffusion Model for Colored Point Cloud Generation (Supplementary Material)

Here, we provide more details of our proposed method, including color point cloud generation from freehand sketch+text and from only freehand sketch. Section 2 presents typical attention maps for sketches and compares the instance scores before and after attention. In Section 3, we also include a more detailed comparison between DPM [2], PVD [4], point-e [3], and STPD.

We also discuss some of the interesting features of our proposed method at the end.

## 1. Color Point Cloud Generation

Returning to the color point cloud generation ability of our proposed model, in this section, we present additional visualizations for 3D generation of the *chair* and *table* categories from freehand sketch and text in point cloud generation (shown in Figure 1). Specifically, only *chair* and *table* categories have color descriptions to support the sketch-text based coloring. We ensure that the edge detection sketches match hand drawn ones. Our conditional method produces colored shapes that are more faithful to the input sketch-text descriptions. Notice how our method pays more attention to the shape details apparent in the sketch (e.g., see the table legs). Note that for better visualization, we up-sample each point cloud to 4096 points.

## 2. Typical Attention Maps Related to Sketches

We also present additional attention maps generated by our attention-capsule network. In this section, we introduce the Instance Score as the evaluation metric. For the instance score on sketch (ISS), we count the number of useful pixels and express it as a percentage of all sketch pixels. For the instance score on capsule (ISC), we compute the highest attention score of each dimension to all capsules and use their average as the evaluation metric.

All visualized sketch and attention score maps are reported in Table 2 and Figure 3. Note that, across all 24 sketches, the average number of useful pixels representing the instance is only 1.07% of the whole sketch. After our capsule-attention module, the attention weight of the instance improves to 58.68% for better diffusion.

## 3. More Point Cloud Generation Results

We further test our STPD model for generating object shapes and categories from sketches. In this section, we focus only on shape and do not include appearance diffusion. Figure 5 shows our results. For each generation, we input the sketch to generate the 3D shape for all four categories (*airplane*, *chair*, *car*, and *table*) used in our experiments. Since there is no text description associated with 3D shape to support color diffusion, we manually assigned colors for better visualization. The numerous generation results demonstrate that our STPD has a good ability to generalize and generate different 3D shapes.

## 4. Detailed Comparisons

We provide a more detailed visualized comparison between PVD [4], DPM [2], Point-E [3], and our STPD. PVD is trained without color. However, this approach has some drawbacks, such as not being able to capture the full range of colors present in a given image.

DPM has been extended to simultaneously capture color and geometry. However, the additional dimension for color may bias the shape and result in higher variance when simply applying the dimension extension on the sketch-text dataset.

Point-E is demonstrated here from a well-trained large dataset and performs great. However, there are two main problems with Point-E: 1) the shape generation does not strictly follow the conditions (e.g., conditions with black metallic arms but not generating the arms) and 2) the color is not accurately related to the shape parts as conditioned (e.g., showing black metallic feet but displaying gray feet).

In contrast, our method, STPD, shows faithful and fine results from the conditions.

## 5. Additional Discussions

In this section, we will discuss three features of our proposed method: (a) generating novel shapes, (b) the effects of slight changes in the sketch, and (c) conflicting sketch and text.

Our STPD method can learn the variance of parts from the training set. Given novel changes (such as back holes

Figure 1. Visualization of the generated objects from free-hand sketch and text. Our model generates better geometry and has the additional functionality of generating colored point clouds. Being conditional, our method produces colored shapes that are more faithful to the input sketch-text descriptions.

and hollow backrest sketch) that do not exist in the training shapes, the STPD learns these changes from the whole shapes and modifies the generated shapes to make the results faithful to the condition. However, we only trained the diffusion model on some categories of ShapeNet [1], which may limit its ability to deal with totally strange objects. It would be an interesting direction to train a more generalized model across large datasets using our extending methods, like the Stable-diffusion did in the 2D area.

We also tested some slight changes in the sketch to show that the results are faithful to the conditions. However, our vanilla STPD model is trained on a small dataset with sketch-text conditions, and we did not consider conflicts during training. Thus, we tested conflicting input sketches and text (such as a sketch with no armrest and text with an armrest). During the diffusion process, the model some-

times followed the sketch and sometimes the text. We believe that the well-trained BERT mask would randomly mask the conflict information in each modality.

## References

[1] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*, 2015. 2

[2] Shitong Luo and Wei Hu. Diffusion probabilistic models for 3d point cloud generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2837–2845, 2021. 1, 4

[3] Alex Nichol, Heewoo Jun, Prafulla Dhariwal, Pamela Mishkin, and Mark Chen. Point-e: A system for generat-

Table 1. Instance Score Comparison before and after capsule attention module, which are denoted as Instance Score on Sketch (ISS%) and Instance Score on Capsule (ISC%).

| Index | 1,1 | 1,2 | 1,3 | 2,1 | 2,2 | 2,3 | 3,1 | 3,2 | 3,3 | 4,1 | 4,2 | 4,3 | Average |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ISS | 1.38 | 1.01 | 1.78 | 1.45 | 1.93 | 0.97 | 0.53 | 1.34 | 1.29 | 0.81 | 1.01 | 1.11 | |
| ISC | 66.59 | 60.80 | 71.33 | 68.57 | 59.85 | 62.07 | 49.06 | 66.74 | 44.39 | 59.04 | 68.87 | 52.75 | |
| Index | 5,1 | 5,2 | 5,3 | 6,1 | 6,2 | 6,3 | 7,1 | 7,2 | 7,3 | 8,1 | 8,2 | 8,3 | |
| ISS | 0.66 | 1.07 | 0.87 | 1.16 | 1.52 | 0.67 | 0.73 | 0.77 | 0.73 | 0.88 | 1.05 | 0.97 | 1.07 |
| ISC | 60.69 | 57.45 | 65.76 | 62.84 | 55.19 | 65.52 | 40.03 | 57.59 | 44.51 | 57.50 | 52.76 | 58.39 | 58.68 |



Figure 2. Visualization of attention score maps from sketch to capsule. For attention score maps, the darker capsule components pay more attention with our attention-routing module in the final feature extraction. The lighter parts denote less attention. We report the Instance score in Table 2 to show the improvement of the attention.

| Text | Sketch | PVD | DPM-extended | Point-E(text) | STPD |
|------|--------|-----|--------------|---------------|------|



*Coffee table, with straight lines and board in olive green*

*An office chair that has the back open with a row of verticals slots. The two arm rests are just enough for the arms to rest. The seating on the office chair is brown and the base of the chair has 5 arms with each on having wheels for rolling around.*

*rectangular shaped wooden chair with blue and green in color with back rest and arms rest are provided.*

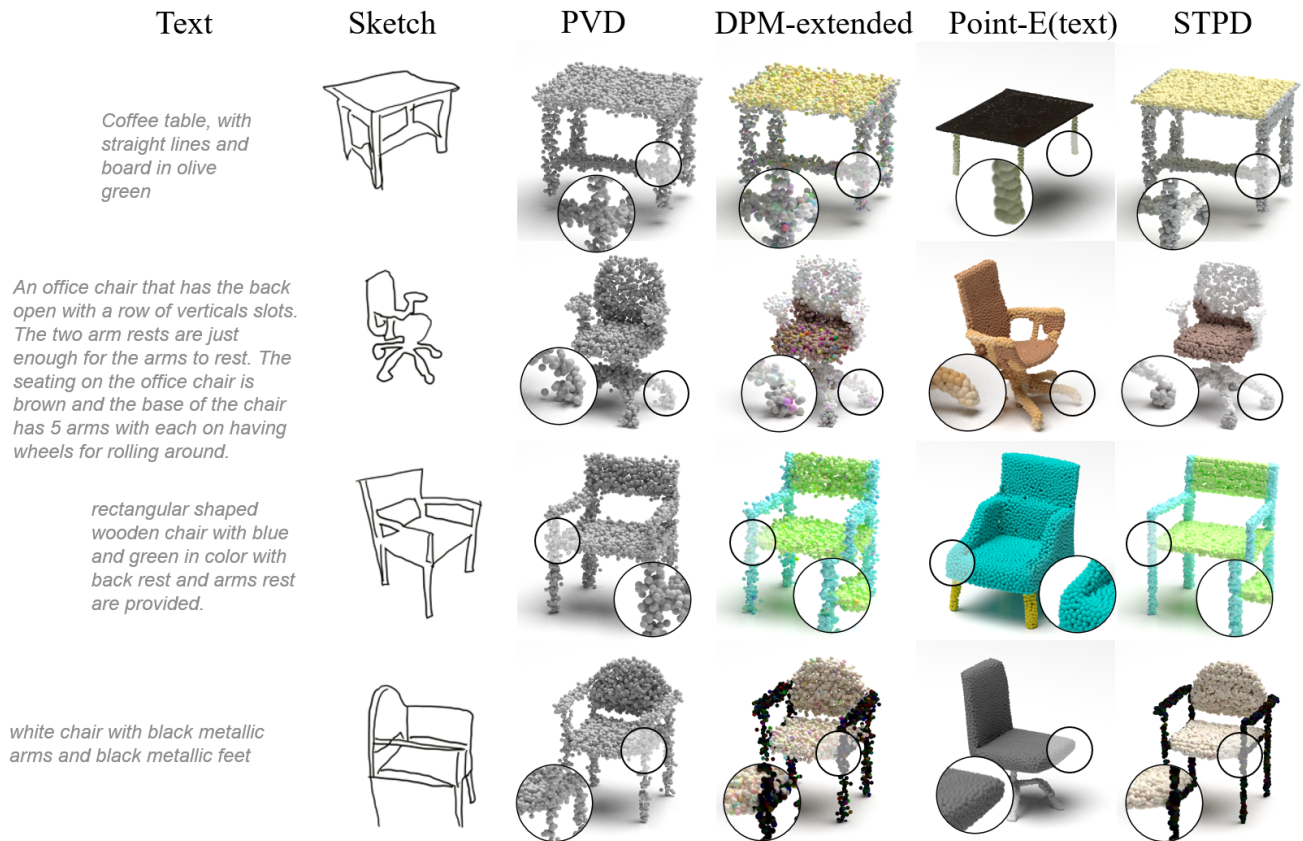*white chair with black metallic arms and black metallic feet*

Figure 3. Visualization comparison Between PVD [4], DPM [2], Point-E [3], and our STPD. We enlarge the local details to demonstrate that our STPD outperforms other baselines in terms of geometry and color generation. Our method produces faithful results from the input conditions.

ing 3d point clouds from complex prompts. *arXiv preprint arXiv:2212.08751*, 2022. 1, 4

[4] Linqi Zhou, Yilun Du, and Jiajun Wu. 3d shape generation and completion through point-voxel diffusion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5826–5835, 2021. 1, 4
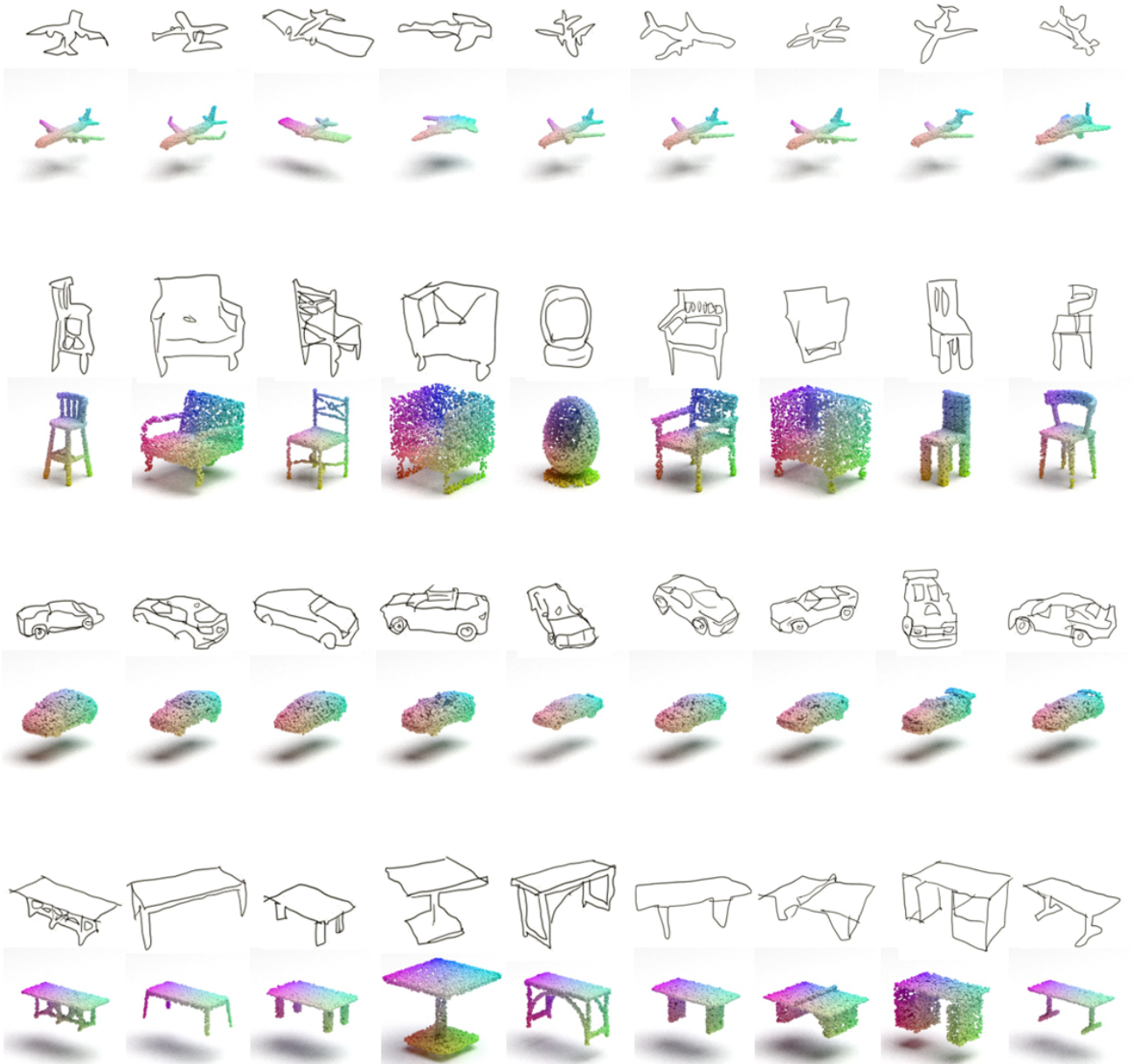
Figure 4. Visualization of point cloud generations for *airplane*, *chair*, *car*, *table* category from free-hand sketches. Our model achieves good generalization ability for generating different shapes that are faithful to the condition provided by the respective sketch.



Figure 5. Feature discussions. (a) generating novel shapes, and (b) effects of slight changes in the sketch, and (c) conflicting sketch & text. The main change part of sketch is highlighted with a green box.