

Blended Point Cloud Diffusion for Localized Text-guided Shape Editing

Supplementary Material

001 Contents

002	1. Implementation Details	1
003	1.1. Training	1
004	1.2. Evaluation	1
005	2. Additional Experiments	2
006	2.1. Comparison to Semantic Editing Paradigms	2
007	2.2. Ablation Study for t_r Values	2
008	2.3. Multi Category Training	3
009	2.4. Per-Category Evaluation	3
010	2.5. Extended Qualitative Results	3
011	2.6. Comparison to RePaint	4
012	2.7. Shape Completion Comparison	6

013	3. Limitations	7
-----	----------------	---

014 We refer readers to the interactive visualizations at [index.html](#). In this document, we provide implementation de-
 015 tails (Section 1), additional experiments and comparisons
 016 (Section 2) and a discussion of limitations (Section 3).

018 1. Implementation Details

019 1.1. Training

020 We trained our model for each category using object-
 021 specific ShapeTalk and l-ShapeTalk subsets. To evaluate
 022 generalization, we also trained a unified model across all
 023 three categories (Chair, Table, and Lamp). For all models,
 024 we used a batch size of 6 and a learning rate of 11×10^{-4} .
 025 The number of epochs, iterations, and training hours for
 026 each category are detailed in Table 1. To encourage the net-
 027 work to rely more on structural guidance rather than text, we
 028 dropped the text guidance with a probability of 0.5, replac-
 029 ing the textual prompt with an empty string. As explained in
 030 Section 3.2.2 of the main paper, we also replaced the con-
 031 ditional point cloud input with the target point cloud with
 032 a probability of 0.1 to support our Inversion-Free Coordi-
 033 nate Blending mechanism. All training was conducted on a
 034 single RTX A5000 GPU (24GB VRAM) for each model.

035 **Data.** To generate the partial point clouds used as guidance
 036 during training and to construct the l-ShapeTalk dataset,
 037 we used the baseline Llama 3 [5] model provided by [un-](#)
 038 [sloth](#). Specifically, we instructed Llama 3 with the follow-
 039 ing prompt for chairs:

040 ### Instruction: Return a single word
 041 using only one of the following options.
 042 Options: back, leg, arm, seat, unknown.

Model	Iterations	Epochs	Training Time (hours)
Chair	13.5 M	250	119
Table	12.1 M	250	107
Lamp	8.2 M	250	74
Airplane	0.9 M	250	52
Guitar	0.4 M	250	31
Knife	0.3 M	250	25
Cap	0.2 M	250	18
Skateboard	0.1 M	250	16
Unified	13.5 M	100	98

Table 1. **Training time.** We report the number of iterations, epochs and the overall training time for each category.

Input: What part of a chair does the
 043 next utterance describe?
 044 If none of the parts are described in
 045 the utterance return unknown. Utterance:
 046 it has larger height.
 047 #### Response: unknown (EOT-token)
 048

We adjusted the prompt for each category, instructing
 049 the language model to extract appropriate part names spe-
 050 cific to that category. To enhance accuracy, we fine-tuned
 051 the model for this task using 200 manually labeled exam-
 052 ples. Samples where the model returned “unknown” were
 053 excluded from the training set, leading to the creation of the
 054 l-ShapeTalk dataset.
 055

For evaluation, we applied the same method but removed
 056 the option to return “unknown”, requiring the language
 057 model to extract the part name it deemed most suitable to
 058 convey the text prompt.
 059

The part name in each l-ShapeTalk sample was provided
 060 as input to a segmentation model to extract binary masks.
 061 For segmentation, we used the PyTorch implementation of
 062 PointNet by [ailia-models](#). These binary masks were in-
 063 verted and then multiplied element-wise with the target
 064 point cloud to generate the guidance partial point cloud.
 065

1.2. Evaluation

ShapeTalk prompts often describe relationships between
 067 objects (e.g., “it has a taller backrest”), which makes them
 068 unsuitable for evaluating individual instances. To address
 069 this, we use Llama 3 to translate ShapeTalk prompts into
 070 more descriptive prompts, such as converting “it has a
 071 taller backrest” to “a chair with a tall backrest”.
 072

We use two CLIP-based [9] metrics, CLIP_{Sim} and
 CLIP_{Dir} , to evaluate edit fidelity. For the CLIP image en-
 073 coder, we rendered a single image for each point cloud from
 074 a consistent viewpoint.
 075

077 CLIP_{Sim} evaluates the similarity between the final output
078 and its textual description. We encoded the rendered
079 image of the output point cloud and the descriptive prompt
080 using their respective CLIP encoders and calculated the co-
081 sine similarity between the resulting encodings.

082 CLIP_{Dir} assesses the semantic direction in CLIP space,
083 capturing the relationship between both the guidance and
084 output shapes. We encoded the rendered images of the input
085 and output point clouds, along with the descriptive prompt
086 and a simple prompt describing the general object shape
087 (e.g., "a chair"). The differences between the two text en-
088 codings and the two image encodings were calculated, and
089 their cosine similarity was used as the metric.

090 All experiments were conducted using the "openai/clip-
091 vit-base-patch32" model, accessed through the [Hugging](#)
092 [Face](#) library.

093 To calculate the rest of the metrics used in the quanti-
094 tative analysis we used the [evaluate_change_it_3d.py](#) script
095 available in [ChangeIt3D's github repo](#).

096 **User Study.** As detailed in Section 4.2, our user study con-
097 sisted of two forms, A and B, each containing 15 ques-
098 tions and answered by 60 different users. In each ques-
099 tion, users were shown a text prompt and an input point
100 cloud, then asked to select one of the editing results gen-
101 erated by ChangeIt3D, Spice-E, and our method, displayed
102 in random order. Both the prompts and input point clouds
103 were randomly selected from the l-ShapeTalk test set. The
104 users were instructed:

105 *"Your task is to select the target object that best repre-
106 sents the prompt while maintaining the shape of the source
107 object as closely as possible. In other words, choose the
108 most suitable target object that effectively embodies the
109 editing of the source object based on the prompt. IMPOR-
110 TANT - If none of the target objects align with the prompt,
111 select the target object that best preserves the shape of the
112 source object."*

113 **Baselines.** We used the official [Changeit3D implementa-](#)
114 [tion](#) to access the ShapeTalk dataset and the pre-trained
115 model weights. The script [evaluate_change_it_3d.py](#) was
116 employed to reproduce their results. Similarly, we used the
117 [Spice-E implementation](#) for obtaining their result. The out-
118 puts were converted into point clouds using marching cubes
119 and random point sampling.

120 2. Additional Experiments

121 To better view 3D results, we recommended viewing the
122 supplemental HTML page, which includes fly-through vi-
123 sualizations demonstrating the quality of our results from
124 multiple views.

125 2.1. Comparison to Semantic Editing Paradigms

In this section, we complement our main paper's compar-
126 isons with a qualitative analysis of alternative approaches
127 to text-guided semantic editing of 3D shapes.

Image Editing and Single View Reconstruction. Recent
129 advancements in image editing and 3D reconstruction meth-
130 ods suggest another potential paradigm: perform image
131 editing on the rendered image of a 3D shape and then use
132 single-view reconstruction to generate an edited 3D shape.
133 We tested this approach by using InstructPix2Pix [12] in
134 conjunction with One-2-3-45++ [6]. Specifically, we ren-
135 dered an image of the input shape, provided it along with an
136 editing prompt to InstructPix2Pix, and used the edited im-
137 age as input to One-2-3-45++ to reconstruct the edited 3D
138 shape. As shown in Figure 3, InstructPix2Pix fail to per-
139 form fine-grained shape editing. To address this, we also
140 fine-tuned InstructPix2Pix using images rendered from the
141 ShapeTalk dataset, as also presented in Figure 3. While
142 the fine-tuned model produces higher-quality results, it still
143 fails to localize edits effectively, especially compared to our
144 method. For single-view reconstruction, shapes generated
145 from InstructPix2Pix results using One-2-3-45++ are shown
146 in Figure 1. Beyond InstructPix2Pix's difficulty with pre-
147 cise localization, One-2-3-45++ introduces slight defects in
148 the reconstructed shapes, such as an asymmetrical backrest
149 (second row from the top) and slightly lopsided legs (third
150 row from the top). However, we note that these defects are
151 minor. The overall high visual quality of the results demon-
152 strates promise for future research in this direction.

Optimization-Based Editing. A widely used approach for
154 textual editing of 3D shapes is Score Distillation Sampling
155 (SDS) [8], which employs inference-time optimization to
156 modify a 3D representation based on a text prompt, using
157 a pre-trained generator as a prior. A qualitative comparison
158 with the SDS based methods Fantasia3D [3] and Vox-E [10]
159 is presented in Figure 3. As shown, while Fantasia3D gen-
160 erates results that resemble the input shapes, it often fails to
161 follow the fine-grained instructions in the editing prompts
162 (e.g., the back is not rounded in the first row, and the chair
163 in the second row lacks four legs). Vox-E also often does
164 not follow the editing instruction (the chair on the third row
165 does not have noticeably thinner legs, chair on the fourth
166 row does not seem to have a taller backrest), but also often
167 fails to correctly maintain the identity of the input objects
168 (adding legs to the chairs on the first and last rows). Addi-
169 tionally, as is common with many SDS-based methods, both
170 methods tend to exhibit noise.

171 2.2. Ablation Study for t_r Values

The parameter t_r controls the balance between coordinate
172 blending steps and steps dedicated solely to shape recon-
173 struction. Higher t_r values increase the number of coor-
174 dinates used in the coordinate blending step, which leads
175 to smoother surfaces but slower inference times.

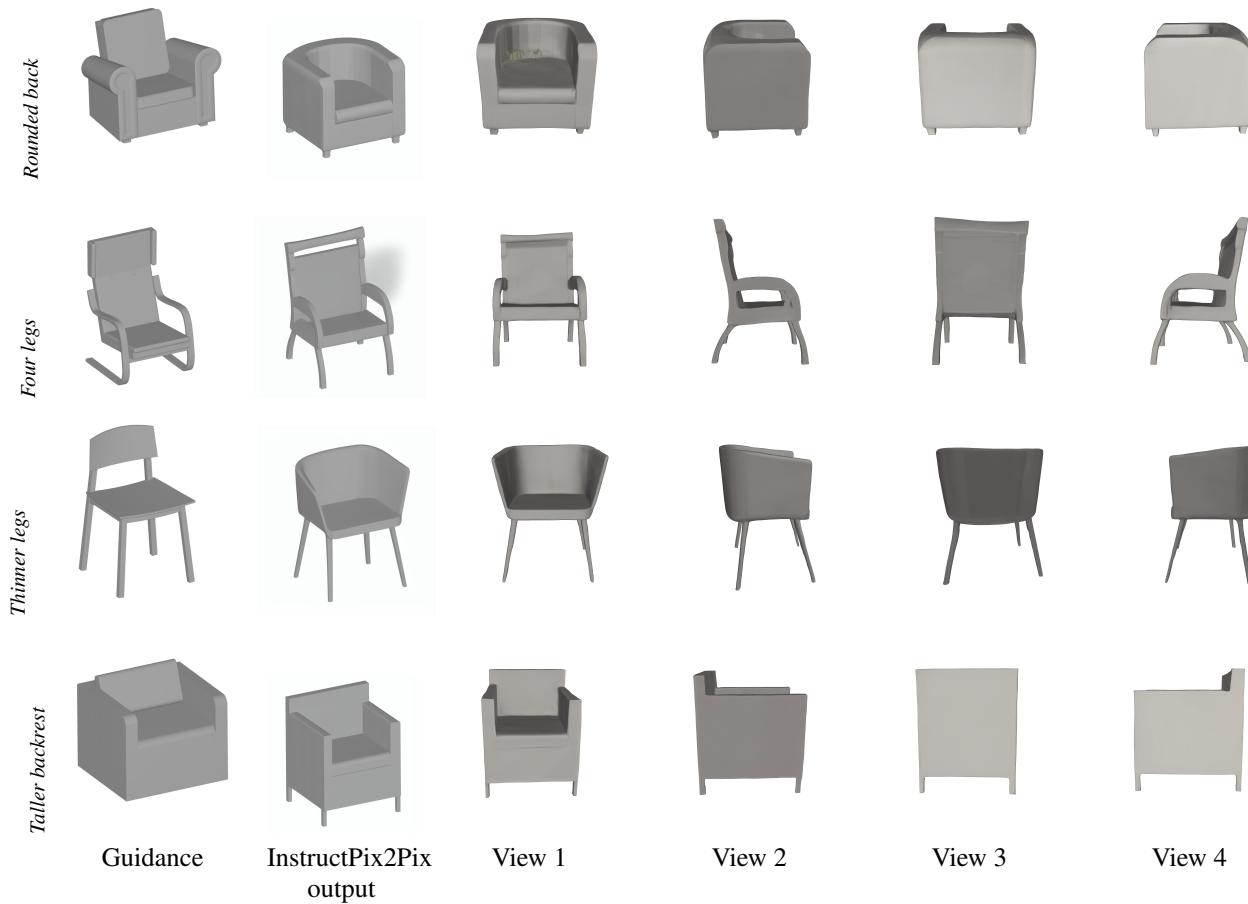


Figure 1. Image Editing Followed by Single-View 3D Reconstruction. We first use InstructPix2Pix to edit the rendered image and then apply One-2-3-45++ for single-view reconstruction. In addition to InstructPix2Pix’s challenges with precise localization, One-2-3-45++ introduces minor defects in the reconstructed shapes.

176 dinate blending steps, providing greater editing freedom
 177 but at the cost of identity preservation. Conversely, lower
 178 t_r values improve identity preservation while reducing the
 179 model’s ability to adhere closely to the text prompt.

180 Figure 2 illustrates the effect of different t_r values. As
 181 t_r decreases, the output aligns more closely with the input
 182 point cloud. While higher t_r values allow greater editing
 183 freedom, they can result in inconsistencies with the masked
 184 guidance point cloud. We set $t_r = 20$ in our experiments as
 185 it offers a good balance between identity preservation and
 186 edit fidelity.

187 **2.3. Multi Category Training**

188 Table 2 demonstrates that our technique is not restricted to
 189 single-category training. We trained our model on a uni-

190 fied dataset comprising all three categories combined. The
 191 results show that our method handles multiple categories
 192 effectively, with the unified model producing results com-
 193 parable to those of the single-category models.

194 **2.4. Per-Category Evaluation**

195 Table 3 provides the evaluation results for each of the three
 196 object categories used for quantitative evaluation. The per-
 197 category scores align closely with the overall average scores
 198 reported in the main paper.

199 **2.5. Extended Qualitative Results**

200 In Figure 4 we present qualitative results from multi-
 201 ple shapeNet categories (*Guitar*, *Airplane*, *Chair*, *Lamp*,
 202 *Sword*, *Hat*, *Skateboard* and *Table*), demonstrating our

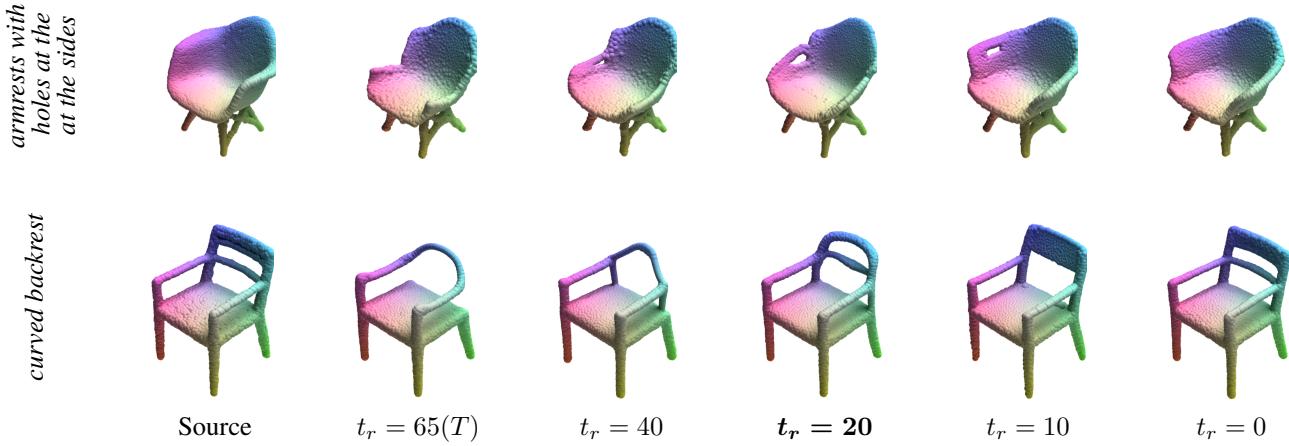


Figure 2. **Ablation Study for t_r Values.** Higher t_r values increase the number of coordinate blending steps, providing better editing freedom but at the cost of inferior identity preservation. Conversely, lower t_r values improve identity preservation while reducing the model’s ability to follow the textual description.

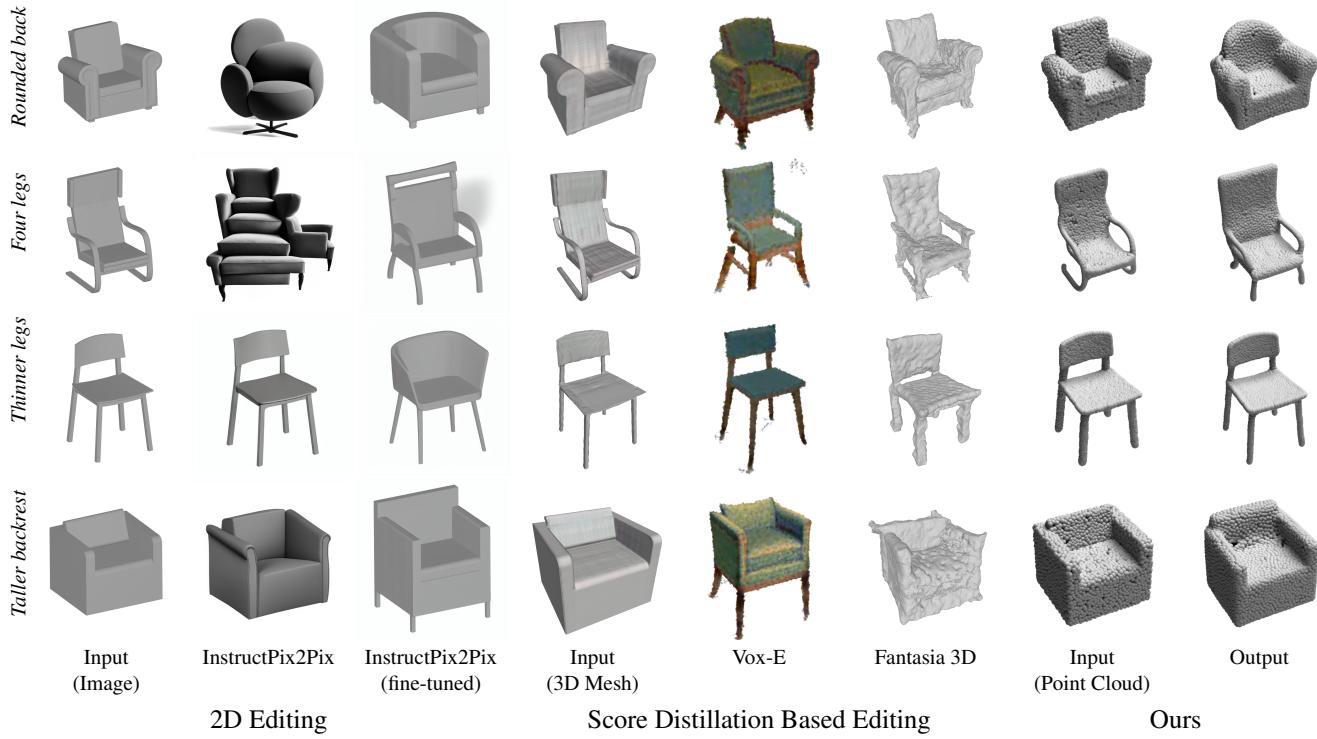


Figure 3. **Qualitative comparison.** We compare our method’s outputs to those of the image editing method InstructPix2Pix [2] as well as the score distillation sampling optimization based 3D editing works Fantasia 3D [3] and Vox-E [10]. As illustrated above, our method outperforms these baselines in terms of edit fidelity, identity preservation and overall visual quality.

203 method’s ability to make meaningful fine grained edits
204 across a wide variety of shape types.

205 2.6. Comparison to RePaint

206 The free form inpainting method RePaint [7] has had a major impact on the field of diffusion based image inpainting.
207

208 This work introduced an inference time algorithm which,
209 somewhat similarly to our coordinate blending algorithm,
210 blends noisy versions of the “known” regions of the image
211 with the predicted denoised versions of the inpainted region
212 according to an input binary mask. RePaint also proposes
213 to resample noise and repeat the process for a given num-

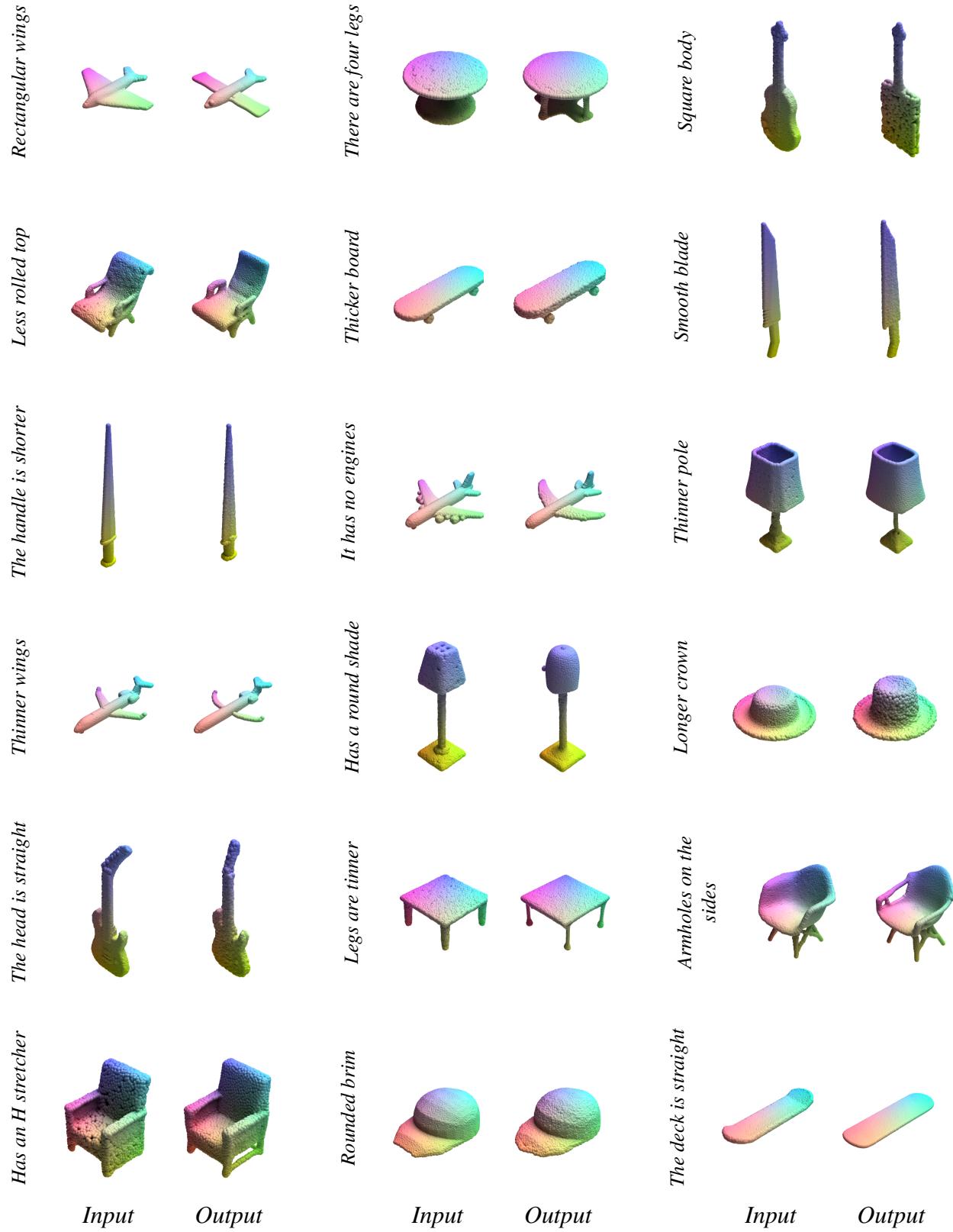


Figure 4. **Results Gallery.** Above we show results over various object categories including *chair*, *table*, *lamp*, *airplane*, *cap*, *guitar*, *skateboard* and *knife*.

Metric	Shapetalk						l-Shapetalk					
	CLIP _{Sim} ↑	CLIP _{Dir} ↑	GD↓	CD↓	FPD↓	l-GD↓	CLIP _{Sim} ↑	CLIP _{Dir} ↑	GD↓	CD↓	FPD↓	l-GD↓
Unified Model	0.26	0.00	0.32	0.03	43.05	0.68	0.26	0.00	0.36	0.09	138.26	0.82
Ours	0.26	0.01	0.34	0.05	33.64	0.78	0.27	0.01	0.29	0.04	13.51	0.55

Table 2. **Comparison to Unified Model.** Our technique performs effectively across multiple categories training, as evidenced by the comparable results.

Metric	Shapetalk						l-Shapetalk						
	CLIP _{Sim} ↑	CLIP _{Dir} ↑	GD↓	CD↓	FPD↓	l-GD↓	CLIP _{Sim} ↑	CLIP _{Dir} ↑	GD↓	CD↓	FPD↓	l-GD↓	
<i>Chair</i>	Changeit3D	0.22	-0.02	0.48	0.05	127.35	0.79	0.22	-0.02	0.68	0.07	156.11	0.87
	Spice-E	0.26	0.00	1.78	0.18	527.33	1.01	0.26	0.00	1.78	0.22	653.16	0.93
	Ours	0.28	0.01	0.27	0.02	14.13	0.74	0.28	0.01	0.24	0.01	5.87	0.54
<i>Table</i>	Changeit3D	0.22	-0.03	0.55	0.16	111.27	1.00	0.22	-0.04	0.66	0.13	141.12	1.04
	Spice-E	0.25	-0.01	1.85	0.40	489.32	0.88	0.26	-0.04	2.85	0.40	621.43	1.07
	Ours	0.27	0.01	0.33	0.03	18.96	0.71	0.26	0.01	0.33	0.03	11.50	0.51
<i>Lamp</i>	Changeit3D	0.19	-0.01	0.93	0.33	310.44	0.66	0.19	-0.02	1.01	0.38	354.63	0.75
	Spice-E	0.23	-0.02	1.88	0.14	153.41	0.94	0.23	-0.02	2.16	0.16	188.56	0.95
	Ours	0.24	0.00	0.41	0.09	67.82	0.88	0.26	0.01	0.31	0.07	23.18	0.60

Table 3. **Per-Category Evaluation**. We compare the performance of ChangeIt3D [1] and Spice-E [11] against ours over the three object categories in the ShapeTalk and l-ShapeTalk datasets.

Metric	CLIP _{Sim} ↑	CLIP _{Dir} ↑	GD↓	CD↓	FPD↓	l-GD↓
RP _{r=1,j=1}	0.22	-0.02	0.57	0.12	69.11	0.72
RP _{r=10,j=1}	0.19	-0.02	0.62	0.15	79.33	0.55
RP _{r=1,j=10}	0.19	-0.03	0.63	0.14	63.48	0.71
RP _{r=10,j=10}	0.21	-0.03	0.72	0.13	64.17	0.57
Ours	0.27	0.01	0.29	0.04	13.51	0.55

Table 4. **Quantitative Evaluation** against RePaint (RP) with different resampling (r) and jumping (j) values. Note that this implementation uses our InPaint-E model and our reconstructed noise on the non-edit regions as directly using Point-E and random noise resulted in highly noisy (and uninformative) outputs.

ber of iterations to further refine this inpainting process. By contrast, in addition to other core differences our work dedicates a significant portion of the inference process to reconstructing the full input shape before starting the inpainting process, as well as operating on a specific model tailored for this task instead of a more general text-to-3D model. To test the significance of some of these design choices we conducted a quantitative comparison against a baseline which resembles RePaint in its function. Specifically, in this baseline inpainting is performed at every inference step ($t_r = T$) and the edit region is initialized with random noise. This baseline also incorporates the “resampling” and “jumping” mechanisms introduced in RePaint. Unlike RePaint however, we used our reconstructed noise for the non-edit region and operated on Inpaint-E, as directly using Point-E and random noise resulted in highly noisy (and uninformative) outputs.

The results of this comparison are presented in Table 4 and clearly shows that our method outperforms this baseline across all metrics.

2.7. Shape Completion Comparison

Shape, or part, completion is a longstanding task involving completing a shape that is missing one or more parts in a way that maintains plausibility and optionally aligns with a text prompt. In Figure 5, we illustrate that these methods are not well suited for fine-grained shape editing as they are inherently blind to the missing parts, focusing on a comparison with SDFusion [4]. Note that we demonstrate this in a slightly different setting in comparison to the one addressed in our work—*unconditional* generation, and not *text-guided* part completion—as their official codebase does not include text-guided part completion.

Nonetheless, we performed a qualitative comparison in which we compared our results (text-guided) against unconditional part completion results of SDFusion. These results show that this method’s ability to maintain identity is somewhat limited, even outside of the edit region (thicker legs on the chair in the second row, as well as adding arms to it). As SDFusion’s part completion in this case is not guided by text it is somewhat hard to judge its quality. However, it is evident that the completed parts often don’t match the general identity of the original shape particularly well (office chair type backrest in row 2, huge elaborate legs in row 1)

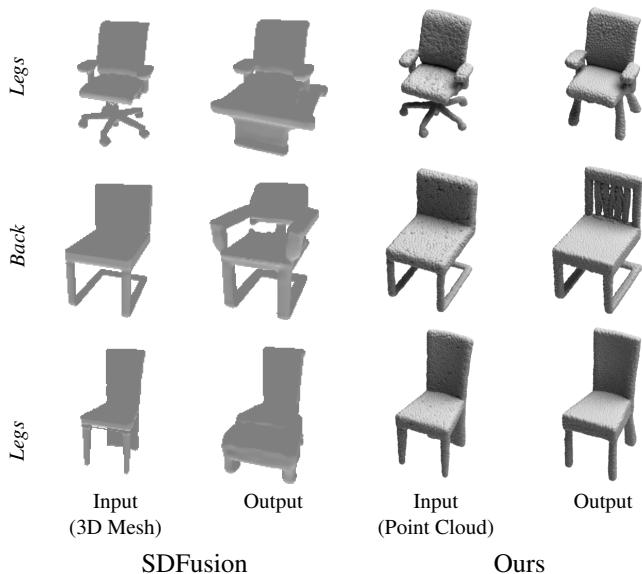


Figure 5. **Qualitative comparison against SDFusion [4] unconditional part completion.** We compare our method’s outputs against the unconditional part completion results of the SDFusion [4] baseline. In each row we task the methods with completing (in SDFusion’s case) or editing¹ (in our case) a different part. As these results show, SDFusion’s ability to preserve identity across all regions of the shape is limited, in comparison to our approach that performs localized fine-grained editing of 3D shapes.

257

compared to our method.

258

3. Limitations

259
260
261
262
263
264
265
266
267
268
269
270
271

While our method performs well in most scenarios, it has certain limitations. First, our inpainting-based approach restricts the ability to generate entirely new objects or parts. BlendedPC is specifically designed for localized editing and struggles with global shape transformations. Second, as a supervised learning approach, the method’s generalizability is constrained by the object categories in the training dataset. This limits performance when editing shape categories not encountered during training. These limitations present exciting opportunities for future research, such as developing techniques for generating object parts, enabling global shape modifications, and improving cross-category generalization.

272

References

273
274
275
276
277

- [1] Panos Achlioptas, Ian Huang, Minhyuk Sung, Sergey Tulyakov, and Leonidas Guibas. Changeit3d: Language-assisted 3d shape edits and deformations. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, page 6, 2022. 6

¹The text guidance provided for our method (from top to bottom): *four legs, spindle backrest, four legs*.

- [2] Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18392–18402, 2023. 4
- [3] Rui Chen, Yongwei Chen, Ningxin Jiao, and Kui Jia. Fantasia3d: Disentangling geometry and appearance for high-quality text-to-3d content creation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 22246–22256, 2023. 2, 4
- [4] Yen-Chi Cheng, Hsin-Ying Lee, Sergey Tulyakov, Alexander G Schwing, and Liang-Yan Gui. Sdfusion: Multimodal 3d shape completion, reconstruction, and generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4456–4465, 2023. 6, 7
- [5] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024. 1
- [6] Minghua Liu, Ruoxi Shi, Linghao Chen, Zhuoyang Zhang, Chao Xu, Xinyue Wei, Hansheng Chen, Chong Zeng, Jiayuan Gu, and Hao Su. One-2-3-45++: Fast single image to 3d objects with consistent multi-view generation and 3d diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10072–10083, 2024. 2
- [7] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11461–11471, 2022. 4
- [8] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022. 2
- [9] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 1
- [10] Etaï Sella, Gal Fiebelman, Peter Hedman, and Hadar Averbuch-Elor. Vox-e: Text-guided voxel editing of 3d objects. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 430–440, 2023. 2, 4
- [11] Etaï Sella, Gal Fiebelman, Noam Atia, and Hadar Averbuch-Elor. Spice-e: Structural priors in 3d diffusion using cross-entity attention. In *ACM SIGGRAPH 2024 Conference Papers*, pages 1–11, 2024. 6
- [12] Jiale Xu, Xintao Wang, Yan-Pei Cao, Weihao Cheng, Ying Shan, and Shenghua Gao. Instructp2p: Learning to edit 3d point clouds with text instructions. *arXiv preprint arXiv:2306.07154*, 2023. 2