# Towards Implicit Text-Guided 3D Shape Generation

Zhengzhe Liu[1]    Yi Wang[2]    Xiaojuan Qi[3*]    Chi-Wing Fu[1*]

[1]The Chinese University of Hong Kong    [2]Shanghai AI Laboratory    [3]The University of Hong Kong

{zzliu,cwfu}@cse.cuhk.edu.hk    wangyi@pjlab.org.cn    xjqi@eee.hku.hk

## Abstract

*In this work, we explore the challenging task of generating 3D shapes from text. Beyond the existing works, we propose a new approach for text-guided 3D shape generation, capable of producing high-fidelity shapes with colors that match the given text description. This work has several technical contributions. First, we decouple the shape and color predictions for learning features in both texts and shapes, and propose the word-level spatial transformer to correlate word features from text with spatial features from shape. Also, we design a cyclic loss to encourage consistency between text and shape, and introduce the shape IMLE to diversify the generated shapes. Further, we extend the framework to enable text-guided shape manipulation. Extensive experiments on the largest existing text-shape benchmark [11] manifest the superiority of this work. The code and the models are available at* https://github.com/liuzhengzhe/Towards-Implicit-Text-Guided-Shape-Generation.

## 1. Introduction

3D shape creation has a wide range of applications, *e.g.*, CAD, games, animations, computational design, augmented reality, etc. Significant progress has been made in recent years by exploiting neural networks and generative models to learn to produce 3D shapes. Yet, existing works [8, 13, 14, 23, 35, 36, 48, 50, 73, 80, 82] focus mostly on generating the overall shapes, whereas the more recent ones [12, 15, 24, 45, 60, 83, 84] attempt to generate shapes with more details.

In this work, we are interested in the challenging task of *text-guided 3D shape generation*—Given a sentence, *e.g.*, "A comfortable red color chair with four legs," we aim to develop a method to automatically generate a 3D shape that follows the text description; see Figure 1 (a) for our example results. This research direction has great potential for efficient 3D shape production, say by taking user speech/text input to guide or condition the process of generating 3D shapes. By this means, we can assist users to readily generate and edit 3D models for diverse applications.
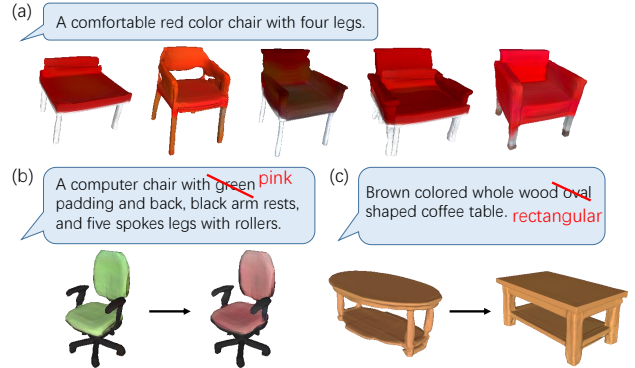
---

*: Corresponding authors



Figure 1. (a) Chairs of different structures and appearances generated by our method from the same given sentence. Our method also allows text-based manipulation in color (b) and in shape (c).

While many methods [41, 62–65, 69, 76, 77, 81, 86, 87] have been developed for generating 2D images from text, the task of generating 3D shapes from text is rather under-explored. Chen *et al*. [11] generate 3D shapes from natural language descriptions by learning joint text and shape embeddings, but the performance and visual quality are highly limited by the low-resolution 3D representations. Another very recent work [34] leverages semantic labels to guide the shape generation, but it requires predefined semantic labels and cannot directly deal with natural language inputs.

To enhance 3D shape generation from text, we propose a new solution by leveraging the implicit representation [14, 49, 55] to predict an occupancy field. Yet, several inherited challenges have not been addressed in the early works for properly adopting the implicit representation for the text-to-shape task. First, the above works generate shapes typically without colors, which are crucial in text-guided 3D shape generation, since text descriptions often contain colors; we empirically found that directly predicting shape and color with a single implicit decoder often lead to shape distortion and color blur. Second, text contains a large amount of spatial-relation information, *e.g.*, "a wooden table on a metal base." Still, spatial-relation local features are ignored in existing works, since the implicit decoder generally considers only the global feature from the auto encoder as input [14]. Third, the generated shapes are not all consistent with the

input texts, largely due to the semantic gap between text and 3D shape and also the lack of effective learning constraints. Last, text-to-shape generation is inherently one-to-many, *i.e.*, diverse results may match the same input text. Yet, the existing regression-based approach outputs only a single shape.

This work presents a new approach for high-fidelity text-guided 3D shape generation. First, we decouple the shape and color predictions for feature learning in both texts and shapes to improve the generation fidelity; this strategy also aids the text-guided shape manipulation. Also, we introduce a word-level spatial transformer to learn to correlate the word features with the spatial domain in shapes. In addition, we design a cyclic loss to encourage the consistency between the generated 3D shape and the input text. Further, we propose a novel style-based latent shape-IMLE generator for producing diversified shapes from the same given text. Last, we extend the framework for text-guided 3D shape manipulation with a two-way cyclic loss. As shown in Figure 1 (b), we may modify the original text and our framework can produce new colored shapes according to the edited text, while keeping the other attributes unchanged.

Extensive experiments on the largest existing text-shape dataset [11] demonstrate the superiority of our approach over the existing works, both qualitatively and quantitatively.

## 2. Related Work

**Text-to-image generation.** Remarkable progress has been made for generating images from text [41, 42, 62–64, 76, 81, 86, 87]. Recently, approaches [57, 65, 68, 69, 77, 79, 85] based on the unconditional GAN [7, 37, 38] were also proposed.

Compared with text-to-image, it is more challenging to generate 3D shapes from texts. First, unlike 2D images, 3D shapes are unstructured and irregular without well-defined grid structures. Also, the text-to-shape task requires a comprehensive prediction of the whole 3D shape, while the text-to-image task addresses image generation, which is a projection of the 3D shape. Further, there are plenty of large-scale image datasets [47, 52, 75] to support text-to-image. Yet, as far as we know, the largest dataset for text-to-shape was proposed in [11], which has $75k$ texts and $15k$ shapes of $128^3$ resolution. The lack of large-scale and high-quality training data makes the text-to-shape task even harder.

**3D shape representations, generation, and manipulation.** Unlike images, 3D shapes can be represented as, *e.g.*, voxel grids [19, 25], point clouds [2, 61], and meshes [22]. Also, various methods [31, 39, 40, 46, 70] have been proposed for generating and manipulating shapes for different 3D representations. Yet, the generated shapes are limited by the resolution and quality of the training set. To generate shapes of arbitrary resolution, recent works [14, 17, 48, 49, 55] start to explore implicit functions, which in fact have been used in many tasks, *e.g.*, single-view reconstruction [45, 51, 80], 3D scene reconstruction [30, 36, 59], and 3D texture gen-

eration [18, 53, 54]. In existing works, a typical approach is to leverage an auto-encoder (AE) to adopt to multiple 3D generation tasks and map the input modalities into the AE's learned feature space, *e.g.*, single-image 3D reconstruction [14, 80], point-cloud-based shape generation [8, 16], and 3D completion [78].

Following the above works, a straightforward approach for text-guided 3D shape generation is to map the text feature into the AE's feature space then adopt an implicit decoder to generate the 3D shape. This simple approach, however, has several drawbacks, as discussed in Section 1.

Recently, several works make it possible to manipulate implicit 3D shapes [20, 28, 33, 88] using a reference box or reference points as guidance. Yet, none of them enables 3D shape manipulation with natural language descriptions.

**3D shape generation from text.** A series of works are proposed to address the tasks on texts and 3D shapes, including learning the text-shape correspondence [3], cross-modal retrieval [27, 72], shape-to-text generation [26], text-guided shape composition [32], and 3D object localization [10].

As far as we are aware of, there are only few works [11, 34] that address the challenging text-to-shape task. Chen *et al*. [11] propose to directly predict colored voxels with adversarial learning on top of a jointly-learned text-shape embedding. Though plausible shapes can be produced, the shape resolution and texture quality are still far from being satisfactory. Also, the generated shapes may not be consistent with the input texts due to the large semantic gap between text and shape. Jahan *et al*. [34] propose a semantic-label guided shape generation approach; however, it can only take one-hot semantic keywords as input and the generated shapes are also unsatisfying in quality, without color and texture.

This work presents a new framework, capable of generating high-fidelity 3D shapes with good semantic correspondence between the text and shape. Also, our framework enables text-guided 3D shape manipulation for both shape and color, outperforming the existing works by a large margin, as demonstrated in the experiments.

**Diversified generation.** Besides GANs, IMLE (Implicit Maximum Likelihood Estimation) is another approach to aid multi-modal generation, *e.g.*, super-resolution [43], semantic-layout-guided image synthesis [44], image decompression [58], and shape completion [5]. Compared with GANs, IMLE mitigates the mode collapse of GANs and boosts the result diversity. In this work, we leverage IMLE for generating multiple shapes from the same text input.

## 3. Methodology

### 3.1. Overview

Given text $\mathbf{T}$, we aim to generate high-quality 3D shape $\mathbf{S}$ with colors, following the description of $\mathbf{T}$. To generate high-quality results, we exploit the implicit occupancy
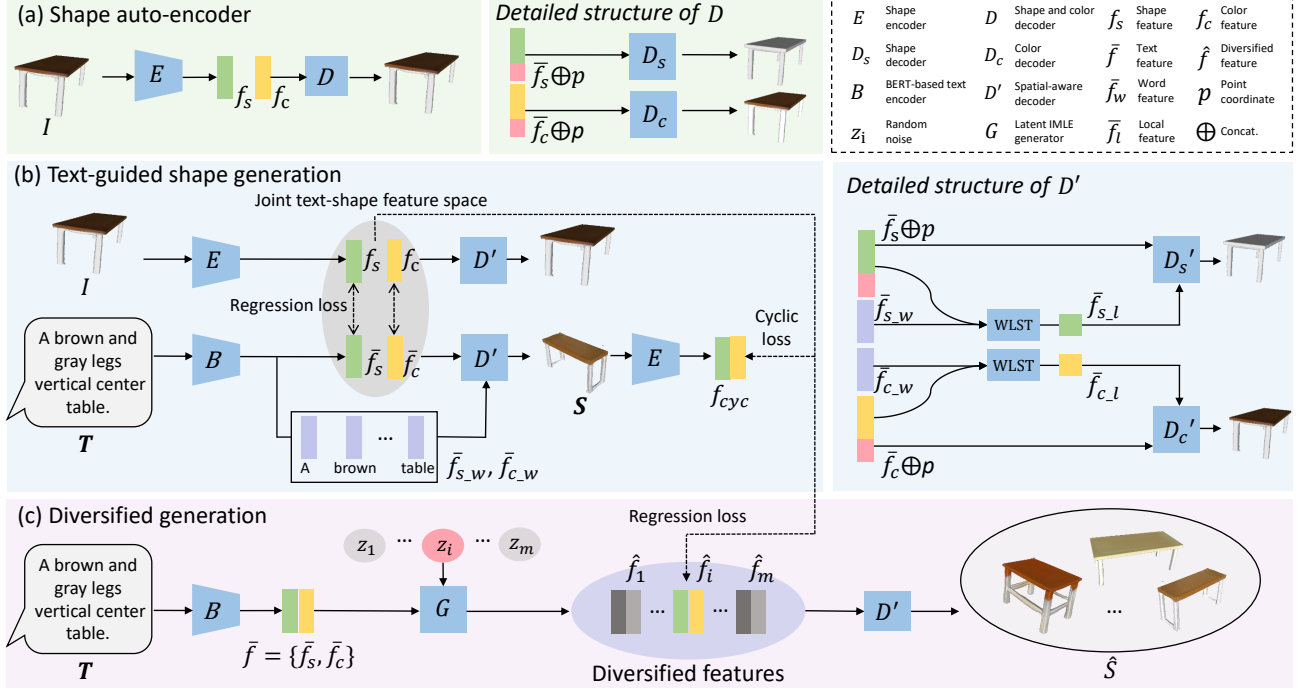
Figure 2. Overview of our text-guided shape generation framework, which has three major parts. (a) First, the shape auto-encoder $\{E, D\}$ extracts shape feature $f_s$ and color feature $f_c$ from the input 3D shape $I$. (b) We then learn to generate the 3D shape in a text-guided manner with the word-level spatial transformer (WLST) and the cyclic consistency loss $f_{cyc}$. (c) Further, we generate diversified 3D shapes from the same given text by adopting a style-based latent shape generator $G$. We only need (c) during the inference.

representation other than the explicit voxel/point/mesh representations to characterize shapes with color. Specifically, the predicted shape with color is denoted as $\mathbf{S} \in \mathbb{R}^{N \times (1+3)}$, including the shape $\in \mathbb{R}^{N \times 1}$(a set of occupancy values in the voxels) and the color $\in \mathbb{R}^{3N}$ (the associated set of RGB values), respectively, where $N$ is the number of sample points, concerning the generation quality.

Our framework consists of a text encoder $B$, feature generator $G$, spatial aware decoder $D'$, and shape encoder $E$. Its overall architecture is given in Figure 2. In inference, $B$ extracts text feature $\bar{f} = \{\bar{f}_s, \bar{f}_c\}$ from text $\mathbf{T}$ (where $\bar{f}_s$ and $\bar{f}_c$ are the shape and color portions of $\bar{f}$, respectively), and $G$ produces multiple instances of such feature $\{\hat{f}_i\}_{i=1}^m$ based on $\bar{f}$ conditioned on various random vectors $\{z_i\}$. Then, $D'$ generates diverse shapes $\{\mathbf{S}_i \in \mathbb{R}^{N \times (1+3)}\}_{i=1}^m$ with color.

The model training of our method is non-trivial. We train the overall framework in *three* stages (see again Figure 2): (a) shape auto-encoder, (b) text-guided shape generation, and (c) diversified shape generation with IMLE. Specifically,

- First, as shown in Figure 2 (a), we train shape encoder $E$ and implicit decoder $D$. As shown in the top middle, unlike existing works [14, 34] that ignore colors in the shape generation, $D$ composes of $D_s$ and $D_c$ that account for the decoding of shape and color, respectively, when $D$ predicts the output shape.

- Then, we adopt BERT-based text encoder $B$ [21] to help

extract text feature $\bar{f} = \{\bar{f}_s, \bar{f}_c\}$ and word-level feature $\bar{f}_w = \{\bar{f}_{s\_w}, \bar{f}_{c\_w}\}$ from input text $\mathbf{T}$ (see Figure 2 (b)), and map $\bar{f}$ into the joint text-shape feature space to reduce the domain gap between the text and the shape. Further, we propose the spatial-aware decoder $D'$ to leverage local feature $\bar{f}_l$ extracted by the word-level spatial transformer (WLST), which explicitly correlates the spatial and word features to improve the fidelity of $\mathbf{S}$. Also, we formulate cyclic loss $L_{cyc}$ to encourage the consistency between shape $\mathbf{S}$ and text $\mathbf{T}$.

- Lastly, we propose to adopt style-based shape generator $G$ that conditions on a set of random noise vectors $\{z_i\}_{i=1}^m$ to enable diversified 3D shape generation with feature $\hat{f}_i$, as shown in Figure 2 (c).

In the following, we will detail each component of the framework and the associated losses.

## 3.2. Shape Auto-Encoder

We extend the auto-encoder in [14] to jointly reconstruct the shape and color. As shown in Figure 2 (a), our shape auto-encoder aims to map the input voxel-based shape $I \in \mathbb{R}^{64 \times 64 \times 64}$ into a compact feature space. Specifically, encoder $E$ [14] extracts the shape and color features $f = \{f_s, f_c\}$ from $I$, whereas decoder $D$ reconstructs the shape and color through $D_s$ and $D_c$, respectively. Inside
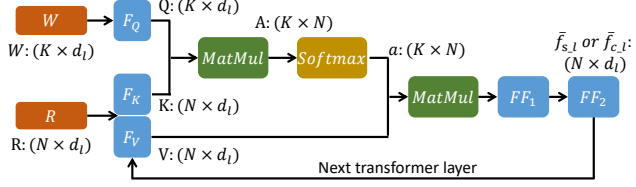
Figure 3. The Word-Level Spatial Transformer architecture. $F_Q$, $F_K$, and $F_V$ are fully-connected layers, whereas $FF_1$ and $FF_2$ are feed-forward networks. The Layer Normalization [6] is omitted.
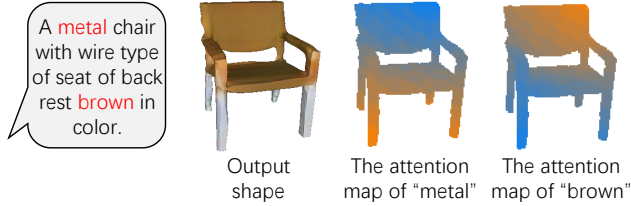


Figure 4. Visualizing the attention map $A$ for the words "metal" and "brown". Warmer colors indicate stronger correlation.

$D$, we concatenate a sample (or query) point coordinate $p = (x, y, z)$ with each feature vector ($f_s$ or $f_c$) as input to $D_s$ or $D_c$. $D_s$ and $D_c$ have the same architecture with seven fully-connected and leaky-ReLU layers, except in the last layer, $D_s$ outputs a single occupancy value and $D_c$ outputs three values for RGB color, both at the sample point $p$.

The shape auto-encoder is trained to reconstruct the shape and color of the input shape with an $L_2$ regression:

$$L_{ae} = \lambda_s \Sigma_p ||D_s(f_s \oplus p) - I(p)||_2^2$$
$$+ \lambda_c \Sigma_{k \in \{R,G,B\}} \Sigma_p ||D_c(f_c \oplus p)[k] - I(p)[k]||_2^2 \mathbb{1}(I(p)),$$
(1)

where $I(p)$ and $I(k, p)$ denote the ground-truth occupancy and color values, respectively, at point $p$; $\oplus$ denotes concatenation; $\mathbb{1}$ is an indicator function of value 1 if $p$ is inside the input shape, and 0, otherwise; and $\lambda_s$ and $\lambda_c$ are weights for the shape and color reconstructions, respectively.

### 3.3. Text-Guided Shape Generation

As shown in Figure 2 (b), the text-guided shape generation network consists of three modules: shape encoder $E$, BERT-based text encoder $B$, and spatial-aware decoder $D'$. With $E$ and $D'$ ($D'_s$ and $D'_c$) initialized by the corresponding components in the shape auto-encoder, our goal here is to train the whole network end-to-end to obtain $B$ and $D'$.

**Text encoder $B$.** We employ the BERT structure [21] to build text encoder $B$ for extracting text feature $\bar{f}$ from input text $\mathbf{T}$ and mapping $\bar{f}$ to the joint text-shape feature space.

**Spatial-aware decoder $D'$.** $D'$ aims to transform text feature $\bar{f}$ to the predicted shape $\mathbf{S}$ with color. Instead of simply using the trained implicit decoder $D$, we construct the spatial-aware decoder $D'$ with the word-level spatial transformer (WLST). In short, we take the local features from WLST to improve the spatial correlation implied from $\mathbf{T}$.

The right side of Figure 2(b) shows the architecture of the spatial-aware decoder $D'$. First, we concatenate $\bar{f}_s$ and $p$ and transform the result $\{\bar{f}_s \oplus p\} \in \mathbb{R}^{N \times (d+3)}$ using a fully-connected layer, where $N$ is the number of sample points for shape reconstruction and $d$ is the channel dimension of $\bar{f}_s$. Then, we transform the word-level BERT features $\{\bar{f}_w\} \in \mathbb{R}^{K \times d_B}$ (where $K$ is the number of words in input text and $d_B$ is the channel dimension of each word feature $\bar{f}_w$) from $B$ using a fully-connected layer. The transformed spatial and word features are denoted as $R \in \mathbb{R}^{N \times d_l}$ and $W \in \mathbb{R}^{K \times d_l}$, respectively, where $R_i \in \mathbb{R}^{d_l}$ is the $i^{th}$ row in $R$ that corresponds to the $i^{th}$ sample point and $W_j \in \mathbb{R}^{d_l}$ is the $j^{th}$ row in $W$ that corresponds to the $j^{th}$ word in input text. Importantly, we formulate the WLST to learn the correlation between $\{R_i\}$ and $\{W_j\}$; see the next paragraph for the details. After that, $D'_s$ takes the global feature $\bar{f}_s$, sample point coordinate $p_i$, and local feature $\bar{f}_{s\_l,i}$ from WLST as inputs to predict the occupancy value at $p_i$ for shape reconstruction.

Figure 3 shows the architecture of the WLST. With the spatial features $R$ and word features $W$, we first establish an attention map $A$ to explicitly correlate each word feature $W_j$ with each sample point $p_i$ given the shape feature $\bar{f}_s$; see Figure 4 for example visualizations of $A$, revealing how it captures the spatial regions in a shape for different words in the input text. Next, we use the $softmax$ function to process $A$ to generate the normalized attention matrix $a$. The output local shape feature $\bar{f}_{s\_l,i}$ of point $p_i$ is the weighted aggregation of the word-level features $W_j$ across the whole input text. Hence, our WLST can be formulated as

$$\bar{f}_{s\_l,i} = \Sigma_j softmax(\frac{F_Q(W_j)F_K(R_i)}{\sqrt{d_l}})F_V(R_i), \quad (2)$$

where $F_Q$, $F_K$, and $F_V$ are fully-connected layers; see Figure 3 for the architecture of the WLST. Similarly, $D'_c$ also leverages a WLST for extracting local color feature $\bar{f}_{c\_l}$.

With the WLST, we can extend the implicit decoder $D$ to take into account the extra local feature $\bar{f}_l = \{\bar{f}_{s\_l}, \bar{f}_{c\_l}\}$ (see Figure 2), which is produced by explicitly learning the correlation between the word-level spatial descriptions and the 3D shape. Hence, we can make every single word in the input text accessible to the shape decoder and enhance the fidelity (or local details) of the generated shape.

**Cyclic consistency loss.** To reduce the semantic gap between the text and shape, we propose a cyclic consistency loss to encourage the consistency between input text $\mathbf{T}$ and output shape $\mathbf{S}$ from $D'$. To form a cycle, we first grid-sample $64 \times 64 \times 64$ points to use $D'$ to generate $\mathbf{S}$, and utilize encoder $E$ from the trained shape auto-encoder to extract features $f_{cyc}$ from $\mathbf{S}$; see Figure 2(b). Then, we define the cyclic consistency loss to operate on the semantic meaningful feature space instead of the low-level occupancy or color values, such that it can regularize the shape generation
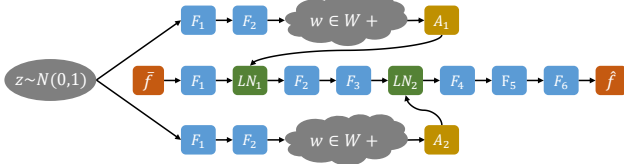
Figure 5. The architecture of our shape-IMLE generator. Inspired by StyleGAN [38], we map random noise $z$ to latent space W+ [1] to control the generator through adaptive Layer Normalization [6] ($A_1$ and $A_2$) at the first and third fully-connected layers.

in a closed loop by encouraging the high-level features $f_{cyc}$ to be similar to $f = \{f_s, f_c\}$ from the shape encoder.

To reduce the memory consumption and training time, we firstly grid-sample $16 \times 16 \times 16$ points to form a low-resolution voxelized shape $S_l$, then tri-linearly upsample $S_l$ to $\mathbf{S}$ of the same resolution as $I$ ($64 \times 64 \times 64$).

**Network training.** Initialized with the shape auto-encoder, we train the text-guided shape generation network end-to-end with the shape auto-encoder loss $L'_{ae}$ on $D'$,

$$
\begin{aligned}
L'_{ae} = &\lambda_s \Sigma_p || D'_s(f_s, p, \bar{f}_{s\_l}, R_{s\_i}) - I(p) ||_2^2 \\
&+ \lambda_c \Sigma_{k \in \{R,G,B\}} \Sigma_p || D'_c(f_c, p, \bar{f}_{c\_l})[k] - I(p)[k] ||_2^2 \mathbb{1}(I(p)),
\end{aligned}
\tag{3}
$$

$$
L_{reg} = \lambda_r || \bar{f} - f ||_2^2,
\tag{4}
$$

$$
\text{and } L_{cyc} = \lambda_{cyc} || f_{cyc} - f ||_2^2,
\tag{5}
$$

where $\lambda_s$, $\lambda_c$, $\lambda_r$, and $\lambda_{cyc}$ are weights.

### 3.4. Diversified 3D Shape Generation

To enable diversified 3D shape generation for the same input text, we propose a style-based latent shape-IMLE generator $G$, namely shape IMLE, which operates in the latent space; see Figure 2(c). Taking text feature $\bar{f} = \bar{f}_s \oplus \bar{f}_c$ from **B** as input, $G$ generates $\{\hat{f}_i = \hat{f}_{s,i} \oplus \hat{f}_{c,i}\}_{i=1}^m$ conditioned on a set of random vectors $Z = \{z_i\}_{i=1}^m$. Different from GANs, which encourage the generated samples to be similar to the target data, IMLE inversely encourages each target data to have a similar generated sample to avoid mode collapse [44]. Figure 5 shows the architecture of the shape IMLE $G$.

For the training of $G$, it is optimized as follows:

$$
\min_{\theta} \mathbb{E}_Z \left[ \min_{k \in \{1,...,m\}} d(G_\theta(\bar{f}, z_k), f) \right]
\tag{6}
$$

where $\theta$ denotes the weights of generator $G$; $d(\cdot, \cdot)$ is a distance metric; and $z_k \sim N(0, 1)$.

With each input $\bar{f}$, we randomly sample $m$ random noise vectors $\{z_i\}$ to generate $m$ different outputs $\{\hat{f}_i\}$. Among them, the one that is most similar to the ground truth $f$, say $\hat{f}_k$, is trained to be closer to $f$ with an $L_2$ regression. So, we can encourage every ground truth $f$ to have a similar generated sample to avoid the mode-collapse issue in GANs,

while promoting diversified shape generation [43, 44]. We train the shape IMLE $G$ with all the other modules $E, B, D'$ frozen (see Figure 2) using an $L_2$ loss on $\hat{f}_k = G(\bar{f}, z_k)$:

$$
L_G = \min_{k \in \{1,...,m\}} || G(\bar{f}, z_k), f ||_2^2.
\tag{7}
$$

During the inference, we feed every feature of $\{\hat{f}_1, \ldots, \hat{f}_m\}$ into $D$ for generating diversified shapes, without using the ground truth $f$ to select the nearest $\hat{f}_k$.

### 3.5. Text-Guided Shape Manipulation

Next, we extend our framework for text-guided shape manipulation, *i.e.*, to generate shape $\dot{S}$ that matches text $\mathbf{T_2}$ that is slightly modified from original text $\mathbf{T_1}$ by replacing/inserting/removing one or a few words, with other attributes unchanged for the same random noise $z$.

Taking shape manipulation (with color unchanged) as an example, we may directly feed feature $\hat{f}_2 = \{\hat{f}_{2s}, \hat{f}_{2c}\}$ from the edited text to $D'$ to generate the new edited shape. Yet, it could cause drastic changes in the unedited region and colors (Figure 7(b)). Considering the decoupled shape and color features, we may mix $\hat{f}_{2s}$ from edited text and $\hat{f}_{1c}$ from original text as input to $D'$. This simple approach ensures the consistency of the unedited attributes but shape and color may not well align with the edited shape (Figure 7(c)), since $\hat{f}_{2s}$ and $\hat{f}_{1c}$ actually come from different texts.

To encourage shape-color alignment, we propose to feed shape feature $\bar{f}_{2s}$ (extracted from text $\mathbf{T_2}$) and color feature $\bar{f}_{1c}$ (extracted from text $\mathbf{T_1}$) to $G_3$ to predict the manipulated feature $\hat{f_{2s}}, \hat{f}_{1c}$. Then, we can feed $\hat{f_{2s}}, \hat{f}_{1c}$ to $D'$ to produce the edited shape $\dot{S}$. Yet, this approach could still lead to certain changes in the unedited attributes (Figure 7(d)). Figure 6 shows our full framework further with the two-way cyclic loss, *i.e.*, $L_{cyc\_c}$ and $L_{cyc\_s}$. Here, we use shape encoder $E$ to extract manipulated feature $\dot{f} = \{\dot{f}_s, \dot{f}_c\}$ from $\dot{S}$ and formulate $L_{cyc\_s}$ for shape consistency ($\dot{f}_s$ and $\hat{f}_{2s}$) and $L_{cyc\_c}$ for color consistency ($\dot{f}_c$ and $\hat{f}_{1c}$). Then, we can formulate the overall loss:

$$
\begin{aligned}
L_{mani} = &(|| \dot{f}_s - \hat{f}_{2s} ||_2^2 + || \dot{f}_c - \hat{f}_{1c} ||_2^2) \mathbb{1}(\text{IoU}(I_1, I_2) > t) \\
&+ L_{G_1} + L_{G_2},
\end{aligned}
\tag{8}
$$

where the first term is the two-way cyclic consistency loss, which takes effect only when the Intersection over Union (IoU) between the associated ground-truth shapes $I_1$ and $I_2$ is larger than threshold $t$. The last two terms fine-tune the shape IMLE for a diversified generation (see Eq. (7)).

To train the framework, we initialize its weights from shape IMLE then finetune $G$ using $L_{mani}$ with all other modules $E, B, D'$ frozen. Also, we randomly sample two un-paired texts $\mathbf{T_1}, \mathbf{T_2}$ to simulate the original and edited texts. With the two-way cyclic loss, the shape IMLE can learn to generate edited shapes with other attributes unchanged,
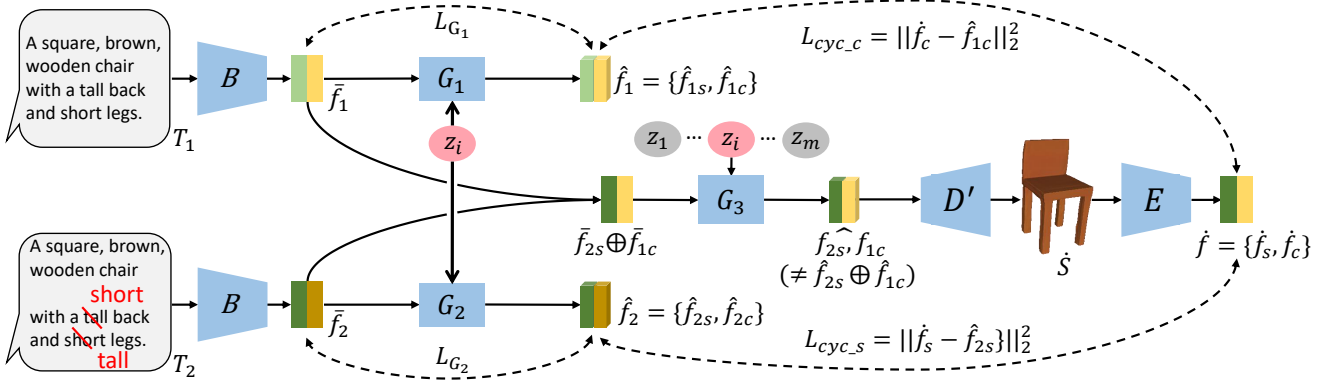
Figure 6. Overview of our text-guided shape manipulation framework (with color unchanged). Given two pieces of text $\mathbf{T}_1, \mathbf{T}_2$, shape IMLE $G_1$ and $G_2$ use the same random noise $z_i$ to generate shapes. $G_3$ takes $\{\bar{f}_{2s}, \bar{f}_{1c}\}$ and $z_i$ as input to generate shape $\dot{S}$ with feature $\{\dot{f}_s, \dot{f}_c\}$ (encoded by $E$), such that $\dot{f}_s$ and $\dot{f}_c$ should be similar to $\hat{f}_{2s}$ and $\hat{f}_{1c}$, respectively. Hence, we propose a two-way cyclic loss ($L_{cyc\_c}$ and $L_{cyc\_s}$) to encourage shape consistency between $\dot{S}$ and $\mathbf{T}_2$, and color consistency between $\dot{S}$ and $\mathbf{T}_1$. $G_1, G_2, G_3$ share the same weights.
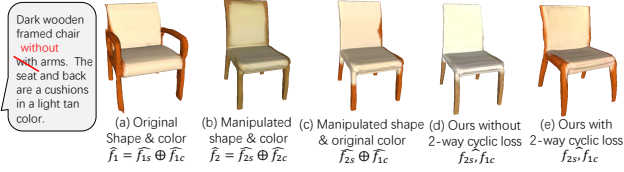


Figure 7. (a) The original shape from the unedited text. (b) The shape from the edited text. It shows that even editing just a color-unrelated word may influence the generated color. (c) Replacing $\hat{f}_{2c}$ with the original color feature $\hat{f}_{1c}$ can cause misalignment between generated shape and color. (d) Our approach without the two-way cyclic loss, the unedited attributes may still change. (e) Our full approach with the two-way cyclic loss produces an edited shape that better preserves the unedited attributes.



Figure 8. Our text-guided shape and color manipulation results.

while better aligning the shape and color. Please see the supplementary material for the details on the color manipulation framework. Besides Figures 1(b,c) and 7(a,e), Figure 8 shows two more text-guided manipulation results.

# 4. Experiments

## 4.1. Dataset and Implementation Details

Our approach is evaluated on the largest text-shape dataset *ShapeNet 3D models with natural language descriptions* [11]. The dataset contains $15,038$ shapes from the table and chair classes of ShapeNet [9]; $75,344$ natural language descriptions, 16.3 words per description on average, and $8,147$ unique words in the whole dataset [11].

We implement our framework in PyTorch [56]. To train

the shape auto-encoder, we sample $4,096$ points with the strategy in [14] and train the network for 500 epochs in $16^3$ resolution, then continue the training for another 500 epochs in $32^3$ resolution with learning rate $1e^{-4}$. For text-guided shape generation, we train the network end-to-end for 200 epochs, then fine-tune it end-to-end in $64^3$ resolution for another 200 epochs. For diversified shape generation, we train the shape IMLE for 100 epochs with learning rate $1e^{-3}$ and the other network modules frozen. Lastly, we fine-tune the shape IMLE for another 100 epochs with the two-way cyclic consistency loss to enable manipulation. We set hyperparameters $d, d_l, \lambda_s, \lambda_c, \lambda_{reg}, \lambda_{cyc}$, and $t$ as 256, 32, 2, 1, 1, 0.005, and 0.01, respectively, using a small validation set.

## 4.2. Comparison with the Existing Works

We compare our method with two existing works [11,34] (see also Section 2) on text-guided shape generation.

For a fair comparison with [11], we transform our generated results into voxels in the same resolution as [11], *i.e.*, $32^3$. Also, we follow its train/val/test ($80\%/10\%/10\%$) split and its evaluation metrics, *i.e.*, IoU, EMD, IS, and Acc (Err=1-Acc), and directly compare our results with the numbers in [11]. Table 1 reports the results, showing that our method outperforms [11] for all evaluation metrics, manifesting its effectiveness. Note that "IS" ranges $[0, 2]$, as it is built upon a two-category classification model, so both methods (1.96 *vs.* 1.97) already achieve satisfying performance in this respective. The qualitative comparisons in Figure 9 also demonstrate the superiority of our approach, which is able to generate much better shapes and colors (see Figure 9 (b, c)) in comparison with [11] (see Figure 9 (a)).

The other work [34] focuses on generating shapes from phrase descriptions; see the left side of Figure 10 (a). Since its setting is very different from ours, we only compare with it qualitatively. To do so, we first prepare sentence

Table 1. Quantitative comparisons with the existing work [11].

| Method | IoU (↑) | IS (↑) | EMD (↓) | Err (↓) |
|---|---|---|---|---|
| Text2Shape [11] | 9.64 | 1.96 | 0.4443 | 2.63 |
| Ours | **12.21** | **1.97** | **0.2071** | **2.52** |



Figure 9. Results by Text2shape [11] (a) *vs*. ours (b,c) *vs*. GT (d).
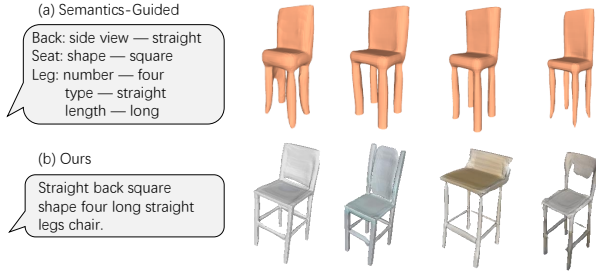


Figure 10. Results generated by [34] (a) *vs*. ours (b).

descriptions that match the phrase descriptions in [34] and then use our model to generate 3D results. Comparing the results shown in Figures 10 (a) and (b), we can see that our model is able to generate more diverse chairs that match the input description ("square shape, long straight leg"), while having varying colors and higher fidelity; please see also the supplementary material for more comparison results.

### 4.3. Ablation Studies

We conduct extensive ablation studies to validate the effectiveness of the key components in "text-guided shape generation" and "diversified generation." To measure the diversity and quality of the generated shapes, we formulate two new metrics, PS and FPD, based on Inception Score (IS) [66] and Fréchet Inception Distance (FID) [29]; please see the supplementary material for the details. To evaluate the text-shape consistency, we adopt R-Precision [81]. To reduce the training time, we train all models in $32^3$ resolution.

**Text-guided shape generation.** We evaluate the effectiveness of the following major components in this module (Figure 2 (a,b)): joint training with a pre-trained auto-encoder (AE), decoupled shape-color decoder (DSCD), WLST module (WLST), and cyclic loss (CL). Please refer to the supple-

Table 2. Ablation studies on text-guided shape generation.

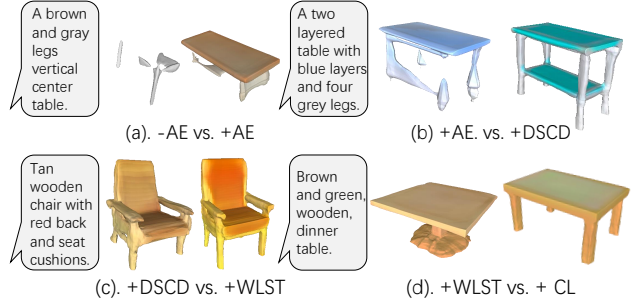| Method | IoU (↑) | PS (↑) | FPD (↓) |
|---|---|---|---|
| Without AE | 0.03 | 1.01±0.00 | 67.37 |
| +AE | 12.04 | 2.95±0.03 | 35.05 |
| further +DSCD | 12.00 | 3.16±0.04 | 31.09 |
| further +WLST | 12.24 | 3.21±0.05 | **30.34** |
| further +CL (full) | **12.33** | **3.26±0.06** | 30.80 |



Figure 11. Qualitative ablation studies on shape generation.

Table 3. Ablation studies on diversified shape generation.

| Method | PS (↑) | FPD (↓) | R-Precision (↑) |
|---|---|---|---|
| Latent GAN | 3.31 ± 0.02 | 30.70 | 21.20 ± 0.11 |
| FC IMLE | 2.93 ± 0.02 | 29.53 | 25.97 ± 0.09 |
| Shape IMLE | 3.39 ± 0.02 | 29.65 | 27.60 ± 0.39 |
| further +WLST | 3.39 ± 0.03 | 28.41 | 34.37 ± 0.09 |
| +WLST+CL (full) | **3.45 ± 0.02** | **27.26** | **40.71 ± 0.10** |

mentary material for the details of each setup.

Quantitative and qualitative results are shown in Table 2 and Figure 11, respectively. Note that all models in this setting achieve satisfying R-Precision ($> 98\%$ except "Without AE"), so we report R-Precision only in the next "Diversified generation" setting. First, auto-encoder joint training (AE) is crucial for model convergence. Without AE, the baseline approach fails to converge, leading to unreasonable results (see Figure 11 (a)) of very low quality. Second, decoupling shape and color in the decoder structure (DSCD) improves both PS and FPD by a large margin, manifesting its effectiveness in promoting high-fidelity and diversified synthesis. This is also verified in the qualitative comparison shown in Figure 11 (b). Third, empowered by the word-level correlation, we can enrich the local details; see "red back and seat cushion" in Figure 11 (c). Lastly, cyclic loss (CL) improves the consistency between the generated shape and input text; see the visual comparison in Figure 11 (d). Note that both WLST and CL benefit the model more in the "Diversified generation" component to be detailed below.

**Diversified generation.** Next, we evaluate the major modules for diversified generation (Figure 2 (c)). First, we re-
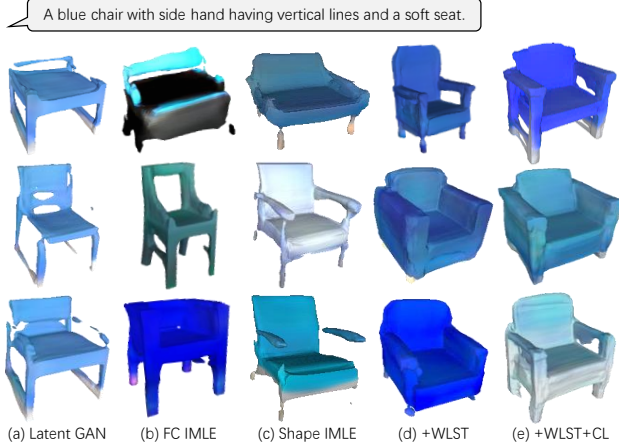
Figure 12. Qualitative ablation studies on diversified generation.

place the style-based shape IMLE with two different components for shape generation: a Latent GAN and a fully-connected IMLE (FC IMLE). Besides, we explore models without the proposed WLST module and cyclic loss (CL), which benefit diversified shape generation. Please refer to the supplementary material for the details of each setup.

Quantitative and qualitative results are shown in Table 3 and Figure 12, respectively. Note that this setting focuses on shape diversity and quality, so we do not adopt "IoU," which measures the shape similarity to the ground truths. In comparison with "Latent GAN," the IMLE model can synthesize diversified colors and avoid generating collapsed invalid shapes (see Figure 12 (a) *vs.* (b)), while attaining better quantitative results. Further, the proposed style-based generator ("shape IMLE") consistently improves on all metrics and yields higher quality shapes with better completeness in comparison with "FC IMLE," as shown in Figure 12 (c). Lastly, the WLST module and cyclic loss further help improves the generation fidelity and text-shape consistency by a large margin as shown in the last two rows of Table 3, manifesting their effectiveness (see Figure 12 (d,e)).

### 4.4. Text-Guided Shape and Color Manipulation

More text-guided manipulation results are shown in Figures 13 and 14, in addition to Figures 1(b,c), 7(a,e), and 8. Thanks to our two-way cyclic loss, our model enables text-guided modification of colors and shapes in the generated results, while trying to keep the other attributes intact. For instance, we are able to modify a "square" table to become "circular," while keeping the other irrelevant regions unchanged, *e.g.*, the legs of the chair; see Figure 13 (a). If we change the word "pink" to "blue," only the associated parts in the shape are changed accordingly; see Figure 14 (a). More comparisons with the existing work [11] and further ablation study on manipulation can be found in the supplementary material.
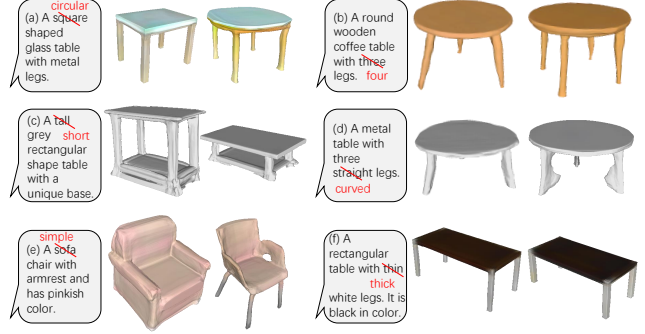


Figure 13. Our text-guided shape manipulation results. We can manipulate (a) the shape of a table, (b) number of legs, (c) height, (d) shape of legs, (e) structure, (f) thickness of legs, and so on.



Figure 14. Our text-guided color manipulation results. We can (c) manipulate the color indirectly using the related word (*e.g.*, "blood" for red), (d) manipulate the material, (e,f) adjust the color brightness using words such as "bright" and "dark", *etc*.

## 5. Conclusion

We have presented a novel framework capable of generating diversified 3D shapes with colors from text descriptions, while allowing flexible text-guided manipulations. Besides the framework, we propose to decouple the shape and color predictions for learning both shape and color features from texts and design the word-level spatial transformer to explicitly correlate words with spatial locations to enhance the local details. Also, we develop the cyclic consistency loss to enhance the text-shape consistency and introduce the style-based shape-IMLE generator for diversifying the generated shapes. Further, we extend the framework for text-guided shape manipulation with the novel two-way cyclic loss. Extensive experimental studies manifest the effectiveness of our framework. Limitation analysis and future works are elaborated in the supplementary material.

# Supplementary Material

## A. Evaluation Metrics

This section introduces the evaluation metrics employed in the experiments. Below, we first introduce the metrics we formulated/extended from the existing ones for the evaluations, and then introduce other metrics from Text2Shape [11].

- **PS** and **FPD**

  Point Score (PS) and Fréchet Point Distance (FPD) measure shape diversity and quality.

  Existing works [46, 67] often utilize the Fréchet Point Distance (FPD) to evaluate the quality of the generated 3D shapes. However, such metric cannot account for color, which is one of the important characteristics in our results that is not present in the previous works. To jointly evaluate the shape and color in the generated results, we formulate PS and extend FPD for shape and color evaluations based on Inception Score (IS) [71] and Fréchet Inception Distance (FID) [29].

  PS measures the KL-divergence between the conditional probabilities of the generated shapes and their marginal probabilities. On the other hand, FPD measures the Wasserstein distance between the distribution of the generated shapes and that of the real samples. The mentioned probabilities are inferred from a pre-trained classification network (*e.g.*, Inception v3 [71] on ImageNet for image generation). In our case, PS and FPD are built upon a newly-trained PointNet [61], since no existing 3D classification network can simultaneously consider both shape and color as far as we know. Specifically, we train a classification-based PointNet on ScanObjectNN [74] for 200 epochs, with a validation classification accuracy of 84.85%.

- **IoU**

  Intersection over Union (IoU) measures the similarity to the ground truth. We evaluate the IoU between the generated shape and ground truth, by measuring the similarity of the occupancy between them.

- **R-precision**

  We adopt **R-precision** [81] to measure the consistency between the generated shape **S** and input text **T**. Specifically, we extract shape and text features using $E$ and $B$, respectively, then evaluate R-Precision three times with different random seeds to reduce the randomness.

- Metrics in Text2Shape [11]

  Chen *et al*. [11] adopt four metrics, including **IoU**, **EMD**, **IS**, and **Acc** (**Err**=1-Acc), for evaluating their

results. IoU and EMD measure the shape and color similarity between the generated shape and ground truth, respectively. IS measures the diversity and quality of the generated shapes, and Err (Acc) measures the quality. Please refer to [11] for the details. For a pair comparison, we train a classification model using the official code released by the author of [11] to evaluate our IS and Acc.

## B. Details of the Baselines

### B.1. Text-Guided Shape Generation

We create the following baselines to evaluate the key modules of our text-guided shape generation framework, including the auto-encoder (AE), decoupled shape-color decoder (DSCD), WLST module (WLST), and cyclic loss (CL).

(i) "Without AE." In this setting, the network is composed of text encoder $B$ and decoder $D$ but without shape encoder $E$. It is optimized to regress the the occupancy and color values of the target shape $I$ with an $L_2$ loss. It serves as our preliminary baseline for text-guided shape generation, where $B$ maps the text **T** to a latent space and $D$ reconstructs the shape and color.

(ii) "+AE." The network is composed of the auto-encoder $E$, $D_{share}$, and text encoder $B$. $E$ encodes the input shape $I$ as a joint shape-color feature $f_{share}$ and $D$ reconstructs the shape and color of $I$. $B$ extracts the text feature $\bar{f}_{share}$ and minimizes the mean squared difference between $\bar{f}_{share}$ and $f_{share}$. It serves as a baseline to directly adopt the auto-encoder-based approach [14] to our task, as introduced in Section 1 of the main paper, so this baseline demonstrates the necessity of the auto-encoder in our approach.

(iii) "Further +DSCD." In this setting, we replace the shared decoder $D_{share}$ with a pair of decoupled shape-color decoders $D = \{D_s, D_c\}$ and replace shared features $f_{share}$ and $\hat{f}_{share}$ in (ii) with a pair of decoupled features $f = \{f_s, f_c\}$ and $\hat{f} = \{\hat{f}_s, \hat{f}_c\}$, respectively. This baseline manifests the effectiveness of the decoupled shape-color decoder (DSCD) in our framework.

(iv) "Further +WLST." Based on (iii), the spatial-aware decoder $D'$ with the WLST is adopted for text-guided shape generation in place of $D$. This baseline manifests the effectiveness of WLST.

(v) "Further +CL" (our full model). Model (iv) is trained with an additional cyclic loss (which is Eq.(5) in the main paper), whereas Model (v) is our full model. Comparing between model (iv) and model (v) verifies the applicability of the cyclic loss.
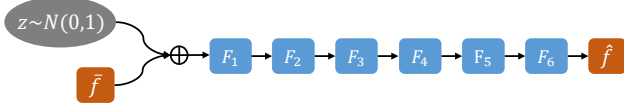
Figure 15. The FC generator architecture in FC IMLE.

## B.2. Diversified Generation

To evaluate the core modules for diversified shape generation, We compare our Shape IMLE with two other approaches: (i) Latent GAN and (ii) fully-connected IMLE (FC IMLE). Besides, we evaluate the performance gain of our proposed WLST and cyclic loss (CL) for diversified shape generation. For each baseline, we generate three different samples for each text with random noises $z_1$ to $z_3$.

(i) "Latent-GAN." We adopt Latent-GAN [2, 4] conditioned on the input text to generate diversified results. We adopt our style-based latent shape-IMLE generator $G$ (Figure 5 in the main paper) as generator and a small network with three fully-connected layers as discriminator $D_{\text{latent}}$. We train the generator and discriminator iteratively using adversarial training with all the other modules frozen.

(ii) "FC IMLE." In this setting, we introduce the IMLE framework for diversified generation. As shown in Figure 15, the IMLE generator $G_{simple}$ is composed of six fully-connected layers that take $\bar{f} \oplus z$ as input. This baseline aims to show the superiority of IMLE.

(iii) "Shape IMLE." In place of FC-IMLE in (ii), we adopt the style-based shape-IMLE generator $G$ shown in Figure 5 in the main paper. This baseline manifests the effectiveness of our shape-IMLE generator $G$.

(iv) "+WLST" and "+CL." Similar to "+WLST" and "+CL" in Section B.1, we again evaluate their capability on improving the diversified generation.

## C. Results of Text-Guided Shape Generation

### C.1. Comparison with Existing Works

In this section, we show more results on comparing our method with [11] and [34]. As shown in Figure 16, our approach is able to generate shapes with higher fidelity compared with [11]. Also, our results are more consistent to the input texts. As shown in Figure 16 (e) on the bottom left of the figure, our approach is able to create a folding chair following the text description, where [11] can only output a regular chair.

The most recent work [34] takes only pre-defined semantic labels as inputs, unlike our approach, which can take
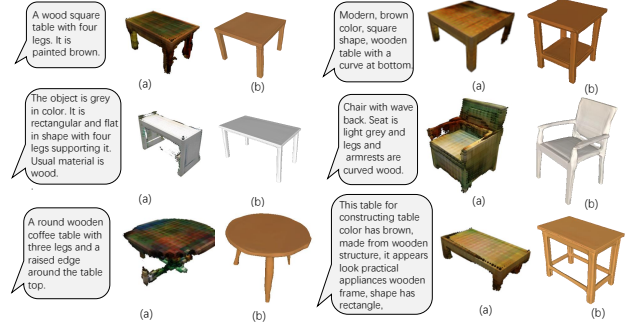


Figure 16. Additional text-guided generation results compared with Text2Shape [11].
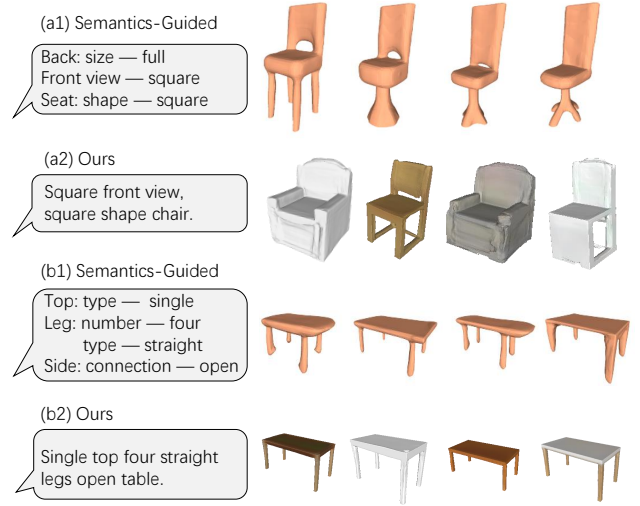


Figure 17. Additional text-guided generation results compared with [34].

natural language as inputs. As shown in Figure 17, our approach can generate more diversified shapes that better match the input text description ("square shape, square view" in Figure 17 (a1, a2)), compared with [34].

### C.2. Additional Generation Results

Further, we show more text-guided shape generation results in Figures 18. These results again manifest the superiority of our approach on diversity, fidelity, and text-shape consistency, demonstrating the capability of our method over the previous ones.

## D. Text-Guided Shape Manipulation

### D.1. Color Manipulation Framework

In this section, we introduce our framework for text-guided color manipulation with shape unchanged. As shown in Figure 19, we feed shape feature $\bar{f}_{1s}$ (extracted from $\mathbf{T}_1$) and color feature $\bar{f}_{2c}$ (extracted from $\mathbf{T}_2$) to $G_3$ to predict

Figure 18. Our text-guided shape generation results.

the manipulated feature $\hat{f_{1s}}, \hat{f}_{2c}$, then feed it to $D'$ to produce the edited shape $\dot{S}$. We then extract manipulated feature $\dot{f} = \{\dot{f}_s, \dot{f}_c\}$ from $\dot{S}$ using $E$ and use the two-way cyclic loss, *i.e.*, $L_{cyc\_s}$ to encourage shape consistency ($\dot{f}_s$ and $\hat{f}_{1s}$) and $L_{cyc\_c}$ to encourage color consistency ($\dot{f}_c$ and $\hat{f}_{2c}$).

Similar to the shape manipulation framework shown in the main paper, we train our color manipulation framework using the following loss:

$$
\begin{aligned}
L_{mani}^{c} = (||\dot{f}_s - \hat{f}_{1s}||_2^2 + ||\dot{f}_c - \hat{f}_{2c}||_2^2)\mathbb{1}(\text{IoU}(I_1, I_2) > t) \\
+ L_{G_1} + L_{G_2},
\end{aligned}
\tag{9}
$$

where the terms have the same definition as Eq.(8) in the main paper.

## D.2. Comparison with the Existing Work

In this section, we compare our method with [11] on shape manipulation capability. As shown in Figure 20, inserting or editing words related to the color attribute leads to undesirable changes in the other attributes, as shown in the results produced by [11], *e.g.*, the shape of the chair back and the table leg, whereas our approach is able to better preserve the shapes (geometries and structures).

## D.3. Ablation Studies

In this section, we evaluate the different strategies for text-guided manipulation quantitatively. To measure the quality

of the manipulated shapes, We adopt PS and FPD as the evaluate metrics. To further evaluate the consistency before and after the manipulation, we calculate R-Precision$_1$ based on $\dot{f}$ (feature from the manipulated shape) and $\hat{f}_1$ (feature from the original text), and assess R-Precision$_2$ based on $\dot{f}$ (feature from the manipulated shape) and $E(D'(\hat{f}_1))$ (feature from the generated shape by the original text). We build a small dataset containing 50 pairs of original and manipulated texts for the evaluation.

(i) Baseline 1. As shown in Figure 7(b) of the main paper, we directly feed the feature from the edited text $\hat{f}_2 = \{\hat{f}_{2s}, \hat{f}_{2c}\}$ to our generation framework. It is the primitive baseline for manipulation because it adopts no mechanism for consistency preserving, but it serves as the upper bound of the shape manipulation quality, since it directly adopts our generation framework (Figure 2 in the main paper) to produce the result without any constraints on the manipulation consistency.

(ii) Baseline 2. As shown in Figure 7(c) of the main paper, the shape is generated by a mixture of $\hat{f}_1$ and $\hat{f}_2$. Specifically, for the shape manipulation, we feed $\hat{f}_{2s} \oplus \hat{f}_{1c}$ to $D'$; and for the color manipulation, we use $\hat{f}_{1s} \oplus \hat{f}_{2c}$.

(iii) Baseline 3. As shown in Figure 7(d) of the main paper, we feed a mixture of $\bar{f}_1$ and $\bar{f}_2$ to $G$ to boost the shape-color alignment. For the shape manipulation, we feed $\bar{f}_{2s} \oplus \bar{f}_{1c}$ to $G$ to derive $\hat{f}_{2s}, \hat{f}_{1c}$; and for the color manipulation, we predict $\hat{f}_{1s}, \hat{f}_{2c}$ from $\hat{f}_{1s} \oplus \hat{f}_{2c}$ with $G$.

(iv) Our full model (Ours). Built upon (iii), we further incorporate the two-way cyclic loss shown in Eq.(8) of the main paper for shape manipulation, and Eq. (9) in this supplementary document for color manipulation.

"Baseline 1" generates a new shape using the edited text without considering what the original shape is. As shown in Table 4, despite of the best diversity and quality it achieves, the lowest R-Precisions indicate the unsatisfying consistency before and after the manipulation (see Figure 7 (b) in the main paper).

On the other hand, "Baseline 2" and "Baseline 3" attain better consistency at the expense of the generation quality, and our manipulation framework with the two-way cyclic loss is able to achieve the best consistency before and after the manipulation, while having better generation quality compared with both "Baseline 2" and "Baseline 3," even being close to "Baseline 1."

## E. Alternative Training Strategy

In this section, we discuss an optional training strategy. Specifically, we jointly train the shape auto-encoder and text encoder $E, D'$ and $B$ end-to-end, instead of following
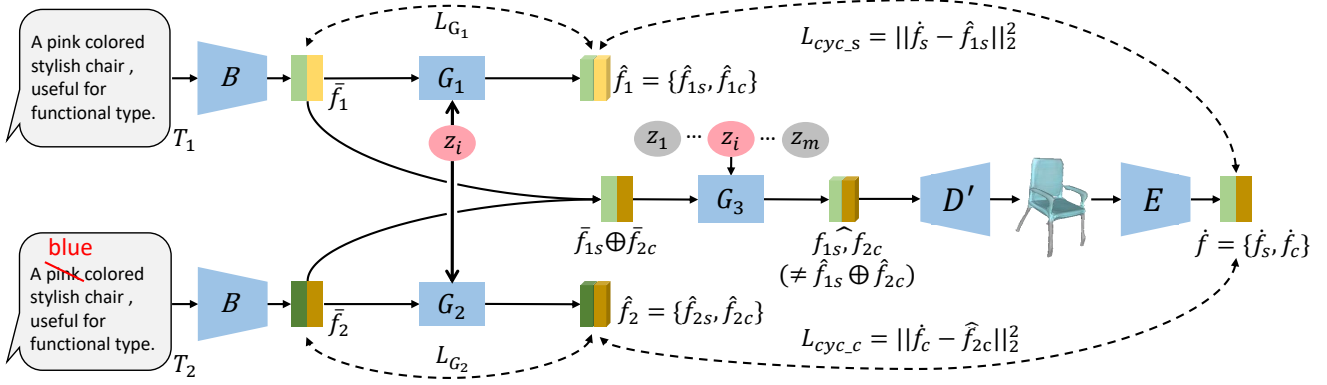
Figure 19. Overview of our text-guided color manipulation framework (with shape unchanged). Given two pieces of text $\mathbf{T}_1, \mathbf{T}_2$, shape IMLE $G_1$ and $G_2$ use the same random noise $z_i$ for shape generation. $G_3$ takes $\{\bar{f}_{1s}, \bar{f}_{2c}\}$ and $z_i$ as input to generate shape $\dot{\mathbf{S}}$ with feature $\{\dot{f}_s, \dot{f}_c\}$ (encoded by $E$), such that $\dot{f}_s$ and $\dot{f}_c$ should be similar to $\hat{f}_{1s}$ and $\hat{f}_{2c}$, respectively. To this end, we propose a two-way cyclic loss to encourage the shape consistency between $\dot{\mathbf{S}}$ and $\mathbf{T}_1$, and the color consistency between $\dot{\mathbf{S}}$ and $\mathbf{T}_2$. $G_1, G_2, G_3$ share the same weights.
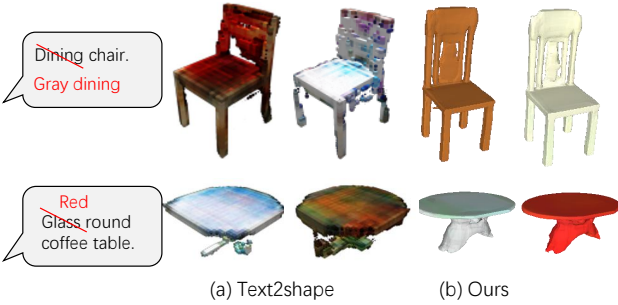


Figure 20. Text-guided manipulation results on comparing our method with Text2Shape [11].

Table 4. Ablation studies on text-guided manipulation.

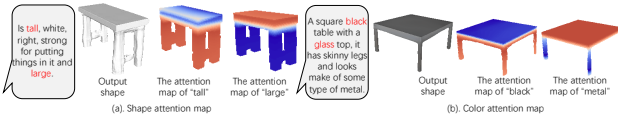| Method | PS ($\uparrow$) | FPD ($\downarrow$) | R-Precision$_1$ ($\uparrow$) | R-Precision$_2$ ($\uparrow$) |
|---|---|---|---|---|
| Baseline 1 | **2.80 ± 3.03** | **31.70** | 36.00 ± 2.83 | 48.67 ± 0.94 |
| Baseline 2 | 2.73 ± 0.39 | 33.77 | 43.33 ± 0.94 | 52.00 ± 3.26 |
| Baseline 3 | 2.75 ± 0.48 | 35.74 | 42.66 ± 3.77 | 56.67 ± 2.49 |
| Ours | 2.76 ± 0.53 | 32.03 | **58.00 ± 2.82** | **67.33 ± 1.89** |



Figure 21. Attention map of the end-to-end training strategy.

the training strategy presented in the main paper that first trains $E, D$, and then jointly trains $E, D', B$. This strategy includes fewer training steps, but needs much more training time because $B$ continuously optimized in the whole training process. This training strategy achieves comparable performance as presented in the main paper, and can generate attention maps that are more consistent with the semantic meaningful shape parts as shown in Figure 21.
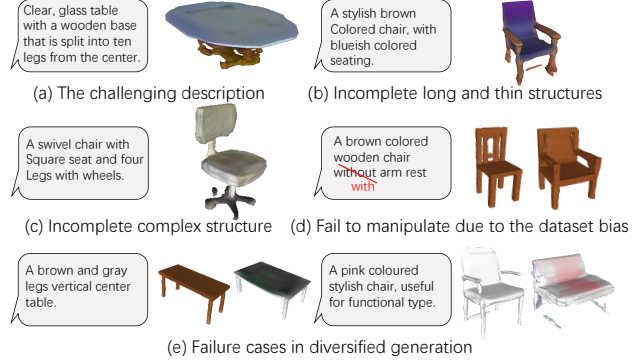


Figure 22. Illustrating the limitations of our current approach.

## F. Limitations and Future Work

Our current approach still has some limitations. First, text-guided shape generation is a very challenging task as discussed in Section 2 of the main paper. For example, some attributes, such as "ten legs" in Figure 22 (a), are extremely challenging to generate. Our framework fails to generate a shape that faithfully follows such description. Second, shapes with long, thin, and fine structures may get distorted or become noisy, as shown in Figure 22 (b, c). To address this issue, we plan to explore more recent 3D implicit representations [15, 24, 45]. Also, the manipulation performance is limited by the inherent bias of the dataset. If we add an armrest to the chair in Figure 22 (d), the manipulated chair will have become wider than the original one. Such a result is partially due to the dataset bias, since armed chairs are typically wider (like sofa) than those without armrests in the dataset. To resolve it, the manipulation process requires a topological understanding of the shapes. In addition, our metrics have some limitations. On the one hand, IoU may not be a good metric for text-to-shape generation task, because a slight difference in height/position between the generated shapes and GT

shape can cause a low IoU; particularly, this is beyond the representative ability of a small piece of text. On the other hand, PS and FPD cannot fully reflect the generative quality, because these two metrics are based on the ScanObjectNN dataset [74], which has a large domain gap from our training dataset ShapeNet. In other words, better PS and FPD simply indicate that the generated shapes are more similar to the ScanObjectNN shapes, not necessarily meaning better quality. Also, there is a trade-off between the diversity and fidelity in our diversified generation. When stronger noise added to encourage the diversity, the quality of some generations cannot be ensured, and some are inconsistent with the text description as shown in Figure 22 (e). It gets more serious when the text is long and contains descriptions on shape details. Last, our approach needs paired text-shape data for training, so we temporarily only explore shapes of table and chair, since the largest existing dataset [11] only provides samples of these two categories. In the future, we will plan to explore zero-shot text-guided shape generation to extend the applicability of this work.

# References

[1] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2StyleGAN: How to embed images into the StyleGAN latent space? In *ICCV*, 2019.

[2] Panos Achlioptas, Olga Diamanti, Ioannis Mitliagkas, and Leonidas J. Guibas. Learning representations and generative models for 3D point clouds. In *ICML*, 2018.

[3] Panos Achlioptas, Judy E. Fan, Robert X.D. Hawkins, Noah D. Goodman, and Leonidas J. Guibas. Learning to refer to 3D objects with natural language. 2018.

[4] Martin Arjovsky and Léon Bottou. Towards principled methods for training generative adversarial networks. *ICLR*, 2017.

[5] Himanshu Arora, Saurabh Mishra, Shichong Peng, Ke Li, and Ali Mahdavi-Amiri. Shape completion via IMLE. *arXiv preprint arXiv:2106.16237*, 2021.

[6] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.

[7] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale GAN training for high fidelity natural image synthesis. *ICLR*, 2019.

[8] Ruojin Cai, Guandao Yang, Hadar Averbuch-Elor, Zekun Hao, Serge Belongie, Noah Snavely, and Bharath Hariharan. Learning gradient fields for shape generation. In *ECCV*, 2020.

[9] Angel X. Chang, Thomas Funkhouser, Leonidas J. Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiong Xiao, Li Yi, and Fisher Yu. ShapeNet: An Information-Rich 3D Model Repository. Technical Report arXiv:1512.03012 [cs.GR], 2015.

[10] Dave Zhenyu Chen, Angel X. Chang, and Matthias Nießner. ScanRefer: 3D object localization in RGB-D scans using natural language. In *ECCV*, 2020.

[11] Kevin Chen, Christopher B. Choy, Manolis Savva, Angel X. Chang, Thomas Funkhouser, and Silvio Savarese. Text2Shape: Generating shapes from natural language by learning joint embeddings. In *ACCV*, 2018.

[12] Zhiqin Chen, Vladimir G. Kim, Matthew Fisher, Noam Aigerman, Hao Zhang, and Siddhartha Chaudhuri. DECOR-GAN: 3D shape detailization by conditional refinement. In *CVPR*, 2021.

[13] Zhiqin Chen, Andrea Tagliasacchi, and Hao Zhang. BSP-Net: Generating compact meshes via binary space partitioning. In *CVPR*, 2020.

[14] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. In *CVPR*, 2019.

[15] Zhang Chen, Yinda Zhang, Kyle Genova, Sean Fanello, Sofien Bouaziz, Christian Hane, Ruofei Du, Cem Keskin, Thomas Funkhouser, and Danhang Tang. Multiresolution deep implicit functions for 3d shape representation. In *ICCV*, 2021.

[16] Julian Chibane, Thiemo Alldieck, and Gerard Pons-Moll. Implicit functions in feature space for 3D shape reconstruction and completion. In *CVPR*, 2020.

[17] Julian Chibane, Aymen Mir, and Gerard Pons-Moll. Neural unsigned distance fields for implicit function learning. In *NeurIPS*, 2020.

[18] Julian Chibane and Gerard Pons-Moll. Implicit feature networks for texture completion from partial 3D data. In *ECCV*, 2020.

[19] Christopher B. Choy, Danfei Xu, JunYoung Gwak, Kevin Chen, and Silvio Savarese. 3D-R2N2: A unified approach for single and multi-view 3D object reconstruction. In *ECCV*, 2016.

[20] Yu Deng, Jiaolong Yang, and Xin Tong. Deformed implicit field: Modeling 3D shapes with learned dense correspondence. In *CVPR*, 2021.

[21] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. *NAACL-HLT*, 2018.

[22] Yutong Feng, Yifan Feng, Haoxuan You, Xibin Zhao, and Yue Gao. MeshNet: Mesh neural network for 3D shape representation. In *AAAI*, 2019.

[23] Lin Gao, Jie Yang, Tong Wu, Yu-Jie Yuan, Hongbo Fu, Yu-Kun Lai, and Hao Zhang. SDM-NET: Deep generative network for structured deformable mesh. *ACM TOG (SIGGRAPH Asia)*, 2019.

[24] Kyle Genova, Forrester Cole, Avneesh Sud, Aaron Sarna, and Thomas Funkhouser. Local deep implicit functions for 3D shape. In *CVPR*, 2020.

[25] Rohit Girdhar, David F. Fouhey, Mikel Rodriguez, and Abhinav Gupta. Learning a predictable and generative vector representation for objects. In *ECCV*, 2016.

[26] Zhizhong Han, Chao Chen, Yu-Shen Liu, and Matthias Zwicker. ShapeCaptioner: Generative caption network for 3D shapes by learning a mapping from parts detected in multiple views to sentences. In *ACM MM*, 2020.

[27] Zhizhong Han, Mingyang Shang, Xiyang Wang, Yu-Shen Liu, and Matthias Zwicker. Y2seq2seq: Cross-modal representation learning for 3D shape and text by joint reconstruction and prediction of view and word sequences. In *AAAI*, 2019.

[28] Zekun Hao, Hadar Averbuch-Elor, Noah Snavely, and Serge Belongie. DualSDF: Semantic shape manipulation using a two-level representation. In *CVPR*, 2020.

[29] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. GANs trained by a two time-scale update rule converge to a local Nash equilibrium. *NIPS*, 2017.

[30] Jiahui Huang, Shi-Sheng Huang, Haoxuan Song, and Shi-Min Hu. DI-Fusion: Online implicit 3D reconstruction with deep priors. In *CVPR*, 2021.

[31] Le Hui, Rui Xu, Jin Xie, Jianjun Qian, and Jian Yang. Progressive point cloud deconvolution generation network. In *ECCV*, 2020.

[32] Faria Huq, Nafees Ahmed, and Anindya Iqbal. Static and animated 3D scene generation from free-form text descriptions. *arXiv preprint arXiv:2010.01549*, 2020.

[33] Moritz Ibing, Isaak Lim, and Leif Kobbelt. 3D shape generation with grid-based implicit functions. In *CVPR*, 2021.

[34] Tansin Jahan, Yanran Guan, and Oliver van Kaick. Semantics-guided latent space exploration for shape generation. In *COMPUT GRAPH FORUM*, 2021.

[35] Chiyu Jiang, Jingwei Huang, Andrea Tagliasacchi, Leonidas J. Guibas, et al. ShapeFlow: Learnable deformations among 3D shapes. *NeurIPS*, 2020.

[36] Chiyu Jiang, Avneesh Sud, Ameesh Makadia, Jingwei Huang, Matthias Nießner, Thomas Funkhouser, et al. Local implicit grid representations for 3D scenes. In *CVPR*, 2020.

[37] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of GANs for improved quality, stability, and variation. *ICLR*, 2018.

[38] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *CVPR*, 2019.

[39] Hyeongju Kim, Hyeonseung Lee, Woo Hyun Kang, Joun Yeop Lee, and Nam Soo Kim. SoftFlow: Probabilistic framework for normalizing flow on manifolds. *NeurIPS*, 2020.

[40] Roman Klokov, Edmond Boyer, and Jakob Verbeek. Discrete point flow networks for efficient point cloud generation. In *ECCV*, 2020.

[41] Bowen Li, Xiaojuan Qi, Thomas Lukasiewicz, and Philip H. S. Torr. Controllable text-to-image generation. *NeurIPS*, 2019.

[42] Bowen Li, Xiaojuan Qi, Thomas Lukasiewicz, and Philip H. S. Torr. ManiGAN: Text-guided image manipulation. In *CVPR*, 2020.

[43] Ke Li, Shichong Peng, Tianhao Zhang, and Jitendra Malik. Multimodal image synthesis with conditional implicit maximum likelihood estimation. *IJCV*, 2020.

[44] Ke Li, Tianhao Zhang, and Jitendra Malik. Diverse image synthesis from semantic layouts via conditional IMLE. In *ICCV*, 2019.

[45] Manyi Li and Hao Zhang. D$^2$IM-Net: Learning detail disentangled implicit fields from single images. *CVPR*, 2021.

[46] Ruihui Li, Xianzhi Li, Ka-Hei Hui, and Chi-Wing Fu. SP-GAN: sphere-guided 3D shape generation and manipulation. *ACM TOG (SIGGRAPH)*, 2021.

[47] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: Common objects in context. In *ECCV*, 2014.

[48] Shi-Lin Liu, Hao-Xiang Guo, Hao Pan, Peng-Shuai Wang, Xin Tong, and Yang Liu. Deep implicit moving least-squares functions for 3D reconstruction. In *CVPR*, 2021.

[49] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3D reconstruction in function space. In *CVPR*, 2019.

[50] Kaichun Mo, Paul Guerrero, Li Yi, Hao Su, Peter Wonka, Niloy Mitra, and Leonidas J. Guibas. StructureNet.: Hierarchical graph networks for 3D shape generation. *ACM TOG (SIGGRAPH Asia)*, 2019.

[51] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Differentiable volumetric rendering: Learning implicit 3D representations without 3D supervision. In *CVPR*, 2020.

[52] Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In *2008 Sixth Indian Conference on Computer Vision, Graphics & Image Processing*, 2008.

[53] Michael Oechsle, Lars Mescheder, Michael Niemeyer, Thilo Strauss, and Andreas Geiger. Texture fields: Learning texture representations in function space. In *ICCV*, 2019.

[54] Michael Oechsle, Michael Niemeyer, Christian Reiser, Lars Mescheder, Thilo Strauss, and Andreas Geiger. Learning implicit surface light fields. In *3DV*, 2020.

[55] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. In *CVPR*, 2019.

[56] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. PyTorch: An imperative style, high-performance deep learning library. *NeurIPS*, 2019.

[57] Or Patashnik, Zongze Wu, Eli Shechtman, Daniel Cohen-Or, and Dani Lischinski. StyleCLIP: Text-driven manipulation of StyleGAN imagery. *ICCV*, 2021.

[58] Shichong Peng and Ke Li. Generating unobserved alternatives. *ICLR*, 2021.

[59] Songyou Peng, Michael Niemeyer, Lars Mescheder, Marc Pollefeys, and Andreas Geiger. Convolutional occupancy networks. *ECCV*, 2020.

[60] Omid Poursaeed, Matthew Fisher, Noam Aigerman, and Vladimir G. Kim. Coupling explicit and implicit surface representations for generative 3D modeling. In *ECCV*, 2020.

[61] Charles R. Qi, Hao Su, Kaichun Mo, and Leonidas J. Guibas. PointNet: Deep learning on point sets for 3D classification and segmentation. In *CVPR*, 2017.

[62] Tingting Qiao, Jing Zhang, Duanqing Xu, and Dacheng Tao. MirrorGAN: Learning text-to-image generation by redescription. In *CVPR*, 2019.

[63] Scott Reed, Zeynep Akata, Xinchen Yan, Lajanugen Logeswaran, Bernt Schiele, and Honglak Lee. Generative adversarial text to image synthesis. In *ICML*, 2016.

[64] Scott E. Reed, Zeynep Akata, Santosh Mohan, Samuel Tenka, Bernt Schiele, and Honglak Lee. Learning what and where to draw. *NIPS*, 2016.

[65] Robin Rombach, Patrick Esser, and Björn Ommer. Network-to-network translation with conditional invertible neural networks. *NeurIPS*, 2020.

[66] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training GANs. *NIPS*, 2016.

[67] Dong Wook Shu, Sung Woo Park, and Junseok Kwon. 3D point cloud generative adversarial network based on tree structured graph convolutions. In *ICCV*, 2019.

[68] Douglas M. Souza, Jônatas Wehrmann, and Duncan D. Ruiz. Efficient neural architecture for text-to-image synthesis. In *IJCNN*, 2020.

[69] David Stap, Maurits Bleeker, Sarah Ibrahimi, and Maartje ter Hoeve. Conditional image generation and manipulation for user-specified content. *CVPRW*, 2020.

[70] Yongbin Sun, Yue Wang, Ziwei Liu, Joshua Siegel, and Sanjay Sarma. PointGrow: Autoregressively learned point cloud generation with self-attention. In *WACV*, 2020.

[71] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *CVPR*, 2016.

[72] Chuan Tang, Xi Yang, Bojian Wu, Zhizhong Han, and Yi Chang. Part2Word: Learning joint embedding of point clouds and text by matching parts to words. *arXiv preprint arXiv:2107.01872*, 2021.

[73] Edgar Tretschk, Ayush Tewari, Vladislav Golyanik, Michael Zollhöfer, Carsten Stoll, and Christian Theobalt. PatchNets: Patch-based generalizable deep implicit 3D shape representations. In *ECCV*, 2020.

[74] Mikaela Angelina Uy, Quang-Hieu Pham, Binh-Son Hua, Thanh Nguyen, and Sai-Kit Yeung. Revisiting point cloud classification: A new benchmark dataset and classification model on real-world data. In *ICCV*, 2019.

[75] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The Caltech-UCSD Birds-200-2011 Dataset. Technical Report CNS-TR-2011-001, California Institute of Technology, 2011.

[76] Hao Wang, Guosheng Lin, Steven Hoi, and Chunyan Miao. Cycle-consistent inverse GAN for text-to-image synthesis. *ACM MM*, 2021.

[77] Zixu Wang, Zhe Quan, Zhi-Jie Wang, Xinjian Hu, and Yangyang Chen. Text to image synthesis with bidirectional generative adversarial network. In *ICME*, 2020.

[78] Rundi Wu, Yixin Zhuang, Kai Xu, Hao Zhang, and Baoquan Chen. PQ-NET: A generative part seq2seq network for 3D shapes. In *CVPR*, 2020.

[79] Weihao Xia, Yujiu Yang, Jing-Hao Xue, and Baoyuan Wu. TediGAN: Text-guided diverse face image generation and manipulation. In *CVPR*, 2021.

[80] Qiangeng Xu, Weiyue Wang, Duygu Ceylan, Radomir Mech, and Ulrich Neumann. DISN: Deep implicit surface network for high-quality single-view 3D reconstruction. *NeurIPS*, 2019.

[81] Tao Xu, Pengchuan Zhang, Qiuyuan Huang, Han Zhang, Zhe Gan, Xiaolei Huang, and Xiaodong He. AttnGAN: Fine-grained text to image generation with attentional generative adversarial networks. In *CVPR*, 2018.

[82] Jie Yang, Kaichun Mo, Yu-Kun Lai, Leonidas J. Guibas, and Lin Gao. DSG-Net: Learning disentangled structure and geometry for 3D shape generation. *ACM TOG (SIGGRAPH Asia)*, 2021.

[83] Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Ronen Basri, and Yaron Lipman. Multiview neural surface reconstruction by disentangling geometry and appearance. *NeurIPS*, 2020.

[84] Wang Yifan, Shihao Wu, Cengiz Oztireli, and Olga Sorkine-Hornung. Iso-Points: Optimizing neural implicit surfaces with hybrid representations. In *CVPR*, 2021.

[85] Mingkuan Yuan and Yuxin Peng. Bridge-GAN: Interpretable representation learning for text-to-image synthesis. *IEEE TCSVT*, 2019.

[86] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiaolei Huang, and Dimitris N. Metaxas. StackGAN: Text to photo-realistic image synthesis with stacked generative adversarial networks. In *ICCV*, 2017.

[87] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiaolei Huang, and Dimitris N. Metaxas. StackGAN++: Realistic image synthesis with stacked generative adversarial networks. *IEEE TPAMI*, 2018.

[88] Zerong Zheng, Tao Yu, Qionghai Dai, and Yebin Liu. Deep implicit templates for 3D shape representation. In *CVPR*, 2021.