

Sketch-A-Shape: Zero-Shot Sketch-to-3D Shape Generation

Aditya Sanghi Pradeep Kumar Jayaraman[‡] Arianna Rampini[‡] Joseph Lambourne
 Hooman Shayani Evan Atherton Saeid Asgari Taghanaki
 Autodesk Research

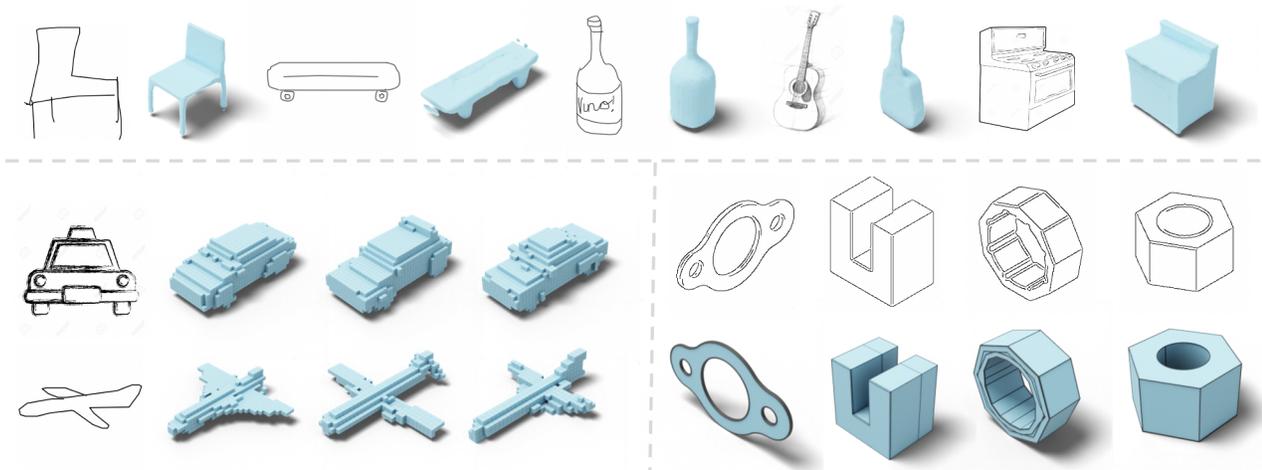


Figure 1: Sketch-A-Shape is a zero-shot sketch-to-3D generative model. Here we show how our method can generalize across voxel, implicit, and CAD representations and synthesize consistent 3D shapes from a variety of inputs ranging from casual doodles to professional sketches with different levels of *ambiguity*.

Abstract

Significant progress has recently been made in creative applications of large pre-trained models for downstream tasks in 3D vision, such as text-to-shape generation. This motivates our investigation of how these pre-trained models can be used effectively to generate 3D shapes from sketches, which has largely remained an open challenge due to the limited sketch-shape paired datasets and the varying level of abstraction in the sketches. We discover that conditioning a 3D generative model on the features (obtained from a frozen large pre-trained vision model) of synthetic renderings during training enables us to effectively generate 3D shapes from sketches at inference time. This suggests that the large pre-trained vision model features carry semantic signals that are resilient to domain shifts, *i.e.*, allowing us to use only RGB renderings, but generalizing to sketches at inference time. We conduct a comprehensive set of experiments investigating different design factors and demonstrate the effectiveness of our straightforward approach for

generation of multiple 3D shapes per each input sketch regardless of their level of abstraction without requiring any paired datasets during training.

1. Introduction

Throughout history, humans have used drawings and other visual representations to communicate complex ideas, concepts, and information.

As hand-drawn sketches have a high level of abstraction, they allow unskilled artists or even young children to convey semantic information about 3D objects [21], while providing trained professionals with a way to quickly express important geometric and stylistic details of a 3D design. The ability to create 3D models which can capture the essence of simple doodles while accurately reproducing 3D shapes described by concept design sketches, will make 3D modelling more accessible to the general public, while allowing designers to rapidly explore many different design

ideas and create virtual models that more accurately reflect the shape, size, and characteristics of real-world objects and environments.

Previous studies have endeavored to employ deep learning techniques in generating 3D shapes from sketches [39, 20, 16, 35], yet there are several limitations that hinder their widespread application. Firstly, there is a lack of (Sketch, 3D shape) paired data at a large scale which forces most methods to be trained on synthetic datasets or data collected on only few categories. Even when a small number of categories of paired sketch-shape data has been collected [35], current methods fail to generalize to different levels of abstractions in the sketches, ranging from casual doodles to detailed professional drawings. Finally, most of the present methods incorporate strong inductive biases, such as view information [79], differentiable rendering [20] and depth estimation [39, 16], thereby constraining their generalizability across 3D representations.

To overcome the challenge of limited availability of paired data, a potential solution is to use prior knowledge encoded in large pre-trained image-text models. Recently, these large pre-trained models have been successfully applied to the 3D domain in creative ways, such as guiding the optimization of differentiable 3D representations [28, 56, 38] or to generate 3D shapes from text prompts using interchangeability of text-image embeddings [61, 62], or using them for representation learning [78, 63].

In this paper, we introduce a straightforward yet effective approach called Sketch-A-Shape, for generating 3D shapes from sketches in a zero-shot setting using pre-trained vision models. Our method is based on the idea that 3D shape rendering features derived from large-scale pre-trained models (such as CLIP [57] and DINOv2 [53]) possess robust local semantic signals that can withstand domain shifts from renderings to sketches. In Sketch-A-Shape, we first train a VQ-VAE to acquire shape embeddings. Following this, a masked transformer is trained to model the distribution of shape embeddings conditioned on local semantic features from an image encoder that is pre-trained and frozen. During inference, the masked transformer is conditioned on local semantic features of the sketch instead, in order to produce the 3D shape. Our findings suggest that with some architectural design choices, this straightforward method enables us to generate several 3D shapes that can generalize across sketches of varying complexities.

To sum up, we make the following contributions:

- We propose Sketch-A-Shape, the first zero-shot approach for sketch-to-3D generation, leveraging large-scale pre-trained models to outdo the need of paired sketch-3D dataset.
- We experimentally show the generalization capability of our method among various datasets (section 4) with

different levels of sketch abstraction, going from simple doodles to professional sketches.

- We conduct thorough experiments to examine the different components of Sketch-A-Shape that contribute to the success of zero-shot shape generation via sketch.

2. Related Work

3D Generative Models. Significant progress has been made in the field of generative models for the creation of 3D shapes in various formats such as voxels [70, 9, 66, 30], CAD [71, 29, 36, 72], implicit representations [42, 7, 54], meshes [48, 19], and point clouds [49, 3, 75, 37, 74, 34]. Recent research on 3D generative models has focused primarily on the development of generative models based on VQ-VAE [47, 8, 73, 77, 62, 72], GAN[4, 65, 17], or diffusion models [80, 40, 50, 26]. The present study concentrates on connecting the sketch modality with 3D shapes across three different 3D representations: voxels, CAD, and implicit representation. Although our approach is based on VQ-VAE, it can be easily extended to GAN or diffusion-based generative models.

3D Zero-Shot Learning. Large pre-trained language and 2D vision models have been creatively used in several downstream 3D vision tasks. Initial works focused on using vision-text models such as CLIP [57] for 3D shape generation using text [61], optimizing nerfs [28], deforming meshes [44], stylizing meshes [46] and animating avatars [25]. More recently, text-to-image models such as Stable Diffusion [59] and Imagen [60], have been used for text-to-shape generation [56, 38], single-view reconstruction [76], and adding texture to 3D shapes [43]. To the best of our knowledge, our work is the first to explore zero-shot 3D shape generation from sketches by leveraging a pre-trained model.

3D Shape Generation from Sketch. Several supervised learning methods have been used to generate 3D shapes from sketches. Works such as [39] use a neural net to estimate depth and normals from a set of viewpoints for a given sketch, which are then integrated into a 3D point cloud. [10] proposes to use a CNN to predict the initial shape and then refine the shape using novel viewpoints using another neural network. Another work [31] represent the 3D shape and its occluding contours in a joint VAE latent space during training which enables them to retrieve a sketch during inference and generate a 3D shape. Sketch2Mesh [20] uses an encoder-decoder architecture to represent and refine a 3D shape to match the target external contour using a differentiable render. Methods such as [69, 79] employ a domain adaption network between unpaired sketch and rendered image data to boost performance on abstract hand-drawn sketches. To address the ambiguity problem of sketches, [79] introduces an additional encoder-decoder to

extract separate view and shape sketch features, while [16] proposes a sketch translator module to fully exploit the spatial information in a sketch and generate suitable features for 3D shape prediction. Recently, [35] trains a diffusion model for generation of 3D point clouds conditioned on sketches using a multi-stage training, and fine-tuning technique. However, we take the novel approach of not training on paired shape-sketch data at all and instead rely on the robustness of the local semantic features from a frozen large pre-trained image encoder such as CLIP.

3. Method

Our approach strives to generate 3D shapes from sketches of different complexities, solely employing synthetic renderings, and without the requirement of a paired dataset of sketches and shapes. The training data consists of 3D shapes, each denoted by \mathbf{S} , which can be represented as a voxel grid, implicit (e.g. occupancy), or CAD, and their r multi-view renderings ($\mathbf{I}^{1:r}$). Our approach involves two training stages. In the first stage, the shapes are transformed into discrete sequences of indices (shape embeddings), denoted by \mathbf{Z} , using a discrete auto-encoder framework [52]. In the second stage, the distribution of these indices is modeled using a transformer-based generative model that is conditioned on features of shape renderings obtained from a frozen pre-trained model. These shape rendering features are a grid of local features from the frozen pre-trained model which are converted into a sequence of local features then conditioned to the transformer through a cross-attention mechanism. During inference, we use an iterative decoding scheme [6] to generate the shape indices iteratively based on features of the sketch. Once the shape indices are generated, we can use the decoder of the auto-encoder to generate the 3D shape. The overall process is illustrated in Figure 2 .

3.1. Stage 1: Training Discrete Auto-encoder

In the first stage, we use an auto-encoder framework to capture the shape distribution into a compressed sequence of discrete indices (shape embeddings) among various modalities. To achieve this, we opt for the Vector Quantized Variational Auto-encoder (VQ-VAE) architecture [52] which efficiently models the 3D shape in a compressed latent space, circumventing posterior collapse and enabling the generation of high-quality 3D shapes. The dataset of 3D shapes \mathbf{S} , are transformed using an encoder, $E(\cdot)$, into a sequence of discrete indices \mathbf{Z} , pointing to a shape dictionary, which their distributions are then modeled in stage 2 using a transformer-based generative model. This is shown below:

$$\mathbf{Z} = VQ(E(\mathbf{S})), \quad \mathbf{S}' = D(\mathbf{Z}) \quad (1)$$

The shape \mathbf{S}' is then generated from \mathbf{Z} using a decoder, $D(\cdot)$, with a reconstruction loss $L_{rec}(S, S')$. We also use the commitment loss [52] to encourage encoder output $E(\cdot)$ commits to an embedding in the shape dictionary, and exponential moving average strategy [58, 81] to encourage dictionary entries to gradually be pulled toward the encoded features. When dealing with voxel representation, we leverage a 3D convolution based on the ResNet architecture [23] for both the encoder and decoder network. Whereas with implicit representation, we rely on a ResNet-based encoder and an up-sampling process for the decoder that generates a higher resolution volume, which is then queried locally to obtain the final occupancy [42, 55]. In the case of CAD representation, we use the SkexGen VQ-VAE architecture [72]. More details of the architectures are provided in the supplementary material.

3.2. Stage 2: Masked Transformer

The goal of stage 2 is to train a prior model which can effectively generate shape indices conditioned on a sketch at inference time. We achieve this by modelling the sequence of discrete indices (shape embedding \mathbf{Z}), produced from stage 1, using a conditional generative model. We use a bi-directional transformer [6] based network which is conditioned on the features of the synthetic 3D renderings using a cross-attention mechanism. During training, we randomly mask a fraction of shape indices with a special mask token [6] to produce \mathbf{Z}^{mask} . The training objective then becomes how to fully unmask the masked indices using the help of the provided conditional information. The training objective is to minimize:

$$L_{mask} = - \mathbb{E}_{\mathbf{Z}, \mathbf{C} \in D} [\log p(\mathbf{Z} | \mathbf{Z}^{mask}, \mathbf{C})] \quad (2)$$

Here, \mathbf{C} represents the conditional information from a given shape \mathbf{S} which are obtained from the multi-view image renderings of the 3D shape. At each iteration, we randomly sample a view to render an image of the 3D shape, which is then converted to local features using a locked pre-trained model. The choice of pre-trained model is an important design criteria which we investigate thoroughly in Section 4.4, and find that using large models trained on diverse data produces the most robust semantic local features which allow domain shift from synthetic renderings to sketches.

The local features sequence can be obtained from different parts of the pre-trained network, which we investigate in Section 4.5. Our findings indicate that utilizing the feature grid output of the deeper layers in the pre-trained models yields better results. This is because deeper layers generate more semantic features, and the grid structure of the feature preserves its local properties. We convert this grid into a sequence using a mapping network comprising of several MLP layers. Finally, we take features ob-

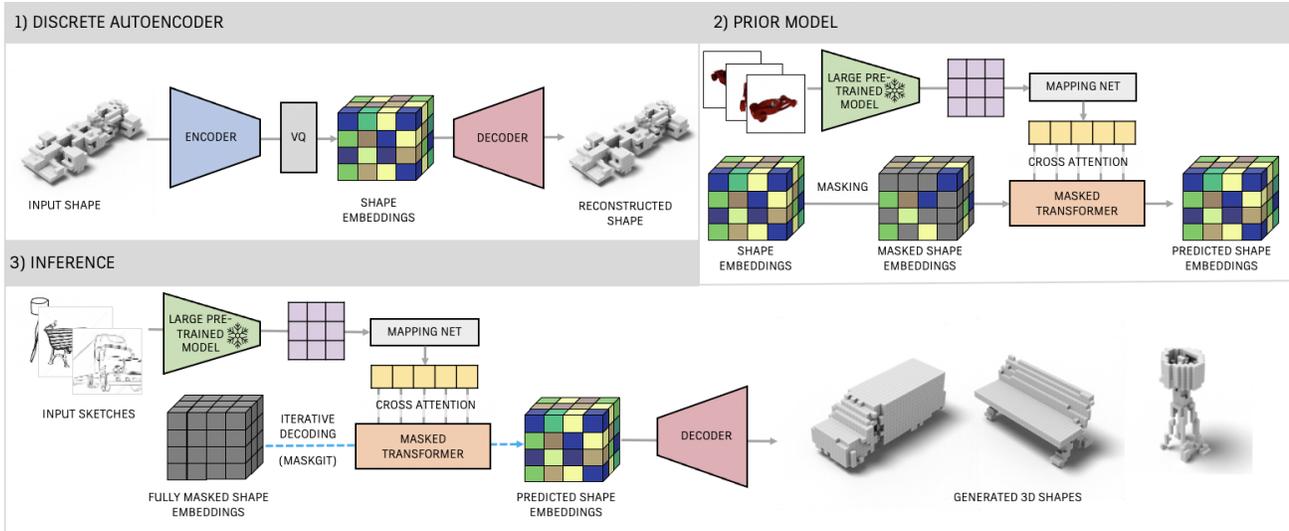


Figure 2: An overview of the Sketch-A-Shape method. The top row shows the two training stages. The bottom row shows how 3D shape is generated via sketch at inference time.

tained and add learnable positional encoding before applying cross-attention with the shape indices’ features at each transformer layer. The choice of conditioning is also an important design choice which we discuss in Section 4.6. Additionally, we replace the local features sequence with a null embedding sequence 5% of the time to allow for classifier-free guidance during inference.

3.3. Inference

During the generation phase, we first convert the sketch into a sequence of local features using the same frozen pre-trained model utilized during training. These local features are semantically robust and serve as the conditioning query for the transformer. We employ an iterative decoding scheme with a cosine schedule, similar to the one proposed in Mask-GIT [6]. The process begins with a completely masked set of indices, which are gradually unmasked in parallel using the conditional information provided by the local features sequence from the sketch. At each time step, the transformer predicts the complete unmasked shape sequence, of which a specific fraction of the highest confidence masked tokens are accepted. These selected tokens are designated as unmasked for the remaining steps, while the rest of the tokens are reset to masked, except for the already unmasked tokens from the previous steps. For each time step, we also apply classifier-free guidance [24] with a guidance scale of 3. This process continues until all the tokens are unmasked. Finally, the completely unmasked tokens are converted into the 3D object using the shape decoder trained in stage 1. It is worth noting that we can restart the same process multiple times to generate differ-

ent 3D shapes for the same sketch query.

4. Experiments

In this section, we present the results of our experiments evaluating the accuracy and fidelity of the generated output produced by our model. We conducted each experiment three times for each metric and reported the average result for each. The experimental setup details are provided in the supplementary material with additional results that may be of interest.

Training Dataset. Our experimentation utilizes two subsets of the ShapeNet(v2) dataset [5]. The first subset, ShapeNet13, consists of 13 categories from ShapeNet, which were also employed in previous studies [9, 42]. In line with Sketch2Model [79], we adopt the same train/val/test partition. The second subset, ShapeNet55, includes all 55 categories of ShapeNet and we follow the same split as [62]. We use the DeepCAD [71] dataset to train our CAD model.

Evaluation Sketch Dataset. One advantage of our method is that it’s not trained on paired (shape, sketch) datasets. Therefore, to comprehensively evaluate its performance, we tested it on various sketch datasets that range from professional to non-expert sketches. Specifically, we utilized the ShapeNet-Sketch dataset [79], which comprises 1300 free-hand sketches across ShapeNet13. In addition, we employed the ImageNet-Sketch dataset [67], which contains 50 sketch images for 1000 ImageNet classes obtained from Google, encompassing a range of professional to non-expert sketches. Moreover, we utilized the TU-Berlin Sketch dataset [13], which includes 20,000 non-expert sketches of

250 object categories. Lastly, QuickDraw Dataset [1, 21] is a collection of 50 million drawings across 345 categories, contributed by players of the game *Quick, Draw!* [32]. ImageNet-Sketch, TU-Berlin Sketch, and QuickDraw datasets also lack ground-truth 3D models, and we only utilized the categories of ShapeNet for these datasets. To evaluate our CAD model we use synthetic edge map sketches but don't train the model using edge maps as augmentation.

Evaluation Metrics. To evaluate our method on different sketch datasets we use two metrics: classification accuracy and human evaluation which are outlined below.

1. **Classifier Accuracy.** As we are dealing with sketch data that lacks ground-truth 3D models, we use the Accuracy (Acc) metric to ensure that the generated shape for a given sketch corresponds to its category. To achieve this, we employ a pre-trained shape classifier, as implemented in [61, 62]. We use this metric for all datasets: ImageNet-Sketch [67], TU-Berlin [13], ShapeNet-Sketch [79], and QuickDraw [1]. We refer to this metric as IS-Acc, TU-Acc, SS-Acc, and QD-Acc, respectively. As our method can generate multiple shape per sketch query, we report the mean across 5 sampled shapes for a given sketch query.
2. **Human Perceptual Evaluation.** We also use Amazon SageMaker Ground Truth and crowd workers from the Mechanical Turk workforce [45] to evaluate how well our generated 3D models preserve important geometric and stylistic details from the sketches.

4.1. Qualitative Results

In Figure 3, we visualize sample generated 3D shapes in different representations such as voxel, implicit, and CAD from sketches of different domains. As shown, our method performs reasonably well on different types of sketches (from simple to professional drawings), particularly when there is ambiguity (such as the view angle of drawings) given the nature of 2D sketches.

4.2. Human Perceptual Evaluation

In addition to generating shapes in the same broad object category as abstract hand drawn sketches, our method is also able to incorporate geometric and stylistic details from a sketch or concept design into the final 3D model. To demonstrate this quantitatively, we run a human perceptual evaluation using Amazon SageMaker Ground Truth and crowd workers from the Mechanical Turk workforce [45]. We evaluate 691 generated models, conditioned on sketches from TU-Berlin [13], ShapeNet-Sketch [79], ImageNet-Sketch [67] and QuickDraw [21].

The human evaluation is posed as a two-alternative forced choice study [14]. The crowd workers are shown images with a sketch on the left hand side and two images

of generated 3D models on the right hand side. An example is shown in Figure 4. One of the generated models was conditioned on the sketch shown, while the other was conditioned on a randomly selected sketch from the same object category. The crowd workers are asked the question "Which of the 3D models on the right hand side best matches the sketch on the left hand side?". The study is designed to measure the extent to which humans perceive our generated 3D models as preserving the shape and stylistic details presented in the sketch, as opposed to simply creating a model from the same object category.

We show each image to 7 independent crowd workers and count the number of images for which 4 or more of them correctly identify the 3D model which was conditioned on the sketch. The results are shown in Table 1. On average 71.1% of our generated 3D models are correctly identified by a majority of the crowd workers. We note that the sketches in TU-Berlin and ShapeNet-Sketch give rise to generations which were easier for the crowd workers to identify, with 74.9% and 73.1% being selected correctly. While these sketches often have a high level of abstraction, they communicate enough detail about the shape for our method to create distinctive 3D models which humans can identify. While ImageNet-Sketch contains superior artwork, often with shading, shadows and other cues to the 3D nature of the shapes, many of the pictures contain full scenes with backgrounds and additional superfluous details. This makes the generation of single objects more challenging, which is reflected by the fact that only 68.1% are correctly identified by the crowd workers. We note qualitatively that in cases where shaded sketches do not contain backgrounds or additional clutter the generated results look better, indicating the utility of our method for quickly generating 3D models from concept designs. The sketches in the QuickDraw dataset are sourced from the from the online game *Quick, Draw!* [32], in which contributors are asked to draw a shape in less than 20 seconds. QuickDraw is the most abstract and noisy sketch dataset, with many of the sketches being drawn with a computer mouse. While our method typically generates 3D shapes of the correct category, only 67.9% of the generations are correctly identified by the crowd workers.

4.3. Comparison with Supervised Models

As there is currently a lack of zero-shot methodologies for generating shapes from sketches, we compared our results to those of a supervised approach called Sketch2Model [79], which was trained on a dataset of paired sketch-3D shapes. We evaluated both methods using our classifier accuracy metric, and the results are presented in Table 2. Our model was not exposed to any sketch-3D pairings during training, but it displays superior generation capabilities compared to Sketch2Model across different datasets.

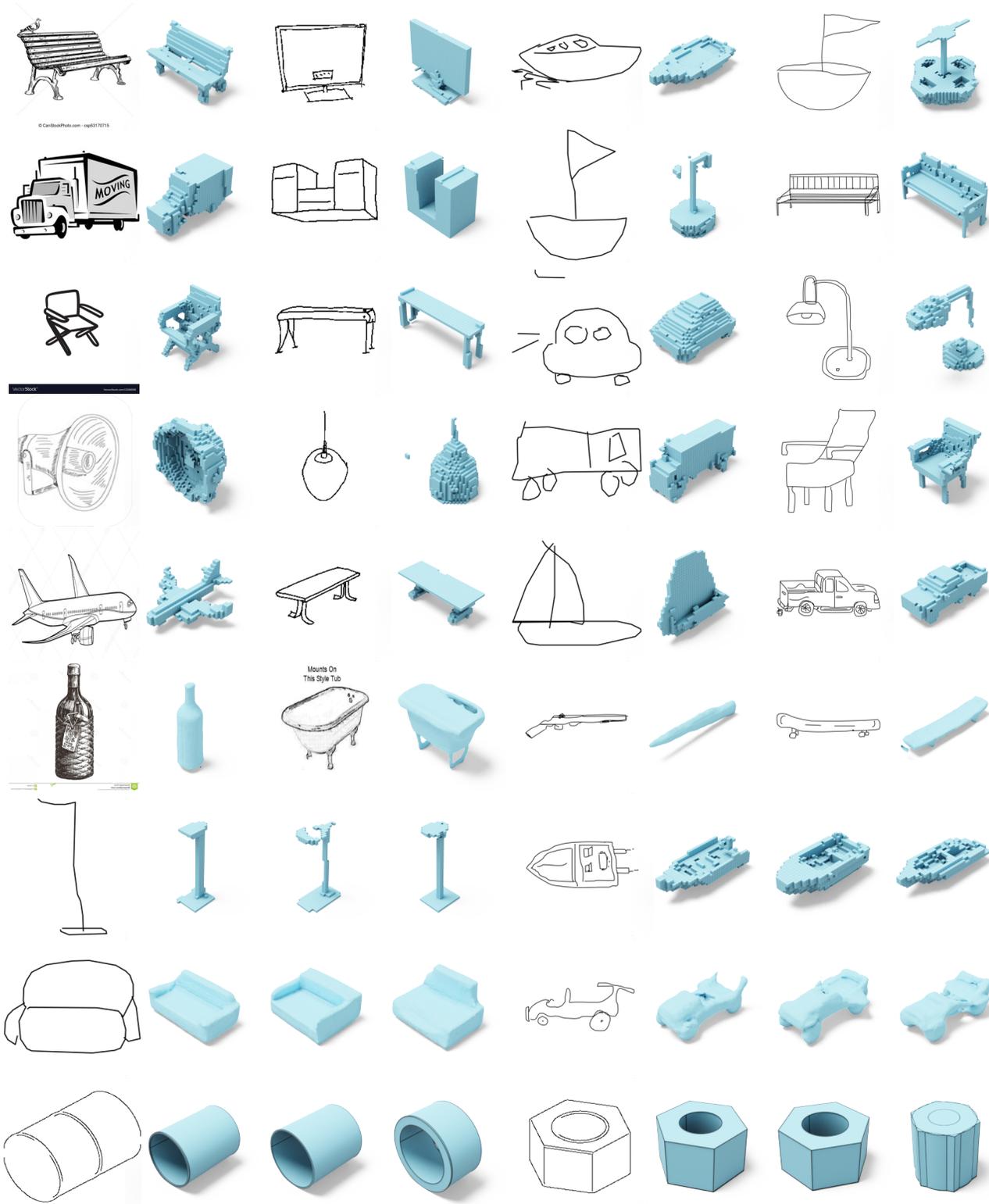


Figure 3: Generated 3D shapes from sketches of different domains. Rows 1–5, and 6 refer to voxel, and implicit representations, respectively. The last three rows show different 3D samples generated per sketch in voxel, implicit, and CAD formats.



Figure 4: An example of the images shown to the Mechanical Turk crowd workers. The question is “Which of the 3D models on the right hand side best matches the sketch on the left hand side?” and the options are “Top” and “Bottom”.

Dataset	% correctly identified
All	71.1%
TU-Berlin	74.9%
ShapeNet-Sketch	73.1%
ImageNet-Sketch	68.1%
QuickDraw	67.9%

Table 1: Results of the human perceptual evaluation by dataset. We show the percentage of generated models for which the majority of the 7 crowd workers reviewing each image correctly identified the 3D model conditioned on the sketch.

Method	QD-Acc \uparrow	TU-Acc \uparrow	SS-Acc \uparrow	IS-Acc \uparrow
S2M [79]	27.4	19.8	26.0	12.0
Ours	58.8	81.5	79.7	74.2

Table 2: Classification accuracy comparison with the supervised method Sketch2Model.

We attribute this difference in performance to several factors. Firstly, we believe that Sketch2Model may be more effective for single-category training rather than for the 13 categories in the ShapeNet dataset. Additionally, because Sketch2Model is a supervised method, it was not exposed to out-of-distribution sketches during training, which may have caused its performance to deteriorate. We provide further details and qualitative comparison with Sketch2Model and other supervised methods in the supplementary material.

4.4. Investigating Pre-Trained Models

This section involves an extensive study of several pre-trained models that are open-sourced and trained on different datasets. The results are present in Table 3. There are 3 major things we investigate through this experiment. First, we investigate the importance of utilizing local grid features of pre-trained models. Typically, pre-trained models possess a global projection vector that is employed for downstream tasks like classification. We compare the efficacy of conditioning our generative model with the global projection vector (row 1) versus the local grid features (row 2). Our findings demonstrate that leveraging local grid features yields better performance compared to the global projection vector for most of the datasets. Furthermore, even from a visual standpoint, we observe that local grid features preserve more local details. It is worth noting that these accuracies are further improved by utilizing classifier-free guidance (CFG), as illustrated in row 3.

Next, we investigate the role of size of pre-trained models and find that increasing the size of the same class of pre-trained model, despite being trained on the same data, results in better zero-shot performance. This phenomenon is evident in the case of the ViT-based [12] CLIP model, where upgrading from the B-32 model to the L-14 model yields a significant improvement in performance. This trend is also observed in the ResNet-based [23] models. Interestingly, it is worth mentioning that the ResNet-based [23] models perform worse than their corresponding ViT-based [12] CLIP models. This could be attributed to the ResNet models’ emphasis on high-frequency, textural features [18].

Finally, we explore how different datasets impact the training of these models. Our findings indicate that the model’s performance remains comparable when trained on extensive datasets such as LAION-2B [64], DINOv2 Dataset [53] or OpenAI internal dataset [57]. However, when we reduce the dataset size significantly, such as in the case of the masked autoencoder [22] trained on 400 times less data from ImageNet [11], its performance significantly declines. Despite being trained on the reconstruction objective, we believe that the masked autoencoder’s performance drop is primarily due to the significantly reduced dataset size, as it still performs reasonably well on this task. Additionally, it is important to highlight that language supervision is unnecessary to acquire resilient features from extensive pre-trained models, as demonstrated by the outcomes of DINOv2.

4.5. Accuracy across Different Layers

In this experiment, we explore the optimal layer of the vision transformer (L-14 model) from which we extract the local conditioning features. Table 4 summarizes our findings. We note that as we delve deeper into the vision transformer architecture, the features extracted from the deeper

Resolution	CFG	Network	Dataset	QD-Acc \uparrow	TU-Acc \uparrow	SS-Acc \uparrow	IS-Acc \uparrow
1 x 512	×	B-32 [57]	OpenAI [57]	36.65	61.14	62.86	55.96
50 x 768	×	B-32 [57]	OpenAI [57]	37.85	63.25	63.78	52.79
50 x 768	✓	B-32 [57]	OpenAI [57]	38.86	65.86	67.36	49.19
197 x 768	✓	B-16 [57]	OpenAI [57]	38.47	71.66	70.72	61.10
257 x 1024	✓	L-14 [57]	OpenAI [57]	55.45	77.15	74.53	69.06
144 x 3072	✓	RN50x16 [57]	OpenAI [57]	34.61	70.81	58.82	59.00
196 x 4096	✓	RN50x64 [57]	OpenAI [57]	46.93	73.79	59.41	64.19
257 x 1024	✓	Open-L-14 [27]	LAION-2B [64]	54.63	77.60	69.03	68.35
256 x 1024	✓	DINO-L-14 [53]	DINOv2 [53]	39.73	71.12	72.10	55.94
197 x 1024	✓	MAE-L [22]	ImageNet [11]	19.31	30.52	38.79	26.65
257 x 1280	✓	MAE-H [22]	ImageNet [11]	18.70	31.63	37.47	31.42

Table 3: Ablations on using different pre-trained models. Resolution indicates the size of the local features obtained from pre-trained model whereas CFG stands for classifier free guidance.

L-14 Layer	QD-Acc \uparrow	TU-Acc \uparrow	SS-Acc \uparrow	IS-Acc \uparrow
1	17.16	22.21	24.16	11.13
7	16.93	25.03	33.24	12.52
13	29.43	51.50	59.64	40.98
19	54.03	75.77	74.45	63.83
24	55.45	77.15	74.53	69.06

Table 4: Accuracy across 24 layers of L-14 architecture.

layers contain more significant semantic information leading to higher accuracy. Moreover, this indicates that the model maintains the positional correlation between patches instead of treating them as global information repositories as visually we can see local semantic generation.

4.6. Design Choices for Conditioning

Table 5 presents our investigation into the impact of the mapping network’s size and the attention mechanism used for conditioning the image features to the transformer. Our results show that incorporating a mapping layer does enhance the model’s performance, with the optimal number of MLP layers being two. Furthermore, our findings suggest that cross-attention with a learnable positional embedding is the most effective conditioning mechanism, as evidenced by the deteriorated performance on removing positional embedding or using self attention as shown in the last two rows of the table.

4.7. Effect of Augmentation

In our final investigation, we explore whether the addition of data augmentation improves the accuracy of shape generation across datasets. The results are summarized in Table 6. We make two noteworthy observations. Firstly, even without data augmentation, our method performs relatively well, indicating the robustness of pre-trained mod-

M	Att. type	QD-Acc \uparrow	TU-Acc \uparrow	SS-Acc \uparrow	IS-Acc \uparrow
0	cross	55.45	77.15	74.53	69.06
1	cross	54.69	77.91	77.19	69.17
2	cross	57.52	79.41	78.18	70.53
3	cross	52.10	76.85	76.85	68.03
2	no pos	55.81	78.75	78.02	71.05
2	self	49.33	76.97	75.42	69.07

Table 5: Ablation on design choices for conditioning. M is the number of MLP layers in mapping network (0 meaning just a linear layer). Note, self stands for self attention, cross for cross attention with positional embedding and no pos stands for cross attention with no positional embedding.

Augmentation	QD-Acc \uparrow	TU-Acc \uparrow	SS-Acc \uparrow	IS-Acc \uparrow
No Aug	57.52	79.41	78.18	70.53
Affine (A)	61.52	78.96	77.95	74.19
Color (C)	58.37	78.03	77.15	71.24
Canny (CN)	53.15	79.04	80.48	68.81
CN + A + C	58.96	81.48	79.71	74.20

Table 6: Effect of different augmentations on the shapes.

els. Secondly, different types of augmentations have a more significant impact on certain datasets than others. For instance, affine transformation significantly enhances the performance of QuickDraw and ImageNet Sketch, while canny edge augmentation is more effective for the ShapeNet Sketch dataset. Consequently, we decide to train a network with all augmentations and find that, on balance across datasets, it performs the best.

5. Conclusion

In this paper, we demonstrate how a 3D generative model conditioned on local features from a pre-trained large-scale image model such as CLIP can be used to generate 3D shapes from sketches. We show how this method can generate multiple shapes for different abstraction of sketches and can be applied to multiple 3D representations. Future work will involve training on much larger and diverse 3D shape datasets and consequently testing on different styles of sketches and levels of details.

References

- [1] The quick, draw! dataset. <https://github.com/googlecreativelab/quickdraw-dataset>.
- [2] Sketch2model: View-aware 3d modeling from single free-hand sketches. <https://github.com/bennyguo/sketch2model>.
- [3] Panos Achlioptas, Olga Diamanti, Ioannis Mitliagkas, and Leonidas Guibas. Learning representations and generative models for 3d point clouds. In *International conference on machine learning*, pages 40–49. PMLR, 2018.
- [4] Eric R Chan, Connor Z Lin, Matthew A Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas J Guibas, Jonathan Tremblay, Sameh Khamis, et al. Efficient geometry-aware 3d generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16123–16133, 2022.
- [5] Angel X. Chang, Thomas Funkhouser, Leonidas 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, 2015. cite arxiv:1512.03012.
- [6] Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T Freeman. Maskgit: Masked generative image transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11315–11325, 2022.
- [7] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5939–5948, 2019.
- [8] An-Chieh Cheng, Xueting Li, Sifei Liu, Min Sun, and Ming-Hsuan Yang. Autoregressive 3d shape generation via canonical mapping. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part III*, pages 89–104. Springer, 2022.
- [9] 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 *European conference on computer vision*, pages 628–644. Springer, 2016.
- [10] Johanna Delanoy, Mathieu Aubry, Phillip Isola, Alexei A Efros, and Adrien Bousseau. 3d sketching using multi-view deep volumetric prediction. *Proceedings of the ACM on Computer Graphics and Interactive Techniques*, 1(1):1–22, 2018.
- [11] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [12] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [13] Mathias Eitz, James Hays, and Marc Alexa. How do humans sketch objects? *ACM Trans. Graph. (Proc. SIGGRAPH)*, 31(4):44:1–44:10, 2012.
- [14] G.T. Fechner. *Elemente der Psychophysik*. Number pt. 1 in *Elemente der Psychophysik*. Breitkopf und Härtel, 1860.
- [15] Rao Fu, Xiao Zhan, Yiwen Chen, Daniel Ritchie, and Srinath Sridhar. Shapecrafter: A recursive text-conditioned 3d shape generation model, 2022.
- [16] Chenjian Gao, Qian Yu, Lu Sheng, Yi-Zhe Song, and Dong Xu. Sketchsampler: Sketch-based 3d reconstruction via view-dependent depth sampling. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part I*, pages 464–479. Springer, 2022.
- [17] Jun Gao, Tianchang Shen, Zian Wang, Wenzheng Chen, Kangxue Yin, Daiqing Li, Or Litany, Zan Gojcic, and Sanja Fidler. Get3d: A generative model of high quality 3d textured shapes learned from images. In *Advances In Neural Information Processing Systems*, 2022.
- [18] Amin Ghiasi, Hamid Kazemi, Eitan Borgnia, Steven Reich, Manli Shu, Micah Goldblum, Andrew Gordon Wilson, and Tom Goldstein. What do vision transformers learn? a visual exploration. *arXiv preprint arXiv:2212.06727*, 2022.
- [19] Thibault Groueix, Matthew Fisher, Vladimir G Kim, Bryan C Russell, and Mathieu Aubry. A papier-mâché approach to learning 3d surface generation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 216–224, 2018.
- [20] Benoit Guillard, Edoardo Remelli, Pierre Yvernay, and Pascal Fua. Sketch2mesh: Reconstructing and editing 3d shapes from sketches. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13023–13032, 2021.
- [21] David Ha and Douglas Eck. A neural representation of sketch drawings. In *International Conference on Learning Representations*, 2018.
- [22] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16000–16009, 2022.
- [23] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [24] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.

- [25] Fangzhou Hong, Mingyuan Zhang, Liang Pan, Zhongang Cai, Lei Yang, and Ziwei Liu. Avatarclip: Zero-shot text-driven generation and animation of 3d avatars. *arXiv preprint arXiv:2205.08535*, 2022.
- [26] Ka-Hei Hui, Ruihui Li, Jingyu Hu, and Chi-Wing Fu. Neural wavelet-domain diffusion for 3d shape generation. In *SIGGRAPH Asia 2022 Conference Papers*, pages 1–9, 2022.
- [27] Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hananeh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. Openclip, July 2021. If you use this software, please cite it as below.
- [28] Ajay Jain, Ben Mildenhall, Jonathan T Barron, Pieter Abbeel, and Ben Poole. Zero-shot text-guided object generation with dream fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 867–876, 2022.
- [29] Pradeep Kumar Jayaraman, Joseph G Lambourne, Nishkrit Desai, Karl DD Willis, Aditya Sanghi, and Nigel JW Morris. Solidgen: An autoregressive model for direct b-rep synthesis. *arXiv preprint arXiv:2203.13944*, 2022.
- [30] Danilo Jimenez Rezende, SM Eslami, Shakir Mohamed, Peter Battaglia, Max Jaderberg, and Nicolas Heess. Unsupervised learning of 3d structure from images. *Advances in neural information processing systems*, 29, 2016.
- [31] Aobo Jin, Qiang Fu, and Zhigang Deng. Contour-based 3d modeling through joint embedding of shapes and contours. In *Symposium on interactive 3D graphics and games*, pages 1–10, 2020.
- [32] J. Jongejan, H. Rowley, T. Kawashima, J. Kim, , and N. Fox-Gieg. The quick, draw! - a.i. experiment. <https://quickdraw.withgoogle.com/>.
- [33] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [34] Wei-Jan Ko, Hui Huang, Yu-Liang Kuo, Chen-Yi Chiu, Li-Heng Wang, and Wei-Chen Chiu. Rpg: Learning recursive point cloud generation. *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 544–551, 2021.
- [35] Di Kong, Qiang Wang, and Yonggang Qi. A diffusion-refinement model for sketch-to-point modeling. In *Proceedings of the Asian Conference on Computer Vision*, pages 1522–1538, 2022.
- [36] Joseph George Lambourne, Karl Willis, Pradeep Kumar Jayaraman, Longfei Zhang, Aditya Sanghi, and Kamal Rahimi Malekshan. Reconstructing editable prismatic cad from rounded voxel models. In *SIGGRAPH Asia 2022 Conference Papers*, pages 1–9, 2022.
- [37] Chun-Liang Li, Manzil Zaheer, Yang Zhang, Barnabas Poczos, and Ruslan Salakhutdinov. Point cloud gan. *arXiv preprint arXiv:1810.05795*, 2018.
- [38] Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. Magic3d: High-resolution text-to-3d content creation. *arXiv preprint arXiv:2211.10440*, 2022.
- [39] Zhaoliang Lun, Matheus Gadelha, Evangelos Kalogerakis, Subhransu Maji, and Rui Wang. 3d shape reconstruction from sketches via multi-view convolutional networks. In *2017 International Conference on 3D Vision (3DV)*, pages 67–77. IEEE, 2017.
- [40] Shitong Luo and Wei Hu. Diffusion probabilistic models for 3d point cloud generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2837–2845, 2021.
- [41] Sarah Boyes Maxseiner. *Sketch Quality Prediction Using Transformers*. PhD thesis, Virginia Tech, 2023.
- [42] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [43] Gal Metzer, Elad Richardson, Or Patashnik, Raja Giryes, and Daniel Cohen-Or. Latent-nerf for shape-guided generation of 3d shapes and textures. *arXiv preprint arXiv:2211.07600*, 2022.
- [44] Oscar Michel, Roi Bar-On, Richard Liu, Sagie Benaim, and Rana Hanocka. Text2mesh: Text-driven neural stylization for meshes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13492–13502, 2022.
- [45] Abhishek Mishra. Machine learning in the aws cloud: Add intelligence to applications with amazon sagemaker and amazon rekognition, 2019.
- [46] Shailesh Mishra and Jonathan Granskog. Clip-based neural neighbor style transfer for 3d assets. *arXiv preprint arXiv:2208.04370*, 2022.
- [47] Paritosh Mittal, Yen-Chi Cheng, Maneesh Singh, and Shubham Tulsiani. Autosdf: Shape priors for 3d completion, reconstruction and generation. *arXiv preprint arXiv:2203.09516*, 2022.
- [48] Charlie Nash, Yaroslav Ganin, S. M. Ali Eslami, and Peter Battaglia. PolyGen: An autoregressive generative model of 3D meshes. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 7220–7229. PMLR, 13–18 Jul 2020.
- [49] Charlie Nash and Christopher KI Williams. The shape variational autoencoder: A deep generative model of part-segmented 3d objects. In *Computer Graphics Forum*, volume 36, pages 1–12. Wiley Online Library, 2017.
- [50] Alex Nichol, Heewoo Jun, Prafulla Dhariwal, Pamela Mishkin, and Mark Chen. Point-e: A system for generating 3d point clouds from complex prompts. *arXiv preprint arXiv:2212.08751*, 2022.
- [51] Alex Nichol, Heewoo Jun, Prafulla Dhariwal, Pamela Mishkin, and Mark Chen. Point-e: A system for generating 3d point clouds from complex prompts, 2022.
- [52] Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learning. *arXiv preprint arXiv:1711.00937*, 2017.
- [53] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez,

- Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.
- [54] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [55] Songyou Peng, Michael Niemeyer, Lars Mescheder, Marc Pollefeys, and Andreas Geiger. Convolutional occupancy networks. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part III 16*, pages 523–540. Springer, 2020.
- [56] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022.
- [57] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021.
- [58] Ali Razavi, Aaron Van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images with vq-vae-2. *Advances in neural information processing systems*, 32, 2019.
- [59] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022.
- [60] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022.
- [61] Aditya Sanghi, Hang Chu, Joseph G Lambourne, Ye Wang, Chin-Yi Cheng, Marco Fumero, and Kamal Rahimi Malekshah. Clip-forge: Towards zero-shot text-to-shape generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18603–18613, 2022.
- [62] Aditya Sanghi, Rao Fu, Vivian Liu, Karl Willis, Hooman Shayani, Amir Hosein Khasahmadi, Srinath Sridhar, and Daniel Ritchie. Textcraft: Zero-shot generation of high-fidelity and diverse shapes from text. *arXiv preprint arXiv:2211.01427*, 2022.
- [63] Kristofer Schlachter, Benjamin Ahlbrand, Zhu Wang, Ken Perlin, and Valerio Ortenzi. Zero-shot multi-modal artist-controlled retrieval and exploration of 3d object sets. In *SIGGRAPH Asia 2022 Technical Communications*, SA '22, New York, NY, USA, 2022. Association for Computing Machinery.
- [64] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *arXiv preprint arXiv:2210.08402*, 2022.
- [65] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. Graf: Generative radiance fields for 3d-aware image synthesis. *Advances in Neural Information Processing Systems*, 33:20154–20166, 2020.
- [66] Maxim Tatarchenko, Alexey Dosovitskiy, and Thomas Brox. Octree generating networks: Efficient convolutional architectures for high-resolution 3d outputs. In *Proceedings of the IEEE international conference on computer vision*, pages 2088–2096, 2017.
- [67] Haohan Wang, Songwei Ge, Zachary Lipton, and Eric P Xing. Learning robust global representations by penalizing local predictive power. In *Advances in Neural Information Processing Systems*, pages 10506–10518, 2019.
- [68] Jiayun Wang, Jierui Lin, Qian Yu, Runtao Liu, Yubei Chen, and Stella X Yu. 3d shape reconstruction from free-hand sketches. *arXiv preprint arXiv:2006.09694*, 2020.
- [69] Lingjing Wang, Cheng Qian, Jifei Wang, and Yi Fang. Un-supervised learning of 3d model reconstruction from hand-drawn sketches. In *Proceedings of the 26th ACM international conference on Multimedia*, pages 1820–1828, 2018.
- [70] Jiajun Wu, Chengkai Zhang, Tianfan Xue, Bill Freeman, and Josh Tenenbaum. Learning a probabilistic latent space of object shapes via 3d generative-adversarial modeling. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016.
- [71] Rundi Wu, Chang Xiao, and Changxi Zheng. Deepcad: A deep generative network for computer-aided design models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6772–6782, 2021.
- [72] Xiang Xu, Karl DD Willis, Joseph G Lambourne, Chin-Yi Cheng, Pradeep Kumar Jayaraman, and Yasutaka Furukawa. Skexgen: Autoregressive generation of cad construction sequences with disentangled codebooks. *arXiv preprint arXiv:2207.04632*, 2022.
- [73] Xingguang Yan, Liqiang Lin, Niloy J Mitra, Dani Lischinski, Daniel Cohen-Or, and Hui Huang. Shapeformer: Transformer-based shape completion via sparse representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6239–6249, 2022.
- [74] Guandao Yang, Xun Huang, Zekun Hao, Ming-Yu Liu, Serge Belongie, and Bharath Hariharan. Pointflow: 3d point cloud generation with continuous normalizing flows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4541–4550, 2019.
- [75] Yaoqing Yang, Chen Feng, Yiru Shen, and Dong Tian. Foldingnet: Point cloud auto-encoder via deep grid deformation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 206–215, 2018.
- [76] Yufei Ye, Shubham Tulsiani, and Abhinav Gupta. Shelf-supervised mesh prediction in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8843–8852, 2021.
- [77] Biao Zhang, Matthias Nießner, and Peter Wonka. 3dilig: Irregular latent grids for 3d generative modeling. *arXiv preprint arXiv:2205.13914*, 2022.

- [78] Renrui Zhang, Ziyu Guo, Wei Zhang, Kunchang Li, Xupeng Miao, Bin Cui, Yu Qiao, Peng Gao, and Hongsheng Li. Pointclip: Point cloud understanding by clip. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8552–8562, 2022.
- [79] Song-Hai Zhang, Yuan-Chen Guo, and Qing-Wen Gu. Sketch2model: View-aware 3d modeling from single free-hand sketches. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6012–6021, 2021.
- [80] Linqi Zhou, Yilun Du, and Jiajun Wu. 3d shape generation and completion through point-voxel diffusion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5826–5835, 2021.
- [81] Adrian Łańcucki, Jan Chorowski, Guillaume Sanchez, Richard Marxer, Nanxin Chen, Hans J. G. A. Dolfing, Sameer Khurana, Tanel Alumäe, and Antoine Laurent. Robust training of vector quantized bottleneck models, 2020.

Supplementary Material

Category	% correctly identified
loudspeaker	68.0%
airplane	58.5%
bench	80.0%
car	69.4%
boat	64.3%
monitor	70.0%
lamp	89.1%
table	80.0%
cabinet	83.3%
sofa	71.4%
gun	65.0%
chair	83.6%
phone	52.5%

Table 7: Results of the human perceptual evaluation by class. We show the percentage of generations the majority of the 7 crowd workers reviewing each image correctly identified the 3D model generated by conditioning on the sketch.

A. Human Perceptual Evaluation

In Subsection 4.2 (main paper) we show the results of the human perceptual study broken down according to the dataset which the target sketch came from. In addition, the results can be broken down based on object category of the target sketch, as shown in Table 7. We see a wide range of performance across the different object categories, with “lamps” being correctly identified 89.1% of the time, while phones are identified just 52.5% of the time, little better than random. Categories which perform well, such as chair, table and sofa, tend to have distinctive shapes which are easy to communicate with sketches. Categories like airplane and gun produce good models, but these are not distinctive enough for the human evaluators to distinguish the correct 3D model from a random model of in the same category. Lower performance on these categories may also relate to the difficulty of drawing objects of these types. We believe having texture can further improve the human perceptual results.

As each generated model is rated by 7 individual crowd workers, we can count the number of raters who correctly identified the generated model, giving us a “shape recognizably score” from 0 to 7. In Figure 5 we show examples from selected categories with the highest and lowest “shape recognizably scores”. For “airplane” category the least recognizable model appears to be in the wrong category, due to the unusual orientation of the sketch. The most and least recognizable sketch in the “bench” category both

Method	Type	IOU \uparrow
Sketch2Mesh [20]	Supervised	0.195
Sketch2Model [79]	Supervised	0.205
Sketch2Point [68]	Supervised	0.163
SketchSampler [16]	Supervised	0.244
ours	Zero-shot	0.292

Table 8: IOU Results Comparison. It can be seen that our method outperforms supervised methods.

come from the Imagenet-Sketch dataset. The sketch for the most recognizable model contains a single bench while the sketch for the least recognizable model also contains background elements like trees and a lamppost. For the “gun” category the most recognizable model is actually from a sketch which looks nothing like a gun. The least recognizable model is a generic gun which does not closely follow the shape of the sketch. The Figure shows how the human evaluation is measure the ability of our method to generate distinctive shapes reflecting the geometry of the sketches as well as general generation quality.

B. Comparison with Supervised Methods

B.1. Quantitative comparison

We evaluate the quality of generated shapes on the ShapeNet-Sketch dataset [79] using Intersection over Union (IOU) with 32^3 voxel shapes, as shown in [41]. This is the only dataset among the four we assessed that includes ground truth 3D voxels. We compare our results to those of other supervised methods presented in Table 8, exactly as in [41]. Our generative model generates 5 shapes based on a given sketch query in the ShapeNet-Sketch dataset and averages the IOU results. Although our method is not trained on any sketch data, it outperforms the supervised baseline. This indicates that the pre-trained model’s learned features are effective in enabling our method to generate 3D shapes using sketches in a zero-shot manner.

B.2. Qualitative comparison

We additionally provide a qualitative comparison with Sketch2Model [79] and SketchSampler [16]. For this comparison, we considered diverse sketches with different levels of abstraction in the same classes of ShapeNet from four datasets: TU-Berlin [13], ShapeNet-Sketch [79], ImageNet-Sketch [67], and QuickDraw [1]. Implementation details can be found in Appendix C. Results are in Figure 6.

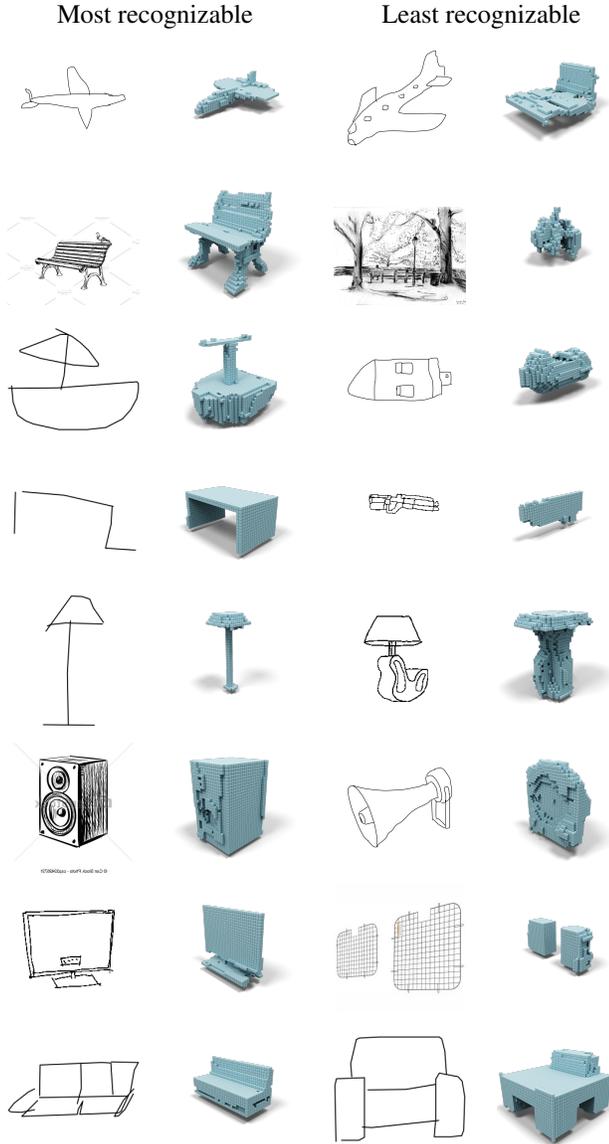


Figure 5: Generated results which were most and least recognizable for the human evaluators.

We can see that Sketch2Model reconstructed meshes often grasp the overall sketch shape, but they appear too smooth and lack geometric details. This method was originally intended for a single category scenario, as presented in the paper. However, this is often impractical.

Similarly, SketchSampler fails to generalize to abstract or out-of-distribution sketches. The resulting point clouds present artifacts and outliers, especially in the direction of the point of view of the sketch (shapes proportion are only preserved when the point clouds are seen from this point of view). Unlike our approach, SketchSampler is designed for

professional sketches only, with reliable shapes and fine-grained details. Thus, it cannot deal with sketches with significant deformation or only expressing conceptual ideas, like the ones in QuickDraw [1].

C. Architecture and Experiment Details

Training Details. We use the Adam Optimizer [33] with a fixed learning rate of $1e-4$ for training. The network is trained for 300 epochs during Stage 1 and for 250 epochs during Stage 2. We do not employ any learning rate scheduler during training. We train the 32^3 voxel model solely on the ShapeNet13 dataset, while the Implicit model is trained on the ShapeNet55 subset. The CAD model is trained on the DeepCAD dataset [72]. This is done to demonstrate the versatility and adaptability of our method to different datasets.

Stage 1 Details. For both the Implicit VQ-VAE and 32^3 VQ-VAE we use a codebook size of 512, a grid size of 8^3 and embedding dimensions of size 64. We employ the ResNet architecture for the 32^3 VQ-VAE, for both the encoder and decoder. In the case of Implicit VQ-VAE, we use the ResNet architecture for the encoder whereas we use a decoder that produces a higher resolution volume, which is then queried locally to obtain the final occupancy [42, 55, 15]. The pretrained VQ-VAE from SkexGen [72] is used for the CAD representation which is composed of three Transformer encoders and decoders for the topology, geometry and extrusions of a CAD model. The models output $4 + 2 + 4 = 10$ codes, with a total codebook size of 1000.

Stage 2 Details. For Stage 2, we utilize a bidirectional Transformer with 8 attention blocks, 8 attention heads, and a token size of 256. We use 24 renderings [9] for both the ShapeNet13 and ShapeNet55 experiments. During inference, we run the Transformer for 15 steps with classifier-free guidance, and the scale parameter is set to 3. The CLIP ViT-L/14 model is employed in all experiments, except in Table 3 of the main paper, where we conduct an ablation study over different pre-trained models. For all experiments, except Table 4, we incorporate cross-attention with learnable positional embedding and a mapping network consisting of 2 layers of MLP. We do not apply any augmentation for the quantitative experiments, except for the results presented in Table 6 and Table 2 of the main paper. For the CAD results, we used a CLIP ViT-B/32 model. **Sketch2Model.** The authors of Sketch2Model released ShapeNet-Synthetic as the training dataset [79], which consists of synthetic sketches of objects from 13 categories from ShapeNet. These objects have been rendered from 20 different views. For training Sketch2Model, we used the official implementation provided in [2], along with the recommended hyperparameters. This implementation uses a step-type learning rate policy, beginning from $1e - 4$ and

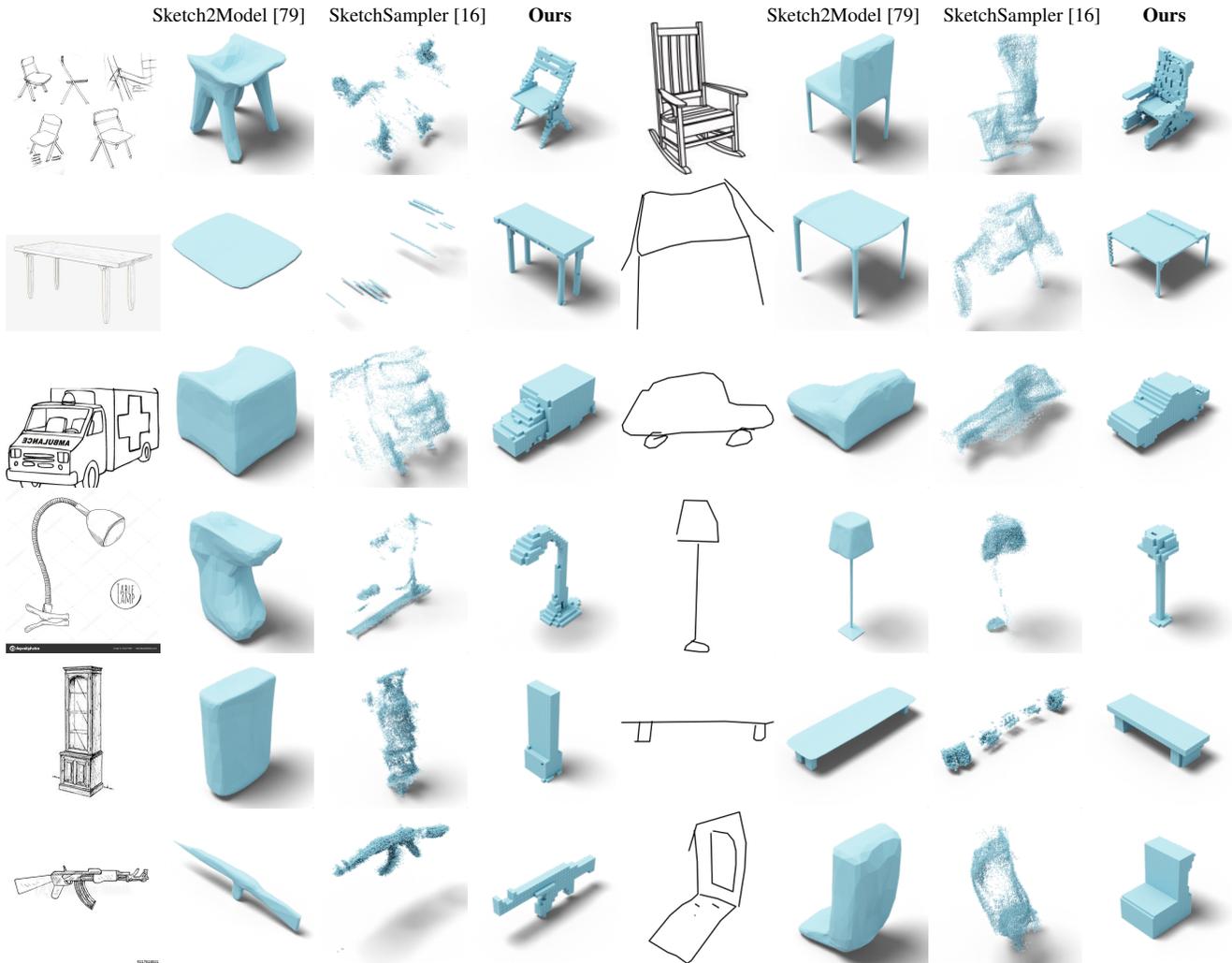


Figure 6: Generated shapes with different methods. We can see that, compared to our method, Sketch2Model meshes lack geometric details and appear smoothed, while SketchSampler pointclouds present artifacts and outliers.

decreasing by 0.3 every 800 epochs, and trains for 2000 epochs with the Adam optimizer. We trained the model on all 13 categories of ShapeNet-Synthetic using the same training/test split of the original paper.

SketchSampler. This method employs as training dataset Synthetic-LineDrawing [16], a paired sketch-3D dataset based on 3D models from ShapeNet. In our experiments, we used the official implementation, cited in the original paper [16]. In particular, we used the pre-trained model released by the authors, and pre-processed the input sketches to be in the same format of Synthetic-LineDrawing ones.

D. Comparison with Point-E

Furthermore, we conducted a comparison between our work and Point-E [51], as illustrated in the table provided below (Row 1). The results clearly demonstrate the supe-



Figure 7: Results of our method on natural images

rior performance of our method, indicating the merit of our design choices.

Method	QD-Acc \uparrow	TU-Acc \uparrow	SS-Acc \uparrow	IS-Acc \uparrow
Point-E	12.6	40.1	43.2	18.9
Ours (CLIP)	58.8	81.5	79.7	74.2
Ours (DINOv2)	39.7	71.1	72.1	55.9

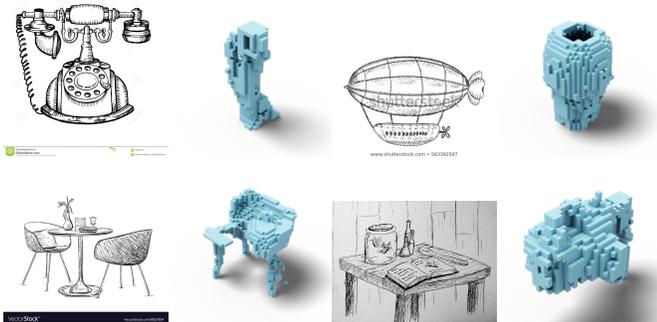


Figure 8: Failure cases of our method

E. Natural Images Results

We explored the applicability of our method to natural images, as it is robust to domain shifts between renderings and sketches. The outcomes are depicted in Figure 7, indicating the efficacy of our method in generating natural images, including those with backgrounds. We believe that this discovery would be of interest to the Single View Reconstruction community.

F. Failure Cases

This section demonstrates the limitations of our method, as illustrated in Figure 8. The outcomes reveal that our method encounters difficulties in generalizing to shapes that are beyond those present in ShapeNet13, as depicted in the first row. Furthermore, our method also faces challenges when dealing with sketches that depict multiple shapes, as shown in the second row. Lastly, our method experiences difficulties in accurately reproducing the local details of shapes, which we consider to be an intriguing direction for future work.

G. Societal Impact

The societal impact of Sketch-to-3D technology can be significant in various fields such as architecture, product design, gaming, and entertainment. With the help of Sketch-to-3D technology, designers and architects can create realistic 3D models quickly and efficiently, reducing the overall time and cost of the design process. However, it is important to note that the widespread adoption of Sketch-to-3D technology could also lead to job displacement in certain industries. As with any technological advancement, it is crucial to consider the potential social and economic impacts and work towards ensuring a smooth transition for workers and

communities affected by such changes.

H. Future Work

We aim to concentrate on expanding this method to handle bigger 3D datasets for our future work. Additionally, we think that enhancing the Stage 1 VQ-VAE can aid in preserving the local features of the 3D shape. Lastly, an intriguing avenue to explore would be to combine sketch with text conditioning, resulting in a more adaptable generative model.

I. More Qualitative Results

Additional results are provided in Figure 9 and Figure 10.

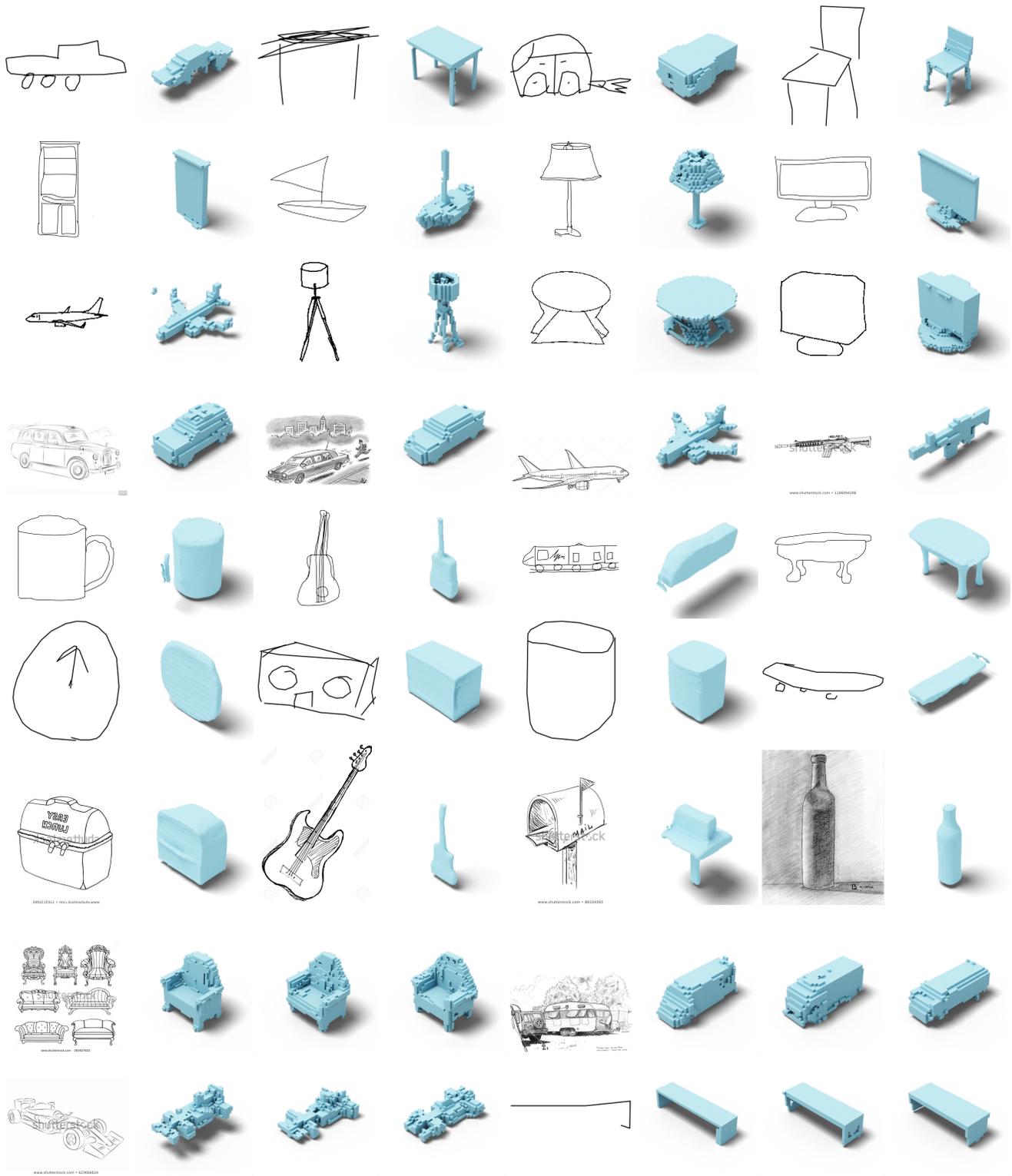


Figure 9: Additional outcomes of our 32^3 voxel model on ShapeNet13 and implicit model on ShapeNet55 are presented. The initial row displays the outcomes of 32^3 voxel model from QuickDraw, the second row from TU Berlin, the third row from ShapeNet Sketch, and the fourth row from Imagenet Sketch. The fifth row displays the outcomes of implicit model from QuickDraw, the sixth row from TU Berlin and the seventh row from Imagenet Sketch. The last two rows exhibit several results obtained from a single sketch.

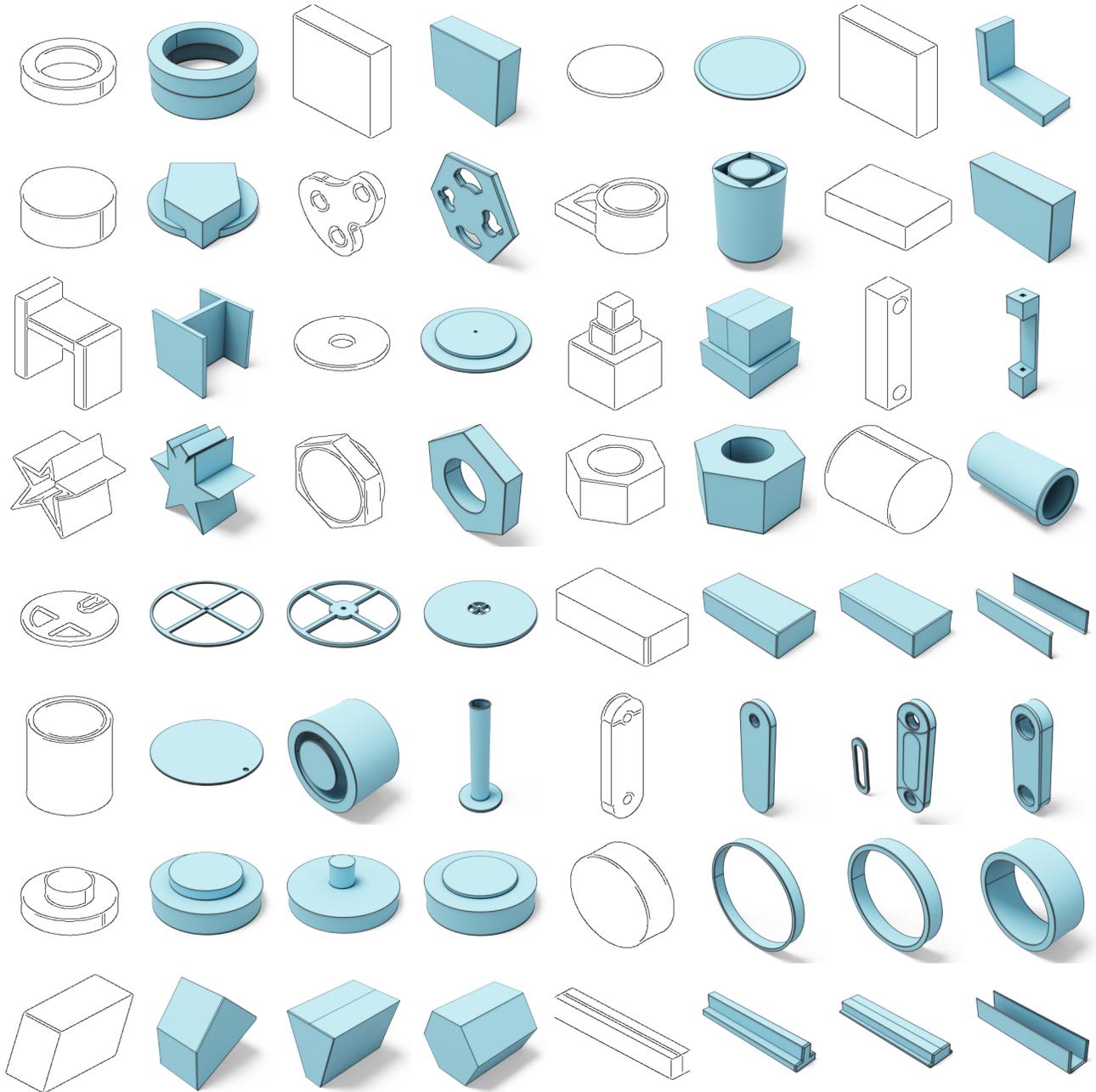


Figure 10: More results from our sketch-CAD model. The last four rows show multiple generation results for each sketch.