

Sculpt3D: Multi-View Consistent Text-to-3D Generation with Sparse 3D Prior

Cheng Chen^{1,2}, Xiaofeng Yang¹, Fan Yang¹, Chengzeng Feng¹, Zhoujie Fu¹,
Chuan-Sheng Foo^{2,3}, Guosheng Lin¹, Fayao Liu²

¹Nanyang Technological University

²Institute for Infocomm Research A*STAR, Singapore

³Centre for Frontier AI Research, A*STAR, Singapore

cheng021@fudan.edu.cn, gslin@ntu.edu.sg fayao.liu@gmail.com

Abstract

Recent works on text-to-3d generation show that using only 2D diffusion supervision for 3D generation tends to produce results with inconsistent appearances (e.g., faces on the back view) and inaccurate shapes (e.g., animals with extra legs). Existing methods mainly address this issue by retraining diffusion models with images rendered from 3D data to ensure multi-view consistency while struggling to balance 2D generation quality with 3D consistency. In this paper, we present a new framework Sculpt3D that equips the current pipeline with explicit injection of 3D priors from retrieved reference objects without re-training the 2D diffusion model. Specifically, we demonstrate that high-quality and diverse 3D geometry can be guaranteed by keypoints supervision through a sparse ray sampling approach. Moreover, to ensure accurate appearances of different views, we further modulate the output of the 2D diffusion model to the correct patterns of the template views without altering the generated object’s style. These two decoupled designs effectively harness 3D information from reference objects to generate 3D objects while preserving the generation quality of the 2D diffusion model. Extensive experiments show our method can largely improve the multi-view consistency while retaining fidelity and diversity. Our project page is available at: <https://stellarcheng.github.io/Sculpt3D/>.

1. Introduction

There has been growing research attention towards text-to-3d generation. Compared to image generation, the data available for 3D generation is less in quantity and lower in quality. Thus, many studies [19, 28, 41] have begun to generate 3D objects using 2D text-to-image models [9, 30] as supervision to leverage their strong priors learned from billions of real images.

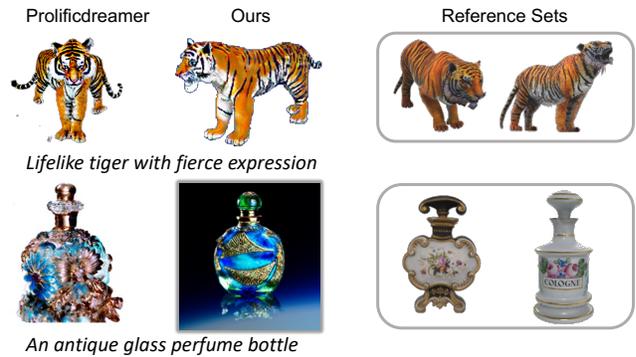


Figure 1. Comparison of objects generated by our method and ProlificDreamer. We retain the 2D model’s capability to produce high-fidelity objects and adaptively learn 3D information from reference templates retrieved from external datasets.

These methods mainly contain two steps: the first step is to continuously sample images from different views of a randomly initialized 3D representation (e.g. NeRF [26], DMTet [33]). The second step uses a 2D diffusion model to individually judge whether each image is a high-quality image that conforms to the text description. Compared to 2D image generation, 3D generation not only requires producing high-quality images for each individual viewpoint but also needs to create plausible shapes and appearances as a whole 3D object. Thus, a high-quality 2D generative model and a mechanism that can accurately provide 3D priors are two keys to achieving decent 3D generation results. Since early works [19, 28, 41] mainly use the sole 2D diffusion model as supervision, they tend to produce inaccurate shapes (shape ambiguity) and appearances that are inconsistent across viewpoints (appearance ambiguity), as shown in Figure 1 left, where examples include incomplete bottles, tigers with multiple legs, and tails.

Recently, some efforts have been made to expand the 3D datasets [10]. Following this, there were immediate attempts to retrain 2D diffusion models on these 3D datasets



Figure 2. Our methods can generate high-fidelity objects with decent shapes using various text prompts. The model adaptively incorporates information from the reference shape displayed on the left, resulting in the creation of objects that range from moderately resembling to substantially diverging from the reference shape. Please find more video results in the supplementary materials.



Figure 3. As shown in the first row, our method can generate diverse 3D objects given the same reference shape. The second row also shows the diverse results generated by randomly selecting reference objects from the top five retrieved samples. All templates are marked as gray and shown in the corner.

to learn 3D information [21, 22, 35]. Although these methods have made impressive progress, they require expensive training costs to re-train the large-scale models, and training 2D diffusion models on rendered images often degrades the

model’s generation quality [21, 35] learned on large real image dataset. Similar situations have arisen in the NLP field. As language models grow huger, it becomes increasingly difficult to inject new information by retraining the models,

thus researchers start to explicitly introduce external knowledge through retrieval augmentation [13, 36]. Motivated by these developments, we design a retrieval mechanism to explicitly supervise the 3D geometry and appearance using retrieved templates without re-training the 2D diffusion model on rendered 3D data.

Explicitly constraining the geometry [37, 39] presents an intrinsic challenge: strong constraints on the shape may make the generated results closely resemble the template, while too lax constraints may fail to ensure a reasonable shape [38]. To adaptively learn the 3D shape information from the template, we exploit the geometric creative capabilities of the 2D diffusion model during volume rendering to enable creative point growth and pruning during the optimization process. Specifically, we design a sparse ray sampling method to selectively discard points, supervising only a minimal number of keypoints that can describe the overall structure, thereby greatly enhancing the 2D diffusion model’s freedom in imaginative shape generation. Moreover, we update the template by pruning and generating new points in areas of low and high NeRF output density, respectively, guided by the diffusion model’s confidence. Since we directly supervise the NeRF without making modifications to the 2D diffusion models, our method can fully preserve the generative quality of the diffusion models while ensuring a decent 3D shape. The generated examples showcased in Figure 2 demonstrate that our method is capable of generating photo-realistic objects that adapt to the template shape, with the diffusion model determining the degree of similarity to the template. In cases where users desire results significantly different from the initial template, we further devised a re-retrieval mechanism that corrects the retrieval results through the generated shape to make full use of the external 3D dataset.

The aforementioned design enables our model to generate diverse and accurate 3D objects in most cases. However, we also observed that the diffusion model may still generate appearances that are inconsistent across views despite the shape being accurate. For instance, it may produce the appearance of an animal’s face at the back or side view, even when the geometry of the face is not generated there. Thus, we further utilize the template’s appearance information to refine the generated objects.

Considering the appearance of generated objects often differs from the template, our challenge here is to correct only the inaccurate aspects of the object’s appearance without altering its style and geometry. Fortunately, recent advances in image controlling [27, 43, 45] enable users to easily modify various attributes of an image, such as style, content, and geometry, in a decoupled manner by training lightweight image adapters. Given the fact that the template always provides accurate guidance on view-specific patterns, like which view the eyes and nose should appear.

We utilize a unified image adapter to first adapt the template to the generated object’s style and then use the adapted image to align the generated erroneous appearances with the correct patterns. As we only modulate the generated patterns for each view without limiting the generated structure, our method only requires four sparse template views to supervise the 3D space partitioned according to four standard orientations. To summarize, our key contributions are:

- We introduce Sculpt3D which explicitly integrates 3D shape and appearance information for multi-view consistent text-to-3d generation while maintaining the high-quality generation capabilities of the 2D diffusion model.
- We enable creative point growth and pruning during the 2D diffusion and 3D geometry co-supervision process, which hones 2D diffusion’s ability to produce shapes that are both accurate and creative. We further use the appearance pattern information of the template to modulate the output of the diffusion model for resolving appearance ambiguities.
- Extensive experiments show that our method is able to significantly improve the multi-view consistency of text-to-3d generation while retaining generalizability.

2. Related work

Large Scale Text-to-Image Diffusion Model. With the tremendous progress in large-scale generative models, a surge of methods [9, 30] have been proposed to perform various types of text-to-image Generation and Editing. To further enhance the generative capabilities of large models, various methods [16, 31] have been proposed to integrate external control signals into these models. ControlNet [45] fine-tunes the Stable Diffusion [30] models to enable more conditional inputs like edge maps, segmentation maps, keypoints, etc. Similar to ControlNet, T2I-Adapter [27] and IPAdapter [43] introduce lightweight adapters for different conditions, providing additional conditional control and supporting the simultaneous use of multiple conditions for one generation. Using external knowledge to augment models has recently drawn attention in both NLP and visual models [2–4, 6, 17]. In image synthesis, Re-Imagen [6] retrieves semantic neighbors to improve the grounding of the diffusion models to real-world knowledge. RDM [2] empowers smaller models with external memory to achieve high-fidelity image generation results. Inspired by these approaches, we utilize template appearances as references to modulate the diffusion process, ensuring the generated images align with the intended viewpoints.

Learning 3D from 2D Diffusion Prior. Pioneering works like Dreamfusion [28] and SJC [40] demonstrate the possibility of supervising NeRF to generate 3D objects using only 2D diffusion. Although their advancements are groundbreaking, the results they produced are somewhat blurry. Subsequent researchers approach the challenge from

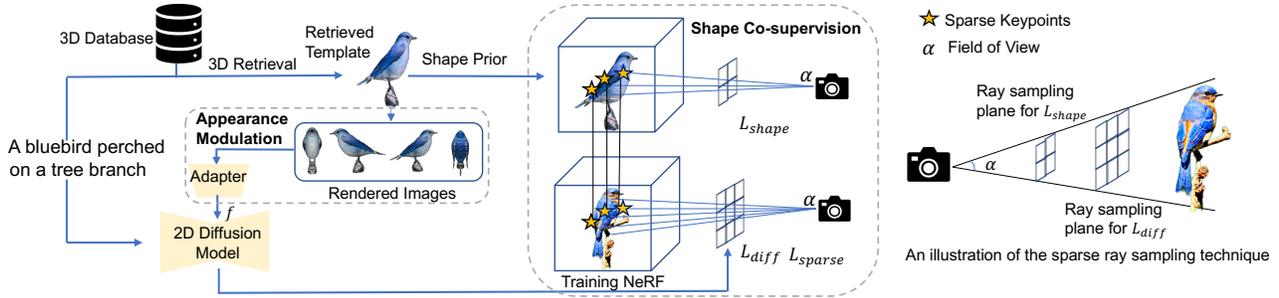


Figure 4. Given a text prompt, we retrieve the most semantically matching samples from an external 3D database. With the retrieved object, we sparsely select the keypoints of the reference shape to co-supervise the geometry with 2D diffusion model. The appearance of the reference object is also used to modulate the 2D diffusion to avoid appearance ambiguity.

various perspectives [7, 8, 42, 44]. Specifically, Magic3D [19] improves both the speed and quality by introducing DMTet [33]. Latent-NeRF [25] seeks to optimize NeRF from an implicit space perspective. Fantasia3D [5] separates geometry and texture modeling to better learn the details of 3D objects. Prolificdreamer [41] introduces the VSD loss to learn the variational distribution of 3D scenes, significantly improving the generation quality. There’s also a line of works focused on image-to-3D generation. For instance, several works [11, 24] use image-conditioned diffusion as a prior to enhancing the generation of unseen viewpoints.

With the release of the large 3D dataset [10], recent methods [22, 23, 34, 35] have attempted to fine-tune 2D diffusion with 3D data. Among them, Zero 1-2-3 [21] introduces camera parameters as conditions to predict images from arbitrary angles relative to the input image. MVDream [35] proposes 3D self-attention to further enhance the generation. Syncdreamer [22] synchronizes the multiview diffusion model to produce multiple new viewpoint images simultaneously.

Different from previous works, our Sculpt3D explicitly explores the 3D priors from reference samples to enhance both the generated shape and appearance without retraining the diffusion model.

3. Approach

As shown in Figure 4, our method uses retrieved templates to provide shape and appearance priors for shape co-supervision and appearance modulation. The details of these components will be given in the following sections.

3.1. Revisiting 2D Diffusion for 3D Generation

Dreamfusion [28] introduces a Score Distillation Sampling (SDS) loss to perform text-to-3d generation. The loss is designed for distilling knowledge from 2D diffusion models to train a 3D representation. Specifically, given a NeRF model $g(\theta)$ which can produce image x at arbitrary camera poses, SDS provides the gradient direction to update θ such that all rendered images are pushed to the high probability

density regions conditioned on the text embedding y under the diffusion prior. The SDS computes the gradient as:

$$\nabla_{\theta} \mathcal{L}_{\text{SDS}}(\phi, x = g(\theta)) = \mathbb{E}_{t, \epsilon} [w(t) (\epsilon_{\phi}(z_t; y, t) - \epsilon) \frac{\partial x}{\partial \theta}], \quad (1)$$

where $w(t)$ is a weighting function, z_t is the noised latent of image x at timestep t , and ϵ_{ϕ} is the denoising network of Stable Diffusion. As aforementioned, while the SDS loss can effectively train the NeRF model, its generated outputs often suffer from oversaturation and are lacking in detail.

To address these issues, Prolificdreamer [41] introduces the VSD Loss. The VSD Loss incorporates the LoRA [16] model to further fit the variational distribution of the 3D scene produced by the training NeRF. It then computes the difference between a pre-trained diffusion model and the LoRA model to guide the NeRF, which is formulated as:

$$\nabla_{\theta} \mathcal{L}_{\text{VSD}}(\theta) \triangleq \mathbb{E}_{t, \epsilon, c} [w(t) \cdot (\epsilon_{\text{pretrain}}(x_t, t, y) - \epsilon_{\phi}(x_t, t, c, y)) \cdot \frac{\partial g(\theta, c)}{\partial \theta}]. \quad (2)$$

In the formula, ϵ_{ϕ} represents the score of a noisy rendered image predicted by the LoRA model, and c is the camera parameter corresponding to the rendered view. We recommend readers refer to Prolificdreamer’s [41] original paper for more details. The VSD loss can effectively improve the fidelity of the generated samples, thus we use it as the 2D diffusion prior by default.

In our experiments, we found that although the VSD loss is able to produce detailed results, its outputs still suffer from inaccurate shapes and appearances. To address these challenges, we next introduce our method of equipping the current pipeline with retrieval capability to explicitly inject 3D priors in the following sections.

3.2. 3D and 2D Co-supervised 3D Generation

Based on previous observations, we now turn to illustrate how to use 3D prior when doing text-guided 3D generation. In our setup, 3D priors can be obtained either by user input

or retrieved from external datasets. Recent advances in representation learning suggest that by scaling up 3D representations, it is accessible to align the CLIP [29] space with 3D data, thereby enabling the retrieval of the most semantically matching objects in a 3D database using natural language. In the experiment, we use the recently released OpenShape [20] model which scales up the 3D backbone to align with CLIP as our 3D retrieval module. NeRF is chosen as our 3D representation due to its flexibility in modifying prior shapes.

To inject the 3D prior, we initially used the 3D template shape to directly initialize the volume density of NeRF. Specifically, inspired by Latent-NeRF [25], we constrain the density of each point during NeRF training. Here, the density label of each point is calculated from the winding number [1] of the normalized template. If the winding numbers show that a point is inside the 3D template shape, we set the density label of that point to 1. Conversely, the point’s label is set to 0. After obtaining an accurate initialized shape, we continue training NeRF using 2D diffusion as supervision. However, we find that the 2D diffusion tends to destroy the initial object shape and converge to a distorted shape. Similar results are also observed by [14, 32]. They find that continuing to modify a well-trained NeRF using either SDS or VSD loss will destroy the initially well-learned 3D representation. In order to effectively generate the correct geometry, we supervise NeRF using both 2D diffusion and 3D shapes. Which is formulated as:

$$L_{co} = L_{diff} + \lambda L_{shape}. \quad (3)$$

As aforementioned, too tight supervision on the shape will make the generated result too similar to the template, and sometimes it even produces an incorrect appearance due to discrepancies between the diffusion prior and the template shape prior. To effectively use 3D prior, we next introduce our shape learning method.

3.2.1 3D Prior Guided Shape Learning

To allow the diffusion model to adaptively learn the 3D prior, we introduce a sparse ray sampling technique to selectively supervise a small number of keypoints that roughly describe the object’s shape. Specifically, every time when randomly sampling a view to train NeRF, 2D diffusion is utilized to supervise all rays to learn the correct RGB and density of each point in the 3D space. At the same time, we maintain the field of view (FOV) unchanged and proportionally reduce the width and height of the ray sampling plane by a factor of N for shape supervision. In this way, as shown in Figure 4, the sampled rays are much sparser, roughly depicting the 3D object shape and providing direct shape guidance. Since the shape constraint only provides a correct sparse prior, diffusion can freely unleash its gen-

erative capabilities in the unconstrained space. The shape supervision loss is defined as a binary cross-entropy loss:

$$L_{shape} = -\frac{1}{|\mathcal{R}|} \sum_{o \in \mathcal{R}} [s_o \log d_o + (1 - s_o) \log(1 - d_o)], \quad (4)$$

where \mathcal{R} denotes the set of keypoints, s_o denotes the density of the keypoint o , and d_o is the NeRF output density.

Considering the co-supervision of 3D shapes and 2D diffusion, two types of conflicts may arise: 2D diffusion might tend to either prune certain points existing in the template or generate points that are not present in the template. To leverage the creativity of the diffusion model to drive the model’s generation, we default the shape supervision scale λ in equation 3 to 0.1. With this configuration, as the diffusion loss scale is larger, points prone to pruning by the diffusion model will have their density optimized towards 0. Conversely, points inclined to be generated will be optimized towards 1. To accelerate the removal of these unnecessary points and the growth of new ones, we further impose a sparsity loss to enforce the generated points’ density to be either zero or one, and the points with a density of zero will be pruned and no longer supervised. As shown in the results Figure 2, our technique can effectively remove unwanted points and generate new points to create new shapes. The sparsity loss is defined as follows:

$$L_{sparse} = \frac{1}{|\mathcal{T}|} \sum_{o \in \mathcal{T}} [\log(d_o) + \log(1 - d_o)], \quad (5)$$

where \mathcal{T} denotes the set of all points. The final loss we used is $L = L_{co} + L_{sparse}$.

The aforementioned design can effectively assist the model in generating new shapes. When the user wishes to significantly increase the downsampling factor N to create objects that differ greatly from the template, we further design a re-retrieval mechanism. Specifically, we extract the initially generated shape representation and use it to retrieve matching shapes in the top 100 objects retrieved by text. This allows for further utilization of 3D datasets to find the reference shape that best matches the structure generated by diffusion.

3.2.2 3D Appearance Modulated 2D Diffusion Prior

Couple shape guidance with 2D diffusion prior can effectively help the model to correctly understand the 3D world, thereby producing correct generation results. However, in our experiments, we also find that the model still cannot infer the correct appearance even when the shape is entirely accurate in some hard cases. To explicitly guide the model to generate the correct appearance for each view, we design an optional technique that uses the appearance of the template as a semantic reference to modulate the diffusion process in hard cases.

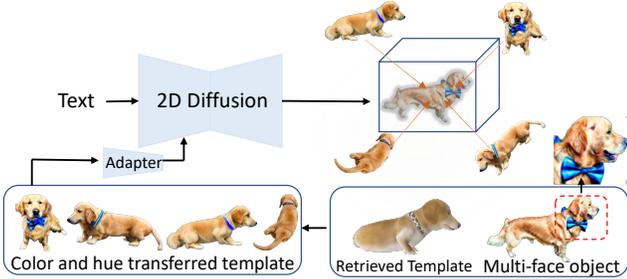


Figure 5. Illustration of the appearance modulation. Four canonical views of the templates are transferred to the generated object’s style to modulate the 2D diffusion.

Specifically, recent works [27, 43] demonstrate that various external signals can be applied to control the output of diffusion by training a lightweight adapter. The adapter can utilize images as prompts to generate results with semantic patterns similar to the reference image. In order to correct the appearance of the generated object without affecting its overall style, we first utilize the adapter to convert the template objects to match the hue and color distribution of the generated one. This can be simply achieved by using the color distribution of the generated object as a condition. Since the converted template view contains accurate view-specific patterns, it is used as the image prompt together with the text to align the diffusion generation results with the correct semantics pattern. As LoRA is designed to fit the scene distribution of the trained NeRF in VSD loss, the diffusion prior coupled with the image adapter can be formulated as:

$$\nabla_{\theta} \mathcal{L}_{diff}(\theta) \triangleq \mathbb{E}_{t, \varepsilon, c} \left[\omega(t) \cdot \left(\epsilon_{\text{pretrain}}(x_t, f, t, y) - \epsilon_{\phi}(x_t, t, c, y) \right) \cdot \frac{\partial g(\theta, c)}{\partial \theta} \right], \quad (6)$$

where f denotes the image features extracted by the adapter. As shown in figure 5, since the semantic pattern of the image is constant within a certain observation range, our method only requires 4 sparse template images corresponding to 4 canonical view spaces.

4. Experiments

To comprehensively assess the effectiveness of our method, we use the text descriptions provided by T3Bench [15] for testing, which contains 100 text prompts covering various types of single objects. Additionally, we use ChatGPT to generate 40 different prompts for testing, including both common everyday objects and some imaginative objects.

4.1. Implementation Details

Our method is built on the implementation from Threestudio [12]. All experiments are conducted on an NVIDIA

A6000 GPU. The model used for 3D retrieval is OpenShape [20], an open-world retrieval model trained using multiple ensemble datasets. By default, the shape ranked first in the retrieval results is used as the reference shape in the experiments. The scale of constraint on shape λ is consistently set to 0.1. The sparse keypoints selection factor N is set to 8, meaning 2D diffusion supervises 3D points on rays sampled from a 512×512 space, while geometry supervision is applied to 3D points on rays sampled from a 64×64 image space. The initial shape is obtained at the 5000 training step when performing shape retrieval. The version of diffusion used in the experiments is Stable Diffusion 1.5. The image adapter used in appearance learning is the publicly available pre-trained T2I-Adapter [27].

4.2. Results of Sculpt3D

We show the generated results of Sculpt3D in Figure 2, including the generated results and the corresponding reference shapes shown on the left. The results demonstrate that our method can generate objects with accurate geometry using various text descriptions while maintaining the ability of 2D diffusion to produce highly realistic appearances. It can be observed that our method can adaptively learn geometry information from the template. Some generated objects resemble the reference shapes, while others show significant differences. Moreover, Figure 3 also showcases Sculpt3D’s capability to produce diverse results. We illustrate this through two sets of examples: one where the same template is used to generate multiple times and another where various templates are randomly chosen from the top five retrieved results. The results are shown in the first and second rows respectively. It can be seen that even when using the same template, Sculpt3D can produce remarkably different results due to the sufficient creative freedom allowed in the generation process. Additionally, since the external 3D dataset contains different samples that conform to the same semantic description, randomly selecting from the retrieval results also effectively yields diverse generative results.

4.3. Comparison to Baselines

We compare our method with five baselines, DreamFusion [28], Latent-NeRF [25], Magic3D [19], Fantasia3D [5], and ProlificDreamer [41]. We use the implementation from Threestudio [12] for all these baselines. As shown in Figure 6, previous methods struggle to generate shapes with reasonable geometry and high-quality appearances, while our method is capable of simultaneously producing objects with good geometry, higher fidelity, and more details.

4.4. Quantitative Evaluation

To quantitatively evaluate the text-to-3d method, T3Bench [15] designs two metrics based on multi-view CLIP score

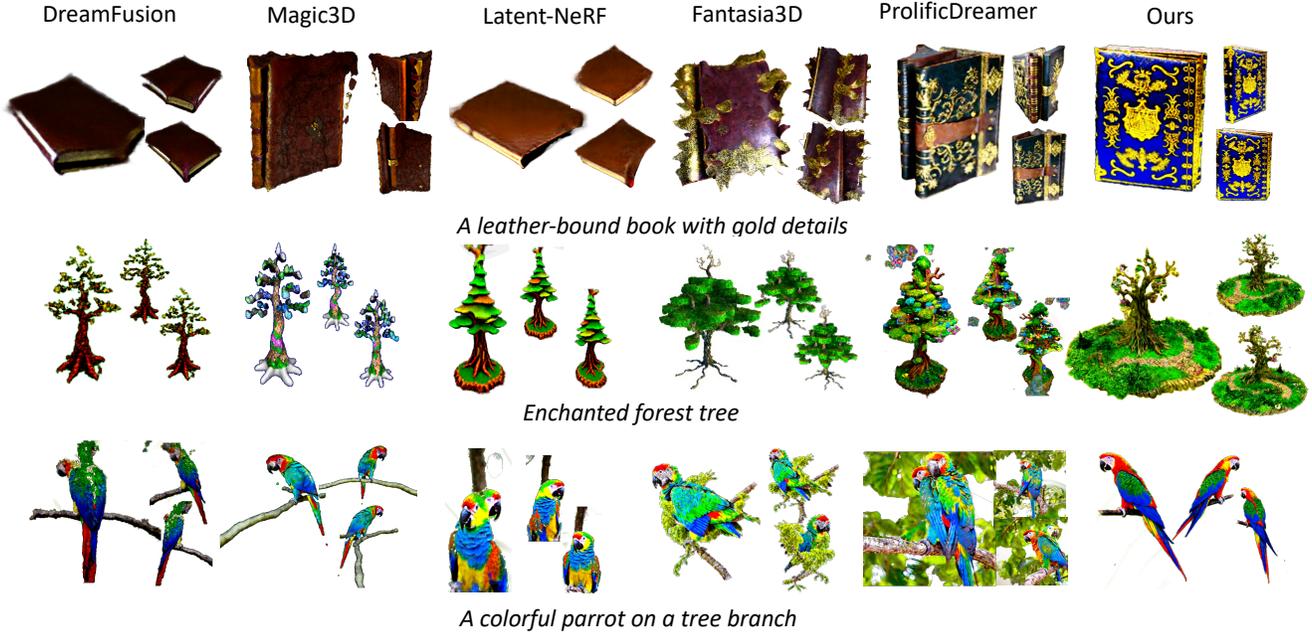


Figure 6. Comparison with baselines. Our method can generate objects with decent shapes, which not only have high fidelity and rich details but also maintain 3D consistency.

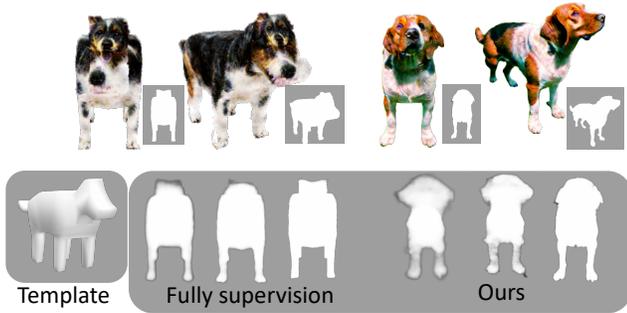


Figure 7. Comparisons with fully shape supervision, the density changes during training are depicted in the gray box below. It is observable that full supervision tends to result in the generated shapes closely resembling the template, thereby leading to a loss in diversity. Moreover, the discrepancy between 2D priors and 3D shape priors could potentially result in inaccuracies in the appearance of the shapes.

and GPT-4 evaluation to assess the generated object’s quality and alignment using 100 prompts. As our work focuses

Table 1. Quantitative comparisons with the baselines.

Methods	Quality	Alignment	Cons. Rate
Dreamfusion	24.9	24.0	34%
Latent-NeRF	34.2	32.0	30%
Magic3D	38.7	35.3	38%
Fantasia3D	29.2	23.5	26%
ProlificDreamer	51.1	47.8	32%
Ours	53.6	49.3	76%



Figure 8. We also showcase the generation results when no matching samples are retrieved. Even though the retrieval model failed to find samples that fully match the semantics, our method is still capable of effectively absorbing useful information from the template to produce correct results.

on generating 3D content with multi-view consistency, we further follow previous work [18] to evaluate 3D consistent rate. Specifically, we randomly select 50 prompts from T3Bench, and manually identify and count 3D inconsistencies (e.g., multiple faces, legs, and other distorted shapes.) of each method. The consistent rate is then determined by dividing the number of 3D consistent objects by the total generated results. As shown in Table 1, our method significantly improves the multi-view consistency rate while sur-

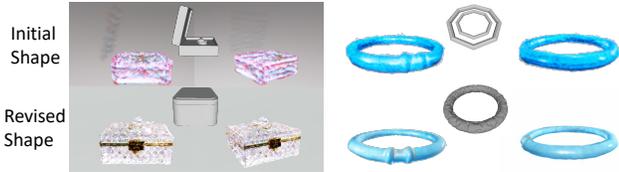


Figure 9. Illustration of the effectiveness of re-retrieval using generated shapes. When users desire results significantly different from the initial template, we can utilize the outputs from the rough generation stage, as shown in the first row, to re-retrieve more accurate reference shapes, demonstrated in the second row.

passing the baseline in both quality and alignment metrics.

4.5. Ablation Study

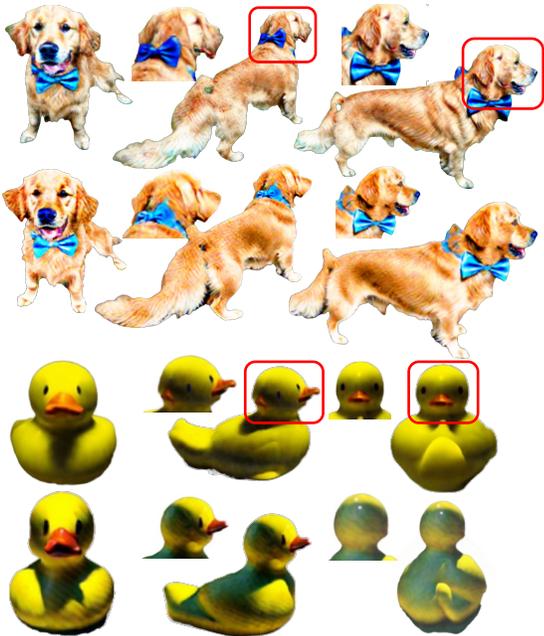


Figure 10. Multi-view inconsistency (appearance) issue. Each column shows a single view of the object. **1st, 3rd rows:** Despite the accuracy of the shape, there is a potential for ambiguity in appearance (highlighted in the red box) that may still arise. **2nd, 4th rows:** Our appearance modulation method can effectively correct this type of appearance confusion.

Effectiveness of Shape Learning. As shown in Figure 7, we compare the progression of NeRF density changes and the final generated results of our method with the fully shape supervision method proposed in Latent-NeRF [30], which is designed to make minor modifications to a given shape surface. The result shows that when applying their constraining coupled with VSD on the high-resolution NeRF, the model converges quickly to the reference shape, excessively limiting the diffusion model’s creativity and resulting in an unnatural shape. Furthermore, the gap be-

tween the diffusion model’s prior and the template’s geometry may cause the model to generate incorrect appearances in inappropriate locations (e.g., multiple dog faces). In cases where the reference shapes are not accurate, it becomes increasingly crucial to grant the diffusion model a suitable level of creative freedom. In Figure 8, we show the generated results when the mismatching templates are retrieved. It is evident that our model can still effectively generate correct results by adding and pruning points based on the inaccurate template shape. We also demonstrate situations where further reducing the supervised points. Figure 9 illustrates the scenario when the sparse selection factor is set to 16. In this scenario, the model may make significant modifications to the initial reference shape, such as removing the reference shape’s lid or combining two rings of a bracelet into one. Our re-retrieval method still can effectively utilize the initial generated shape to retrieve the most matching sample from the candidate set, thus accelerating the creation of the final high-quality object.

Effectiveness of Appearance Modulation. Though the correct shape is guaranteed by shape supervision, some challenging cases also show appearance ambiguities. As shown in the first row of Figure 10, our baseline model can generate the correct shape of a yellow rubber duck, but behind its head, despite the absence of corresponding shapes for the mouth and eyes, it still generates the appearance of a duck’s face at the back view. By using the template with the correct pattern to refine erroneous appearances, our method can effectively adapt the generated objects to the correct appearance without changing their overall style. In the case of the golden retriever, which presents a difficult pose, our baseline model initially generates an extra face on the side by mistake. Through refinement, we are able to effectively correct the appearance mismatch for this challenging pose.

5. Conclusion & Limitation

Conclusion: In this paper, we propose Sculpt3D which explicitly utilizes the 3D shape and appearance information from the retrieved template to aid text-to-3D generation. Sculpt3D is capable of performing creative point growth and pruning within a framework of sparse geometry constraints, thus enabling flexible and accurate shape generation. Moreover, we use the correct pattern information from the template’s appearance to fix ambiguities in the generated object’s appearances without changing their style. Sculpt3D enhances multi-view consistency in a manner that explicitly supervises the 3D representation, thereby preserving the generative capabilities of 2D diffusion. Experiments on text-to-3d benchmarks show the effectiveness of the proposed model, and more extensive ablation studies further confirm the generalizability of our method.

Limitation: While our method has shown promising performance, we also note some limitations. Since we ex-

PLICITLY supervise the geometry, it is difficult for our method to correctly generate when the initial retrieved shape exceeds the prior of the 3D dataset. Early in our development, we experimented with generating an initial 3D shape without constraints and then using that shape to retrieve matching reference objects. Unfortunately, this approach often fails to produce reasonable reference objects due to the limitations of existing generation methods.

Acknowledgements. This research work is supported by the Agency for Science, Technology and Research (A*STAR) under its MTC Programmatic Funds (Grant No. M23L7b0021).

References

- [1] Gavin Barill, Neil G Dickson, Ryan Schmidt, David IW Levin, and Alec Jacobson. Fast winding numbers for soups and clouds. *ACM Transactions on Graphics (TOG)*, 37(4): 1–12, 2018. [5](#)
- [2] Andreas Blattmann, Robin Rombach, Kaan Oktay, Jonas Müller, and Björn Ommer. Retrieval-augmented diffusion models. *Advances in Neural Information Processing Systems*, 35:15309–15324, 2022. [3](#)
- [3] Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al. Improving language models by retrieving from trillions of tokens. In *International conference on machine learning*, pages 2206–2240. PMLR, 2022.
- [4] Cheng Chen, Zhenshan Tan, Qingrong Cheng, Xin Jiang, Qun Liu, Yudong Zhu, and Xiaodong Gu. Utc: A unified transformer with inter-task contrastive learning for visual dialog. In *Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition*, pages 18103–18112, 2022. [3](#)
- [5] Rui Chen, Yongwei Chen, Ningxin Jiao, and Kui Jia. Fantasia3d: Disentangling geometry and appearance for high-quality text-to-3d content creation. *arXiv preprint arXiv:2303.13873*, 2023. [4](#), [6](#), [12](#)
- [6] Wenhui Chen, Hexiang Hu, Chitwan Saharia, and William W Cohen. Re-imagen: Retrieval-augmented text-to-image generator. *arXiv preprint arXiv:2209.14491*, 2022. [3](#)
- [7] Yiwen Chen, Zilong Chen, Chi Zhang, Feng Wang, Xiaofeng Yang, Yikai Wang, Zhongang Cai, Lei Yang, Huaping Liu, and Guosheng Lin. Gaussianeditor: Swift and controllable 3d editing with gaussian splatting, 2023. [4](#)
- [8] Yiwen Chen, Chi Zhang, Xiaofeng Yang, Zhongang Cai, Gang Yu, Lei Yang, and Guosheng Lin. It3d: Improved text-to-3d generation with explicit view synthesis. *arXiv preprint arXiv:2308.11473*, 2023. [4](#)
- [9] Florinel-Alin Croitoru, Vlad Hondru, Radu Tudor Ionescu, and Mubarak Shah. Diffusion models in vision: A survey. *T-PAMI*, 2023. [1](#), [3](#)
- [10] Matt Deitke, Ruoshi Liu, Matthew Wallingford, Huong Ngo, Oscar Michel, Aditya Kusupati, Alan Fan, Christian Laforte, Vikram Voleti, Samir Yitzhak Gadre, et al. Objaverse-xl: A universe of 10m+ 3d objects. *arXiv preprint arXiv:2307.05663*, 2023. [1](#), [4](#), [11](#)
- [11] Congyue Deng, Chiyu Jiang, Charles R Qi, Xinchun Yan, Yin Zhou, Leonidas Guibas, Dragomir Anguelov, et al. Nerdi: Single-view nerf synthesis with language-guided diffusion as general image priors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20637–20647, 2023. [4](#)
- [12] Yuan-Chen Guo, Ying-Tian Liu, Ruizhi Shao, Christian Laforte, Vikram Voleti, Guan Luo, Chia-Hao Chen, Zi-Xin Zou, Chen Wang, Yan-Pei Cao, and Song-Hai Zhang. threestudio: A unified framework for 3d content generation. <https://github.com/threestudio-project/threestudio>, 2023. [6](#)
- [13] Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. Retrieval augmented language model pre-training. In *International conference on machine learning*, pages 3929–3938. PMLR, 2020. [3](#)
- [14] Ayaan Haque, Matthew Tancik, Alexei A Efros, Aleksander Holynski, and Angjoo Kanazawa. Instruct-nerf2nerf: Editing 3d scenes with instructions. *arXiv preprint arXiv:2303.12789*, 2023. [5](#)
- [15] Yuze He, Yushi Bai, Matthieu Lin, Wang Zhao, Yubin Hu, Jenny Sheng, Ran Yi, Juanzi Li, and Yong-Jin Liu. T3bench: Benchmarking current progress in text-to-3d generation. *arXiv preprint arXiv:2310.02977*, 2023. [6](#)
- [16] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021. [3](#), [4](#)
- [17] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474, 2020. [3](#)
- [18] Weiyu Li, Rui Chen, Xuelin Chen, and Ping Tan. Sweetdreamer: Aligning geometric priors in 2d diffusion for consistent text-to-3d. *arXiv preprint arXiv:2310.02596*, 2023. [7](#)
- [19] 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. In *CVPR*, 2023. [1](#), [4](#), [6](#), [12](#)
- [20] Minghua Liu, Ruoxi Shi, Kaiming Kuang, Yinhao Zhu, Xuanlin Li, Shizhong Han, Hong Cai, Fatih Porikli, and Hao Su. Openshape: Scaling up 3d shape representation towards open-world understanding. *arXiv preprint arXiv:2305.10764*, 2023. [5](#), [6](#)
- [21] Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick. Zero-1-to-3: Zero-shot one image to 3d object. In *ICCV*, 2023. [2](#), [4](#)
- [22] Yuan Liu, Cheng Lin, Zijiao Zeng, Xiaoxiao Long, Lingjie Liu, Taku Komura, and Wenping Wang. Syncdreamer: Generating multiview-consistent images from a single-view image. *arXiv preprint arXiv:2309.03453*, 2023. [2](#), [4](#)
- [23] Xiaoxiao Long, Yuan-Chen Guo, Cheng Lin, Yuan Liu, Zhiyang Dou, Lingjie Liu, Yuexin Ma, Song-Hai Zhang,

- Marc Habermann, Christian Theobalt, et al. Wonder3d: Single image to 3d using cross-domain diffusion. *arXiv preprint arXiv:2310.15008*, 2023. 4
- [24] Luke Melas-Kyriazi, Iro Laina, Christian Rupprecht, and Andrea Vedaldi. Realfusion: 360deg reconstruction of any object from a single image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8446–8455, 2023. 4
- [25] Gal Metzer, Elad Richardson, Or Patashnik, Raja Giryes, and Daniel Cohen-Or. Latent-nerf for shape-guided generation of 3d shapes and textures. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12663–12673, 2023. 4, 5, 6, 12
- [26] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020. 1
- [27] Chong Mou, Xintao Wang, Liangbin Xie, Jian Zhang, Zhonggang Qi, Ying Shan, and Xiaohu Qie. T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models. *arXiv preprint arXiv:2302.08453*, 2023. 3, 6, 11
- [28] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. In *ICLR*, 2023. 1, 3, 4, 6
- [29] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 5
- [30] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 1, 3, 8
- [31] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22500–22510, 2023. 3
- [32] Ruizhi Shao, Jingxiang Sun, Cheng Peng, Zerong Zheng, Boyao Zhou, Hongwen Zhang, and Yebin Liu. Control4d: Dynamic portrait editing by learning 4d gan from 2d diffusion-based editor. *arXiv preprint arXiv:2305.20082*, 2023. 5
- [33] Tianchang Shen, Jun Gao, Kangxue Yin, Ming-Yu Liu, and Sanja Fidler. Deep marching tetrahedra: a hybrid representation for high-resolution 3d shape synthesis. *Advances in Neural Information Processing Systems*, 34:6087–6101, 2021. 1, 4
- [34] Ruoxi Shi, Hansheng Chen, Zhuoyang Zhang, Minghua Liu, Chao Xu, Xinyue Wei, Linghao Chen, Chong Zeng, and Hao Su. Zero123++: a single image to consistent multi-view diffusion base model. *arXiv preprint arXiv:2310.15110*, 2023. 4
- [35] Yichun Shi, Peng Wang, Jiangleong Ye, Mai Long, Kejie Li, and Xiao Yang. Mvdream: Multi-view diffusion for 3d generation. *arXiv preprint arXiv:2308.16512*, 2023. 2, 4
- [36] Kurt Shuster, Spencer Poff, Moya Chen, Douwe Kiela, and Jason Weston. Retrieval augmentation reduces hallucination in conversation. *arXiv preprint arXiv:2104.07567*, 2021. 3
- [37] Chaoyue Song, Jiacheng Wei, Ruibo Li, Fayao Liu, and Guosheng Lin. 3d pose transfer with correspondence learning and mesh refinement. *arXiv preprint arXiv:2109.15025*, 2021. 3
- [38] Chaoyue Song, Tianyi Chen, Yiwen Chen, Jiacheng Wei, Chuan Sheng Foo, Fayao Liu, and Guosheng Lin. Moda: Modeling deformable 3d objects from casual videos, 2023. 3
- [39] Chaoyue Song, Jiacheng Wei, Ruibo Li, Fayao Liu, and Guosheng Lin. Unsupervised 3d pose transfer with cross consistency and dual reconstruction. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–13, 2023. 3
- [40] Haochen Wang, Xiaodan Du, Jiahao Li, Raymond A Yeh, and Greg Shakhnarovich. Score jacobian chaining: Lifting pretrained 2d diffusion models for 3d generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12619–12629, 2023. 3
- [41] Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. Prolificdreamer: High-fidelity and diverse text-to-3d generation with variational score distillation. *arXiv preprint arXiv:2305.16213*, 2023. 1, 4, 6, 12
- [42] Xiaofeng Yang, Yiwen Chen, Cheng Chen, Chi Zhang, Yi Xu, Xulei Yang, Fayao Liu, and Guosheng Lin. Learn to optimize denoising scores for 3d generation: A unified and improved diffusion prior on nerf and 3d gaussian splatting. *arXiv preprint arXiv:2312.04820*, 2023. 4
- [43] Hu Ye, Jun Zhang, Sibio Liu, Xiao Han, and Wei Yang. Ip-adapter: Text compatible image prompt adapter for text-to-image diffusion models. *arXiv preprint arXiv:2308.06721*, 2023. 3, 6
- [44] Chi Zhang, Yiwen Chen, Yijun Fu, Zhenglin Zhou, Gang Yu, Billz Wang, Bin Fu, Tao Chen, Guosheng Lin, and Chunhua Shen. Styleavatar3d: Leveraging image-text diffusion models for high-fidelity 3d avatar generation. *arXiv preprint arXiv:2305.19012*, 2023. 4
- [45] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3836–3847, 2023. 3

Sculpt3D: Multi-View Consistent Text-to-3D Generation with Sparse 3D Prior: Supplementary Material

A. Supplementary Materials

We have prepared supplementary materials, including a document and a video, to provide a more comprehensive understanding of our Sculpt3D. In the document, we discuss the technical details of our implementation in Sec. B. Moreover, we present additional examples and comparisons in Sec. C to demonstrate the performance of our method. Furthermore, we have prepared a video that showcases the results and comparisons of Sculpt3D.

B. Technical Details

Implementation details. This section provides more implementation details of our experiments.

- In our experiments, we observe that most objects in Objectverse [10] are aligned with the observer’s frontal view. Thus we only normalized the vertices and centers of the templates without manually adjusting their poses. In addition to the 100 prompts provided by T3bench, we also use ChatGPT to generate 40 additional prompts, including common objects as well as some unusual and special items. The prompts we used with ChatGPT are as follows. *I am utilizing a text-to-3D model to generate various 3D objects. Please create 40 prompts for me, including both everyday common objects and some unusual and special items.*
- When performing appearance modulation, three adapters are utilized to correct erroneous patterns without altering the style of the generated objects. The adapters we used are a spatial color palette adapter, a structure adapter, and an image adapter. Specifically, since T2I-Adapter [27] supports combining multiple adapters to utilize complementary ability between different adapters, we use sketches extracted from the template objects as structure pattern conditions and the hue and color distribution of the generated object as spatial color conditions. Both the sketch and spatial color palette extraction model are the default models from the T2I-Adapter. This approach effectively retains the pattern information of the template while transferring it to the color distribution of the generated object. When the color distribution of the template is transferred, it is used as the image condition to modulate the diffusion process using the equation 6.

C. More Results

Loss balancing ablation In addition to the ablation studies of the shape learning provided in Sec. 4.5, we conduct another ablation study to discuss the effectiveness of sparse ray sampling, which is described in Sec. 3.2.1. We first remove sparse ray sampling and keep the value of λ in equation 3 as 0.1 to evaluate the effectiveness of sparse ray sampling.

As shown in Figure C.1, the results show that removing sparse ray sampling causes the generated objects to closely resemble the template, due to the geometric constraints being uniformly applied to all points. For example, the folds in the hat closely match those in the template, and the back cover of the water gun doesn’t close. As shown in the third column of Figure C.1, by implementing sparse ray sampling our method can generate imaginative and reasonable geometry under the guidance of the reference shape.



Figure C.1. Ablation on the sparse ray sampling strategy.

For the choice of λ in equation 3, we study the effect of it by applying sparse ray sampling with λ values of 1, 0.1, and 0.01. The results are shown in Figure C.2. It’s evident that even at $\lambda = 1$, our sparse sampling approach is able to provide sufficient flexibility for the model to learn new shapes. Compared to the results with $\lambda = 1$, setting λ as 0.1 can further increase the geometry freedom in the generated results. For instance, the shape of the straw hat is obviously changed. When set λ as 0.01, the model can create significantly new geometries, but it may produce undesirable outcomes. Therefore, we default λ as 0.1 in our experiments.

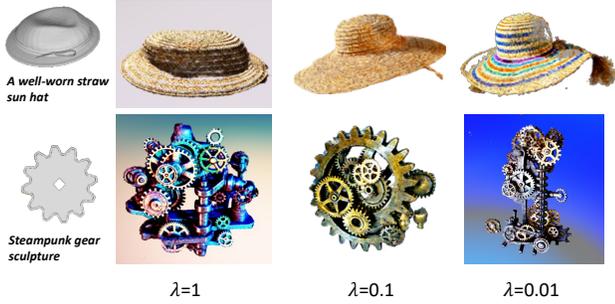


Figure C.2. Ablation on the shape co-supervision value λ in equation 3.

Details of comparison with baselines. To further validate the effectiveness of the sparse prior scheme in sculpt3d, we conduct two additional experiments. We first study the effectiveness of the sketch shape loss proposed by Latent-NeRF [25]. They propose it to allow the model to make slight changes in the template’s surface, the loss is formulated as:

$$L_{\text{Sketch-Shape}} = CE(\alpha_{\text{NeRF}}(p), \alpha_{\text{GT}}(p)) \cdot (1 - e^{-\frac{d^2}{2\sigma_S}}), \quad (7)$$

where α_{NeRF} and α_{GT} are the NeRF occupancy and template shape’s occupancy, respectively. The loss is applied to all points, d represents the distance of a point p from the surface, and σ_S is a hyperparameter that controls how lenient the loss is. A higher σ_S value means a more relaxed constraint to the surface of the generated object. Since their method operates with the SDS loss at a low resolution of 64 rendering, for a more comprehensive comparison, we use their code to conduct experiments in their 64 setting and combine it with the VSD loss to train at a higher resolution of 512 rendering. To fully utilize the new geometry generation capability of their method, we employ the maximum value of σ_S , 1.2, as used by them in all experiments.

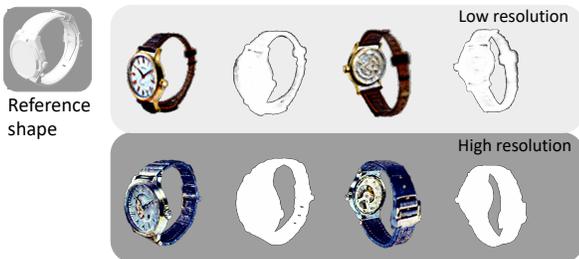


Figure C.3. Ablation of the strategy proposed by [25] in both low and high resolution rendering.

The experimental results are illustrated in Figure C.3, where we showcase the generated outcomes at two different resolutions along with their corresponding occupancy. It is observable that at lower resolutions, their method is

able to alter surfaces, like thinning watch straps. However, at higher resolutions, their approach struggles to change the object’s shape, which results in the generated geometries closely resembling the templates. Additionally, it is noted that their method of relaxing surface constraints often leaves residual artifacts on the surface. This is evident in the occupancy results of the watch straps in the first row, stemming from an incomplete removal of surface density.

Comparison with mesh initiation. In the main text, we mentioned that directly using the template’s shape to initialize NeRF’s density can not guarantee a satisfactory shape. Unlike our approach, Fantasia3D uses a mesh-based DM Tet as a 3D representation, thus supporting initialization with an initial mesh. To more comprehensively verify the role of geometric initialization, we also used our template to initialize Fantasia3D. As shown in Figure C.4, the results show that simple initialization is hard to ensure the subsequent learning direction of the model. Despite the model being initialized by a reasonable shape, it still produces unsatisfactory outcomes, such as the distorted shapes of birds and books. This further underscores the necessity of employing geometry and 2D co-supervision during the learning process.



Figure C.4. Ablation on the mesh initiation strategy.

More comparisons. Here, we showcase more multi-view examples generated by our method in Figure C.5.

Furthermore, we also provide more comparisons with baselines in Figure C.6 and Figure C.7. To compare with the best results demonstrated by the baseline methods, we follow the previous works [5, 19, 41] to directly copy the figures from the corresponding papers for comparisons.

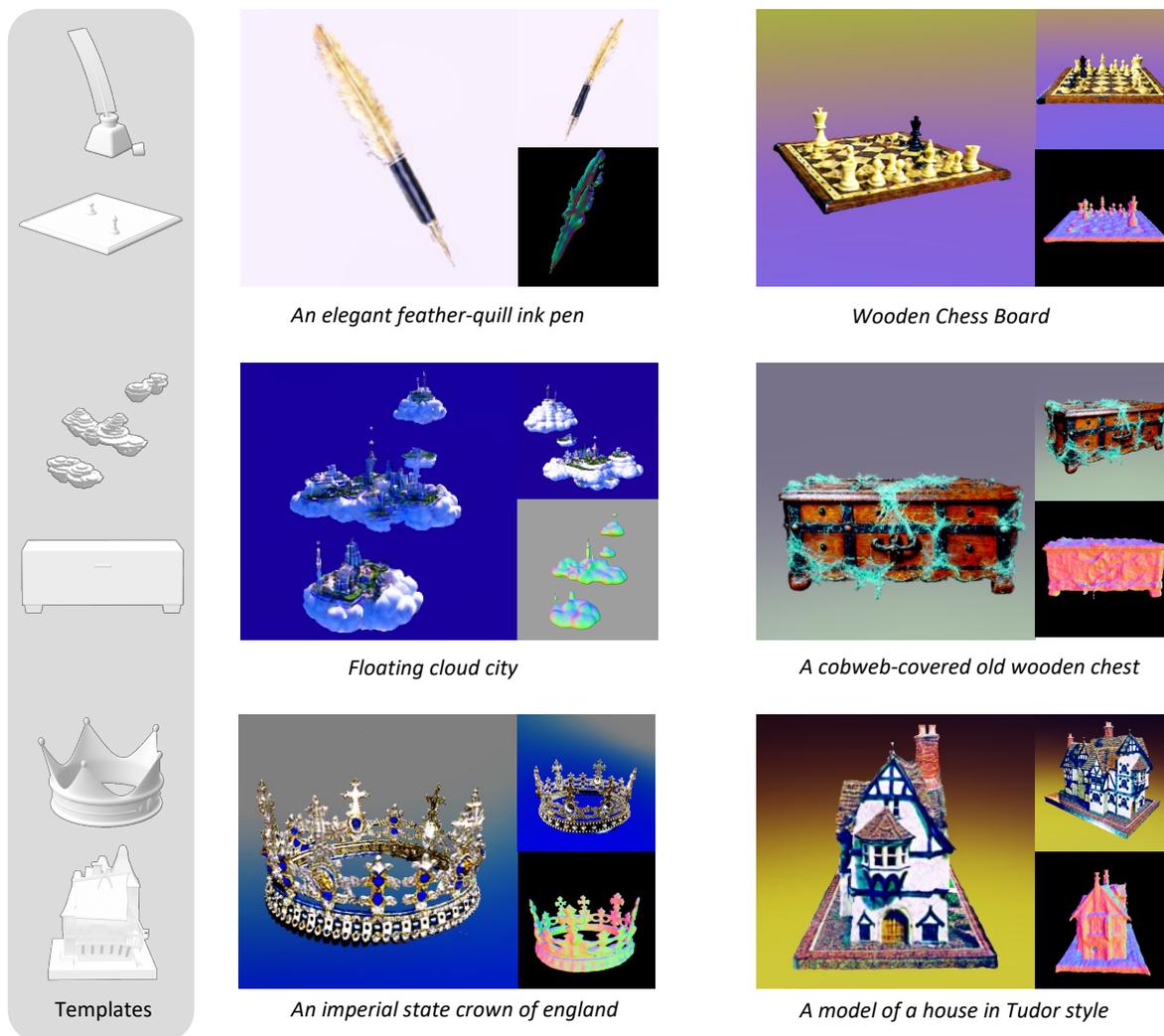


Figure C.5. More multi-view examples generated by our method, the retrieved templates are shown on the left.

A small saguaro cactus planted in a clay pot.



Ours



ProlificDreamer



Magic3D



Dreamfusion

A car made out of sushi.



Ours



ProlificDreamer



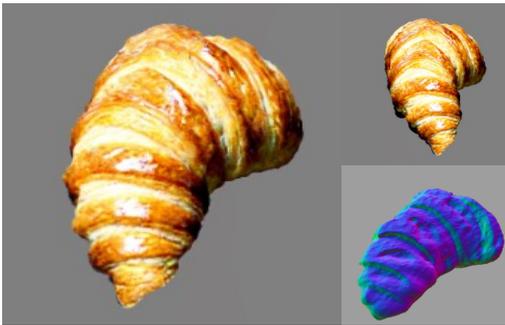
Fantasia3D



Magic3D

Figure C.6. Additional examples for qualitative comparison with baselines.

A delicious croissant.



Ours



ProlificDreamer



Fantasia3D



Dreamfusion

Figure C.7. Additional examples for qualitative comparison with baselines.