

Make-Your-3D: Fast and Consistent Subject-Driven 3D Content Generation

Fangfu Liu, Hanyang Wang, Weiliang Chen, Haowen Sun, and Yueqi Duan*

Tsinghua University

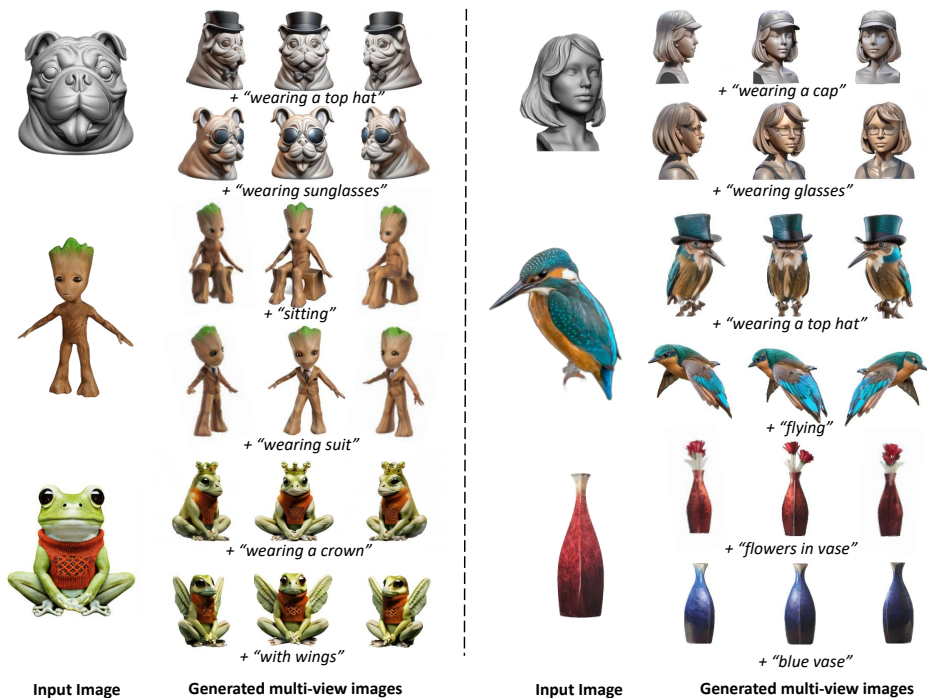


Fig. 1: Make-Your-3D can personalize 3D contents from only a **single image** of a subject with text-driven modifications within only **5 minutes**.

Abstract. Recent years have witnessed the strong power of 3D generation models, which offer a new level of creative flexibility by allowing users to guide the 3D content generation process through a single image or natural language. However, it remains challenging for existing 3D generation methods to create subject-driven 3D content across diverse prompts. In this paper, we introduce a novel 3D customization method, dubbed **Make-Your-3D** that can personalize high-fidelity and consistent 3D content from only a single image of a subject with text description within 5 minutes. Our key insight is to harmonize the distributions of a multi-view diffusion model and an identity-specific 2D generative model, aligning them with the distribution of the desired

* Corresponding author.

3D subject. Specifically, we design a co-evolution framework to reduce the variance of distributions, where each model undergoes a process of learning from the other through identity-aware optimization and subject-prior optimization, respectively. Extensive experiments demonstrate that our method can produce high-quality, consistent, and subject-specific 3D content with text-driven modifications that are unseen in subject image. Project page: <https://liuff19.github.io/Make-Your-3D/>.

Keywords: 3D Generation · Personalization · Fast Speed

1 Introduction

Subject-driven customization has emerged as a prominent aspect within the field of generative models [40, 44, 60, 66, 68], providing a wide range of multimedia applications [61, 70]. It aims to synthesize the individual subject in diverse contexts and styles while retaining high fidelity to subject-specific identities. Despite the great progress of personalization in text-to-image (T2I) [3, 8, 44, 45, 67] and text-to-video (T2V) [32, 42, 58, 60, 63] models, the exploration of customized 3D generation remains relatively limited.

Driven by the advancements in neural 3D representations [21, 34], extensive datasets [5, 9, 10], and the diffusion-based generative models [16, 43], recent works [26, 29, 37, 47, 52, 57] have demonstrated high-quality automated 3D generation from text or image prompt. Although text or image prompts allow for some degree of control over the generated 3D asset, it remains challenging to produce high-fidelity and subject-specific 3D content with text-driven modifications such as novel colors, poses, or attributes that are unseen in any of the input subject images. Enabling the creation of subject-specific 3D assets through flexible text controls would greatly simplify the workflow for artists and streamline 3D acquisition processes. One notable attempt for subject-driven 3D generation is DreamBooth3D [40], which combines a personalizing model (DreamBooth [44]) and a text-to-3D model (DreamFusion [37]) with two-stage tuning for DreamBooth to optimize NeRF [34] representation. However, it has two inherent limitations: (a) time-consuming optimization of NeRF representation with the two fine-tuning stages in DreamBooth and (b) requirements of multiple subject-specific images as input, which significantly limits the range of applications.

In this paper, we propose **Make-Your-3D**, a novel co-evolution framework for fast and consistent subject-driven 3D content generation. Specifically, given only a single casual subject image, we can generate subject-specific 3D assets that align with text-driven modifications as contextualization within 5 minutes, which is **36× faster** than DreamBooth3D [40]. As shown in Fig. 1, we personalize 3D content with geometric consistency and strong appearance identity preservation from given subjects while also respecting the variations (*e.g.*, sitting or wearing suit) provided by input text prompts.

For Make-Your-3D, we draw inspiration from recent advancements in personalized models [67] and multi-view diffusion models [29]. Despite the customization capability of personalized models and the 3D consistency of multi-view diffu-

sion models, there remains a domain gap between the targeted subject and these two models, particularly when the subject is unseen in the training data [44, 47]. Therefore, the key idea of our method is to harmonize the distribution of the identity-specific 2D generative model and multi-view diffusion model, aligning them with the distribution of the desired 3D subject. Specifically, we design a co-evolution framework to reduce the variance of distributions, where each model undergoes a process of learning from the other through identity-aware optimization and subject-prior optimization, respectively. Given a casual image captured from a subject, we first lift it to a 3D space through a multi-view diffusion model and capture its multi-views to optimize the 2D personalized model with an identity-enhancement process, which imposes enhanced identity-aware information into the 2D model. Next, we apply the original 2D personalized model to multi-views of the subject with modified text description and obtain more diverse images of the subject. Then we optimize the multi-view diffusion model with such various subject images in the subject-prior optimization process, infusing the subject-specific prior into the multi-view diffusion model. Finally, we cascade the two optimized models to process the single input image of the targeted subject and generate consistent subject-driven 3D results.

We conduct extensive experiments on the dataset used in DreamBooth3D [44] and open-vocabulary wild images captured from a subject with different styles. We also conduct a user study to evaluate the subject and prompt fidelity in our synthesized 3D results. The experiments validate that our method is capable of producing vivid and high-fidelity 3D assets with strong adherence to the given subject while highly respecting the contextualization in the input text prompts. Compared to DreamBooth3D [40], our method not only surpasses in terms of quality, resolution, and consistency, but also shows a remarkable $36\times$ speed improvement for efficiency. Unlike DreamBooth3D [40], our approach takes only a single wild image as input, eliminating the need for 3-6 carefully selected images of the same subject. Fig. 1 shows sample results of Make-Your-3D on different subjects and contextualizations, indicating that our co-evolution framework promotes the capability of generating subject-specific 3D assets for both models.

2 Related Works

Text-to-3D Generation. With exciting breakthroughs emerging in the image and video generation [4, 17, 41, 64, 71], there has been growing interest in 3D content generation [25, 28], particularly in text-to-3D generation [24]. One approach is to utilize extensive data [5, 9, 10] to train 3D generative models [6, 15, 22, 30, 31, 36, 48, 69, 73], akin to Text-to-Image (T2I) generation. However, due to constraints by the scale and quality of paired text-3D data, these methods are often limited to specific object categories and may exhibit a perceived lack of realism. Another line of text-to-3D is pioneered by DreamFusion [37], which employed a Score Distillation Sampling (SDS) loss to optimize a parametric 3D representation guided by the pre-trained 2D diffusion model [46]. Following works concentrate on improving 3D consistency, novel view quality, and genera-

tion speed through strategies such as incorporating 3D priors [27, 47, 55], crafting a tailored optimization strategy [26, 50, 57], and selecting more expressive and efficient representations [7, 52]. However, only text is not informative enough to express complex scenes or concepts, which can be a hindrance to 3D content creation. Moreover, it is often unattainable to create contextually diverse 3D assets that precisely align with the user-desired objects.

Image-to-3D Generation. Given the rich information embedded in images, numerous studies [11, 29, 33, 52, 54, 59, 65, 72] have explored to generate 3D content from a single image. Early attempts integrated the input image into the optimization pipeline by creating loss based on predicted depth [11, 65] or object masks [33], comparing them with the rendered image. Magic123 [38] designs a two-stage coarse-to-fine framework for high-quality image-to-3D generation, employing textual inversion to ensure the generation of object-preserving geometry and textures. DreamGaussian [52] and Repaint123 [72] leverage a more efficient Gaussian splatting representation [21], significantly improving the optimization speed. Wonder3D [29] uses a 2D diffusion model to generate multi-view normal maps with color images and applies a geometry-aware normal fusion algorithm for direct surface extraction. However, despite the capability of these approaches to generate 3D content from a single image, their excessive reliance on images to maintain 3D consistency results in a lack of diversity in the generated 3D content, sometimes resembling more of a 3D reconstruction task. While HarmonyView [59] introduces the concept of harmonizing consistency and diversity, the displayed results often fall short of delivering satisfactory diversity, let alone achieving subject-driven customization. Different from their methods, our work is dedicated to reconstructing the concept of the provided object rather than the input image, thereby preserving the generated diversity.

Subject-Driven Content Creation. An increasing number of works [14, 19, 20, 23, 40, 44] are focusing on subject-driven generation, enabling users to personalize the generated content for specific subjects and concepts. Dreambooth [44] fine-tunes the 2D diffusion model and expands the model’s language-vision dictionary with rare tokens using multiple images, achieving personalized text-to-image generation. IP-Adapter [67] realizes controllable generation by incorporating image prompt in the text-to-image models via the design of a lightweight decoupled cross-attention mechanism. VideoBooth [20] injects image prompts into the text-to-video model (T2V) in a coarse-to-fine manner, achieving customized content generation for videos. Despite remarkable success in personalizing T2I and T2V models, they do not generate 3D assets or afford any 3D control. The first attempt in 3D subject-driven generation is DreamBooth3D [40], which proposes a simple pipeline utilizing DreamBooth [44] for 3D subject-driven generation. However, its generation is constrained by DreamBooth [44], requiring heavy fine-tuning stages with multiple subject images over 3 hours, which limits the range of applications. In contrast, our method can achieve fast subject-driven 3D content generation from only a single image of a subject within 5 minutes.

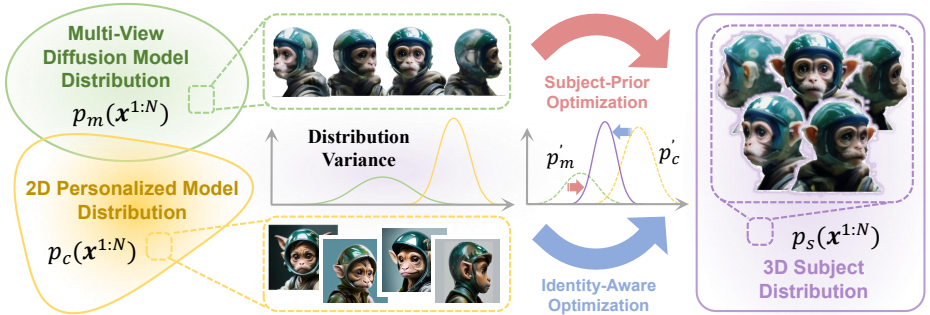


Fig. 2: Distribution variance between the wild subject and pre-trained models. Taking a *monkey* image and text prompt “with elf ears” as input, the pre-trained 2D personalized model and multi-view diffusion model generate images out of the distribution of desired ones, *i.e.*, the specific monkey with elf ears. To solve the problem, we carefully design a co-evolution framework including subject-prior and identity-aware optimization to harmonize the distributions and achieves desired 3D assets.

3 Method

In this section, we introduce our framework, *i.e.*, Make-Your-3D, for fast and consistent subject-driven 3D content generation. Our goal is to align output 3D assets with the distribution of the desired subject. To this end, we first review the scheme of diffusion model (Sec. 3.1), which is the basis of our pre-trained multi-view and personalized models. Then we analyze the distribution variance and optimization target of our method (Sec. 3.2). We further introduce our co-evolution framework, including identity-aware optimization (Sec. 3.3) and subject-prior optimization (Sec. 3.4). Finally, we present our mesh extraction process (Sec. 3.5). An overview of our framework is depicted in Fig. 3.

3.1 Preliminaries

Diffusion Models [16, 49] are probabilistic generative models that comprise two processes: (a) a forward process that gradually adds Gaussian noise to the data following a T steps Markov chain and (b) a denoising process that generates samples from the Gaussian distribution. Let $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ represent the real data under the additional condition \mathbf{c} , and let $t \in [0, T]$ denote the time step of the diffusion process. The training objective of the diffusion model ϵ_θ , which predicts noise, is calculated as the following variational bound:

$$\mathcal{L}_{diff} = \mathbb{E}_{\mathbf{x}_0, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \mathbf{c}, t} \|\epsilon - \epsilon_\theta(\mathbf{x}_t, \mathbf{c}, t)\|^2, \quad (1)$$

where $\mathbf{x}_t = \alpha_t \mathbf{x}_0 + \sigma_t \epsilon$ is the noisy data at step t , and α_t, σ_t are fixed sequence of the noise schedule. For the conditional diffusion models, classifier-free guidance [18] is often employed as a prevailing method. In the sampling stage, the

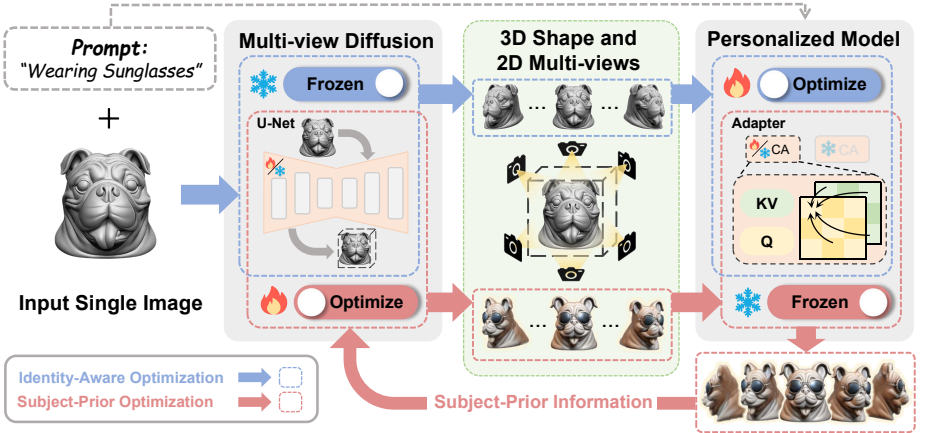


Fig. 3: The overall framework of our proposed Make-Your-3D. Our framework includes identity-aware optimization of 2D personalized model and subject-prior optimization of multi-view diffusion model to approximate subject distribution. The identity-aware optimization (Sec. 3.3) lifts input image to 3D space through a frozen multi-view diffusion model and optimizes the 2D personalized model via multi-views. The subject-prior optimization (Sec. 3.4) adopts diverse images from frozen personalized model to infuse the subject-specific prior into the multi-view diffusion model.

prediction noise is computed as a combination of conditional model $\epsilon_{\theta}(\mathbf{x}_t, \mathbf{c}, t)$ and unconditional model $\epsilon_{\theta}(\mathbf{x}_t, \mathbf{c})$:

$$\hat{\epsilon}_{\theta}(\mathbf{x}_t, \mathbf{c}, t) = w\epsilon_{\theta}(\mathbf{x}_t, \mathbf{c}, t) + (1 - w)\epsilon_{\theta}(\mathbf{x}_t, t), \quad (2)$$

where w is the guidance scale to adjust the alignment with condition \mathbf{c} . In our study, we utilize the open-source 2D personalized model [67] and the multi-view diffusion model [29], which are both built upon the Stable-Diffusion model [43], to achieve fast and consistent subject-driven 3D generation.

3.2 The Distribution of 3D Subject

Powered by recent advances in personalizing text-to-image models [44, 45, 67] and image-to-3D models [29, 38, 56], an intuitive idea to achieve customized 3D generation is to naively combine these methods. However, as shown in Fig. 2, it fails to yield satisfactory subject-specific 3D assets due to distribution variance between the wild subject and the pre-trained models [12, 44, 67]. To explore how to approximate the subject domain, we first propose that the distribution of the 3D subject denoted as $q_s(\mathbf{z})$, can be modeled as a joint distribution as:

$$q_s(\mathbf{z}) = p_s(\mathbf{x}^{1:N} | \mathcal{I}, y) = p_s(\mathbf{x}^1 | \mathcal{I}, y) \cdot p_s(\mathbf{x}^{2:N} | \mathcal{I}, y), \quad (3)$$

where $\mathbf{x}^{1:N}$ are 2D multi-view color images observed from 3D subject conditioned on subject image \mathcal{I} and text-driven modification y . As we have the pre-trained

2D customized model $p_c(\mathbf{x}^1|\mathcal{I}, y)$ and multi-view diffusion model $p_m(\mathbf{x}^{1:N}|\mathcal{I}) = \prod_{n=1}^N p_m(\mathbf{x}^n|\mathcal{I})$, our key insight involves optimizing both p_c and p_m to p'_c and p'_m , which closely align with the distribution of p_s respectively, *i.e.*,

$$p'_c(\mathbf{x}^1|\mathcal{I}, y) \approx p_s(\mathbf{x}^1|\mathcal{I}, y), p'_m(\mathbf{x}^{2:N}|\mathcal{I}) \approx p_s(\mathbf{x}^{2:N}|\mathcal{I}). \quad (4)$$

Given that the condition y is independent of \mathbf{x}^n in the multi-view diffusion model, we can ultimately approximate $q_s(\mathbf{z})$ with the optimized models p'_c and p'_m , *i.e.*, $q_s(\mathbf{z}) \approx p'_c \cdot p'_m$. Finally, we can formulate the joint distribution as a Markov chain within the diffusion scheme (omit the symbol y and \mathcal{I} for simplicity):

$$p'_c(\mathbf{x}_T^1) p'_m(\mathbf{x}_T^{2:N}) \prod_t p'_m(\mathbf{x}_{t-1}^1 | \mathbf{x}_t^1) p'_m(\mathbf{x}_{t-1}^{2:N} | \mathbf{x}_t^{2:N}), \quad (5)$$

where $p'_c(\mathbf{x}_T^1)$ and $p'_m(\mathbf{x}_T^{2:N})$ are Gaussian noises. Our key problem is to characterize the distribution $p'_c \rightarrow p_s$ and $p'_m \rightarrow p_s$ so that we can sample from this Markov chain to generate 3D assets in the subject distribution. Inspired by the derivation above, we carefully design a co-evolution framework, including identity-aware optimization for p_c and subject-prior optimization for p_m .

3.3 Identity-Aware Optimization

To mimic the appearance of subjects from images and synthesize novel renditions of them in different contexts, DreamBooth [44] finetunes all the parameters of the T2I model with 3-6 captured subject images. However, such a tuning strategy that relies on several images, is inefficient and constrained to scenarios where only one image can serve as input. In contrast, we use only a single subject image as input and choose a more efficient adapter-based T2I model (*i.e.*, IP-Adapter [67]) as our 2D personalized model p_c . Despite the user-friendly appeal of using a single input image, the 2D personalized model [67] suffers from distribution variance [12] between the subject and the training data, leading to outputs that do not resemble the subject as shown in Fig. 7, *i.e.*, $p_c(\mathbf{x}^1|\mathcal{I}, y) \neq p_s(\mathbf{x}^1|\mathcal{I}, y)$. To approximate subject distribution and enhance awareness of identity, we first leverage a multi-view diffusion model [29] p_m to generate the multi-views $\mathbf{x}^{(1:N)}$ with view-direction aware prompt $y^{(1:N)}$ given the input subject image \mathcal{I} and text-driven prompt y . Then we apply augmentations to $\mathbf{x}^{(1:N)}$ and process them through the pretrained CLIP image encoder [39] \mathcal{F} and get $\mathcal{F}(\mathbf{x}^{(1:N)})$. Finally, in the adapter module of the 2D personalized model, we use $\mathcal{F}(\mathbf{x}^{(1:N)})$ and $y^{(1:N)}$ to optimize the parameters of image cross-attention layer while freezing the original UNet model and text cross-attention modules. We follow the similar training objective to obtain p'_c in Eq. 1 with the condition $\mathbf{c} = \{y, \mathcal{F}(\mathbf{x})\}$. The empirical analysis in Sec. 4 demonstrates that our multi-view-based, identity-aware optimization effectively narrows the gap between p_c and the subject domain p_s .

3.4 Subject-Prior Optimization

Powered by the 3D consistency of neural radiance fields, DreamBooth3D [40] distills the fine-tuned DreamBooth to generate 3D assets via score distillation

sampling (SDS) [37]. However, this framework suffers from low resolution and time-consuming optimization for per-sample training from scratch, as shown in Fig. 5, 9, limiting the practical usage. In this work, we optimize a more efficient multi-view diffusion framework based on Wonder3D [29] to approximate $p_s(\mathbf{x}^{2:N}|\mathcal{I})$ while better achieving fast and high-fidelity personalized 3D generation. Given the multi-views $\mathbf{x}^{(1:N)}$ from subject image \mathcal{I} as discussed in Sec. 3.3, we process them through the original 2D personalized model with text-driven modification. Then we obtain diverse outputs $\tilde{\mathbf{x}}^{(1:N)}$ from multi-views, which coarsely adhere to the driven text and subject style with strong subject knowledge prior. In addition, we further exploit the subject geometry prior represented by normal maps $\tilde{\mathbf{n}}^{(1:N)}$ inferred from $\tilde{\mathbf{x}}^{(1:N)}$ by using the off-the-shelf single-view estimator [13]. Finally, we optimize the cross-domain self-attention module in the UNet framework based on multi-view diffusion model [29] to incorporate the subject-specific prior knowledge in the views of 3D distribution. Our objective function consists of two terms: (a) an image diffusion term for subject-prior enhancement and (b) a parameter preservation term for maintaining multi-view ability, which is computed as:

$$\mathcal{L}_{prior} = \mathbb{E}_{\mathbf{x}_0, \mathbf{n}_0, \epsilon, \mathbf{c}_n, \mathbf{c}_i, t} \|\epsilon - \epsilon_\theta(\mathbf{x}_t, \mathbf{c}_n, \mathbf{c}_i, t)\|^2 + \lambda \frac{\|\theta - \theta_0\|_1}{N_\theta}, \quad (6)$$

where $\mathbf{c}_i, \mathbf{c}_n$ are the condition of the subject image with corresponding normal maps, θ_0 is the initial parameter of original multi-view diffusion, N_θ is the number of parameters, and λ is a balancing parameter set to 1. Grounded by visualization studies in Sec. 4, our subject-prior optimization strategy can successfully impose subject prior knowledge into multi-view diffusion model, which aligns $p'_m(\cdot|\mathcal{I}) \rightarrow p_s(\cdot|\mathcal{I})$ by Eq. 6 and yields more desired and consistent subject-driven 3D assets.

3.5 Subject-Driven Mesh Extraction

As done in previous co-evolution framework (*i.e.*, identity-aware optimization and subject-prior optimization), we have aligned the 2D personalized model p'_c and the multi-view diffusion model p'_m with the subject distribution p_s . Given the subject image \mathcal{I} and text modification y , we first cascade the two optimized models to process them and obtain subject-driven multiview color images $\hat{\mathbf{x}}^{(1:N)}$ with respect to Eq. 3. From $\hat{\mathbf{x}}^{(1:N)}$, we apply a recent U-Net based Gaussian model pretrained in LGM [51] to predict 3D Gaussians. Next, we train an efficient NeRF (*i.e.*, Instant-NGP [35]) by using the rendered images from 3D Gaussians, and then convert the NeRF to polygonal meshes [53]. More details can be found in our supplementary materials. With adequately optimized implementation in identity-aware optimization (~ 1 min), subject-driven optimization (~ 3 min), and mesh conversion (~ 1 min), our framework can understand the visual subject in a reference image by approximating the subject distribution, and fast produce high-fidelity, consistent unseen personalized 3D content driven by text modification.

4 Experiments

In this section, we conduct extensive experiments to evaluate our subject-driven 3D content generation framework Make-Your-3D, and show the comparison results against DreamBooth3D [40]. We first present our qualitative results in multi-views and comparisons with baselines [40, 47] in various applications (*e.g.*, stylization, accessorization) (Sec. 4.2). Then we report the quantitative results with a user study (Sec. 4.3). Finally, we carry out more open settings and ablation studies to further verify the efficacy of our framework design (Sec. 4.4). Please refer to the supplementary materials for more visualizations, comparisons, and detailed analysis.

4.1 Experiment Setup

Implementation Details. In our framework implementation, we choose IP-Adapter [67] as our 2D personalized model backbone and apply the learning rates of $1e-4$ with a 0.01 weight decay to the image cross-attention layers with our multi-view based identity-aware optimization. On the other hand, we employ Wonder3D [29] as our multi-view diffusion model and utilize diverse images generated by the original 2D personalized model to subject its U-Net module to subject-prior optimization, with a $5e-5$ learning rates and a $1e-2$ weight decay. Notably, for each subject image, it takes only 5 minutes to complete all optimization stages on a single NVIDIA RTX3090 (24GB) GPU, which is far more efficient than 3 hours tuning on 4 core TPUv4 used in DreamBooth3D [40]. We use a fixed 30 iterations to optimize the personalized model. For the multi-view diffusion model, we use around 100 iterations in subject-prior optimization across different objects. To reconstruct 3D geometry, our method is built on the instant-NGP [35] based Gaussian reconstruction method [53].

Baselines and Metrics. We extensively compare our method with two baselines: DreamBooth3D [40] and an implementation for multi-view dreambooth in MVDream [47]. Since the two baseline methods [40, 47] do not have released related code, their results are obtained by downloading from their project pages. For metrics, we mainly show our results with notable comparisons through visualization. Following [37, 40], we evaluate our approach with the CLIP R-Precision metric in CLIP ViT-B/16, ViT-B/32, and ViT-L-14 models. We also conduct a user study to further demonstrate the subject-driven fidelity, prompt fidelity, consistency, and overall quality of our method.

4.2 Qualitative Results

Visual Results of Make-Your-3D. Fig. 4 shows sample visual results of our method across different subjects with customized text prompts. The results demonstrate high-fidelity and consistent 3D generation with Make-Your-3D for

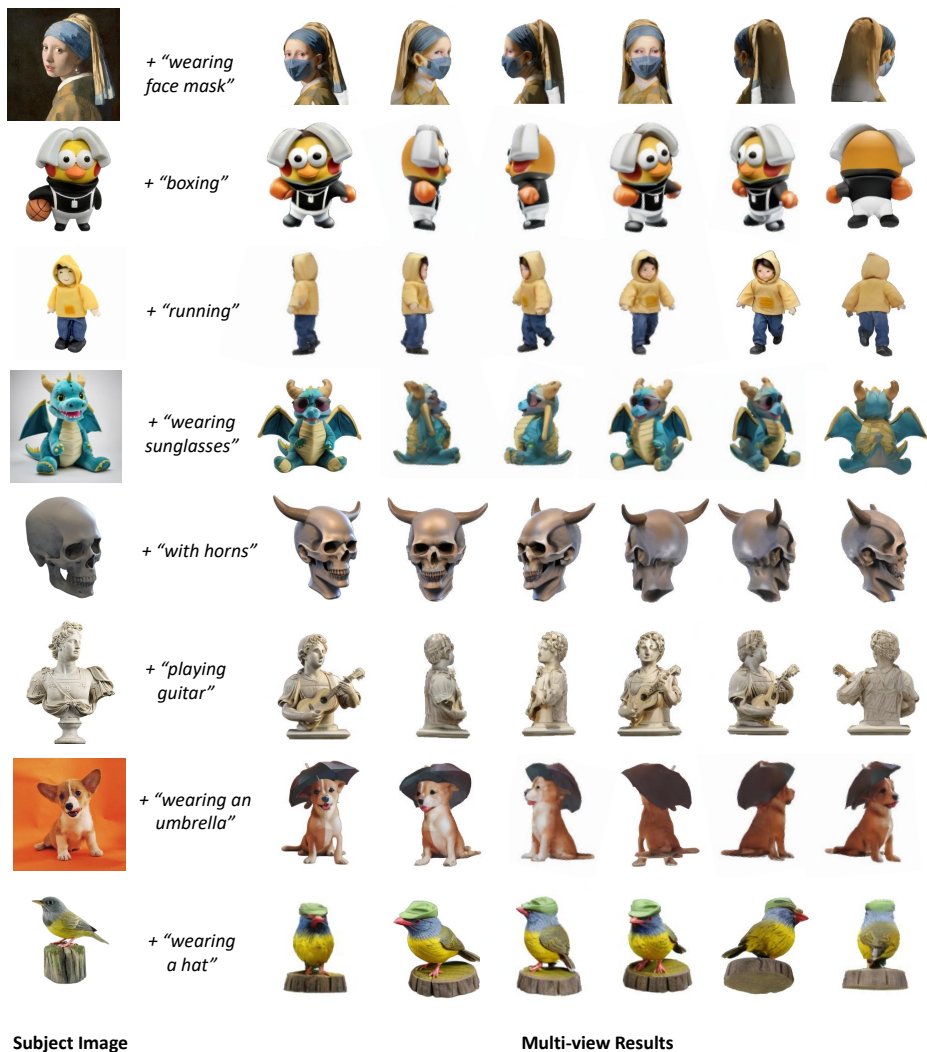


Fig. 4: Visual results of Make-Your-3D on different subjects with customized text inputs. The multi-view results demonstrate that our method can generate 3D assets with high-fidelity, 3D consistency, subject preservation, and faithfulness to the text prompts.

open-vocabulary wild input subject images, achieving faithful alignment respecting the context in the input text prompt.

Qualitative Comparisons. We compare our method with DreamBooth3D [40] in various applications shown in Fig. 5, including color editing, accessorization, stylization, and motion modification. We observe that DreamBooth3D only generates coarse results with multiple images as input over 3 hours. In contrast,



Fig. 5: The qualitative comparisons with DreamBooth3D. We use the same text prompt and only one of the input images as in DreamBooth3D. Notice ours perform better on the object details with less input images.

our Make-Your-3D produces higher quality subject-driven 3D results with compelling object details from only **a single** subject image within only 5 minutes, which is **36 \times** faster than DreamBooth3D. We also conduct comparisons with a recent multi-view DreamBooth implementation in MVDream [47] in Fig. 6. The results further indicate that our method can not only achieve great 3D consistency but also better preserve subject identity without overfitting to data bias in terms of generated styles shown in mutli-view DreamBooth [47].

4.3 Quantitative Results

Table 1 shows the average CLIP R-Precision over 160 evenly spaces azimuth renders at a fixed elevation of 40 degrees, following the same setting in DreamBooth3D [40] for fair-

Table 1: Quantitative comparisons on rendered images with text prompts using different CLIP retrieval models.

	ViT-B/16 \uparrow	ViT-B/32 \uparrow	ViT-L-14 \uparrow
DreamBooth3D [40]	0.783	0.710	0.797
MV DreamBooth [47]	0.805	0.735	0.813
Make-Your-3D (Ours)	0.817	0.764	0.826

ness. Results clearly demonstrate higher scores for Make-Your-3D, indicating better 3D consistency and text-prompt alignment of our results. For user study, we render 360-degree videos of subject-driven 3D models and show each volun-

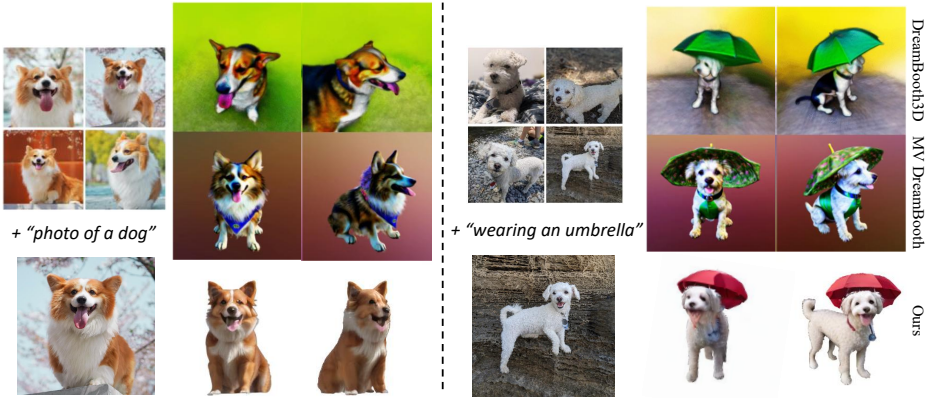


Fig. 6: The multi-view qualitative comparisons. We are able to generate more realistic objects and achieve better subject preservation with multi-view consistency.



Fig. 7: Ablation study of our method. We ablate the design choices of identity-aware optimization and subject-prior optimization.

teer with five samples of rendered video from a random method. They can rate in four aspects: 3D consistency, subject fidelity, prompt fidelity, and overall quality on a scale of 1-10, with higher scores indicating better performance. We collect results from 30 volunteers shown in Table 2. We find our method is significantly preferred by users over these aspects.

4.4 Ablation Study and Discussion

We carry out ablation studies on the design of Make-Your-3D framework in Fig. 7 to verify the effectiveness of our co-evolution framework. Specifically, we perform ablation on identity-aware optimization and subject-prior optimization. The results reveal that the omission of any of two elements leads to a degradation in



Fig. 8: More personalization results for humans. Given a customized description and a face image, we can generate high-quality attributes (*e.g.*, beard, clothes) for the 3D character according to various contexts.

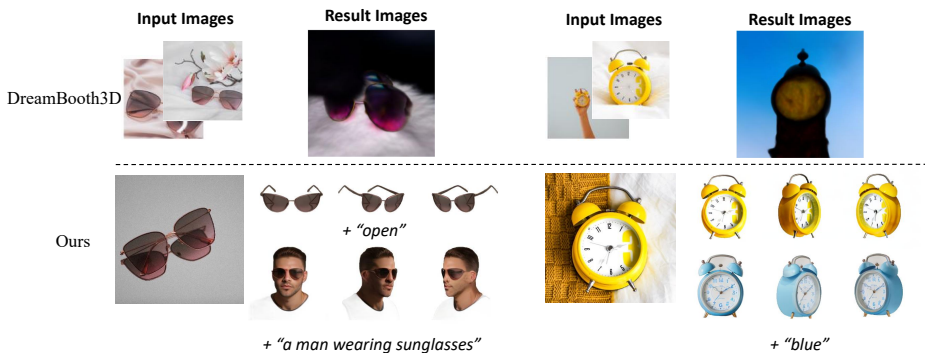


Fig. 9: Comparisons with the failure cases in DreamBooth3D [40]. As DreamBooth3D fails to reconstruct thin object structures like sunglasses and suffers from limited view variation, Our method has made significant improvements in fine details of thin objects and fast 3D personalization from a single subject image.

terms of subject-driven fidelity. Notably, the absence of identity-aware optimization leads to worse subject preservation and consistency. The lack of subject-prior optimization results in less plausible multi-view rendering, especially in cases where the back view lacks informative subject-prior guidance. This illustrates the effectiveness of our overall framework (Fig. 3) that can approximate subject

distribution and have great identity-specific preservation. Moreover, our method is robust in various open-vocabulary settings from wild web images and achieves high-quality results in failure cases of DreamBooth3D [40]

Table 2: Quantitative comparison results of DreamBooth3D [40], multi-view DreamBooth [47] and our Make-Your-3D on the multi-view consistency, subject fidelity, prompt fidelity, and overall quality score in a user study, rated on a range of 1-10, with higher scores indicating better performance.

Method	Multi-view Consistency	Subject Fidelity	Prompt Fidelity	Overall Quality
DreamBooth3D [40]	6.05	6.42	6.89	5.33
MV DreamBooth [47]	8.76	7.55	6.73	7.59
Make-Your-3D (Ours)	9.01	8.91	8.70	9.05

shown in Fig. 9. Driven by our co-evolution framework, we can serve more applications such as human personalization shown in Fig. 8, where we can change their attributes like hair, clothes, and more. These surprising results further support the effectiveness of the co-evolution framework in our Make-Your-3D and present the great potential for subject-driven customization. More impressive results on different applications can be found in our supplementary materials.

5 Conclusion

In this paper, we have proposed Make-Your-3D, a method for fast and consistent subject-driven 3D content generation. To approximate the distribution of the 3D subject, we introduce a novel co-evolution framework. This includes an identity-aware optimization for 2D personalized model and a subject-specific optimization for multi-view diffusion model, through which each model adapts and improves the other’s capacity to capture the subject-driven identity. Therefore, our method bridges the distribution variance from the 3D subject, achieving high-fidelity, multi-view coherent, and subject-specific 3D assets that faithfully adhere to the contextualization in text guidance (*e.g.*, playing guitar, boxing, etc.). Notably, we only need a single subject image as input and produce per 3D result within 5 minutes, 36× faster than DreamBooth3D [40]. Extensive qualitative and quantitative experiments verify the effectiveness and efficiency of our co-evolution framework on 3D content personalization and demonstrate the potential for a wide range of applications.

Limitations and Future Work. Although our Make-Your-3D allows for high-quality 3D personalization and demonstrates better performance than previous work, the quality still seems to be limited to the backbone itself based on Stable Diffusion v1.5. The larger diffusion model such as SDXL [1] will further improve our performance. In future work, we are interested in exploring the 3D scene-level personalization which is a more challenging and complex task. We hope that our Make-Your-3D will pave the way for future advancements, as we believe this technology of subject-driven 3D generation may have a disruptive effect on various sectors, including advertising, entertainment, fashion, and more.

References

1. stable-diffusion-xl-base-1.0. <https://huggingface.co/stabilityai/stable-diffusion-xl-base-1.0>, accessed: 2023-08-29 **14**
2. Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F.L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., et al.: Gpt-4 technical report. arXiv preprint arXiv:2303.08774 (2023) **21, 22, 23**
3. Avrahami, O., Hertz, A., Vinker, Y., Arar, M., Fruchter, S., Fried, O., Cohen-Or, D., Lischinski, D.: The chosen one: Consistent characters in text-to-image diffusion models. arXiv preprint arXiv:2311.10093 (2023) **2**
4. Cao, H., Tan, C., Gao, Z., Xu, Y., Chen, G., Heng, P.A., Li, S.Z.: A survey on generative diffusion models. *IEEE Transactions on Knowledge and Data Engineering* pp. 1–20 (2024). <https://doi.org/10.1109/TKDE.2024.3361474> **3**
5. Chang, A.X., Funkhouser, T., Guibas, L., Hanrahan, P., Huang, Q., Li, Z., Savarese, S., Savva, M., Song, S., Su, H., et al.: Shapenet: An information-rich 3d model repository. arXiv preprint arXiv:1512.03012 (2015) **2, 3**
6. Chen, H., Gu, J., Chen, A., Tian, W., Tu, Z., Liu, L., Su, H.: Single-stage diffusion nerf: A unified approach to 3d generation and reconstruction (2023) **3**
7. Chen, R., Chen, Y., Jiao, N., Jia, K.: Fantasia3d: Disentangling geometry and appearance for high-quality text-to-3d content creation (2023) **4**
8. Chen, W., Hu, H., Li, Y., Ruiz, N., Jia, X., Chang, M.W., Cohen, W.W.: Subject-driven text-to-image generation via apprenticeship learning. *Advances in Neural Information Processing Systems* **36** (2024) **2**
9. Deitke, M., Liu, R., Wallingford, M., Ngo, H., Michel, O., Kusupati, A., Fan, A., Laforte, C., Voleti, V., Gadre, S.Y., et al.: Objaverse-xl: A universe of 10m+ 3d objects. *Advances in Neural Information Processing Systems* **36** (2024) **2, 3**
10. Deitke, M., Schwenk, D., Salvador, J., Weihs, L., Michel, O., Vanderbilt, E., Schmidt, L., Ehsani, K., Kembhavi, A., Farhadi, A.: Objaverse: A universe of annotated 3d objects. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 13142–13153 (2023) **2, 3**
11. Deng, C., Jiang, C.M., Qi, C.R., Yan, X., Zhou, Y., Guibas, L., Anguelov, D.: Nerdi: Single-view nerf synthesis with language-guided diffusion as general image priors (2022) **4**
12. Du, C., Li, Y., Qiu, Z., Xu, C.: Stable diffusion is unstable. *Advances in Neural Information Processing Systems* **36** (2024) **6, 7**
13. Eftekhari, A., Sax, A., Malik, J., Zamir, A.: Omnidata: A scalable pipeline for making multi-task mid-level vision datasets from 3d scans. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 10786–10796 (2021) **8**
14. Gal, R., Alaluf, Y., Atzmon, Y., Patashnik, O., Bermano, A.H., Chechik, G., Cohen-Or, D.: An image is worth one word: Personalizing text-to-image generation using textual inversion (2022) **4**
15. Gupta, A., Xiong, W., Nie, Y., Jones, I., Oğuz, B.: 3dgen: Triplane latent diffusion for textured mesh generation (2023) **3**
16. Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models. *Advances in neural information processing systems* **33**, 6840–6851 (2020) **2, 5**
17. Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models (2020) **3**
18. Ho, J., Salimans, T.: Classifier-free diffusion guidance. arXiv preprint arXiv:2207.12598 (2022) **5**

19. Huang, N., Zhang, T., Yuan, Y., Chen, D., Zhang, S.: Customize-it-3d: High-quality 3d creation from a single image using subject-specific knowledge prior (2024) [4](#)
20. Jiang, Y., Wu, T., Yang, S., Si, C., Lin, D., Qiao, Y., Loy, C.C., Liu, Z.: Videobooth: Diffusion-based video generation with image prompts (2023) [4](#)
21. Kerbl, B., Kopanas, G., Leimkühler, T., Drettakis, G.: 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics* **42**(4) (2023) [2](#), [4](#)
22. Kim, S.W., Brown, B., Yin, K., Kreis, K., Schwarz, K., Li, D., Rombach, R., Torralba, A., Fidler, S.: Neuralfield-ldm: Scene generation with hierarchical latent diffusion models (2023) [3](#)
23. Kumari, N., Zhang, B., Zhang, R., Shechtman, E., Zhu, J.Y.: Multi-concept customization of text-to-image diffusion (2023) [4](#)
24. Li, C., Zhang, C., Waghware, A., Lee, L.H., Rameau, F., Yang, Y., Bae, S.H., Hong, C.S.: Generative ai meets 3d: A survey on text-to-3d in aigc era. *arXiv preprint arXiv:2305.06131* (2023) [3](#)
25. Li, X., Zhang, Q., Kang, D., Cheng, W., Gao, Y., Zhang, J., Liang, Z., Liao, J., Cao, Y.P., Shan, Y.: Advances in 3d generation: A survey (2024) [3](#)
26. Lin, C.H., Gao, J., Tang, L., Takikawa, T., Zeng, X., Huang, X., Kreis, K., Fidler, S., Liu, M.Y., Lin, T.Y.: Magic3d: High-resolution text-to-3d content creation. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 300–309 (2023) [2](#), [4](#)
27. Liu, F., Wu, D., Wei, Y., Rao, Y., Duan, Y.: Sherpa3d: Boosting high-fidelity text-to-3d generation via coarse 3d prior (2023) [4](#)
28. Liu, J., Huang, X., Huang, T., Chen, L., Hou, Y., Tang, S., Liu, Z., Ouyang, W., Zuo, W., Jiang, J., et al.: A comprehensive survey on 3d content generation. *arXiv preprint arXiv:2402.01166* (2024) [3](#)
29. Long, X., Guo, Y.C., Lin, C., Liu, Y., Dou, Z., Liu, L., Ma, Y., Zhang, S.H., Habermann, M., Theobalt, C., et al.: Wonder3d: Single image to 3d using cross-domain diffusion. *arXiv preprint arXiv:2310.15008* (2023) [2](#), [4](#), [6](#), [7](#), [8](#), [9](#), [19](#)
30. Lorraine, J., Xie, K., Zeng, X., Lin, C.H., Takikawa, T., Sharp, N., Lin, T.Y., Liu, M.Y., Fidler, S., Lucas, J.: Att3d: Amortized text-to-3d object synthesis (2023) [3](#)
31. Luo, S., Hu, W.: Diffusion probabilistic models for 3d point cloud generation (2021) [3](#)
32. Ma, Z., Zhou, D., Yeh, C.H., Wang, X.S., Li, X., Yang, H., Dong, Z., Keutzer, K., Feng, J.: Magic-me: Identity-specific video customized diffusion. *arXiv preprint arXiv:2402.09368* (2024) [2](#)
33. Melas-Kyriazi, L., Rupperecht, C., Laina, I., Vedaldi, A.: Realfusion: 360deg reconstruction of any object from a single image (2023) [4](#)
34. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM* **65**(1), 99–106 (2021) [2](#)
35. Müller, T., Evans, A., Schied, C., Keller, A.: Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)* **41**(4), 1–15 (2022) [8](#), [9](#), [19](#)
36. Nichol, A., Jun, H., Dhariwal, P., Mishkin, P., Chen, M.: Point-e: A system for generating 3d point clouds from complex prompts (2022) [3](#)
37. Poole, B., Jain, A., Barron, J.T., Mildenhall, B.: Dreamfusion: Text-to-3d using 2d diffusion (2022) [2](#), [3](#), [8](#), [9](#), [19](#)
38. Qian, G., Mai, J., Hamdi, A., Ren, J., Siarohin, A., Li, B., Lee, H.Y., Skorokhodov, I., Wonka, P., Tulyakov, S., Ghanem, B.: Magic123: One image to high-quality 3d object generation using both 2d and 3d diffusion priors (2023) [4](#), [6](#)

39. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: International conference on machine learning. pp. 8748–8763. PMLR (2021) [7](#)
40. Raj, A., Kaza, S., Poole, B., Niemeyer, M., Ruiz, N., Mildenhall, B., Zada, S., Aberman, K., Rubinstein, M., Barron, J., et al.: Dreambooth3d: Subject-driven text-to-3d generation. arXiv preprint arXiv:2303.13508 (2023) [2](#), [3](#), [4](#), [7](#), [9](#), [10](#), [11](#), [13](#), [14](#), [20](#)
41. Reed, S., Akata, Z., Yan, X., Logeswaran, L., Schiele, B., Lee, H.: Generative adversarial text to image synthesis. In: International conference on machine learning. pp. 1060–1069. PMLR (2016) [3](#)
42. Ren, Y., Zhou, Y., Yang, J., Shi, J., Liu, D., Liu, F., Kwon, M., Shrivastava, A.: Customize-a-video: One-shot motion customization of text-to-video diffusion models (2024) [2](#)
43. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 10684–10695 (2022) [2](#), [6](#)
44. Ruiz, N., Li, Y., Jampani, V., Pritch, Y., Rubinstein, M., Aberman, K.: 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. pp. 22500–22510 (2023) [2](#), [3](#), [4](#), [6](#), [7](#), [19](#)
45. Ruiz, N., Li, Y., Jampani, V., Wei, W., Hou, T., Pritch, Y., Wadhwa, N., Rubinstein, M., Aberman, K.: Hyperdreambooth: Hypernetworks for fast personalization of text-to-image models. arXiv preprint arXiv:2307.06949 (2023) [2](#), [6](#)
46. Saharia, C., Chan, W., Saxena, S., Li, L., Whang, J., Denton, E., Ghasemipour, S.K.S., Ayan, B.K., Mahdavi, S.S., Lopes, R.G., Salimans, T., Ho, J., Fleet, D.J., Norouzi, M.: Photorealistic text-to-image diffusion models with deep language understanding (2022) [3](#)
47. Shi, Y., Wang, P., Ye, J., Long, M., Li, K., Yang, X.: Mvdream: Multi-view diffusion for 3d generation (2023) [2](#), [3](#), [4](#), [9](#), [11](#), [14](#)
48. Shue, J.R., Chan, E.R., Po, R., Ankner, Z., Wu, J., Wetzstein, G.: 3d neural field generation using triplane diffusion (2022) [3](#)
49. Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., Ganguli, S.: Deep unsupervised learning using nonequilibrium thermodynamics. In: International conference on machine learning. pp. 2256–2265. PMLR (2015) [5](#)
50. Sun, J., Zhang, B., Shao, R., Wang, L., Liu, W., Xie, Z., Liu, Y.: Dreamcraft3d: Hierarchical 3d generation with bootstrapped diffusion prior (2023) [4](#)
51. Tang, J., Chen, Z., Chen, X., Wang, T., Zeng, G., Liu, Z.: Lgm: Large multi-view gaussian model for high-resolution 3d content creation (2024) [8](#), [19](#)
52. Tang, J., Ren, J., Zhou, H., Liu, Z., Zeng, G.: Dreamgaussian: Generative gaussian splatting for efficient 3d content creation. arXiv preprint arXiv:2309.16653 (2023) [2](#), [4](#)
53. Tang, J., Zhou, H., Chen, X., Hu, T., Ding, E., Wang, J., Zeng, G.: Delicate textured mesh recovery from nerf via adaptive surface refinement. arXiv preprint arXiv:2303.02091 (2023) [8](#), [9](#), [19](#)
54. Tang, J., Wang, T., Zhang, B., Zhang, T., Yi, R., Ma, L., Chen, D.: Make-it-3d: High-fidelity 3d creation from a single image with diffusion prior (2023) [4](#)
55. Tang, S., Zhang, F., Chen, J., Wang, P., Furukawa, Y.: Mvdifusion: Enabling holistic multi-view image generation with correspondence-aware diffusion (2023) [4](#)

56. Wang, P., Shi, Y.: Imagedream: Image-prompt multi-view diffusion for 3d generation. arXiv preprint arXiv:2312.02201 (2023) [6](#)
57. Wang, Z., Lu, C., Wang, Y., Bao, F., Li, C., Su, H., Zhu, J.: Prolificdreamer: High-fidelity and diverse text-to-3d generation with variational score distillation (2023) [2](#), [4](#)
58. Wei, Y., Zhang, S., Qing, Z., Yuan, H., Liu, Z., Liu, Y., Zhang, Y., Zhou, J., Shan, H.: Dreamvideo: Composing your dream videos with customized subject and motion. arXiv preprint arXiv:2312.04433 (2023) [2](#)
59. Woo, S., Park, B., Go, H., Kim, J.Y., Kim, C.: Harmonyview: Harmonizing consistency and diversity in one-image-to-3d (2023) [4](#)
60. Wu, J.Z., Ge, Y., Wang, X., Lei, S.W., Gu, Y., Shi, Y., Hsu, W., Shan, Y., Qie, X., Shou, M.Z.: Tune-a-video: One-shot tuning of image diffusion models for text-to-video generation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 7623–7633 (2023) [2](#)
61. Wu, J., Gan, W., Chen, Z., Wan, S., Lin, H.: Ai-generated content (aigc): A survey. arXiv preprint arXiv:2304.06632 (2023) [2](#)
62. Wu, T., Yang, G., Li, Z., Zhang, K., Liu, Z., Guibas, L., Lin, D., Wetzstein, G.: Gpt-4v (ision) is a human-aligned evaluator for text-to-3d generation. arXiv preprint arXiv:2401.04092 (2024) [21](#)
63. Xing, J., Xia, M., Liu, Y., Zhang, Y., Zhang, Y., He, Y., Liu, H., Chen, H., Cun, X., Wang, X., et al.: Make-your-video: Customized video generation using textual and structural guidance. arXiv preprint arXiv:2306.00943 (2023) [2](#)
64. Xing, Z., Feng, Q., Chen, H., Dai, Q., Hu, H., Xu, H., Wu, Z., Jiang, Y.G.: A survey on video diffusion models (2023) [3](#)
65. Xu, D., Jiang, Y., Wang, P., Fan, Z., Wang, Y., Wang, Z.: Neurallift-360: Lifting an in-the-wild 2d photo to a 3d object with 360deg views (2023) [4](#)
66. Yang, L., Zhang, Z., Song, Y., Hong, S., Xu, R., Zhao, Y., Zhang, W., Cui, B., Yang, M.H.: Diffusion models: A comprehensive survey of methods and applications. ACM Computing Surveys **56**(4), 1–39 (2023) [2](#)
67. Ye, H., Zhang, J., Liu, S., Han, X., Yang, W.: Ip-adapter: Text compatible image prompt adapter for text-to-image diffusion models. arXiv preprint arXiv:2308.06721 (2023) [2](#), [4](#), [6](#), [7](#), [9](#)
68. Zeng, Y., Lu, Y., Ji, X., Yao, Y., Zhu, H., Cao, X.: Avatarbooth: High-quality and customizable 3d human avatar generation. arXiv preprint arXiv:2306.09864 (2023) [2](#)
69. Zhang, B., Tang, J., Niessner, M., Wonka, P.: 3dshape2vecset: A 3d shape representation for neural fields and generative diffusion models (2023) [3](#)
70. Zhang, C., Zhang, C., Zhang, M., Kweon, I.S.: Text-to-image diffusion model in generative ai: A survey. arXiv preprint arXiv:2303.07909 (2023) [2](#)
71. Zhang, C., Zhang, C., Zhang, M., Kweon, I.S.: Text-to-image diffusion models in generative ai: A survey (2023) [3](#)
72. Zhang, J., Tang, Z., Pang, Y., Cheng, X., Jin, P., Wei, Y., Ning, M., Yuan, L.: Repaint123: Fast and high-quality one image to 3d generation with progressive controllable 2d repainting (2023) [4](#)
73. Zhao, Z., Liu, W., Chen, X., Zeng, X., Wang, R., Cheng, P., Fu, B., Chen, T., Yu, G., Gao, S.: Michelangelo: Conditional 3d shape generation based on shape-image-text aligned latent representation (2023) [3](#)

A Additional Implementation Details

In this section, we provide additional details about the implementation of our subject-driven 3D generation framework **Make-Your-3D**.

A.1 Hyperparameter Settings

In our subject-driven optimization, We retain the optimizer settings and ϵ -prediction strategy from the pretrained process, with a 0.9 adam β_1 , a 0.999 adam β_2 , a $1e - 2$ adam weight decay, and a $1e - 8$ adam ϵ . During the optimization, we use a reduced image size of 256×256 .

A.2 Training Details

The time consumption for the three stages in our framework is as follows: identity-aware optimization takes ~ 1 minute, subject-driven optimization takes ~ 3 minutes, and mesh conversion takes approximately ~ 1 minute. The mesh conversion module is designed to be adaptable to a variety of mesh extraction models [29, 51, 53]. For the sake of improved efficiency, we have chosen to utilize the LGM [51], which is capable of generating 3D Gaussians of objects within a mere 7 seconds. After that, we train an efficient NeRF (*i.e.*, Instant-NGP [35]) by using the rendered images from 3D Gaussians, and then convert the NeRF to polygonal meshes [53]. Specifically, we train two hash grids to reconstruct the geometry and appearance from Gaussian renderings. Please refer to LGM [51] for more details of generated Gaussians. With adequately optimized implementation, it takes extra 1 minute to perform this Gaussians to NeRF to mesh conversion. Our codes for implementation will be available upon acceptance.

A.3 Designing Prompts for One-Shot Personalization

Our goal is to let the diffusion model’s be deeply aware of a new subject’s identity. As mentioned in DreamBooth [44], a common way to personalize a diffusion model is to “implant” a new (*unique identifier, subject*) pair into the model’s “dictionary” and label input multi-view images of the subject "a [identifier] [class noun]", where [identifier] is a unique identifier (*e.g.*, “xxy5syt00”) linked to the subject and [class noun] is a coarse class descriptor of the subject (*e.g.*, cat, dog, watch, etc.). However, this simple method loses all the position information in multi-view images, which is strongly important for identity awareness. To fully leverage the concealed information, we propose a novel (*unique identifier, subject, direction*) pair inspired by view-dependent prompting in DreamFusion [37] and automatically label each image "a [identifier] [class noun], [direction]", where "direction" is one of the six directions (*e.g.*, front, back, left, etc.) for the corresponding multi-view image. By employing this design, we ensure that positional cues are incorporated into our identity-aware optimization process. We utilize the resulting six matched pairs for the optimization, further enhancing the model’s ability to capture and comprehend the subject’s identity.

B Additional Results

To further demonstrate the effectiveness and impressive visualization results of our Make-Your-3D, we conducted more experiments including additional comparison results against DreamBooth3D [40] and visual results (*e.g.*, multi-view images, textured meshes, normals, etc.).

B.1 More Qualitative Comparisons

Fig. 10 demonstrates additional qualitative comparisons with DreamBooth3D [40]. We observed that DreamBooth3D tends to generate coarse results for limited subjects, which lack sufficient identity consistency and also suffer from overfitting issues. For example, the red backpack in Fig. 10 exhibits three small signs at the right bottom, which are not preserved by DreamBooth during its generation process. On the other hand, our proposed method successfully maintains this distinct feature while generating high-resolution outputs.

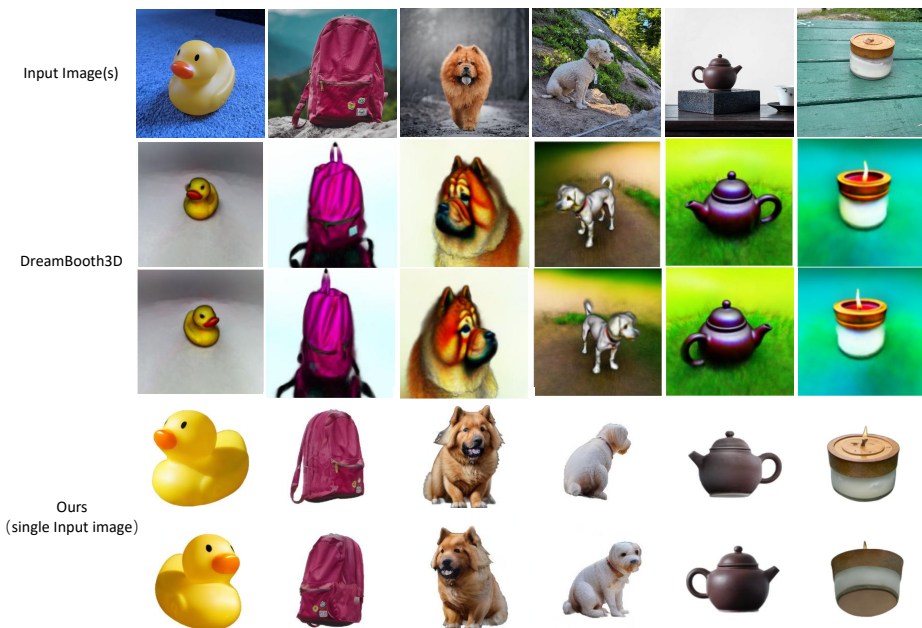


Fig. 10: More qualitative comparisons with DreamBooth3D on subjects of DreamBooth dataset. Notice ours perform better on the object’s identity consistency with less input images.

B.2 More Visual Results

Fig. 11, 16 shows multiple views of the assets rendered for the subjects with different text prompts. Fig. 17 provides additional results with associated nor-

mals and textured meshes to demonstrate the 3D consistency of our results on a variety of customized subjects.

B.3 More Human Personalization Results

Fig. 14, 15 presents additional personalized examples of human faces generated by our co-evolution framework. Our approach enables modification of attributes such as hairstyle, clothing, and more, as well as the capability to alter expressions, makeup, and styling. These remarkable results provide further evidence of the effectiveness of the co-evolution framework in our Make-Your-3D application and demonstrate its broad potential for various areas, such as customizing 3D characters, virtual reality, online clothing try-on, and beyond.

C GPT4-V for 3D Evaluation

We choose a recent automatic and versatile evaluation metric GPTEval3D [62] based on GPT-4Vision (GPT-4V) [2] for additional pairwise comparison. For the two 3D assets, we render them from four or nine viewpoints. These two images will be concatenated together before passing into GPT-4V along with the text instructions. GPT-4V will return a decision of which of the two 3D assets is better according to the instruction. As shown in Fig. 12, We evaluate our results in three main criteria: text-asset alignment, 3D plausibility and texture details.

D Ethical Statement

We confirm that all images used in this paper for research and publication have been obtained and used in a manner compliant with ethical standards. The individuals depicted in these images have given consent for their use, or the images are sourced from publicly available datasets and were used in accordance with the terms of use and permissions. Furthermore, the publication and use of these images do not pose any societal or ethical harm. We have taken necessary precautions to ensure that the research presented in this paper respects individual rights, including the right to privacy and the fundamental principles of ethical research conduct.

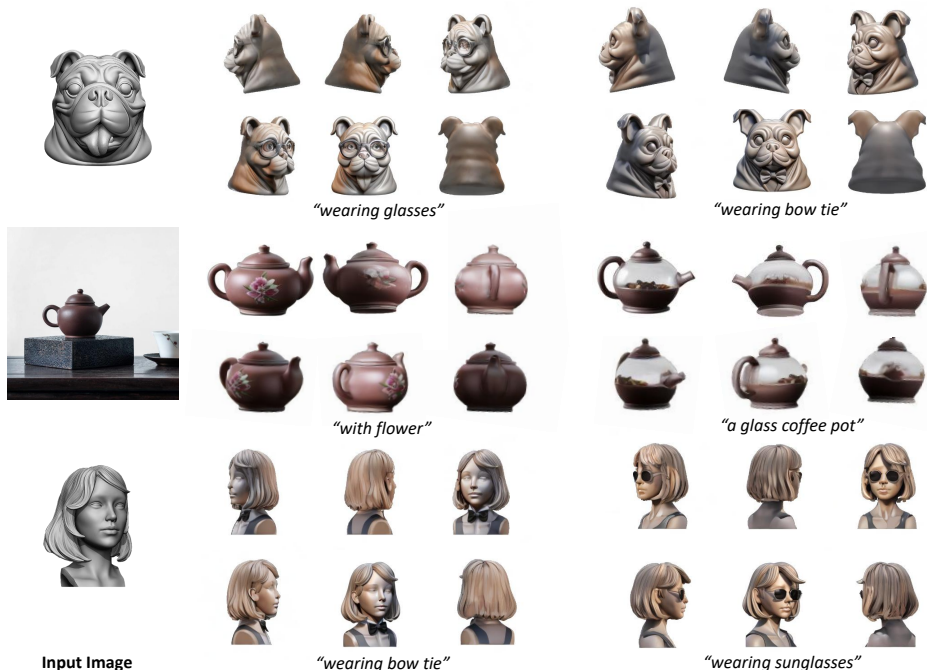


Fig. 11: More personalization results for a subject with different text inputs.

ours

Dreambooth3D

“A dog jumping”

Text-Asset Alignment: The left model ...in a jumping posture, with a fairly neutral expression and typical canine anatomy. The right model ...mid-jump, but with an exaggerated facial expression and a more dynamic pose...but the left one has a more realistic demeanor.

3D Plausibility: The left dog has a smooth and continuous body structure, with limbs proportionately sized and positioned. The right dog, ... has somewhat exaggerated proportions and a less coherent body structure, particularly around the hind legs and neck.

Final answer: neutral neutral

ours

Dreambooth3D

“A red bag”

3D Plausibility: The left model ...with proportional sizing and recognizable parts such as pockets, zippers, and straps that maintain their form across different angles. The right model ...some warping and less consistent detailing.

Low-Level Texture Details: The left model shows more defined textures, with clear fabric patterns, seams, and finer details like the brand logo and zipper pullers being visible. The right model has some textural details, but they appear more blurred and less distinct, particularly in the frontal view.

Final answer: left left left

Fig. 12: Examples of the analysis by GPT-4V [2]. Given two 3D assets, we ask GPT-4V to compare them on various aspects and provide an explanation. We find that GPT-4V’s preference closely aligns with that of humans in our user study.



Fig. 13: More examples of the analysis by GPT-4V [2].

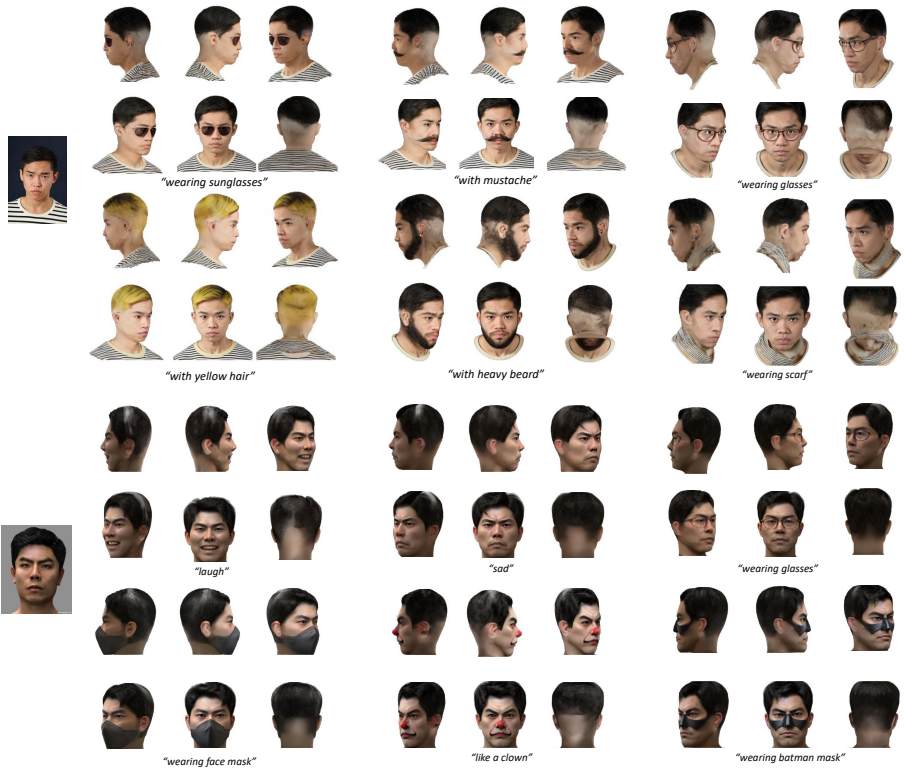


Fig. 14: More personalization results for human faces.



Fig. 15: More personalization results for one person with multiple customized text inputs.

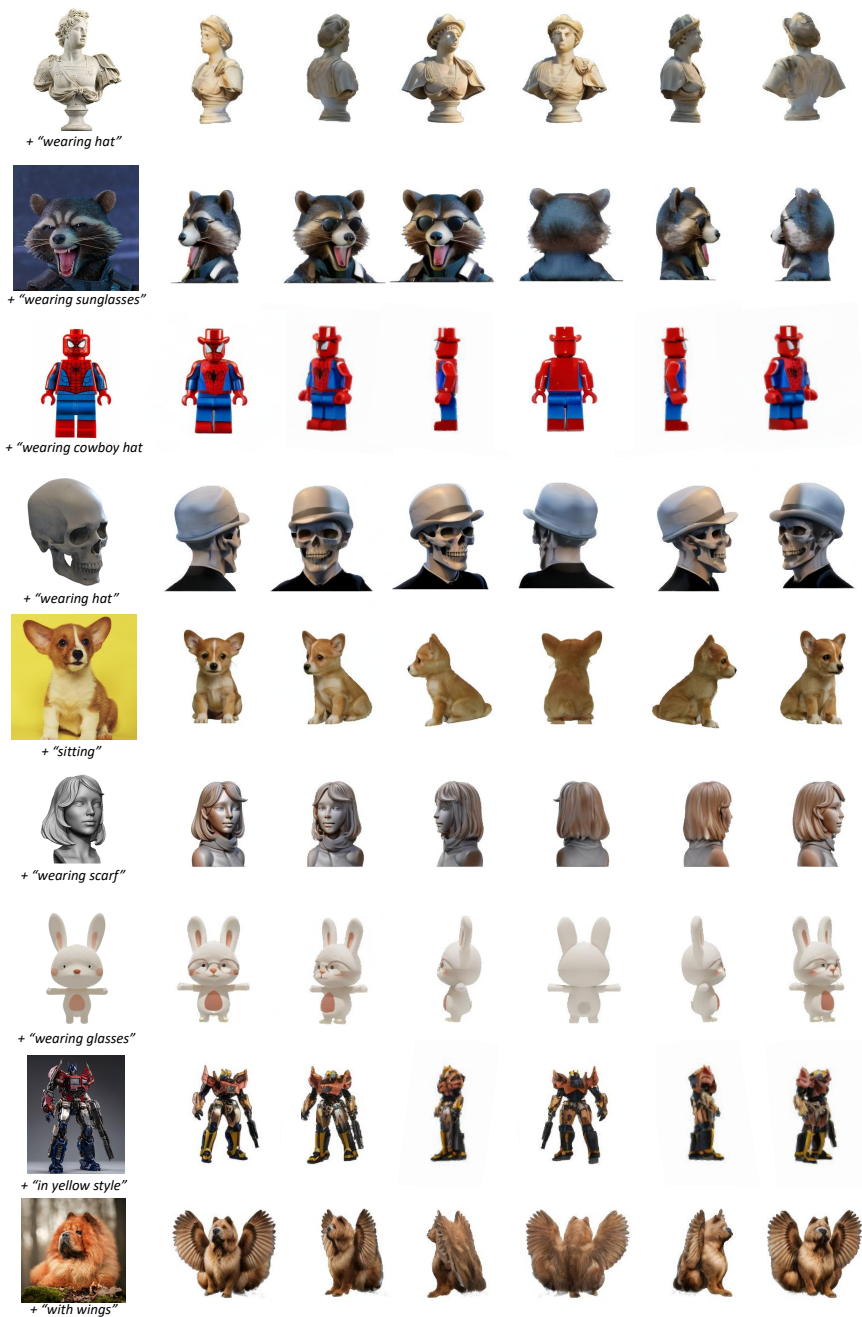


Fig. 16: More visual results of Make-Your-3D on different subjects with customized text inputs.

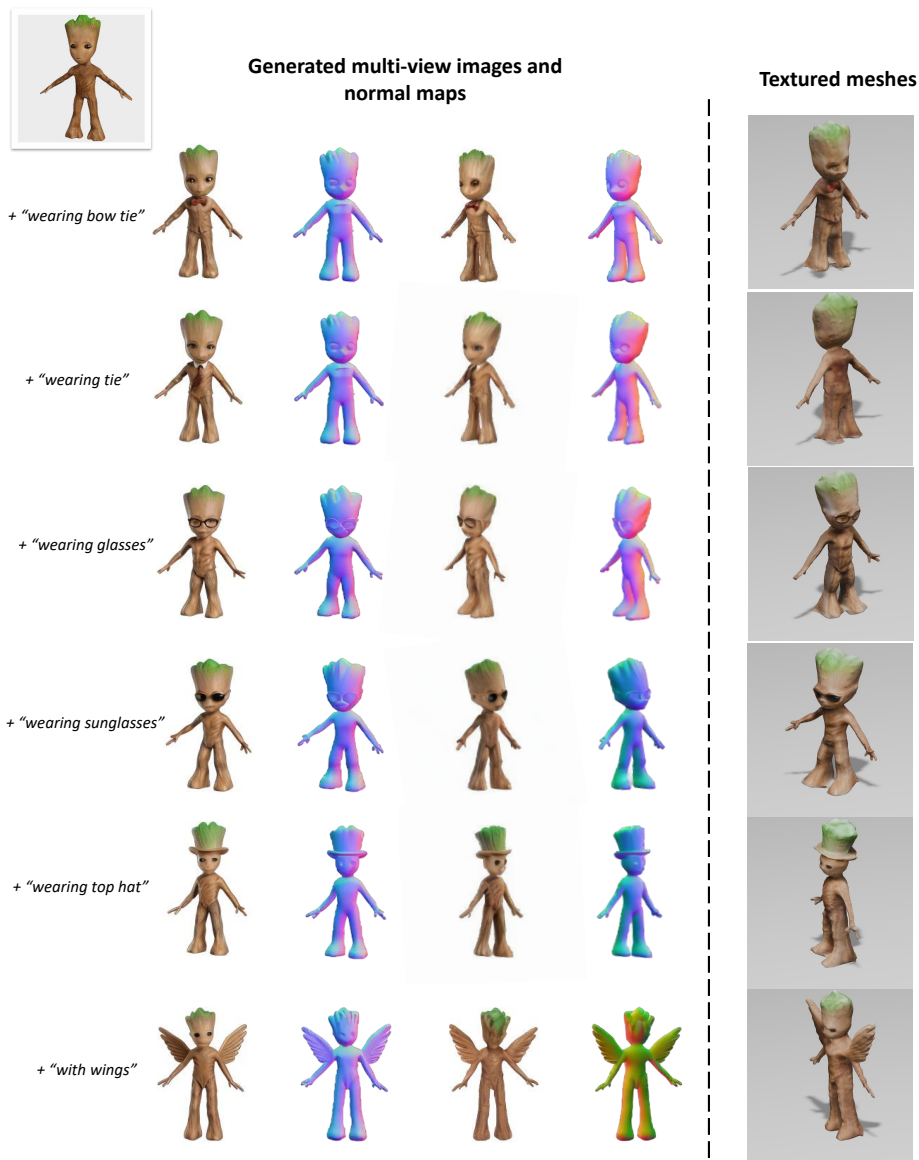


Fig. 17: Textured meshes and normal maps on a subject "Gelute" with various customized text inputs.