

DreamAvatar: Text-and-Shape Guided 3D Human Avatar Generation via Diffusion Models

Yukang Cao^{1*} Yan-Pei Cao^{2*} Kai Han^{1†} Ying Shan² Kwan-Yee K. Wong¹

¹The University of Hong Kong

²ARC Lab, Tencent PCG

Abstract

We present *DreamAvatar*, a text-and-shape guided framework for generating high-quality 3D human avatars with controllable poses. While encouraging results have been produced by recent methods on text-guided 3D common object generation, generating high-quality human avatars remains an open challenge due to the complexity of the human body’s shape, pose, and appearance. We propose *DreamAvatar* to tackle this challenge, which utilizes a trainable NeRF for predicting density and color features for 3D points and a pre-trained text-to-image diffusion model for providing 2D self-supervision. Specifically, we leverage SMPL models to provide rough pose and shape guidance for the generation. We introduce a dual space design that comprises a canonical space and an observation space, which are related by a learnable deformation field through the NeRF, allowing for the transfer of well-optimized texture and geometry from the canonical space to the target posed avatar. Additionally, we exploit a normal-consistency regularization to allow for more vivid generation with detailed geometry and texture. Through extensive evaluations, we demonstrate that *DreamAvatar* significantly outperforms existing methods, establishing a new state-of-the-art for text-and-shape guided 3D human generation. Project page: <https://yukangcao.github.io/DreamAvatar/>.

1. Introduction

The creation of 3D graphical human models has received great attention in recent years due to its wide-ranging applications in fields such as film making, video games, AR/VR, and human-robotic interaction. Traditional methods for building such complex 3D models require thousands of man-hours of trained artists and engineers [8, 9], making the process both time-consuming and highly expert-dependent. With the development of deep learning methods, we have witnessed the emergence of promising methods [39, 5, 46]

which can reconstruct 3D human models from monocular images. These techniques, however, still face challenges in fully recovering details from the input images and rely heavily on the training dataset. To tackle these challenges and simplify the modeling process, incorporating generative models into 3D human avatar modeling has recently received increasing attention from the research community. This approach has the potential to alleviate the need for large 3D datasets and facilitate easier and more accessible 3D human avatar modeling.

To leverage the potential of 2D generative image models for 3D content generation, recent approaches [32, 18, 22] have utilized pre-trained text-guided image diffusion models to optimize Neural Radiance Fields (NeRFs) [24]. DreamFusion [32] introduces a novel Score Distillation Sampling (SDS) strategy to self-supervise the optimization process and achieves promising results. However, human bodies, which are the primary focus of this paper, are widely regarded as complex due to their articulated structures such as the head, arms, thighs, and feet, each capable of posing in various ways. As a result, while DreamFusion [32] and subsequent methods (e.g., Magic3D [18]) produce impressive results, they lack the proper constraints to enforce consistent 3D human structures and often struggle to generate detailed textures for 3D human avatars. More recently, Latent-NeRF [22] introduces a sketch-shape loss based on the 3D shape guidance, but it still faces challenges in generating reasonable results for human bodies. TEXTure [35] proposes generating a texture map for a given 3D shape, but it may suffer from the self-occlusion problem and cannot generate novel geometric structures.

In this paper, we present *DreamAvatar*, a novel framework for generating high-quality 3D human avatars from text prompts and shape priors. Inspired by previous works [32, 18], *DreamAvatar* employs a trainable NeRF as the base representation for predicting density and color features for each 3D point. Coupled with a pre-trained text-to-image diffusion model [36], *DreamAvatar* can be trained to generate 3D avatars using 2D self-supervision. *The key innovation of DreamAvatar lies in three main aspects. Firstly, we leverage SMPL models to provide shape priors, which*

* Equal contribution

† Corresponding author



Figure 1. Results of DreamAvatar. DreamAvatar can generate high-quality geometry and texture for any type of human avatar, based on a simple text description.

yield rough pose and shape guidance for the generation process. *Secondly*, we introduce a *dual space design* comprising a canonical space and an observation space, which are related by a learnable deformation field through the NeRF. The two spaces are jointly optimized, allowing for the transfer of well-optimized texture and geometry from the canonical space to the target posed avatar, and providing more constraints on the canonical model. The idea behind this design is to select a canonical pose that minimizes self-occlusions and is easy for image diffusion models to generate. *Thirdly*, to enable a more vivid generation with detailed geometry such as dresses, cloaks, accessories, and mecha, based on the text guidance, we further exploit a *normal-consistency regularization*, inspired by [42]. Specifically, we extend the MLP in the NeRF to predict a normal vector for each 3D point and encourage the agreement between the predicted normal and the gradient of the density field.

We extensively evaluate our method and compare it with previous methods. *Our DreamAvatar significantly outperforms existing methods and can generate high-quality 3D human avatars with text-consistent geometry and geometry-consistent texture.* Our code will be made publicly available.

2. Related Work

Text-guided 2D image generation. Recently, the CLIP model [33] (Contrastive Language-Image Pre-training) was proposed with the aim of classifying images and text by mapping them to a shared feature space. However, this model is not consistent with the way human perceives language, and it may not fully capture the intended meanings. With the improvements in neural networks and text-image datasets, the diffusion model has been introduced to handle more complex semantic concepts [2, 22, 34, 38]. Follow-up methods are designed to improve computational efficiency, for instance, utilizing a cascade of super-resolution models [2, 38] or sampling from a low-resolution latent space and decoding the latent features into high-resolution images [22]. DreamBooth [37] fine-tunes the diffusion model for certain subjects, while ControlNet [50] and T2I-Adapter [26] propose controlling the pre-trained diffusion model with additional information. However, text-to-3D generation remains a challenge due to the computational cost and the lack of text-3D paired datasets.

Text-guided 3D content generation. Text-guided 3D content generation methods have emerged based on the success of text-guided 2D image generation. Earlier works, such as CLIP-forge [40], generate objects by learning a normalizing flow model from textual descriptions, but these methods are computationally expensive. DreamField [15], CLIP-mesh [17], AvatarCLIP [14], Text2Mesh [23], and Dream3D [47] rely on a pre-trained image-text model [33] to optimize the underlying 3D representation (NeRF or

mesh). Although these approaches eliminate the need for 3D training data through the use of pre-trained vision-language models, they fail to produce convincing 2D renderings.

Recently, DreamFusion [32] proposes score distillation sampling based on the pre-trained diffusion model [38] to enable text-guided 3D generation. Magic3D [18] improves it by introducing a coarse-to-fine pipeline to generate fine-grained text-guided 3D textured mesh. Point-E [27] optimizes the point cloud based on the diffusion model. Latent-NeRF [22] improves training efficiency by directly optimizing the latent features. TEXTure [35] applies a depth-diffusion model [36] to generate texture maps for a given 3D mesh. Despite their promising performance, these methods still struggle in generating text-guided 3D human avatars due to the inherent challenges of this task.

3D human generative models. 3D generative methods based 3D voxel grids [11, 12, 20, 45], point clouds [1, 21, 25, 48, 49, 52], meshes [51] are commonly represented and often require expensive and limited 3D datasets. In recent years, various methods [28, 19, 31, 44, 7] have been proposed to utilize the neural rendering technique NeRF and train on 2D human videos for novel view synthesis. Following their steps, EG3D [6] and GNARF [3] propose a tri-plane representation that uses GANs for 3D generation from latent codes. ENARF-GAN [29] extend NARF [28] to human representation. Meanwhile, EVA3D [13] and HumanGen [16] are proposed to generate human radiance fields directly from the 2D StyleGAN-Human [10] dataset. Although these methods have produced convincing results, they do not have the ability to “dream” or generate new subjects that have not been seen during training.

3. Methodology

Here, we introduce our text-and-shape guided generative network, DreamAvatar, which utilizes a trainable NeRF and a pre-trained text-to-image stable diffusion model [36] to generate 3D human avatars with controllable poses. DreamAvatar incorporates two modelling spaces which are related by a learnable deformation field and jointly optimized through the trainable NeRF module (see Fig. 2). In the following subsections, we first provide the preliminaries that underpin our method in Sec. 3.1. Next, we delve into the details of our method and discuss: (1) the density field based on SMPL models used to evolve geometry, (2) the dual spaces related by a learnable deformation field, and (3) the normal-consistency loss that enables the network to generate text-consistent geometry and geometry-consistent texture in Sec. 3.2.

3.1. Preliminaries

Text-guided 3D generation methods Recent text-guided 3D generation models [32, 18, 22] showcase promising re-

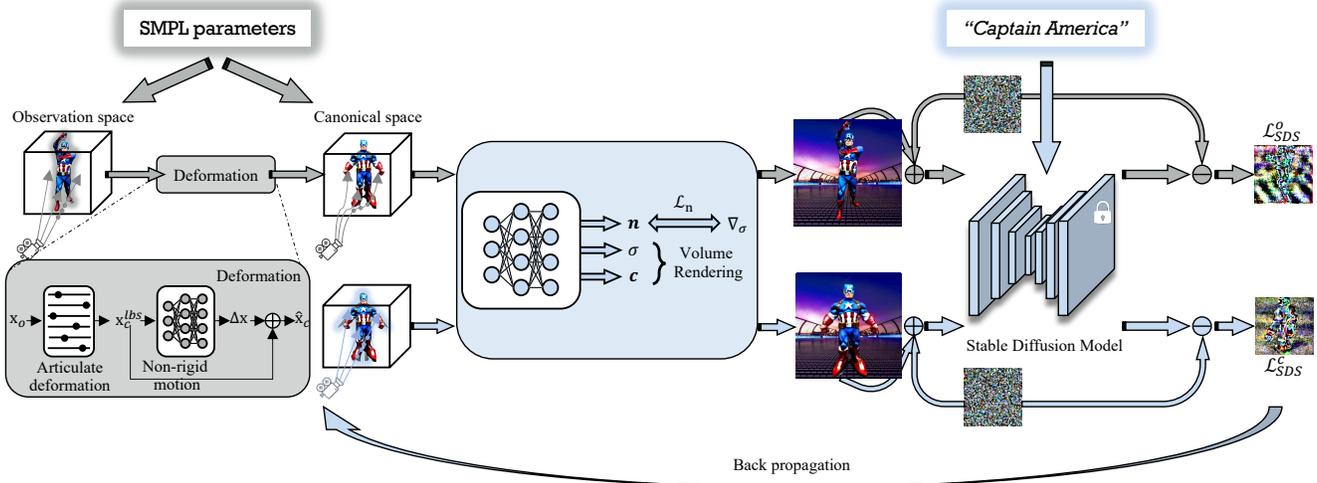


Figure 2. Overview of DreamAvatar. Our network takes as input a text prompt and SMPL parameters to optimize a trainable NeRF model via a pre-trained denoising stable diffusion model. At the core of our network are two modeling spaces that are related by an SMPL-based learnable deformation field, which robustly controls the pose and transfers the high-quality geometry and texture from the canonical space to the observation space.

sults by incorporating three fundamental components:

(1) *NeRF model* that represents a 3D scene via an implicit function:

$$F_\theta(\gamma(\mathbf{x})) \mapsto (\sigma, \mathbf{c}), \quad (1)$$

where \mathbf{x} is a 3D point that will be processed by a grid frequency encoder $\gamma(\cdot)$ [24], σ is the density value and \mathbf{c} represents color. Generally, the implicit function $F_\theta(\cdot)$ is implemented as an MLP with trainable parameters θ .

(2) *Volume Rendering* technique that effectively renders a 3D scene onto a 2D image. For each image pixel, the rendering is done by casting a ray \mathbf{r} from the pixel location into the 3D scene and sampling 3D points μ_i along \mathbf{r} . Density and color of the sampled points are predicted by F_θ . The color C of each image pixel is then calculated by:

$$C(\mathbf{r}) = \sum_i W_i \mathbf{c}_i, \quad W_i = \alpha_i \prod_{j < i} (1 - \alpha_j) \quad (2)$$

where $\alpha_i = 1 - e^{(-\sigma_i \|\mu_i - \mu_{i+1}\|)}$. Note that \mathbf{c} and C are defined as RGB color in [32, 18], and as 4D latent features in [22]. We follow [22] for better training efficiency and higher resolution.

(3) *Score Distillation Sampling (SDS)* derived on text-guided diffusion models ϕ [36, 38]. We employ a pre-trained stable diffusion model [36] with the learned denoising function $\epsilon_\phi(x_t; y, t)$. Here x_t denotes the noisy image at noise level t , and y is the text embedding. Given a rendered image from the NeRF model, we add random noise ϵ to obtain a noisy image x , and the self-supervise SDS loss for optimizing the NeRF model will minimize the difference between the predicted noise $\epsilon_\phi(x_t; y, t)$ and added noise ϵ :

$$\nabla_\theta \mathcal{L}_{SDS}(\phi, F_\theta) = \mathbb{E}_{t, \epsilon} [w(t) (\epsilon_\phi(x_t; y, t) - \epsilon) \frac{\partial x}{\partial \theta}], \quad (3)$$

where $w(t)$ represents a weighting function that depends on the noise level t . See the right part of Fig. 2 for a better view of the process.

SMPL [4, 30] 3D parametric human model It builds a 3D human shape by 6,890 body vertices. Formally, by assembling pose parameters ξ and shape parameters β , we can obtain the 3D SMPL human model by:

$$T_P(\beta, \xi) = \bar{T} + B_S(\beta; \mathcal{S}) + B_P(\xi; \mathcal{P}), \quad (4)$$

$$M(\beta, \xi) = \text{LBS}(T_P(\beta, \xi), J(\beta), \xi, \mathcal{W}), \quad (5)$$

where $T_P(\cdot)$ represents the non-rigid deformation from canonical model \bar{T} by using shape blend shape function B_S , and pose blend shape function B_P . \mathcal{S} and \mathcal{P} are the principal components of vertex displacements. $\text{LBS}(\cdot)$ denotes the linear blend skinning function, corresponding to articulated deformation. It poses $T_P(\cdot)$ based on the parameters ξ and joint locations $J(\beta)$, by using the blend weights \mathcal{W} , individually for each body vertex:

$$\mathbf{v}_o = \mathcal{G} \cdot \mathbf{v}_c, \quad \mathcal{G} = \sum_{k=1}^K w_k \mathcal{G}_k(\xi, j_k), \quad (6)$$

where \mathbf{v}_c is an SMPL vertex under the canonical pose, \mathbf{v}_o denotes the corresponding vertex under the observed pose, w_k is the skinning weight, $\mathcal{G}_k(\xi, j_k)$ is the affine deformation that transforms the k -th joint j_k from the canonical space to observed space, and K is the number of neighboring joints.

Unlike the original SMPL which defines ‘‘T-pose’’ as the canonical model, here, we use ‘‘A-pose’’ as the canonical model which is a more natural human rest pose and makes it easier to generate geometry-consistent texture (see Fig. 2).

3.2. DreamAvatar

As illustrated in Fig. 2, our proposed framework takes as input a text prompt and SMPL parameters, defining the target shape and pose in the observation space. DreamAvatar conducts Score Distillation Sampling (SDS)-based optimization [32] in both observation space and canonical space simultaneously, and learns a deformation field to relate the two spaces. To represent the observation space and canonical space, we utilize an extended neural radiance field where the density, latent feature [22], and normal direction of each sample point can be queried and optimized. We utilize the input SMPL parameters to handle different body parts separately and derive reasonable initial density values in each space. Like Latent-NeRF [22], we employ a pre-trained stable diffusion model to provide self-supervisions. Besides, our framework includes a normal-consistency loss between the predicted normal and the normal calculated by finite difference from the estimated density field, improving the text-geometry consistency in the generated human avatar.

SMPL-derived density field We propose to make our NeRF model evolve from the density field derived from input SMPL models. Specifically, for a 3D point \mathbf{x}_c in the canonical space, we first calculate its signed distance d to the SMPL surface, and then convert it to a density value $\bar{\sigma}_c$:

$$\bar{\sigma}_c = \max(0, \text{softplus}^{-1}(\tau)), \quad \tau = \frac{1}{a} \text{sigmoid}(-d/a),$$

where $\text{sigmoid}(x) = 1/(1 + e^{-x})$, $\text{softplus}^{-1}(x) = \log(e^x - 1)$, and a is a predefined hyper-parameter [47], which is set to 0.001 in our experiments. For a point \mathbf{x}_o in observation space (the upper branch in Fig. 2), we first obtain its $\bar{\sigma}_o$ based on the observed SMPL model and then deform it back to canonical space.

Deformation field Inspired by HumanNeRF [44], we employ a deformation field to map a point \mathbf{x}_o from the observation space to a corresponding point $\hat{\mathbf{x}}_c$ in the canonical space. Specifically, we divide it into two parts: (1) articulated deformation that applies the inverse transformation of SMPL linear blend skinning $LBS(\cdot)$ (Sec. 3.1), and (2) non-rigid motion implemented as an MLP to learn the corrective offset:

$$\hat{\mathbf{x}}_c = \mathbf{x}_c^{lbs} + \text{MLP}_{\theta_{NR}}(\gamma(\mathbf{x}_c^{lbs})), \quad \mathbf{x}_c^{lbs} = \mathcal{G}^{-1} \cdot \mathbf{x}_o, \quad (7)$$

where \mathcal{G} is obtained from the observed SMPL vertex closest to \mathbf{x}_o .

Dual spaces For a 3D point \mathbf{x}_c in canonical space and \mathbf{x}_o in the observation space, we calculate their density σ_c, σ_o and latent color feature $\mathbf{c}_c, \mathbf{c}_o$ by:

$$F(\mathbf{x}_c, \bar{\sigma}_c) = F_\theta(\gamma(\hat{\mathbf{x}}_c)) + (\bar{\sigma}_c, \mathbf{0}) \mapsto (\sigma_c, \mathbf{c}_c), \quad (8)$$

$$F(\mathbf{x}_o, \bar{\sigma}_o) = F_\theta(\gamma(\hat{\mathbf{x}}_c)) + (\bar{\sigma}_o, \mathbf{0}) \mapsto (\sigma_o, \mathbf{c}_o), \quad (9)$$

where $\hat{\mathbf{x}}_c$ is the corresponding point of \mathbf{x}_o in canonical space, and $\bar{\sigma}_c, \bar{\sigma}_o$ are the density values inferred from corresponding SMPL parameters. The introduction of the dual space design serves two purposes. Firstly, it enables us to define a canonical pose that minimizes self-occlusions in 3D and is easier for image diffusion models to generate, which benefits both 2D generation and 3D optimization. Secondly, it facilitates the transfer of optimized texture and geometry from the canonical space to the target pose, while also providing additional constraints on the canonical model.

Background scene Following DreamFusion [32], we implement an environment map MLP that takes the positionally-encoded ray direction as input and similarly predicts the 4D latent color feature $C_{bg}(\mathbf{r})$ for each rendered ray. We will then composite the previously acquired latent color feature $C(\mathbf{r})$ for each rendered ray on top of this background feature with the accumulated alpha value:

$$C'(\mathbf{r}) = C(\mathbf{r}) + (1 - \sum_i W_i)C_{bg}(\mathbf{r}), \quad (10)$$

where the weight W_i follows the definition in Eq. (2).

Training objectives Our experiments show that our setup for evolving the density field derived from SMPL models effectively controls the pose. However, it might also limit the model from producing text-consistent geometry and geometry-consistent texture, as shown in Fig. 8, since it constrains the model to the SMPL shape. To overcome this limitation, we draw inspiration from Ref-NeRF [42] and adopt a normal-consistency loss between the gradient of the density field ($\nabla\sigma$) and the predicted normal direction, by expanding the MLP to predict a 3D normal vector \mathbf{n} for each point:

$$F_\theta(\gamma(\mathbf{x})) \mapsto (\sigma, \mathbf{c}, \mathbf{n}), \quad (11)$$

$$\mathcal{L}_n = b \|\nabla\sigma - \mathbf{n}\|^2, \quad b = \|1 - e^{-\sigma}\| \quad (12)$$

It will unleash the network from the strong SMPL-derived density field, allowing the generation of text-consistent geometry and high-quality texture.

Combining the normal-consistency loss with the SDS loss, the learning objective to train DreamAvatar can be written as:

$$\mathcal{L} = \lambda_{SDS} * (\mathcal{L}_{SDS}^c + \mathcal{L}_{SDS}^o) + \lambda_n * \mathcal{L}_n, \quad (13)$$

where the two SDS losses are defined over canonical space and observation space respectively:

$$\mathcal{L}_{SDS}^c = \mathbb{E}_{t, \epsilon}[w(t)(\epsilon_\theta(x_t \sim c; y, t) - \epsilon) \frac{\partial x}{\partial \theta}], \quad (14)$$

$$\mathcal{L}_{SDS}^o = \mathbb{E}_{t, \epsilon}[w(t)(\epsilon_\theta(x_t \sim o; y, t) - \epsilon) \frac{\partial x}{\partial \theta}]. \quad (15)$$

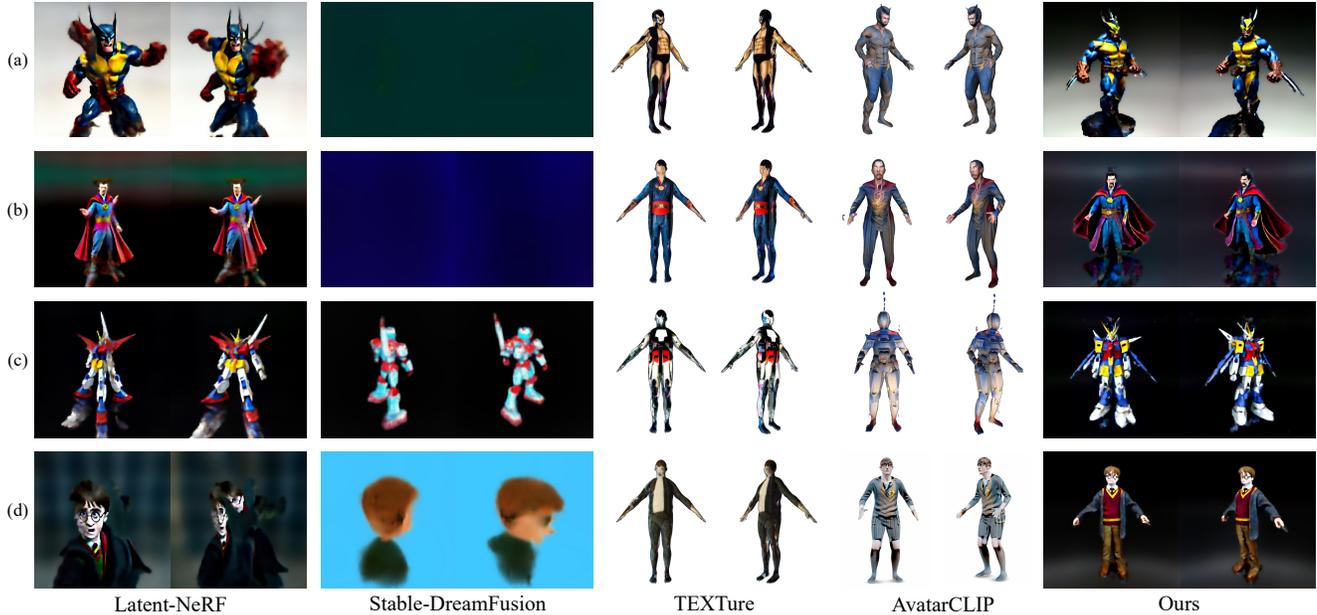


Figure 3. Comparison with existing text-to-3D methods in the canonical space. While other methods struggle or fail to generate reasonable results, our method can perform much better with high-quality geometry and texture. (a): “Wolverine, marvel character”, (b): “Doctor Strange”, (c) “Mobile suit Gundam”, (d) “Harry Potter”.



Figure 4. Avatar generation with different poses. Our method can handle and control the 3D generation with any pose.

4. Experiments

We now validate the effectiveness and capability of our proposed framework from a variety of text prompts and provide comparisons with existing text-guided 3D generation methods given the same text prompts.

Implementation details We follow [22] to implement the NeRF [24] and Stable Diffusion components of our method. For the latter, we use the HuggingFace diffusers, Ver.1-4. Considering there is no official implementation for score distillation sampling, we adopt it from the open-source [41]. Typically, for each text prompt, we train our network for 10,000 iterations, which takes around 2 hours on one single NVIDIA RTX 2080Ti GPU.

Baseline methods As there is no public official implementation for DreamFusion [32] and Magic3D [18], we mainly compare our method with four SOTA methods: (1)

Stable-DreamFusion, which is the open-source from [41]. It re-implements DreamFusion [32] but also adopts several improvements proposed by Magic3D, *e.g.*, training an additional MLP to predict per-point normal instead of calculating the finite difference; (2) Latent-NeRF [22], which uses latent diffusion model [36] to improve training efficiency and utilizes 3D prior to supervising the training; (3) TEXTure [35], which applies the text-guided depth diffusion model to generate the texture map for provided 3D shape; (4) AvatarCLIP [14], which employs CLIP vision-language model and NeUS [43] as the backbone for the text-guided human generation.

4.1. Qualitative Evaluations

Avatar generation with different styles In Fig. 1, we provide a diverse set of 3D human avatars generated by our DreamAvatar. We can consistently observe high-quality ge-

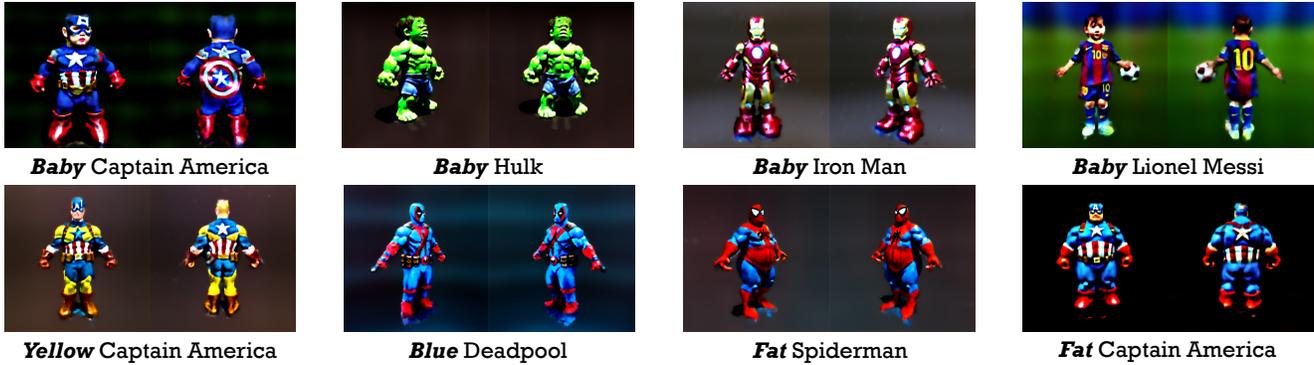


Figure 5. Text manipulation over the avatar generations. Our methods can generate the avatars in different styles by editing the text prompt.

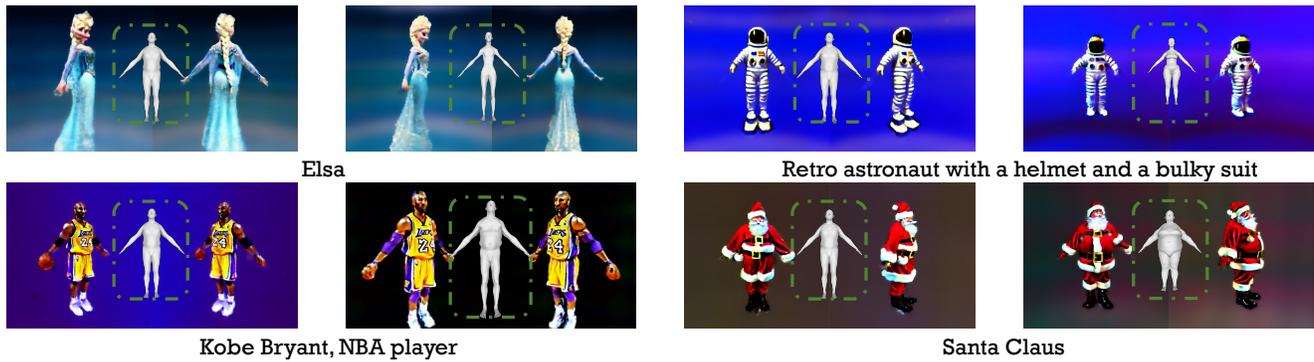


Figure 6. Shape modification via SMPL shape parameters. Our method can generate “thin”, “short”, “tall”, and “fat” 3D avatars based on the input SMPL shape parameters. In the middle is the employed SMPL model. We rescale them with the same ratio during the training.

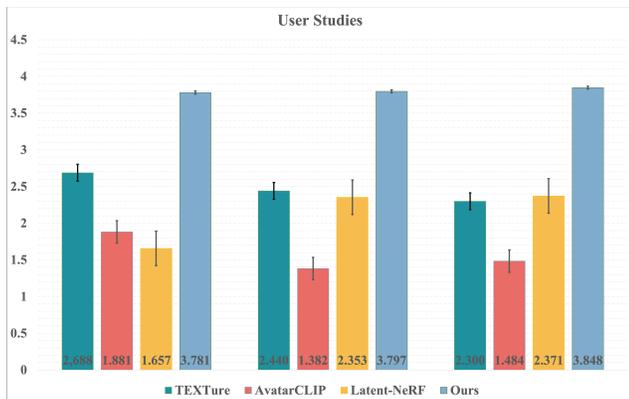


Figure 7. User studies on rotated 3D human avatars.

ometry and texture from all these examples from different viewpoints. It demonstrates the capability of our method to generate different types of avatars, *e.g.*, human beings, superheroes, comic characters, painted avatars, robots, etc.

Comparison with SOTA methods We provide qualitative comparisons with existing SOTA methods in Fig. 3. For a fair comparison, we use the same 3D canonical SMPL model for Latent-NeRF, TEXTure, and Ours, while using the 3D model provided by AvatarCLIP for better results from their method. We can observe that our method con-

sistently achieves notably better results in terms of both geometry and texture. Meanwhile, our method has the capability to generate text-consistent 3D geometry, geometry-consistent texture, and maintain the specified pose from the 3D shape prior. See more comparisons in supplementary.

Avatar generation under different poses We validate the effectiveness of our method for generating 3D human avatars in various poses, which is not achievable by other existing methods (see Fig. 4) due to the absence of the shape prior. Furthermore, our DreamAvatar can maintain high-quality texture and geometry for extreme poses, *e.g.*, self-occlusion pose and crouching, thanks to our dual space design. See supplementary for more results.

Text manipulation on avatar generation We further explore the capabilities of DreamAvatar by editing the text prompt for controlled generation (see Fig. 5), *e.g.*, by specifying age, color, and body type. Our method can generate faithful avatars fulfilling the text with additional descriptive information.

Shape modification via SMPL shape parameters We further demonstrate the possibility of our method to generate different size of 3D human avatars, *e.g.*, thin, short, tall, fat, by editing the SMPL shape parameters in Fig. 6.



Figure 8. Analysis of the setup for SMPL-derived density $\bar{\sigma}_c, \bar{\sigma}_o$, and normal-consistency loss \mathcal{L}_n . I: Ours, II: Ours without $\bar{\sigma}_c, \bar{\sigma}_o$, III: Ours without \mathcal{L}_n .

4.2. User Studies

To further validate the quality of the generation results, we conduct user studies to compare four different SOTA methods: Latent-NeRF, TEXTure, AvatarCLIP, and our DreamAvatar. We ask 23 volunteers to rank these methods in terms of (1) geometry quality, (2) texture quality, (3) consistency with the text. We randomly select 30 generated results (in the form of rendered rotating videos) and ask each volunteer to score the methods from 1 (worst) to 4 (best) for each example and each aspect. We provide the final rates in Fig. 7. Our method achieves the best rank among those three aspects, and the standard errors demonstrate that our method can consistently and robustly perform.

4.3. Further Analysis

Effectiveness of SMPL-derived density $\bar{\sigma}_c, \bar{\sigma}_o$ We ablate the NeRF models optimized without the SMPL-derived density $\bar{\sigma}_c, \bar{\sigma}_o$ as the basis for density prediction. We find that without using $\bar{\sigma}_c, \bar{\sigma}_o$, (1) the generated avatars are of low-quality geometry and texture with strong outliers and (2) the generated shapes are not constrained to be reasonable human bodies and not view-consistent (see Fig. 8).

Effectiveness of normal consistency loss \mathcal{L}_n We disable the normal-related loss term \mathcal{L}_n , and present the comparisons in Fig. 8. As can be seen, the generation quality, especially the detail of geometry and texture, notably decreases, validating the effectiveness of the introduced normal-consistency regularization (see Fig. 8).

Effectiveness of dual space design To validate our design, we experiment with two degenerated versions of our framework: (1) only the canonical space (Fig. 9 (a)), (2) only the observation space without deformation field (Fig. 9 (b, c)). We clearly see that none of these degenerated de-

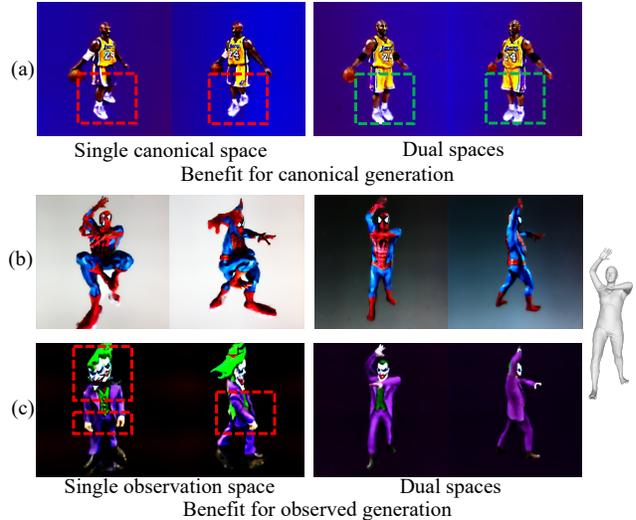


Figure 9. Analysis of dual space design. (a): single canonical space; (b), (c): single observation space w/o the deformation field. Prompt: (a) “Kobe Bryant, NBA player”, (b) “A standing spider-man”, (c) “Joker”.

signs can perform as well as our dual space design. In other words, our dual space design can improve the robustness of the generation in canonical space, have the capability to control the generated pose, and ensure the consistency between geometry and texture.

5. Conclusions

In this paper, we have introduced DreamAvatar, an effective framework for text-and-shape guided 3D human avatar generation. In DreamAvatar, we propose to leverage the parametric SMPL model to provide shape prior, guiding the generation with a rough pose and shape. We also propose a dual space design, allowing the generation to be well-constrained in the canonical space and the observation space through a learnable deformation field. Additionally, we also introduce a normal-consistency loss to facilitate the generation with text-consistent geometry and geometry-consistent texture. Extensive experiments show that our method has achieved state-of-the-art 3D human avatar generation.

Limitations Although DreamAvatar can generate sensationally high-quality clothing geometry and texture, the 3D mesh extracted from NeRF is still relatively coarse, limiting the performance. Moreover, The quality of face texture is sometimes random, resulting less realistic face generation.

Societal impact As the generation matures for both geometry and texture for human avatars, it could decrease the workload and open the gate for the next generation of the metaverse. However, there is still a danger of nefarious use of this technology to generate plausible renderings of individuals. We encourage the usage and research to be performed in an open and transparent way.

Acknowledgement This work is partially supported by Hong Kong Research Grant Council - Early Career Scheme (Grant No. 27208022) and HKU Seed Fund for Basic Research.

References

- [1] Panos Achlioptas, Olga Diamanti, Ioannis Mitliagkas, and Leonidas Guibas. Learning representations and generative models for 3d point clouds. In *International conference on machine learning*, 2018. 3
- [2] Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Karsten Kreis, Miika Aittala, Timo Aila, Samuli Laine, Bryan Catanzaro, et al. ediffi: Text-to-image diffusion models with an ensemble of expert denoisers. *arXiv preprint arXiv:2211.01324*, 2022. 3
- [3] Alexander W. Bergman, Petr Kellnhofer, Wang Yifan, Eric R. Chan, David B. Lindell, and Gordon Wetzstein. Generative neural articulated radiance fields. In *Advances in Neural Information Processing Systems*, 2022. 3
- [4] Federica Bogo, Angjoo Kanazawa, Christoph Lassner, Peter Gehler, Javier Romero, and Michael J Black. Keep it smpl: Automatic estimation of 3d human pose and shape from a single image. In *European Conference on Computer Vision*, 2016. 4
- [5] Yukang Cao, Guanying Chen, Kai Han, Wenqi Yang, and Kwan-Yee K Wong. Jiff: Jointly-aligned implicit face function for high quality single view clothed human reconstruction. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2022. 1
- [6] Eric R. Chan, Connor Z. Lin, Matthew A. Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas Guibas, Jonathan Tremblay, Sameh Khamis, Tero Karras, and Gordon Wetzstein. Efficient geometry-aware 3D generative adversarial networks. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2022. 3
- [7] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 3
- [8] Alvaro Collet, Ming Chuang, Pat Sweeney, Don Gillett, Dennis Evseev, David Calabrese, Hugues Hoppe, Adam G. Kirk, and Steve Sullivan. High-quality streamable free-viewpoint video. *ACM Transactions on Graphics*, 2015. 1
- [9] Massimiliano Favalli, Alessandro Fornaciai, Ilaria Isola, Simone Tarquini, and Luca Nannipieri. Multiview 3d reconstruction in geosciences. *Computers Geosciences*, 2012. 1
- [10] Jianglin Fu, Shikai Li, Yuming Jiang, Kwan-Yee Lin, Chen Qian, Chen Change Loy, Wayne Wu, and Ziwei Liu. Stylegan-human: A data-centric odyssey of human generation. In *European Conference on Computer Vision*, 2022. 3
- [11] Matheus Gadelha, Subhransu Maji, and Rui Wang. 3d shape induction from 2d views of multiple objects. In *International Conference on 3D Vision*, 2017. 3
- [12] Philipp Henzler, Niloy J Mitra, and Tobias Ritschel. Escaping plato’s cave: 3d shape from adversarial rendering. In *International Conference on Computer Vision*, 2019. 3
- [13] Fangzhou Hong, Zhaoxi Chen, Yushi LAN, Liang Pan, and Ziwei Liu. EVA3d: Compositional 3d human generation from 2d image collections. In *International Conference on Learning Representations*, 2023. 3
- [14] Fangzhou Hong, Mingyuan Zhang, Liang Pan, Zhonggang Cai, Lei Yang, and Ziwei Liu. Avatarclip: Zero-shot text-driven generation and animation of 3d avatars. *ACM Transactions on Graphics*, 2022. 3, 6, 12, 17
- [15] Ajay Jain, Ben Mildenhall, Jonathan T. Barron, Pieter Abbeel, and Ben Poole. Zero-shot text-guided object generation with dream fields. 2022. 3
- [16] Suyi Jiang, Haoran Jiang, Ziyu Wang, Haimin Luo, Wenzheng Chen, and Lan Xu. Humangen: Generating human radiance fields with explicit priors. *arXiv preprint arXiv:2212.05321*, 2022. 3
- [17] Nasir Mohammad Khalid, Tianhao Xie, Eugene Belilovsky, and Popa Tiberiu. Clip-mesh: Generating textured meshes from text using pretrained image-text models. 2022. 3
- [18] Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. Magic3d: High-resolution text-to-3d content creation. *arXiv preprint arXiv:2211.10440*, 2022. 1, 3, 4, 6
- [19] Lingjie Liu, Marc Habermann, Viktor Rudnev, Kripasindhu Sarkar, Jiatao Gu, and Christian Theobalt. Neural actor: Neural free-view synthesis of human actors with pose control. *ACM SIGGRAPH Asia*, 2021. 3
- [20] Sebastian Lunz, Yingzhen Li, Andrew Fitzgibbon, and Nate Kushman. Inverse graphics gan: Learning to generate 3d shapes from unstructured 2d data. *arXiv preprint arXiv:2002.12674*, 2020. 3
- [21] Shitong Luo and Wei Hu. Diffusion probabilistic models for 3d point cloud generation. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2021. 3
- [22] Gal Metzger, Elad Richardson, Or Patashnik, Raja Giryes, and Daniel Cohen-Or. Latent-nerf for shape-guided generation of 3d shapes and textures. *arXiv preprint arXiv:2211.07600*, 2022. 1, 3, 4, 5, 6, 12, 14, 17
- [23] Oscar Michel, Roi Bar-On, Richard Liu, Sagie Benaim, and Rana Hanocka. Text2mesh: Text-driven neural stylization for meshes. In *IEEE Conference on*

- Computer Vision and Pattern Recognition*, June 2022. 3
- [24] 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 *European Conference on Computer Vision*, 2020. 1, 4, 6
- [25] Kaichun Mo, Paul Guerrero, Li Yi, Hao Su, Peter Wonka, Niloy Mitra, and Leonidas Guibas. StructureNet: Hierarchical graph networks for 3d shape generation. *ACM Transactions on Graphics*, 2019. 3
- [26] Chong Mou, Xintao Wang, Liangbin Xie, Jian Zhang, Zhongang 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
- [27] Alex Nichol, Heewoo Jun, Prafulla Dhariwal, Pamela Mishkin, and Mark Chen. Point-e: A system for generating 3d point clouds from complex prompts. *arXiv preprint arXiv:2212.08751*, 2022. 3
- [28] Atsuhiko Noguchi, Xiao Sun, Stephen Lin, and Tatsuya Harada. Neural articulated radiance field. In *International Conference on Computer Vision*, 2021. 3
- [29] Atsuhiko Noguchi, Xiao Sun, Stephen Lin, and Tatsuya Harada. Unsupervised learning of efficient geometry-aware neural articulated representations. In *European Conference on Computer Vision*, 2022. 3
- [30] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed AA Osman, Dimitrios Tzionas, and Michael J Black. Expressive body capture: 3d hands, face, and body from a single image. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 4
- [31] Sida Peng, Yuanqing Zhang, Yinghao Xu, Qianqian Wang, Qing Shuai, Hujun Bao, and Xiaowei Zhou. Neural body: Implicit neural representations with structured latent codes for novel view synthesis of dynamic humans. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2021. 3
- [32] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. DreamFusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022. 1, 3, 4, 5, 6
- [33] 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*, 2021. 3
- [34] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022. 3
- [35] Elad Richardson, Gal Metzer, Yuval Alaluf, Raja Giryes, and Daniel Cohen-Or. Texture: Text-guided texturing of 3d shapes. *arXiv preprint arXiv:2302.01721*, 2023. 1, 3, 6, 12, 17
- [36] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2022. 1, 3, 4, 6
- [37] 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. 2022. 3
- [38] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022. 3, 4
- [39] Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao Li. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In *International Conference on Computer Vision*, 2019. 1
- [40] Aditya Sanghi, Hang Chu, Joseph G Lambourne, Ye Wang, Chin-Yi Cheng, Marco Fumero, and Kamal Rahimi Malekshani. Clip-forge: Towards zero-shot text-to-shape generation. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2022. 3
- [41] Jiayang Tang. A pytorch implementation of the text-to-3d model dreamfusion, powered by the stable diffusion text-to-2d model. <https://github.com/ashawkey/stable-dreamfusion>. 6, 12, 17
- [42] Dor Verbin, Peter Hedman, Ben Mildenhall, Todd Zickler, Jonathan T Barron, and Pratul P Srinivasan. Ref-nerf: Structured view-dependent appearance for neural radiance fields. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2022. 3, 5
- [43] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *Conference on Neural Information Processing Systems*, 2021. 6
- [44] Chung-Yi Weng, Brian Curless, Pratul P Srinivasan, Jonathan T Barron, and Ira Kemelmacher-Shlizerman. Humannerf: Free-viewpoint rendering of moving people from monocular video. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2022. 3, 5
- [45] Jiajun Wu, Chengkai Zhang, Tianfan Xue, Bill Freeman, and Josh Tenenbaum. Learning a probabilistic latent space of object shapes via 3d generative-adversarial modeling. *Advances in Neural Information Processing Systems*, 29, 2016. 3

- [46] Yuliang Xiu, Jinlong Yang, Xu Cao, Dimitrios Tzionas, and Michael J. Black. ECON: Explicit Clothed humans Obtained from Normals. In *IEEE-Learning Transferable Visual Models From Natural Language Supervision Conference on Computer Vision and Pattern Recognition*, 2023. 1
- [47] Jiale Xu, Xintao Wang, Weihao Cheng, Yan-Pei Cao, Ying Shan, Xiaohu Qie, and Shenghua Gao. Dream3d: Zero-shot text-to-3d synthesis using 3d shape prior and text-to-image diffusion models. *arXiv preprint arXiv:2212.14704*, 2022. 3, 5
- [48] Guandao Yang, Xun Huang, Zekun Hao, Ming-Yu Liu, Serge Belongie, and Bharath Hariharan. Point-flow: 3d point cloud generation with continuous normalizing flows. In *International Conference on Computer Vision*, 2019. 3
- [49] Xiaohui Zeng, Arash Vahdat, Francis Williams, Zan Gojcic, Or Litany, Sanja Fidler, and Karsten Kreis. Lion: Latent point diffusion models for 3d shape generation. In *Advances in Neural Information Processing Systems*, 2022. 3
- [50] Lvmin Zhang and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. *arXiv preprint arXiv:2302.05543*, 2023. 3
- [51] Yuxuan Zhang, Wenzheng Chen, Huan Ling, Jun Gao, Yinan Zhang, Antonio Torralba, and Sanja Fidler. Image gans meet differentiable rendering for inverse graphics and interpretable 3d neural rendering. 2021. 3
- [52] Linqi Zhou, Yilun Du, and Jiajun Wu. 3d shape generation and completion through point-voxel diffusion. In *International Conference on Computer Vision*, 2021. 3

A. Implementation details

To implement our method, we utilize the grid frequency encoder $\gamma(\cdot)$ from open-source Stable DreamFusion [41], which maps the input $\mathbf{x} \in \mathbb{R}^3$ onto a higher-frequency dimension, resulting in $\gamma(\mathbf{x}) \in \mathbb{R}^{32}$. The multilayer perceptron (MLP) within our NeRF model is comprised of three layers with dimensions [32, 64, 64, 3+1+4], where the channel ‘3’, ‘1’, and ‘4’ correspond to predicted normals, density values, and latent color features, respectively. A similar MLP architecture, [32, 64, 64, 3], is applied to address the non-rigid motion. The learning rate and λ_n are set to 0.001 and 0.0005, respectively, to ensure effective training of our approach.

Training complexity The reported results are achieved after being trained for 10,000 epochs on a single 2080Ti GPU for approximately 2 hours. From our experiments, Latent-NeRF [22] requires roughly 1 hour for 10,000 epochs, AvatarCLIP [14] takes around 5 hours for 30,000 epochs, and the open-source Stable DreamFusion [41] needs approximately 1.6 hours for 10,000 epochs. TEXTure [35] takes less than 2 minutes for 10 epochs to obtain the results, as it only predicts the texture map without evolving the geometry.

B. Further Analysis

B.1. Effects of non-rigid motion

The deformation field that we have implemented serves a crucial role in connecting the canonical and observation spaces. It comprises two distinct components: articulate deformation and non-rigid motion. In order to showcase the effectiveness of our design, we disable the non-rigid motion and present the comparisons in Fig. 10. The results indicate that non-rigid motion can be instrumental in eliminating outliers and improving the robustness of pose control.

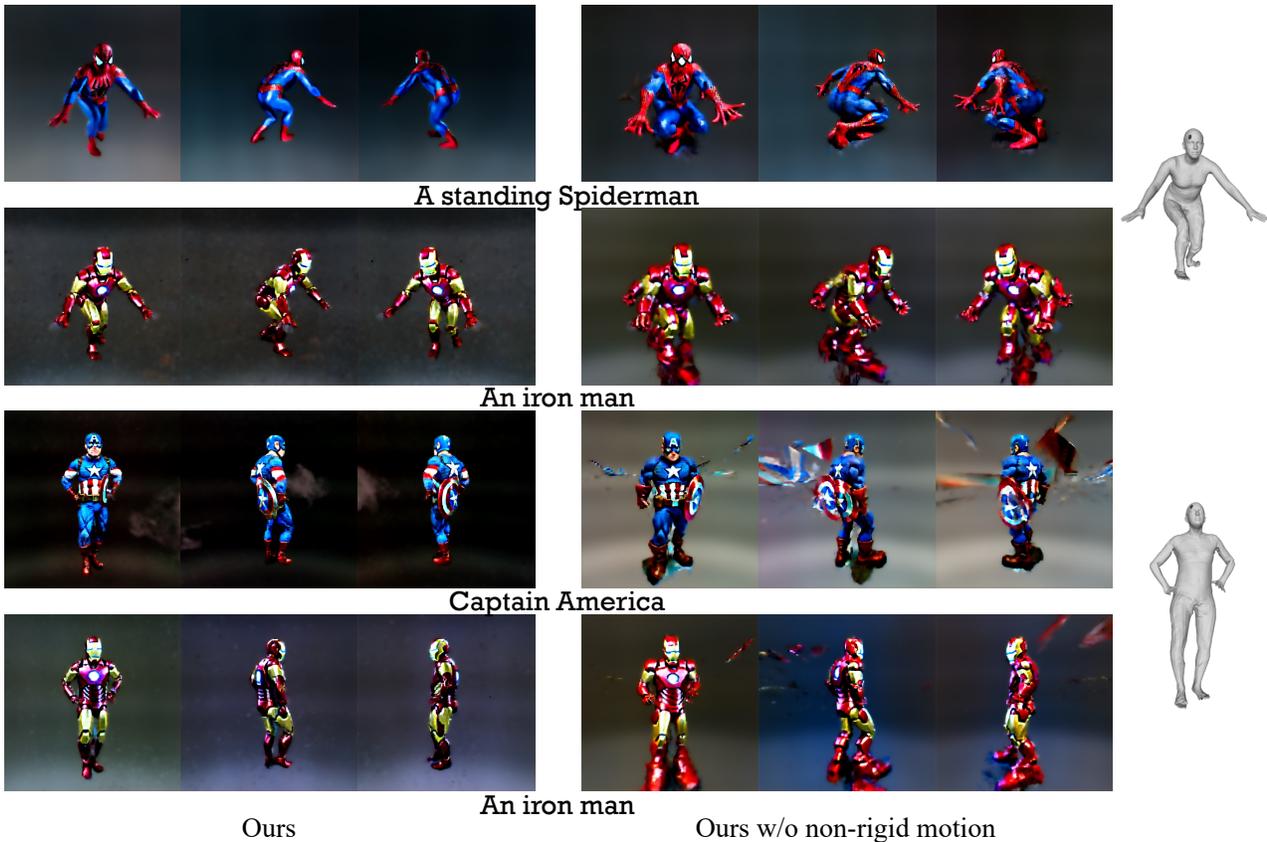


Figure 10. Effects of non-rigid motion. Our non-rigid motion component can eliminate the outliers for robust pose control.

B.2. Effects of dual spaces

In addition to the analysis in the main paper, we here offer more analysis of our dual-space design. Firstly we compare the setups between ‘Dual spaces’ and ‘Observation space w/ deformation field’ in Fig. 11. Based on these comparisons, we can conclude that our network employing the dual space design can enhance the texture quality, especially for the challenging poses with self-occlusion, and improve the robustness for pose control.

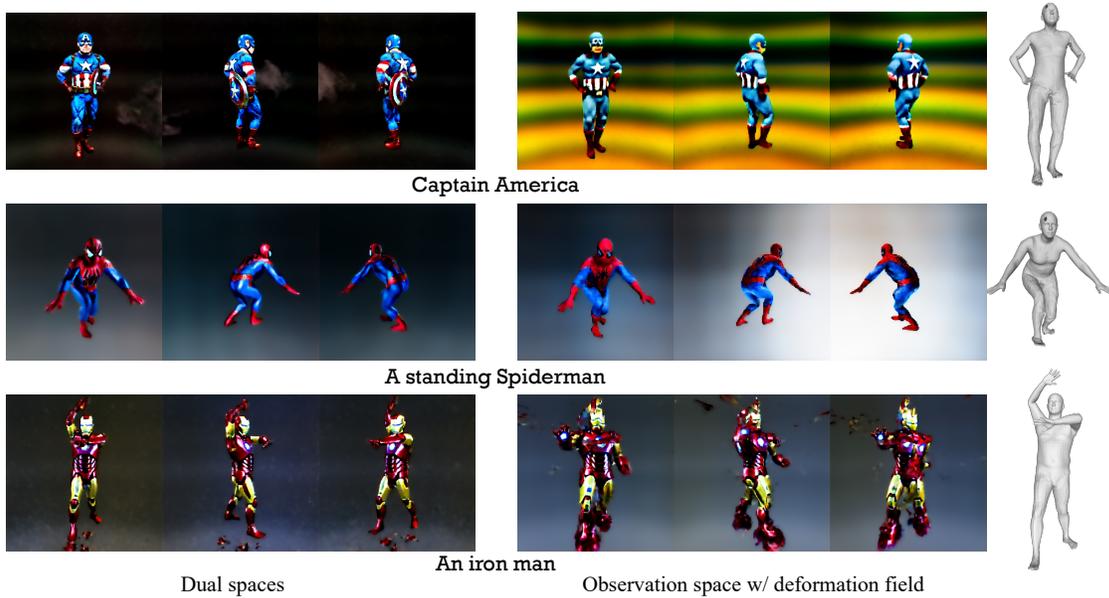


Figure 11. Comparison between our dual space design and the single observation space design with deformation field. Our dual space design largely improves the texture quality and robustness.

In Fig. 12, we then provide more comparison between ‘Dual spaces’ and ‘Single canonical space’ to further demonstrate that our dual-space design will improve the robustness of the generation in canonical space.

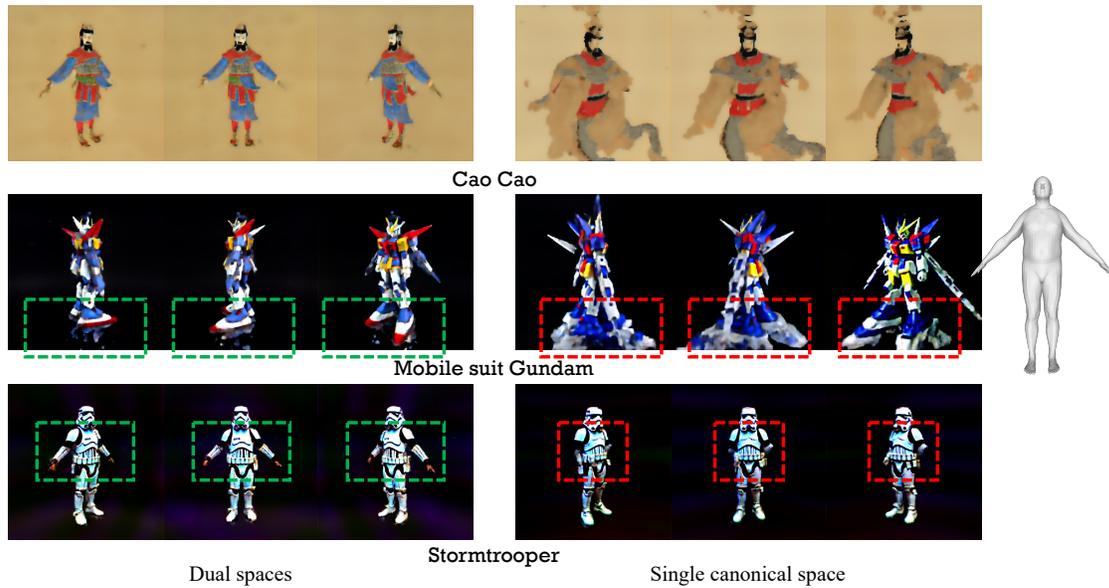


Figure 12. Comparison between our dual space design and the single canonical space design. Our dual space design improves the robustness of the generation in the canonical space.

B.3. Effects of Stable Diffusion model

Following Latent-NeRF [22], we adopt a color representation in latent space instead of directly predicting RGB values in 3-dimensional space. More precisely, we estimate the $64 \times 64 \times 4$ latent color features, which are subsequently decoded to produce RGB images at a resolution of $512 \times 512 \times 3$. This design enables the generation of higher-resolution textures and significantly enhances training efficiency See Fig. 13.

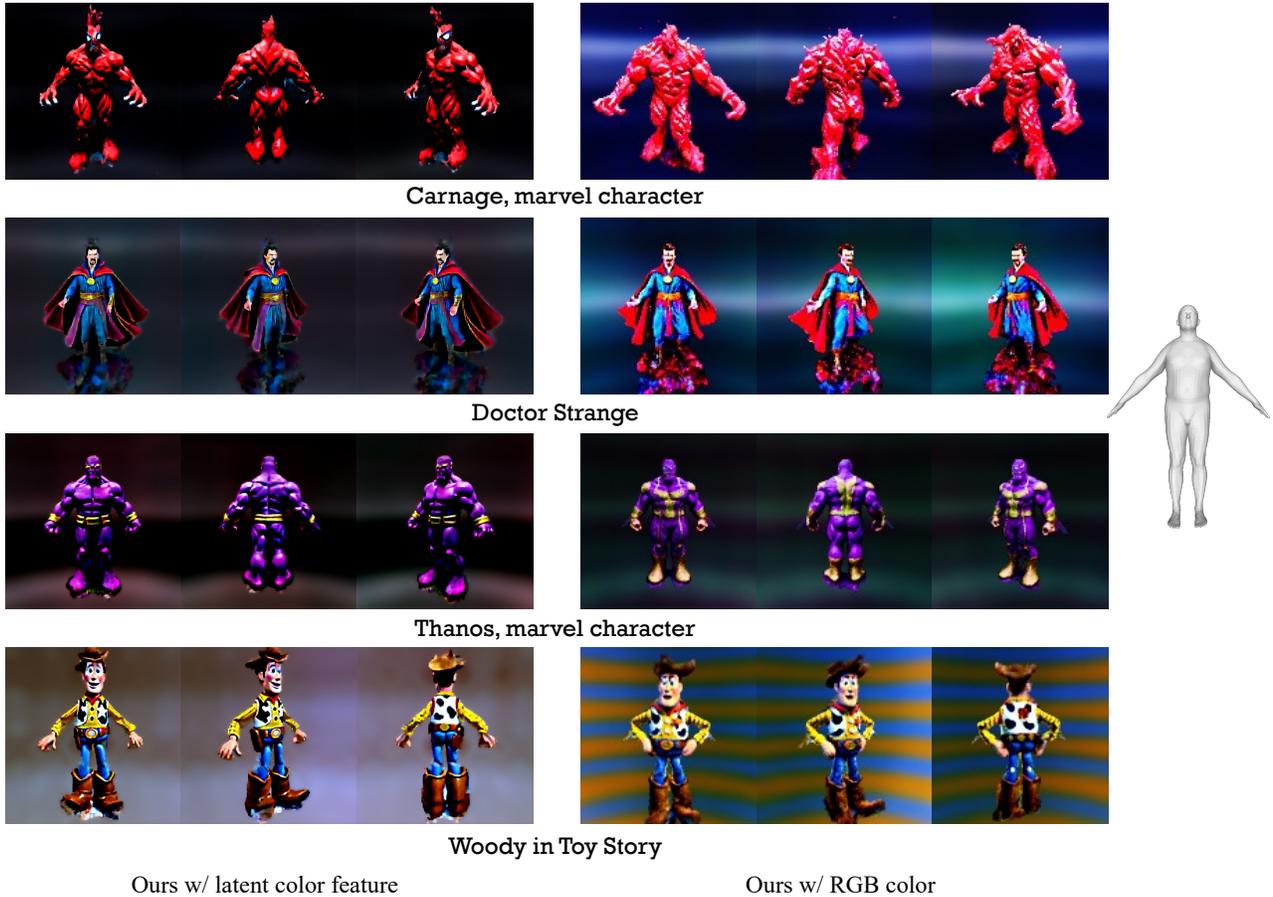


Figure 13. Comparison between ‘ours with latent color feature’ and ‘ours with RGB color’. The latent color feature is significantly more effective in generating high-resolution textures.

C. Qualitative comparisons

C.1. Comparison on different poses

We further compare our method with Latent-NeRF and TEXTure, which also have the ability to take in 3D priors at various poses. The results of this comparison, presented in Fig. 14-Fig. 16, indicate that our proposed method can effectively handle varied poses, while TEXTure is susceptible to self-occlusion issues and Latent-NeRF faces challenges in perceiving 3D information. Note that TEXTure only predicts the texture without evolving the geometry.

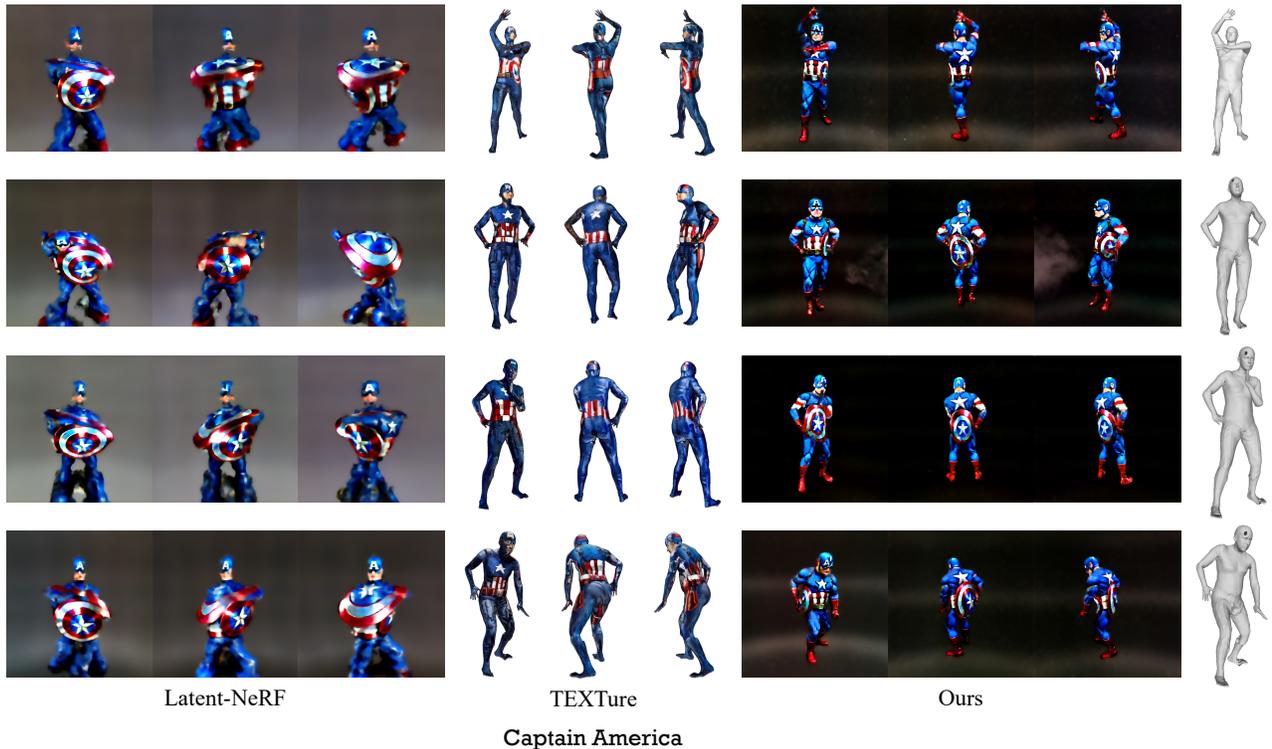


Figure 14. Qualitative comparison between our method and SOTA methods at different poses.

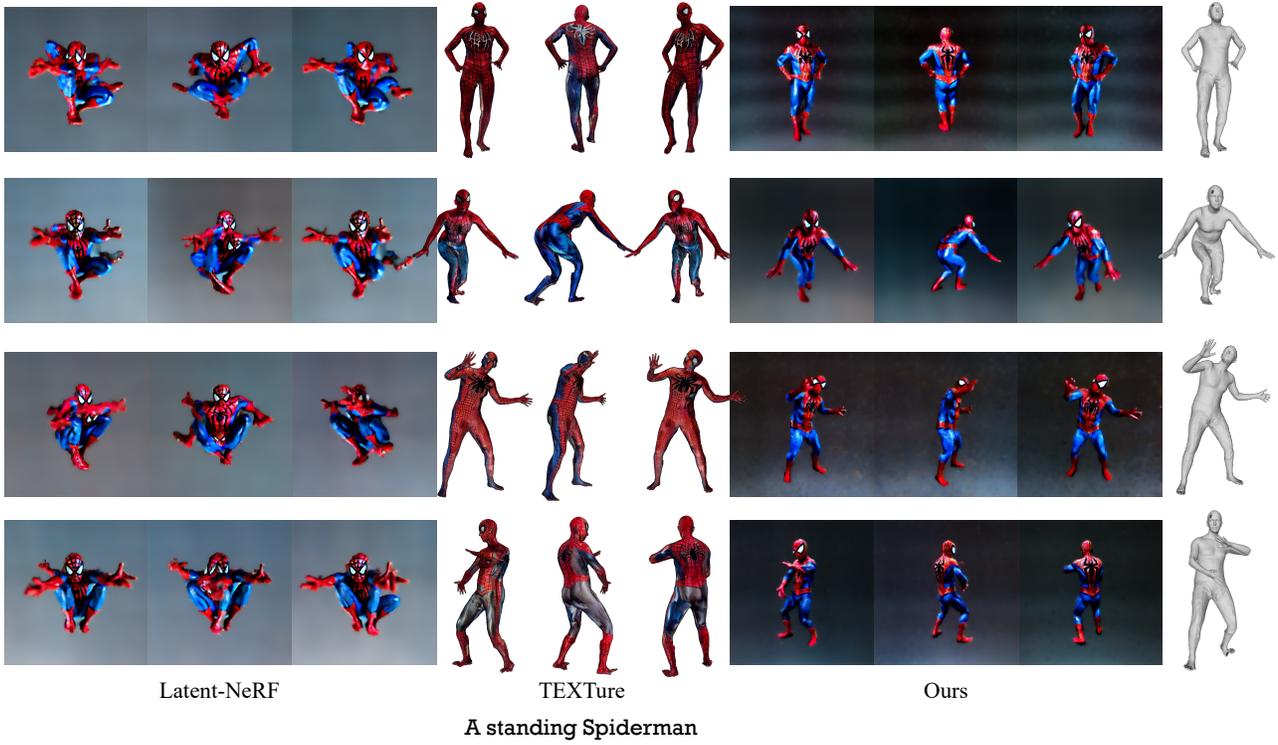


Figure 15. Qualitative comparison between our method and SOTA methods at different poses.

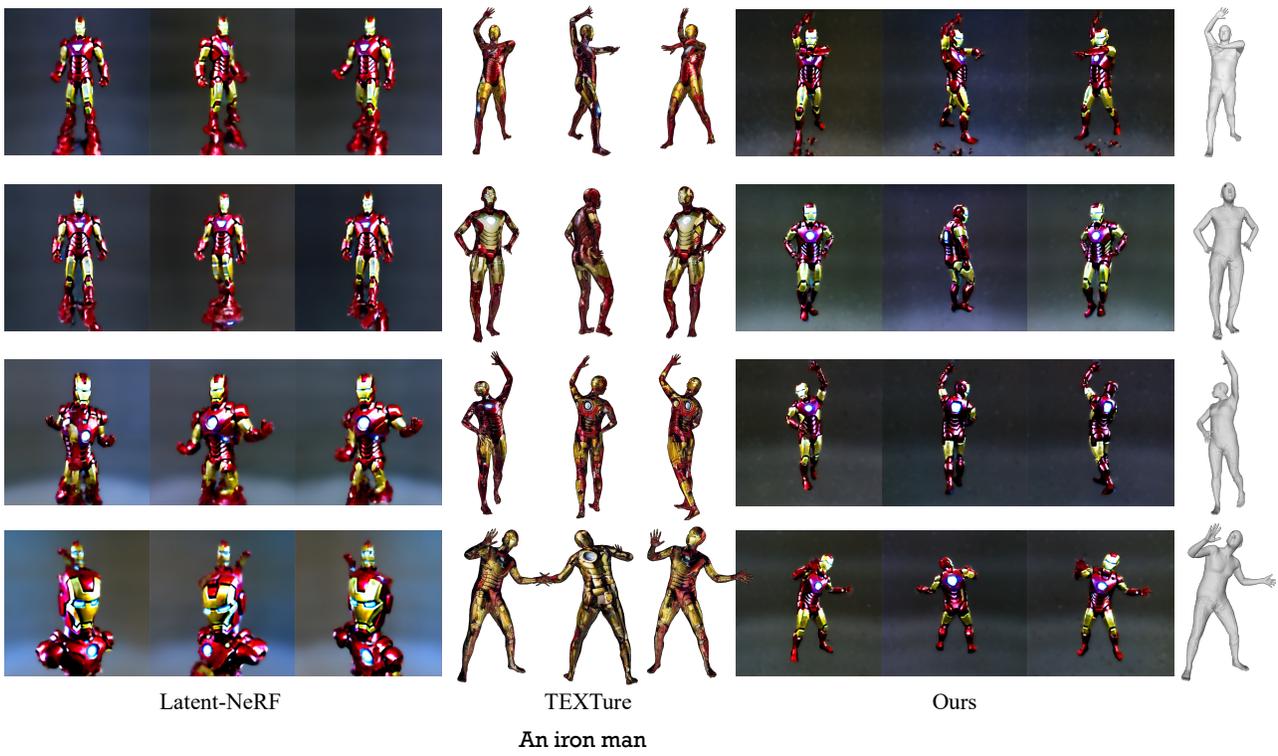


Figure 16. Qualitative comparison between our method and SOTA methods at different poses.

C.2. More comparison with SOTA methods

Here, we provide more comparison with SOTA methods, including Latent-NeRF [22], TEXTure [35], and Avatar-CLIP [14]. DreamFusion [41] is not compared due to the absence of the official implementation. Our method consistently outperforms the other competing methods in terms of geometry, texture, and consistency with text (see Fig. 17 - Fig. 19).

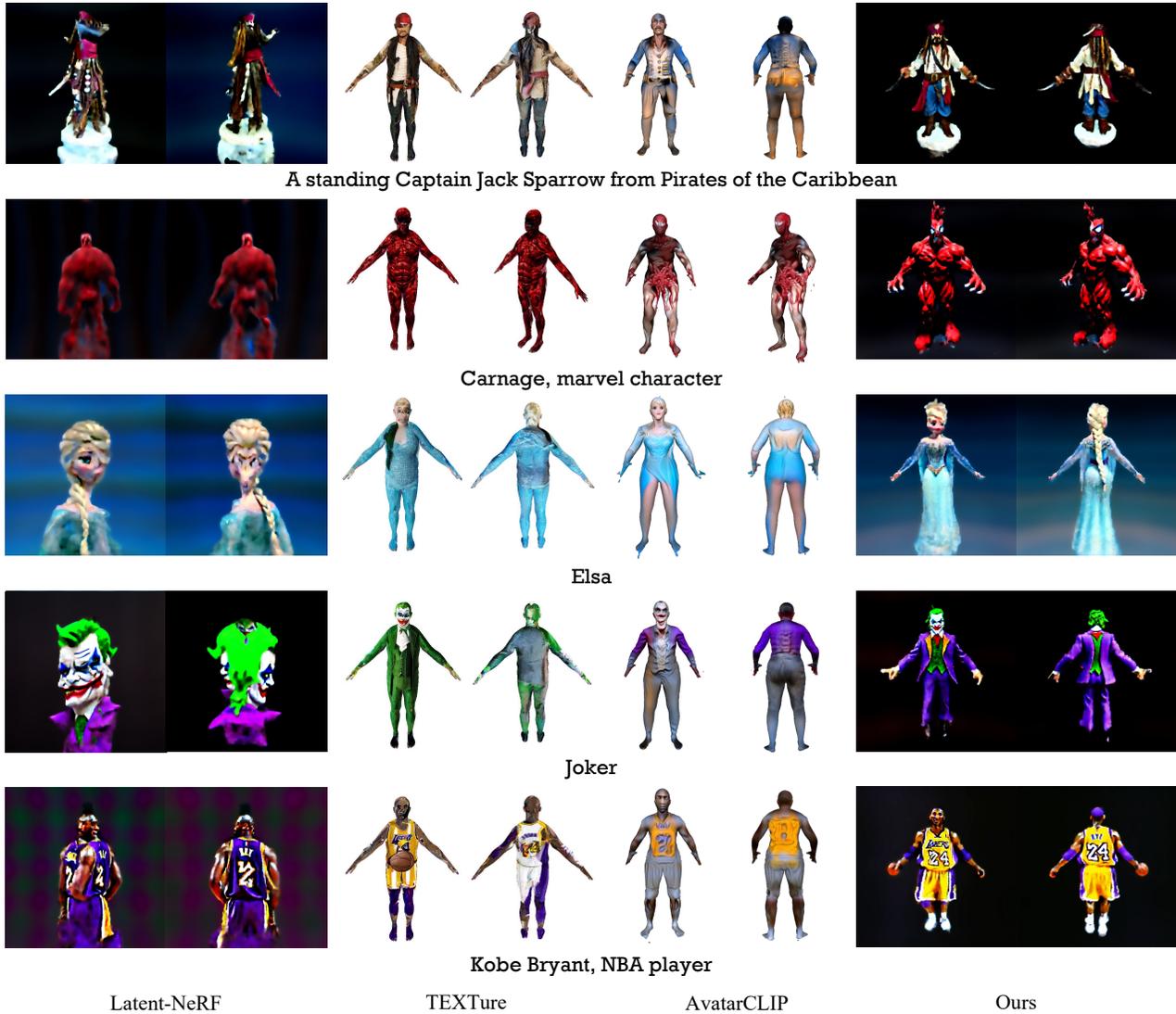
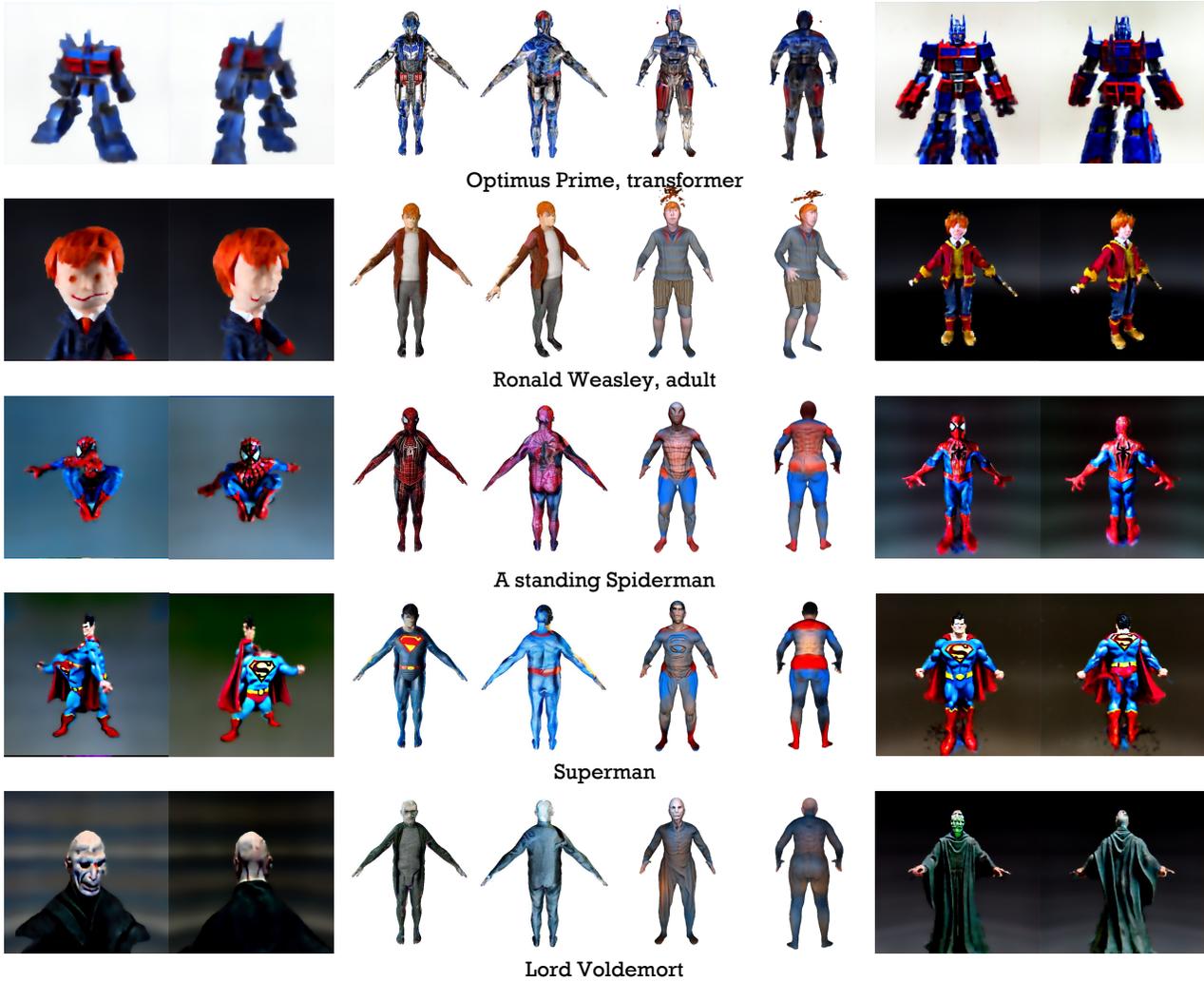


Figure 17. Qualitative comparison between our method and SOTA methods with various prompts.



Optimus Prime, transformer

Ronald Weasley, adult

A standing Spiderman

Superman

Lord Voldemort

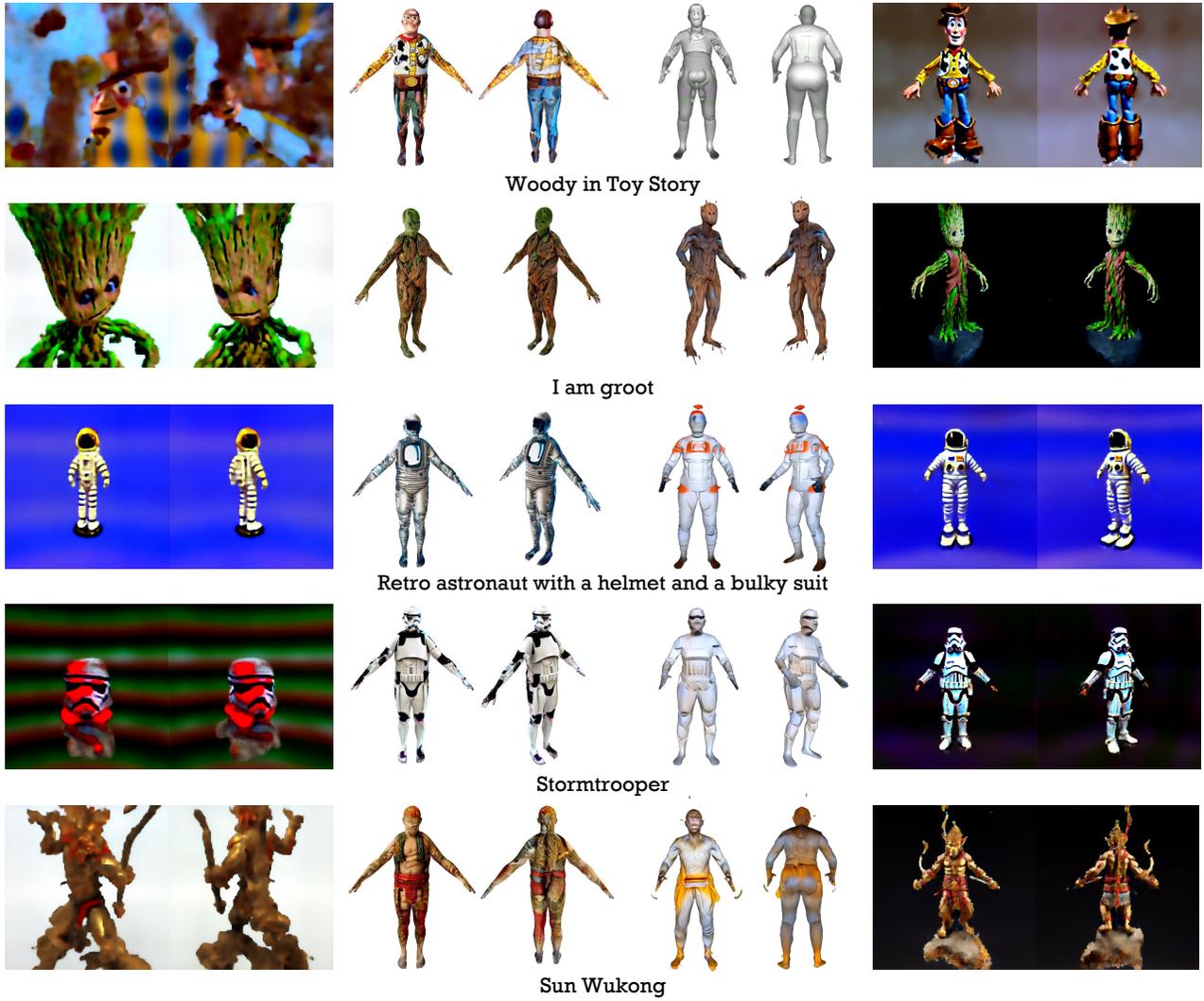
Latent-NeRF

TEXTure

AvatarCLIP

Ours

Figure 18. Qualitative comparison between our method and SOTA methods with various prompts.



Latent-NeRF

TEXTure

AvatarCLIP

Ours

Figure 19. Qualitative comparison between our method and SOTA methods with various prompts.