

DreamFace: Progressive Generation of Animatable 3D Faces under Text Guidance

LONGWEN ZHANG*, ShanghaiTech University and Deemos Technology Co., Ltd., China

QIWEI QIU*, ShanghaiTech University and Deemos Technology Co., Ltd., China

HONGYANG LIN*, ShanghaiTech University and Deemos Technology Co., Ltd., China

QIXUAN ZHANG, ShanghaiTech University and Deemos Technology Co., Ltd., China

CHENG SHI, ShanghaiTech University, China

WEI YANG, Huazhong University of Science and Technology, China

YE SHI, ShanghaiTech University, China

SIBEI YANG, ShanghaiTech University, China

LAN XU, ShanghaiTech University and Shanghai Engineering Research Center of Intelligent Vision and Imaging, China

JINGYI YU, ShanghaiTech University and Shanghai Engineering Research Center of Intelligent Vision and Imaging, China



Fig. 1. DreamFace generates personalized 3D physically-based facial assets under text guidance, which are compatible with the existing CG pipeline, with desired shapes, textures, and fine-grained animations for realistic rendering. Go to DreamFace project page <https://sites.google.com/view/dreamface>, watch our video at <https://youtu.be/yCuvzgGMvPM> and experience DreamFace online at <https://hyperhuman.top> !

Emerging Metaverse applications demand accessible, accurate and easy-to-use tools for 3D digital human creations in order to depict different cultures and societies as if in the physical world. Recent large-scale vision-language advances pave the way for novices to conveniently customize 3D content.

*Equal contributions.

Authors' addresses: Longwen Zhang, ShanghaiTech University and Deemos Technology Co., Ltd., Shanghai, China, zhanglw2@shanghaitech.edu.cn; Qiwei Qiu, ShanghaiTech University and Deemos Technology Co., Ltd., Shanghai, China, qiugw@shanghaitech.edu.cn; Hongyang Lin, ShanghaiTech University and Deemos Technology Co., Ltd., Shanghai, China, linhy@shanghaitech.edu.cn; Qixuan Zhang, ShanghaiTech University and Deemos Technology Co., Ltd., Shanghai, China, zhangqx1@shanghaitech.edu.cn; Cheng Shi, ShanghaiTech University, Shanghai, China, shicheng2022@shanghaitech.edu.cn; Wei Yang, Huazhong University of Science and Technology, Wuhan, China, weiyangcs@hust.edu.cn; Ye Shi, ShanghaiTech University, Shanghai, China, shiye@shanghaitech.edu.cn; Sibe Yang, ShanghaiTech University, Shanghai, China, yangsb@shanghaitech.edu.cn; Lan Xu, ShanghaiTech University and Shanghai Engineering Research Center of Intelligent Vision and Imaging, Shanghai, China, xulan1@shanghaitech.edu.cn; Jingyi Yu, ShanghaiTech University and Shanghai Engineering Research Center of Intelligent Vision and Imaging, Shanghai, China, yujingyi@shanghaitech.edu.cn.

However, the generated CG-friendly assets still cannot represent the desired facial traits for human characteristics. In this paper, we present DreamFace, a progressive scheme to generate personalized 3D faces under text guidance. It enables layman users to naturally customize 3D facial assets that are compatible with CG pipelines, with desired shapes, textures and fine-grained animation capabilities. From a text input to describe the facial traits, we first introduce a coarse-to-fine scheme to generate the neutral facial geometry with a unified topology. We employ a selection strategy in the CLIP embedding space to generate coarse geometry, and subsequently optimize both the detailed displacements and normals using Score Distillation Sampling (SDS) from the generic Latent Diffusion Model (LDM). Then, for neutral appearance generation, we introduce a dual-path mechanism, which combines the generic LDM with a novel texture LDM to ensure both the diversity and textural specification in the UV space. We also employ a two-stage optimization to perform SDS in both the latent and image spaces to significantly provide compact priors for fine-grained synthesis. It also enables learning the mapping from the compact latent space into physically-based textures (diffuse albedo, specular intensity, normal maps, etc.). Our

generated neutral assets naturally support blendshapes-based facial animations, thanks to the unified geometric topology. We further improve the animation ability with personalized deformation characteristics. To this end, we learn the universal expression prior in a latent space with neutral asset conditioning using the cross-identity hypernetwork, we subsequently train a neural facial tracker from video input space into the pre-trained expression space for personalized fine-grained animation. Extensive qualitative and quantitative experiments validate the effectiveness and generalizability of DreamFace. Notably, DreamFace can generate realistic 3D facial assets with physically-based rendering quality and rich animation ability from video footage, even for fashion icons or exotic characters in cartoons and fiction movies.

CCS Concepts: • Computing methodologies → Computer graphics.

Additional Key Words and Phrases: Text-Driven Generation, 3D Digital Humans, Physically-based Facial Assets

1 INTRODUCTION

Digital human face models should equally reflect the diversity of the human experience in order to depict different cultures, societies, and environments that make up our physical world. A successful 3D digital human creation tool will hence allow users to create and customize models that follow their own identities and physical characteristics. These include skin color, hair types, facial shapes, expressions and appearance with rich facial traits. In addition to being realistic and believable, the produced assets also need to match specific themes or concepts, e.g., from as complicated as the movie script and literature, or as simple as textual descriptions by novice users. It has been in high demand for a variety of applications for feature films, game productions, and most recently, immersive experiences in the Metaverse.

However, creating these believable 3D human characters is not for anyone. Early attempts [Alexander et al. 2009; Debevec et al. 2000] generally require expensive apparatus and immense artistic expertise, and hence are limited to celebrities for feature film productions. It has been a long journey for the graphics community to democratize the accessible use of 3D facial assets to the mass crowd, equipped with powerful neural generative techniques, from variational autoencoders (VAEs) [Kingma and Welling 2014], generative adversarial networks (GANs) [Goodfellow et al. 2020] to the latest Diffusion Models [Ho et al. 2020]. Seminal generation approaches such as StyleGAN2 [Karras et al. 2020] have successfully produced highly realistic 2D facial rendering almost indistinguishable from real photos. Various methods [Chan et al. 2022; Deng et al. 2022; Or-El et al. 2022] further combine 2D GANs with 3D implicit representations, achieving 3D-aware facial synthesis with detailed hairstyles, textures or expressions. Yet, these rendering-based schemes are still difficult to integrate seamlessly into the existing CG production pipeline. A few recent works [Li et al. 2020b,a] combine production-ready facial assets into the VAE or GAN-based generation framework. But they still are not able to produce a high degree of diversity, partially due to the limited training set of 3D or paired data and the limited flexibility of the generation process. Hence, the vision of creating facial assets from novices' text prompt remains far-reaching. Fortunately, recent large-scale vision-language models [Radford et al. 2021; Rombach et al. 2022] pave the way to lower the barrier to entry for novices further. Using the advanced

diffusion models, recent attempts have achieved huge success for zero-shot 2D image content creation from text prompts [Ramesh et al. 2022; Saharia et al. 2022], or even for 3D generation[Jain et al. 2022; Lin et al. 2022; Metzer et al. 2022; Poole et al. 2022]. A few works utilize the pre-trained CLIP model [Radford et al. 2021] to generate animatable full-body avatars [Hong et al. 2022b] or facial texture maps [Aneja et al. 2022]. Despite the diversity, the generated assets from the above methods still cannot represent the desired facial traits for human characteristics, e.g., lacking detailed facial geometry, fine-grained animation or physically-based textures.

In this paper, we present DreamFace, a progressive scheme to generate personalized 3D faces with text guidance. As shown in Fig. 1, with only prompt controls, DreamFace enables layman users to customize 3D facial assets with the desired shape and physically-based textures, as well as empowered animation capabilities. The resulting assets (mesh, texture, normal maps, etc.) are also compatible with the existing graphics production pipeline, significantly democratizing the use of generated facial assets. In particular, we demonstrate that DreamFace not only benefits the CG production industry but also incentivizes numerous innovative applications.

At the core of our approach is the organic integration of the recent large-scale vision-language advances with dynamic physically-based facial assets. For high-quality and more controllable generation, we design DreamFace in a progressive framework, which consists of three sequential modules: geometry generation, physically-based texture diffusion, and animation empowerment. From a text input to describe the facial traits of desired avatar, we first introduce a coarse-to-fine scheme to generate the corresponding neutral facial geometry with the topology structure from ICT-FaceKit [Li et al. 2020a]. In the coarse stage, we select the optimal coarse geometry from a diverse candidate pool randomly pre-sampled from the shape space of ICT-FaceKit, by comparing relative matching scores between the prompt and the rendered geometry images in the CLIP embedding space. We then perform fine-grained detail carving on top of the coarse geometry regarding vertex displacements and detailed normal maps in the tangent space. Analogous to DreamFusion [Poole et al. 2022], we learn the detailed displacements and normal maps using the Score Distillation Sampling technique (SDS) to compute gradients from the generic Latent Diffusion Model (LDM), i.e., Stable Diffusion [Rombach et al. 2022]. Then, our physically-based appearance generation aims to predict the neural facial assets that are consistent with both the predicted geometry and text prompt. Specifically, we adopt a dual-path mechanism to utilize two kinds of diffusion models: one generic LDM for diverse generation ability from general prompt inputs and a novel texture LDM to ensure the textural specifications in the UV space. For training our texture LDM, we augment existing UV texture datasets with our physically-based one and subsequently adopt a prompt tuning strategy to compensate for the difference between various datasets. For efficient appearance generation, we further introduce a two-stage optimization to perform Score Distillation Sampling (SDS) in both the latent and image spaces, where the latent one significantly provides compact priors for fine-grained synthesis. Note that both SDS processes are enhanced with the generic and texture LDMs through our dual-path design. To obtain physically-based assets, we further learn the mapping from the compact latent space into

diffuse albedo, specular intensity, and normal maps, followed by a super-resolution module to generate 4K textures for high-quality rendering.

For animating the generated neutral geometry and appearance, the brute-force approach would adopt the default blendshapes provided by the parametric model ICT-FaceKit [Li et al. 2020a], since our assets share the same geometric topology with ICT-FaceKit. One could utilize existing face trackers to obtain the corresponding expression parameters from video footage so as to seamlessly animate our facial assets. We further introduce an enhanced animation scheme accompanied by our neutral assets to maintain more personalized motion characteristics than the general blendshapes. We first follow the recent work [Cao et al. 2022] to adopt the cross-identity hypernetwork, so as to learn a universal prior for modeling the expression space of the generated facial assets. We adopt the U-Net architecture to transform the neutral asset (in terms of geometry images) into the generated facial mesh under various expressions. With this universal prior with neutral asset conditioning, we train a video facial tracker as an encoder from the image space into the pre-trained expression space, enabling video-based personalized animation. Finally, we showcase the capability of DreamFace for generating realistic 3D facial assets with rich animation and physically-based rendering quality, even for fashion icons or exotic characters in cartoons and fiction movies. With our DreamFace, only through textural guidance, even novices can naturally create the human characters they have in mind and easily customize the creations for novel effects like aging and virtual makeup, and even further animate the creations using in-the-wild video footage.

To summarize, our main contributions include:

- We present DreamFace, a novel generation scheme to bridge recent vision-language models with animatable and physically-based facial assets, with progressive learning to disentangle geometry, appearance and animation ability.
- We introduce a dual-path appearance generation design to combine a novel texture diffusion model with the pre-trained one, equipped with a two-stage optimization in both the latent and image spaces.
- We demonstrate the animation ability of the generated facial assets using blendshapes or an empowered personalized scheme, and further showcase the applications of DreamFace to design human characters naturally.

2 RELATED WORKS

Face modeling and generation. The generation of 3D faces has seen rapid progress in recent years. Many research efforts have attempted to generate realistic facial modeling, which draws a solid foundation for generation research. Parametric models using Principal Component Analysis (PCA) were first used to express facial components containing geometry and textures, especially 3DMM model [Blanz and Vetter 1999]. More work extended the expressiveness of parametric models to represent the facial details faithfully [Cao et al. 2013; Li et al. 2017; Ranjan et al. 2018]. GAN [Goodfellow et al. 2020] methods demonstrated their strong generative ability to produce high-fidelity results. Seminal generation approaches such

as StyleGAN2 [Karras et al. 2020] generate the face image of a vast abundance. Extend to 3D generation, a series of methods [Chan et al. 2022; Deng et al. 2022; Liu et al. 2022; Or-El et al. 2022; Tewari et al. 2020] focused on embedding 3D priors within the GAN framework to generate 3D faces consistent with known attributes. Such generative methods model 3D faces with varieties of representation, including explicit representation based on mesh [Abrevaya et al. 2019; Gecer et al. 2021; Lattas et al. 2021; Li et al. 2020a; Liao et al. 2020] and implicit representation based on SDF, such as StyleSDF [Or-El et al. 2022], or based on NeRF [Mildenhall et al. 2021], like Head-NeRF [Hong et al. 2022a] and MofaNeRF [Zhuang et al. 2022]. Although implicit representation methods achieve remarkable detail quality and visual results, they are hard to be integrated into the existing CG production pipeline. On the other hand, even though explicit representation methods enjoy good compatibility, the over-smoothed geometry cannot faithfully reproduce detailed features. In this paper, we design DreamFace, a highly compatible scheme that can generate highly realistic and detailed facial assets. Specifically, DreamFace applies the explicit face mesh with the same geometric topology as ICT-FaceKit [Li et al. 2020a] to utilize displacement and tangent space normal maps to achieve fine-grained rendering.

Texture generation. Generating faces with mesh representation always requires high-quality facial textures. There is a large corpus of research works in the field of generative models for UV textures based on self-supervised methods [Fukamizu et al. 2019; Gecer et al. 2020, 2019, 2021; Lee et al. 2020]. Though these self-supervised methods reached relatively satisfactory results, their performance is still far from ultra-realistic rendering. In contrast, supervised methods [Bao et al. 2021; Lattas et al. 2020, 2021; Li et al. 2020a] adopted image-to-image translation frameworks [Isola et al. 2017; Wang et al. 2018] to generate another reflectance textures from a single albedo map, ensuring high-level correspondence between the generated textures. However, they still have difficulty assembling a large-scale dataset that always requires a photometric capture system [Ghosh et al. 2011; Ma et al. 2007]. Different texture standards and data quality further hinder joint training across different datasets. In addition, most of these methods are based on GANs that are neither scalable nor stable on text-conditional generation tasks. Recent advances show the unprecedented generative abilities and compositional power of Diffusion models [Ho et al. 2020; Sohl-Dickstein et al. 2015], which have not yet been developed in texture generation. To this end, we propose to train the diffusion model for producing high-quality physically-based textures with a novel prompt tuning strategy to bridge the uneven quality across all training data.

Text to 3D generation. In recent years, there has been a great deal of interest in text-driven image generation [Mansimov et al. 2016; Reed et al. 2016]. Leverage on the powerful vision-language representations enriched by CLIP [Radford et al. 2021], DALL-E2 [Ramesh et al. 2022] further demonstrates zero-shot text-driven image synthesis capabilities by excessively scaling up training data size. Many works, such as VQGAN-CLIP [Crowson et al. 2022] and GLIDE [Nichol et al. 2021], guided the generative image models by embedding the distance with the latent code of text prompts encoded by the CLIP

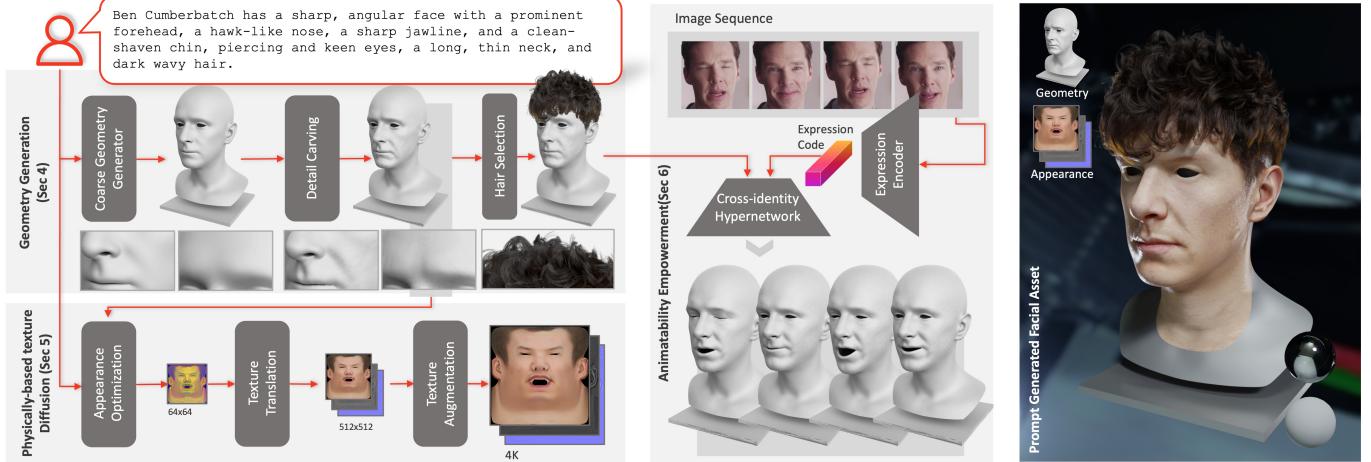


Fig. 2. The overview of DreamFace. Our pipeline mainly includes three modules, including geometry generation (Sec. 4), physically-based texture diffusion (Sec. 5), and animatability empowerment (Sec. 6). Given textual guidance, DreamFace is able to generate facial assets that closely resemble the described characteristics in terms of shape and appearance. Our approach is consistent with industry standards in computer graphics production and is able to achieve photo-realistic results when driven and rendered.

model. StyleCLIP [Patashnik et al. 2021] followed by StyleGAN-Nada [Gal et al. 2022] combined the generative power of StyleGAN and the expressive power of CLIP, enabling the image generation to have remarkable diversity from various domains. The state-of-the-art text-to-image generation method [Rombach et al. 2022] uses the Latent Diffusion Model that shifts the diffusion process to the latent space using pretrained autoencoders. This approach makes the training more stable and feasible to train on larger datasets to achieve a more remarkable richness. In contrast to the well-developed 2D image generation task, 3D generation needs more effort to generate high-quality and divisible 3D models. Sanghi et al. [2022] attempted to directly generate 3D objects by training a autoencoder with 3D representations. Further combination of the vision-language model with the rendered image from differentiable rendering on mesh representation [Chen et al. 2022; Khalid et al. 2022; Michel et al. 2022] or NeRF representation [Hong et al. 2022b; Jain et al. 2022] gives more progress to the generation and editing of 3D objects. DreamFusion [Poole et al. 2022] pioneered the introduction of the diffusion model for the supervision of generated objects but suffered for a long generation time. Magic3D [Lin et al. 2022] accelerates the DreamFusion [Poole et al. 2022] by optimizing the initial coarse model generated by low-resolution diffusion prior, using a sparse 3D hash grid structure [Müller et al. 2022]. Latent-NeRF [Metzger et al. 2022] utilized LDM to directly optimize NeRF on latent space and proved the effectiveness of optimizing latent textures on the mesh. The above methods either use CLIP or the Diffusion model to generate stunning results but still suffer from Two-faces Janus problems, resulting in insufficient quality. To handle this problem, we develop a novel dual-path optimization scheme that leverages the learned texture prior to guiding the Score Distillation Sampling and still retains powerful text-guided generation capability.

Facial Animation. Unlike implicit representation, explicit representation supports generating expressions directly by generic expression blendshapes. Some face-tracking techniques [Apple 2023; Feng et al. 2021; Fyffe et al. 2015; Somepalli et al. 2021] can provide a lightweight capture solution for facial animation with a single RGB camera input. To further enhance the inaccurate result of generic facial animation, many research efforts made huge progress. Laine et al. [2017] trained a characteristic neural network with registered 4D scans. Lombardi et al. [2018] encodes geometry and view-dependent texture information with a VAE framework that enables video-driven animation from VR Headset cameras. To allow for cross-identity retargeting, Moser et al. [2021] applied image-to-image translation and extracted a common representation between the input video and rendered CG sequence to predict blendshape weights. The state-of-the-art neural animation method [Zhang et al. 2022] proposed a production-level animation pipeline that generates high-quality facial geometry for person-specific performance. A very recent work Cao et al. [2022] proposed a Universal Prior Model (UPM), which can produce personalized expressions from unseen identities. Inspired by this work, we train a cross-identity hypernetwork together with an image expression encoder to enable personalized animation for our generated facial assets by novice users.

3 OVERVIEW

Here we introduce DreamFace, a progressive generation scheme to marry large-scale vision-language models with personalized facial assets that are compatible with CG engines. As shown in Fig. 2, from only prompt controls, our DreamFace disentangles the generation framework into three cooperated modules, named geometry generation, physically-based texture diffusion, and animation empowerment, respectively.

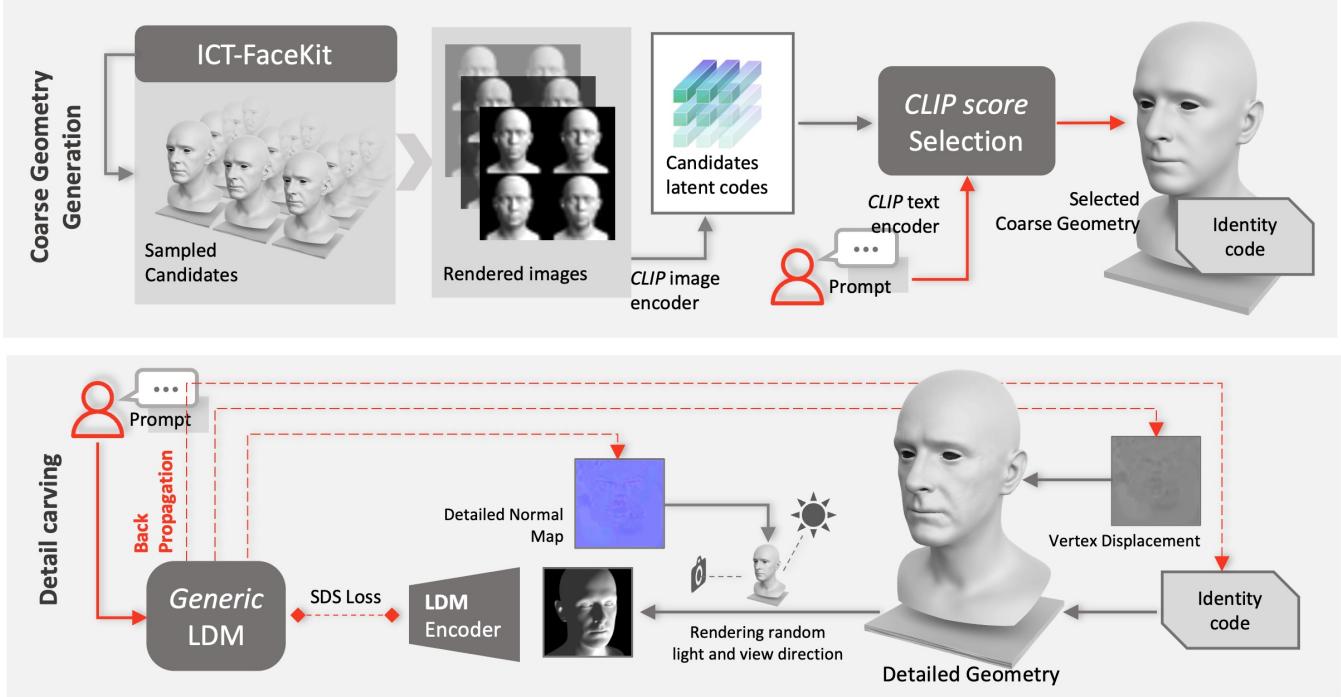


Fig. 3. Geometry generation pipeline. Given the input prompt, we utilize the CLIP model to select the coarse geometry candidates with the highest matching score. Next, we employ a generic LDM to perform SDS on the rendered images under random view and lighting conditions. This allows us to add facial details to the coarse geometry via vertex displacement and detailed normal map, resulting in a highly detailed geometry.

For geometry generation, we adopt a coarse-to-fine scheme to obtain neutral geometry with fine-grained facial traits and topology structure from the parametric model ICT-FaceKit [Li et al. 2020a] (Sec. 4). In the coarse stage, we select an optimal one from a diverse candidate pool pre-sampled from ICT-FaceKit by calculating the matching scores based on CLIP model [Radford et al. 2021]. Then, we perform fine-grained detail carving on top of the coarse geometry by applying the Score Distillation Sampling technique [Poole et al. 2022] to optimize both the detailed displacements and normal maps in the tangent space. We also adopt an efficient prompt-based hair selection to increase realism.

We further generate a physically-based appearance consistent with both the predicted neutral geometry and text prompt (Sec. 5). We introduce a dual-path mechanism with two latent diffusion models (LDMs), including a generic LDM for diversifying the generated results from arbitrary prompt inputs, as well as a novel texture LDM in the UV space. The texture LDM is trained using augmented UV texture datasets with physically-based ones, as well as prompt tuning for the compensation of data variance. We also introduce a two-stage optimization to perform the Score Distillation Sampling (SDS) in both the latent and image spaces, which are enhanced through our dual-path design. It enables well-sculpted latent space with compact priors for efficient and fine-grained synthesis. Finally, we learn the mapping from the latent space into physically-based assets and subsequently adopt a super-resolution to generate 4K textures.

For animation, our generated neutral assets naturally support existing facial trackers to obtain the expression parameters, benefiting from the by-default blendshapes and topology from ICT-FaceKit. We further empower the animation ability of our assets to model personalized deformation characteristics (Sec. 6). We employ the cross-identity hypernetwork [Cao et al. 2022] to learn the universal expression prior in a latent space, where a U-Net architecture is adopted to transform the neutral asset (as geometry images) into the generated facial meshes with diverse expressions. On top of such expression prior with neutral asset conditioning, we trained a neural video facial tracker, as an encoder from the image space into the pre-trained expression space, achieving video-based personalized fine-grained animation.

4 GEOMETRY GENERATION

From a user prompt \mathcal{P} that describes the facial characteristics, we first introduce a coarse-to-fine scheme to generate the corresponding neutral facial geometry with the topology structure from ICT-FaceKit, which consists of 14062 vertices and 28068 faces. We propose a prompt-based selection framework that chooses the optimal coarse geometry from diverse candidates randomly sampled from ICT-FaceKit shape space, which has the best CLIP matching score with user prompt. We then perform fine-grained detail carving on top of the coarse geometry in terms of vertex displacements and detailed normal maps in the tangent space. Similar to DreamFusion [Poole et al. 2022], We rely on the Score Distillation Sampling

(SDS) loss of the pretrained generic LDM, i.e., Stable Diffusion [Rombach et al. 2022], for guiding the details carving. We also use the pre-trained LDM image autoencoder from Stable Diffusion with \mathcal{E} as encoder and \mathcal{D} as decoder, which converts a 512×512 image to 64×64 latent code and back. In the following, we first describe the process of generating coarse geometry and then introduce the facial details carving process based on the coarse geometry.

4.1 Coarse geometry generation

As depicted in Fig. 3, the coarse geometry generation process involves randomly sampling candidates from the shape space of ICT-FaceKit [Li et al. 2020a]. We then render the front and left/right 3/4 views of the selected face geometry. And then, we use the CLIP model [Radford et al. 2021] to extract features from the images and select the one that best matches the input user prompt. We use the geometry corresponding to the best match as our coarse head mesh.

For sampling candidates from ICT-FaceKit, we obtain a shape space with $|\beta| = 100$ basis. The formulate of a certain shape from the shape space is:

$$\mathbf{T} = T(\boldsymbol{\beta}) = \bar{\mathbf{T}} + \sum_i \beta_i S_i, \quad (1)$$

where $\bar{\mathbf{T}}$ is the mean face, S_i is the shape components, and \mathbf{T} are corresponding generated head mesh. The coarse geometry candidates are sampled from the shape space via choose $\boldsymbol{\beta}$ from a multivariate normal distribution $\boldsymbol{\beta} \sim \mathcal{N}(0, 1)$. Then, the sampled candidates share a similar distribution that covers the span of the possible facial shapes in ICT-FaceKit.

To find the best match from candidate geometries, we first render the front and left/right 3/4 views of the selected face geometry under 10 directional lightings from different angles. This process generates 30 images in total for each candidate. We then project both images and text prompt to CLIP embeddings and calculate their matching scores. Specifically, given a face candidate i with parameter $\boldsymbol{\beta}_i$. We render the corresponding mesh $T(\boldsymbol{\beta}_i)$ using a mesh renderer $\mathcal{R}_m(\cdot)$ with the predefined 3 camera poses $c \in \mathcal{I}_m$ and 10 lightings $l \in \mathcal{J}_m$. The rendering process then is: $\mathbf{I}_i^{c,l} = \mathcal{R}_m(\boldsymbol{\beta}_i, c, l)$, where \mathbf{I} denotes the rendered images. Then we embed the images using the CLIP image encoder $\mathcal{E}_{\text{vision}}$ and average the generated latent codes for each candidate geometry (to guide the CLIP image encoder to focus more on geometry instead of appearances), which can be formulated as follows:

$$e_i = \mathbb{E}_{c,l} [\mathcal{E}_{\text{vision}}(\mathcal{R}_m(\boldsymbol{\beta}_i, c, l))], \quad (2)$$

where $\mathbb{E}(\cdot)$ represents the expected value computation. The text embedding then is simply obtained through encoding the user prompt as $e_t = \mathcal{E}_{\text{text}}(\mathcal{P})$, where $\mathcal{E}_{\text{text}}$ is the text encoder of CLIP. Instead of calculating correlations between e_i and e_t , we compute in a relative way following AvatarClip [Hong et al. 2022b] as:

$$s = \lambda_d s_d + \lambda_r s_r, \quad \text{where } s_d = \tilde{e}_i \cdot \tilde{e}_t, \quad s_r = \Delta \tilde{e}_i \cdot \Delta \tilde{e}_t \quad (3)$$

and \tilde{x} represents the normalized value of x , $\Delta e_i = e_i - \bar{e}_i$, $\Delta e_t = e_t - \bar{e}_t$, and \bar{e}_i , \bar{e}_t are the anchor embeddings, i.e., that of the anchor text, i.e., “the face”, and the mean mesh $\bar{\mathbf{T}}$. We use the combined value s as the final matching score. Finally, we select the candidate with the maximum matching score as the coarse geometry \mathbf{T}^* , where the corresponding identity code is $\boldsymbol{\beta}^*$.

4.2 Detail carving

We then carve details on the selected coarse geometry \mathbf{T}^* by optimizing additional vertex displacements and detailed normal maps in the tangent space. Recall that the coarse mesh is selected from a set of randomly generated geometries using parametric representations. The selected coarse mesh is over-smoothing and deviates from the text prompt. To increase realism and make the geometry more closely match the input prompt, we develop a detail carving process. We model the face details using two components, i.e., the vertex displacement \mathcal{V}_d and geometric detailed tangent space normal map \mathcal{N}_d , on top of the coarse mesh \mathbf{T}^* . The vertex displacement imposes geometric correction to \mathbf{T}^* , while the geometric detailed tangent space normal map provides more details during rendering, such as deep wrinkles. As illustrated in Fig. 3, we use the generic LDM, i.e., Stable Diffusion, to add facial details to the coarse geometry with prompt guidance. Specifically, the vertex displacement and tangent space normal map takes effect during the face rendering process, and hence gradient of the generic LDM with rendered faces as input can be passed using the Score Distillation Sampling technique for learning \mathcal{V}_d and \mathcal{N}_d .

With the selected coarse geometry, our detailed geometry is formulated as:

$$\mathbf{T}^\dagger = \mathbf{T}^* + \mathcal{V}_d \odot \mathbf{n}(\mathbf{T}^*), \quad (4)$$

where $\mathbf{n}(\cdot) \in \mathbb{R}^{3 \times N}$ represents the vertex normal, \odot represents element-wise multiplication. Then with the corrected mesh with vertex displacement added, we can render its images given a camera pose c and lighting direction l :

$$\mathbf{I} = \mathcal{R}_m(\mathbf{T}^\dagger, \mathcal{V}_d, \mathcal{N}_d, c, l), \quad (5)$$

where \mathcal{R}_m is the differentiable render using tangent space normal maps [Laine et al. 2020; Munkberg et al. 2022]. We can write the SDS loss of generic LDM on rendered image \mathbf{I} as follows:

$$\nabla_{x_d} \mathcal{L}_{\text{SDS}}(\mathbf{I}) \triangleq \mathbb{E}_{t,\epsilon} \left[w(t) (\epsilon_\phi(z_t^d; t, \mathcal{P}) - \epsilon) \frac{\partial z_t^d}{\partial \mathbf{I}} \frac{\partial \mathbf{I}}{\partial x_d} \right] \quad (6)$$

where $x_d = [\mathcal{V}_d, \mathcal{N}_d, \boldsymbol{\beta}^*]$ are the optimizing parameters, $z^d = \mathcal{E}(\mathbf{I})$ is the encoded image using LDM image encoder, $w(t)$ is a weighting depending on discrete time step t , ϵ_ϕ is the denoiser of generic LDM with classifier-free guidance. Notice we also refine the shape parameters $\boldsymbol{\beta}^*$, denoted as $\boldsymbol{\beta}^\dagger$.

During the training process, for each rendering we randomly sample a camera pose c and lighting direction l . The sampling space of the camera pose is constrained to an arc with endpoints from the left 45° to the right 45° w.r.t. the front-facing view. The lighting directions l are purely random within the hemisphere in the front part of the face. Besides the SDS loss, we further add regularization terms to ensure the rationality of generated details. The additional regularization losses include:

$$\begin{aligned} \mathcal{L}_{\text{sha}} &= \|\boldsymbol{\beta}^\dagger - \boldsymbol{\beta}^*\|_2^2, & \mathcal{L}_{\text{geo}} &= \text{Laplacian}(\mathbf{T}^\dagger, \mathbf{T}^*), \\ \mathcal{L}_{\text{map}} &= \|\Delta \mathcal{N}_d\|_2^2 + \|\nabla \mathcal{N}_d\|_2^2, \end{aligned} \quad (7)$$

where $\text{Laplacian}(\cdot)$ represents the Laplacian smooth loss between two meshes, and \mathcal{L}_{map} regularizes both the gradient and divergence of the detailed normal map for smoothing. The final optimization

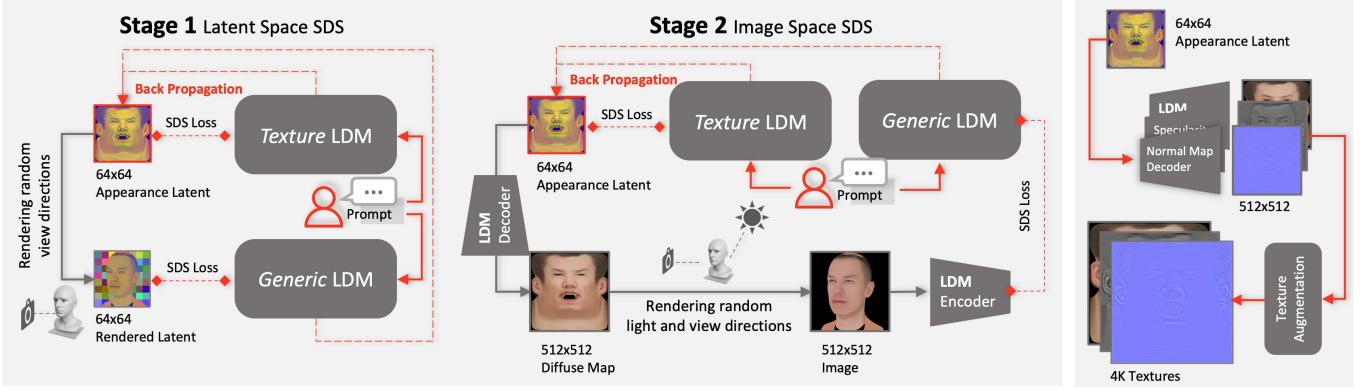


Fig. 4. The overview of physically-based texture diffusion. To generate detailed and realistic textures that match the input prompt, DreamFace performs Dual-path SDS on textures with the use of both a generic LDM and a texture LDM, in both the latent space and image space. By jointly optimizing using two LDMs, we are able to generate high-quality diffuse texture maps that match the input prompt and are consistent with UV unwrapping. An additional texture translation and augmentation module are also included to generate all physically-based textures with high resolution, suitable for rendering.

objective is defined as follows:

$$\mathcal{L}_{carve} = \mathcal{L}_{SDS} + \mathcal{L}_{sha} + \mathcal{L}_{geo} + \mathcal{L}_{map}, \quad (8)$$

where the corresponding weights of each term are ignored for clear presentation. Following DreamFusion [Poole et al. 2022], we uniformly random the discrete time step $t \sim \text{Uniform}(0.02t_{\max}, 0.98t_{\max})$ when optimizing. The detail carving process results in the generation of a detailed geometry that closely matches the input prompt and effectively conveys the described facial characteristics. This detailed geometry is of high quality and is suitable for use in the subsequent appearance generation process.

4.3 Hair selection

In addition to generating the facial geometry, our pipeline also includes the generation of a realistic hairstyle that matches the input prompt. Similar to the geometry generation process, we use CLIP to select the hairstyle candidates with the highest matching score with the input prompt. Our hair dataset comprises 16 hairstyles created by professional artists. We first select the hairstyle that best matches the prompt and then select the predefined hair color by rendering the hair on the head of the geometry. This results in a detailed facial asset with a corresponding hairstyle that closely matches the input prompt.

5 PHYSICALLY-BASED TEXTURE DIFFUSION

Appearance is a critical aspect of animatable neural facial assets, as it allows people to recognize and distinguish faces in a single glance. This section aims to generate an appearance, controlled by the diffuse, specularity and normal maps in texture space in addition to the detailed geometry, that closely matches the input prompt. Beyond effectively conveying the described characteristics of the input prompt, the consistency of texture maps with UV unwrapping is necessary to ensure compatibility with the existing computer graphics production pipeline. Therefore, we propose a dual-path mechanism to jointly optimize texture using two kinds of diffusion models, one generic diffusion model for diverse generation ability

from general prompt inputs, and a novel texture one to ensure the textural specifications in the UV space, to predict the neutral facial assets that are consistent with both the predicted geometry and text prompt. Moreover, for efficient appearance generation, inspired by Latent-NeRF [Metzler et al. 2022], we further introduce a two-stage optimization to perform Score Distillation Sampling (SDS) in both the latent and image spaces, where the latent one significantly provides compact priors for fine-grained synthesis. In the following sections, we will first introduce how to pretrain a diffusion model in the texture space (Sec. 5.1) and then present our efficient two-stage optimization scheme with dual-path enhancement (Sec. 5.2). Finally, we generate physically-based textures with high resolution in Sec. 5.3.

5.1 Learning diffusion model in texture space

The synthesis of image from user-defined text prompt has significantly progressed via text-conditional diffusion models (DMs) [Ho et al. 2022; Nichol et al. 2021; Saharia et al. 2022], where DMs break down the generation problem into a sequential denoising process. Using DMs for text-conditional generation is much easier compared with other deep generative models, as DMs can be easily scaled up and trained stably on billions of image-text pairs to model complex distributions of them. Therefore, we build up our texture generator on the DM pretrained on large-scale image-text pairs. Moreover, denoising at the texture resolution results in unbearable computational costs. We follow the Latent Diffusion Model (LDM) technique [Rombach et al. 2022] to conduct the diffusion process in latent space, where an autoencoder with \mathcal{E} to encode the textures into latent codes, and \mathcal{D} to restore into textures from latent codes.

However, one remaining issue is that the LDM pretrained on natural images is unaware of texture specifications that define how the parts in texture correspond to semantics on face geometry. To enforce the texture specifications and meanwhile preserve the generating ability, we first collect a diverse UV texture dataset and then train our texture LDM by finetuning the pretrained LDM on the dataset to supervise the conformity of texture specifications.

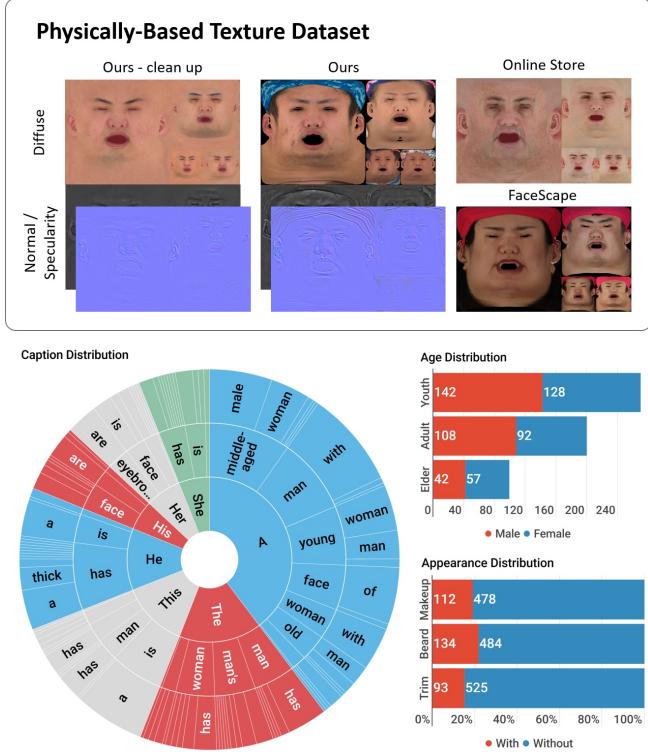


Fig. 5. Our collected physically-based texture dataset. The upper figure demonstrates several samples from several sources, which are made under different standards. The bottom left figure demonstrates the distribution of annotated text prompts. The bottom right bars illustrate the distribution of textures.

5.1.1 Data Collection. We collect the UV texture dataset from multiple sources, including our captured facial scans using the multi-view photometric capture system, public datasets [Yang et al. 2020], and textures from commercial datasets [3DScanStore 2023]. By combining data from diverse datasets, the proposed dataset comprises a diverse range of skin tones, ages, and genders. However, the data merged from the other various datasets were acquired using different standards and exhibited different forms, including variations in UV unwrapping, artistic modifications, and lighting conditions. To ensure consistency and appropriateness for our research purposes, our team of expert artists and researchers performed extensive unification and annotating. A more detailed analysis of the dataset is illustrated in Fig. 5. Our unified dataset can be formulated as $\mathcal{U} = \{U_d, U_s, U_n\}$, where U_d, U_s, U_n represent the diffuse map, specularity map and normal map, respectively. Next, to collect text prompts \mathcal{T} of textures, we first render the textures with geometries using our mesh and texture render and then ask annotators to generate the corresponding text descriptions following specific rules. For the textures without geometries, we render them with a template human face mesh.

Non-face region masking. In addition to texture-text pairs $(\mathcal{U}, \mathcal{T})$, we additionally collect non-face region masks B of textures, which

include hairs, caps, markers, holes, missing edges, etc. These non-face regions in textures distract the texture LDM from learning facial appearances and make the training less stable. Hence, we combine the face color detection model [Chaves-González et al. 2010] with a texture component mask to accurately extract the mask of the non-face region. We condition the texture LDM on masks B , in which the non-face regions are marked as zero and the rest as one, to let the LDM be aware of the non-face regions. We will describe this process in detail in the following section.

5.1.2 Prompt Tuning. A sizable portion of diffuse maps in our dataset containing unwanted lighting effects and deviating from strict diffuse maps is one of the biggest challenges in producing diffuse maps for use in practical applications. Many of the data in our collection were not captured using a system like Light Stage [Ma et al. 2007], which employs polarized patterns to get rid of specularity and produce pure diffuse lighting. As a result, our dataset of diffuse maps contains a proportion of textures that exhibit undesired lighting effects and irregularities that are not suitable for our texture generation pipeline. However, we aim for the texture LDM to learn from a diverse range of textures while ensuring that generated textures during inference do not contain these unwanted artifacts. To address this issue, we first identify two dual domains: those originating from the desired domain, where lighting and specularity have been removed, and those from the undesired domains, where lighting has not been removed, denoted as Ω_d and Ω_u , respectively. Then we provide a novel Prompt Tuning method to overcome this problem and guarantee that our trained texture LDM can create diffuse maps inside the desired domain Ω_d only but can be trained on all the texture-text pairs in both desired and undesired domains $\{\Omega_d, \Omega_u\}$.

As shown in Fig. 6, one straightforward idea is to design two extra domain-specific prompts to indicate the textures in one particular domain explicitly, such as “ \mathcal{T} in the desired domain” and “ \mathcal{T} in the undesired domain”. However, the characteristics of the domain may not be accurately captured through the text encoder, i.e., the CLIP model here, from the handcraft prompts defined based on human understanding, which will be evaluated in Sec. 7.3.1. In addition, identifying appropriate handcraft prompts via prompt engineering is time-consuming and unstable, which is not practical for our needs. Therefore, instead of using the handcraft text prompts, we utilize prompt tuning to learn continuous text prompts, i.e., the word embedding vectors C_d and C_u , to represent the two domains, respectively. Then, the domain-specific text prompts $C = \{C_d, C_u\}$ are combined with the collected ones \mathcal{T} to represent and form the domain-specific texture-prompt pairs $(U_d, Q(\mathcal{T}, C))$, where $Q(\cdot, \cdot)$ is the function that tokenizes the prompts \mathcal{T} and concatenates them with the corresponding domain-specific text prompts. In particular, the domain-specific prompt C_d is concatenated to the prompt belonging to the Ω_d , and vice versa.

To learn a LDM to generate diffuse maps U_d , the modified learning objective conditioned on prompts $Q(\mathcal{T}, C)$ and the non-face mask B is formulated as follows:

$$\mathcal{L}_{\text{LDM}}(\epsilon_\theta, C) = \mathbb{E}_{(U_d, \mathcal{T}), t} \left[\|\epsilon_\theta(\mathcal{E}(U_d)_t, t, \mathcal{E}_{\text{text}}(Q(\mathcal{T}, C)), B) - \epsilon\|_2^2 \right], \quad (9)$$

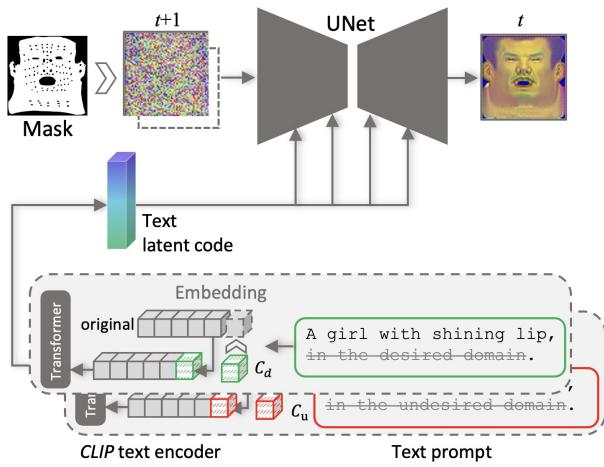


Fig. 6. The overview of our Texture LDM training pipeline. Our approach utilizes two methods to generate high-quality diffuse maps: (1) Prompt Tuning, instead of handcraft domain-specific text prompts, two domain-specific continuous text prompts C_d and C_u are combined with corresponding text prompt, which will be optimized during U-Net denoiser training to avoid unstable and time-consuming prompt engineering for handcraft prompt generation. (2) Non-face region masking, the denoising process of LDM will be additionally conditioned on a non-face region mask to ensure that the generated diffuse map is free of any undesired elements.

where the domain-specific text prompts C and the LDM U-Net denoiser ϵ_θ are the optimizing parameters, $\mathcal{E}(\cdot)$ is the encoder of the LDM, and $\mathcal{E}_{\text{text}}(\cdot)$ is the CLIP text encoder to encode the prompts $Q(\mathcal{T}, C)$. To condition the LDM on the non-face mask B , we scale B to the same size as the latent code and concatenate it to the noisy latent code at each time step for denoising.

During inference, the trained texture LDM can generate diffuse maps in the Ω_d by concatenating the learned C_d with user-defined text prompts. Similarly, it can create high-quality diffuse maps free of any undesired elements by offering a mask filled with ones.

5.2 Two-stage Dual-path Appearance Optimization

Following previous works [Avrahami et al. 2022; Metzer et al. 2022; Turner 2022] that explored the spatial consistencies of super-pixels in latent space, we conduct SDS in two-stage to optimize the texture in latent space first and then refine them in image space. This two-stage optimization scheme allows computation efficiency while retaining the high quality of texture generation.

Dual-path Optimization. Unlike previous 3D generation works [Lin et al. 2022; Metzer et al. 2022; Poole et al. 2022; Vahdat et al. 2021] that only apply one diffusion model to conduct SDS, we propose to combine a generic LDM and a texture LDM to generate textures in a dual-path optimization scheme, where the generic LDM retains the diverse generation ability from general input prompts, while a texture LDM aims to ensure the textural specifications in the UV space. Specifically, our dual-path optimization performs SDS using both the generic LDM, i.e., the Stable Diffusion [Rombach et al. 2022], and

our pretrained texture LDM (see Sec. 5.1) simultaneously. Although the generic LDM is capable of generating plausible textures in image space according to general text prompts, its produced textures are inconsistent with geometry, which leads the facial elements to deviate from their corresponding positions in geometry. Our dual-path optimization scheme solves this problem by additionally employing the pretrained texture LDM on the basis of the generic one to ensure the generated texture follows the UV specifications. We apply the novel dual-path optimization in both optimization stages in the latent space and the image space, respectively, enabling the generalization of textures and consistency in UV space. In the first stage, we utilize dual-path optimization to generate textures in the latent space that provides compact priors for fine-grained synthesis. Similar to the first stage, in the second stage, we also apply the dual-path optimization to enforce UV map specifications of texture and retain generalization ability, but with a detailed normal map and random lighting applied to enhance appearance details in the image space while disentangling lighting from diffuse map.

In the following sections, we introduce the design detail of our dual-path optimization scheme in latent space SDS (Sec. 5.2.1) and image space SDS (Sec. 5.2.2), respectively.

5.2.1 Latent space SDS. Inspired by Latent-NeRF [Metzer et al. 2022] exploring the spatial consistencies of super-pixels and performing the denoising process efficiently, we perform both texture rendering and SDS on latent space directly via our dual-path optimization scheme. Specifically, given the generic LDM denoiser ϵ_ϕ and the texture LDM denoiser ϵ_θ , we perform the SDS simultaneously but operate in different modalities, including rendered latent code for the generic LDM and texture latent code for the texture LDM. We denote the texture latent code to be learned as $z \in \mathbb{R}^{64 \times 64 \times 4}$, and the rendered latent code of z with geometry and camera as $z^r \in \mathbb{R}^{64 \times 64 \times 4}$. The rendered latent code z^r is generated using differentiable mesh renderer \mathcal{R}_a with texture using random background augmentation similar to previous methods [Hong et al. 2022b; Khalid et al. 2022], which could be formulated as follows:

$$z^r = \mathcal{R}_a(T^\dagger, z, c), \quad (10)$$

where T^\dagger is the detailed geometry generated from Sec. 4.2, and c conforms to previous defined camera distribution also in Sec. 4.2. Then, the SDS loss of both LDMs could be formulated as follows:

$$\begin{aligned} \nabla_z \mathcal{L}_{\text{SDS}}^g &= \mathbb{E}_{t,\epsilon} \left[w_r(t)(\epsilon_\phi(z_t^r; t, \mathcal{P}) - \epsilon) \frac{\partial z^r}{\partial z} \right], \\ \nabla_z \mathcal{L}_{\text{SDS}}^\tau &= \mathbb{E}_{t,\epsilon} \left[w_\tau(t)(\epsilon_\theta(z_t; t, \mathcal{P}') - \epsilon) \right], \end{aligned} \quad (11)$$

where t is the uniform discrete time step shared by both LDMs, \mathcal{P} is the input user prompt, and \mathcal{P}' is the augmented prompt for texture LDM with specifically designed keyword appended, as discussed in Sec. 5.1.2. We compute $\partial z^r / \partial z$ from the differential renderer \mathcal{R}_a in Eqn. 10. The final optimization objective then is the combination of the two SDS losses as follows:

$$\mathcal{L}_{\text{tex}} = \lambda_g \mathcal{L}_{\text{SDS}}^g + \lambda_\tau \mathcal{L}_{\text{SDS}}^\tau, \quad (12)$$

where λ_g and λ_τ are corresponding weights that require deliberate design which we will further discuss in Sec. 7.3.3 and supplementary video. We observe that following DreamFusion [Poole et al. 2022]

to uniformly sample the time step and pass the gradient leads to fragility in optimization. In contrast, we vary the t as in the denoising process, i.e., we decrease t with equal intervals from t_{\max} to 0 as in a real denoising process. Finally, by performing SDS under our dual-path optimization scheme, we obtain a learned texture latent code z that matches the user-defined prompt and conforms to the texture specification.

5.2.2 Image space SDS. Our proposed dual-path optimization scheme based on a pair of generic and texture LDMs allows efficient learning of well-formed texture latent code in the first stage. In the second stage, we apply the dual-path optimization scheme in the image space to further obtain a more detailed diffuse map. Specifically, we perform SDS of the generic LDM path on rendered RGB images and utilize the LDM autoencoder \mathcal{E} and \mathcal{D} to map between the latent space and the image space, while keeping the SDS process for the texture LDM path the same as in the first stage. The optimization scheme in the image space SDS follows the basic framework of that in the latent space SDS process, but we apply detailed normal maps and random lighting to enhance appearance details while disentangling lighting from diffuse maps during general LDM SDS. Specifically, instead of rendering the latent code z directly into z' , we convert the 64×64 texture latent code z into a 512×512 diffuse texture map using the LDM decoder \mathcal{D} . Then, we render the geometry T^\dagger using the decoded texture into an 512×512 image under random camera pose c and lighting l with detailed normal map applied, and the rendering process is defined as follows:

$$\mathbf{I} = \mathcal{R}_a(T^\dagger, \mathcal{N}_d, \mathcal{D}(z), c, l). \quad (13)$$

We then use the LDM encoder \mathcal{E} to encode the rendered image \mathbf{I} into latent code $z^{Hr} = \mathcal{E}(\mathbf{I})$ and perform the SDS process of the generic LDM. In particular, SDS loss of the generic LDM is

$$\nabla_z \mathcal{L}_{\text{SDS}}^g = \mathbb{E}_{t,\epsilon} \left[w_r(t) (\epsilon_\phi(z_t^{Hr}; t, \mathcal{P}) - \epsilon) \frac{\partial z^{Hr}}{\partial \mathbf{I}} \frac{\partial \mathbf{I}}{\partial \mathcal{D}(z)} \frac{\partial \mathcal{D}(z)}{\partial z} \right] \quad (14)$$

according to the Leibniz rule. And the SDS process of the texture SDS path is the same with the Eq. 11. After the dual-path optimization in image space, we are able to obtain a finetuned texture latent code z^H that not only has a higher level of details but also conforms to our texture space prior. Finally, the z^H can be decoded into the high-quality diffuse map via the decoder \mathcal{D} .

5.3 Physically-based textures generation

Beyond high-quality texture optimized in Sec. 5.2, an indispensable component of high-fidelity facial assets is the physically-based textures including high-resolution diffuse, specularity and normal maps, which enable photo-realistic rendering using existing CG production pipeline. There is a strong correlation between components of physically-based textures as proposed by Li et al. [2020a]. Hence we can infer the specularity and normal maps from the diffuse texture with an image-to-image translation technique, and we re-use the latent space of the LDM autoencoder for its compactness and efficiency.

Texture Translation. Recall that in the dual path optimization process, we obtain a fine-tuned texture latent code z^H , which could

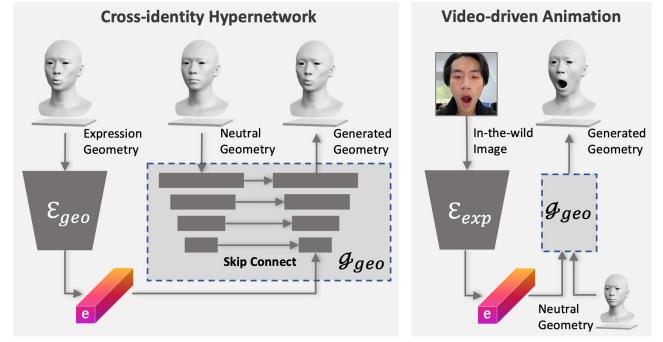


Fig. 7. Overview of animatability empowerment. First, we train a geometry generator to learn a latent space of expression where the decoder is extended to condition on neutral geometry. Then, an expression encoder is further trained to extract expression features from RGB images. Thus, we are able to generate personalized animations by conditioning on the given neutral geometry using monocular RGB images.

be decoded to diffuse map through \mathcal{D} . To further obtain the specularity and normal maps, we develop an image-to-image translation technique inspired by the previous works [Li et al. 2020b]. Instead of training an image-to-image translation module from scratch, we manage to utilize the compactness of the existing LDM autoencoder latent space. Besides the pretrained LDM encoder \mathcal{E} and decoder \mathcal{D} , we train another two specific decoders $\mathcal{D}_s, \mathcal{D}_n$ that decode the existing texture latent code z^H into specularity map and normal map, respectively. Both decoders are trained using our physically-based texture dataset, consisting of 370 tuples of diffuse maps U_d , specularity maps U_s , and normal maps U_n , collected by the photometric multi-stereo capture system. We then encode diffuse maps into latent codes using \mathcal{E} as $u_d = \mathcal{E}(U_d)$. And the corresponding learning objective is formulated as follows:

$$\mathcal{L}_{\text{tex}} = \ell(\mathcal{D}_s(u_d), U_s) + \ell(\mathcal{D}_n(u_d), U_n), \quad (15)$$

where $\ell(\cdot)$ represents the same loss terms described in Stable Diffusion [Rombach et al. 2022].

Texture Augmentation. After generating specularity and normal maps from corresponding diffuse map input, we further enhance their quality and upscale them to 4K resolution to add pore-level details while preserving the identity information. Specifically, we first fine-tune the face restoration network RestoreFormer [Wang et al. 2022] on our dataset at 512×512 resolution to enhance the facial details. Then we refine the super-resolution model Real-ESRGAN [Wang et al. 2021] using our high-resolution textures to further generate 4096×4096 physically-based textures, which are essential for photo-realistic rendering. The final normal texture synthesizes pore-level details on the face and is also represented in tangent space. We integrate it with the geometric detailed tangent space normal map \mathcal{N}_d during the rendering process.

6 ANIMATABILITY EMPOWERMENT

In addition to generating fine face geometry and physically-based textures, our framework also empowers the animatability of the generated facial asset. While our asset directly supports traditional

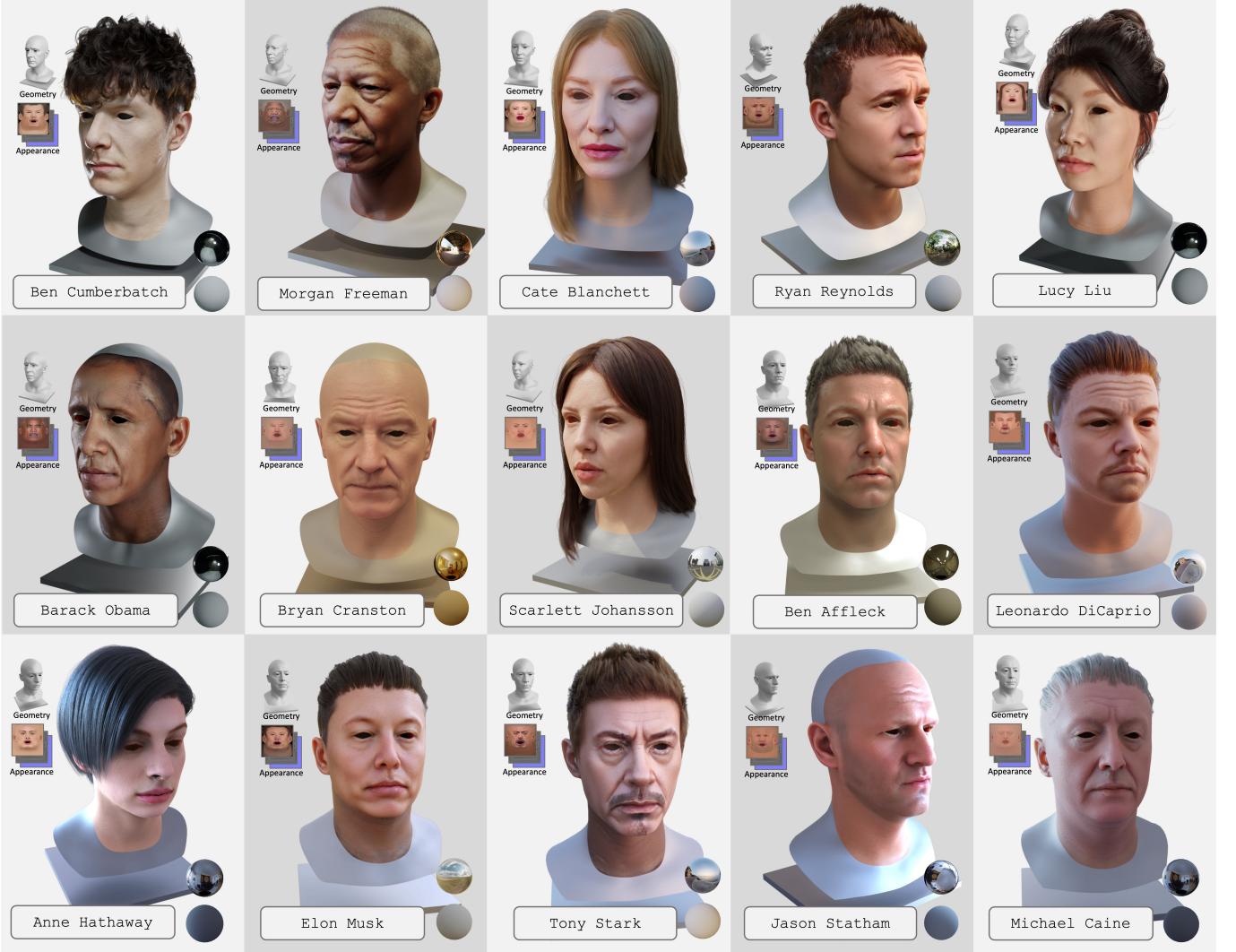


Fig. 8. Generated facial assets of celebrities. Our approach generates facial assets of celebrities that capture their personalized characteristics and achieve a high degree of resemblance. By generating physically-based textures, our facial assets achieve photo-realistic results using the modern CG rendering pipeline.

blendshape-based animation due to its consistent geometric topology, one could utilize existing facial performance capture techniques [Apple 2023; Li et al. 2017; Moser et al. 2021] to obtain the corresponding expression blendshape parameters from images and videos to animate our generated facial asset.

Beyond generic animatability using blendshapes, we explore introducing an enhanced animation scheme to produce personalized expressions while preserving the unique characteristics of our generated facial asset and creating more realistic and believable animations.

In contrast with [Laine et al. 2017; Lombardi et al. 2018; Zhang et al. 2022] which learn person-specific facial animations from extensive facial captures for one identity, we develop a cross-identity geometry hypernetwork that generates person-specific expressions and serves as a universal prior for expression space. Inspired by

Cao et al. [2022], we condition the geometry generator on neutral geometry and train the hypernetwork using numerous geometries with various expressions.

Cross-identity hypernetwork. Our geometry hypernetwork consists of a geometry expression encoder \mathcal{E}_{geo} and a geometry generator \mathcal{G}_{geo} , where the geometry generator is basically a U-Net conditioned on neutral geometry without expressions, as exhibited in Fig. 7. The geometry expression encoder encodes different facial expressions into a unified expression latent code, while the geometry generator uses this code, along with the neutral geometry of a specific identity, to generate the desired facial geometry with corresponding expressions. The forward process could be formulated as



Fig. 9. Generated facial assets from descriptions. Our approach generates facial assets that faithfully match the characteristics described in the prompts. Through our animatability empowerment, the generated facial assets can be animated using a single RGB image and rendered photo-realistically in modern CG pipelines.

follows:

$$\begin{aligned} z_e &= \mathcal{E}_{\text{geo}}(\mathbf{G}), \quad \tilde{\mathbf{G}} = \mathcal{G}_{\text{geo}}(\mathbf{G}_0, z_e), \\ \mathcal{L}_{\text{recon}} &= \|\tilde{\mathbf{G}} - \mathbf{G}\|_2^2, \end{aligned} \quad (16)$$

where \mathbf{G} is the input geometry, z_e is the expression code, \mathbf{G}_0 is the neutral geometry of \mathbf{G} , $\tilde{\mathbf{G}}$ is the generated geometry, and $\mathcal{L}_{\text{recon}}$ is the difference between the generated geometry and original geometry corresponding to z_e as training loss. We train this geometry hypernetwork in a self-supervised manner, using a dataset of geometries with various expressions and identities. Once trained, the geometry generator \mathcal{G}_{geo} is capable of producing geometries with the desired expressions for a specific identity using the unified expression latent code.

Video-driven animation. With the geometry generator \mathcal{G}_{geo} well trained, we additionally train an image expression encoder \mathcal{E}_{exp} to extract unified expression latent code from RGB images. We freeze the geometry generator in this process, and train \mathcal{E}_{exp} under the

supervision of both real images and randomly rendered images from geometries \mathbf{G} and textures \mathbf{A} in the dataset, which could be formulated as:

$$\begin{aligned} \hat{z}_e &= \mathcal{E}_{\text{exp}}(\mathcal{R}_a(\mathbf{G}, \mathbf{A})), \quad \hat{\mathbf{G}} = \mathcal{G}_{\text{geo}}(\mathbf{G}_0, \hat{z}_e), \\ \mathcal{L}_{\text{exp}} &= \|\hat{\mathbf{G}} - \mathbf{G}\|_2^2, \end{aligned} \quad (17)$$

where \mathcal{R}_a is a mesh and texture renderer that renders under random augmentations, and \mathcal{L}_{exp} is the difference between ground-truth geometry with expression and the generated geometries. After training, \mathcal{E}_{exp} is able to extract expression latent code from in-the-wild videos, and \mathcal{G}_{geo} could further produce personalized animations for our generated facial asset.

Dataset. In order to train our geometry generator with both generalizability to various identities and accurate capturing of fine-grained expression details, we capture a dataset that includes a large number of static scans of different expressions and identities, as well as several dynamic performance sequences from various performers.



Fig. 10. Generation out of distribution. The upper row shows the rendering results from the differentiable renderer, and the lower row shows the corresponding diffuse maps. Our framework faithfully reveals the facial characteristics of characters, even if they are not present in our texture dataset, for example, the pink nose of *Na'vi* and the metallic patterns of *Black Panther*'s mask. In addition, our texture LDM serves as a robust prior, ensuring that the generated facial components share a consistent UV space.

Our dataset consists of 38400 registered mesh and 614400 images from 300 identities across different genders, ages, and ethnicities.

Network training. The geometry encoder \mathcal{E}_{geo} consists of 7 sequential convolution layers with kernel size of 5 and stride of 2 that has [32,64,128,256,512,512,512] intermediate channels and a linear layer that maps to a latent space with 256 dimensions. The geometry generator \mathcal{G}_{geo} is implemented as a U-Net with similar architecture as proposed by Cao et al. [2022], which down-samples input 5 times with [32,64,128,256,256] intermediate channels. The image encoder \mathcal{E}_{exp} shares the same architecture with \mathcal{E}_{geo} . The input and output geometry are formulated as geometry map of size 256×256, as described in previous works [Li et al. 2020b,a; Zhang et al. 2022]. The input image is cropped to the facial region and down-sampled to 256×256. The first stage of training the geometry hypernetwork \mathcal{E}_{geo} , \mathcal{D}_{geo} takes 5 days to converge. And the second stage training the image encoder \mathcal{E}_{exp} takes 48 hours to converge.

All training is done using Pytorch with AdaBelief optimizer and learning rate 5e-5, using a single Nvidia A6000 GPU.

7 EXPERIMENTS

In this section, we present the experiment results of DreamFace for generating animatable neural facial assets. We first introduce the implementation details and showcase a gallery of generated high-quality assets by DreamFace, highlighting the wide scope of applications that our approach enables. We then provide a detailed evaluation of modules in our pipeline, including geometry, appearance and animation modules, both qualitatively and quantitatively. At last, we compare with the state-of-the-art methods, followed by a comprehensive user study.



Fig. 11. Texture editing using prompts and sketches. By directly using our trained texture LDM with a prompt, one can achieve global editing effects such as aging and makeup. By further combining masks or sketches, one can create various effects such as tattoos, beards, and birthmarks.

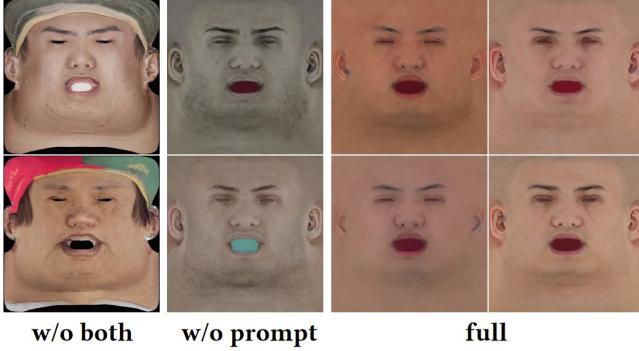


Fig. 12. Video-driven animation of generated facial assets. For each case, we show the input RGB images (left) and the personalized driven results (right). Our framework provides each generated facial asset with personalized expressions from a single image.

7.1 Implementation details

We rely on the generic LDM and texture LDM for geometry generation and texture diffusion. In our implementation, we use the pre-trained LDM checkpoint *stable-diffusion-v1-5* from RunwayML [RunwayML 2022] as our generic LDM, and use the CLIP checkpoint *clip-vit-large-patch14* from OpenAI [OpenAI 2022]. While for the texture LDM, we duplicate a generic LDM and finetune it as described in 5.1. We train the texture LDM on two Nvidia A6000 GPUs for 150 epochs which takes about 12 hours.

During the geometry generation process, we sample one million candidates from ICT-FaceKit according to our pre-defined distribution and then perform 300 steps of detail carving. In the texture diffusion stage, we perform 200 steps of dual-path optimization in latent space and 200 steps in image space. By our two-stage design, the SDS-based generation process is very efficient, enabling the generation of a high-quality facial asset within 5 minutes on a single Nvidia A6000 GPU.



Training setting	KID↓	
	Full dataset	Desired domain
w/o m&p	0.1282	0.2737
w/o prompt	0.2467	0.1068
full	0.2125	0.0578

Fig. 13. Qualitative and quantitative comparison with different variants on learning texture LDM. With our two careful designs, the texture LDM produces textures with the best quality, where the lighting and unwanted elements are well removed, as illustrated in the upper figure, and also achieves the lowest KID in desired texture domain as shown in the lower table.

7.2 Generation Results

DreamFace is a novel framework for generating realistic 3D facial assets that are not only true-to-life in appearance, but also rich in animation capabilities. Our approach allows for the creation of highly detailed and personalized characters, from fashion icons to exotic creatures from fiction and film, all through the use of simple textual prompts.

By only providing a simple prompt of celebrity or general description, one can create realistic 3D facial assets that highly resemble the described characteristics, as demonstrated in Fig. 8 and Fig. 9, respectively. With DreamFace, one can even generate faces of fashion icons or unreal humanoid characters from fiction, movies, or even dreams in mind while retaining high recognition, as illustrated in Fig. 10. By leveraging both the generic LDM and the texture LDM, DreamFace generates animatable neural facial assets that are compatible with modern computer graphics pipelines and achieve photo-realistic rendering, utilizing the robust prior knowledge learned by these generative models through natural language prompts.

One can continue to customize its facial appearance for novel effects with additional textual guidance or even hand-made sketches, as illustrated in Fig. 11. Our texture LDM directly supports denoising the existing texture under prompt guidance, such as adding wrinkles or makeup. By mixing the weight of trained texture LDM with a generic LDM, the mixed texture LDM is capable of generating desired patterns from noise or hand-made sketches in the masked areas that seamlessly blend with original diffuse maps (refer to stable-diffusion-webui [AUTOMATIC1111 2022] for detail implementation).

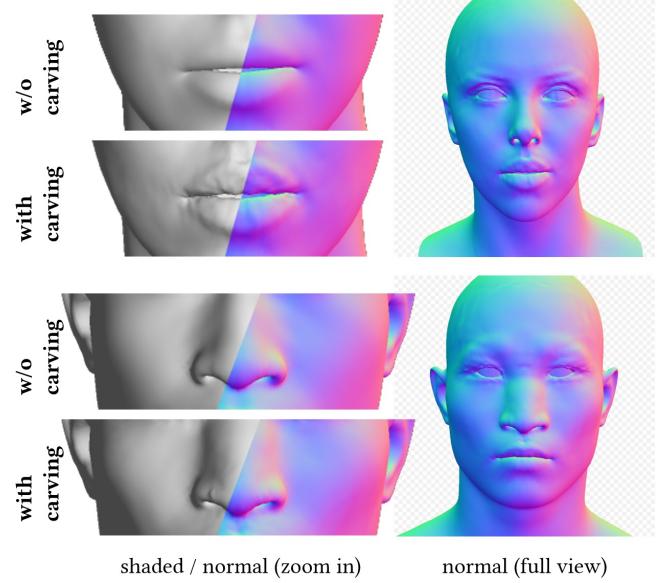


Fig. 14. Qualitative comparison on detail carving. For each case, we compare coarse geometry and detailed geometry with a zoomed-in view (left) and a full view (right). The upper case uses the prompt "the face of Scarlett Johansson" whose mouth is well-carved, demonstrating her characteristic traits. The lower case uses the prompt "the face with a Na'vi style nose in the movie Avatar" where a wide nose bridge resembles the indigenous alien species. This comparison shows the effectiveness of our detail carving in capturing specific traits and features of different characters.



Fig. 15. Qualitative comparison on appearance latent code generation. For each case, we compare two variations with our pipeline together with rendering results using our differentiable renderer. The upper case has the prompt "the face of Anthony Hopkins" and the lower case has the prompt "the face of a young lady with exquisite makeup". Compared to variations without SDS in latent space and using texture LDM, our pipeline produces the best diffuse maps with precise, well-defined components in consistent UV space.

Empowered by DreamFace, one can easily animate his creation using in-the-wild video footage directly with nuanced performance captured. Given a generated facial asset with geometry T^\dagger , we directly inference its geometry from a input RGB image V by $\hat{G} =$

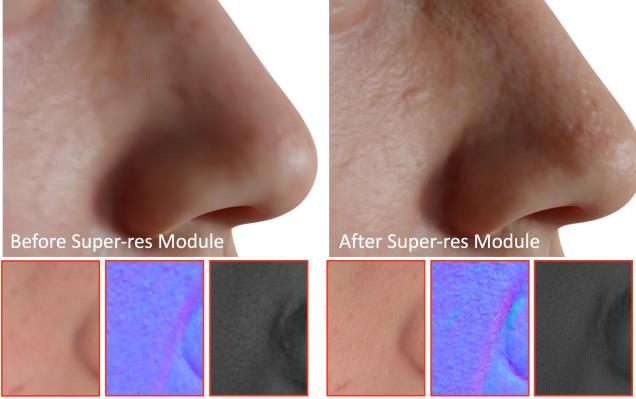


Fig. 16. Qualitative comparison of physically-based texture augmentation. We show the directly decoded textures from generated texture latent code (left) and the high-resolution physically-based textures (right).

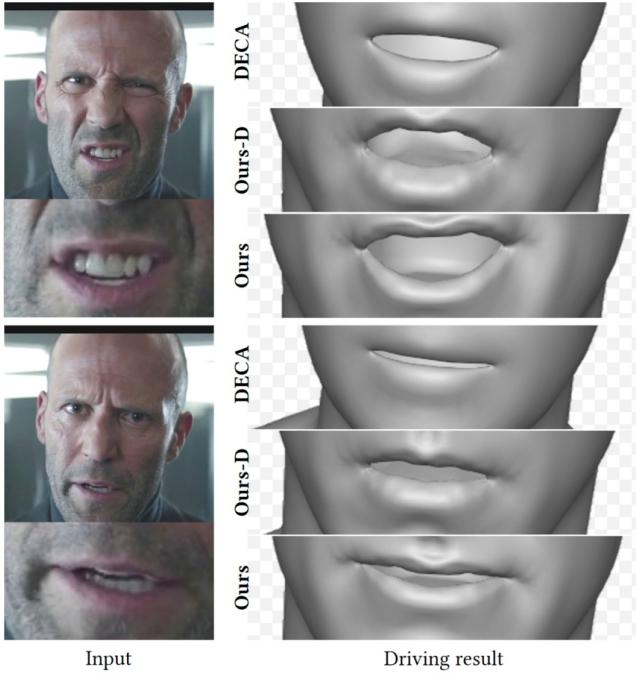


Fig. 17. Qualitative comparison of expression generation with the previous generalized method. Compared to DECA, which uses generic blendshapes for expression control, our framework provides fine-grained expression details while carefully capturing nuanced performance, such as “left sneering” (upper).

$\mathcal{G}_{\text{geo}}(\mathbf{T}^\dagger, \mathcal{E}_{\text{exp}}(\mathbf{V}))$. As illustrated in Fig. 12, our animation pipeline empower animatability by generating personalized expressions and faithfully restoring the detailed performance. With DreamFace, even novices can naturally create characters straight from their imagination, ready to be driven and rendered in stunning detail.



Fig. 18. Qualitative and quantitative comparison between prompt controlled generation methods. We show several generated results of each method (upper) and quantitative results (lower). Compared to previous methods, our method preserves more details on appearance while keeping a fine-grained geometry. In addition, our method achieves the best text-matching score and has the least running time.

7.3 Evaluation

7.3.1 Evaluation of texture LDM. In this section, we evaluate the quality of the generated textures produced by our texture LDM. Specifically, we employ Kernel Inception Distance (KID) to measure the performance of texture generation without classifier-free guidance on our full dataset and textures only from our desired domain. To evaluate the performance of our approach, we conduct experiments with three different settings: **w/o m&p**, **w/o prompt**, and **full**. **w/o m&p** denotes the naive LDM training pipeline without prompt tuning or masking. **w/o prompt** denotes the variant without prompt tuning, and **full** denotes our pipeline with both designs. We generate 1000 samples for each setting and calculate KID 50 times with a batch size equal to the minimum available size of each pair. As illustrated in Fig. 13, **naive** results in lighting remaining and unwanted elements like hat and black area. **w/o prompt** removes the unwanted elements like black area but may produce textures with inconsistent skin tones. Our **full** generates desired

Case	Resemble the specific character?
Celebrity	72.3%
Out of distribution	71.6%

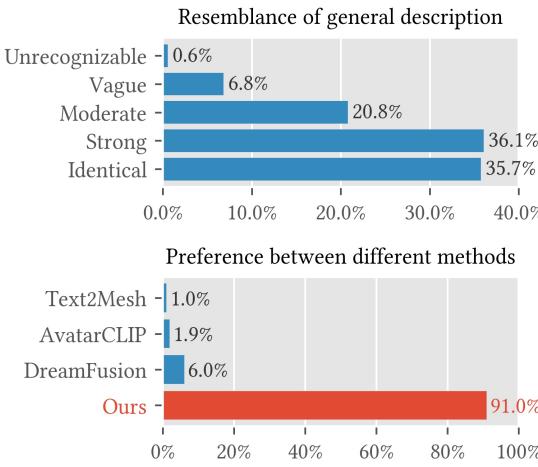


Fig. 19. Quantitative results of user study. The upper table illustrates our high resemblance to specific characters, while the lower figures demonstrate the high user ratings on matching general descriptions and the extremely high user preference for our method compared to other methods.

textures with high quality, which demonstrates the effectiveness of our designs.

7.3.2 Detail carving. Detail carving, as proposed in Sec. 4.2, is a key step of our framework that further reveals characteristics at the geometry level. Let **w/o carving** denote the coarse geometry and let **with carving** denote our detailed geometry. As illustrated in Fig. 14, different carved geometries demonstrate the capability of our approach both in capturing characteristic traits of a specific individual and generating distinct facial features even across species.

7.3.3 Ablation study of texture diffusion. Our proposed texture latent code generation pipeline, as outlined in Sec.5.2.1, utilizes both SDS in the latent space and joint optimization using a texture LDM to generate high-quality diffuse maps. Latent space SDS produces texture latent code more efficiently due to low resolution, while only using image space SDS takes much longer time (about 15 minutes, which is 3-4x times longer compared to latent space SDS) to generate satisfactory results. To evaluate the performance of our pipeline, we compare it to two variations: one without the use of SDS in the latent space or texture LDM (similar to DreamFusion [Poole et al. 2022]), denoted as **w/o l&t**, and another one without the use of texture LDM (similar to Latent-NeRF [Metzger et al. 2022]), denoted as **w/o texture**. Our full pipeline, denoted as **full**, includes both SDS and texture LDM. As shown in Fig.15, the results demonstrate that the **w/o l&t** pipeline produces diffuse maps with a lot of noise and artifacts, **w/o texture** pipeline produces diffuse maps where components are drifting in UV space and lighting exists, while the

full pipeline gives the best results that meet the standard used in physically-based rendering.

7.3.4 Texture augmentation. We demonstrate the effectiveness of texture augmentation in Fig. 16. After generating satisfactory texture latent code, our texture translation and augmentation module further provide detailed physically-based textures, including the diffuse map, specularity map and normal map in high resolution, which are essential for photo-realistic rendering.

7.3.5 Comparison of animation generation. We compare our animation method that generates personalized expressions from images to the previous method DECA [Feng et al. 2021] that infers parameters of general expression blendshapes. To adapt FLAME [Li et al. 2017] blendshapes to our asset, we transfer the expression parameters (blendshape weights and jaw rotation) to control our generated facial asset with the rig that same as FLAME. Let **DECA**, **Ours-D**, **Ours** denote the original results from DECA, the transferred expression using our asset, and our enhanced animation scheme, respectively. As illustrated in Fig. 17, compared to DECA, our framework provides personalized expressions and more accurately restores the detailed expression in the video.

7.3.6 Comparison with state-of-the-art methods. We further present a comprehensive comparison of our proposed method with state-of-the-art prompt-controlled generation techniques, including Text2Mesh [Michel et al. 2022], AvatarCLIP [Hong et al. 2022b] and DreamFusion [Poole et al. 2022]. Our goal is to evaluate the quality and efficiency of our approach in comparison to other methods. To conduct the comparison, we generate 10 different characters using each method. For Text2Mesh, we use the official implementation. For AvatarCLIP, we use the generated results from the official website. For DreamFusion, we use a re-implementation version [Tang 2022] that uses Stable Diffusion. For all these approaches, we use the prompt "the realistic face of SOMEONE" for the generation. To quantitatively evaluate the generated results, we calculate the CLIP score by computing the cosine similarity of image features and text features using the prompt "the realistic face of SOMEONE". Additionally, we also measure the running time of each method. As illustrated in Fig. 18, our method demonstrates its capability to generate high-quality realistic assets efficiently and outperforms others in terms of both CLIP score and running time.

7.3.7 User study. Finally, we conducted a comprehensive user study to evaluate the performance of our generated facial assets in terms of their ability to match the given prompts and their level of resemblance to specific characters or general descriptions. We recruited 180 volunteers to participate in the study, which consisted of three main evaluations: matching specific characters, matching general descriptions, and user preference across different methods. For the evaluation of specific characters, we generated 27 different samples that included both celebrities and out-of-distribution characters and asked volunteers to rate the level of the resemblance of the generated results to the specific characters. For the evaluation of general descriptions, we generated 10 different cases with diversity and asked volunteers to score the conformity of the results with the descriptions on a five-point scale. Additionally, we compared our method with Text2Mesh, AvatarCLIP, and DreamFusion using

10 different prompts, and asked volunteers to select their preferred method. As shown in Fig. 19, our method achieved high levels of resemblance and is significantly more preferred than the other three methods.

7.4 Limitation

As a brave attempt, DreamFace presents a novel approach for the progressive generation, editing, and animation of neural facial assets using simple text prompts. Our framework demonstrates a high level of realism and resemblance in the generated assets, even for novice users. However, it is important to note that there are limitations to our approach. One limitation is that while we have incorporated hairstyle generation into our framework, the generation of full facial components such as eyes and mouth interiors is currently not possible. The modeling and rendering of eyes, in particular, presents significant technical challenges that are yet to be fully addressed. Another limitation is that as a framework based on prompt-conditioned diffusion models, DreamFace is constrained by the capabilities of these models. For example, while GAN-based methods have the ability to invert a given sample easily, the inversion of diffusion models has not been fully explored. While the performance of DreamFace is influenced by the state-of-the-art in diffusion model research, further advancements in this area may improve our results. Lastly, while our framework has demonstrated strong animation capabilities through the native support of blendshapes and an enhanced animation scheme, there is still potential for further research in the area of prompt-based animation control. The generation of facial motion, including lively expressions and nuanced performance, is a particularly interesting and challenging problem that remains to be solved.

Potential ethical implications. As a text-driven generation method, DreamFace operates on the output of pre-trained large-scale vision-language models such as CLIP and Stable Diffusion, and is subject to any biases present in their training data, like gender or racial preferences. However, it is essential to note that these biases are not unique to DreamFace, and are inherent in any generation method that relies on these pre-trained models. Additionally, the ease of use and high quality of the generated assets may also raise concerns about the potential for misuse, such as creating fake videos or impersonating individuals. It is important for future work to consider and address these ethical issues in the development of text-driven generation methods by developing methods to mitigate biases in pre-trained large-scale vision-language models and ensuring that the data used for training these models are diverse, representative, and carefully reviewed.

8 CONCLUSION

We have presented DreamFace, to progressively generate personalized 3D faces that are compatible with CG engines from only text-prompt controls. Through DreamFace, even novices can naturally create 3D human characters with desired shapes and textures they have in mind, easily customize the creations for novel effects like aging and virtual makeup, and even further animate the

creations using in-the-wild video footage. Specifically, our coarse-to-fine scheme efficiently generates high-quality geometry, including the coarse geometry with a unified topology and the nuanced displacement and normal details. For appearance generation, our dual-path mechanism organically combines two latent diffusion models, i.e., a generic LDM and a texture LDM, achieving diverse and consistent results with the textural specification. Our two-stage optimization in both the latent and image spaces achieves highly efficient and fine-grained synthesis, enabling the mapping from the compact latent space to physically-based textures. Our neutral assets naturally support blendshapes-based facial animations, while the accompanied neural animation scheme further provides personalized fine-grained animation from only video input. Extensive experimental results and user studies have demonstrated the effectiveness of DreamFace. We believe that our approach renews the generation of accessible 3D facial assets with physically-based rendering quality and rich animation ability in the neural and prompt-interaction era. It not only benefits the CG production industry but also incentivizes numerous innovative applications for VR/AR and the emerging Metaverse.

REFERENCES

- 3DScanStore. 2023. 3DScanStore. <https://www.3dscanstore.com>.
- Victoria Fernández Abrevaya, Adnane Boukhayma, Stefanie Wuhrer, and Edmond Boyer. 2019. A decoupled 3d facial shape model by adversarial training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 9419–9428.
- Oleg Alexander, Mike Rogers, William Lambeth, Matt Chiang, and Paul Debevec. 2009. The Digital Emily Project: Photoreal Facial Modeling and Animation. In *ACM SIGGRAPH 2009 Courses* (New Orleans, Louisiana) (*SIGGRAPH ’09*). Association for Computing Machinery, New York, NY, USA, Article 12, 15 pages. <https://doi.org/10.1145/1667239.1667251>
- Shivangi Aneja, Justus Thies, Angela Dai, and Matthias Nießner. 2022. ClipFace: Text-guided Editing of Textured 3D Morphable Models. In *ArXiv preprint arXiv:2212.01406*.
- Apple. 2023. ARKit - Face Tracking. <https://developer.apple.com/documentation/arkit/arfaceanchor>.
- AUTOMATIC1111. 2022. stable-diffusion-webui. <https://github.com/AUTOMATIC1111/stable-diffusion-webui>.
- Omri Avrahami, Ohad Fried, and Dani Lischinski. 2022. Blended Latent Diffusion. *arXiv preprint arXiv:2206.02779* (2022).
- Linchao Bao, Xiangkai Lin, Yajing Chen, Haoxian Zhang, Sheng Wang, Xuefei Zhe, Di Kang, Haozhi Huang, Xinwei Jiang, Jue Wang, Dong Yu, and Zhengyou Zhang. 2021. High-Fidelity 3D Digital Human Head Creation from RGB-D Selfies. *ACM Transactions on Graphics* (2021).
- Volker Blanz and Thomas Vetter. 1999. A morphable model for the synthesis of 3D faces. In *Proceedings of the 26th annual conference on Computer graphics and interactive techniques*. 187–194.
- Chen Cao, Tomas Simon, Jin Kyu Kim, Gabe Schwartz, Michael Zollhoefer, Shun-Suke Saito, Stephen Lombardi, Shih-En Wei, Danielle Belko, Shouo-I Yu, Yaser Sheikh, and Jason Saragih. 2022. Authentic Volumetric Avatars from a Phone Scan. *ACM Trans. Graph.* 41, 4, Article 163 (jul 2022), 19 pages. <https://doi.org/10.1145/3528223.3530143>
- Chen Cao, Yanlin Weng, Shun Zhou, Yiyi Tong, and Kun Zhou. 2013. Facewarehouse: A 3d facial expression database for visual computing. *IEEE Transactions on Visualization and Computer Graphics* 20, 3 (2013), 413–425.
- Eric R Chan, Connor Z Lin, Matthew A Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas J Guibas, Jonathan Tremblay, Sameh Khamis, et al. 2022. Efficient geometry-aware 3D generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 16123–16133.
- Jose M. Chaves-González, Miguel A. Vega-Rodríguez, Juan A. Gómez-Pulido, and Juan M. Sánchez-Pérez. 2010. Detecting Skin in Face Recognition Systems: A Colour Spaces Study. 20, 3 (may 2010), 806–823. <https://doi.org/10.1016/j.dsp.2009.10.008>
- Yongwei Chen, Rui Chen, Jiabao Lei, Yabin Zhang, and Kui Jia. 2022. TANGO: Text-driven Photorealistic and Robust 3D Stylization via Lighting Decomposition. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- Katherine Crowson, Stella Biderman, Daniel Kornis, Dashiell Stander, Eric Hallahan, Louis Castricato, and Edward Raff. 2022. Vqgan-clip: Open domain image generation and editing with natural language guidance. In *European Conference on Computer Vision*. Springer, 88–105.

- Paul Debevec, Tim Hawkins, Chris Tchou, Haarm-Pieter Duiker, Westley Sarokin, and Mark Sagar. 2000. Acquiring the Reflectance Field of a Human Face (*SIGGRAPH ’00*). ACM Press/Addison-Wesley Publishing Co., USA, 145–156. <https://doi.org/10.1145/344779.344855>
- Yu Deng, Jiaolong Yang, Jianfeng Xiang, and Xin Tong. 2022. GRAM: Generative Radiance Manifolds for 3D-Aware Image Generation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- Yao Feng, Haiwen Feng, Michael J. Black, and Timo Bolkart. 2021. Learning an Animatable Detailed 3D Face Model from In-The-Wild Images. *ACM Transactions on Graphics (Proc. SIGGRAPH)* 40, 8. <https://doi.org/10.1145/3450626.3459936>
- Kentaro Fukamizu, Masaaki Kondo, and Ryuichi Sakamoto. 2019. Generation High resolution 3D model from natural language by Generative Adversarial Network. *arXiv preprint arXiv:1901.07165* (2019).
- Graham Fyffe, Andrew Jones, Oleg Alexander, Ryosuke Ichikari, and Paul Debevec. 2015. Driving High-Resolution Facial Scans with Video Performance Capture. *ACM Trans. Graph.* 34, 1, Article 8 (dec 2015), 14 pages. <https://doi.org/10.1145/2638549>
- Rinon Gal, Or Patashnik, Haggai Maron, Amit H. Bermano, Gal Chechik, and Daniel Cohen-Or. 2022. StyleGAN-NADA: CLIP-Guided Domain Adaptation of Image Generators. *ACM Trans. Graph.* 41, 4, Article 141 (jul 2022), 13 pages. <https://doi.org/10.1145/3528223.3530164>
- Baris Gecer, Alexandros Lattas, Stylianos Ploumpis, Jiankang Deng, Athanasios Pappaioannou, Stylianos Moschoglou, and Stefanos Zafeiriou. 2020. Synthesizing coupled 3d face modalities by trunk-branch generative adversarial networks. In *European conference on computer vision*. Springer, 415–433.
- Baris Gecer, Stylianos Ploumpis, Irene Kotsia, and Stefanos Zafeiriou. 2019. GANFIT: Generative Adversarial Network Fitting for High Fidelity 3D Face Reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Baris Gecer, Stylianos Ploumpis, Irene Kotsia, and Stefanos P Zafeiriou. 2021. Fast-GANFIT: Generative Adversarial Network for High Fidelity 3D Face Reconstruction. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2021).
- Abhijeet Ghosh, Graham Fyffe, Borom Tunwattanapong, Jay Busch, Xueming Yu, and Paul Debevec. 2011. Multiview face capture using polarized spherical gradient illumination. In *Proceedings of the 2011 SIGGRAPH Asia Conference*. 1–10.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2020. Generative Adversarial Networks. *Commun. ACM* 63, 11 (oct 2020), 139–144. <https://doi.org/10.1145/3422622>
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems* 33 (2020), 6840–6851.
- Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim Salimans. 2022. Cascaded Diffusion Models for High Fidelity Image Generation. *J. Mach. Learn. Res.* 23 (2022), 47–1.
- Fangzhou Hong, Mingyuan Zhang, Liang Pan, Zhongang Cai, Lei Yang, and Ziwei Liu. 2022b. AvatarCLIP: Zero-Shot Text-Driven Generation and Animation of 3D Avatars. *ACM Transactions on Graphics (TOG)* 41, 4 (2022), 1–19.
- Yang Hong, Bo Peng, Haiyao Xiao, Ligang Liu, and Juyong Zhang. 2022a. Headnerf: A real-time nerf-based parametric head model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* 20374–20384.
- Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. 2017. Image-to-Image Translation with Conditional Adversarial Networks. *CVPR* (2017).
- Ajay Jain, Ben Mildenhall, Jonathan T Barron, Pieter Abbeel, and Ben Poole. 2022. Zero-shot text-guided object generation with dream fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 867–876.
- Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. 2020. Analyzing and Improving the Image Quality of StyleGAN. In *Proc. CVPR*. Nasir Mohammad Khalid, Tianhao Xie, Eugene Belilovsky, and Popa Tiberiu. 2022. CLIP-Mesh: Generating textured meshes from text using pretrained image-text models. (December 2022).
- Diederik P. Kingma and Max Welling. 2014. Auto-Encoding Variational Bayes. In *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14–16, 2014, Conference Track Proceedings*.
- Samuli Laine, Janne Hellsten, Tero Karras, Yeongho Seol, Jaakko Lehtinen, and Timo Aila. 2020. Modular Primitives for High-Performance Differentiable Rendering. *ACM Transactions on Graphics* 39, 6 (2020).
- Samuli Laine, Tero Karras, Timo Aila, Antti Herva, Shunsuke Saito, Ronald Yu, Hao Li, and Jaakko Lehtinen. 2017. Production-level facial performance capture using deep convolutional neural networks. In *Proceedings of the ACM SIGGRAPH/Eurographics symposium on computer animation*. 1–10.
- Alexandros Lattas, Stylianos Moschoglou, Baris Gecer, Stylianos Ploumpis, Vasileios Triantafyllou, Abhijeet Ghosh, and Stefanos Zafeiriou. 2020. AvatarMe: Realistically Renderable 3D Facial Reconstruction ‘In-the-Wild’. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Alexandros Lattas, Stylianos Moschoglou, Stylianos Ploumpis, Baris Gecer, Abhijeet Ghosh, and Stefanos Zafeiriou. 2021. AvatarMe++: Facial shape and brdf inference with photorealistic rendering-aware gans. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 44, 12 (2021), 9269–9284.
- Myunggi Lee, Wonwoong Cho, Moonheum Kim, David Inouye, and Nojun Kwak. 2020. Styleuv: Diverse and high-fidelity uv map generative model. *arXiv preprint arXiv:2011.12893* (2020).
- J. Li, Z. Kuang, Y. Zhao, M. He, and H. Li. 2020b. Dynamic Facial Asset and Rig Generation from a Single Scan. *ACM Transactions on Graphics (TOG)* (2020).
- Rui long Li, Karl Bladin, Yajie Zhao, Chinmay Chinara, Owen Ingraham, Pengda Xiang, Xinglei Ren, Pratusha Prasad, Bipin Kishore, Jun Xing, et al. 2020a. Learning formulation of physically-based face attributes. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 3410–3419.
- Tianye Li, Timo Bolkart, Michael J. Black, Hao Li, and Javier Romero. 2017. Learning a model of facial shape and expression from 4D scans. *ACM Transactions on Graphics (Proc. SIGGRAPH Asia)* 36, 6 (2017), 194:1–194:17. <https://doi.org/10.1145/3130800.3130813>
- Yiyi Liao, Katja Schwarz, Lars Mescheder, and Andreas Geiger. 2020. Towards unsupervised learning of generative models for 3d controllable image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 5871–5880.
- Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. 2022. Magic3D: High-Resolution Text-to-3D Content Creation. *arXiv preprint arXiv:2211.10440* (2022).
- Yuchen Liu, Zhixin Shu, Yijun Li, Zhe Lin, Richard Zhang, and SY Kung. 2022. 3d-fm gan: Towards 3d-controllable face manipulation. In *European Conference on Computer Vision*. Springer, 107–125.
- Stephen Lombardi, Jason Saragih, Tomas Simon, and Yaser Sheikh. 2018. Deep appearance models for face rendering. *ACM Transactions on Graphics (TOG)* 37, 4 (2018), 1–13.
- Wan-Chun Ma, Tim Hawkins, Pieter Peers, Charles-Felix Chabert, Malte Weiss, Paul E Debevec, et al. 2007. Rapid Acquisition of Specular and Diffuse Normal Maps from Polarized Spherical Gradient Illumination. *Rendering Techniques* 2007, 9 (2007), 10.
- Elman Mansimov, Emilio Parisotto, Jimmy Ba, and Ruslan Salakhutdinov. 2016. Generating Images from Captions with Attention. In *ICLR*.
- Gal Metzer, Elad Richardson, Or Patashnik, Raja Giryes, and Daniel Cohen-Or. 2022. Latent-NeRF for Shape-Guided Generation of 3D Shapes and Textures. *arXiv e-prints* (2022), arXiv-2211.
- Oscar Michel, Roi Bar-On, Richard Liu, Sagie Benaim, and Rana Hanocka. 2022. Text2Mesh: Text-Driven Neural Stylization for Meshes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 13492–13502.
- Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. 2021. Nerf: Representing scenes as neural radiance fields for view synthesis. *Commun. ACM* 65, 1 (2021), 99–106.
- Lucio Moser, Chinyu Chien, Mark Williams, Jose Serra, Darren Hendler, and Doug Roble. 2021. Semi-Supervised Video-Driven Facial Animation Transfer for Production. *ACM Trans. Graph.* 40, 6, Article 222 (dec 2021), 18 pages. <https://doi.org/10.1145/3478513.3480515>
- Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. 2022. Instant Neural Graphics Primitives with a Multiresolution Hash Encoding. *ACM Trans. Graph.* 41, 4, Article 102 (July 2022), 15 pages. <https://doi.org/10.1145/3528223.3530127>
- Jacob Munkberg, Jon Hasselgren, Tianchang Shen, Jun Gao, Wenzheng Chen, Alex Evans, Thomas Müller, and Sanja Fidler. 2022. Extracting Triangular 3D Models, Materials, and Lighting From Images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8280–8290.
- Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. 2021. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741* (2021).
- OpenAI. 2022. CLIP-ViT-Large-Patch14. <https://huggingface.co/openai/clip-vit-large-patch14>.
- Roy Or-El, Xuan Luo, Mengyi Shan, Eli Shechtman, Jeong Joon Park, and Ira Kemelmacher-Shlizerman. 2022. StyleSDF: High-Resolution 3D-Consistent Image and Geometry Generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 13503–13513.
- Or Patashnik, Zongze Wu, Eli Shechtman, Daniel Cohen-Or, and Dani Lischinski. 2021. Styleclip: Text-driven manipulation of stylegan imagery. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2085–2094.
- Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. 2022. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988* (2022).
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning Transferable Visual Models From Natural Language Supervision. In *Proceedings of the 38th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 139)*, Marina Meila and Tong Zhang (Eds.). PMLR, 8748–8763. <https://proceedings.mlr.press/v139/radford21a.html>

- Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. 2022. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125* (2022).
- Anurag Ranjan, Timo Bolkart, Soubhik Sanyal, and Michael J Black. 2018. Generating 3D faces using convolutional mesh autoencoders. In *Proceedings of the European Conference on Computer Vision (ECCV)*. 704–720.
- Scott Reed, Zeynep Akata, Xinchen Yan, Lajanugen Logeswaran, Bernt Schiele, and Honglak Lee. 2016. Generative adversarial text to image synthesis. In *International conference on machine learning*. PMLR, 1060–1069.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10684–10695.
- RunwayML. 2022. Stable Diffusion v1.5. <https://huggingface.co/runwayml/stable-diffusion-v1-5>.
- 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, Tim Salimans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. 2022. Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding. <https://doi.org/10.48550/ARXIV.2205.11487>
- Aditya Sanghi, Hang Chu, Joseph G. Lambourne, Ye Wang, Chin-Yi Cheng, Marco Fumero, and Kamal Rahimi Malekshan. 2022. CLIP-Forge: Towards Zero-Shot Text-To-Shape Generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 18603–18613.
- Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. 2015. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*. PMLR, 2256–2265.
- Meghana Rao Somepalli, M.D. Sai Charan, S Shruthi, and Suja Palaniswamy. 2021. Implementation of Single Camera Markerless Facial Motion Capture using Blendshapes. In *2021 IEEE International Conference on Computation System and Information Technology for Sustainable Solutions (CSITSS)*. 1–6. <https://doi.org/10.1109/CSITSS54238.2021.9683460>
- Jiaxiang Tang. 2022. Stable-dreamfusion: Text-to-3D with Stable-diffusion. <https://github.com/ashawkey/stable-dreamfusion>.
- Ayush Tewari, Mohamed Elgharib, Gaurav Bharaj, Florian Bernard, Hans-Peter Seidel, Patrick Pérez, Michael Zollhofer, and Christian Theobalt. 2020. Stylerig: Rigging stylegan for 3d control over portrait images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6142–6151.
- Kevin Turner. 2022. Decoding Latents to RGB Without Upscaling. <https://discuss.huggingface.co/t/decoding-latents-to-rgb-without-upscaling/23204/2>.
- Arash Vahdat, Karsten Kreis, and Jan Kautz. 2021. Score-based generative modeling in latent space. *Advances in Neural Information Processing Systems* 34 (2021), 11287–11302.
- Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. 2018. High-Resolution Image Synthesis and Semantic Manipulation with Conditional GANs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Xintao Wang, Liangbin Xie, Chao Dong, and Ying Shan. 2021. Real-ESRGAN: Training Real-World Blind Super-Resolution with Pure Synthetic Data. In *International Conference on Computer Vision Workshops (ICCVW)*.
- Zhouxia Wang, Jiawei Zhang, Runjian Chen, Wenping Wang, and Ping Luo. 2022. RestoreFormer: High-Quality Blind Face Restoration from Undegraded Key-Value Pairs. (2022).
- Haotian Yang, Hao Zhu, Yanru Wang, Mingkai Huang, Qiu Shen, Ruigang Yang, and Xun Cao. 2020. Facescape: a large-scale high quality 3d face dataset and detailed riggable 3d face prediction. In *Proceedings of the ieee/cvf conference on computer vision and pattern recognition*. 601–610.
- Longwen Zhang, Chuxiao Zeng, Qixuan Zhang, Hongyang Lin, Ruixiang Cao, Wei Yang, Lan Xu, and Jingyi Yu. 2022. Video-Driven Neural Physically-Based Facial Asset for Production. *ACM Trans. Graph.* 41, 6, Article 208 (nov 2022), 16 pages. <https://doi.org/10.1145/3550454.3555445>
- Yiyu Zhuang, Hao Zhu, Xusen Sun, and Xun Cao. 2022. Mofanerf: Morphable facial neural radiance field. In *European Conference on Computer Vision*. Springer, 268–285.