

UVMMap-ID: A Controllable and Personalized UV Map Generative Model

Weijie Wang^{1,3,6*}, Jichao Zhang^{2*†}, Chang Liu¹, Xia Li³, Xingqian Xu⁴, Humphrey Shi⁵, Nicu Sebe¹, Bruno Lepri⁶

¹University of Trento, Trento Italy, ²Ocean University of China, Qingdao, China, ³ETH Zürich, Zürich, Switzerland,

⁴PicsArt AI Research, Atlanta, United States, ⁵Georgia Institute of Technology, Atlanta, United States

⁶Fondazione Bruno Kessler, Trento, Italy

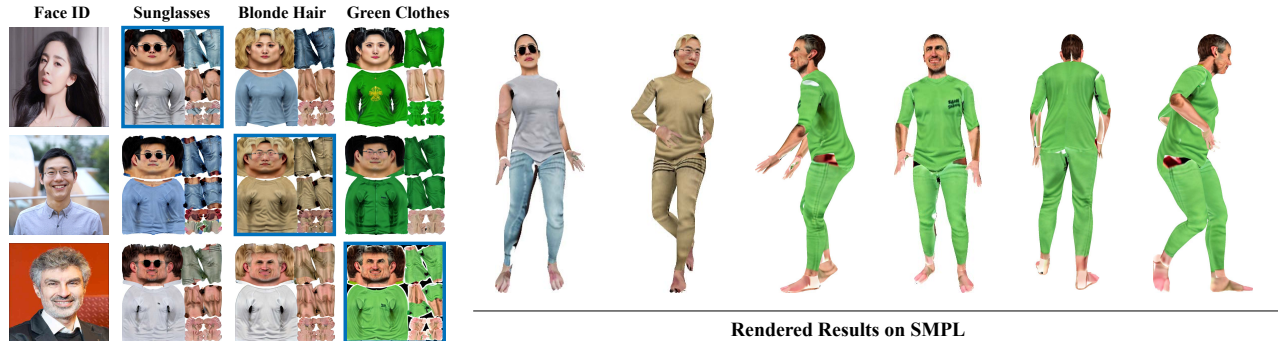


Figure 1: Our method can synthesize high-quality textures while enabling a controllable and personalized generation with the given text prompts and Face ID (Left). The textures can be directly applied to SMPL meshes [29] (Right).

ABSTRACT

Recently, diffusion models have made significant strides in synthesizing realistic 2D human images based on provided text prompts. Building upon this, researchers have extended 2D text-to-image diffusion models into the 3D domain for generating human textures (UV Maps). However, some important problems about UV Map Generative models are still not solved, i.e., how to generate personalized texture maps for any given face image, and how to define and evaluate the quality of these generated texture maps. To solve the above problems, we introduce a novel method, UVMMap-ID, which is a controllable and personalized UV Map generative model. Unlike traditional large-scale training methods in 2D, we propose to fine-tune a pre-trained text-to-image diffusion model which is integrated with a face fusion module for achieving ID-driven customized generation. To support the finetuning strategy, we introduce a small-scale attribute-balanced training dataset, including high-quality textures with labeled text and Face ID. Additionally, we introduce some metrics to evaluate the multiple aspects of the textures. Finally, both quantitative and qualitative analyses demonstrate the effectiveness of our method in controllable and

The first two authors contribute equally, and Jichao Zhang is the corresponding author.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

ACM MM, 2024, Melbourne, Australia

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

personalized UV Map generation. Code is publicly available via <https://github.com/twowwj/UVMMap-ID>.

CCS CONCEPTS

• **Do Not Use This Code** → **Generate the Correct Terms for Your Paper**; *Generate the Correct Terms for Your Paper*; Generate the Correct Terms for Your Paper; Generate the Correct Terms for Your Paper.

KEYWORDS

Generative Model, Diffusion Model, 3D Avatar Generation, Multi-Modal Generation

ACM Reference Format:

Weijie Wang^{1,3,6*}, Jichao Zhang^{2*†}, Chang Liu¹, Xia Li³, Xingqian Xu⁴, Humphrey Shi⁵, Nicu Sebe¹, Bruno Lepri⁶. 2024. UVMMap-ID: A Controllable and Personalized UV Map Generative Model. In *Proceedings of the 32nd ACM International Conference on Multimedia (MM'24)*, October 28–November 1, 2024, Melbourne, Australia. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

The development of 3D human models has garnered significant attention in recent years, owing to its versatile applications across various domains, including filmmaking, video games, augmented reality/virtual reality (AR/VR), and human-robot interaction. Among the myriad tasks essential for crafting digital humans, texture synthesis stands out as a pivotal element in achieving the photorealistic

quality of 3D avatars. However, creating 3D textures in the traditional computer graphics pipeline is time-consuming and labor-intensive. Thus, it is important to utilize generation techniques to design diverse texture maps automatically.

Texture (UV map) generation has been a focus in previous approaches for tasks such as 3D face and human reconstruction. These methods leverage generators from Generative Adversarial Networks (GANs) to estimate textures either in an unsupervised [9, 45, 52, 59] or supervised [24, 25] manner. Subsequently, the texture estimation model is integrated into the avatar fitting stage. Nonetheless, these methods are limited in generating novel textures and need more support for controllable generation.

Large-scale text-to-image diffusion models [36, 38], nowadays, have been proven very effective over cross-model generation tasks, which should mainly attributed to the scalable 2D image-text data pairs along with large-scale parallel computation. Yet we notice that the lack of large-scale 3D texture data makes training high-quality texture generative models quite challenging. Inspired by the pre-trained strategy of DreamBooth, SMPLitex [5] has employed a few texture maps (UV defined by SMPL [29]) to fine-tune a pretrained text-to-image diffusion model. It has been observed that this approach enables the synthesis of texture maps while supporting its foundation text-driven task. However, the inability of SMPLitex to support personalized texture generation poses a significant limitation on their approach, particularly in applications where user customization is crucial. Personalized texture generation enables the tailoring of textures to specific individual preferences, fostering a comprehensive experience in 3D applications, including avatars, VR, and gaming. Besides personalization, evaluating the quality of generated textures within the UV space remains an unresolved challenge, leaving more space for research.

In this paper, we introduce the UVMap-ID method, a UV map generative model that supports ID-driven personalized generation tasks. Specifically, we fine-tune a pretrained text-to-image diffusion model using a small-scale training dataset. In contrast to 2D personalized methods [7, 46, 49, 56] that necessitate large-scale training data in 2D methods, our dataset, which is attribute-balanced (i.e., "Race and Gender"), comprises around 750 image-ID pairs: the textures map with annotated text prompts, the corresponding portrait faces. To enable the ability of ID-driven personalized generation, we extend the stable diffusion with an additional face fusion module. Moreover, we introduce some corresponding metrics to evaluate the quality of generated textures from multiple aspects, i.e., fidelity, structure preservation, ID preservation, and text-image alignment. Remarkably, our model achieves high-quality and diverse texture synthesis within just several hours of training, while also supporting controllable and personalized synthesis with the user-provided image ID.

In summary, our contributions are as follows:

- We are the first to propose a controllable and personalized UV map generative model capable of synthesizing diverse and personalized texture maps.
- We propose an efficient fine-tuning strategy for training an ID-driven extension architecture for StableDiffusion, utilizing only a small-scale training dataset.

- We utilize our method to produce a new dataset, containing around 5k UVMap-ID image pairs, and will make it publicly available. Our small-scale attribute-balanced training dataset, the larger-scale dataset, and metrics for textures play a bridging role in guiding subsequent work in this field.

2 RELATED WORK

UV-Map Generative Model. This model aims to generate diverse textures based on the generative models, such as Generative Adversarial Networks [10], Diffusion Models [13, 43]. Existing works utilize this technique in the 3D face reconstruction with the 3D morphable model (3DMM) [3] or human reconstruction with the SMPL [29]. For face texture generation, GANFIT [9] first uses 10,000 high-resolution textures to train the GAN generator, then takes this GAN generator as the statistical parametric representation of the facial texture in the fitting progress. To avoid the training using the limited numbers and diversity of texture map, StyleUV [25] integrates the 2D image fitting and rendering stages into the adversarial networks. Additionally, some methods focus on contributing the 3D facial UV-texture datasets, such as Facescape [55], and FFHQ-UV [1]. For human texture generation, most of the works learn to recover the full texture from a single human image. The Re-Identification metric as supervised in this task is proposed [45]. To further improve the quality of texture generation, Zhao. et al [59] introduce a consistency learning to enforce the cross-view consistency of texture prediction during training. Texformer [52] introduces the transformer architecture to exploit global information of the input, effectively facilitating higher-quality texture generation. Different from these methods without using any ground-truth 3D textures, Verica. et al [24] non-rigidly registers the SMPL model to thousands of 3D scans, and encodes the appearances as texture maps. And these 3D textures are used to train a texture completed model. However, these mentioned methods cannot support diverse and text-guided texture generation. The most related work to ours is SMPLitex [5]. Motivated by the Dreambooth [37], SMPLitex utilizes a few texture maps to fine-tune the pretrained text-guided diffusion model to enable the textures inpainting and text-guided texture generation task. Compared to SMPLitex, our method supports both text-guided and ID-driven personalized texture generation.

Text-to-3D Avatar Generation. Text-guided 3D content generation has achieved great success with the development of 3D representation methods and generative models. Lots of methods utilize the frozen image-text joint embedding models from CLIP [33] to optimize the underlined 3D representation, such as NeRF [30] where some of them work on generation for general 3D object [18, 31, 40, 50, 54], or human Avatar [14, 16]. The most famous work is Dream Fields [18] which first demonstrated the effectiveness of combining the CLIP model and NeRF representation for 3D object creation, but 3D objects produced by this approach tend to lack realism and accuracy. DreamFusion [32] introduces Score Distillation Sampling (SDS) loss which is based on probability density distillation that enables the use of a pretrained 2D diffusion model as a prior for optimization of a parametric NeRF representation. By using SDS loss instead of CLIP, DreamFusion generates high-quality coherent 3D objects while aligning with the given text prompt. Recently, many similar methods with SDS loss have occurred to

improve text-to-3D results in various aspects, such as enhancing the realism of rendering with detailed geometry [6], solving the multiple-view inconsistency problem [27, 42] or using variational score distillation (VSD) [47] method instead of SDS to improve the fidelity and diversity of 3D content generation. However, high-quality human avatars remain a challenge due to the complexity of the human body’s shape, pose, and appearance. To make the avatar animatable, DreamAvatar [4] and AvatarCraft [19] integrate the SMPL prior into the NeRF or SDF representation with a deformable field. To improve the avatar’s quality and avoid the cartoon-like appearance, DreamHuman [23] uses a spherical harmonics lighting model instead of diffuse reflectance model and additionally optimizes a spherical harmonics coefficients; HumanNorm [17] introduces a normal diffusion model to enhance the diffusion model’s understanding of 3D geometry to further improve the texture and geometry’s quality. More recently, HumanGaussian [28] integrates 3D Gaussian representation instead of NeRF into 3D Human Avatar generation to reduce training time. Compared with these text-to-3D works, we focus on achieving a controllable texture generation but don’t care about the generation of geometry.

Text-Driven Personalized Diffusion Models. Diffusion model [13, 43], is a class of generative modeling in which it iteratively transforms noises to samples simulating the true data distribution. Diffusion models generally outperformed other traditional methods, such as GANs, due to the fact that the output quality has been notably improved across diverse domains. Diffusion models are widely used for text-to-image generation [34, 36, 38], and also stand out supporting more cross-model tasks [2, 35, 53]. One of the foundation works, Stable diffusion [36], applies the diffusion process on latent space, reducing training computation while preserving quality. While other methods, such as Imagen [38] and DALL-E2 [34], generate samples directed over pixel space, have also proven effective. Finetune-wise, DreamBooth [37] and LoRA [15] introduces a subject-driven training approach, enabling text controls, and offers a compelling feature for precise personalizing. Text Inversion [8] and VideoBooth [20] suggest an alternative solution via latent inversion before editing. Another class of methods [7, 46, 48, 49, 51, 56–58, 60] extends the model with additional networks to extract and adopt conditional inputs that guide the generation. Representatively, IP-Adapter [56] introduces a decoupled U-Net that injects conditional hidden features to the original diffusion U-Net, achieving an accurate control from the reference input. Some concurrent 2D methods such as Instant-ID [46], Infinite-ID [49] and SSR-Encoder [58], also attracted lots of attention. In this work, we share goals similar to IP-Adapter and Instant-ID, focusing on 3D human texture rather than 2D generation.

3 METHODS

Given a reference portrait describing the facial appearance (Face ID) of the target individual, our model aims to generate a texture that aligns with the facial appearance of the target person and fits the structure of the UV map defined by SMPL. In this section, we first provide a brief introduction to Denoising Diffusion Probabilistic Models [13] in Section 3.1, laying the foundational framework and network architecture for our method. Subsequently, detailed explanations of design specifics are presented in Section 3.2. Then, we

will explain the pipeline we use to build the dataset in Section 3.3. Finally, we introduce some metrics for UV textures in Section 3.4.

3.1 Preliminary: Denoising Diffusion Probabilistic Models

The denoising diffusion probabilistic models operate by simulating a forward process that adds noise to an image or its latent representation over a series of time steps, transforming them into Gaussian noise. Conversely, the reverse process seeks to recover the original image or latent representation by iterative denoising. This bidirectional process is key to the diffusion models’ ability to generate high-fidelity images. Our work leverages Stable Diffusion (SD), a pretrained generative model that could generate high-quality images from a text prompt. Specifically, given an image x , SD first uses a pretrained autoencoder to encode x into latent: $z = \mathcal{E}(x)$. Then, noise is gradually added to z over a sequence of T steps, transitioning the data distribution from the original data distribution to a Gaussian Noise distribution, and the noise added forward a Markov chain of conditional Gaussian distributions defines the process:

$$q(z_t|z_{t-1}) = \mathcal{N}(z_t; \sqrt{1 - \beta_t}z_{t-1}, \beta_t\mathbf{I}),$$

where β_t is the variance schedule. During training, the denoising u-net ϵ_θ of SD aims to learn to reconstruct the original latent z from the noise, modeled by:

$$p_\theta(z_{t-1}|z_t) = \mathcal{N}(z_{t-1}; \mu_\theta(z_t, t), \sigma_\theta^2(z_t, t)\mathbf{I}),$$

and the learning objective is defined as follows:

$$L(\theta) = \mathbb{E}_{z_t, c, \epsilon, t} [\|\epsilon - \epsilon_\theta(z_t, c, t)\|^2],$$

where c represents text conditional embeddings.

3.2 Fine-Tuning Text-to-Image Models for ID-Driven UV Map Generation

Fig. 2 provides the pipeline of our proposed approach. The initial input to the pipeline consists of random noise and a reference portrait. Our text-to-image model is configured based on the design of SD, employing the same framework and trained weights of SD. Motivated by DreamBooth [37], we propose to utilize the finetuning strategy with a prior preservation loss (Fig. 2 (Left)) applying to text-to-image diffusion architecture integrating with a face fusion module (Fig. 2 (Right)).

3.2.1 Face Fusion Module. To enable Stable Diffusion to accept additional image information, (i.e., the portraits), the previous methods mainly leverages the CLIP image encoder, either directly substituting the CLIP text encoder or through decoupled cross-attention mechanism to separate cross-attention layers for text features and image features [34, 56]. Nevertheless, the CLIP image encoder is constrained by its operation on images of lower resolution, which particularly impacts its efficacy in encoding face images by failing to encapsulate comprehensive details. Moreover, CLIP’s architecture, fundamentally designed to align semantic features between text and images, mainly focuses on high-level feature correspondence. This orientation towards semantic feature matching inadvertently results in a dilution of finer, detailed features during the encoding process, posing a challenge for applications requiring precise detail retention. Hence, we propose to use the face embedding extracted

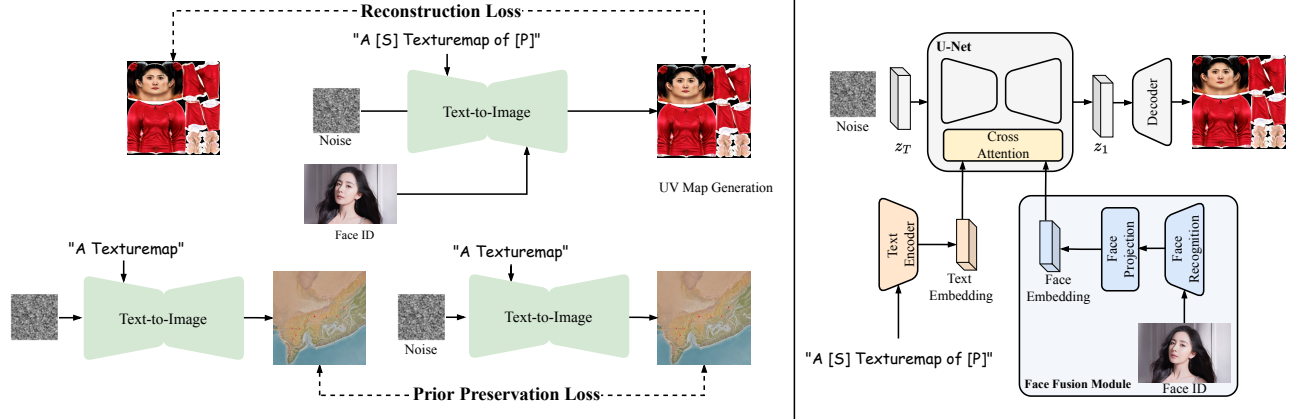


Figure 2: The left side of the figure shows the overview of our proposed pipeline. Given a reference image as face ID, we utilize a pre-trained text-to-image diffusion model, where the input is a combination of a noised UV Map and text prompt of a unique identifier and characteristics of the portrait where "A [S] Texturemap of [P]," where [S] is a unique identifier and [P] represents the race and gender. To maintain the quality of images generated by the pre-trained model and effectively process textual features, we adopt a prior preservation loss. The right side of the figure shows the detailed architecture of our model, where facial information is mapped to the same dimensions as text embeddings through a facial recognition model and face projection layers. Subsequently, we merge facial and textual information via decoupled cross-attention, which is then integrated into the pre-trained text-to-image model.

by the face recognition models and linear projection layers to provide SD with human face information. Also, to preserve the original model's ability to process text information while integrating image information, we adopt the decoupled cross-attention mechanism [56], ensuring a seamless blend of both modalities. Given query feature Z , image feature c_i and the text feature c_t , the output Z' of decoupled cross-attention layers is:

$$Z' = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V + \text{softmax}\left(\frac{Q(K')^T}{\sqrt{d_k}}\right)V',$$

where $Q = ZW_q$, $K = c_tW_k$, $V = c_tW_v$, $K' = c_tW'_k$, $V' = c_tW'_v$, and the W_q , W_k , W_v , W'_k and W'_v are learnable parameters of the projection layers. Similar fusion modules have been utilized in some concurrent 2D methods [46, 49].

3.2.2 Prior Preservation Loss. We observed that when using "UV texture map" as the text prompt, SD often fails to generate any correct UV maps. This is likely because SD is trained on data scraped from the internet, where real UV texture maps are rarely found in the training resources. Also, our goal is to generate images with a small training set (about 750 images in our dataset), each featuring different facial characteristics of individuals, and generating accurate faces has always been a weakness of SD. Additionally, our input incorporates extra face image information, and during fine-tuning, we would like to ensure our model does not lose SD's original capability to correctly process textual information. To this end, we introduced prior preservation loss, as proposed in Dreambooth [37], to ensure the model retains its generalization ability and does not overfit the few-shot examples provided during the personalization process.

However, our objectives differ fundamentally from Dreambooth in two ways. Firstly, Dreambooth targets subject-driven generation, whereas our model aims at generating specific formats of images, the UV texture maps. This leads to a situation where Dreambooth requires re-fine-tuning the entire SD for each subject, while our model, after training, can generate corresponding UV maps for any input face ID. This distinction arises because, in DreamBooth, one unique identifier represents a single unique subject, whereas our unique identifier [S] denotes one unique kind of image structure (UV Map defined by SMPL). Secondly, we added extra facial information [P] to our text prompts during training to further preserve the original capabilities of the text encoder, enabling it to effectively parse attributes such as race and gender. For detailed experiments, please refer to Section 4.4

Formally, the training loss of our model is defined as:

$$L(\theta) = \mathbb{E}_{z_t, c, \epsilon, t} \left[\|\epsilon - \epsilon_\theta(z_t, c, t)\|^2 \right] + \mathbb{E}_{z_t, c', \epsilon, t} \left[\|\epsilon_{\text{pr}} - \epsilon_\theta(z_t, c', t)\|^2 \right],$$

where c' is a fixed conditional text prompt "a texturemap" and ϵ_{pr} is the generate data using the frozen diffusion model with c' .

3.3 Dataset

Training Dataset In this part, we describe the process of constructing our dataset, which is centered around the generation of high-quality and diverse UV texture maps for digital human models. Our approach can be segmented into three stages:

1) **Celebrity Selection:** In the initial phase of our dataset creation, we aimed for a balanced and inclusive representation by employing OpenAI's ChatGPT to generate a list of 150 celebrities. Our selection was structured to include equal representation across three ethnic

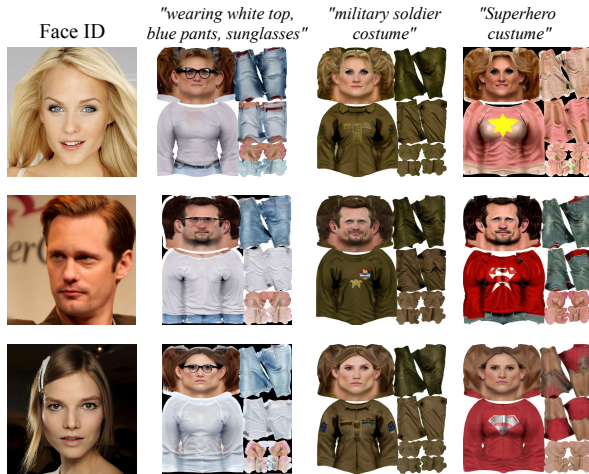


Figure 3: Personalized textures generation results using face IDs from CelebA-HQ dataset.

groups: African American, Asian, and White, with 50 celebrities from each group. To further enhance the diversity and applicability of our dataset, we ensured gender balance within each ethnic category, selecting 25 male and 25 female celebrities. We use celebrities because SMPLitex accepts only text input, and celebrity portraits are readily available. This approach allows us to link names, portraits, and corresponding UV texture maps effectively.

2) UV Texture Map Generation: We employed SMPLitex to generate UV texture maps for each of the selected celebrities. This process resulted in 50 UV texture maps per celebrity, totaling 7,500 initial texture maps.

3) Manual Selection: To ensure the highest quality and relevance for our dataset, we manually reviewed the generated UV texture maps and selected 5 maps per celebrity that best met our predefined criteria. These criteria included clarity, detail accuracy, and representation quality of ethnic features. This manual selection process narrowed our dataset to 750 UV texture maps with 5 UV texture maps per ID.

A New Dataset: CelebA-HQ-UV We utilize our method with personalized generation to produce a new dataset, which contains 5k UVMap-ID pairs. Specifically, we select 5000 high-resolution face images from CelebA-HQ [21] as reference image IDs of our methods. For every ID, our method produces 10 textures and selects 2 by the evaluation of multiple aspects, i.e., the quality of textures, the preservation of UV structure, and the preservation of face ID. Fig. 3 shows some results using three face IDs from CelebA-HQ. We refer to this dataset as CelebA-HQ-UV, and will make it publicly available. Note that we define a list of text prompts for these generations which will be introduced in the supplementary material.

3.4 Metrics

As previously mentioned, assessing the quality of generated textures within the UV space defined by SMPL poses a significant challenge, especially within the scope of our personalized generation task. In this paper, we introduced four metrics to evaluate

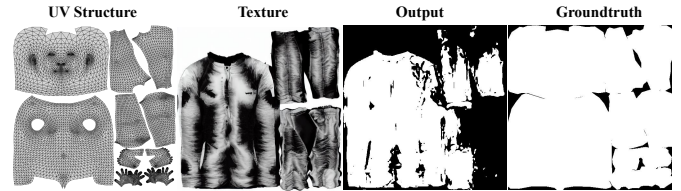


Figure 4: It shows UV structures, textures from SMPLitex, extracted semantic segmentation, and semantic groundtruth from left to right.

the quality of the generated textures from multiple aspects: Inception Scores [39] to evaluate the fidelity and diversity, Semantic Structure Preservation (SSP) to evaluate structure preservation of UV space defined by SMPL [29], Deep Face Recognition (DFR) to evaluate Face ID preservation and CLIP-Text (CLIPT) [20, 48] score to evaluate the text-image alignment.

Inception Score (IS) on UV textures and rendered results The Inception Score (IS) and Fréchet Inception distance [12] are widely utilized metrics for evaluating the diversity and quality of 2D images generated by generative models. FID is a well-established measure that compares the inception similarity score between distributions of generated and real images. One key distinction between IS and FID is that IS is computed solely using fake samples, eliminating the need for real samples in its calculation. Due to the lack of real sample distribution, we employ the IS to directly evaluate the quality of 5000 generated textures rather than FID. We refer to IS on textures of UV space as IS (UV). Additionally, we render these textures into 2D space by applying them to the SMPL Mesh. Subsequently, we utilize IS to evaluate the quality of 5000 rendered human images in 2D space. We refer to this type of IS as IS (R).

Semantic Structure Preservation (SSP) To assess the preservation of UV structures in generated textures, we introduce a novel metric termed Semantic Structure Preservation (SSP). Notably, we have observed instances where the generated textures from SMPLitex [5] may not faithfully retain these underlying structures, as illustrated in Fig. 4. The SSP metric is designed to quantify this preservation. We leverage off-the-shelf human parsing techniques [26] to extract semantic segmentation from the generated images and then compare it with ground truth segmentation (Fig. 4 (right)). We conduct this comparison across a dataset comprising 1000 images and compute the mean difference as the SSP score.

Deep Face Recognition (DFR) To assess the preservation of identity (ID) within textures, a crucial aspect of personalized image generation tasks, we propose employing Deep Face Recognition (DFR) methods to quantify the similarity between generated textures and reference facial images. Specifically, we leverage the off-the-shelf tool [41] to do face recognition between the textures and image ID. We use 10 face IDs, and 100 samples for every ID and report the successful numbers. We refer to this metric as the DFR score which is reported as a measure of the preservation of identity within the generated textures.

CLIP-Text (CLIPT) To measure the alignment of the generated textures and given text prompts, we use the CLIP-Text (CLIPT) score followed by 2D methods [20, 48]. This metric is calculated




	Face ID	<i>"wearing yellow clothes"</i>	<i>"wearing sunglasses"</i>	<i>"wearing white shirt and jeans"</i>	<i>"military soldier costume"</i>	<i>"santa claus costume"</i>
Asian woman						
Asian man						
Asian man						
White woman						
Asian woman						
Asian woman						
Asian woman						
Asian man						
White man						

Figure 5: Our personalized generation results. The 1st column shows reference faces, obtained from the website, and not existing in our training set.

Methods	IS (R) \uparrow	IS (UV) \uparrow	SSP \downarrow	CLIPT \uparrow	DFR \uparrow
SMPLitex [5]	1.46 \pm 0.020	1.95 \pm 0.049	10.45	29.40	62
UVMap-ID	1.78 \pm 0.020	1.89 \pm 0.027	8.46	29.12	792

Table 1: Quantitative results using four metrics: inception scores on rendered images (IS (R)), inception scores on UV maps (IS (UV)), Semantic Structure Preservation (SSP), CLIP Text (CLIPT), Deep Face Recognition (DFR).

using the cosine similarity of the CLIP text embeddings of the given text prompts and CLIP image embeddings of the generated textures. We compute the CLIPT score using 1000 text-prompt pairs.

4 EXPERIMENTS

4.1 Training Details

Our experiments are based on the Realistic_Vision_V4 model, which is further fine-tuned on Stable Diffusion v_1.5 [36], and could produce more photorealistic images. Additionally, we utilize the buffalo_l pre-trained face recognition model from SCRFD [11], and pre-trained projection layers from [56]. The experimental code is developed using the HuggingFace Diffusers library [44]. During training, we fine-tune the entire U-Net, text encoder and face projection layers, and keep the VAE encoder and decoder of Stable Diffusion frozen. The UVMap-ID training is conducted on a single machine equipped with an A40 GPU for 1500 steps, with a batch size of 2. We employ the AdamW optimizer [22] with a fixed learning rate of 1e-6 and a weight decay of 0.01. Our dataset comprises images with a resolution of 512 \times 512, hence we generate images at this resolution during training. In the inference phase, we use a 50-step DDIM sampler [43] and set the classifier-free guidance scale to 7.5.

4.2 Baselines

We take the texture generation model SMPLitex [5] as the baseline. And all results from SMPLitex are produced from their released code and pretrained model. SMPLitex does not support image-driven personalized generation. Thus, we provide image ID's name in the text prompts for SMPLitex, but not for our method.

4.3 Comparisons

Fig. 5 shows diverse personalized texture generation results from our methods. Our reference face IDs (1st column images) are collected from a diverse range of sources on the website, thus encompassing a wide variety of characteristics, including different ethnicities, genders, occupations, levels of fame, and even facial poses. As shown in the 2nd-6th columns of Fig. 5, our generated UV textures effectively preserve the identity features of these reference face IDs, demonstrating the effectiveness and robustness of our methods in personalized generation. Moreover, our method also achieves accurate text-driven controllable generation.

We conducted visualization comparisons with SMPLitex [5], as depicted in Fig. 6. Notably, SMPLitex is not an image-driven method. Therefore, while we utilized some well-known celebrities as image IDs and provided their names in text prompts for SMPLitex, we deliberately omitted this information for our method to ensure

Methods	DFR \uparrow
UVMap-ID w/o "Race and Gender"	436
UVMap-ID w/ "Race and Gender"	792

Table 2: Ablation Study for "Race and Gender" label.

Methods	IS (R) \uparrow	IS (UV) \uparrow	SSP \downarrow	CLIPT \uparrow	DFR \uparrow
UVMap-ID (1)	1.88 \pm 0.028	2.03 \pm 0.039	10.59	29.09	734
UVMap-ID (2)	1.78 \pm 0.020	1.89 \pm 0.027	8.46	29.12	792
UVMap-ID (5)	1.55 \pm 0.017	1.55 \pm 0.084	8.74	29.27	798

Table 3: Ablation studies of Training data. UVMap-ID (N) denotes the number (N) of textures for each ID in the training stage.

a fairer comparison. Remarkably, our results exhibit a higher degree of similarity in face ID preservation compared to SMPLitex, underscoring the superiority of our method in maintaining identity features during personalized texture generation. Moreover, our approach also demonstrates superior structural preservation compared to SMPLitex, as evidenced by the "Jay Chou" row (Top-Right).

Quantitative results using four metrics are shown in Table 1. We observe that SMPLitex achieves better IS (UV) scores than our method. We attribute this to the fact that our approach is image-driven, which means that the provided reference ID constrains the diversity of generated images, a crucial aspect of IS. In contrast, our method achieves a higher IS (R) than SMPLitex. As mentioned, SMPLitex often struggles to preserve UV structures effectively, resulting in unrealistic renderings. The comparison of structure preservation can be validated by our achieved superior SSP score. Moreover, our DFR score significantly outperforms the Baseline, validating that our method achieves better similarity to the target ID in personalized texture generation tasks. Additionally, the high success rate of 837 out of 1000 demonstrates the robustness of our method to reference images. Furthermore, we observe that our CLIPT score is comparable to the baseline, indicating that the "image prompt" generated by our image encoder does not significantly affect the control capability of the text prompt.

4.4 Ablation Studies

"Race and Gender" in prompts As shown in Fig. 7, we analyze the impact of including race and gender labels in prompts during training, assessing how this additional information affects generative model performance. As indicated in Table 2, incorporating race and gender labels significantly enhances the model's DFR score compared to the version without these labels (UVMap-ID w/o "Race and Gender"). This indicates that the facial recognition model we use focuses more on the structural information of the human face, while the label supplements the missing information such as skin color.

Training Data In this part, we explore the impact of varying the number of UV maps used per image ID during training. Our model, UVMap-ID, is evaluated using a consistent training strategy, except that each image ID in the training dataset is processed using 1, 2, or 5 UV maps. These setups are denoted as UVMap-ID (1), UVMap-ID (2), and UVMap-ID (5) respectively.

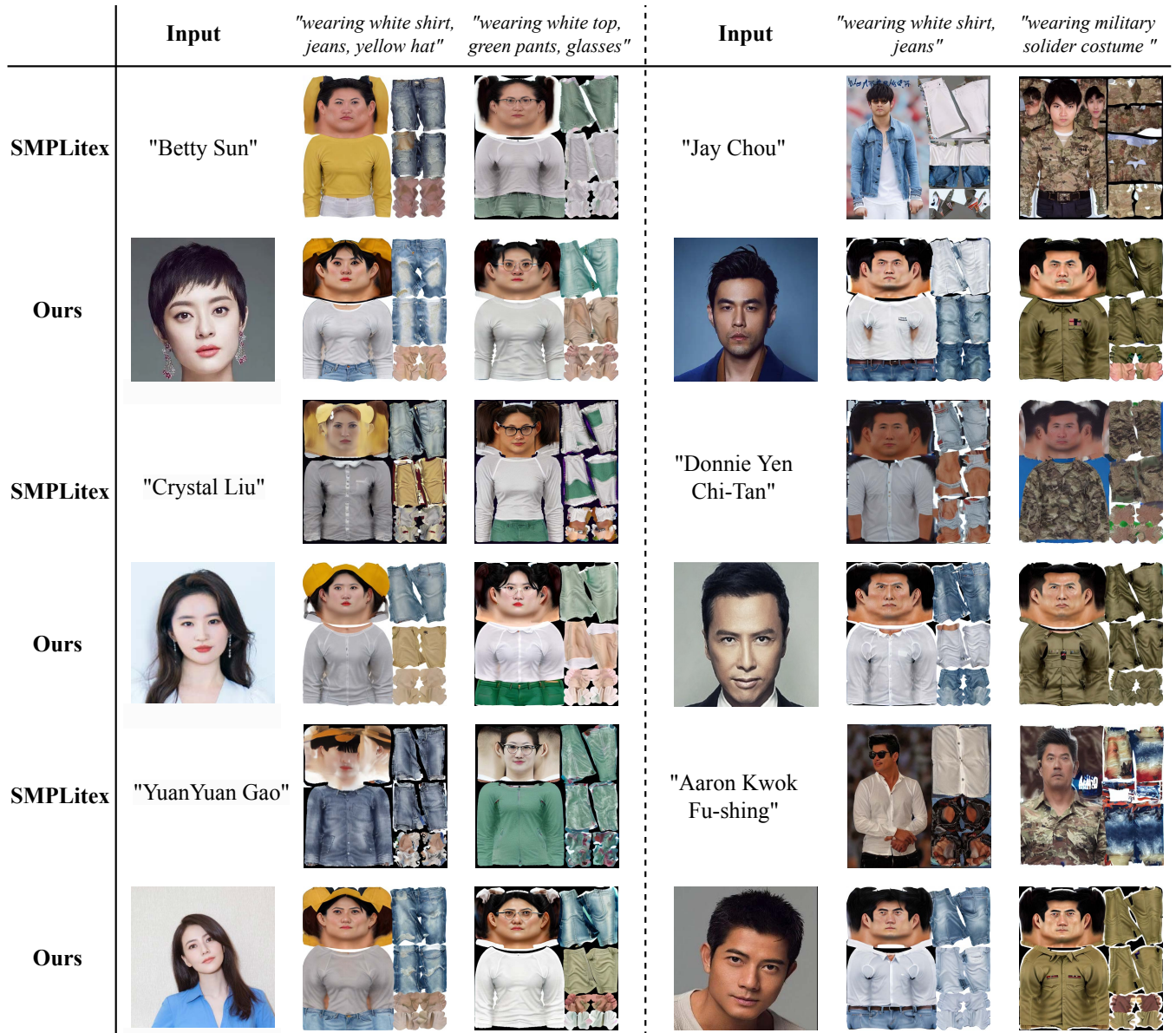


Figure 6: Comparison with SMPLitex [5] results. SMPLitex is not an image ID-driven method. Thus, we provided these celebrities' names in the test prompts for SMPLitex, but not for ours. Taking "Betty Sun" as an example (upper-left corner), the test prompt of SMPLitex is "a texturemap of Betty Sun wearing...", and our test prompt is "a texturemap of Asian woman wearing...". Note that image IDs are not existing in our training data.

Table 3 highlights the performance metrics across these configurations. Based on the results shown in Table 3, we have chosen UVMaP-ID (2) as our base model. This configuration utilizes two UV maps, which provide a diverse dataset sufficient to capture the critical variations in facial features, without overloading the pre-trained model. UVMaP-ID (2) strikes a balance, delivering remarkable realism in image generation while effectively maintaining the identity of reference images.

5 CONCLUSIONS

In this paper, we introduce UVMaP-ID, the first method for ID-driven personalized texture generation. UVMaP-ID takes the StableDiffusion as the backbone and extends it with an additional face fusion module. Moreover, our method is a highly efficient model with only several hours fine-tuning strategy on a small-scale dataset. Additionally, we also explore the evaluation of quality for UV textures and introduce some corresponding metrics. Finally, with user provided face images, our method can automatically create high-quality UV textures with the preservation of face ID while enabling

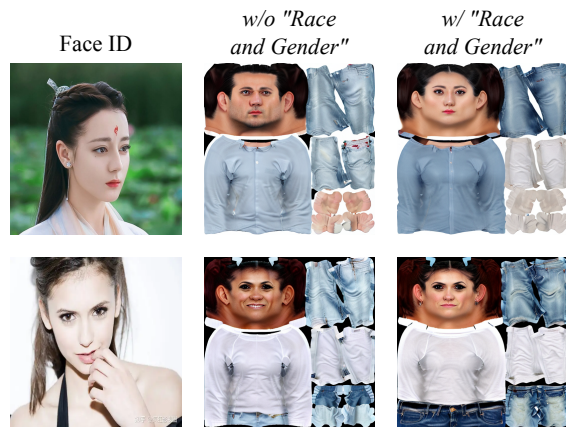


Figure 7: Qualitative ablation studies of between w/o and $w/$ "Race and Gender" labels. The 1st-row results show our full method preserves the "Gender" attribute and the 2nd-row results show our full method preserves the "Race" attribute.

text-driven controls, which is a very available application for 3D avatar creation in compute graphics fields. By using our method, we create a new dataset, CelebA-HQ-UV, comprising textures and face ID pairs. This dataset will be shared with the community to facilitate further research. We desire to explore the interactive editing of textures in the future.

6 ACKNOWLEDGEMENTS

This work has been partially supported by the European Union's Horizon Europe research and innovation program under grant agreement No. 101120237 (ELIAS). Bruno Lepri and Nicu Sebe also acknowledge the support of the PNRR project FAIR - Future AI Research (PE00000013), under the NRRP MUR program funded by the NextGenerationEU.

REFERENCES

- [1] Haoran Bai, Di Kang, Haoxian Zhang, Jinshan Pan, and Linchao Bao. 2023. Ffhquv: Normalized facial uv-texture dataset for 3d face reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 362–371.
- [2] Fan Bao, Shen Nie, Kaiwen Xue, Chongxuan Li, Shi Pu, Yaole Wang, Gang Yue, Yue Cao, Hang Su, and Jun Zhu. 2023. One transformer fits all distributions in multi-modal diffusion at scale. In *International Conference on Machine Learning*. PMLR, 1692–1717.
- [3] 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.
- [4] Yukang Cao, Yan-Pei Cao, Kai Han, Ying Shan, and Kwan-Yee K Wong. 2023. Dreamavatar: Text-and-shape guided 3d human avatar generation via diffusion models. *arXiv preprint arXiv:2304.00916* (2023).
- [5] Dan Casas and Marc Comino Trinidad. 2023. Smlitex: A generative model and dataset for 3d human texture estimation from single image. *BMVC* (2023).
- [6] Rui Chen, Yongwei Chen, Ningxin Jiao, and Kui Jia. 2023. Fantasia3D: Disentangling Geometry and Appearance for High-quality Text-to-3D Content Creation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.
- [7] Siying Cui, Jiankang Deng, Jia Guo, Xiang An, Yongle Zhao, Xinyu Wei, and Ziyong Feng. 2024. IDAdapter: Learning Mixed Features for Tuning-Free Personalization of Text-to-Image Models. <https://api.semanticscholar.org/CorpusID:268537084>
- [8] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. 2023. An image is worth one word: Personalizing text-to-image generation using textual inversion. *ICLR* (2023).
- [9] 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*. 1155–1164.
- [10] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial nets. *NeurIPS* (2014).
- [11] Jia Guo, Jiankang Deng, Alexandros Lattas, and Stefanos Zafeiriou. 2021. Sample and Computation Redistribution for Efficient Face Detection. *arXiv:2105.04714* [cs.CV]
- [12] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. 2017. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems* 30 (2017).
- [13] Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. *Advances in neural information processing systems* 33 (2020), 6840–6851.
- [14] Fangzhou Hong, Mingyuan Zhang, Liang Pan, Zhongang Cai, Lei Yang, and Ziwei Liu. 2022. Avatarclip: Zero-shot text-driven generation and animation of 3d avatars. *arXiv preprint arXiv:2205.08535* (2022).
- [15] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *ICLR 2022* (2021).
- [16] Shuo Huang, Zongxin Yang, Liangting Li, Yi Yang, and Jia Jia. 2023. AvatarFusion: Zero-shot Generation of Clothing-Decoupled 3D Avatars Using 2D Diffusion. In *Proceedings of the 31st ACM International Conference on Multimedia*. 5734–5745.
- [17] Xin Huang, Ruizhi Shao, Qi Zhang, Hongwen Zhang, Ying Feng, Yebin Liu, and Qing Wang. 2024. HumanNorm: Learning Normal Diffusion Model for High-quality and Realistic 3D Human Generation. *CVPR* (2024).
- [18] 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.
- [19] Ruixiang Jiang, Can Wang, Jingbo Zhang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao. 2023. AvatarCraft: Transforming Text into Neural Human Avatars with Parameterized Shape and Pose Control. *ICCV* (2023).
- [20] Yuming Jiang, Tianxing Wu, Shuai Yang, Chenyang Si, Dahua Lin, Yu Qiao, Chen Change Loy, and Ziwei Liu. 2024. VideoBooth: Diffusion-based Video Generation with Image Prompts. *CVPR* (2024).
- [21] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. 2018. Progressive growing of gans for improved quality, stability, and variation. *ICLR* (2018).
- [22] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- [23] Nikos Kolotouros, Thiemo Alldieck, Andrei Zanfir, Eduard Gabriel Bazavan, Mihai Fieraru, and Cristian Sminchisescu. 2023. DreamHuman: Animatable 3D Avatars from Text. *NeurIPS* (2023).
- [24] Verica Lazova, Eldar Insafutdinov, and Gerard Pons-Moll. 2019. 360-degree textures of people in clothing from a single image. In *2019 International Conference on 3D Vision (3DV)*. IEEE, 643–653.
- [25] 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).
- [26] Peike Li, Yunqiu Xu, Yunchao Wei, and Yi Yang. 2020. Self-Correction for Human Parsing. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2020). <https://doi.org/10.1109/TPAMI.2020.3048039>
- [27] Weiyu Li, Rui Chen, Xuelin Chen, and Ping Tan. 2024. SweetDreamer: Aligning Geometric Priors in 2D Diffusion for Consistent Text-to-3D. *ICLR* (2024).
- [28] Xian Liu, Xiaohang Zhan, Jiaxiang Tang, Ying Shan, Gang Zeng, Dahua Lin, Xihui Liu, and Ziwei Liu. 2024. Humanguassian: Text-driven 3d human generation with gaussian splatting. *CVPR* (2024).
- [29] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. 2015. SMPL: A Skinned Multi-Person Linear Model. *ACM Trans. Graphics (Proc. SIGGRAPH Asia)* 34, 6 (Oct. 2015), 248:1–248:16.
- [30] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. 2020. NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis. In *ECCV*.
- [31] Nasir Mohammad Khalid, Tianhao Xie, Eugene Belilovsky, and Tiberiu Popa. 2022. Clip-mesh: Generating textured meshes from text using pretrained image-text models. In *SIGGRAPH Asia 2022 conference papers*. 1–8.
- [32] Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. 2022. DreamFusion: Text-to-3D using 2D Diffusion. *ICLR 2023* (2022).
- [33] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*. PMLR, 8748–8763.
- [34] 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* 1, 2 (2022), 3.
- [35] Joshua Ramette, Josiah Sinclair, Zachary Vendeiro, Alyssa Rudelis, Marko Cetina, and Vladan Vuletić. 2022. Any-to-any connected cavity-mediated architecture

- for quantum computing with trapped ions or rydberg arrays. *PRX Quantum* 3, 1 (2022), 010344.
- [36] 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.
- [37] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. 2023. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 22500–22510.
- [38] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. 2022. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems* 35 (2022), 36479–36494.
- [39] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. 2016. Improved techniques for training gans. *Advances in neural information processing systems* 29 (2016).
- [40] Aditya Sanghi, Hang Chu, Joseph G Lambourne, Ye Wang, Chin-Yi Cheng, Marco Fumero, and Kamal Rahimi Malekshah. 2022. Clip-forge: Towards zero-shot text-to-shape generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 18603–18613.
- [41] Sefik Ilkin Serengil and Alper Ozpinar. 2020. LightFace: A Hybrid Deep Face Recognition Framework. In *2020 Innovations in Intelligent Systems and Applications Conference (ASYU)*. IEEE, 23–27. <https://doi.org/10.1109/ASYU50717.2020.9259802>
- [42] Yichun Shi, Peng Wang, Jianglong Ye, Mai Long, Kejie Li, and Xiao Yang. 2023. Mvdream: Multi-view diffusion for 3d generation. *arXiv preprint arXiv:2308.16512* (2023).
- [43] Jiaming Song, Chenlin Meng, and Stefano Ermon. 2020. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502* (2020).
- [44] Patrick Von Platen, Suraj Patil, Anton Lozhkov, Pedro Cuenca, Nathan Lambert, Kashif Rasul, Mishig Davaadorj, and Thomas Wolf. 2022. Diffusers: State-of-the-art diffusion models.
- [45] Jian Wang, Yunshan Zhong, Yachun Li, Chi Zhang, and Yichen Wei. 2019. Re-identification supervised texture generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 11846–11856.
- [46] Qixun Wang, Xu Bai, Haofan Wang, Zekui Qin, and Anthony Chen. 2024. InstantID: Zero-shot Identity-Preserving Generation in Seconds. *arXiv preprint arXiv:2401.07519* (2024).
- [47] Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. 2023. ProlificDreamer: High-Fidelity and Diverse Text-to-3D Generation with Variational Score Distillation. *NeurIPS* (2023).
- [48] Yuxiang Wei, Yabo Zhang, Zhilong Ji, Jinfeng Bai, Lei Zhang, and Wangmeng Zuo. 2023. ELITE: Encoding Visual Concepts into Textual Embeddings for Customized Text-to-Image Generation. *ICCV* (2023).
- [49] Yi Wu, Ziqiang Li, Heliang Zheng, Chaoyue Wang, and Bin Li. 2024. InfiniteID: Identity-preserved Personalization via ID-semantics Decoupling Paradigm. <https://api.semanticscholar.org/CorpusID:268531420>
- [50] Jiale Xu, Xintao Wang, Weihao Cheng, Yan-Pei Cao, Ying Shan, Xiaoju Qie, and Shenghua Gao. 2023. Dream3d: Zero-shot text-to-3d synthesis using 3d shape prior and text-to-image diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 20908–20918.
- [51] Xingqian Xu, Jiayi Guo, Zhangyang Wang, Gao Huang, Irfan Essa, and Humphrey Shi. 2023. Prompt-Free Diffusion: Taking "Text" out of Text-to-Image Diffusion Models. *arXiv preprint arXiv:2305.16223* (2023).
- [52] Xiangyu Xu and Chen Change Loy. 2021. 3D human texture estimation from a single image with transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*. 13849–13858.
- [53] Xingqian Xu, Zhangyang Wang, Gong Zhang, Kai Wang, and Humphrey Shi. 2023. Versatile diffusion: Text, images and variations all in one diffusion model. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 7754–7765.
- [54] Haibo Yang, Yang Chen, Yingwei Pan, Ting Yao, Zhineng Chen, and Tao Mei. 2023. 3dstyle-diffusion: Pursuing fine-grained text-driven 3d stylization with 2d diffusion models. In *Proceedings of the 31st ACM International Conference on Multimedia*. 6860–6868.
- [55] 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.
- [56] Hu Ye, Jun Zhang, Sibio Liu, Xiao Han, and Wei Yang. 2023. IP-Adapter: Text Compatible Image Prompt Adapter for Text-to-Image Diffusion Models. (2023).
- [57] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. 2023. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 3836–3847.
- [58] Yuxuan Zhang, Yiren Song, Jiaming Liu, Rui Wang, Jinpeng Yu, Hao Tang, Huaxia Li, Xu Tang, Yao Hu, Han Pan, and Zhongliang Jing. 2024. SSR-Encoder: Encoding Selective Subject Representation for Subject-Driven Generation.
- [59] Fang Zhao, Shengcai Liao, Kaihao Zhang, and Ling Shao. 2020. Human parsing based texture transfer from single image to 3D human via cross-view consistency. *Advances in Neural Information Processing Systems* 33 (2020), 14326–14337.
- [60] Shanshan Zhong, Zhongzhan Huang, Weushao Wen, Jinghui Qin, and Liang Lin. 2023. Sur-adapter: Enhancing text-to-image pre-trained diffusion models with large language models. In *Proceedings of the 31st ACM International Conference on Multimedia*. 567–578.