

# DINAR: Diffusion Inpainting of Neural Textures for One-Shot Human Avatars

David Svitov  
Samsung AI Center  
d.svitov@samsung.com

Dmitrii Gudkov  
Samsung AI Center  
d.gudkov@samsung.com

Renat Bashirov  
Samsung AI Center  
r.bashirov@samsung.com

Victor Lempitsky  
Cinemersive Labs  
victor@cinemersivelabs.com

## Abstract

*We present DINAR, an approach for creating realistic rigged fullbody avatars from single RGB images. Similarly to previous works, our method uses neural textures combined with the SMPL-X body model to achieve photo-realistic quality of avatars while keeping them easy to animate and fast to infer. To restore the texture, we use a latent diffusion model and show how such model can be trained in the neural texture space. The use of the diffusion model allows us to realistically reconstruct large unseen regions such as the back of a person given the frontal view. The models in our pipeline are trained using 2D images and videos only. In the experiments, our approach achieves state-of-the-art rendering quality and good generalization to new poses and viewpoints. In particular, the approach improves state-of-the-art on the SnapshotPeople public benchmark.*

## 1. Introduction

The use of fullbody avatars in the virtual and augmented reality applications [34] is one of the drivers behind the recent surge of interest in fullbody photorealistic avatars [42, 4, 40]. Apart from the realism and fidelity of avatars, the ease of acquisition of new personalized avatars is of paramount importance. Towards this end, several works propose methods to restore 3D textured model of a human from a single image [48, 49, 19, 4] but such models require additional efforts to produce rigging for animation. The use of additional rigging methods significantly complicates the process of obtaining an avatar and often restricts the poses that can be handled. At the same time, some of the recent methods use textured parametric models of human body [56, 28] while applying inpainting in the texture space. Current texture-based methods, however, lack photo-realism and rendering quality.

An alternative to using classical RGB textures directly is to use deferred neural rendering [53]. Such approaches make it possible to create human avatars controlled by the parametric model [29, 40]. The resulting avatars are more photo-realistic and easier to animate. However, existing approaches require a video sequence to create an avatar [41]. The StylePeople system [16], which is also based on deferred neural rendering and parametric model, provides an opportunity to create avatars from single images, however the quality of rendering for unobserved body parts is low.

We propose a new method to create photo-realistic animatable human avatars from a single photo. To make the avatars animatable we leverage neural texture approach [53] along with the SMPL-X parametric body model [40]. We propose a new architecture for generating the neural textures, in which the texture comprises both the RGB part explicitly extracted from the input photograph by warping and additional neural channels obtained by mapping the image to a latent vector space and decoding the result into the texture space. As is customary with neural rendering [53, 16, 41], the texture generation is trained in an end-to-end fashion with the rendering network.

To restore the neural texture for unobserved parts of the human body we develop a diffusion model [22]. This approach allows us to obtain photo-realistic human avatars from single images. In the presence of multiple images, we can merge neural textures corresponding to different images while restoring parts that are still missing by diffusion-based inpainting. The use of diffusion for inpainting distinguishes our approach from several previous works [28, 20, 17] including StylePeople [16] that rely on generative adversarial framework [15] to perform inpainting of human body textures. As in other image domains [10, 47], we found that the use of diffusion alleviates problems with mode collapse and allows to obtain plausible samples from complex multi-modal distributions. To the best of our knowledge, we are the first to extend the diffusion models to the task of generating human body textures.



Figure 1: **Animations of one-shot avatars.** We generate avatars of previously unseen people from single images and animate the avatars by changing their SMPL-X poses. Our model produces plausible results for new poses and views, even for complex and loose garments like skirts and can handle complex poses (such as hands near the body) gracefully.

To sum up, our contributions are as follows:

- We propose a new approach for modeling human avatars based on neural textures that combine the RGB and the latent components.
- We adapt the diffusion framework for neural textures and demonstrate that it is capable of inpainting such textures.
- We demonstrate the ability of our system to build realistic animatable avatars from a single photograph.

The proposed approach allows us to obtain a photorealistic animatable person’s avatar from a single image. Specifically, when the input photograph is taken from the front, we observe that the person’s back is restored in a consistent photo-realistic manner. We demonstrate the effectiveness and accuracy of our approach on real-world images from SnapshotPeople [3] public benchmark and images of people in natural poses.

## 2. Related work

**Full body human avatar.** Many avatar systems are based on parametric human body models, of which the most popular are the SMPL [29] body model as well as the SMPL-X model [40] which augments SMPL with facial expressions and hand articulations. Such models represent human body without garments or hair. Approaches based on deferred neural rendering (DNR) [53] or neural radiance fields (NeRF) [33] can be used to add clothing and perform photo-realistic rendering. DNR uses a multi-channel trainable neural texture and a rendering convolutional network

to render the resulting avatars in a realistic way [41, 16]. It makes it easier to animate the resulting avatar. NeRF uses sampling along a ray in implicit space for rendering and allows one to extract accurate and consistent geometry [48, 23, 4].

**Video full body avatar.** The parametric SMPL model is used to obtain a photo-realistic human avatar from a monocular video [2, 3, 54, 41]. The use of a video enables the creation of high-quality avatars due to the large amount of input information. However, the applicability of such approaches is limited by the difficulty of obtaining such videos.

A number of techniques [41, 16, 27] use neural textures to simulate out-of-mesh details in order to create avatars from monocular video. For a more temporally consistent representation of loose clothing, some methods [60, 2, 3] use vertex displacements in addition to the texture. The highest photorealism and temporal consistency in the video avatar task are currently achieved by NeRF-based methods [38, 54, 24]. Such methods use SMPL to translate an avatar to the implicit canonical space. Then render it by integrating along rays from a camera.

**One-shot full body avatar.** One-shot avatar approaches reconstruct human body avatars from a single image. Early works achieved this by inpainting partial RGB textures [56, 28]. These approaches did not allow realistic modeling of avatars with clothing. More recent works on one-shot body modeling relied on implicit geometry and radiance models, which predict occupancy and color with the multi-layer perceptron conditioned on feature vectors extracted from the input image [48, 4, 19]. While this line of work often man-

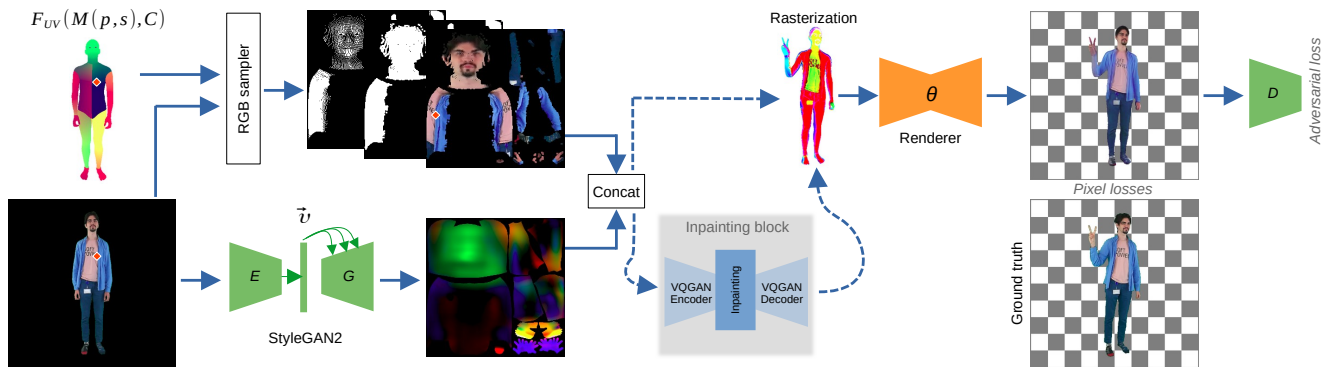


Figure 2: **Overview of our neural texture avatar pipeline:** Given an RGB image and the corresponding SMPL-X model as input, we use a UV-map to sample the RGB texture. Using an encoder and a StyleGAN2 generator, we convert the input image into a neural texture. We rasterize SMPL-X with concatenation of the neural and RGB textures and translate it to the RGB+mask output using the rendering network. The inpainting block replaces missing texture fragments and is described in detail in Section 3.3.

ages to recover intricate geometric details, the realism of the recovered texture in the unobserved parts is usually limited.

The ARCH system [23] uses the rigged canonical space to build avatars suitable for animations. ARCH++ [20] improves the quality of resulting avatars by revisiting the major steps of ARCH. They also solve the challenge of unseen surfaces by synthesizing the back views of a person from the front view. PHORHUM [4] improves the modeling of geometry by adding the reasoning about scene illumination and albedo of the surface. S3F [9] also enables changing the avatar lighting. But it takes advantage of 3D Features connected to the vertices of the parametric model, which allows one to animate the resulting avatar. ICON [55] uses local features to avoid the dependence of the reconstructed geometry on the global pose. Their approach first estimates separate models from each view and then merges the models using SCANimate [50]. The method uses RGB textures applied to reconstructed geometry, which limits the rendering photo-realism.

An alternative approach to getting one-shot avatars is to use avatar generative models as proposed in StylePeople [16]. The authors circumvent the need to reconstruct the unseen parts by exploiting the entanglement in the GANs latent space. Unfortunately, the imperfection of their underlying generative model often leads to implausible appearance of unobserved parts.

**Diffusion models.** Diffusion models [51] are probabilistic models for learning the distribution of  $p(x)$  by gradual denoising of a normally distributed variable. Such denoising corresponds to learning the inverse process for a fixed Markov chain of length  $T$ . The most successful models for image generation [10, 22, 43] use the reweighted variant of the variational lower bound on  $p(x)$ . These models can also be interpreted as denoising autoencoders

$\epsilon_\omega(x_t, t); t = 1, \dots, T$  with shared weights. These autoencoders learn to predict  $x_{t-1}$  with reduced noise level over  $x_t$ . In [22] was demonstrated that such denoising models can be trained with a simplified loss function:

$$L_{DM} = \mathbb{E}_{x, \epsilon \sim \mathcal{N}(0,1), t} [||\epsilon - \epsilon_\omega(x_t, t)||_2^2], \quad (1)$$

where  $t$  is uniformly sampled from  $\{1, \dots, T\}$ . Our work builds on recent advances in inpainting with diffusion models [46, 31, 44]. In particular, we use latent diffusion [43] in our approach, which has been shown to be effective in inpainting of RGB images.

### 3. Method

Our approach has two components: the avatar generation model and the inpainting model. The scheme of our avatar generation model is presented in Fig. 2. The model reconstructs the neural texture from the input image using two pathways and then uses the texturing operation as well as the neural rendering to synthesize images of the avatar. Then we train the inpainting model based on the denoising diffusion probabilistic model (DDPM) [22] on top of the pretrained avatar generation model. This model is used to restore body areas unrepresented in the input image.

#### 3.1. Avatar generation model

Our approach creates a 3D rigged avatars of clothed human using several neural networks and the texturing modules trained together in the end-to-end fashion. The overall architecture takes as input an RGB image and a parametric SMPL-X [37] body model. During training, we fit SMPL-X models to the train images using an implementation of SMPLifyX [37] approach with an additional segmentation loss that allows us to match human silhouettes better.

In more detail, we use an SMPL-X fixed-topology mesh  $M(p, s)$ , driven by sets of pose parameters  $p$  and shape parameters  $s$ . Here the pose parameters  $p$  are used to rig the mesh. Information about person appearance is stored in the  $L$ -channeled neural texture  $T$ . One-to-one correspondence between the mesh  $M$  and the texture  $T$  is set as a UV-map. For SMPL-X, we employ a customized UV-map with a front cut to avoid difficult to inpaint back-view seams. For the mapping of texture  $T$  to the output image we use the rasterizer  $R(T, M, C_{\text{target}})$  depending on the mesh  $M$  and target camera position  $C_{\text{target}}$ . As result, the rasterizer  $R$  produces an image of size  $H \times W \times L$ . Where  $H$  and  $W$  determine the output image size and  $L$  is the texture features size.

We set parameters of the rasterizer  $R$  so that  $H$  and  $W$  correspond to the height and width of the input RGB image  $I_{\text{rgb}}$ . In this case, we can use UV not only to map feature vectors from the neural texture  $T$ , but also to sample color values from the input image  $I_{\text{rgb}}$  into the texture space:  $T_{\text{rgb}}$ . We transfer the color value from  $I_{\text{rgb}}$  to the texture  $T_{\text{rgb}}$  point specified by the UV-map. This RGB texture allows us to explicitly save information about high-frequency details and original colors, which are hard to preserve when mapping the whole image to a vector of limited dimensionality (as discussed below). We additionally apply inpainting of small gaps with averaging neighbor pixels to fill the gaps in  $T_{\text{rgb}}$ . We also save the binary map of sampled pixels to the  $B_{\text{smp}}$  and the map of the sampled and inpainted pixels to the  $B_{\text{fill}}$ .

The main part of the neural texture is  $T_{\text{gen}}$ . It has the number of channels  $L = 16$  and is generated using the encoder-generator architecture  $T_{\text{gen}} = G(E(I_{\text{rgb}}))$ . The encoder  $E$  is based on the StyleGAN2 [25] discriminator architecture and compresses the input image  $I_{\text{rgb}}$  to a vector of dimension 512. The generator  $G$  has the architecture of the StyleGAN2 generator and converts the vector into a  $T_{\text{gen}}$  neural texture with the number of channels  $L = 16$  as in StylePeople [16].

The final neural texture used in our method has a dimension of  $256 \times 256 \times 21$  and consists of the concatenation of: the generated texture  $T_{\text{gen}}$  ( $256 \times 256 \times 16$ ), the texture  $T_{\text{rgb}}$  sampled from the RGB image ( $256 \times 256 \times 3$ ) and the two binary segmentation maps ( $B_{\text{smp}}$  and  $B_{\text{fill}}$ ):

$$T = T_{\text{gen}} \oplus T_{\text{rgb}} \oplus B_{\text{smp}} \oplus B_{\text{fill}}. \quad (2)$$

We note that such an approach with the explicit use of RGB channels as part of the neural texture was originally used in [53].

We use the neural renderer  $\theta(R(T, M, C_{\text{target}}))$  to translate the rasterized image with  $L + 3 + 1 + 1$  channels into  $I_{\text{rend}}$  output RGB image. Neural renderer  $\theta$  has a U-Net [45] architecture with ResNet [18] blocks. We train the renderer  $\theta$  jointly with the neural texture encoder  $E$  and generator  $G$ . Thus, rendering an avatar with a texture  $T$  (Eq. 2) in a new

pose  $p_{\text{target}}$  has the following form:

$$I_{\text{rend}} = \theta(R(T, M(p_{\text{target}}, s), C_{\text{target}})) \quad (3)$$

During training, we minimize the following losses: the difference loss  $L_2$  and the perceptual loss  $LPIPS$  [58] between the rendered  $I_{\text{rend}}$  and ground truth  $I_{\text{GT}}$  images; the Dice loss [52] between the ground truth and the predicted segmentation masks. We calculate  $L_2$  and  $LPIPS$  for the entire image and additionally with a weight of 0.1 for the area with a face, since it has been demonstrated in [13] that the face is very important for human perception. Additionally, we use an adversarial loss to make  $I_{\text{rend}}$  look more realistic and sharp. We use nonsaturating adversarial loss  $Adv$  [15] with the StyleGAN2 discriminator  $D$  with R1-regularization [32]. The overall loss thus has the following form:

$$Loss = \lambda_1 \cdot L_2 + \lambda_2 \cdot LPIPS + \lambda_3 \cdot Dice + \lambda_4 \cdot Adv + \lambda_5 \cdot R1_{\text{reg}}. \quad (4)$$

We describe the choice of hyperparameters  $\lambda_{1..5}$  in Section 4.

Our training is performed on multi-view image sets such as sets of frames of the video (or such as sets of renders of 3D models). The training data and network training details are discussed in Section 4. During training, in each case we take two different frames from the same set, one as the input image  $I_{\text{rgb}}$  and the other as the target image  $I_{\text{GT}}$  (the two thus having different camera parameters as well as two different body pose parameters). This is essential for the pipeline to generalize to new camera positions  $C$  and poses  $p$ . To accomplish that, the renderer and the texture generation modules learn to inpaint small fragments unseen in  $I_{\text{rgb}}$ .

### 3.2. Texture merging

While the texture generator and the renderer learn to compensate for the small amount of the unseen surfaces that may be present in the target view, we have found that such ability is limited to small changes in body pose and/or camera parameters.

The easiest way to obtain an avatar that can be rendered from arbitrary angles is to create it from several images by merging the corresponding neural textures. For that we can use a simple blending scheme (Fig. 3). Specifically, assume the  $N$  input images  $I_{\text{rgb}}^i$  from different view points are given that result in  $N$  neural textures  $T^i$ . These textures naturally cover different areas of the body visible in different input images. To combine these textures, we define the function  $F(T^1 \dots T^N, \lambda^1 \dots \lambda^N)$  shown in Fig. 3. The  $\lambda^i$  is auxiliary information about the visible in  $I_{\text{rgb}}^i$  part of the texture.



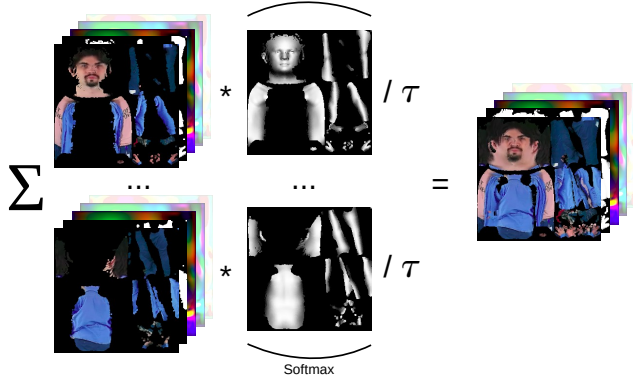


Figure 3: **Texture merging:** We get the full neural texture as the weighted sum of textures from different view points. As weights, we use the angles between the normal vectors and the direction of the camera. For simplicity, only the RGB channels are shown in the figure, but the merging affects all channels.

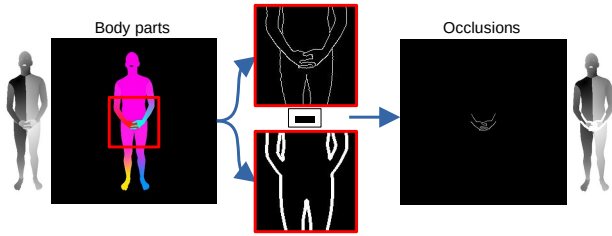


Figure 4: **Oclusions detection.** We use rasterization with the colored body parts texture to detect areas occluded by limbs. The detected areas are masked out of the final UV-render to reduce inaccuracies in RGB texture sampling.

As auxiliary information  $\lambda^i$  we use the angles between the normal vectors of corresponding point of the mesh  $M^i$  and direction vector of the camera. Thus  $\lambda^i$  defines how frontal each texture point is to the camera. We perform texture merging utilizing this information, emphasizing the more frontal pixels for each  $T^i$ . We aggregate textures  $T^i$  using weighted average with weights calculated as  $\vec{w} = softmax(\frac{\lambda}{\tau})$ . The  $\tau$  factor controls the sharpness of edges at the junction of merged textures. The final texture is calculated as:

$$T = F(T^1 \dots T^N, \lambda^1 \dots \lambda^N) = \sum_{i=1}^N T^i \cdot w^i. \quad (5)$$

This technique allows us to get few-shot avatars by merging one-shot avatars for different views. More sophisticated blending schemes such as pyramid blending [36] or Poisson blending [39] can be also used.

### 3.3. The inpainting model

As the final piece of our approach, we train a network that can create complete neural textures from a single photographs. We do that using supervised learning, where we use the incomplete textures based on single photographs as inputs and the combined textures aggregated from multiple views as ground truth.

Since the distribution of plausible (“correct”) complete textures given the input partial texture is usually highly complex and multi-modal, we use the denoising diffusion probabilistic model (DDPM) framework [22] and train a denoising network instead of the direct mapping from the input to the output.

**Reducing texture space.** As described above, in our experiments the neural texture  $T$  has a resolution of  $256 \times 256 \times 21$ . This leads to a huge memory requirements during the diffusion model training. In order to reduce memory consumption and improve the network convergence, we first reduce the neural texture size using VQGAN autoencoder [12] analogously to the reduction of an RGB image size performed in [43]. As demonstrated in Fig. 2 VQGAN is added as an alternative branch for the input of the renderer  $\theta$ . After the pretraining of VQGAN, the pipeline is finetuned end-to-end in order to adapt the renderer to VQGAN decompression artifacts in the neural texture  $T_{res}$ .

We use several loss functions to train the VQGAN autoencoder to more accurately restore neural textures. To improve the visual quality of the avatar after texture decompression, we use the loss functions in the RGB space. We render the avatar as described in (3) with the restored texture  $T_{res}$  and optimize the loss function (4). Also, we use an additional L2 loss in texture space  $\|T - T_{res}\|_2^2$  for additional regularization and preservation of neural texture properties for all views.

**Texture inpainting.** After adding the VQGAN branch to the pipeline, we train the DDPM model in its latent space. Thus, the diffusion model is applied to  $T_c^j = E_{VQ}(T^j)$  with size  $64 \times 64 \times 3$ , obtained after the compression of the single-view based texture  $T^j$  by the VQGAN encoder  $E_{VQ}$ . Index  $j$  in our experiments corresponds to the front view. Following [43] we train DDPM using a U-Net architecture with attention. We condition the denoising model with  $T_c^j \oplus b(B_{fill}^j)$ , where  $b$  is a bilinear resize to the spatial size of  $T_c^j$ . As the input to the model, we feed in the concatenation of the condition and the compressed merged texture  $T_c$  corrupted with normally-distributed noise corresponding to the diffusion timestep  $t$ . The U-Net architecture thus trains to denoise the input  $T_c$  by minimizing the loss (1). This allows  $T_c$  to be iteratively derived from pure Gaussian noise during conditional inference.

As mentioned above, we train diffusion inpainting (Fig. 5) using the merged textures (section 3.2) as ground truth. Here, we generate textures from a dataset of 3D human

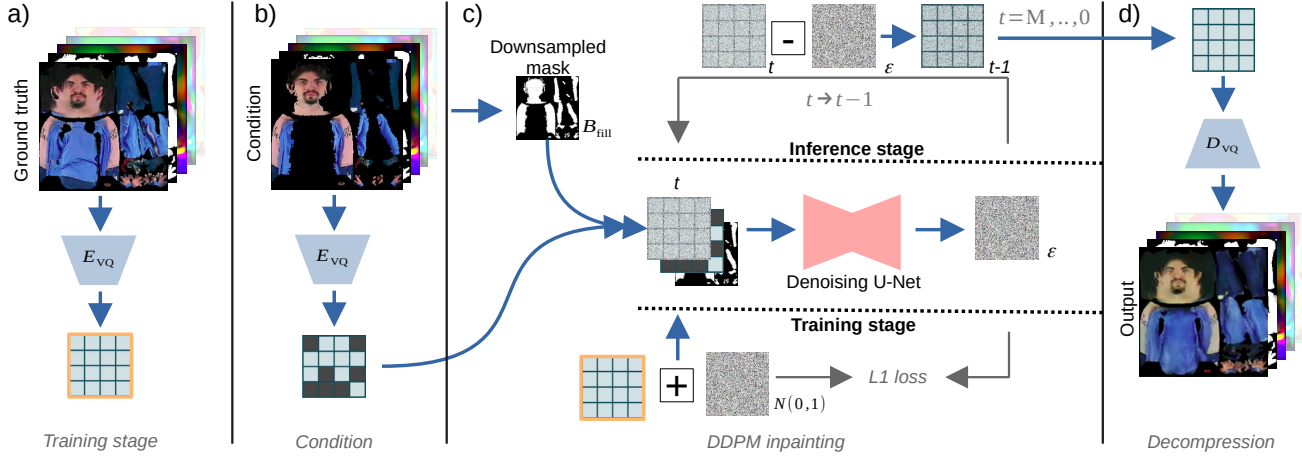


Figure 5: **Texture inpainting:** a) In the training step, we reduce the size of the ground truth and partial (b) neural textures using the VQGAN encoder. c) At train U-Net learns to remove noise from Ground truth texture for a given step  $t$ . At inference, we iteratively apply  $M$  steps of denoising to sample the image. d) The sampled texture is transformed to the original size by the VQGAN decoder.

scans. Multi-view dataset of people photographs with good angle coverage could be used as training data as well, and we chose to use the scan render dataset solely because of its availability to us. The training dataset is discussed in section 4.

### 3.4. Inference

We perform inference in two stages. To get an animated avatar from a photo, we first get a partial texture from the input image, and then we inpaint it with the DDPM model.

To create the partial texture  $T^j$  we inference the model shown in Figure 2 with input image  $I_{rgb}^j$ . We also apply several techniques to improve avatar quality in the inference stage. To enhance the texture details in the visible part, we perform a few optimization steps of RGB channels with gradients from the differentiable renderer for the input image. To reduce the impact of SMPL-X fitting imperfections, we detect areas of human self-occlusion in the input image (Section 3.5) and do not sample the texture along the overlap outline. The restoration of these areas are left for the inpainting process described below.

During the inpainting stage the DDPM restores the whole texture  $T_c$  from the noise conditioned on the partial texture. To do this, we perform  $M$  denoising steps starting with pure noise as shown in Figure 5. We then transform  $T_c$  into the restored full-size texture  $T_{res}$  using the VQGAN decoder  $D_{VQ}$ . To retain all the details from the input view, we then merge the restored texture with the input texture  $T^j$  using the  $B_{fill}^j$  mask:

$$T = T_{res} \cdot (1 - B_{fill}^j) + T^j \cdot B_{fill}^j \quad (6)$$

The resulting texture has all the details visible in the  $I_{rgb}^j$  image, while the parts of a person invisible on  $I_{rgb}^j$  are restored

by the DDPM in the VQGAN latent space. An avatar with a texture  $T$  obtained in this way can be rendered from new view points and in new poses.

### 3.5. Addressing proxy-geometry imperfections

We have found that due to imperfect SMPL-X meshes, pixels are wrongly sampled from one body part to another at the self-occluded areas (e.g. hands in front of the body). This results in implausible renderings (Fig. 7 (b, c)).

To address this issue, we do not sample the RGB texture along the outline of such overlapping body parts. We employ rasterization with a colormap as a texture to find overlapping regions (Fig. 4). We assigned each limb in the colormap a separate color and made the transition between them smooth using a color gradient. This enables us to avoid having seams in the rasterization. We detect edges in the colormap rasterization with Canny algorithm [8]. We then determine a person’s contours employing binarized SMPL-X rasterization. By taking the contours out of edges, we obtain an occlusion map. We use the resulting occlusion map to mask out areas in the UV-render. This enables us to rely on inpainting in later pipeline steps rather than sampling pixels in overlapped areas.

## 4. Experiments

### 4.1. Implementation details

Now we present implementation details and hyperparameters values. Our model is trained on RGB images at the  $512 \times 512$  resolution. First, we train the generator and the renderer with the next losses: the L2 loss, the LPIPS loss, the Dice loss [52] and the adversarial nonsaturating loss [15] with R1 regularization [32]. We use a weighted



Figure 6: **Qualitative comparison on the SnapshotPeople benchmark.** We compare our method with state-of-the-art approaches for creating avatars: PIFu, PHORHUM, ARCH, ARCH++, StylePeople. For each method, we show a front view and a back view to evaluate the quality of back reconstruction. Results for PHORHUM, ARCH and ARCH++ provided by the authors and registered with ground truth images.

sum of losses with the following weights. L2 loss with  $\lambda_1 = 2.2$ ; LPIPS loss with  $\lambda_2 = 1.0$ ; Dice loss with  $\lambda_3 = 1.0$ ; Adversarial loss with  $\lambda_4 = 0.01$ . We use the lazy regularization R1 as proposed in [25] every 16 iterations with the weight  $\lambda_5 = 0.1$ . To calculate LPIPS, we take a random  $256 \times 256$  crop of the  $512 \times 512$  images. We train the generator and the renderer end-to-end for 100,000 steps using the ADAM optimizer [26] with  $2e-3$  learning rate and the batch size of four.

Then we train VQGAN to compress neural textures to  $6 \times 64 \times 64$  tensors consisting of vectors of length six from a trainable dictionary with 8192 elements. We first train only the VQGAN branch for 300,000 steps. Then we finetune the pipeline end-to-end for additional 20,000 steps to reduce the renderer artifacts when processing neural textures after VQGAN. After that, we train the diffusion probabilistic model to restore missing parts of the texture. We use the U-Net architecture with the BigGAN [7] residual blocks for up- and downsampling and with attention layers on three levels of its feature hierarchy. To additionally prevent overfitting, we use dropout with a probability of 0.5 in the residual blocks. We train the diffusion model for 50,000 iterations with AdamW optimizer [30] with a batch size of 128 and a learning rate of  $1.0e-6$ .

## 4.2. Datasets

To train our pipeline we used only 2D images obtained by rendering Texel [1] dataset. We pretrained the neural texture generator and the neural renderer using 2D images of 13,000 people in diverse poses. We have noticed that diverse poses are crucial to train realistically animatable avatars. For each image, we obtained a segmentation mask using Graphonomy [14] segmentation and fit the SMPL-X parametric model using SMPLify-X [37]. We also used a segmentation Dice loss [52] to improve the body shape of fitted SMPL-X.

To train VGQAN and the diffusion model, we used renders from Texel dataset. We acquired 3333 human scans from the Texel dataset. They are people in different clothing with various body shapes, skin tones and genders. We rendered each scan from 8 different views to get a multi-view dataset. We also augmented the renders with camera angle changes and color shifts. Thus, for each person in the dataset, we got 72 images. Note that any images from different views are suitable for training the model, not necessarily obtained from 3D scans.

We qualitatively evaluate our avatars and their animations on AzurePeople [5] dataset (Fig. 1). This dataset contains

Method	Same view				Novel view		
	MS-SSIM $\uparrow$	PSNR $\uparrow$	LPIPS $\downarrow$	DISTS $\downarrow$	DISTS $\downarrow$	ReID $\downarrow$	KID $\downarrow$
PIFu	0.9793	26.2828	0.0404	0.0706	0.1839	0.09769	0.0907
Phorhum	0.9603	24.2112	0.0531	0.0948	0.1564	0.09149	0.0144
ARCH	0.9223	20.6499	0.0732	0.1432	0.2039	0.09575	0.0974
ARCH++	0.9526	22.5729	0.0540	0.0842	0.1750	0.09098	0.0589
StylePeople	0.9610	23.0584	0.0848	0.0852	0.1698	0.13090	0.0530
<b>DINAR (Ours)</b>	0.9687	24.4182	0.0504	0.0703	0.1407	0.07855	0.0133

Table 1: **Metrics comparison on the SnapshotPeople benchmark.** We compared our method not only with other parametric model based approaches (StylePeople), but also with approaches that restore geometry and require using additional methods for rigging (PIFu, Phorhum) or restore geometry in the canonical pose (ARCH, ARCH++).

Method	MS-SSIM $\uparrow$	PSNR $\uparrow$	LPIPS $\downarrow$
RGB texture	0,8825	22,8041	0,1204
Neural texture	0,8632	23,0935	0,1285
Both textures	0,9126	24,1041	0,1005
Both + Inpainting	0,9199	24,8510	0,0904
Both + Inpainting + Occlusion detector	0,9201	24,8619	0,0905
Both + Inpainting + Ocl. + RGB tuning	0,9203	24,9758	0,0901

Table 2: **Quantitative ablation study.** Metrics were measured for images of 15 people with diverse poses and view points.

diverse dressed people standing in natural poses. We also quantitatively evaluate our pipeline on the SnapshotPeople [3] public benchmark. It contains 24 videos of people rotating in A-pose. We select frames with front and back views from each video to measure the accuracy of front and back reconstruction respectively. For each image we get a segmentation mask and SMPL-X fit as described above.

### 4.3. Quantitative results

To numerically evaluate our avatars, we report several metrics (Table 1). For clarity, the table is divided into three sections: PIFu [48], PHORHUM [4] - require additional rigging; ARCH [23], ARCH++ [20] - restore geometry in the canonical space; StylePeople [16], DINAR (ours) - utilize a parametric human model and therefore can be animated without additional procedures. We measured reference-based metrics in the front view avatars for the SnapshotPeople public benchmark, namely: Multi-scale structural similarity (MS-SSIM  $\uparrow$ ), Peak signal-to-noise ratio (PSNR  $\uparrow$ ), Learned Perceptual Image Patch Similarity (LPIPS  $\downarrow$ ) [59]. We found that our method works on par with non-rigged methods.

To evaluate the quality of the back view (and thus the ability to generalize to new views), we report Kernel Inception distance (KID  $\downarrow$ ) [6] measurements. This metric allows one to evaluate generated images quality and is more

suitable for small amounts of data than FID [21]. Our approach results in the highest KID value compared to other methods. To assess identity preservation we measured re-identification score (ReID  $\downarrow$ ) based on FlipReID [35] model for human re-identification. Our method shows the best results in person’s identity preserving between front and back views. To additionally validate the quality of the textures and to measure structural similarity in cases with unaligned ground truth images we measured Deep Image Structure and Texture Similarity (DISTS  $\downarrow$ ) [11]. We provide measurements for front view and back view in the table. Our method produce the most naturally looking avatars from both views.

### 4.4. Qualitative Results

Metric value does not always correlate well with human perception. In Fig 6, we show the qualitative results of our method in comparison with other one-shot avatars methods: PIFu, PHORHUM, ARCH, ARCH++ and StylePeople. The figure shows both front and back views of the avatars. Additional comparisons are shown in Supplementary materials, including results on the THuman2.0 dataset [57].

Overall, our method realistically reconstructs the texture of clothing fabrics on the back (*e.g.* pleats on pants), which boosts the realism of the renders. Using the whole information from the given image allows us not to copy unwanted patterns from front to back (as is commonly done by the pixel-aligned methods while recovering the texture for the back). Using sampled RGB texture as an addition to a neural texture allows us to achieve photo-realistic facial details and preserve high frequency details. We note that PIFu accurately reproduces the color of the avatar and restores the geometry well. However, it does not preserve high-frequency details, which is why avatars suffer from a lack of photo-realism. PHORHUM generates highly photo-realistic avatars but often suffers from color shifts. Another methodological shortcoming of this approach is the absence of a human body prior. Therefore, the model can be over-fitted on training dataset’s human poses, which may lead to incorrect work with unseen poses. Avatars generated by



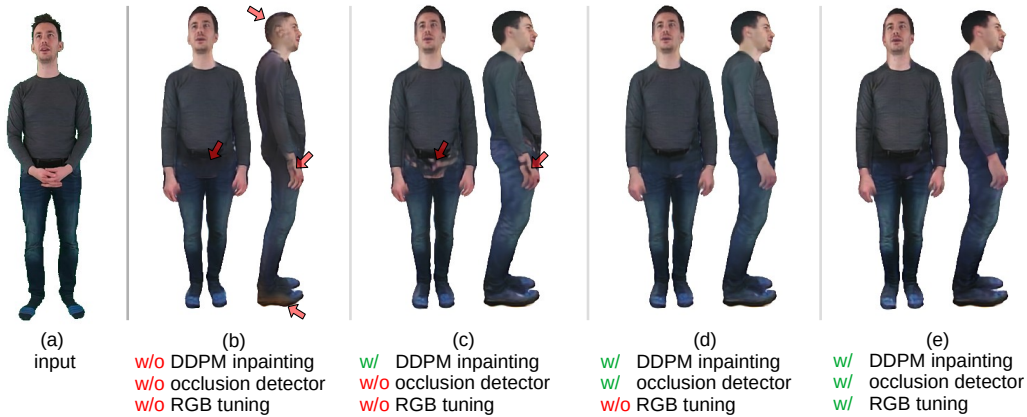


Figure 7: **Ablation study.** b) Avatars without DDPM inpainting are prone to artifacts and lack of realism. c) Disabling the occlusions detector leads to artifacts due to sampling errors at the edge of the overlapped areas. d) Disabling RGB finetuning results in less high-frequency detail and worse color matching. e) The best result is achieved by using all the steps of the pipeline.

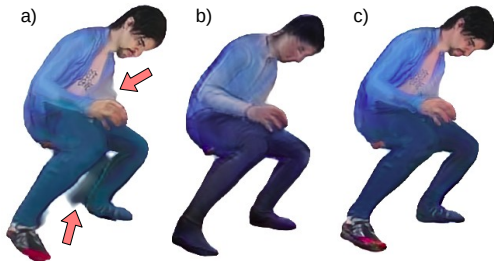


Figure 8: **Texture ablation study.** a) RGB texture; b) Neural texture; c) Both textures.

ARCH contain strong color artifacts and suffer from geometry restoration errors. ARCH++ significantly improves geometry and color quality for the frontal view, but the back view still suffers from color shift and artifacts. StylePeople is based on a parametric human model and can be easily animated without the use of third party methods or additional rigging. However, the coverage of the latent space of their model is limited, which leads to overfitting and poor generalization to unseen people, when performing inference based on a single view.

## 5. Ablation study

We quantitatively (Table 2) and qualitatively (Fig. 7, 8) ablate texture choice and the proposed method steps. For this, we used 15 images from AzurePeople [5] with diverse poses and view points. Eight people’s input images from the front view point, and seven from the back view point.

We have studied the efficiency of using RGB and neural channels in the texture for the one-shot avatar task. We found that using both kinds of channels results in better met-

rics than using them separately. We also studied the impact of using a DDPM inpainting model (Section 3.3) on the final metrics. In scenarios without a model for *Inpainting*, the renderer takes over this function. The *Occlusion detection* step is used to remove RGB texture sampling errors caused by SMPL-X inaccuracy and is described in section 3.5. The final step in the ablation study is *RGB channels fine-tuning* on the input image for a few steps (Section 3.4).

The best metrics are achieved using all the steps of the proposed approach. We also provide a visual comparison of the various steps of the ablation study in the Figure 7.

Figure 8 (a) illustrates how the renderer is unable to obtain information about the out-of-mesh details while using only the RGB texture. This leads to artifacts around the avatar in new poses. In turn, the use of only neural channels leads to the loss of high-frequency information (Fig. 8 (b)). Using both types of textures allows us to pass information about out-of-mesh elements to the renderer and preserve high-frequency details (Fig. 8 (c)).

## 6. Conclusion

We have presented a new approach for modeling human avatars based on neural textures that combines the RGB and the latent components. The RGB components are used to preserve the high frequency details, while the neural components add hair and clothing to the base SMPL-X mesh. Using the parametric SMPL-X model as a basis makes it easy to animate the resulting avatar. Our method restores missing texture parts using an adapted diffusion framework for inpainting such textures. Our method thus creates rigged avatars, while also improving the rendering quality of the unseen body parts when compared to modern non-rigged human model reconstruction methods.

## References

- [1] Texel 3d body model dataset. <https://texel.graphics/texel-3d-body-model-dataset/>. Accessed: 09-11-2022. **7**
- [2] Thiemo Alldieck, Marcus Magnor, Weipeng Xu, Christian Theobalt, and Gerard Pons-Moll. Detailed human avatars from monocular video. In *International Conference on 3D Vision (3DV)*, pages 98–109. IEEE, 2018. **2**
- [3] Thiemo Alldieck, Marcus Magnor, Weipeng Xu, Christian Theobalt, and Gerard Pons-Moll. Video based reconstruction of 3d people models. In *CVPR*, pages 8387–8397, Jun 2018. **2, 8**
- [4] Thiemo Alldieck, Mihai Zanfir, and Cristian Sminchisescu. Photorealistic monocular 3d reconstruction of humans wearing clothing. In *CVPR*, pages 1506–1515, 2022. **1, 2, 3, 8**
- [5] Renat Bashirov, Anastasia Ianina, Karim Iskakov, Yevgeniy Kononenko, Valeriya Strizhkova, Victor Lempitsky, and Alexander Vakhitov. Real-time rgbd-based extended body pose estimation. In *WACV*, pages 2807–2816, 2021. **7, 9**
- [6] Mikołaj Bińkowski, Dougal J. Sutherland, Michael Arbel, and Arthur Gretton. Demystifying MMD GANs. In *ICLR*, 2018. **8**
- [7] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale GAN training for high fidelity natural image synthesis. In *ICLR*, 2019. **7**
- [8] John F. Canny. A computational approach to edge detection. *IEEE TPAMI*, PAMI-8:679–698, 1986. **6**
- [9] Enric Corona, Mihai Zanfir, Thiemo Alldieck, Eduard Gabriel Bazavan, Andrei Zanfir, and Cristian Sminchisescu. Structured 3d features for reconstructing controllable avatars. In *CVPR*, pages 16954–16964, 2023. **3**
- [10] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021. **1, 3**
- [11] Keyan Ding, Kede Ma, Shiqi Wang, and Eero P Simoncelli. Image quality assessment: Unifying structure and texture similarity. *IEEE TPAMI*, 2020. **8**
- [12] Patrick Esser, Robin Rombach, and Björn Ommer. Taming transformers for high-resolution image synthesis. *CVPR*, pages 12868–12878, 2021. **5**
- [13] Anna Frühstück, Krishna Kumar Singh, Eli Shechtman, Niloy J Mitra, Peter Wonka, and Jingwan Lu. Insetgan for full-body image generation. In *CVPR*, pages 7723–7732, 2022. **4**
- [14] Ke Gong, Yiming Gao, Xiaodan Liang, Xiaohui Shen, Meng Wang, and Liang Lin. Graphonomy: Universal human parsing via graph transfer learning. In *CVPR*, 2019. **7**
- [15] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc., 2014. **1, 4, 6**
- [16] Artur Grigorev, Karim Iskakov, Anastasia Ianina, Renat Bashirov, Ilya Zakharkin, Alexander Vakhitov, and Victor Lempitsky. Stylepeople: A generative model of fullbody human avatars. In *CVPR*, pages 5151–5160, 2021. **1, 2, 3, 4, 8**
- [17] Artur Grigorev, Artem Sevastopolsky, Alexander Vakhitov, and Victor Lempitsky. Coordinate-based texture inpainting for pose-guided human image generation. In *CVPR*, pages 12135–12144, 2019. **1**
- [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. **4**
- [19] Tong He, John Collomosse, Hailin Jin, and Stefano Soatto. Geo-pifu: Geometry and pixel aligned implicit functions for single-view human reconstruction. *Advances in Neural Information Processing Systems*, 33:9276–9287, 2020. **1, 2**
- [20] Tong He, Yuanlu Xu, Shunsuke Saito, Stefano Soatto, and Tony Tung. Arch++: Animation-ready clothed human reconstruction revisited. In *ICCV*, pages 11046–11056, 2021. **1, 3, 8**
- [21] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017. **8**
- [22] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020. **1, 3, 5**
- [23] Zeng Huang, Yuanlu Xu, Christoph Lassner, Hao Li, and Tony Tung. Arch: Animatable reconstruction of clothed humans. In *CVPR*, pages 3093–3102, 2020. **2, 3, 8**
- [24] Wei Jiang, Kwang Moo Yi, Golnoosh Samei, Oncel Tuzel, and Anurag Ranjan. Neuman: Neural human radiance field from a single video. In *ECCV*, pages 402–418. Springer, 2022. **2**
- [25] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *CVPR*, pages 8110–8119, 2020. **4, 7**
- [26] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2015. **7**
- [27] Alexey Larionov, Evgeniya Ustinova, Mikhail Sidorenko, David Svitov, Ilya Zakharkin, Victor Lempitsky, and Renat Bashirov. Morf: Mobile realistic fullbody avatars from a monocular video. *arXiv preprint arXiv:2303.10275*, 2023. **2**
- [28] Verica Lazova, Eldar Insafutdinov, and Gerard Pons-Moll. 360-degree textures of people in clothing from a single image. In *International Conference on 3D Vision (3DV)*, pages 643–653. IEEE, 2019. **1, 2**
- [29] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. Smpl: A skinned multi-person linear model. *ACM transactions on graphics (TOG)*, 34(6):1–16, 2015. **1, 2**
- [30] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR*, 2019. **7**
- [31] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. In *CVPR*, pages 11461–11471, 2022. **3**

- [32] Lars Mescheder, Andreas Geiger, and Sebastian Nowozin. Which training methods for gans do actually converge? In *International conference on machine learning*, pages 3481–3490. PMLR, 2018. [4](#), [6](#)
- [33] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021. [2](#)
- [34] Stylianos Mystakidis. Metaverse. *Encyclopedia*, 2(1):486–497, 2022. [1](#)
- [35] Xingyang Ni and Esa Rahtu. Flipreid: Closing the gap between training and inference in person re-identification. *9th European Workshop on Visual Information Processing (EU-VIP)*, pages 1–6, 2021. [8](#)
- [36] Joan M Ogden, Edward H Adelson, James R Bergen, and Peter J Burt. Pyramid-based computer graphics. *RCA engineer*, 30(5):4–15, 1985. [5](#)
- [37] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed A. A. Osman, Dimitrios Tzionas, and Michael J. Black. Expressive body capture: 3d hands, face, and body from a single image. In *CVPR*, 2019. [3](#), [7](#)
- [38] Sida Peng, Junting Dong, Qianqian Wang, Shangzhan Zhang, Qing Shuai, Xiaowei Zhou, and Hujun Bao. Animatable neural radiance fields for modeling dynamic human bodies. In *ICCV*, pages 14314–14323, 2021. [2](#)
- [39] Patrick Pérez, Michel Gangnet, and Andrew Blake. Poisson image editing. In *ACM SIGGRAPH 2003 Papers*, pages 313–318. 2003. [5](#)
- [40] Sergey Prokudin, Michael J Black, and Javier Romero. Smpix: Neural avatars from 3d human models. In *WACV*, pages 1810–1819, 2021. [1](#), [2](#)
- [41] Amit Raj, Julian Tanke, James Hays, Minh Vo, Carsten Stoll, and Christoph Lassner. Anr: Articulated neural rendering for virtual avatars. In *CVPR*, pages 3722–3731, 2021. [1](#), [2](#)
- [42] Edoardo Remelli, Timur Bagautdinov, Shunsuke Saito, Chenglei Wu, Tomas Simon, Shih-En Wei, Kaiwen Guo, Zhe Cao, Fabian Prada, Jason Saragih, et al. Drivable volumetric avatars using texel-aligned features. In *ACM SIGGRAPH Conference Proceedings*, pages 1–9, 2022. [1](#)
- [43] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, pages 10684–10695, 2022. [3](#), [5](#)
- [44] Andres Romero, Angela Castillo, Jose Abril-Nova, Radu Timofte, Ritwik Das, Sanchit Hira, Zhihong Pan, Min Zhang, Baopu Li, Dongliang He, et al. Ntire 2022 image inpainting challenge: Report. In *CVPR*, pages 1150–1182, 2022. [3](#)
- [45] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. [4](#)
- [46] Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion models. In *ACM SIGGRAPH 2022 Conference Proceedings*, pages 1–10, 2022. [3](#)
- [47] Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad Norouzi. Image super-resolution via iterative refinement. *IEEE TPAMI*, 2022. [1](#)
- [48] Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao Li. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In *ICCV*, pages 2304–2314, 2019. [1](#), [2](#), [8](#)
- [49] Shunsuke Saito, Tomas Simon, Jason Saragih, and Hanbyul Joo. Pifuhd: Multi-level pixel-aligned implicit function for high-resolution 3d human digitization. In *CVPR*, pages 84–93, 2020. [1](#)
- [50] Shunsuke Saito, Jinlong Yang, Qianli Ma, and Michael J Black. Scanimat: Weakly supervised learning of skinned clothed avatar networks. In *CVPR*, pages 2886–2897, 2021. [3](#)
- [51] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265. PMLR, 2015. [3](#)
- [52] Carole H Sudre, Wenqi Li, Tom Vercauteren, Sebastien Ourselin, and M Jorge Cardoso. Generalised dice overlap as a deep learning loss function for highly unbalanced segmentations. In *Deep learning in medical image analysis and multimodal learning for clinical decision support*, pages 240–248. Springer, 2017. [4](#), [6](#), [7](#)
- [53] Justus Thies, Michael Zollhöfer, and Matthias Nießner. Deferred neural rendering: Image synthesis using neural textures. *ACM Transactions on Graphics (TOG)*, 38(4):1–12, 2019. [1](#), [2](#), [4](#)
- [54] Chung-Yi Weng, Brian Curless, Pratul P Srinivasan, Jonathan T Barron, and Ira Kemelmacher-Shlizerman. Humanerf: Free-viewpoint rendering of moving people from monocular video. In *CVPR*, pages 16210–16220, 2022. [2](#)
- [55] Yulian Xiu, Jinlong Yang, Dimitrios Tzionas, and Michael J Black. Icon: Implicit clothed humans obtained from normals. In *CVPR*, pages 13286–13296. IEEE, 2022. [3](#)
- [56] Xiangyu Xu and Chen Change Loy. 3d human texture estimation from a single image with transformers. In *ICCV*, pages 13849–13858, 2021. [1](#), [2](#)
- [57] Tao Yu, Zerong Zheng, Kaiwen Guo, Pengpeng Liu, Qionghai Dai, and Yebin Liu. Function4d: Real-time human volumetric capture from very sparse consumer rgbd sensors. In *CVPR*, June 2021. [8](#)
- [58] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. [4](#)
- [59] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, pages 586–595, 2018. [8](#)
- [60] Hao Zhao, Jinsong Zhang, Yu-Kun Lai, Zerong Zheng, Yingdi Xie, Yebin Liu, and Kun Li. High-fidelity human avatars from a single rgb camera. In *CVPR*, pages 15904–15913, 2022. [2](#)