

AvatarCLIP: Zero-Shot Text-Driven Generation and Animation of 3D Avatars

FANGZHOU HONG*, S-Lab, Nanyang Technological University, Singapore

MINGYUAN ZHANG*, S-Lab, Nanyang Technological University, Singapore

LIANG PAN, S-Lab, Nanyang Technological University, Singapore

ZHONGANG CAI, S-Lab, Nanyang Technological University, Singapore and SenseTime Research, China

LEI YANG, SenseTime Research, China

ZIWEI LIU†, S-Lab, Nanyang Technological University, Singapore

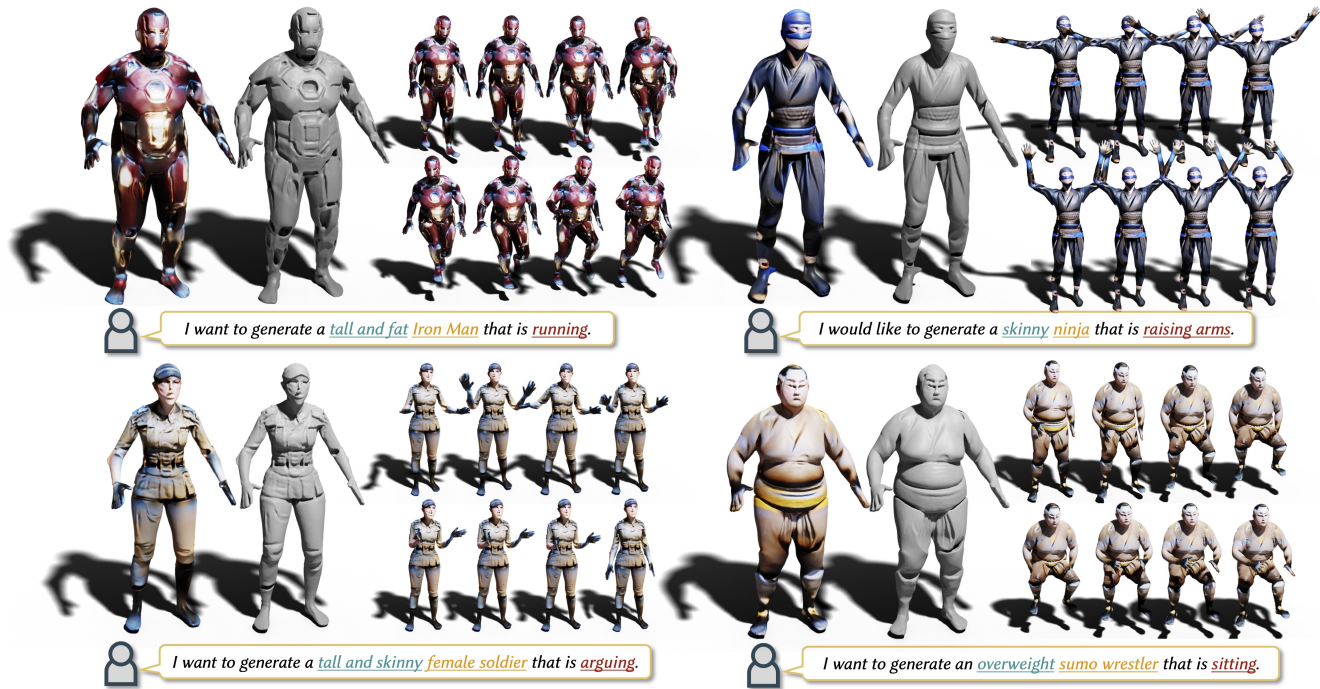


Fig. 1. In this work, we present AvatarCLIP, a novel zero-shot text-driven 3D avatar generation and animation pipeline. Driven by natural language descriptions of the desired **shape**, **appearance** and **motion** of the avatar, AvatarCLIP is capable of robustly generating 3D avatar models with vivid texture, high-quality geometry and reasonable motions.

*Both authors contributed equally to this research.

†corresponding author

Authors' addresses: Fangzhou Hong, fangzhou001@ntu.edu.sg, S-Lab, Nanyang Technological University, Singapore; Mingyuan Zhang, mingyuan001@ntu.edu.sg, S-Lab, Nanyang Technological University, Singapore; Liang Pan, liang.pan@ntu.edu.sg, S-Lab, Nanyang Technological University, Singapore; Zhongang Cai, caizhongang@sensetime.com, S-Lab, Nanyang Technological University, Singapore and SenseTime Research, China; Lei Yang, yanglei@sensetime.com, SenseTime Research, China; Ziwei Liu, ziwei.liu@ntu.edu.sg, S-Lab, Nanyang Technological University, Singapore.

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 ACM 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.

3D avatar creation plays a crucial role in the digital age. However, the whole production process is prohibitively time-consuming and labor-intensive. To democratize this technology to a larger audience, we propose AvatarCLIP, a zero-shot text-driven framework for 3D avatar generation and animation. Unlike professional software that requires expert knowledge, AvatarCLIP empowers layman users to customize a 3D avatar with the desired shape and texture, and drive the avatar with the described motions using solely natural languages. Our key insight is to take advantage of the powerful vision-language model CLIP for supervising neural human generation, in terms of 3D geometry, texture and animation. Specifically, driven by natural language descriptions, we initialize 3D human geometry generation with a shape VAE network. Based on the generated 3D human shapes, a volume rendering model is utilized to further facilitate geometry sculpting and

© 2022 Association for Computing Machinery.
0730-0301/2022/7-ART161 \$15.00
<https://doi.org/10.1145/3528223.3530094>

texture generation. Moreover, by leveraging the priors learned in the motion VAE, a CLIP-guided reference-based motion synthesis method is proposed for the animation of the generated 3D avatar. Extensive qualitative and quantitative experiments validate the effectiveness and generalizability of AvatarCLIP on a wide range of avatars. Remarkably, AvatarCLIP can generate unseen 3D avatars with novel animations, achieving superior zero-shot capability. Codes are available at <https://github.com/hongfz16/AvatarCLIP>.

CCS Concepts: • **Computing methodologies** → **Rendering**; *Animation*; Shape modeling.

Additional Key Words and Phrases: Zero-Shot Generation, Text-Driven Generation, 3D Avatar Generation, 3D Avatar Animation

ACM Reference Format:

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. *ACM Trans. Graph.* 41, 4, Article 161 (July 2022), 19 pages. <https://doi.org/10.1145/3528223.3530094>

1 INTRODUCTION

Creating digital avatars has become an important part in movie, game, and fashion industries. The whole process includes creating shapes of the character, drawing textures, rigging skeletons, and driving the avatar with the captured motions. Each step of the process requires many specialists that are familiar with professional software, large numbers of working hours, and expensive equipment, which are only affordable by large companies. The good news is that recent progress in academia, such as large-scale pre-trained models [Radford et al. 2021] and advanced human representations [Loper et al. 2015] are making it possible for this complicated work to be available to small studios and even reach out to the mass crowd. In this work, we take a step further and propose AvatarCLIP, which is capable of generating and animating 3D avatars solely from natural language descriptions as shown in Fig. 1.

Previously, there are several similar efforts, e.g. avatar generation, and motion synthesis, towards this vision. For avatar generation, several attempts have been made in 2D-GAN-based human generation [Lewis et al. 2021; Sarkar et al. 2021b] and 3D neural human generation [Grigorev et al. 2021; Zhang et al. 2021]. However, most of them either compromise the generation quality and diversity or have limited control over the generation process. Not to mention, most generated avatars cannot be easily animated. As for motion synthesis, impressive progresses have been made in recent years. They are able to generate motions conditioned on action classes [Guo et al. 2020; Petrovich et al. 2021] or motion trajectories [Kania et al. 2021; Ling et al. 2020]. However, they typically require paired data for fully supervised training, which not only limits the richness of the generated motions but also makes the generation process less flexible. There still exists a considerable gap between existing works and the vision of making digital avatar creation manageable for the mass crowd.

To skirt the complicated operations, natural languages could be used as a user-friendly control signal for convenient 3D avatars generation and animation. However, there are no existing high-quality avatar-text datasets to support supervised text-driven 3D avatar generation. As for animating 3D avatars, a few attempts [Ghosh et al. 2021; Tevet et al. 2022] have been made towards text-driven

motion generation by leveraging a motion-text dataset. Nonetheless, restricted by the scarcity of paired motion-text data, those fully supervised methods have limited generalizability.

Fortunately, recent advances in vision-language models pave the way toward zero-shot text-driven generation. CLIP [Radford et al. 2021] is a vision-language pre-trained model trained with large-scale image-text pairs. Through direct supervision on images, CLIP shows great success in zero-shot text-guided image generation [Ramesh et al. 2021]. Inspired by this thread of works, we choose to take advantage of the powerful CLIP to achieve zero-shot text-driven generation and animation of 3D avatars. However, neither the 3D avatar nor motion sequences can be directly supervised by CLIP. Major challenges exist in both static avatar generation and motion synthesis.

For static 3D avatar generation, the challenges lie in three aspects, namely texture generation, geometry generation, and the ability to be animated. Inspired by the recent advances in neural rendering [Tewari et al. 2021], the CLIP supervision can be applied to the rendered images to guide the generation of an implicit 3D avatar, which facilitates the generation of avatar textures. Moreover, to speed up the optimization, we propose to initialize the implicit function based on a template human mesh. At optimization time, we also use the template mesh constraint to control the overall shape of the implicit 3D avatar. To adapt to the modern graphics pipeline and, more importantly, to be able to be animated later, we need to extract meshes from generated implicit representations. Therefore, other than the texture, it is also desirable to generate high-quality geometry. To tackle the problem, the key insight is that when inspecting the geometry of 3D models on the computer, users usually turn off texture shading to get texture-less renderings, which can directly reveal the geometry. Therefore, we propose to randomly cast light on the surface of the implicit 3D avatar to get texture-less renderings, upon which the CLIP supervision is applied. Last but not least, to make the generated implicit 3D avatar animatable, we propose to leverage the recent achievements in parametric human models [Loper et al. 2015]. Specifically, we align and register the generated 3D avatar to a SMPL mesh. So that it can be driven by the SMPL skeletons.

CLIP is only trained with static images and insensitive to sequential motions. Therefore, it is inherently challenging to generate reasonable motion sequences using only the supervision from CLIP. To tackle this problem, we divide the whole process into two stages: 1) generating candidate poses with the guidance of CLIP, and 2) synthesizing smooth and valid motions with the candidate poses as references. In the first stage, a code-book consisting of diverse poses is created by clustering. Poses that match motion descriptions are selected by CLIP from the code-book. These generated poses serve as important clues for the second stage. A motion VAE is utilized in the second stage to learn motion priors, which facilitates the reference-guided motion synthesis.

With careful designs of the whole pipeline, AvatarCLIP is capable of generating high-quality 3D avatars with reasonable motion sequences guided by input texts as shown in Fig. 1. To evaluate our framework quantitatively, comprehensive user studies are conducted in terms of both avatar generation and animation to show

our superiority over existing solutions. Moreover, qualitative experiments are also performed to validate the effectiveness of each component in our framework.

To sum up, our contributions are listed below: **1)** To the best of our knowledge, it is the first text-driven full avatar synthesis pipeline that includes the generation of shape, texture, and motion. **2)** Incorporating the power of large-scale pre-trained models, the proposed framework demonstrates strong zero-shot generation ability. Our avatar generation pipeline is capable of robustly generating animation-ready 3D avatars with high-quality texture and geometry. **3)** Benefiting from the motion VAE, a novel zero-shot text-guided reference-based motion synthesis approach is proposed. **4)** Extensive qualitative and quantitative experiments show that the generated avatars and motions are of higher quality compared to existing methods and are highly consistent with the corresponding input natural languages.

2 RELATED WORK

Avatar Modeling and Generation. For its wide application in industries, human modeling has been thoroughly studied for decades. Driven by a large-scale human body dataset [Pishchulin et al. 2017], SMPL [Loper et al. 2015] and SMPL-X [Pavlakos et al. 2019] are proposed as a parametric human model. For its strong interpretability and compatibility with the modern graphics pipeline, we choose SMPL as the template of our avatar. However, they only provide the ability to model the human body without clothes. Many efforts [Bhatnagar et al. 2020, 2019; Hong et al. 2021; Jiang et al. 2020] have been made on the modeling of clothed humans.

Due to the complexity of clothes and accessories, non-parametric human modelings [Alldieck et al. 2018; Burov et al. 2021; Corona et al. 2021; Habermann et al. 2021; Huang et al. 2020; Mihajlovic et al. 2021; Palafox et al. 2021; Saito et al. 2021; Weng et al. 2019] are proposed to offer more flexibility on realistic human modeling. Moreover, inspired by recent advances in volume rendering, non-parametric human modeling methods based on the neural radiance field (*i.e.* NeRF [Jain et al. 2021b; Mildenhall et al. 2020]) have also been studied [Habermann et al. 2021; Liu et al. 2021; Peng et al. 2021a,b; Xu et al. 2021a; Zhao et al. 2021]. Combining advantages of volume rendering and SDF, NeuS [Wang et al. 2021b] is proposed recently to achieve high-quality geometry and color reconstruction. That justifies our choice of NeuS as base representations of avatars.

With the rapid development of deep learning in recent years, impressive progresses have been shown on 2D image generation. 2D face generation [Fu et al. 2022; Jiang et al. 2021; Karras et al. 2019, 2020] is now a very mature technology for the simplicity of the face structure and large-scale high-quality face datasets [Liu et al. 2015]. Recent works [Han et al. 2018; Jiang et al. 2022; Lewis et al. 2021; Sarkar et al. 2021b] have also demonstrated wonderful results in terms of 2D human body generation and manipulation. Without the knowledge of 3D space, it has always been a difficult task to animate the 2D human body [Sarkar et al. 2021a; Siarohin et al. 2019a,b; Yoon et al. 2021].

The 3D human generation [Chen et al. 2022; Noguchi et al. 2021, 2022] is barely explored until the very recent. Combining the powerful SMPL-X and StyleGAN, StylePeople [Grigorev et al. 2021]

proposes a data-driven animatable 3D avatar generation method. Inspired by the advancements in NeRF-GAN [Chan et al. 2021], a recent work [Zhang et al. 2021] proposes to bring in 3D awareness to traditional 2D generation pipelines.

Motion Synthesis. Serving as one of the most significant parts of animation, motion synthesis has been the research focus of many researchers. Several large-scale datasets [Cai et al. 2022, 2021; Ionescu et al. 2013; Mahmood et al. 2019; Mehta et al. 2017; Varol et al. 2017; von Marcard et al. 2018] provide human motions as sequences of 3D keypoints or SMPL parameters. The rapid advancements of motion datasets stimulate researches on the motion synthesis. Early works apply classical machine learning to unconditional motion synthesis [Ikemoto et al. 2009; Mukai and Kuriyama 2005; Rose et al. 1998]. DVGANs [Xiao Lin 2014], Text2Action [Ahn et al. 2018] and Language2Pose [Ahuja and Morency 2019] generate motions conditioned on short texts using fully annotated data. Action2Motion [Guo et al. 2020] and Actor [Petrovich et al. 2021] condition the motion generation on pre-selected action classes. These methods require large amounts of data [Hong et al. 2022] with annotations of action classes or language descriptions, which limits their applications. On the contrary, our proposed AvatarCLIP can drive the human model by general natural languages without any paired data. Some other works [Aggarwal and Parikh 2021; Li et al. 2021] focus on music-conditioned motion synthesis. Moreover, some works [Bergamin et al. 2019; Won and Lee 2019] focus on the physics-based motion synthesis. They construct motion sequences with physical constrains for more realistic generation results.

Zero-shot Text-driven Generation. Text-to-image synthesis [Mansimov et al. 2015] has long been studied. The ability to zero-shot generalize to unseen categories is first shown by [Reed et al. 2016]. CLIP and DALL-E [Ramesh et al. 2021] further show the incredible text-to-image synthesis ability by excessively scale-up the size of training data. Benefiting from the zero-shot ability of CLIP, many amazing zero-shot text-driven applications [Frans et al. 2021; Patashnik et al. 2021; Xu et al. 2021b] are being developed. Combining CLIP with 3D representations like NeRF or mesh, zero-shot text-driven 3D object generation [Jain et al. 2021a; Jetchev 2021; Michel et al. 2021; Sanghi et al. 2021] and manipulation [Wang et al. 2021a] have also come true in recent months.

3 OUR APPROACH

Recall that our goal is zero-shot text-driven 3D avatar generation and animation, which can be formally defined as follows. The inputs are natural languages, $\text{text} = \{t_{\text{shape}}, t_{\text{app}}, t_{\text{motion}}\}$. The three texts correspond to the descriptions of the desired body shape, appearance and motion. The output is two-part, including **a)** an animatable 3D avatar represented as a mesh $M = \{V, F, C\}$, where V is the vertices, F stands for faces, C represents the vertex colors; **b)** a sequence of poses $\Theta = \{\theta_i\}_{i=1}^L$ comprising the desired motion, where L is the length of the sequence.

3.1 Preliminaries

CLIP. [Radford et al. 2021] is a vision-language pre-trained model trained with large-scale image-text datasets. It consists of an image

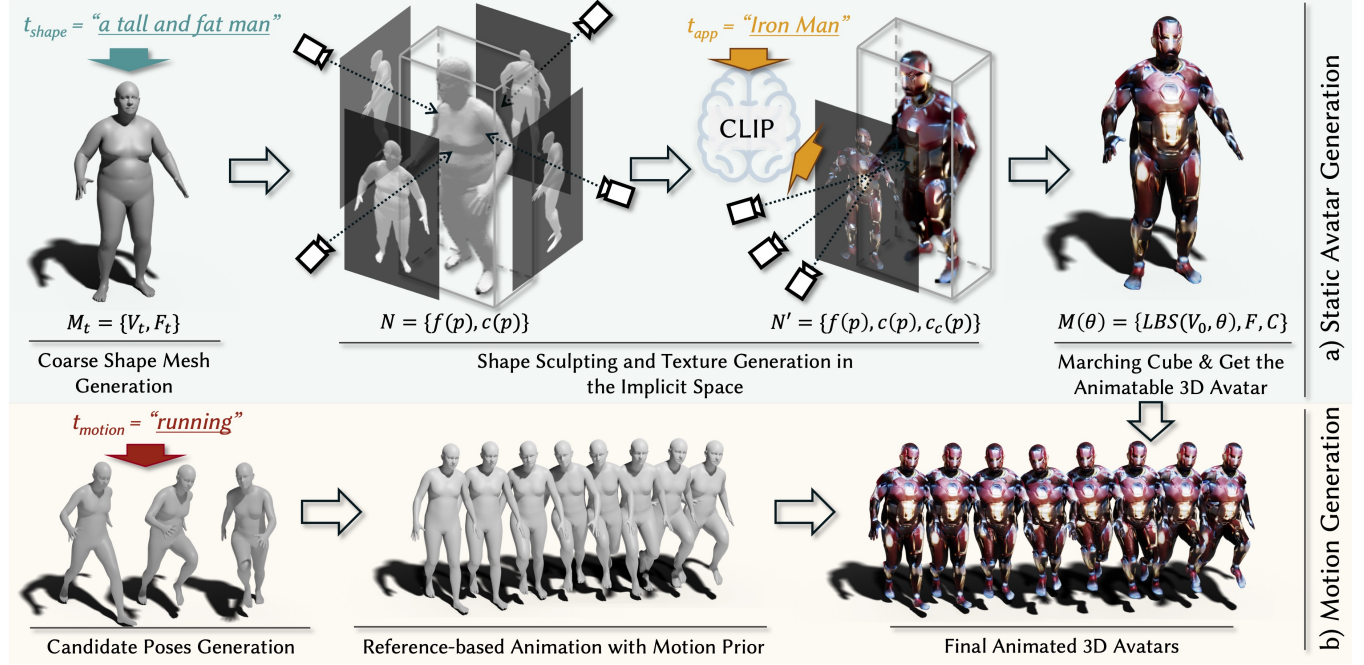


Fig. 2. **An Overview of the Pipeline of AvatarCLIP.** The whole pipeline is divided into two parts: a) Static Avatar Generation; b) Motion Generation. Assume the user want to generate ‘a tall and fat Iron Man that is running’. An animatable avatar is generated guided by t_{shape} = ‘a tall and fat man’ and t_{app} = ‘Iron Man’. Then a motion sequence matching the description t_{motion} = ‘running’ is generated to animate the generated avatar.

encoder E_I and a text encoder E_T . The encoders are trained in the way that the latent codes of paired images and texts are pulled together and unpaired ones are pushed apart. The resulting joint latent space of images and texts enables the zero-shot text-driven generation by encouraging the latent code of the generated content to align with the latent code of the text description. Formally, the CLIP-guided loss function is defined as

$$\mathcal{L}_{clip}(I, T) = 1 - \text{norm}(E_I(I)) \cdot \text{norm}(E_T(T)), \quad (1)$$

where (\cdot) represents the cosine distance. By minimizing $\mathcal{L}_{clip}(I, T)$, the generated image I is encouraged to match the description of the text prompt T .

SMPL. [Loper et al. 2015] is a parametric human model driven by large-scale aligned human surface scans [Pishchulin et al. 2017]. The SMPL model can be formally defined as $M_{SMPL}(\beta, \theta; \Phi)$, where $\beta \in \mathbb{R}^{10}$ controls the body shapes, $\theta \in \mathbb{R}^{24 \times 3}$ contains the axis angles of each joint. In this work, we use SMPL as the template meshes for the initialization of the implicit 3D avatar and the automatic skeleton-binding and animation.

NeuS. [Wang et al. 2021b] proposes a novel volume rendering method that combines the advantages of SDF and NeRF [Mildenhall et al. 2020], leading to high-quality surface reconstruction as well as photo-realistic novel view rendering. For some viewing point o and viewing direction v , the color of the corresponding pixel is accumulated along the ray by

$$C(o, v) = \int_0^\infty w(t)c(p(t), v)dt, \quad (2)$$

where $p(t) = o + vt$ is some point along the ray, $c(p(t), v)$ output the color at point $p(t)$, which is implemented by MLPs. $w(t)$ is a weighting function for the point $p(t)$. $w(t)$ is designed to be unbiased and occlusion-aware, so that the optimization on multi-view images can lead to the accurate learning of a SDF representation. It is defined as

$$w(t) = \frac{\phi_s(f(p(t)))}{\int_0^\infty \phi_s(f(p(u)))du}, \quad (3)$$

where ϕ_s is the logistic density distribution, f is an SDF network. NeuS is used in our work as the representation of the implicit 3D avatar.

3.2 Pipeline Overview

As shown in Fig. 2, the pipeline of our AvatarCLIP can be divided to two parts, *i.e.* static avatar generation and motion generation. For the first part, the natural language description of the shape t_{shape} is used for the generation of a coarse body shape β_t . Together with a pre-defined standing pose θ_{stand} , we get the template mesh $M_t = \{V_t, F_{SMPL}\}$, where vertices $V_t = M(\beta_t, \theta_{stand}; \Phi)$ and faces F_{SMPL} are given by SMPL. M_t is then rendered to multi-view images for the training of a NeuS model N , which is later used as the initialization of our implicit 3D avatar. Guided by the appearance description t_{app} , N is further optimized by CLIP in a shape-preserving way for shape sculpting and texture generation. After that, the target static 3D avatar mesh $M = \{V, F, C\}$ is extracted from N' by the marching cube algorithm [Lorensen and Cline 1987] and aligned with M_t to

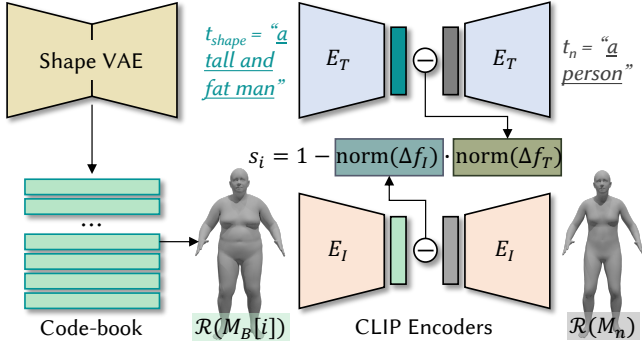


Fig. 3. **Illustration of the Coarse Shape Generation.** A shape VAE is trained to construct a code-book which is used for CLIP-guided query to get a best match for the input text t_{shape} . To introduce the awareness of body attributes like height, a neutral shape M_n and text t_n is defined as the anchor. The relative direction in latent space is used for the CLIP-guided query.

get ready for animation. For the second part, the natural language description of the motion t_{motion} is used to generate candidate poses from a pre-calculated code-book. Then the candidate poses are used as references for the optimization of a pre-trained motion VAE to get the desired motion sequence.

3.3 Static Avatar Generation

3.3.1 Coarse Shape Generation. As the first step of generating avatars, we propose to generate a coarse shape as the template mesh M_t from the input description of the shape t_{shape} . For this step, SMPL is used as the source of possible human body shapes. As shown in Fig. 3, we first construct a code-book by clustering body shapes. Then CLIP is used to extract features of the renderings of body shapes in the code-book and texts to get the best matching body shape. Although the process is straightforward, careful designs are required to make full use of both powerful tools.

For the code-book construction, it is important to sample uniformly to cover most of the possible body shapes. The SMPL shape space $\beta \in \mathbb{R}^{10}$ is learned by the PCA decomposition, which results in non-uniform distribution. Therefore, we propose to train a shape VAE for a more uniform latent shape space. The code-book $B_{\text{shape}} \in \mathbb{R}^{K \times d_s}$ is constructed by K-Means clustering on the latent space of the shape VAE, which can be decoded to a set of body meshes $M_B = \{V_B, F_{\text{SMPL}}\}$, where K is the size of the code-book, d_s is the dimension of the latent code.

For the CLIP-guided code-book query, it is important to design reasonable query scores. We observe that for some attributes like body heights, it is hard to be determined only by looking at the renderings of the body without any reference. Therefore, we propose to set a reference and let CLIP score the body shapes in a relative way, which is inspired by the usage of CLIP in 2D image editing [Patashnik et al. 2021]. We define a neutral state body shape M_n and the corresponding neutral state text t_n as the reference. The scoring of each entry s_i in the code-book is defined as

$$s_i = 1 - \text{norm}(\Delta f_i) \cdot \text{norm}(\Delta f_T), \quad (4)$$

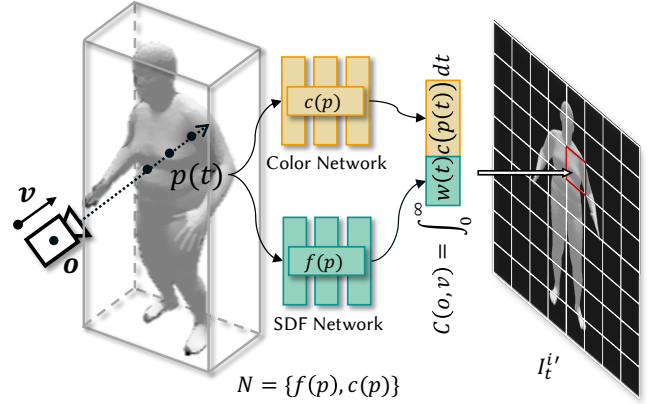


Fig. 4. **Initialization of the Implicit 3D Avatar.** Multi-view renderings of the template mesh M_t is used to optimize a randomly initialized NeuS network N , which is later used as an initialization of the implicit 3D avatar.

where $\Delta f_i = E_I(\mathcal{R}(M_B[i])) - E_I(\mathcal{R}(M_n))$, $\Delta f_T = E_T(t_{\text{shape}}) - E_T(t_n)$, $\mathcal{R}(\cdot)$ denotes a mesh renderer, E_I , E_T represents CLIP image and text encoders. Δf_T is the relative direction guided by the text, which is intended to be aligned with the visual relative shape differences Δf_i in the CLIP latent space. Then the entry with the maximum score among the code-book is retrieved as the coarse shape $M_t = M_B[\text{argmax}_i(s_i)]$.

3.3.2 Shape Sculpting and Texture Generation. The generated template mesh M_t represents a desired coarse naked body shape. To generate high-quality 3D avatars, the shape and texture need to be further sculpted and generated to match the description of the appearance t_{app} . As discussed previously, we choose to use an implicit representation, *i.e.* NeuS, as the base 3D representation in this step for its advantages in both geometry and colors. In order to speed up the optimization and more importantly, control the general shape of the avatar for the convenience of animation, this step is designed as a two-stage optimization process.

The first stage creates an equivalent implicit representation of M_t by optimizing a randomly initialized NeuS N with the multi-view renderings $\{I_{M_t}^i\}$ of the template mesh M_t . Specifically, as shown in Fig. 4, the NeuS $N = \{f(p), c(p)\}$ comprises of two sub-networks. The SDF network $f(p)$ takes some point p as input and outputs the signed distance to its nearest surface. The color network $c(p)$ takes some point p as input and outputs the color at that point. Both $f(p)$ and $c(p)$ are implemented using MLPs. It should be noted that compared to the original implementation of NeuS, we omit the color network's dependency on the viewing direction. The viewing direction dependency is originally designed to model the shading that changes with the viewing angle. In our case, it is redundant to model such shading effects since our goal is to generate 3D avatars with consistent textures, which ideally should be albedo maps. Similar to the original design, N is optimized by a three-part loss function

$$\mathcal{L}_1 = \mathcal{L}_{\text{color}} + \lambda_1 \mathcal{L}_{\text{reg}} + \lambda_2 \mathcal{L}_{\text{mask}}. \quad (5)$$

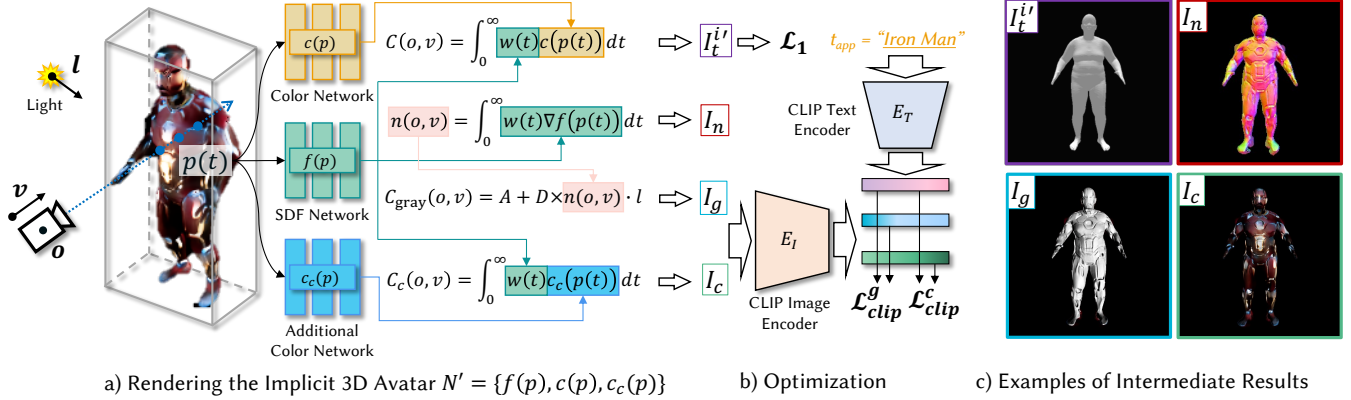


Fig. 5. **Detailed Method of Shape Sculpting and Texture Generation.** An additional color network $c_c(p)$ is appended to the initialized implicit 3D avatar for texture generation, which is illustrated in a). Three types of constraints introduced for the optimization are shown in b), including reconstruction loss \mathcal{L}_1 , CLIP-guided losses \mathcal{L}_{clip}^g and \mathcal{L}_{clip}^c for the geometry sculpting and texture generation, respectively. The sub-figure c) shows examples of intermediate results.

\mathcal{L}_{color} is the reconstruction loss between the rendered images and the ground truth multi-view renderings $\{I_{M_t}^i\}$. \mathcal{L}_{reg} is an Eikonal term [Gropp et al. 2020] to regularize the SDF $f(p)$. \mathcal{L}_{mask} is a mask loss that encourages the network to only reconstruct the foreground object, *i.e.* the template mesh M_t . The resulting NeuS N serves as an initialization for the second stage optimization.

The second stage stylizes the initialized NeuS N from the first stage by the power of CLIP. At the same time, the coarse shape M_t should still be maintained, resulting in two possible solutions. The first one is to fix the weights of $f(p)$, which ensures the final generated shape is the same as M_t . The color network $c(p)$ is optimized to ‘colorize’ the fixed shape. However, as discussed before, M_t is generated from SMPL, which only provides the shape of a naked human body. Purely adding textures on the template shapes, although simple and tractable, is not considered the optimal solution. Because not only textures but also geometry is crucial in the process of avatar creation. To allow the fine-level shape sculpting, we need to fine-tune the weights of $f(p)$. At the same time, to maintain the general shape of the template mesh, $c(p)$ should maintain its original function of reconstructing the template mesh to allow the reconstruction loss during the optimization. This leads to the second solution, where we keep the original two sub-networks $f(p)$ and $c(p)$ and introduce an additional color network $c_c(p)$. Both color networks share the same SDF network $f(p)$ as shown in Fig. 5. $f(p)$ and $c(p)$ are in charge of the reconstruction of the template mesh M_t . $f(p)$ together with $c_c(p)$ are responsible for the stylizing part and comprises the targeting implicit 3D avatar. Formally, the new NeuS model $N' = \{f(p), c(p), c_c(p)\}$ now consists of three sub-networks, where $f(p)$ and $c(p)$ are initialized by the pre-trained N , and $c_c(p)$ is a randomly initialized color network. All three sub-networks are

trainable. Similarly, we ignore the viewing direction dependency in the additional color network.

The second-stage optimization is supervised by a three-part loss function

$$\mathcal{L}_2 = \mathcal{L}_1 + \lambda_3 \mathcal{L}_{clip}^c + \lambda_4 \mathcal{L}_{clip}^g, \quad (6)$$

where \mathcal{L}_1 is the reconstruction loss over NeuS $\{f(p), c(p)\}$ as defined in Eq. 5. \mathcal{L}_{clip}^c and \mathcal{L}_{clip}^g are CLIP-guided loss functions that guide the texture and geometry generation to match the description t_1 using two different types of renderings, which are introduced as follows.

The first type of the rendering is the colored rendering I_c of the NeuS model $\{f(p), c_c(p)\}$, which is calculated by applying Eq. 2 to each pixel of the image. The second type is the texture-less rendering I_g of the same NeuS model. The rendering algorithm used here is ambient and diffuse shading. For each ray $\{o, v\}$, the normal direction $n(o, v)$ of the first surface point it intersects with can be calculated by the accumulation of the gradient of the SDF function at each position along the ray, which is formulated as

$$n(o, v) = \int_0^\infty w(t) \nabla f(p(t)) dt, \quad (7)$$

where $w(t)$ is the weighting function as defined in Eq. 3. To render the texture-less model, a random light direction needs to be sampled. To avoid the condition where the light and camera are at opposite sides of the model and the geometry details of the model cannot be revealed, the light direction is uniformly sampled in a small range around the camera direction. Formally, defining the camera direction in the spherical coordinate system as polar and azimuthal angles $\{\theta_c, \phi_c\}$, then the light direction l is sampled from $\{\theta_c + X_1, \phi_c + X_2\}$, where $X_1, X_2 \sim \mathcal{U}(-\pi/4, \pi/4)$. Since coloring is not needed for the

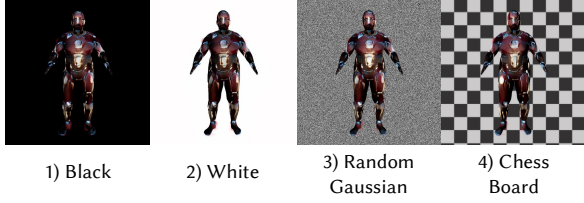


Fig. 6. Examples of four types of random background augmentations.

texture-less rendering, we could simply calculate the gray level of the ray $\{o, v\}$ by

$$C_{gray}(o, v) = A + D \times n(o, v) \cdot l, \quad (8)$$

where $A \sim \mathcal{U}(0, 0.2)$ is randomly sampled from a uniform distribution, $D = 1 - A$ is the diffusion. By applying Eq. 8 to each pixel of the image, we get the texture-less rendering I_g . In practice, we find that random shading on textured renderings I_c benefiting the uniformity of the generated textures, which is discussed later in the ablation study. The random shading on I_c is similar to the texture-less rendering, which is formally defined as

$$C_{shade}(o, v) = A + D \times n(o, v) \cdot l * C(o, v). \quad (9)$$

Then the two types of CLIP-guided loss functions are formally defined as $\mathcal{L}_{clip}^c = \mathcal{L}_{clip}(I_c, t_{app})$, $\mathcal{L}_{clip}^g = \mathcal{L}_{clip}(I_g, t_{app})$.

To render a $H \times W$ image, a total of $H \times W \times Q$ queries need to be performed given Q query times for each ray. Due to the high memory footprint of volume rendering, the rendering resolution $H \times W$ is heavily constrained. In our experiment on one 32GB GPU, the maximum H_{max} and W_{max} are around 110, which is far from the resolution upper bound of 224 provided by CLIP. With the network structures and other hyper-parameters not modified, to increase the rendering resolution, we propose a dilated silhouettes-based rendering strategy based on the fact that the rays not encountering any surface do not contribute to the final rendering while consuming large amounts of memories. Specifically, we could get the rays that have high chances of encountering surfaces by calculating the silhouette of the rendering of the template mesh M_t with the given camera parameter. Moreover, to ensure a proper amount of shape sculpting space is allowed, we further dilate the silhouettes. The rays within the silhouettes are calculated and make contributions to the final rendering. Defining the ratio between the area of the dilated silhouette and the total area of the image as r_s , this rendering strategy dynamically increases the maximum resolution to $H'_{max} \times W'_{max} = H_{max} \times W_{max} / r_s$. Empirically, the resolution can be increased to around 150^2 .

In order to further increase the robustness of the optimization process, three augmentation strategies are proposed: **a)** random background augmentation; **b)** random camera parameter sampling; **c)** semantic-aware prompt augmentation. They are introduced as follows.

Inspired by Dream Fields [Jain et al. 2021a], random background augmentation helps CLIP to focus more on the foreground object and prevents the generation of randomly floating volumes. As shown in Fig. 6, we randomly augment the backgrounds of the renderings to

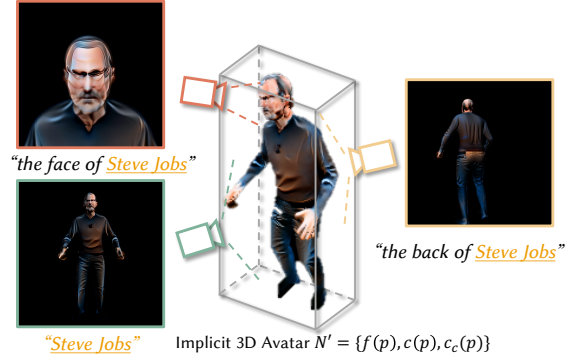


Fig. 7. Examples of two types of prompt augmentations. One for the detailed refinement of the face and the other for the back side of the avatar.

1) pure black background; 2) pure white background; 3) Gaussian noises; 4) Gaussian blurred chessboard with random block sizes.

To prevent the network from finding short-cut solutions that only give reasonable renderings for several fixed camera positions, we randomly sample the camera extrinsic parameters for each optimization iteration in a manually-defined importance sampling way to get more coherent and smooth results. We choose to work with the ‘look at’ mode of the camera, which consists of a look-at point, a camera position, and an up direction. We set the up direction to always align with the up direction of the template body. The look at position is sampled from a Gaussian distribution $X, Y, Z \sim \mathcal{N}(0, 0.1)$, which are then clipped between -0.3 and 0.3 to prevent the avatar from being out of the frame. The camera position is sampled using a spherical coordinate system with the radius R sampled from a uniform distribution $\mathcal{U}(1, 2)$, the polar angle θ_c sampled from a uniform distribution $\mathcal{U}(0, 2\pi)$, the azimuthal angle ϕ_c sampled from a Gaussian distribution $\mathcal{N}(0, \pi/3)$ for the camera to point at the front side of the avatar in the most of the iterations.

So far, no human prior is introduced in the whole optimization process, which may result in generating textures at wrong body parts or not generating textures for important body parts, which will later be shown in the ablation study. To prevent such problems, we propose to explicitly bring in human prior in the optimization process by augmenting the prompts in a semantic-aware manner. For example, as shown in Fig. 7, if $t_{app} = \text{‘Steve Jobs’}$, we would augment t_{app} to two additional prompts $t_{face} = \text{‘the face of Steve Jobs’}$ and $t_{back} = \text{‘the back of Steve Jobs’}$. For every four iterations, the look-at point of the camera is set to be the center of the face to get renderings of the face, which will be supervised by the prompt t_{face} . The ‘face augmentation’ directly supervises the generation of the face, which is important to the quality of the generated avatar since humans are more sensitive to faces. For the second augmented prompt, when the randomly sampled camera points at the back of the avatar, t_{back} is used as the corresponding text to explicitly guide the generation of the back of the avatar.

3.3.3 Make the Static Avatar Animatable. With the generated implicit 3D avatar N' , the marching cube algorithm is performed to

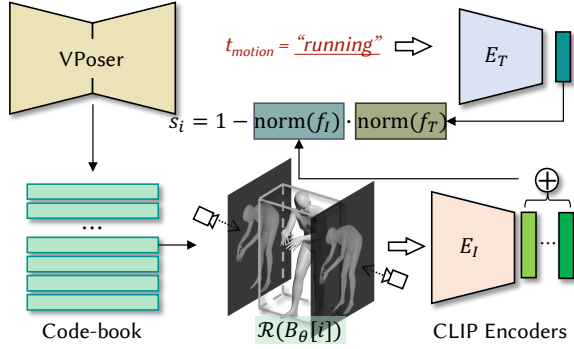


Fig. 8. **Detailed Pipeline of Candidate Poses Generation.** The pre-trained VPoser is first used to build a code-book. Given text description t_{motion} , each pose feature f_i from the code-book is used to calculate the similarity with the text feature f_T , which is used to select Top-K entries as candidate poses.

extract its meshes $M = \{V, F, C\}$ before making it animatable. Naturally, with the careful design of keeping the overall shape unchanged during the optimization, the generated mesh can be aligned with the initial template mesh. Firstly, the nearest neighbor for each vertex in V is retrieved in the vertices of the template mesh M_t . Blend weights of each vertex in V are copied from the nearest vertex in M_t . Secondly, an inverse LBS algorithm is used to bring the avatar's standing pose θ_{stand} back to the zero pose θ_0 . The vertices V are transformed to V_{θ_0} . Finally, V_{θ_0} can be driven by any pose θ using the LBS algorithm. Hence, for any pose, the animated avatar can be formally defined as $M(\theta) = (\text{LBS}(V_{\theta_0}, \theta), F, C)$.

3.4 Motion Generation

Empirically, CLIP is not capable of directly estimating the similarities between motion sequences and natural language descriptions. It also lacks the ability to assess the smoothness or rationality of motion sequences. These two limitations suggest that it is hard to generate motions only using CLIP supervision. Hence, we have to introduce other modules to provide motion priors. However, CLIP has the ability to value the similarity between a rendered human pose and a description. Furthermore, similar poses can be regarded as references for the expected motions. Based on the above observations, we propose a two-stage motion generation process: 1) candidate poses generation guided by CLIP. 2) motion sequences generation using motion priors with candidate poses as references. The details are illustrated as follows.

3.4.1 Candidate Poses Generation. To generate poses consistent with the given description t_{motion} , an intuitive method is to directly optimize the parameter θ in the SMPL model or the latent code of a pre-trained pose VAE (e.g. VPoser [Pavlakos et al. 2019]). However, they can hardly yield reasonable poses due to difficulties during optimization, which are later shown in the experiments. Hence, it is not a wise choice to directly optimize the poses.

As shown in Fig. 8, to avoid the direct optimization, we first create a code-book from AMASS dataset [Mahmood et al. 2019]. To reduce dimensions, VPoser is used to encode poses to $z \in \mathbb{R}^{d_p}$. Then we

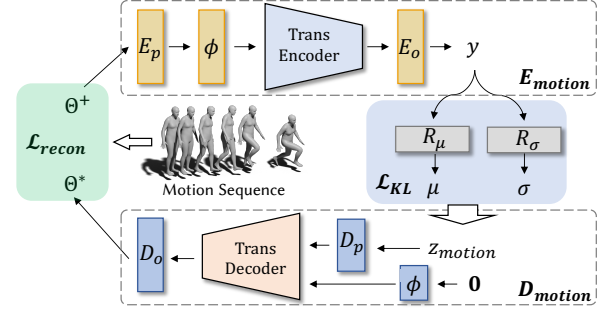


Fig. 9. **Structure of the Motion VAE.** The motion VAE contains three parts: the encoder E_{motion} , the decoder D_{motion} , and a reparameterization module. The reconstruction loss $\mathcal{L}_{\text{recon}}$ and the KL-divergence term \mathcal{L}_{KL} are used for the motion VAE training.

use K-Means to acquire K cluster centroids that form our pose code-book $B_z = \mathbb{R}^{K \times d_p}$. Each element of B_z is then decoded by VPoser, which leads to a set of poses B_θ .

Given the motion description t_{motion} , we calculate the similarity between t_{motion} and each pose $B_\theta[i]$ from the code-book B_θ , which can be defined as

$$s_i = 1 - \text{norm}(E_I(\mathcal{R}(B_\theta[i]))) \cdot \text{norm}(E_T(t_{\text{motion}})). \quad (10)$$

Top-k scores s_i and their corresponding poses $B_\theta[i]$ are selected to construct the candidate pose set S , which serves as references to generate a motion sequence in the next stage.

3.4.2 Reference-Based Animation with Motion Prior. We propose a two-fold method to generate a target motion sequence that matches the motion description t_{motion} . 1) A motion VAE is trained to capture human motion priors. 2) We optimize the latent code of the motion VAE using candidate poses S as references. The details are introduced as follows.

Motion VAE Pre-train. We take inspirations from Actor [Petrovich et al. 2021] to construct the motion VAE. Specifically, as shown in Fig. 9, the motion VAE contains three parts: the motion encoder E_{motion} , the reparameterization module, and the motion decoder D_{motion} . A motion sequence is denoted as $\Theta^+ \in \mathbb{R}^{L \times 24 \times 6}$, where L represents the length of the motion. Each joint is represented by a continuous 6-D tensor [Zhou et al. 2019].

The motion encoder E_{motion} contains a projection layer, a positional embedding layer, several transformer encoder layers and an output layer. Formally, E_{motion} can be defined as $y = E_{\text{motion}}(\Theta^+) = E_o(\text{TransEncoder}(\phi(E_p(\Theta^+))))$, where E_p and E_o are fully connected layers, ϕ represents positional embedding operation [Vaswani et al. 2017]. TransEncoder includes multiple transformer encoder layers [Vaswani et al. 2017].

The reparameterization module yields a Gaussian distribution where $\mu = R_\mu(y)$, $\sigma = R_\sigma(y)$. R_μ, R_σ are fully connected layers to calculate the mean μ and the standard deviation σ of the distribution, respectively. Using the reparameterization trick [Kingma and Welling 2013], a random latent code z_{motion} is sampled under the distribution $\mathcal{N}(\mu, \sigma)$. The latent code is further decoded by D_{motion} , which can be formally defined as $\Theta^* = D_{\text{motion}}(z_{\text{motion}}) =$

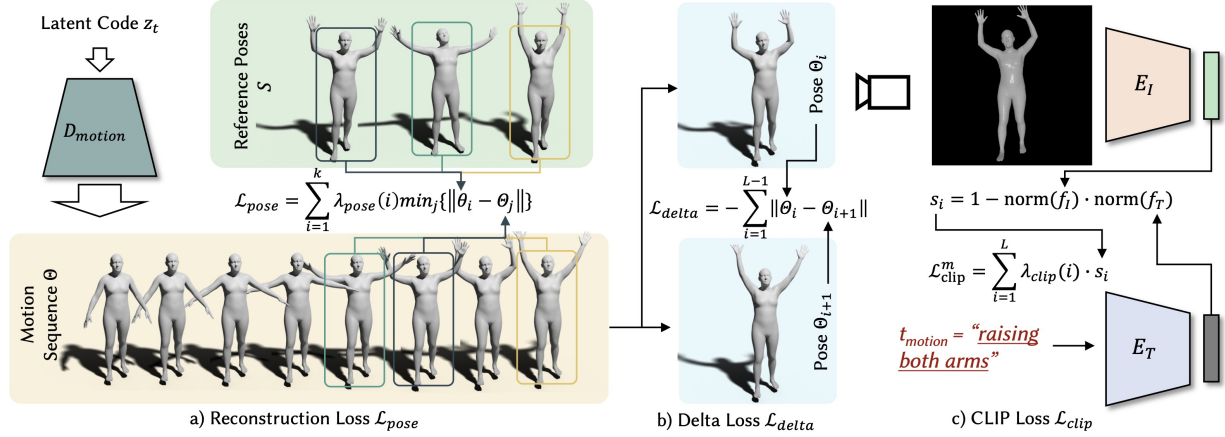


Fig. 10. **Detailed Pipeline of Reference-Based Animation with Motion Prior.** Three constraint terms are designed to optimize the latent code z_t . For the motion sequence Θ decoded by D_{motion} , $\mathcal{L}_{\text{pose}}$ is used to minimize the distance between each candidate pose and the nearest pose in Θ . $\mathcal{L}_{\text{delta}}$ is an adjustable loss item that measures the differences between adjacent poses and is capable of controlling the intensity of motion. $\mathcal{L}_{\text{clip}}^m$ measures the similarity between description t_{motion} and each pose in Θ .

$D_o(\text{TransDecoder}(D_p(z_{\text{motion}}), \phi(\mathbf{0})))$, where D_p and D_o are fully connected layers, $\mathbf{0}$ is a zero-vector. TransDecoder contains several transformer decoder layers. A loss function with two terms is proposed to train the motion VAE, which is defined as

$$\mathcal{L}_{\text{mVAE}} = \lambda_5 \cdot \mathcal{L}_{\text{KL}} + \mathcal{L}_{\text{recon}}, \quad (11)$$

where λ_5 is a hyper-parameter to balance two terms. \mathcal{L}_{KL} computes the KL-divergence to enforce the distribution assumption. $\mathcal{L}_{\text{recon}}$ is the mean squared error between Θ^+ and Θ^* for the reconstruction of the motion sequences.

Optimization on the Motion VAE. With the pre-trained motion VAE, we attempt to optimize its latent code z_t to synthesize a motion sequence $\Theta = D_{\text{motion}}(z_t)$. As shown in Fig. 10, three optimization constraints are proposed as

$$\mathcal{L}_{\text{motion}} = \mathcal{L}_{\text{pose}} + \lambda_6 \mathcal{L}_{\text{delta}} + \lambda_7 \mathcal{L}_{\text{clip}}^m, \quad (12)$$

where λ_6, λ_7 are hyper-parameters to balance the terms. $\mathcal{L}_{\text{pose}}$ is the reconstruction term between the decoded motion sequence Θ and candidate poses S . $\mathcal{L}_{\text{delta}}$ measures the range of the motion to prevent the motion from being overly-smoothed. $\mathcal{L}_{\text{clip}}^m$ encourages each single pose in the motion to match the input motion description. Details of three loss terms are introduced as follows.

Given reference poses $S = \{\theta_1, \theta_2, \dots, \theta_k\}$, the target is to construct a motion sequence that is close enough to these poses. We propose to minimize the distance between θ_i and its nearest frame Θ_j , where $i \in \{1, 2, \dots, k\}$, and $j \in \{1, 2, \dots, L\}$. Note that we assume θ_i is less similar to t_{motion} with larger i . Therefore, we use a coefficient $\lambda_{\text{pose}}(i) = 1 - \frac{i-1}{k}$ to focus on candidate poses with higher similarities. Formally, the reconstruction loss is defined as

$$\mathcal{L}_{\text{pose}} = \sum_{i=1}^k \lambda_{\text{pose}}(i) \min_j \{\|\theta_i - \Theta_j\|\}. \quad (13)$$

Empirically, only using the reconstruction loss $\mathcal{L}_{\text{pose}}$, the generated motion tends to be over-smoothed. To generate motions

with larger motion ranges, we design a motion range term $\mathcal{L}_{\text{delta}}$ to measure the smoothness of adjacent poses,

$$\mathcal{L}_{\text{delta}} = - \sum_{i=1}^{L-1} \|\theta_i - \theta_{i+1}\|, \quad (14)$$

which serves as a penalty term against over-smoothed motions. More intense motions will be generated when increasing λ_6 .

The matching scheme of the reconstruction term $\mathcal{L}_{\text{pose}}$ does not guarantee the ordering of candidate poses. Lack of supervision on the pose orderings would lead to unstable generation results. Furthermore, candidate poses might only contribute to a small part of the final motion sequence, which would lead to unexpected motion pieces. To tackle these two problems, we design an additional CLIP-guided loss term

$$\mathcal{L}_{\text{clip}}^m = \sum_{i=1}^L \lambda_{\text{clip}}(i) \cdot s_i, \quad (15)$$

where s_i is the similarity score between the pose θ_i and text description t_{motion} , which is defined as

$$s_i = 1 - \text{norm}(E_I(\mathcal{R}(\theta_i))) \cdot \text{norm}(E_T(t_{\text{motion}})). \quad (16)$$

$\lambda_{\text{clip}}(i) = \frac{i}{L}$ is a monotonically increasing function so that $\mathcal{L}_{\text{clip}}^m$ gives higher penalty to the later poses in the sequence. With the above CLIP-guided term, the whole motion sequence will be more consistent with t_{motion} . Empirically, we find that we only need to sample a small part of poses in Θ for the calculation of this term, which will speed up the optimization without observable degradation in performance.

4 EXPERIMENTS

4.1 Experimental Setup

4.1.1 Implementation Details. The shape VAE (Sec. 3.3.1) uses a two-layer MLP with a 16 dimension latent space for the encoder

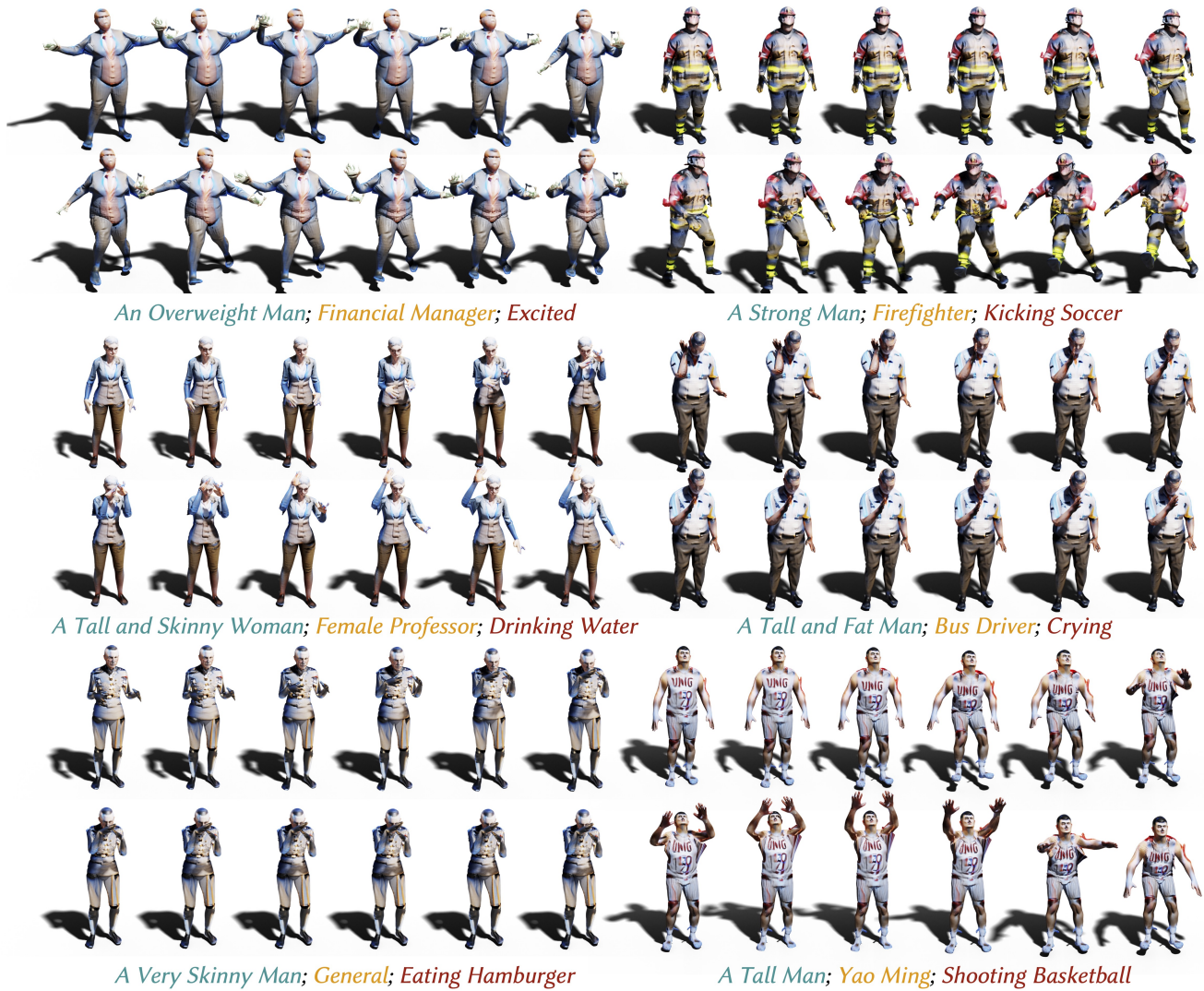


Fig. 11. **Overall Results of AvatarCLIP.** Renderings of several animated 3D avatars are shown in sequence. The corresponding driving texts for *shape*, *appearance* and *motion* are put below the sequences.

and decoder. Then 40,960 latent codes are sampled from a Gaussian distribution and clustered to 2,048 centroids by K-Means which composite the code-book. For the shape sculpting and texture generation (Sec. 3.3.2), we adjust the original NeuS such that the SDF network uses a 6-layer MLP and the color network uses a 4-layer MLP. For each ray, we perform 32 uniform samplings and 32 importance samplings. The Adam algorithm [Kingma and Ba 2014] with a learning rate of 5×10^{-4} is used for 30,000 iterations of optimization.

For candidate pose generation (Sec. 3.4.1), we directly use the pre-trained VPoser [Pavlakos et al. 2019] and use K-Means to acquire 4,096 cluster centroids from the AMASS dataset [Mahmood et al. 2019]. We select top-5 poses for the next stage (Sec.3.4.2). Our proposed Motion VAE has a 256 dimension latent space. The length

of the motion is 60. We train 100 epochs for motion VAE on AMASS dataset. An Adam optimizer with a 5×10^{-4} learning rate is used. During optimizing latent code for t_{motion} , Adam is used for 5,000 iterations of optimization with a 1×10^{-2} learning rate.

4.1.2 Baselines. Though it is the first work to generate and animate 3D avatars in a zero-shot text-driven manner, we design reasonable baseline methods for the evaluation of each part of AvatarCLIP. For the coarse shape generation, we design a baseline method where the shape parameters (*i.e.* the SMPL β and latent code of the shape VAE) are directly optimized by CLIP-guided losses. For the shape sculpting and texture generation, we compare our design with Text2Mesh [Michel et al. 2021]. Moreover, we introduce a

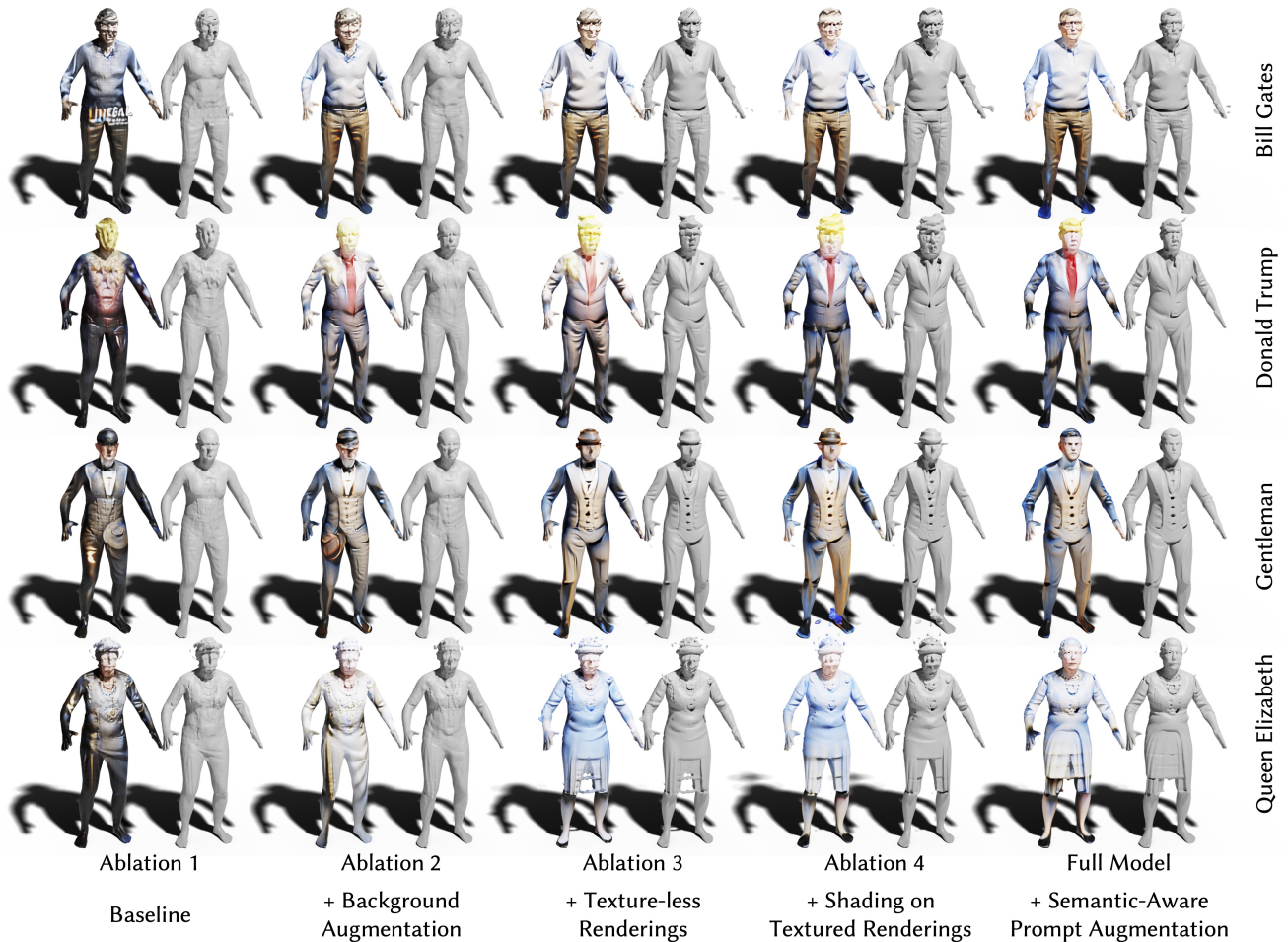


Fig. 12. **Ablation Study on Static Avatar Generation.** Four ablation studies are performed to validate our design choices in avatar generation. Specifically, four ablation settings subsequently add 1) background augmentation, 2) texture-less renderings, 3) shading on textured renderings, 4) semantic-aware prompt augmentation.

NeRF-based baseline by adapting Dream Fields [Jain et al. 2021a], where our ‘additional color network’ design is added to its pipeline to constraint the general shape of the avatar.

For the candidate pose generation, three baseline methods are designed to compare with our method. To illustrate the difficulty of direct optimization, we set two baselines that directly optimize on SMPL parameter θ and latent code z in VPoser. Moreover, inspired by CLIP-Forge [Sanghi et al. 2021], we use Real NVP [Dinh et al. 2016] to get a bi-projection between the normal distribution and latent space distribution of VPoser. This normalization flow network is conditioned on the CLIP features. This method does not need paired pose-text data. As for the second part of the motion generation, we design two baseline methods to compare with: Baseline (i) first sorts candidate poses S by their similarity scores s_i . Then, direct interpolations over the latent codes between each pair of adjacent poses are performed to generate the motion sequence. Baseline (ii)

uses the motion VAE to introduce motion priors into the generative pipeline. But (ii) directly calculates the reconstruction loss without using the re-weighting technique.

4.2 Overall Results

Overall results of the whole pipeline of AvatarCLIP are shown in Fig. 11. Avatars with diverse body shapes along with varied appearances are generated with high quality. They are driven by generated motion sequences that are reasonable and consistent with the input descriptions. In a zero-shot style, AvatarCLIP is capable of generating animatable avatars and motions, making use of the strong prior in pre-trained models. The whole process of avatar generation and animation, which originally requires expert knowledge of professional software, can now be simply driven by natural languages with the help of our proposed AvatarCLIP.

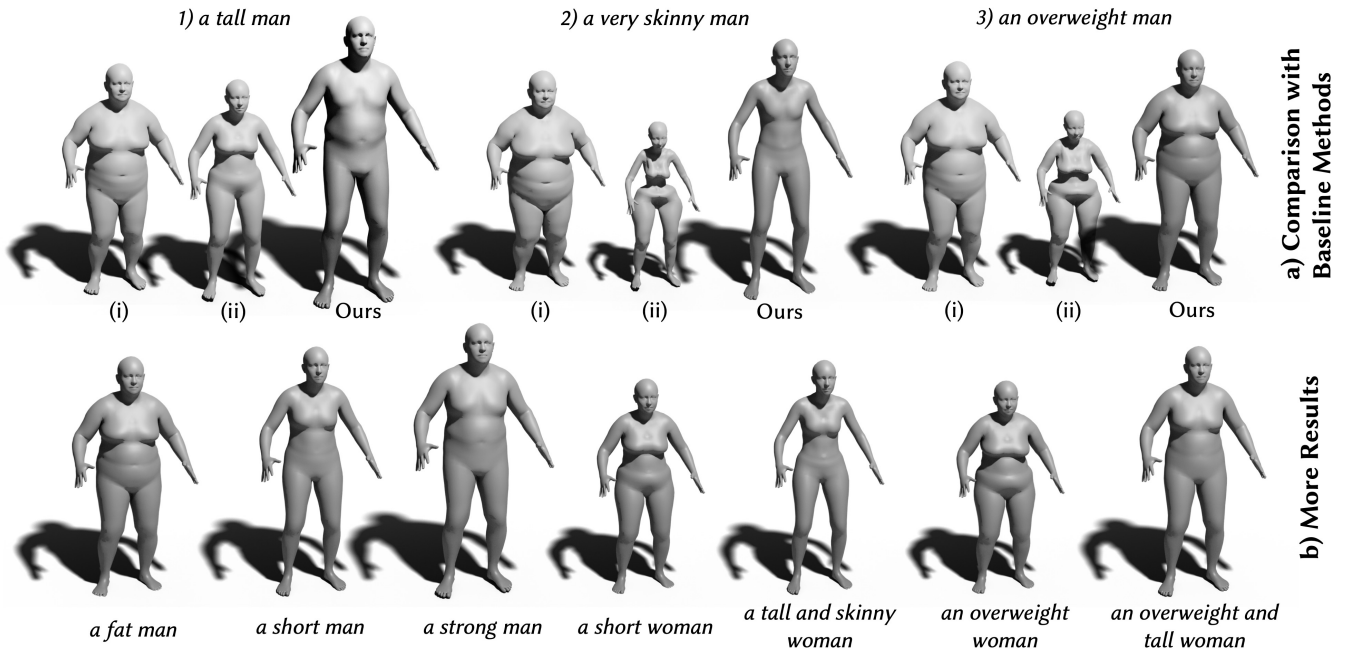


Fig. 13. **Results of the Coarse Shape Generation.** a) The qualitative comparison between our method and two baselines (i) direct optimization over SMPL β space, (ii) direct optimization over VAE latent space. Two baseline methods fail to generate reasonable body shapes. While our method successfully generates body shapes that match the descriptions above. b) More results of the coarse shape generation are shown.

4.3 Experiments on Avatar Generation

4.3.1 Ablation Study. To validate the effectiveness of various designs in the avatar generation module, we perform extensive ablation studies. We ablate the designs of 1) background augmentation; 2) supervision on texture-less renderings; 3) random shading on the textured renderings; 4) semantic-aware prompt augmentation. The ablation settings shown in Fig. 12 are formed by subsequently adding the above four designs to a baseline method where only textured renderings are supervised by CLIP. As shown in the first two columns of Fig. 12, background augmentation has a great influence on the texture generation, without which the textures tend to be very dark. Comparing the second and third columns, adding the supervision on texture-less renderings improves the geometry quality by a large margin. The geometry of ‘Ablation 2’ has lots of random bumps, which make the surfaces noisy. While the geometry of ‘Ablation 3’ is smooth and has detailed wrinkles of garments. As shown by the ‘Ablation 3’ and ‘Ablation 4’, adding random shadings on textured renderings helps the generation of more uniform textures. For example, the ‘Donald Trump’ of ‘Ablation 3’ has a brighter upper body than the lower one, which is improved in ‘Ablation 4’. Without the awareness of human body semantics, the previous four settings cannot generate correct faces for the avatars. The last column, which uses the semantic-aware prompt augmentation, has the best results in terms of the face generation.

4.3.2 Qualitative Results of Coarse Shape Generation. For this part, we design two intuitive baseline methods where direct CLIP supervision is back-propagated to the shape parameters. As shown in Fig. 13 (a), both optimization methods fail to generate body shapes consistent with description texts. Even opposite text guidance (e.g. ‘skinny’ and ‘overweight’) leads to the same optimization direction. In comparison, our method robustly generates reasonable body shapes agreeing with input texts. More diverse qualitative results of our method are shown in Fig. 13 (b).

4.3.3 Qualitative Results of Shape Sculpting and Texture Generation.

Broad Range of Driving Texts. Throughout extensive experiments, our method is capable of generating avatars from a wide range of appearances descriptions including three types: 1) celebrities, 2) fictional characters and 3) general words that describe people, as shown in Fig. 17. As shown in Fig. 17 (a), given celebrity names as the appearance description, the most iconic outfit of the celebrity is generated. Thanks to our design of semantic-aware prompt augmentation, the faces are also generated correctly. For the fictional character generation as illustrated in Fig. 17 (b), avatars of most text descriptions can be correctly generated. Interestingly, for the characters that have accessories with complex geometry (e.g. helmets of ‘Batman’, the dress of ‘Elsa’), the optimization process has the tendency of ‘growing’ new structures out of the template human body. As for the general descriptions, our method can handle very broad ranges including common job names (e.g. ‘Professor’, ‘Doctor’), words that describe people at a certain age (e.g. ‘Teenager’,

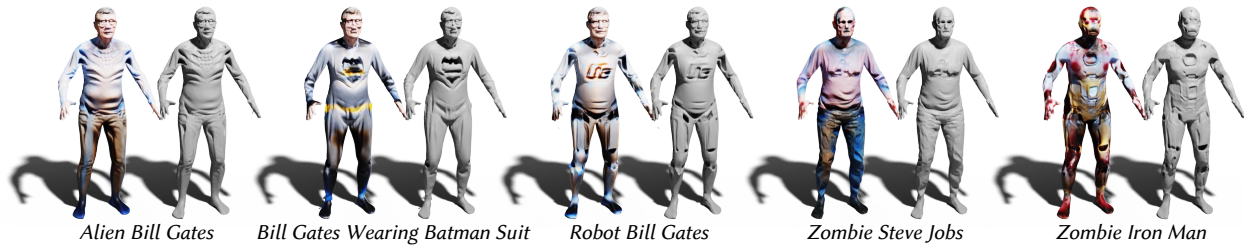


Fig. 14. **Results of Concept Mixing.** We show several concept mixing avatar generation results. Natural blending of different concepts is demonstrated, which further proves the generalizability of our method. The corresponding input texts are listed below the renderings.

‘Senior Citizen’) and other fantasy professions (e.g. ‘Witch’, ‘Wizard’). It can be observed that other than the avatars themselves, their respective iconic objects can also be generated. For example, the ‘Gardener’ grasps flowers and grasses in his hands.

Zero-Shot Controlling Ability. Other than the overall descriptions of the appearance as shown above, our method is also capable of zero-shot controlling at a more detailed level. As shown in Fig. 18, we can control the faces generated on the avatars, e.g. Bill Gates wearing an Iron Man suit, by tuning the semantic-aware prompt augmentation. Moreover, we can control the clothing of the avatar by direct text guidance, e.g. ‘Steve Jobs in white shirt’.

Concept Mixing. Inspired by DALL-E [Ramesh et al. 2021] and Dream Fields [Jain et al. 2021a], one of the most exciting applications of CLIP-driven generation is concept mixing. As shown in Fig. 14, we demonstrate examples of mixing fantasy elements with celebrities. While maintaining the recognizable identities, the fantasy elements blend in with the generated avatars naturally.

Geometry Quality. The critical design of texture-less rendering supervision mainly contributes to the geometry generation. We mainly compare our AvatarCLIP with the adapted Dream Field, which is based on NeRF. As shown in Fig. 19, our method consistently outperforms Dream Fields in terms of geometry quality. Detailed muscle shapes, armor curves, and cloth wrinkles can be generated.

Robustness. Other than the generation quality, we also investigate the robustness of our algorithm compared with the baseline method Text2Mesh. For each method, we use the same five random seeds for five independent runs for the same prompt. As shown in Fig. 20, our method manages to output results with high quality and consistency with the input text. Text2Mesh fails most runs, which shows that the representations of meshes are unstable for optimization, especially with weak supervision.

Failure Cases. Although a wide range of generation results are experimented with and demonstrated above, there exist failure cases in generating loose garments and accessories. For example, as shown in Fig. 17, the dress of ‘Elsa’ and the cloak of ‘Doctor Strange’ are not generated. The geometry of exaggerated hair and beard (e.g. the breaded Forrest Gump) is also challenging to be generated correctly.

This is caused by the reconstruction loss in the optimization process. The ‘growing’ is discouraged and very limited changes of the geometry are allowed.

4.3.4 Quantitative Results. To quantitatively evaluate the results of our avatar generation method, we ask 22 volunteers to perform a user study in terms of 1) the consistency with input texts, 2) texture quality, and 3) geometry quality. We randomly select 8 input texts that describe appearances. They are used for avatar generation by three methods, i.e. Dream Field, Text2Mesh and our AvatarCLIP. For each sample, the volunteers are asked to score the results of three methods from 1 to 5 in terms of the above three aspects. As illustrated in Fig. 15, our method consistently outperforms the other two baseline methods in all three aspects. Moreover, the standard deviations of our method are the lowest among the three methods, which also demonstrates the stable quality of our method.

4.4 Experiments on Motion Generation

4.4.1 Qualitative Results of Motion Generation.

Ablation Study. To evaluate the effectiveness of our design choices in the reference-based animation module (i.e. the proposed three constraint terms, the usage of motion VAE as a motion prior), we compare our method with two baseline methods as shown in Fig. 21. For ‘Brush Teeth’, (i) generates unordered and unrelated pose sequences. (ii) also fails to generate reasonable motions, which is caused by its blind focus on reconstructing the unordered candidate poses. By introducing a re-weighting mechanism, our method not only focuses on reconstruction but also considers the rationality of the generated motion. For ‘Kick Soccer’, the motion sequences from (i) (ii) have no drastic changes when the leg kicks out. $\mathcal{L}_{\text{delta}}$ plays a significant role here to control the intensity of motions. As for ‘Raise Both Arms’, it is supposed to generate a motion sequence from a neutral pose to a pose with raised arms. However, (i) generates a motion sequence that is contrary to the expected result. (ii) introduces several unrelated actions. With the help of $\mathcal{L}_{\text{clip}}^m$, our method is capable of generating motions with the correct ordering and better consistency with the descriptions.

Failure Cases. We also demonstrate some failure cases in Fig. 23. Throughout experiments, we find it hard to precisely control the body parts. For example, we cannot specifically control left or right

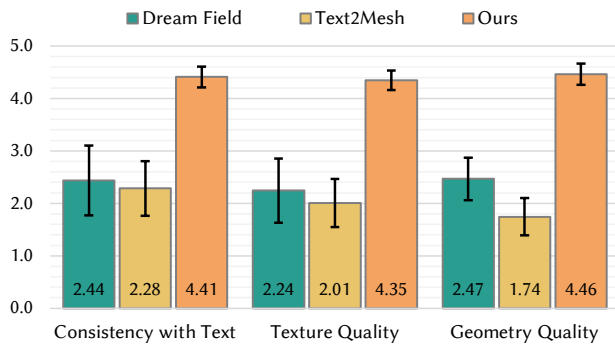


Fig. 15. **User Study on Static Avatar Generation** quantitatively shows the superiority of AvatarCLIP over other two baseline methods in three aspects: 1) consistency with text, 2) texture quality, and 3) geometry quality.

hands, as shown in ‘raising left arm’ case. Moreover, limited by the diversity of candidate poses, it is challenging to generate more complex motions like ‘hugging’, ‘playing piano’ and ‘dribbling’.

4.4.2 Qualitative Results of Candidate Pose Generation.

Comparison with Baseline Methods. Comparisons between different candidate pose generation methods are shown in Fig. 22 (a). Both direct optimization methods (i) (ii) fail to generate reasonable poses, let alone poses that are consistent with the given description. These results suggest that direct optimization on human poses parameters is intractable. Compared with (i) and (ii), conditioned Real NVP (iii) can yield rational poses. However, compared to the proposed solution based on the code-book, the generated poses from (iii) are less reasonable and of lower quality.

Broad Range of Driving Texts. To demonstrate the zero-shot ability of the pose generation method, we experiment with four categories of motion descriptions: 1) abstract emotion descriptions (e.g. ‘tired’ and ‘sad’); 2) common action descriptions (e.g. ‘walking’ and ‘squatting’); 3) descriptions of motions related to body parts (e.g. ‘raising both arms’, ‘washing hands’) 4) descriptions of motions that involve interaction with objects (e.g. ‘shooting basketball’). They are shown in Fig.22 (b).

4.4.3 Quantitative Results. Quantitatively, we evaluate the candidate pose generation and reference-based animation separately. For the candidate pose generation, we ask 58 volunteers to select the candidate poses that are the most consistent with the given text inputs for 15 randomly selected samples. The percentage of the selected times of each method for each text input is used as the scores for counting. As shown in Fig. 16 (a), compared with the baseline methods introduced in Sec. 4.4.2, our method outperforms them by a large margin. For the reference-based animation, we ask 20 volunteers to score the results from 1 to 5 in terms of the consistency with the input texts and the overall quality for 10 randomly selected samples. As shown in Fig. 16 (b), compared with the two baseline methods introduced in Sec. 4.4.1, our method outperforms them by large margins in both consistency and quality.

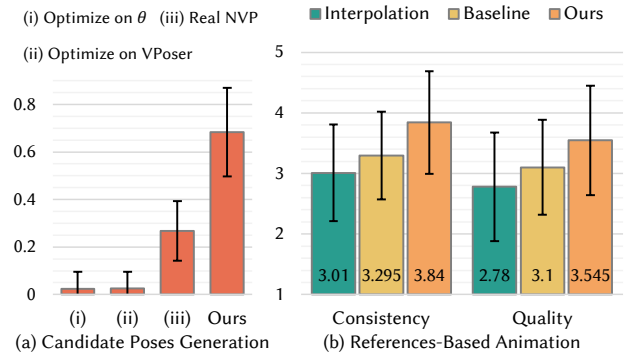


Fig. 16. **Results of User Study on Motion Generation.** For candidate poses generation, our method outperforms both direct optimization and sampling from multi-modality Real NVP. For animation, our method acquires higher scores than both baseline methods noticeably.

5 DISCUSSION

In this work, by proposing AvatarCLIP, we make the originally complex and demanding 3D avatar creation accessible to layman users in a text-driven style. It is made possible with the powerful priors provided by the pre-trained models including the shape/ motion VAE and the large-scale vision-language pre-trained model CLIP. Extensive experiments are conducted to validate the effectiveness of careful designs of our methods.

Limitations. For the avatar generation, limited by the weak supervision and low resolution of CLIP, the results are not perfect if zoomed in. Besides, it is hard to generate avatars with large variations given the same prompt. For the same prompt, the CLIP text feature is always the same. Therefore, the optimization directions are the same, which leads to similar results across different runs. For the motion synthesis, limited by the code-book design, it is hard to generate out-of-distribution candidate poses, which limits the ability to generate complex motions. Moreover, due to the lack of video CLIP, it is difficult to generate stylized motion.

Potential Negative Impacts. The usage of pre-trained models might induce ethical issues. For example, if we let t_{app} = ‘doctor’, the generated avatar is male. If t_{app} = ‘nurse’, the generated avatar is female, which demonstrates the gender bias. We think the problem originates from the large-scale internet data used for CLIP training, which might be biased if not carefully reviewed. Future works regarding the ethical issues of large-scale pre-trained models are required for the zero-shot techniques safe to be used. Moreover, with the democratization of producing avatars and animations, users can easily produce fake videos of celebrities, which might be misused and cause negative social impacts.

ACKNOWLEDGMENTS

This study is supported by NTU NAP, MOE AcRF Tier 2 (T2EP20221-0033), and under the RIE2020 Industry Alignment Fund – Industry Collaboration Projects (IAF-ICP) Funding Initiative, as well as cash and in-kind contribution from the industry partner(s).



Fig. 17. **Results of Static Avatar Generation.** A broad range of driving texts are tested, including celebrity names, frictional character names and general descriptions. The generated avatars are rendered with and without texture for the convenience of observing textures and geometry.



Fig. 18. **Results of More Detailed Controlling.** The left two examples demonstrate the first type of detailed controlling, where the semantic-aware prompt augmentation is adapted. The right three examples further show controls over clothes simply by specifying in the input texts.

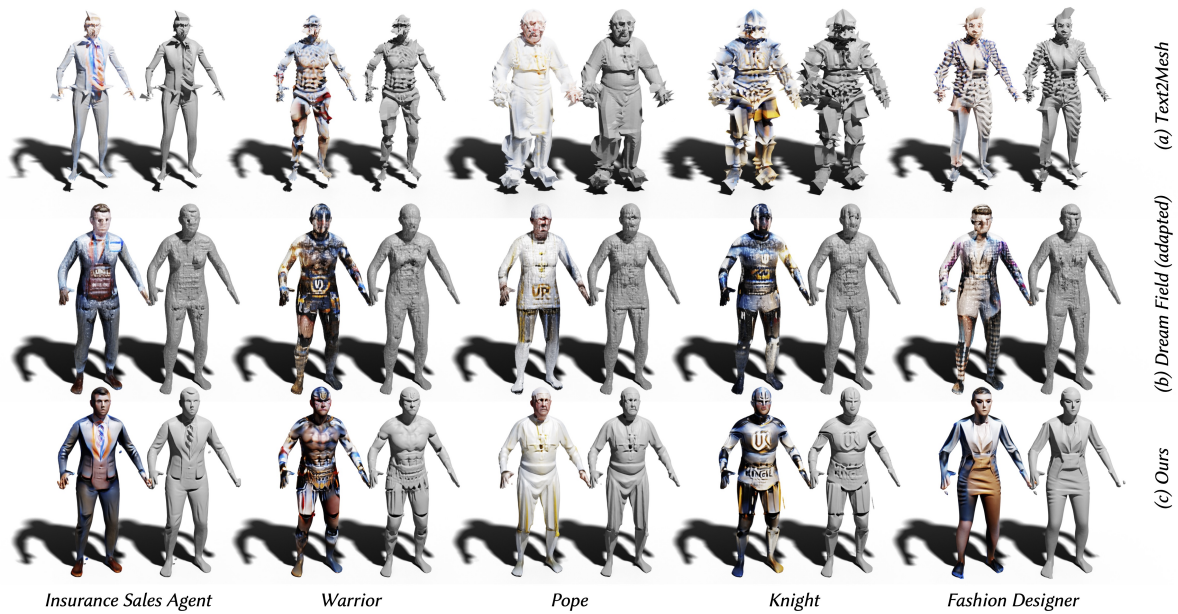


Fig. 19. **Qualitative Comparison with Baseline Methods.** Side-by-side comparisons between our method and (a) Text2Mesh [Michel et al. 2021] (the first line), (b) Dream Field [Jain et al. 2021a] (the second line) are demonstrated. Results of our method clearly show better quality in terms of both geometry and texture.

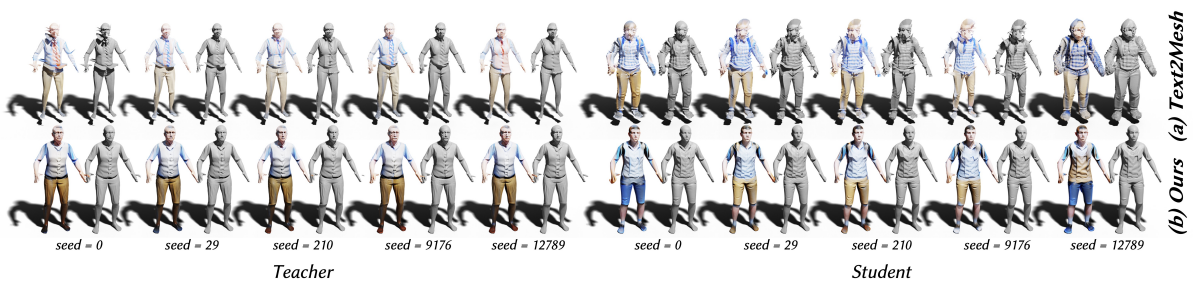


Fig. 20. **Results of Multiple Runs with Different Random Seeds.** By running the optimization multiple times with different random seeds, we observe that our method succeeds all runs, while Text2Mesh shows unstable results.

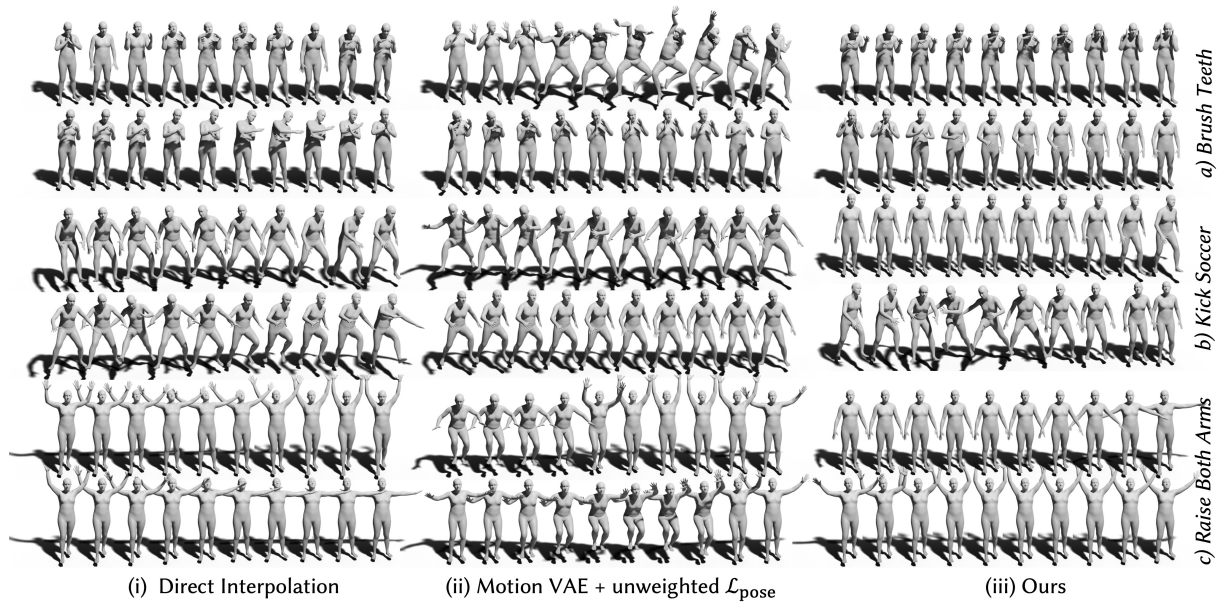


Fig. 21. **Qualitative Ablation of Motion Generation.** Compared with direct interpolation or solely using unweighted $\mathcal{L}_{\text{pose}}$, our proposed method can generate stable and reasonable motion sequences that are consistent with the given description.

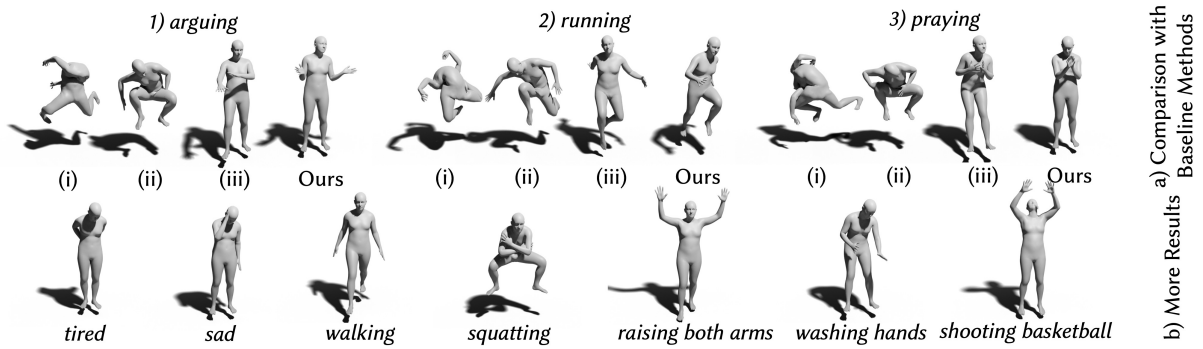


Fig. 22. **Results of the Candidate Pose Generation.** a) compares three baseline methods with ours. Baseline (i) directly optimizes on SMPL parameter θ . Baseline (ii) directly optimizes over the latent space of the VPoser. Both methods can hardly generate reasonable poses. Baseline (iii) utilizes a multi-modality Real NVP, which can generate relatively reasonable poses but still worse than our method by a clear margin. b) demonstrates more results of the candidate poses generation.

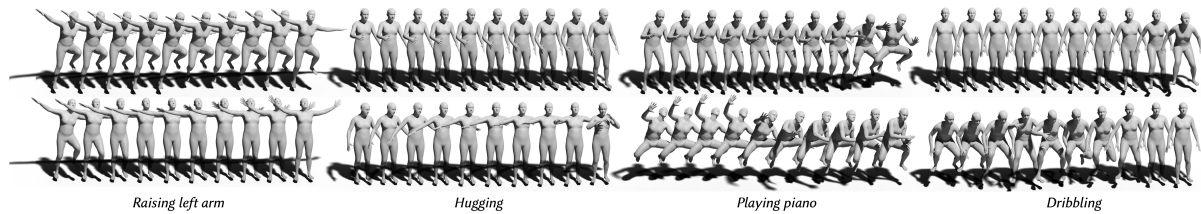


Fig. 23. **Failure Cases of Motion Generation.** More complex motions and more detailed controls over the motions are challenging to generate. Corresponding input texts are listed below.

REFERENCES

- Gunjan Aggarwal and Devi Parikh. 2021. Dance2Music: Automatic Dance-driven Music Generation. *arXiv preprint arXiv:2107.06252* (2021).
- Hyemin Ahn, Timothy Ha, Yunho Choi, Hwiyeon Yoo, and Songhwa Oh. 2018. Text2Action: Generative Adversarial Synthesis from Language to Action. In *2018 IEEE International Conference on Robotics and Automation, ICRA 2018, Brisbane, Australia, May 21–25, 2018*. IEEE, 1–5. <https://doi.org/10.1109/ICRA.2018.8460608>
- Chaitanya Ahuja and Louis-Philippe Morency. 2019. Language2Pose: Natural Language Grounded Pose Forecasting. In *2019 International Conference on 3D Vision (3DV)*, 719–728. <https://doi.org/10.1109/3DV.2019.00084>
- Thiemo Alldieck, Marcus Magnor, Weipeng Xu, Christian Theobalt, and Gerard Pons-Moll. 2018. Video based reconstruction of 3d people models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 8387–8397.
- Kevin Bergamin, Simon Clavet, Daniel Holden, and James Richard Forbes. 2019. DRCon: data-driven responsive control of physics-based characters. *ACM Transactions On Graphics (TOG)* 38, 6 (2019), 1–11.
- Bharat Lal Bhatnagar, Cristian Sminchisescu, Christian Theobalt, and Gerard Pons-Moll. 2020. Combining implicit function learning and parametric models for 3d human reconstruction. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II*. Springer, 311–329.
- Bharat Lal Bhatnagar, Garvita Tiwari, Christian Theobalt, and Gerard Pons-Moll. 2019. Multi-garment net: Learning to dress 3d people from images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5420–5430.
- Andrei Burov, Matthias Nießner, and Justus Thies. 2021. Dynamic surface function networks for clothed human bodies. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 10754–10764.
- Zhongang Cai, Daxuan Ren, Ailing Zeng, Zhengyu Lin, Tao Yu, Wenjia Wang, Xiangyu Fan, Yang Gao, Yifan Yu, Liang Pan, Fangzhou Hong, Mingyuan Zhang, Chen Change Loy, Lei Yang, and Ziwei Liu. 2022. HuMMan: Multi-Modal 4D Human Dataset for Versatile Sensing and Modeling. *arXiv preprint arXiv:2204.13686* (2022).
- Zhongang Cai, Mingyuan Zhang, Jiawei Ren, Chen Wei, Daxuan Ren, Jiatong Li, Zhengyu Lin, Haiyu Zhao, Shuai Yi, Lei Yang, et al. 2021. Playing for 3D Human Recovery. *arXiv preprint arXiv:2110.07588* (2021).
- Eric R Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. 2021. pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5799–5809.
- Xu Chen, Tianjian Jiang, Jie Song, Jinlong Yang, Michael J Black, Andreas Geiger, and Otmarr Hilliges. 2022. gDNA: Towards Generative Detailed Neural Avatars. *arXiv* (2022).
- Enric Corona, Albert Pumarola, Guillem Alenya, Gerard Pons-Moll, and Francesc Moreno-Noguer. 2021. SMPLicit: Topology-aware generative model for clothed people. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 11875–11885.
- Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. 2016. Density estimation using real nvp. *arXiv preprint arXiv:1605.08803* (2016).
- Kevin Frans, LB Soros, and Olaf Witkowski. 2021. Clipdraw: Exploring text-to-drawing synthesis through language-image encoders. *arXiv preprint arXiv:2106.14843* (2021).
- Jianglin Fu, Shikai Li, Yuming Jiang, Kwan-Yee Lin, Chen Qian, Chen-Change Loy, Wayne Wu, and Ziwei Liu. 2022. StyleGAN-Human: A Data-Centric Odyssey of Human Generation. *arXiv preprint arXiv:2204.11823* (2022).
- Anindita Ghosh, Noshaba Cheema, Cennet Oguz, Christian Theobalt, and Philipp Slusallek. 2021. Text-Based Motion Synthesis with a Hierarchical Two-Stream RNN. In *ACM SIGGRAPH 2021 Posters*. 1–2.
- Artur Grigorev, Karim Isakov, Anastasia Ianina, Renat Bashirov, Ilya Zakharkin, Alexander Vakhitov, and Victor Lempitsky. 2021. StylePeople: A Generative Model of Fullbody Human Avatars. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5151–5160.
- Amos Gropp, Lior Yariv, Niv Haim, Matan Atzmon, and Yaron Lipman. 2020. Implicit geometric regularization for learning shapes. *arXiv preprint arXiv:2002.10099* (2020).
- Chuan Guo, Xinxin Zuo, Sen Wang, Shihao Zou, Qingyao Sun, Annan Deng, Minglun Gong, and Li Cheng. 2020. Action2motion: Conditioned generation of 3d human motions. In *Proceedings of the 28th ACM International Conference on Multimedia*. 2021–2029.
- Marc Habermann, Lingjie Liu, Weipeng Xu, Michael Zollhoefer, Gerard Pons-Moll, and Christian Theobalt. 2021. Real-time deep dynamic characters. *ACM Transactions on Graphics (TOG)* 40, 4 (2021), 1–16.
- Xintong Han, Zuxuan Wu, Zhe Wu, Ruichi Yu, and Larry S Davis. 2018. Viton: An image-based virtual try-on network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 7543–7552.
- Fangzhou Hong, Liang Pan, Zhongang Cai, and Ziwei Liu. 2021. Garment4D: Garment Reconstruction from Point Cloud Sequences. In *Advances in Neural Information Processing Systems*.
- Fangzhou Hong, Liang Pan, Zhongang Cai, and Ziwei Liu. 2022. Versatile Multi-Modal Pre-Training for Human-Centric Perception. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- Zeng Huang, Yuanlu Xu, Christoph Lassner, Hao Li, and Tony Tung. 2020. Arch: Animatable reconstruction of clothed humans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 3093–3102.
- Leslie Ikenoto, Okan Arikan, and David Forsyth. 2009. Generalizing Motion Edits with Gaussian Processes. *ACM Trans. Graph.* 28, 1, Article 1 (feb 2009), 12 pages. <https://doi.org/10.1145/1477926.1477927>
- Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. 2013. Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *IEEE transactions on pattern analysis and machine intelligence* 36, 7 (2013), 1325–1339.
- Ajay Jain, Ben Mildenhall, Jonathan T Barron, Pieter Abbeel, and Ben Poole. 2021a. Zero-Shot Text-Guided Object Generation with Dream Fields. *arXiv preprint arXiv:2112.01455* (2021).
- Ajay Jain, Matthew Tancik, and Pieter Abbeel. 2021b. Putting nerf on a diet: Semantically consistent few-shot view synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5885–5894.
- Nikolay Jetchev. 2021. ClipMatrix: Text-controlled Creation of 3D Textured Meshes. *arXiv preprint arXiv:2109.12922* (2021).
- Boyi Jiang, Juyong Zhang, Yang Hong, Jinhao Luo, Ligang Liu, and Hujun Bao. 2020. Bcnet: Learning body and cloth shape from a single image. In *European Conference on Computer Vision*. Springer, 18–35.
- Yuming Jiang, Ziqi Huang, Xingang Pan, Chen Change Loy, and Ziwei Liu. 2021. Talk-to-Edit: Fine-Grained Facial Editing via Dialog. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 13799–13808.
- Yuming Jiang, Shuai Yang, Haonan Qiu, Wayne Wu, Chen Change Loy, and Ziwei Liu. 2022. Text2Human: Text-Driven Controllable Human Image Generation. *ACM Transactions on Graphics (TOG)* 41, 4 (2022), 1–11.
- Kacper Kania, Marek Kowalski, and Tomasz Trzcinski. 2021. TrajeVAE—Controllable Human Motion Generation from Trajectories. *arXiv preprint arXiv:2104.00351* (2021).
- Tero Karras, Samuli Laine, and Timo Aila. 2019. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4401–4410.
- Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. 2020. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8110–8119.
- Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- Diederik P Kingma and Max Welling. 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114* (2013).
- Kathleen M Lewis, Srivatsan Varadharajan, and Ira Kemelmacher-Shlizerman. 2021. Tryongan: Body-aware try-on via layered interpolation. *ACM Transactions on Graphics (TOG)* 40, 4 (2021), 1–10.
- Ruilong Li, Shan Yang, David A Ross, and Angjoo Kanazawa. 2021. Ai choreographer: Music conditioned 3d dance generation with aist++. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 13401–13412.
- Hung Yu Ling, Fabio Zinno, George Cheng, and Michiel Van De Panne. 2020. Character controllers using motion vaes. *ACM Transactions on Graphics (TOG)* 39, 4 (2020), 40–1.
- Lingjie Liu, Marc Habermann, Viktor Rudnev, Kripasindhu Sarkar, Jiatong Gu, and Christian Theobalt. 2021. Neural actor: Neural free-view synthesis of human actors with pose control. *ACM Transactions on Graphics (TOG)* 40, 6 (2021), 1–16.
- Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. 2015. Deep Learning Face Attributes in the Wild. In *Proceedings of International Conference on Computer Vision (ICCV)*.
- Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. 2015. SMPL: A skinned multi-person linear model. *ACM transactions on graphics (TOG)* 34, 6 (2015), 1–16.
- William E Lorensen and Harvey E Cline. 1987. Marching cubes: A high resolution 3D surface construction algorithm. *ACM siggraph computer graphics* 21, 4 (1987), 163–169.
- Naureen Mahmood, Nima Ghorbani, Nikolaus F Troje, Gerard Pons-Moll, and Michael J Black. 2019. AMASS: Archive of motion capture as surface shapes. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5442–5451.
- Elman Mansimov, Emilio Parisotto, Jimmy Lei Ba, and Ruslan Salakhutdinov. 2015. Generating images from captions with attention. *arXiv preprint arXiv:1511.02793* (2015).
- Dushyant Mehta, Helge Rhodin, Dan Casas, Pascal Fua, Aleksandr Sotnychenko, Weipeng Xu, and Christian Theobalt. 2017. Monocular 3d human pose estimation in the wild using improved cnn supervision. In *2017 international conference on 3D vision (3DV)*. IEEE, 506–516.
- Oscar Michel, Roi Bar-On, Richard Liu, Sagie Benaim, and Rana Hanocka. 2021. Text2Mesh: Text-Driven Neural Stylization for Meshes. *arXiv preprint arXiv:2112.03221* (2021).
- Marko Mihajlovic, Yan Zhang, Michael J Black, and Siyu Tang. 2021. LEAP: Learning articulated occupancy of people. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10461–10471.

- 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 *European conference on computer vision*. Springer, 405–421.
- Tomohiko Mukai and Shigeru Kuriyama. 2005. Geostatistical Motion Interpolation. *ACM Trans. Graph.* 24, 3 (jul 2005), 1062–1070. <https://doi.org/10.1145/1073204.1073313>
- Atsuhiko Noguchi, Xiao Sun, Stephen Lin, and Tatsuya Harada. 2021. Neural articulated radiance field. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5762–5772.
- Atsuhiko Noguchi, Xiao Sun, Stephen Lin, and Tatsuya Harada. 2022. Unsupervised Learning of Efficient Geometry-Aware Neural Articulated Representations. *arXiv preprint arXiv:2204.08839* (2022).
- Pablo Palafox, Aljaž Božič, Justus Thies, Matthias Nießner, and Angela Dai. 2021. NPMs: Neural Parametric Models for 3D Deformable Shapes. *arXiv preprint arXiv:2104.00702* (2021).
- Or Patashnik, Zongze Wu, Eli Shechtman, Daniel Cohen-Or, and Dani Lischinski. 2021. Styleclip: Text-driven manipulation of stylegan imagery. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2085–2094.
- Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed AA Osman, Dimitrios Tzionas, and Michael J Black. 2019. Expressive body capture: 3d hands, face, and body from a single image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10975–10985.
- Sida Peng, Junting Dong, Qianqian Wang, Shangzhan Zhang, Qing Shuai, Hujun Bao, and Xiaowei Zhou. 2021a. Animatable neural radiance fields for human body modeling. *arXiv e-prints* (2021), arXiv–2105.
- Sida Peng, Yuanqing Zhang, Yinghao Xu, Qianqian Wang, Qing Shuai, Hujun Bao, and Xiaowei Zhou. 2021b. Neural body: Implicit neural representations with structured latent codes for novel view synthesis of dynamic humans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 9054–9063.
- Mathis Petrovich, Michael J Black, and Gül Varol. 2021. Action-Conditioned 3D Human Motion Synthesis with Transformer VAE. *arXiv preprint arXiv:2104.05670* (2021).
- Leonid Pishchulin, Stefanie Wuhrer, Thomas Helten, Christian Theobalt, and Bernt Schiele. 2017. Building Statistical Shape Spaces for 3D Human Modeling. *Pattern Recognition* (2017).
- 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. *arXiv preprint arXiv:2103.00020* (2021).
- Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. 2021. Zero-shot text-to-image generation. *arXiv preprint arXiv:2102.12092* (2021).
- Scott Reed, Zeynep Akata, Xinchun Yan, Lajanugen Logeswaran, Bernt Schiele, and Honglak Lee. 2016. Generative adversarial text to image synthesis. In *International Conference on Machine Learning*. PMLR, 1060–1069.
- Charles Rose, Michael F. Cohen, and Bobby Bodenheimer. 1998. Verbs and Adverbs: Multidimensional Motion Interpolation. *IEEE Comput. Graph. Appl.* 18, 5 (sep 1998), 32–40. <https://doi.org/10.1109/38.708559>
- Shunsuke Saito, Jinlong Yang, Qianli Ma, and Michael J Black. 2021. SCANimate: Weakly supervised learning of skinned clothed avatar networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2886–2897.
- Aditya Sanghi, Hang Chu, Joseph G Lambourne, Ye Wang, Chin-Yi Cheng, and Marco Fumero. 2021. Clip-forge: Towards zero-shot text-to-shape generation. *arXiv preprint arXiv:2110.02624* (2021).
- Kripasindhu Sarkar, Vladislav Golyanik, Lingjie Liu, and Christian Theobalt. 2021a. Style and pose control for image synthesis of humans from a single monocular view. *arXiv preprint arXiv:2102.11263* (2021).
- Kripasindhu Sarkar, Lingjie Liu, Vladislav Golyanik, and Christian Theobalt. 2021b. HumanGAN: A Generative Model of Humans Images. *arXiv preprint arXiv:2103.06902* (2021).
- Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. 2019a. Animating arbitrary objects via deep motion transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2377–2386.
- Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. 2019b. First order motion model for image animation. *Advances in Neural Information Processing Systems* 32 (2019), 7137–7147.
- Guy Tevet, Brian Gordon, Amir Hertz, Amit H Bermanto, and Daniel Cohen-Or. 2022. MotionCLIP: Exposing Human Motion Generation to CLIP Space. *arXiv preprint arXiv:2203.08063* (2022).
- Ayush Tewari, Justus Thies, Ben Mildenhall, Pratul Srinivasan, Edgar Tretschk, Yifan Wang, Christoph Lassner, Vincent Sitzmann, Ricardo Martin-Brualla, Stephen Lombardi, et al. 2021. Advances in neural rendering. *arXiv preprint arXiv:2111.05849* (2021).
- Gül Varol, Javier Romero, Xavier Martin, Naureen Mahmood, Michael J. Black, Ivan Laptev, and Cordelia Schmid. 2017. Learning from Synthetic Humans. In *CVPR*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*. 5998–6008.
- Timo von Marcard, Roberto Henschel, Michael Black, Bodo Rosenhahn, and Gerard Pons-Moll. 2018. Recovering Accurate 3D Human Pose in The Wild Using IMUs and a Moving Camera. In *European Conference on Computer Vision (ECCV)*.
- Can Wang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao. 2021a. CLIP-NeRF: Text-and-Image Driven Manipulation of Neural Radiance Fields. *arXiv preprint arXiv:2112.05139* (2021).
- Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. 2021b. NeuS: Learning Neural Implicit Surfaces by Volume Rendering for Multi-view Reconstruction. *arXiv preprint arXiv:2106.10689* (2021).
- Chung-Yi Weng, Brian Curless, and Ira Kemelmacher-Shlizerman. 2019. Photo wake-up: 3d character animation from a single photo. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5908–5917.
- Jungdam Won and Jehee Lee. 2019. Learning body shape variation in physics-based characters. *ACM Transactions on Graphics (TOG)* 38, 6 (2019), 1–12.
- Mohamed R. Amer Xiao Lin. 2014. Human Motion Modeling using DVGANs. *arXiv preprint arXiv:1804.10652* (2014).
- Hongyi Xu, Thiemo Alldieck, and Cristian Sminchisescu. 2021a. H-nerf: Neural radiance fields for rendering and temporal reconstruction of humans in motion. *Advances in Neural Information Processing Systems* 34 (2021).
- Mengde Xu, Zheng Zhang, Fangyun Wei, Yutong Lin, Yue Cao, Han Hu, and Xiang Bai. 2021b. A Simple Baseline for Zero-shot Semantic Segmentation with Pre-trained Vision-language Model. *arXiv preprint arXiv:2112.14757* (2021).
- Jae Shin Yoon, Lingjie Liu, Vladislav Golyanik, Kripasindhu Sarkar, Hyun Soo Park, and Christian Theobalt. 2021. Pose-guided human animation from a single image in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 15039–15048.
- Jichao Zhang, Enver Sangineto, Hao Tang, Aliaksandr Siarohin, Zhun Zhong, Nicu Sebe, and Wei Wang. 2021. 3D-Aware Semantic-Guided Generative Model for Human Synthesis. *arXiv preprint arXiv:2112.01422* (2021).
- Fuqiang Zhao, Wei Yang, Jiakai Zhang, Pei Lin, Yingliang Zhang, Jingyi Yu, and Lan Xu. 2021. HumanNeRF: Generalizable Neural Human Radiance Field from Sparse Inputs. *arXiv preprint arXiv:2112.02789* (2021).
- Yi Zhou, Connelly Barnes, Jingwan Lu, Jimei Yang, and Hao Li. 2019. On the continuity of rotation representations in neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5745–5753.