

AVATARGO: ZERO-SHOT 4D HUMAN-OBJECT INTERACTION GENERATION AND ANIMATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Recent advancements in diffusion models have led to significant improvements in the generation and animation of 4D full-body human-object interactions (HOI). Nevertheless, existing methods primarily focus on SMPL-based motion generation, which is limited by the scarcity of realistic large-scale interaction data. This constraint affects their ability to create everyday HOI scenes. This paper addresses this challenge using a zero-shot approach with a pre-trained diffusion model. Despite this potential, achieving our goals is difficult due to the diffusion model’s lack of understanding of “*where*” and “*how*” objects interact with the human body. To tackle these issues, we introduce **AvatarGO**, a novel framework designed to generate animatable 4D HOI scenes directly from textual inputs. Specifically, **1**) for the “*where*” challenge, we propose **LLM-guided contact retargeting**, which employs Lang-SAM to identify the contact body part from text prompts, ensuring precise representation of human-object spatial relations. **2**) For the “*how*” challenge, we introduce **correspondence-aware motion optimization** that constructs motion fields for both human and object models using the linear blend skinning function from SMPL-X. Our framework not only generates coherent compositional motions, but also exhibits greater robustness in handling penetration issues. Extensive experiments with existing methods validate AvatarGO’s superior generation and animation capabilities on a variety of human-object pairs and diverse poses. As the first attempt to synthesize 4D avatars with object interactions, we hope AvatarGO could open new doors for human-centric 4D content creation.

1 INTRODUCTION

The creation of 4D human-object interaction (HOI) holds immense significance across a wide range of industries, including augmented/virtual reality (AR/VR) and game development, as it forms the foundation of the 4D virtual world. Traditionally, developing such models has required extensive human effort and specialized engineering expertise. Fortunately, thanks to the collections of HOI datasets (Li et al., 2023b; Bhatnagar et al., 2022; Jiang et al., 2023a) and the recent advancements in diffusion models (Saharia et al., 2022; Ramesh et al., 2022; Balaji et al., 2022; Stability.AI, 2022; 2023), existing HOI generative techniques (Zhang et al., 2022; 2023; 2024; Shafir et al., 2023; Kapon et al., 2024; Chen et al., 2024a) have exhibited promising capabilities by generating 4D human motions with object interactions from textual inputs. Nonetheless, these methods primarily focus on SMPL-based (Loper et al., 2015; Pavlakos et al., 2019) motion generation, which struggles to capture the realistic appearance of subjects encountered in everyday life. Although InterDreamer (Xu et al., 2024b) has recently proposed to generate text-aligned 4D HOI sequences in a zero-shot manner, their output is still largely constrained by the SMPL model. This highlights a pressing need for more realistic and generalizable methods tailored specifically to model 4D human-object interactive content. We take the initiative and showcase the potential of addressing this challenge by leveraging the 3D generative methods in a zero-shot manner.

In recent times, 3D generative methods (Poole et al., 2022; Tang et al., 2023; Liu et al., 2023c; Lin et al., 2023; Wang et al., 2023d; Cao et al., 2023b; Liao et al., 2023) and Large Language Models (LLMs) (Wu et al., 2023a) have garnered increasing interest. These progressives have led to the development of text-guided 3D compositional generation techniques that are capable of comprehending intricate relations and creating complex 3D scenes incorporating multiple subjects. Notably, GraphDreamer (Gao et al., 2023) utilizes LLMs to construct a graph where nodes represent



081 **Figure 1: Examples of 4D animation results obtained via AvatarGO.** AvatarGO effectively
082 produces diverse human-object compositions with correct spatial correlations and contact areas. It
083 achieves joint animation of humans and objects while avoiding penetration issues.

084
085
086 objects and edges denote their relations. ComboVerse (Chen et al., 2024b) proposes spatial-aware
087 score distillation sampling to amplify the spatial correlation. Subsequent studies (Epstein et al., 2024;
088 Zhou et al., 2024) further explore the potential of jointly optimizing layouts to composite different
089 components.

090 Despite the promising performance demonstrated by existing methods, they encounter two major
091 challenges in generating 4D HOI scenes: 1) *Incorrect contact area*: While LLMs excel at capturing
092 the relationships, optimization with diffusion models faces difficulties in accurately defining the
093 contact area between various objects, particularly those with complex articulated structures like
094 human bodies. Although efforts like InterFusion (Dai et al., 2024) have constructed 2D human-object
095 interaction datasets to retrieve human poses from text prompts, they still encounter challenges in
096 defining the optimal contact body parts for cases outside the training distribution. 2) *Limitations in*
097 *4D compositional animation*: While existing techniques like DreamGaussian4D (Ren et al., 2023)
098 and TC4D (Bahmani et al., 2024) employ video diffusion models (Blattmann et al., 2023; Guo et al.,
099 2023a) to animate 3D static scenes, they often treat the entire scene as one subject during optimization,
100 leading to unrealistic animation results. Despite initiatives like Comp4D (Xu et al., 2024a), which
101 utilize trajectories to animate 3D objects individually, modeling contact between various subjects
remains a challenge.

102 In this paper, we propose **AvatarGO**, a novel framework for compositional 4D avatar generation with
103 object interactions. By taking the text prompts as inputs, we assume that the 3D human and object
104 models as well as the human motion sequences can be individually generated by adopting existing
105 generative techniques (Tang et al., 2023; Liu et al., 2023d; Zhang et al., 2023; 2024). Specifically,
106 we adopt DreamGaussian4D (Ren et al., 2023) as our baseline considering its superior training efficiency
107 and focus on addressing the challenges associated with human-object interactions. To achieve this
objective, AvatarGO integrates two key innovations to learn “where” and “how” the object should

interact with the human body: **1) LLM-guided contact retargeting.** Given the limited availability of human-object interaction images in the 2D dataset used for diffusion model training, it’s difficult to identify the most appropriate contact area between humans and objects. To tackle this issue, we propose leveraging Lang-SAM (lan, 2023) to identify the contact body part from text prompts, which serves as the initialization for the optimization procedure. **2) Correspondence-aware motion optimization.** Building upon the observation that penetration is absent in static composited models, we introduce correspondence-aware motion optimization that leverages SMPL-X as an intermediary to maintain the correspondence between humans and objects when they are animated to a new pose, thus demonstrating greater robustness in handling penetration issues.

We thoroughly assess AvatarGO by compositing diverse pairs of 3D humans and objects and animating them across various motion sequences (see Fig. 1). Our experimental results show that our method excels at identifying optimal contact areas and exhibits greater robustness in handling penetration issues during animation, significantly outperforming existing techniques. We will make our code publicly available.

2 RELATED WORK

3D Content Generation. Leveraging advances in diffusion-based text-to-2D image generation (Saharia et al., 2022; Ramesh et al., 2022; Balaji et al., 2022; Stability.AI, 2022; 2023), DreamFusion introduced Score Distillation Sampling (SDS) to generate 3D content via pre-trained models, utilizing technologies like NeRF (Mildenhall et al., 2020), DMTET (Shen et al., 2021), 3D Gaussian Splatting (Kerbl et al., 2023)). Subsequent research has focused on enhancing output quality (Lin et al., 2023; Chen et al., 2023b; Wang et al., 2023d), controlling generation processes (Metzer et al., 2022; Seo et al., 2023), improving training efficiency (Wang et al., 2023a; Wu et al., 2024; Tang et al., 2023), and extending capabilities on 3D texturing (Richardson et al., 2023; Cao et al., 2023a; Chen et al., 2023a; Tang et al., 2024b). Addressing 3D human body complexity, recent studies (Cao et al., 2023b;c; Liao et al., 2023; Jiang et al., 2023b; Huang et al., 2023b; Kolotouros et al., 2023; Zeng et al., 2023; Huang et al., 2023a) have been proposed for creating controllable 3D human avatars, although these still require significant input-specific training time. The proliferation of large 3D datasets (Deitke et al., 2023; 2024; Wu et al., 2023b) has propelled 3D generation techniques forward. Notably, Zero-1-to-3 (Liu et al., 2023c), Zero123++ (Shi et al., 2023a), and MVDream (Shi et al., 2023b) use 2D diffusion models to generate consistent multi-view images, serving as inputs for efficient 3D model generation tools like SyncDreamer (Liu et al., 2023e), Wonder3D (Long et al., 2023), One-2-3-45 (Liu et al., 2023b;a), UniDream (Liu et al., 2023f), MVDiffusion++ (Tang et al., 2024c), and Make-Your-3D (Liu et al., 2024). Additionally, building on transformer (Vaswani et al., 2017) and image processor advancements (e.g., DINO (Caron et al., 2021; Oquab et al., 2023)), Large Reconstruction Models (Hong et al., 2023; Wang et al., 2023b; Xu et al., 2023; Li et al., 2023a) implement transformer-based architectures to derive 3D tri-plane tokens from image features. 3DTopia (Hong et al., 2024) uses hybrid diffusion priors to produce high-fidelity 3D objects. Meanwhile, methods like LGM (Tang et al., 2024a), CRM (Wang et al., 2024), and GRM (Yinghao et al., 2024) explore various 3D representations for improved performance, such as 3D Gaussian Splatting (Kerbl et al., 2023) and FlexiCube (Shen et al., 2023). Despite these advances, challenges remain in generating complex compositional 3D scenes.

3D Compositional Generation. To address the compositional nature of 3D content, a few efforts have been made recently. *Epstein et al* (Epstein et al., 2024) and GALA3D (Zhou et al., 2024) propose optimizing component layouts for integrated object scenes. ComboVerse (Chen et al., 2024b) introduces spatial-aware score distillation sampling (SSDS) to effectively learn object spatial relations. GraphDreamer (Gao et al., 2023) uses large language models to form graph structures where nodes and edges represent objects and their relationships, respectively, showing promising results. Challenges remain in modeling interactions between humans and objects. InterFusion (Dai et al., 2024) develops a 2D dataset for human-object interactions, enabling text-guided pose retrieval and scene generation. However, this approach lacks precise control over interaction areas and is not readily adaptable to 4D scenarios.

4D Content Generation. Recent advances in video diffusion models and score distillation sampling have spurred a variety of 4D scene generation techniques. Make-A-Video3D (MAV3D) (Singer et al., 2023) utilizes HexPlane features for 4D representations. 4D-fy (Bahmani et al., 2023) and

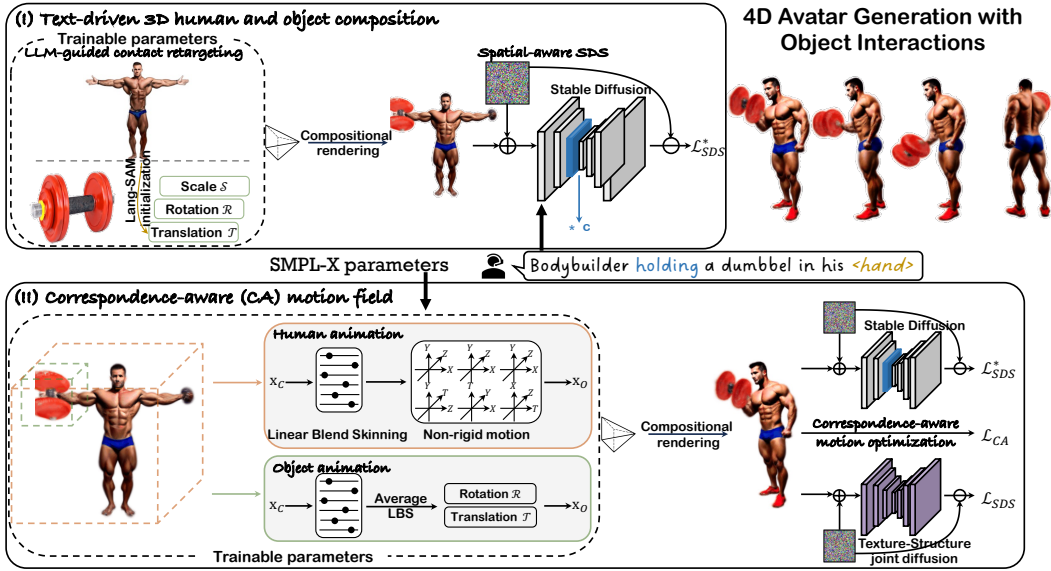


Figure 2: **Overview of AvatarGO.** AvatarGO takes the text prompts as input to generate 4D avatars with object interactions. At the core of our network are: 1) *Text-driven 3D human and object composition* that employs large language models to retarget the contact areas from texts and spatial-aware SDS to composite the 3D models. 2) *Correspondence-aware motion optimization* which jointly optimizes the animation for humans and objects. It effectively maintains the spatial correspondence during animation, addressing the penetration issues.

DreamGaussian4D (Ren et al., 2023) employ multi-stage optimization pipelines to transform static 3D into dynamic 4D scenes. Dream-in-4D (Zheng et al., 2023) allows for personalized 4D generation using image guidance, while Consistent4D (Jiang et al., 2023c) uses video inputs with RIFE (Huang et al., 2022) and a super-resolution module for scene creation. 4DGen (Yin et al., 2023) and AnimatableDreamer (Wang et al., 2023c) focus on controllable motion generation via driving videos. More recently, Comp4D (Xu et al., 2024a) and TC4D (Bahmani et al., 2024) have introduced trajectory-based approaches for creating 4D compositional scenes. While these technologies show promise, they often struggle to produce 4D avatars that effectively interact with objects. Although GAvatar (Yuan et al., 2023) excels in 4D human animation, its object interaction capabilities are limited.

3 METHODOLOGY

Given a generated 3D avatar and a specific 3D object, AvatarGO generates compositional 4D avatars with object interactions based on text instructions. In the subsequent sections, we first introduce the preliminaries (in Sec. 3.1), including static 3D content generation and parametric human model SMPL-X. Next, we will describe the key components of AvatarGO, including (1) text-driven 3D human and object composition (in Sec. 3.2), and (2) correspondence-aware motion optimization for achieving synchronized human and object animation (in Sec. 3.3). The overview of AvatarGO is shown in Fig. 2.

3.1 PRELIMINARIES

3D Model Generation. Recently, DreamGaussian (Tang et al., 2023) showcases promising results with largely improved training efficiency by incorporating two major components:

(1) *3D Gaussian Splatting (3DGS).* 3DGS (Kerbl et al., 2023) directly defines the 3D space through a set of Gaussians parameterized by their 3D position μ , opacity α , anisotropic covariance Σ , and spherical harmonic coefficients sh . The sh term is used to capture the view-dependent appearance of

the scene and Σ can be decomposed to:

$$\Sigma = RSS^T R^T, \quad (1)$$

where R is the rotation matrix expressed by a quaternion $q \in \mathbf{SO}(3)$, and S is the scaling matrix, represented by a 3D vector s . Essentially, each Gaussian centered at point (mean) μ is defined as:

$$G(\mathbf{x}, \mu) = e^{-\frac{1}{2}(\mathbf{x}-\mu)^T \Sigma^{-1}(\mathbf{x}-\mu)}, \quad (2)$$

where \mathbf{x} is the 3D query point.

For rendering the 3D Gaussians onto the 2D image space, 3DGS incorporates a tile-based rasterizer and point-based α -blend rendering. Specifically, the color $C(u)$ of a pixel u can be calculated as:

$$C(u) = \sum_{i \in N} T_i c_i \alpha_i \mathcal{SH}(sh_i, v), \quad T_i = G(\mathbf{x}, \mu_i) \prod_{j=1}^{i-1} (1 - \alpha_j G(\mathbf{x}, \mu_j)), \quad (3)$$

where T represents the transmittance, \mathcal{SH} denotes the spherical harmonic function, and v indicates the viewing direction. By optimizing the Gaussian attributes $\{G : \mu, q, s, \sigma, c\}$ and dynamically adjusting the density of 3D Gaussians (*i.e.*, densifying and pruning), DreamGaussian achieves high-quality generations from either textual or visual inputs.

(2) *Score Distillation Sampling (SDS)*. Starting with the latent feature z extracted from a 3DGS rendering x , SDS introduces random noise ϵ to z , yielding a noisy latent variable z_t . This variable is then processed by a pre-trained denoising function $\epsilon_\phi(z_t; y, t)$ to estimate the added noise. The SDS loss then calculates the difference between predicted and added noise, with its gradient calculated by:

$$\nabla_\theta \mathcal{L}_{\text{SDS}}(\phi, g(\theta)) = \mathbb{E}_{t, \epsilon \sim \mathcal{N}(0,1)} \left[w(t) (\epsilon_\phi(z_t; y, t) - \epsilon) \frac{\partial z}{\partial \mathbf{x}} \frac{\partial \mathbf{x}}{\partial \theta} \right], \quad (4)$$

where y denotes the text embedding, $w(t)$ weights the loss from noise level t . We do not apply the mesh extraction and texture optimization proposed in DreamGaussian to obtain the 3D models.

SMPL-X (Loper et al., 2015; Pavlakos et al., 2019). With pose parameter θ , shape parameter β , and expression parameter ϕ as inputs, SMPL-X maps the canonical model to the observation space:

$$M(\beta, \theta, \phi) = \text{LBS}(\mathbf{T}(\beta, \theta, \phi), J(\beta), \theta, \mathcal{W}), \quad (5a)$$

$$\mathbf{T}(\beta, \theta, \phi) = \mathbf{T} + B_s(\beta) + B_e(\phi) + B_p(\theta), \quad (5b)$$

where M denotes the function defining the mesh model of a human body, and \mathbf{T} represents the transformed vertices. \mathcal{W} stands for blend weights, B_s , B_e , and B_p are functions respectively for shape, expression, and pose blend shapes. $\text{LBS}(\cdot)$ indicates the linear blend skinning function that poses each body vertex of SMPL-X according to:

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

where \mathbf{v}_c and \mathbf{v}_o represent SMPL-X vertices under the canonical pose and observation space, respectively. w_k is the skinning weight, $\mathcal{G}_k(\theta, j_k)$ is the affine deformation that maps the k -th joint j_k from the canonical space to observation space, and K denotes the number of neighboring joints.

3.2 TEXT-DRIVEN 3D HUMAN AND OBJECT COMPOSITION

With the help of DreamGaussian (Tang et al., 2023), we efficiently generate the 3D avatar G_h and the 3D object G_o individually based on 3DGS and SDS (discussed in Sec. 3.1). We noticed that even with manual adjustments, such as rescaling and rotating the 3D objects, it’s difficult to directly rig the generated 3D human and object models accurately (see Appx. F). Therefore, we strive to seamlessly composite G_h and G_o based on the text prompt in this stage. Specifically, the Gaussian attributes of G_h and G_o would be optimized, as well as three trainable global parameters of G_o , including rotation $\mathcal{R} \in \mathbb{R}^4$, scaling factor $\mathcal{S} \in \mathbb{R}$, and the translation matrix $\mathcal{T} \in \mathbb{R}^3$:

$$\mathbf{X}_{G_o} := \mathcal{S} \cdot (\mathbf{X}_{G_o} \cdot \mathcal{R} + \mathcal{T}), \quad (7)$$

where \mathbf{X}_{G_o} is the set of static Gaussian points.

270 However, solely utilizing SDS for optimization could frequently lead to disproportionate relationships
 271 and erroneous contact areas (see Fig. 3). This issue can be attributed to two potential factors: (1) the
 272 absence of emphasis on words describing human-object interaction, which decreases the model’s
 273 ability to comprehend the relationships between humans and objects; (2) the complexity inherent in
 274 human subjects, posing challenges for the diffusion model to identify the most suitable contact areas
 275 (see Sec. 4.3).

276 **Spatial-aware SDS (SSDS).** Following ComboVerse (Chen et al., 2024b), we employ SSDS to
 277 facilitate the compositional 3D generation between the human and the object. Specifically, SSDS
 278 augments the SDS with a spatial relationship between the human and the object by scaling the
 279 attention maps of the designated tokens $\langle \text{token}^* \rangle$ with a constant factor c (where $c > 1$):

$$280 \text{ATT} := \begin{cases} c \cdot \text{ATT}_{\langle \text{token} \rangle}, & \text{if } \langle \text{token} \rangle = \langle \text{token}^* \rangle, \\ \text{ATT}_{\langle \text{token} \rangle}, & \text{otherwise.} \end{cases} \quad (8)$$

283 Here, $\langle \text{token}^* \rangle$ corresponds to the tokens encoding the human-object interaction term, such as
 284 $\langle \text{'holding'} \rangle$, which can be identified through Large Language Models (LLMs) or specified by
 285 the user. Consequently, the spatial-aware SDS loss can be written as:

$$286 \nabla_{\theta} \mathcal{L}_{\text{SSDS}}(\phi^*, g(\theta)) = \mathbb{E}_{t, \epsilon \sim \mathcal{N}(0,1)} \left[w(t) (\epsilon_{\phi^*}(z_t; y, t) - \epsilon) \frac{\partial z}{\partial x} \frac{\partial x}{\partial \theta} \right], \quad (9)$$

289 where ϕ^* denotes the pre-trained denoising function with the adjusted attention maps.

290 **LLM-guided Contact Retargeting.** While spatial-aware SDS could benefit in understanding spatial
 291 correlations, it still faces difficulties in identifying the most appropriate contact area (See Fig. 3),
 292 which serves as a key component for human-object interaction. According to our studies (see Appx. E
 293 for visualization), the diffusion model struggles to accurately estimate contacts, even in the 2D images
 294 generated for human-object interaction. To tackle this issue, we propose leveraging Lang-SAM (Ian,
 295 2023) to identify the contact area from text prompts. Specifically, starting from the 3D human model
 296 G_h , we render it from a frontal viewpoint to produce the image I . This image, alongside textual
 297 inputs, undergoes Lang-SAM model to derive 2D segmentation masks \mathcal{M} :

$$298 \text{LangSAM}(I, \langle \text{body-part} \rangle) \rightarrow \mathcal{M}, \quad (10)$$

300 where $\langle \text{body-part} \rangle$ represents the text describing the human body part, such as $\langle \text{'hand'} \rangle$.
 301 Subsequently, we back-project the 2D segmentation labels onto the 3D Gaussians via inverse rendering
 302 (Chen et al., 2023c). Specifically, for each pixel u on the segmentation maps, we update the mask
 303 value (0 or 1) back to the Gaussians via:

$$304 w_i = \sum_{i \in \mathcal{N}} o_i(u) \times T_i(u) \times \mathcal{M}(u), \quad (11)$$

306 where w_i represents the weight of the i -th Gaussian, \mathcal{N} is the collection of Gaussians that can be
 307 projected onto the pixel u . $o(\cdot)$, $T(\cdot)$, and $\mathcal{M}(\cdot)$ respectively denote the opacity, transmittance, and
 308 segmentation mask value. Following the weight updates, we assess whether a Gaussian corresponds
 309 to the segmented region of the human body part by comparing its weight against a pre-defined
 310 threshold a . We then initialize the translation parameter \mathcal{T} according to:

$$311 \mathcal{T} = (\mathbf{w}^T * \boldsymbol{\mu}) / \sum \mathbf{w}, \quad (12)$$

313 where $\mathbf{w} = \{w_1, \dots, w_N | w_i = 0/1\} \in \mathbb{R}^{N \times 1}$, $\boldsymbol{\mu} = \{\mu_1, \dots, \mu_N\} \in \mathbb{R}^{N \times 3}$, and N is the number of
 314 Gaussain points within the human model G_h .

316 3.3 CORRESPONDENCE-AWARE MOTION FIELD

318 Following the compositional integration of 3D humans and objects, animating them synchronously
 319 presents an additional challenge owing to potential penetration issues. This problem stems from the
 320 absence of a well-defined motion field for the object. To this end, we establish the motion fields for
 321 both human and object models using the linear blend skinning function from SMPL-X (as in Eq. 6),
 322 and propose a correspondence-aware motion optimization aimed at optimizing the trainable global
 323 parameters of the object model, *i.e.*, rotation (\mathcal{R}) and translation (\mathcal{T}), to improve robustness against
 penetration issues between humans and objects.

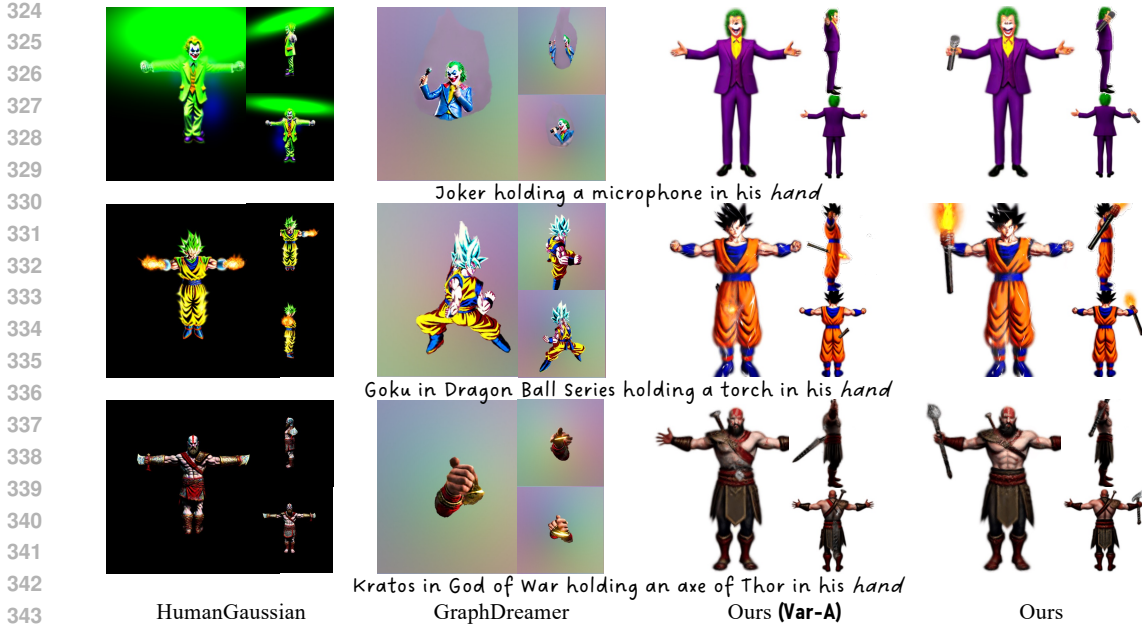


Figure 3: Comparisons on 3D compositional generations.

Human Animation. Given the motion sequence, we first construct a deformation field, which consists of two components: (1) articulated deformation utilizing the SMPL-X linear blend skinning function $LBS(\cdot)$, and (2) non-rigid motion learning the offset based on HexPlane features (Cao & Johnson, 2023), to deform the point \mathbf{x}_c from the canonical space to \mathbf{x}_o in the observation space:

$$\mathbf{x}_o = \mathcal{G} \cdot \mathbf{x}_c + \text{MLP}(F(\mathbf{x}_c, \mathbf{t})), \tag{13}$$

where $F(\cdot)$ denotes the HexPlane-based feature extraction network, and \mathbf{t} indicates the timestamp. We derive \mathcal{G} from the closet canonical SMPL-X vertex to \mathbf{x}_c .

Object Animation. Similar to the human animation, we calculate the deformation matrix \mathcal{G}_c for each Gaussian point \mathbf{x} within the object model G_o based on its closest canonical SMPL-X vertex. Given our experimental definition of 3D objects as rigid bodies, we then compute their average to establish the intermediate motion field for the object:

$$\mathbf{X}_o = \mathcal{G}'_c \cdot \mathbf{X}_c, \quad \mathcal{G}'_c = \frac{\sum_{i \in [1, M]} \mathcal{G}_{c_i}}{M}, \tag{14}$$

where $\mathbf{X}_o = \{\mathbf{x}_{o_1}, \dots, \mathbf{x}_{o_M}\}$, $\mathbf{X}_c = \{\mathbf{x}_{c_1}, \dots, \mathbf{x}_{c_M}\}$, and M is the total number of Gaussian points within G_o . Although animating the object directly using SMPL-X linear blend skinning may seem like a simple solution, it can result in penetration issues between the human and the object (see Fig. 5). This challenge arises primarily from the absence of proper constraints to maintain the correspondence between these two models.

Correspondence-aware Motion Optimization. Drawing insight from the fact that our method is robust in handling penetration issues in static composited models across various scenarios, we propose a correspondence-aware motion optimization to preserve the correspondence between human and object, thereby addressing the penetration problem. Specifically, we extend the above motion field (Eq. 14) to include two additional trainable parameters \mathcal{R} and \mathcal{T} :

$$\mathbf{X}_o := \mathbf{X}_o \cdot \mathcal{R} + \mathcal{T}. \tag{15}$$

where \mathbf{X}_o is obtained in Eq. 14. Rather than naively optimizing the parameters via SDS, we propose a novel correspondence-aware training objective that leverages SMPL-X as an intermediary to maintain the correspondence between human and object when they are animated to a new pose:

$$\mathcal{L}_{CA} = \text{MSE}(\mathcal{G}_c, \mathcal{G}_o), \quad \mathcal{G}_c = \{\mathcal{G}_{c_0}, \dots, \mathcal{G}_{c_M}\}, \quad \mathcal{G}_o = \{\mathcal{G}_{o_0}, \dots, \mathcal{G}_{o_M}\} \tag{16}$$

where \mathcal{G}_{c_i} and \mathcal{G}_{o_i} is respectively derived based on \mathbf{x}_{c_i} , \mathbf{x}_{o_i} and their corresponding SMPL-X models.



Figure 4: **Comparisons on 4D avatar animation with object interactions.** ‘*’ indicates that HumanGaussian directly employs the SMPL LBS function for animation.

In addition to our correspondence-aware loss, we also incorporate the spatial-aware SDS as in Eq. 9 and the texture-structure joint SDS from HumanGaussian (Liu et al., 2023d) to enhance the overall quality:

$$\begin{aligned} \nabla_{\theta} \mathcal{L}_{\text{SDS}}(\phi, g(\theta)) = & \lambda_1 \cdot \mathbb{E}_{t, \epsilon \sim \mathcal{N}(0,1)} \left[w(t) (\epsilon_{\phi}(z_{x_t}; y, t) - \epsilon_x) \frac{\partial z_x}{\partial x} \frac{\partial x}{\partial \theta} \right] \\ & + \lambda_2 \cdot \mathbb{E}_{t, \epsilon \sim \mathcal{N}(0,1)} \left[w(t) (\epsilon_{\phi}(z_{d_t}; y, t) - \epsilon_d) \frac{\partial z_d}{\partial d} \frac{\partial d}{\partial \theta} \right], \end{aligned} \quad (17)$$

where λ_1 and λ_2 are hyper-parameters to balance the impact of structural and textural losses, while d denotes the depth renderings.

The overall loss function to optimize the 4D animative scene is then given by:

$$\mathcal{L} = \lambda_{CA} \cdot \mathcal{L}_{CA} + \lambda_{\text{SDS}} \cdot \mathcal{L}_{\text{SDS}} + \lambda_{\text{SSDS}} \cdot \mathcal{L}_{\text{SSDS}}, \quad (18)$$

where λ_{CA} , λ_{SDS} , and λ_{SSDS} represents weights to balance the respective losses.

4 EXPERIMENTS

We now validate the effectiveness and capability of our proposed framework to animate various 3D avatar-object pairs with different poses and provide comparisons with existing 3D and 4D compositional generation methods.

Implementation Details. We follow DreamGaussian4D (Ren et al., 2023) to implement the 3D Gaussian Splatting (Kerbl et al., 2023) and the HexPlane (Cao & Johnson, 2023) in our method. We utilize the pre-trained Texture-Structure joint diffusion model from HumanGaussian (Liu et al., 2023d) and version 2.1 of Stable Diffusion (Stability.AI, 2022) to respectively calculate the SDS and spatial-aware SDS in our implementation. Typically, for each 3D avatar-object pair, we train the 3D stage with a batch size of 16 for 400 epochs, and the 4D stage with a batch size of 10 for 400 epochs. The training takes around 10 minutes for the 3D stage and 20 minutes for the 4D stage on a single NVIDIA A100 GPU. We use Adam (Kingma & Ba, 2015) optimizer for back-propagation. Additional implementation details can be found in the Appx. B.

Comparison Methods for 3D Static Generation. We first compare the 3D static generation results with HumanGaussian (Liu et al., 2023d) and GraphDreamer (Gao et al., 2023). Since ComboVerse (Chen et al., 2024b) lacks an official code release and relies on image inputs, we compare static AvatarGO with an alternative variant, *i.e.*, “Ours (Var-A)”, by only using the spatial-aware score distillation sampling (SSDS) in ComboVerse to composite 3D humans and avatars. We cannot compare with GALA3D as their source code is not publicly accessible.

Comparison Methods for 4D Animation. Since there are no specific methods tailored for 4D avatar animation with object interactions, we assess AvatarGO’s efficacy against three recent 4D generation techniques (i.e., DreamGaussian4D (Ren et al., 2023), HumanGaussian (Liu et al., 2023d), and TC4D (Bahmani et al., 2024)), as well as one alternative variant “Ours (Var-B)”. To implement **Var-B**, we utilize human hand motion sequences as trajectories to guide the transformation of 3D objects and follow Comp4D to integrate the video diffusion model to compute SDS. Because InterDreamer (Xu et al., 2024b) and InterFusion (Dai et al., 2024) have not released their code, we could not include their results for comparison. See more motivation for designing “Ours (Var-A)” and “Ours (Var-B)” in Appx. C.

4.1 QUALITATIVE EVALUATIONS

4D Avatar Generation with Object Interaction. In Fig. 1, we present a diverse collection of avatar-object pairs that are animated to different poses. These renderings consistently showcase high-fidelity results from various viewpoints. Thanks to our proposed LLM-guided contact retargeting and correspondence-aware motion optimization, our method can deliver appropriate human-object interactions and demonstrate superior robustness to the penetration issues.

Comparison on 3D Generation. We provide qualitative comparisons with existing methods on 3D generation in Fig. 3. We can observe: 1) without the aid of LLMs, HumanGaussian struggles to determine the spatial correlations between humans and objects; 2) Despite using graphs to establish relationships, GraphDreamer is confused by the meaningful contact, resulting in unsatisfactory outcomes. 3) Optimizing \mathcal{R} , \mathcal{S} , and \mathcal{T} with only SSDS is inadequate to move the object to the correct area. Conversely, AvatarGO consistently outperforms with precise human-object interactions.

Comparisons on 4D Animation. In Fig. 4, we compare our 4D animation results with SOTA methods. We take the rendering from our 3D compositions stage as the input for DreamGaussian4D. The following observations can be made: 1) Even with human-object interaction images, DreamGaussian4D, which employs video diffusion models, struggles with animating the composited scene. 2) Direct animation via SMPL LBS function, as in HumanGaussian, tends to yield unsmooth results, especially for the arms. 3) TC4D faces similar issues as the DreamGaussian4D. Meanwhile, it treats the entire scene as a single entity, lacking both local and large-scale motions for individual objects. 4) One may think applying trajectory to objects seems like a simple solution (as in Comp4D). However, as seen in “Ours (Var-B)”, it can disrupt spatial correlations between humans and objects. These points further validate the necessity of AvatarGO. Our method can consistently deliver superior results with correct relationships and better robustness to penetration issues. See the Appx. A, H, J, K for more comparisons.

4.2 QUANTITATIVE EVALUATIONS

CLIP-based Metrics. We use CLIP-based metrics (CLIP-Score (CLIP-S), CLIP image similarity (CLIP-Image), and CLIP Directional Similarity (CLIP-DS) (Brooks et al., 2023; Gal et al., 2022)) with CLIP-ViT-L/14 model. Among them, CLIP-S measures the similarity between texts and their corresponding models, CLIP-Image denotes the similarity between compositional models and human models, and CLIP-DS represents the alignment between changes in text captions (e.g., “Iron Man” to “Iron Man holding an axe of Thor in his hand”) and corresponding changes in images. Through Tab. 1, our method maintains the human identity in the composited scenes (see CLIP-Image). Note that “Ours (Var-A)” and GraphDreamer is slightly better for this metric as they struggle to do the composition (see Fig. 3). Meanwhile, “Ours” and “Ours (static)” consistently achieve better results than HumanGaussian and other variants, further affirming the objective superiority of AvatarGO.

User Studies We further conduct user studies to compare with DreamGaussian4D, HumanGaussian, TC4D, and “Ours (Var-A)”. 24 Volunteers rated these methods independently based on seven criteria from 1 (worst) to 5 (best): (1) Level of penetration; (2) Accuracy of the relative scale between humans and objects; (3) Accuracy

Table 1: Quantitative Evaluation.

	GraphDreamer	TC4D	HumanGaussian	Ours (Var-A)	Ours (Var-B)	Ours (static)	Ours
CLIP-Image \uparrow	98.44	89.50	83.93	97.88	92.11	93.45	92.20
CLIP-S \uparrow	8.09	19.84	23.69	25.36	30.57	32.27	32.84
CLIP-DS \uparrow	1.71	15.28	4.71	0.91	25.90	33.80	28.03

Table 2: User studies.

	Dream Gaussian4D	Human Gaussian	TC4D	Ours (Var-B)	Ours
Level of penetration \uparrow	1.267	1.084	4.236	1.537	4.872
Accuracy of relative scale \uparrow	1.183	1.092	4.308	3.947	4.788
Accuracy of contact \uparrow	1.654	1.137	4.412	2.137	4.802
Motion quality \uparrow	1.321	2.156	1.947	1.673	4.592
Motion amount \uparrow	2.118	3.781	1.517	4.159	4.934
Text alignment \uparrow	2.047	1.918	4.515	2.462	4.767
Overall Performance \uparrow	3.467	1.633	4.033	2.033	4.869



Figure 5: **Analysis of correspondence-aware motion field.**

of contact; (4) Motion quality; (5) Motion amount; (6) Text alignment; (7) Overall performance. Detailed results have been presented in Tab. 2. Key observations include: 1) Both DreamGaussian4D and HumanGaussian have difficulty providing satisfactory outcomes for human-object interaction (HOI) scenes. 2) Although TC4D performs well with HOI generations, it only produces global motions, leading to less optimal motion quality and quantity compared to our method. Our final design consistently delivers superior results for all seven criteria, outperforming the other methods across the board.

4.3 ABLATION STUDIES

Analysis of LLM-guided Contact Retargeting. We first conduct evaluations to validate the efficacy of employing Lang-SAM for retargeting the accurate contact area. See Fig. 3. By comparing “Ours (Var-A)” and Ours, we can conclude that without Lang-SAM, the model struggles to produce correct human-object interaction in the 3D compositional generation.

Analysis of Correspondence-aware Motion Field. In Fig. 5, we first compare our proposed training objectives \mathcal{L}_{CA} with two alternative strategies: 1) “SDF distance loss” which minimizes the change of signed distance field (SDF) between objects and humans when they are animated to a new pose, and 2) “SDF label loss” that supervise the label of SDF instead. These comparisons demonstrate the effectiveness of our proposed method for maintaining spatial correlations during the animations. Additionally, we validate our model’s design by further comparing it with two variants: 1) “w/o $\mathcal{R}, \mathcal{T}, \mathcal{L}_{CA}$ ” which disables the trainable parameters \mathcal{R}, \mathcal{T} , Eq. 15, and our proposed loss \mathcal{L}_{CA} . This setting represents the scenario where the object is moved directly with the contact point. and 2) “w/o \mathcal{L}_{CA} ” which trains the animation network solely with SDS loss ($\mathcal{L}_{SDS}^*, \mathcal{L}_{SDS}$). These comparisons underscore the necessity of these components in achieving 4D animation with better robustness to the penetration issues.

Analysis of Spatial-aware SDS.

We finally assess the effectiveness of spatial-aware SDS (SSDS) and present the results in Fig. 6. Notably, we observe that SSDS plays a crucial role in preventing the optimization of $\mathcal{R}, \mathcal{S}, \mathcal{T}$ from vanishing during 3D compositional generation. Additionally, there is a drop in the quality of the animated avatars when disabling SSDS.

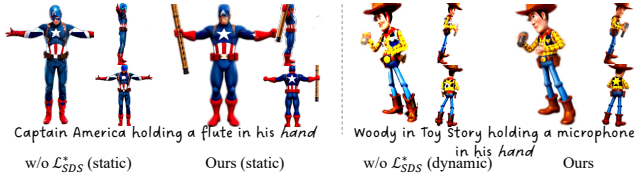


Figure 6: **Analysis of Spatial-aware SDS.**

5 CONCLUSIONS

In this paper, we have introduced AvatarGO, the first attempt for text-guided 4D avatar generation with object interactions. Within AvatarGO, we proposed to employ large language model for comprehending the most suitable contact area between humans and objects. We also presented a novel correspondence-aware motion optimization that utilizes SMPL-X as an intermediary to enhance the model’s resilience to penetration issues when animating 3D humans and objects into new poses. Extensive evaluations demonstrated that our method has achieved high-fidelity 4D animations across diverse 3D avatar-object pairs and poses, surpassing current state-of-the-arts by a large margin.

Limitations. While opening new doors for human-centric 4D content generation, we acknowledge AvatarGO has certain limitations: 1) Our pipeline operates under the assumption of rigid-body dynamics for 3D objects, making it unsuitable for animating non-rigid content such as flags; 2) our method presumes that continuous contact between objects and avatars, making it challenges for tasks like “Dribbling the basketball,” where the human and object inevitably disconnect at certain points. Nevertheless, our current approach does not cover all possible scenarios, it effectively handles continuous contact and rigid connections, which are commonly encountered in real-world applications.

REFERENCES

- 540
541
542 Language segment anything. <https://github.com/paulguerrero/lang-sam.git>,
543 2023. 3, 6
- 544 Sherwin Bahmani, Ivan Skorokhodov, Victor Rong, Gordon Wetzstein, Leonidas Guibas, Peter
545 Wonka, Sergey Tulyakov, Jeong Joon Park, Andrea Tagliasacchi, and David B Lindell. 4d-fy:
546 Text-to-4d generation using hybrid score distillation sampling. *arXiv preprint arXiv:2311.17984*,
547 2023. 3
- 548
549 Sherwin Bahmani, Xian Liu, Yifan Wang, Ivan Skorokhodov, Victor Rong, Ziwei Liu, Xihui Liu,
550 Jeong Joon Park, Sergey Tulyakov, Gordon Wetzstein, et al. Tc4d: Trajectory-conditioned text-to-
551 4d generation. *arXiv preprint arXiv:2403.17920*, 2024. 2, 4, 9, 19
- 552 Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Karsten Kreis, Miika Aittala,
553 Timo Aila, Samuli Laine, Bryan Catanzaro, et al. ediffi: Text-to-image diffusion models with an
554 ensemble of expert denoisers. *arXiv preprint arXiv:2211.01324*, 2022. 1, 3
- 555
556 Bharat Lal Bhatnagar, Xianghui Xie, Ilya Petrov, Cristian Sminchisescu, Christian Theobalt, and
557 Gerard Pons-Moll. Behave: Dataset and method for tracking human object interactions. In *IEEE*
558 *Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, jun 2022. 1
- 559
560 Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik
561 Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling
562 latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023. 2
- 563
564 Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image
565 editing instructions. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2023. 9
- 566
567 Ang Cao and Justin Johnson. Hexplane: A fast representation for dynamic scenes. In *Proceedings of*
568 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 130–141, 2023. 7, 8,
18
- 569
570 Tianshi Cao, Karsten Kreis, Sanja Fidler, Nicholas Sharp, and Kangxue Yin. Textfusion: Synthesizing
571 3d textures with text-guided image diffusion models. In *Proceedings of the IEEE/CVF International*
572 *Conference on Computer Vision*, pp. 4169–4181, 2023a. 3
- 573
574 Yukang Cao, Yan-Pei Cao, Kai Han, Ying Shan, and Kwan-Yee K Wong. Dreamavatar: Text-and-
575 shape guided 3d human avatar generation via diffusion models. *arXiv preprint arXiv:2304.00916*,
2023b. 1, 3
- 576
577 Yukang Cao, Yan-Pei Cao, Kai Han, Ying Shan, and Kwan-Yee K Wong. Guide3d: Create 3d avatars
578 from text and image guidance. *arXiv preprint arXiv:2308.09705*, 2023c. 3
- 579
580 Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and
581 Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the*
IEEE/CVF international conference on computer vision, pp. 9650–9660, 2021. 3
- 582
583 Dave Zhenyu Chen, Yawar Siddiqui, Hsin-Ying Lee, Sergey Tulyakov, and Matthias Nießner.
584 Text2tex: Text-driven texture synthesis via diffusion models. *arXiv preprint arXiv:2303.11396*,
2023a. 3
- 585
586 Rui Chen, Yongwei Chen, Ningxin Jiao, and Kui Jia. Fantasia3d: Disentangling geometry and
587 appearance for high-quality text-to-3d content creation. *arXiv preprint arXiv:2303.13873*, 2023b.
588 3
- 589
590 Rui Chen, Mingyi Shi, Shaoli Huang, Ping Tan, Taku Komura, and Xuelin Chen. Taming diffusion
591 probabilistic models for character control. In *SIGGRAPH*, 2024a. 1
- 592
593 Yiwen Chen, Zilong Chen, Chi Zhang, Feng Wang, Xiaofeng Yang, Yikai Wang, Zhongang Cai, Lei
Yang, Huaping Liu, and Guosheng Lin. Gaussianeditor: Swift and controllable 3d editing with
gaussian splatting. *arXiv preprint arXiv:2311.14521*, 2023c. 6

- 594 Yongwei Chen, Tengfei Wang, Tong Wu, Xingang Pan, Kui Jia, and Ziwei Liu. Comboverse:
595 Compositional 3d assets creation using spatially-aware diffusion guidance. *arXiv preprint*
596 *arXiv:2403.12409*, 2024b. 2, 3, 6, 8
- 597 Sisi Dai, Wenhao Li, Haowen Sun, Haibin Huang, Chongyang Ma, Hui Huang, Kai Xu, and
598 Ruizhen Hu. Interfusion: Text-driven generation of 3d human-object interaction. *arXiv preprint*
599 *arXiv:2403.15612*, 2024. 2, 3, 9
- 600 Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig
601 Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of anno-
602 tated 3d objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
603 *Recognition*, pp. 13142–13153, 2023. 3
- 604 Matt Deitke, Ruoshi Liu, Matthew Wallingford, Huong Ngo, Oscar Michel, Aditya Kusupati, Alan
605 Fan, Christian Laforte, Vikram Voleti, Samir Yitzhak Gadre, et al. Objaverse-xl: A universe of
606 10m+ 3d objects. *Advances in Neural Information Processing Systems*, 36, 2024. 3
- 607 Dave Epstein, Ben Poole, Ben Mildenhall, Alexei A Efros, and Aleksander Holynski. Disentangled
608 3d scene generation with layout learning. *arXiv preprint arXiv:2402.16936*, 2024. 2, 3
- 609 Rinon Gal, Or Patashnik, Haggai Maron, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or.
610 Stylegan-nada: Clip-guided domain adaptation of image generators. *ACM Transactions on*
611 *Graphics (TOG)*, 2022. 9
- 612 Gege Gao, Weiyang Liu, Anpei Chen, Andreas Geiger, and Bernhard Schölkopf. Graphdreamer:
613 Compositional 3d scene synthesis from scene graphs. *arXiv preprint arXiv:2312.00093*, 2023. 1,
614 3, 8, 19, 22
- 615 Xun Guo, Mingwu Zheng, Liang Hou, Yuan Gao, Yufan Deng, Chongyang Ma, Weiming Hu,
616 Zhengjun Zha, Haibin Huang, Pengfei Wan, et al. I2v-adapter: A general image-to-video adapter
617 for video diffusion models. *arXiv preprint arXiv:2312.16693*, 2023a. 2
- 618 Yuan-Chen Guo, Ying-Tian Liu, Ruizhi Shao, Christian Laforte, Vikram Voleti, Guan Luo, Chia-
619 Hao Chen, Zi-Xin Zou, Chen Wang, Yan-Pei Cao, and Song-Hai Zhang. threestudio: A unified
620 framework for 3d content generation. [https://github.com/threestudio-project/](https://github.com/threestudio-project/threestudio)
621 [threestudio](https://github.com/threestudio-project/threestudio), 2023b. 17, 18
- 622 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
623 recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2016. 18
- 624 Fangzhou Hong, Jiaxiang Tang, Ziang Cao, Min Shi, Tong Wu, Zhaoxi Chen, Tengfei Wang, Liang
625 Pan, Dahua Lin, and Ziwei Liu. 3dtopia: Large text-to-3d generation model with hybrid diffusion
626 priors. *arXiv preprint arXiv:2403.02234*, 2024. 3
- 627 Yicong Hong, Kai Zhang, Jiuxiang Gu, Sai Bi, Yang Zhou, Difan Liu, Feng Liu, Kalyan Sunkavalli,
628 Trung Bui, and Hao Tan. Lrm: Large reconstruction model for single image to 3d. *arXiv preprint*
629 *arXiv:2311.04400*, 2023. 3
- 630 Xin Huang, Ruizhi Shao, Qi Zhang, Hongwen Zhang, Ying Feng, Yebin Liu, and Qing Wang.
631 Humannorm: Learning normal diffusion model for high-quality and realistic 3d human generation.
632 *arXiv preprint arXiv:2310.01406*, 2023a. 3
- 633 Yukun Huang, Jianan Wang, Ailing Zeng, He Cao, Xianbiao Qi, Yukai Shi, Zheng-Jun Zha, and
634 Lei Zhang. Dreamwaltz: Make a scene with complex 3d animatable avatars. *arXiv preprint*
635 *arXiv:2305.12529*, 2023b. 3
- 636 Zhewei Huang, Tianyuan Zhang, Wen Heng, Boxin Shi, and Shuchang Zhou. Real-time intermediate
637 flow estimation for video frame interpolation. In *European Conference on Computer Vision*, pp.
638 624–642. Springer, 2022. 4
- 639 Nan Jiang, Tengyu Liu, Zhexuan Cao, Jieming Cui, Zhiyuan Zhang, Yixin Chen, He Wang, Yixin
640 Zhu, and Siyuan Huang. Full-body articulated human-object interaction. In *Proceedings of the*
641 *IEEE/CVF International Conference on Computer Vision*, pp. 9365–9376, 2023a. 1

- 648 Ruixiang Jiang, Can Wang, Jingbo Zhang, Menglei Chai, Mingming He, Dongdong Chen, and Jing
649 Liao. Avatarcraft: Transforming text into neural human avatars with parameterized shape and pose
650 control. *arXiv preprint arXiv:2303.17606*, 2023b. 3
- 651 Yanqin Jiang, Li Zhang, Jin Gao, Weimin Hu, and Yao Yao. Consistent4d: Consistent 360 $\{\backslash\deg\}$
652 dynamic object generation from monocular video. *arXiv preprint arXiv:2311.02848*, 2023c. 4
- 653 Roy Kapon, Guy Tevet, Daniel Cohen-Or, and Amit H Bermano. Mas: Multi-view ancestral sampling
654 for 3d motion generation using 2d diffusion. In *Proceedings of the IEEE/CVF Conference on*
655 *Computer Vision and Pattern Recognition*, pp. 1965–1974, 2024. 1
- 656 Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting
657 for real-time radiance field rendering. *ACM Transactions on Graphics (ToG)*, 42(4):1–14, 2023. 3,
658 4, 8
- 659 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International*
660 *Conference on Learning Representations*, 2015. 8
- 661 Nikos Kolotouros, Thiemo Alldieck, Andrei Zanfir, Eduard Gabriel Bazavan, Mihai Fieraru,
662 and Cristian Sminchisescu. Dreamhuman: Animatable 3d avatars from text. *arXiv preprint*
663 *arXiv:2306.09329*, 2023. 3
- 664 Jiahao Li, Hao Tan, Kai Zhang, Zexiang Xu, Fujun Luan, Yinghao Xu, Yicong Hong, Kalyan
665 Sunkavalli, Greg Shakhnarovich, and Sai Bi. Instant3d: Fast text-to-3d with sparse-view generation
666 and large reconstruction model. *arXiv preprint arXiv:2311.06214*, 2023a. 3
- 667 Jiaman Li, Jiajun Wu, and C Karen Liu. Object motion guided human motion synthesis. *ACM Trans.*
668 *Graph.*, 42(6), 2023b. 1
- 669 Tingting Liao, Hongwei Yi, Yuliang Xiu, Jiaxiang Tang, Yangyi Huang, Justus Thies, and Michael J
670 Black. Tada! text to animatable digital avatars. *arXiv preprint arXiv:2308.10899*, 2023. 1, 3
- 671 Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten
672 Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. Magic3d: High-resolution text-to-3d content
673 creation. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2023. 1, 3
- 674 Fangfu Liu, Hanyang Wang, Weiliang Chen, Haowen Sun, and Yueqi Duan. Make-your-3d: Fast and
675 consistent subject-driven 3d content generation. *arXiv preprint arXiv:2403.09625*, 2024. 3
- 676 Minghua Liu, Ruoxi Shi, Linghao Chen, Zhuoyang Zhang, Chao Xu, Xinyue Wei, Hansheng Chen,
677 Chong Zeng, Jiayuan Gu, and Hao Su. One-2-3-45++: Fast single image to 3d objects with
678 consistent multi-view generation and 3d diffusion. *arXiv preprint arXiv:2311.07885*, 2023a. 3
- 679 Minghua Liu, Chao Xu, Haian Jin, Linghao Chen, Zexiang Xu, Hao Su, et al. One-2-3-45: Any single
680 image to 3d mesh in 45 seconds without per-shape optimization. *arXiv preprint arXiv:2306.16928*,
681 2023b. 3
- 682 Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick.
683 Zero-1-to-3: Zero-shot one image to 3d object. In *Proceedings of the IEEE/CVF International*
684 *Conference on Computer Vision*, pp. 9298–9309, 2023c. 1, 3
- 685 Xian Liu, Xiaohang Zhan, Jiayang Tang, Ying Shan, Gang Zeng, Dahua Lin, Xihui Liu, and Ziwei
686 Liu. Humangaussian: Text-driven 3d human generation with gaussian splatting. *arXiv preprint*
687 *arXiv:2311.17061*, 2023d. 2, 8, 9, 19, 22, 23
- 688 Yuan Liu, Cheng Lin, Zijiao Zeng, Xiaoxiao Long, Lingjie Liu, Taku Komura, and Wenping Wang.
689 Syncdreamer: Generating multiview-consistent images from a single-view image. *arXiv preprint*
690 *arXiv:2309.03453*, 2023e. 3
- 691 Zexiang Liu, Yangguang Li, Youtian Lin, Xin Yu, Sida Peng, Yan-Pei Cao, Xiaojuan Qi, Xiaoshui
692 Huang, Ding Liang, and Wanli Ouyang. Unidream: Unifying diffusion priors for relightable
693 text-to-3d generation. *arXiv preprint arXiv:2312.08754*, 2023f. 3

- 702 Xiaoxiao Long, Yuan-Chen Guo, Cheng Lin, Yuan Liu, Zhiyang Dou, Lingjie Liu, Yuexin Ma,
703 Song-Hai Zhang, Marc Habermann, Christian Theobalt, et al. Wonder3d: Single image to 3d using
704 cross-domain diffusion. *arXiv preprint arXiv:2310.15008*, 2023. 3
- 705 Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. SMPL:
706 A skinned multi-person linear model. *ACM Trans. Graphics, Asia*, 2015. 1, 5
- 707 Gal Metzer, Elad Richardson, Or Patashnik, Raja Giryes, and Daniel Cohen-Or. Latent-nerf for
708 shape-guided generation of 3d shapes and textures. *arXiv preprint arXiv:2211.07600*, 2022. 3
- 709 Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and
710 Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European
711 Conference on Computer Vision*, 2020. 3
- 712 Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov,
713 Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning
714 robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 3
- 715 Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed A. A. Osman, Dimitrios
716 Tzionas, and Michael J. Black. Expressive body capture: 3d hands, face, and body from a single
717 image. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 1, 5
- 718 Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d
719 diffusion. In *International Conference on Learning Representations*, 2022. 1
- 720 Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-
721 conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022. 1, 3
- 722 Jiawei Ren, Liang Pan, Jiayang Tang, Chi Zhang, Ang Cao, Gang Zeng, and Ziwei Liu. Dreamgaus-
723 sian4d: Generative 4d gaussian splatting. *arXiv preprint arXiv:2312.17142*, 2023. 2, 4, 8, 9, 17,
724 19, 23
- 725 Elad Richardson, Gal Metzer, Yuval Alaluf, Raja Giryes, and Daniel Cohen-Or. Texture: Text-guided
726 texturing of 3d shapes. *arXiv preprint arXiv:2302.01721*, 2023. 3
- 727 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed
728 Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al.
729 Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint
730 arXiv:2205.11487*, 2022. 1, 3
- 731 Junyoung Seo, Wooseok Jang, Min-Seop Kwak, Jaehoon Ko, Hyeonsu Kim, Junho Kim, Jin-Hwa
732 Kim, Jiyoung Lee, and Seungryong Kim. Let 2d diffusion model know 3d-consistency for robust
733 text-to-3d generation. *arXiv preprint arXiv:2303.07937*, 2023. 3
- 734 Yonatan Shafir, Guy Tevet, Roy Kapon, and Amit H Bermano. Human motion diffusion as a
735 generative prior. *arXiv preprint arXiv:2303.01418*, 2023. 1
- 736 Tianchang Shen, Jun Gao, Kangxue Yin, Ming-Yu Liu, and Sanja Fidler. Deep marching tetrahedra:
737 a hybrid representation for high-resolution 3d shape synthesis. In *Advances in Neural Information
738 Processing Systems*, 2021. 3
- 739 Tianchang Shen, Jacob Munkberg, Jon Hasselgren, Kangxue Yin, Zian Wang, Wenzheng Chen, Zan
740 Gojcic, Sanja Fidler, Nicholas Sharp, and Jun Gao. Flexible isosurface extraction for gradient-based
741 mesh optimization. *ACM Transactions on Graphics (TOG)*, 42(4):1–16, 2023. 3
- 742 Ruoxi Shi, Hansheng Chen, Zhuoyang Zhang, Minghua Liu, Chao Xu, Xinyue Wei, Linghao Chen,
743 Chong Zeng, and Hao Su. Zero123++: a single image to consistent multi-view diffusion base
744 model. *arXiv preprint arXiv:2310.15110*, 2023a. 3
- 745 Yichun Shi, Peng Wang, Jianglong Ye, Long Mai, Kejie Li, and Xiao Yang. Mvdream: Multi-view
746 diffusion for 3d generation. *arXiv:2308.16512*, 2023b. 3
- 747 Uriel Singer, Shelly Sheynin, Adam Polyak, Oron Ashual, Iurii Makarov, Filippos Kokkinos, Naman
748 Goyal, Andrea Vedaldi, Devi Parikh, Justin Johnson, et al. Text-to-4d dynamic scene generation.
749 *arXiv preprint arXiv:2301.11280*, 2023. 3
- 750
751
752
753
754
755

- 756 Stability.AI. Stable diffusion. [https://stability.ai/blog/](https://stability.ai/blog/stable-diffusion-public-release)
757 [stable-diffusion-public-release](https://stability.ai/blog/stable-diffusion-public-release), 2022. 1, 3, 8
758
- 759 Stability.AI. Stability AI releases DeepFloyd IF, a powerful text-to-image model
760 that can smartly integrate text into images. [https://stability.ai/blog/](https://stability.ai/blog/deepfloyd-if-text-to-image-model)
761 [deepfloyd-if-text-to-image-model](https://stability.ai/blog/deepfloyd-if-text-to-image-model), 2023. 1, 3
- 762 Jiaxiang Tang, Jiawei Ren, Hang Zhou, Ziwei Liu, and Gang Zeng. Dreamgaussian: Generative
763 gaussian splatting for efficient 3d content creation. *arXiv preprint arXiv:2309.16653*, 2023. 1, 2, 3,
764 4, 5
- 765 Jiaxiang Tang, Zhaoxi Chen, Xiaokang Chen, Tengfei Wang, Gang Zeng, and Ziwei Liu. Lgm:
766 Large multi-view gaussian model for high-resolution 3d content creation. *arXiv preprint*
767 *arXiv:2402.05054*, 2024a. 3
768
- 769 Jiaxiang Tang, Ruijie Lu, Xiaokang Chen, Xiang Wen, Gang Zeng, and Ziwei Liu. Intex: Interactive
770 text-to-texture synthesis via unified depth-aware inpainting. *arXiv preprint arXiv:2403.11878*,
771 2024b. 3
- 772 Shitao Tang, Jiacheng Chen, Dilin Wang, Chengzhou Tang, Fuyang Zhang, Yuchen Fan, Vikas
773 Chandra, Yasutaka Furukawa, and Rakesh Ranjan. Mvdifffusion++: A dense high-resolution
774 multi-view diffusion model for single or sparse-view 3d object reconstruction. *arXiv preprint*
775 *arXiv:2402.12712*, 2024c. 3
776
- 777 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
778 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing*
779 *systems*, 30, 2017. 3
- 780 Peihao Wang, Zhiwen Fan, Dejia Xu, Dilin Wang, Sreyas Mohan, Forrest Iandola, Rakesh Ranjan,
781 Yilei Li, Qiang Liu, Zhangyang Wang, et al. Steindreamer: Variance reduction for text-to-3d score
782 distillation via stein identity. *arXiv preprint arXiv:2401.00604*, 2023a. 3
- 783 Peng Wang, Hao Tan, Sai Bi, Yinghao Xu, Fujun Luan, Kalyan Sunkavalli, Wenping Wang, Zexiang
784 Xu, and Kai Zhang. Pf-lrm: Pose-free large reconstruction model for joint pose and shape
785 prediction. *arXiv preprint arXiv:2311.12024*, 2023b. 3
786
- 787 Xinzhou Wang, Yikai Wang, Junliang Ye, Zhengyi Wang, Fuchun Sun, Pengkun Liu, Ling Wang, Kai
788 Sun, Xintong Wang, and Bin He. Animatable-dreamer: Text-guided non-rigid 3d model generation
789 and reconstruction with canonical score distillation. *arXiv preprint arXiv:2312.03795*, 2023c. 4
- 790 Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. Prolific-
791 dreamer: High-fidelity and diverse text-to-3d generation with variational score distillation. *arXiv*
792 *preprint arXiv:2305.16213*, 2023d. 1, 3
- 793 Zhengyi Wang, Yikai Wang, Yifei Chen, Chendong Xiang, Shuo Chen, Dajiang Yu, Chongxuan Li,
794 Hang Su, and Jun Zhu. Crm: Single image to 3d textured mesh with convolutional reconstruction
795 model. *arXiv preprint arXiv:2403.05034*, 2024. 3
796
- 797 Tianyu Wu, Shizhu He, Jingping Liu, Siqi Sun, Kang Liu, Qing-Long Han, and Yang Tang. A brief
798 overview of chatgpt: The history, status quo and potential future development. *IEEE/CAA Journal*
799 *of Automatica Sinica*, 10(5):1122–1136, 2023a. 1
- 800 Tong Wu, Jiarui Zhang, Xiao Fu, Yuxin Wang, Jiawei Ren, Liang Pan, Wayne Wu, Lei Yang, Jiaqi
801 Wang, Chen Qian, et al. Omniobject3d: Large-vocabulary 3d object dataset for realistic perception,
802 reconstruction and generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision*
803 *and Pattern Recognition*, pp. 803–814, 2023b. 3
- 804 Zike Wu, Pan Zhou, Xuanyu Yi, Xiaoding Yuan, and Hanwang Zhang. Consistent3d: Towards
805 consistent high-fidelity text-to-3d generation with deterministic sampling prior. *arXiv preprint*
806 *arXiv:2401.09050*, 2024. 3
807
- 808 Dejia Xu, Hanwen Liang, Neel P Bhatt, Hezhen Hu, Hanxue Liang, Konstantinos N Plataniotis,
809 and Zhangyang Wang. Comp4d: Llm-guided compositional 4d scene generation. *arXiv preprint*
arXiv:2403.16993, 2024a. 2, 4, 18

- 810 Sirui Xu, Ziyin Wang, Yu-Xiong Wang, and Liang-Yan Gui. Interdreamer: Zero-shot text to 3d
811 dynamic human-object interaction. *arXiv preprint arXiv:2403.19652*, 2024b. 1, 9
812
- 813 Yinghao Xu, Hao Tan, Fujun Luan, Sai Bi, Peng Wang, Jiahao Li, Zifan Shi, Kalyan Sunkavalli,
814 Gordon Wetzstein, Zexiang Xu, et al. Dmv3d: Denoising multi-view diffusion using 3d large
815 reconstruction model. *arXiv preprint arXiv:2311.09217*, 2023. 3
- 816 Yuyang Yin, Dejia Xu, Zhangyang Wang, Yao Zhao, and Yunchao Wei. 4dgen: Grounded 4d content
817 generation with spatial-temporal consistency. *arXiv preprint arXiv:2312.17225*, 2023. 4
818
- 819 Xu Yinghao, Shi Zifan, Yifan Wang, Chen Hansheng, Yang Ceyuan, Peng Sida, Shen Yujun, and
820 Wetzstein Gordon. Grm: Large gaussian reconstruction model for efficient 3d reconstruction and
821 generation, 2024. 3
- 822 Ye Yuan, Xueting Li, Yangyi Huang, Shalini De Mello, Koki Nagano, Jan Kautz, and Umar
823 Iqbal. Gavatar: Animatable 3d gaussian avatars with implicit mesh learning. *arXiv preprint*
824 *arXiv:2312.11461*, 2023. 4
- 825 Yifei Zeng, Yuanxun Lu, Xinya Ji, Yao Yao, Hao Zhu, and Xun Cao. Avatarbooth: High-quality and
826 customizable 3d human avatar generation. *arXiv preprint arXiv:2306.09864*, 2023. 3
827
- 828 Lvmin Zhang and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models.
829 *arXiv preprint arXiv:2302.05543*, 2023. 19
- 830 Mingyuan Zhang, Zhongang Cai, Liang Pan, Fangzhou Hong, Xinying Guo, Lei Yang, and Ziwei
831 Liu. Motiondiffuse: Text-driven human motion generation with diffusion model. *arXiv preprint*
832 *arXiv:2208.15001*, 2022. 1
- 833
- 834 Mingyuan Zhang, Xinying Guo, Liang Pan, Zhongang Cai, Fangzhou Hong, Huirong Li, Lei Yang,
835 and Ziwei Liu. Remodiffuse: Retrieval-augmented motion diffusion model. In *Proceedings of the*
836 *IEEE/CVF International Conference on Computer Vision*, pp. 364–373, 2023. 1, 2
- 837 Mingyuan Zhang, Daisheng Jin, Chenyang Gu, Fangzhou Hong, Zhongang Cai, Jingfang Huang,
838 Chongzhi Zhang, Xinying Guo, Lei Yang, Ying He, et al. Large motion model for unified
839 multi-modal motion generation. *arXiv preprint arXiv:2404.01284*, 2024. 1, 2
840
- 841 Yufeng Zheng, Xueting Li, Koki Nagano, Sifei Liu, Otmar Hilliges, and Shalini De Mello. A unified
842 approach for text-and image-guided 4d scene generation. *arXiv preprint arXiv:2311.16854*, 2023.
843 4
- 844 Xiaoyu Zhou, Xingjian Ran, Yajiao Xiong, Jinlin He, Zhiwei Lin, Yongtao Wang, Deqing Sun,
845 and Ming-Hsuan Yang. Gala3d: Towards text-to-3d complex scene generation via layout-guided
846 generative gaussian splatting. *arXiv preprint arXiv:2402.07207*, 2024. 2, 3
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

864 A VIDEO RESULTS

865
866 To better visualize the generated results, we offer an improved demonstration of our method through
867 rotated videos in the supplementary materials. To access this demonstration, please open the file
868 named “**index.html**” provided in the supplementary.
869

870 B IMPLEMENTATION DETAILS

871
872 Our network is built upon the official implementation of DreamGaussian4D (Ren et al., 2023) and
873 Threestudio (Guo et al., 2023b) (an open-source 3D generative project).
874

875 To ensure easy reproducibility, we first include all the hyperparameters for our 3D composition stage
876 in Tab. 3.
877

878 Table 3: **Hyper-parameters of AvatarGO - 3D composition stage.**

Camera setting	Camera distance range	2.
	Radius	2.0
	Elevation range	(-30, 30)
	FoV range	49.1
Render setting	Resolution for 0-120 epochs	(128, 128)
	Resolution for 120-240 iters	(256, 256)
	Resolution for 240-400 iters	(512, 512)
Diffusion setting	Guidance scale	7.5
	t range	(0.01, 0.97)
	Minimal step percent	0.01
	Maximal step percent	0.97
	$\omega(t)$	$\sqrt{\alpha_t}(1 - \alpha_t)$
Initialization	Rotation \mathcal{R}	<code>torch.normal(mean=[0.5, 0.5, 0.5, 0.5], std=0.1)</code>
	Translation \mathcal{T}	0.0
	Scale \mathcal{S}	<code>torch.normal(mean=1.0, std=0.3)</code>
Learning rate	Rotation \mathcal{R}	0.005
	Translation \mathcal{T}	0.005
	Scale \mathcal{S}	0.005
LLM-guided contact retargeting	threshold a	1e-7
Training objectives	λ_{DS}^*	1.0
Hardware	GPU	1 × NVIDIA A100 (80GB)

Table 4: Hyper-parameters of AvatarGO - 4D animataion stage.

Camera setting	Camera distance range	2.
	Radius	2.0
	Elevation range	(-30, 30)
	FoV range	49.1
Render setting	Resolution for 0-120 epochs	(128, 128)
	Resolution for 120-240 iters	(256, 256)
	Resolution for 240-400 iters	(512, 512)
Diffusion setting to calculate \mathcal{L}_{SDS}^*	Guidance scale	7.5
	t range	(0.01, 0.97)
	Minimal step percent	0.01
	Maximal step percent	0.97
	$\omega(t)$	$\sqrt{\alpha_t(1 - \alpha_t)}$
Diffusion setting to calculate \mathcal{L}_{SDS}	Guidance scale	7.5
	Guidance rescale	0.75
	t range	(0.02, 0.98)
	Minimal step percent	0.02
	Maximal step percent	0.98
	gradient clip	[0, 1.5, 2.0, 1000]
	gradient clip pixel	True
	gradient clip threshold	1.0
$\omega(t)$	$\sqrt{\alpha_t(1 - \alpha_t)}$	
Initialization	Rotation \mathcal{R}	[-0.16, -0.16, -0.16, 0.5]
	Translation \mathcal{T}	0.0
Learning rate	Rotation \mathcal{R}	0.001
	Translation \mathcal{T}	0.001
Training objectives	λ_{CA}	1e+3
	λ_{SDS}^*	1.0
	λ_{SDS}	1.0
Hardware	GPU	1 × NVIDIA A100 (80GB)

In the 4D animation stage, we apply HexPlane (Cao & Johnson, 2023) to produce features from point position \mathbf{x}_c and timestamp \mathbf{t} , followed by an MLP to predict the offset for Gaussian attributes, i.e., point location \mathbf{x} , scaling matrix s , rotation matrix R . Specifically, the HexPlane encoder lifts the inputs to a higher frequency dimension $F((\mathbf{x}_c, \mathbf{t})) \in \mathbb{R}^{128}$, while the MLP is set to the default in DreamGaussian4D with ResNet (He et al., 2016).

To further ensure easy reproducibility, we first include all the hyperparameters for our 4D animation stage in Tab. 4. The other hyper-parameters are set to be the default of DreamGaussian4D (Guo et al., 2023b).

C MORE EXPLANATION ON DESIGNING “OURS (VAR-A)” AND “OURS (VAR-B)”

“Ours (VAR-A)”: This is a version where we have disabled the Lang-SAM initialization in our 3D static compositional generation. Comparing this with our final method shows that without assistance from Lang-SAM, the diffusion model struggles to accurately interpret human-object images.

“Ours (VAR-B)”: While Comp4D (Xu et al., 2024a) separates 3D scenes into two components and applies trajectories to one component for compositional 4D generation, it leaves the other component static. This method is not suitable for our scenarios where both humans and objects are dynamic. Therefore, we design “Ours (Var-B)” by adopting the Comp4D strategy: allowing the object to follow a trajectory while the human moves independently. Specifically, we replace our correspondence-aware motion supervision, as defined in Eq. 16, with SDS supervision strategy via the video diffusion model used in Comp4D. Comparing this approach with our final method demonstrates that our correspondence-aware motion supervision more effectively preserves the relationship between humans and objects throughout the animation process.

972 D TRAINING COMPLEXITY

973
974 In our study, our results, detailed in both the main paper and the Appendix, involve training the 3D
975 stage for 400 epochs on a single NVIDIA A100 GPU, taking approximately 10 minutes. Similarly,
976 the 4D stage requires roughly 20 minutes of training on the same GPU. To compare with other
977 methods: 1) In the experiments for 3D compositional generation, HumanGaussian (Liu et al., 2023d)
978 demands approximately 2 hours to complete 3600 epochs; GraphDreamer (Gao et al., 2023) adopts
979 a two-stage training approach, with the coarse stage taking roughly 3 hours for 10000 epochs and
980 the fine stage requiring around 6 hours for 20000 epochs. 2) Additionally, in our experiments with
981 4D animation, DreamGaussian4D (Ren et al., 2023) completes training of their 3-stage network in
982 around 10 minutes; TC4D (Bahmani et al., 2024) demands approximately 1 hour for the first stage
983 over 10000 epochs, 3 hours for the second stage over 20000 epochs, and roughly 30 hours for the
984 third stage over 30000 epochs.

985 E 2D HUMAN-OBJECT INTERACTION IMAGE GENERATION

986 Because of the limited availability of human-object interaction images within the 2D dataset utilized
987 for training diffusion models, existing models encounter challenges in accurately capturing the spatial
988 dynamics and contact between humans and objects. This limitation is evident in Figure 7, where
989 we noticed that during the process of 2D image generation, the diffusion model would struggle to
990 create such images. This inadequacy significantly hampers the ability of diffusion models to generate
991 realistic 3D human-object interactions.
992
993



1006 **Figure 7: Example generation of human-object interaction images.** Images generated by pose-
1007 conditioned ControlNet Zhang & Agrawala (2023)
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

F DIRECT RIGGING OF 3D OBJECT AND HUMAN MODELS

We conducted experiments by directly positioning the 3D objects in a reasonable position relative to the humans. As shown in Fig. 8, without further adjustments such as rescaling or rotating, the relationships between humans and objects are not accurately depicted. Penetration issues will also exist in some examples. Even with manual adjustments, such as rescaling and rotating the 3D objects, significant human effort is required, and the interactions between humans and objects still lack accuracy. For instance, Fig. 8 illustrates that humans frequently appear with open hands, which fails to convincingly "hold" the objects and significantly undermines the user experience.

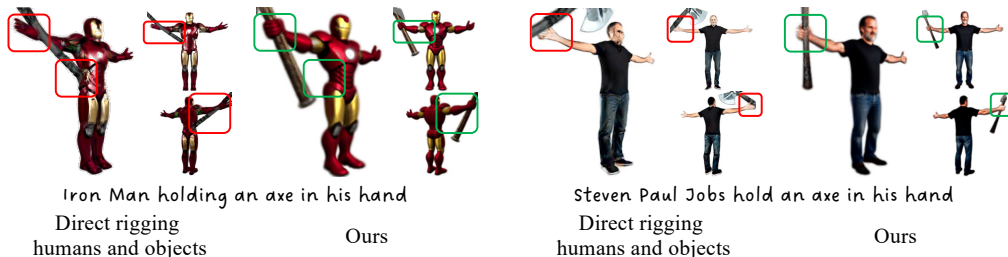


Figure 8: Evaluation by directly rigging humans and objects

G ANALYSIS BY DETERMINING THE ANIMATION OF OBJECT BY ONLY THE CONTACT PART

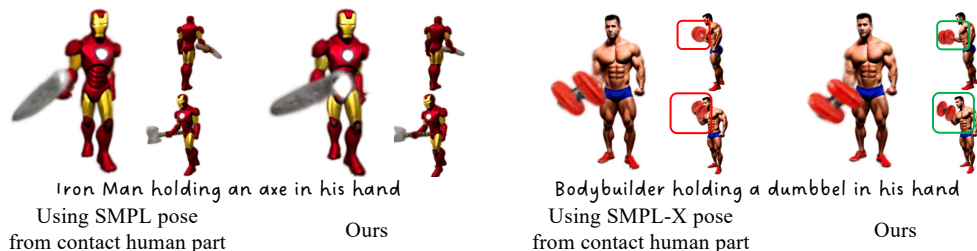


Figure 9: Evaluation by using SMPL-X pose from contact human part

We conducted experiments using the contact part of the human body to determine the object’s motion. The results are shown in Fig. 9. We found that this approach works well when the object is positioned far from the body, but it can encounter penetration issues when the object is close to the body (see "Bodybuilder holding a dumbbell in his hand"). We will incorporate this discussion into the updated paper.

H COMPARISONS WITH AVATARCRAFT, DREAMWALTZ AND DREAMAVATAR

In Fig. 10, we provide qualitative comparisons with AvatarCraft, DreamWaltz, and DreamAvatar. We observed that AvatarCraft and DreamAvatar are highly constrained by the SMPL prior model, making it difficult for them to create human models with effective object interactions. While DreamWaltz can generate some object interactions, these interactions are often inaccurate. Additionally, DreamWaltz has trouble maintaining proper interactions throughout the animation, as presented in Fig. 11.



Figure 10: Qualitative comparisons with DreamWaltz, AvatarCraft, and DreamAvatar



Figure 11: Evaluation on DreamWaltz's animated results

I SOCIETAL IMPACT.

The progress in 4D avatar generation with object interactions holds promise for numerous AR/VR applications, yet also raises concerns regarding potential misuse, such as creating misleading or nonexistent human-object pairings. We advocate for responsible research and deployment, promoting openness and transparency in practices to mitigate any potential negative consequences.

1134 J MORE COMPARISONS ON 3D GENERATION
 1135

1136 We provide more qualitative comparisons with HumanGaussian (Liu et al., 2023d), Graph-
 1137 Dreamer (Gao et al., 2023), and ‘Ours (Var-A)’ in Fig. 12. These results serve to reinforce
 1138 the claims made in Sec. 4 of the main paper, providing further evidence of the superior performance
 1139 of AvatarGO in compositing 3D human and object models.

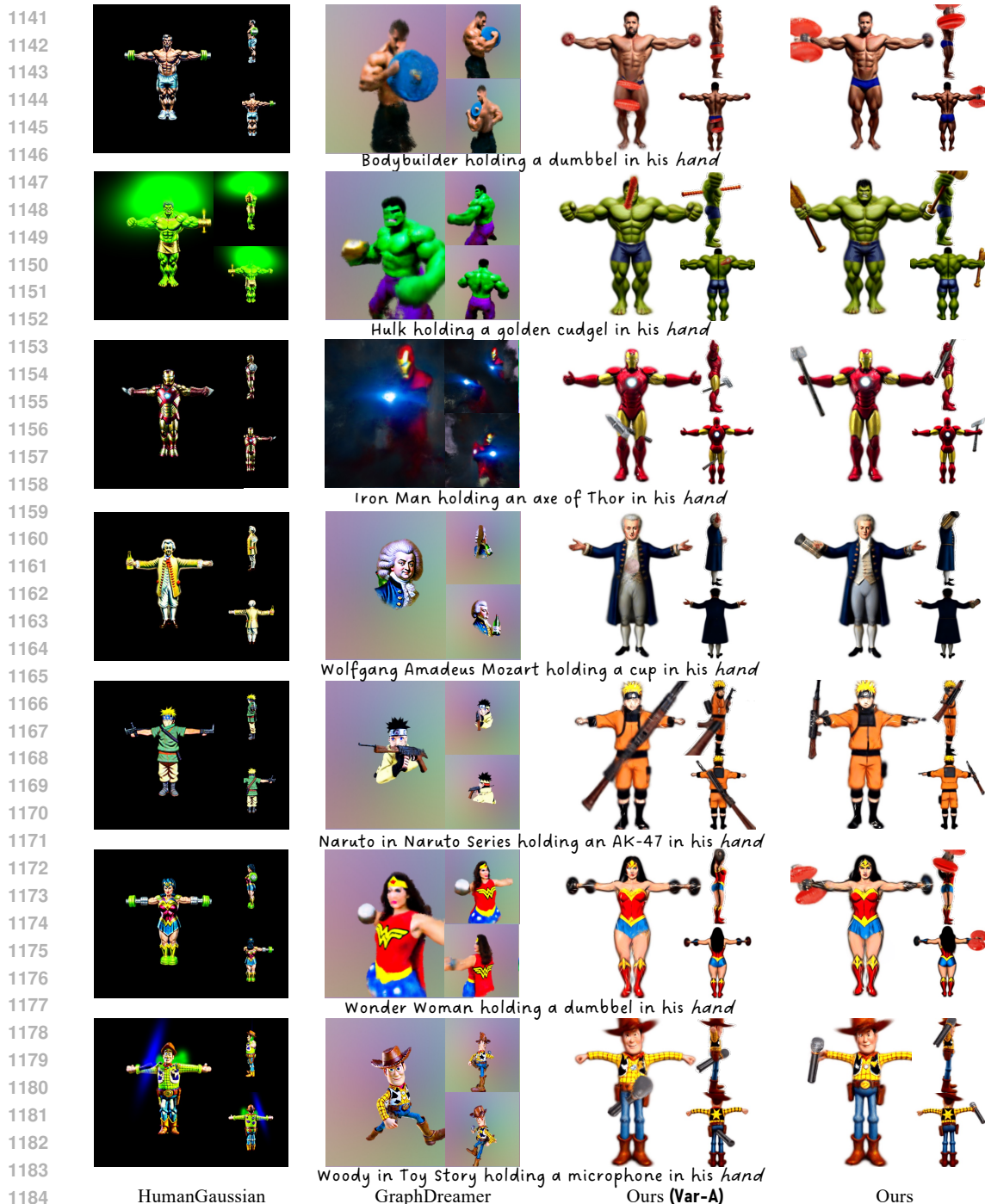


Figure 12: Comparisons on 3D compositional generations.

K MORE COMPARISONS ON 4D ANIMATION

We further provide more qualitative comparisons of 4D animation with DreamGaussian4D (Ren et al., 2023), HumanGaussian (Liu et al., 2023d), and “Ours (Var-B)”. The results can be found in Fig. 13. These comparisons further demonstrate the superiority of AvatarGO in maintaining the spatial correlation during animations and in addressing the penetration issues.

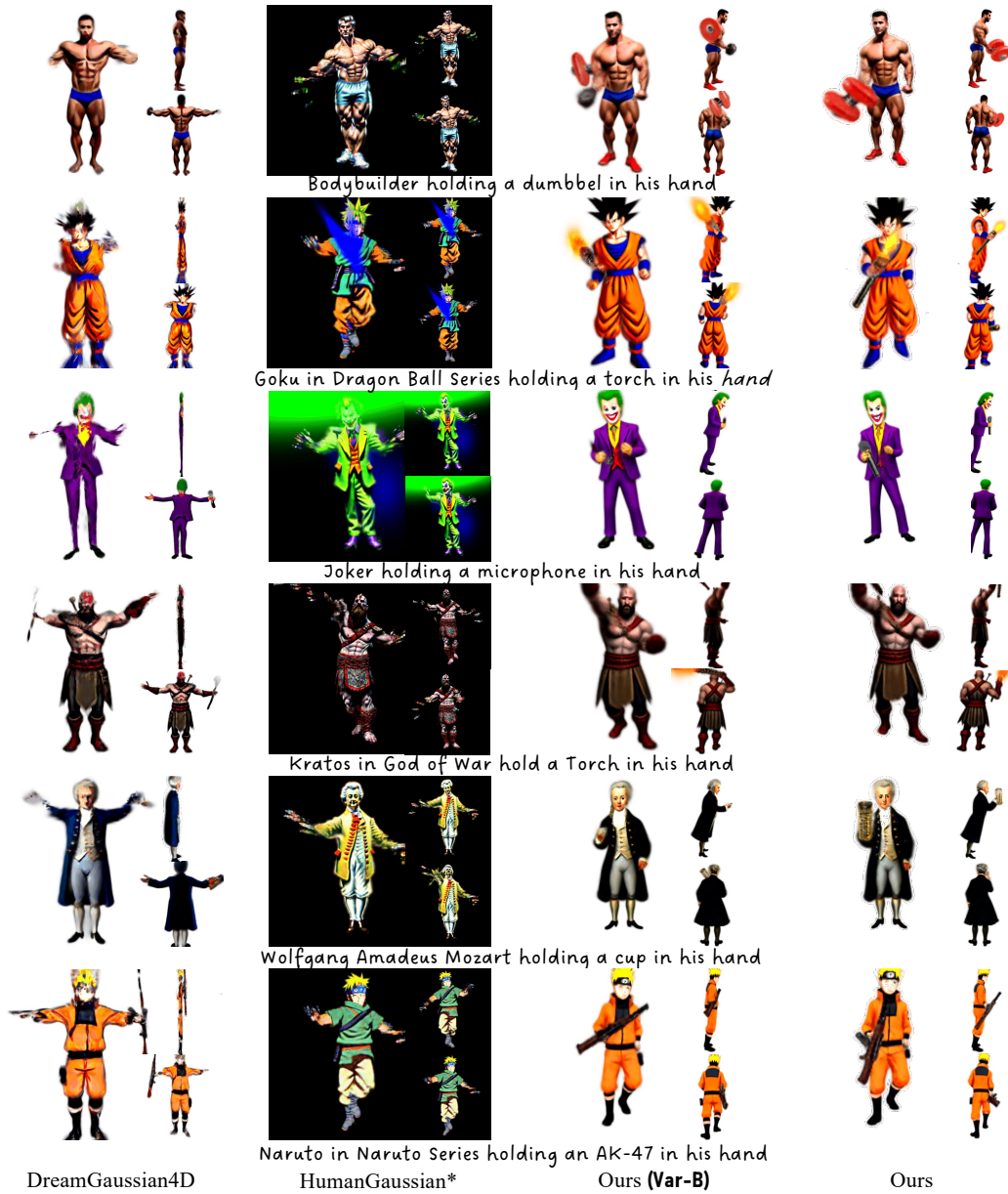


Figure 13: Comparisons on 4D animation.