

# DreamHOI: Subject-Driven Generation of 3D Human-Object Interactions with Diffusion Priors

Hanwen Zhu<sup>1,2†</sup> Ruining Li<sup>1\*</sup> Tomas Jakob<sup>1\*</sup>  
<sup>1</sup>University of Oxford <sup>2</sup>Carnegie Mellon University  
 thomaszh@cs.cmu.edu {ruining, tomj}@robots.ox.ac.uk  
[DreamHOI.github.io](https://DreamHOI.github.io)

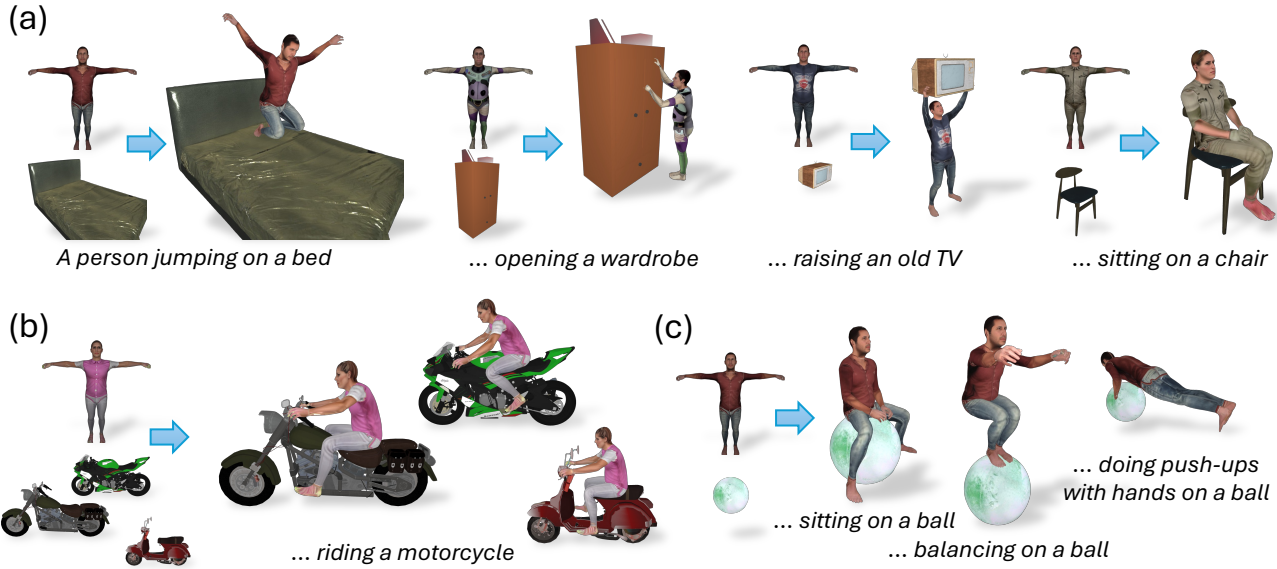


Figure 1. **Generated Human-Object Interactions by DreamHOI.** (a-c) DreamHOI takes as inputs a skinned human model, an object mesh, and a textual description of the intended interaction between them. It then poses the human model to create realistic interactions. (b) Given the same interaction description, the generated pose naturally conforms to the intricacies of the input object to be interacted with. (c) Given a fixed object, the generated poses vary faithfully to different intended interactions.

## Abstract

We present *DreamHOI*, a novel method for zero-shot synthesis of human-object interactions (HOIs), enabling a 3D human model to realistically interact with any given object based on a textual description. The complexity of this task arises from the diverse categories and geometries of real-world objects, as well as the limited availability of datasets that cover a wide range of HOIs. To circumvent the need for extensive data, we leverage text-to-image diffusion models trained on billions of image-caption pairs. We optimize the articulation of a skinned human mesh using *Score Distillation Sampling* (SDS) gradients obtained from these models, which predict image-space edits. However, directly backpropagating image-space gradients into com-

plex articulation parameters is ineffective due to the local nature of such gradients. To overcome this, we introduce a dual implicit-explicit representation of a skinned mesh, combining (implicit) neural radiance fields (NeRFs) with (explicit) skeleton-driven mesh articulation. During optimization, we transition between implicit and explicit forms, grounding the NeRF generation while refining the mesh articulation. We validate our approach through extensive experiments, demonstrating its effectiveness in generating realistic HOIs. The code can be found on the project page at <https://DreamHOI.github.io/>.

## 1. Introduction

The goal of this work is to create a method that can make existing 3D human models realistically interact with any given

<sup>†</sup>Work done while at University of Oxford.

<sup>\*</sup>Equal advising.

3D object, conditioned on a textual description of their interaction. For instance, given a 3D motorcycle and “riding” interaction, we aim to deform the 3D human model so that it realistically appears to ride the given motorcycle (Fig. 1). This task allows for automated population of virtual 3D environments with humans that naturally interact with objects, which has significant implications for industries such as movie and video game production as well as product advertisement.

The challenge in achieving open-world synthesis of human-object interaction (HOI) stems from the fact that the target deformation of the human model is influenced not only by the specified interaction but also by the actual geometry of the object. The riding posture on a cruiser motorcycle typically differs from that on a scooter or sports motorcycle, as shown in Figure 1b; moreover, individual cruiser motorcycles have distinct geometries, each necessitating specific deformations of the character to accommodate them. To train a model for this task in a *supervised* manner requires triplets of 3D objects, interactions, and target human poses. While a few such datasets have been recently collected [1, 18, 20, 49, 63, 67], methods trained on them [8, 36, 59] are constrained by the limited object and interaction coverage in these datasets. Constructing even larger datasets of immensely diverse real-world objects which humans can interact with is exceptionally difficult and expensive.

To bypass tedious data collection, we propose to distill the guidance for HOI synthesis from off-the-shelf text-to-image diffusion models that have been trained on large-scale image-caption pairs. These models can act as a “critic” and provide image-space gradients through Score Distillation Sampling (SDS) [38]. These gradients indicate how to modify the image to better align with the given text prompt and can be used to optimize the actual image parameters through backpropagation. This approach has proven effective in text-to-3D generation, where the underlying 3D representation, such as NeRF [32], is rendered differentially to obtain images for a text-to-2D diffusion model to “critique”. The local “edits” proposed by the 2D model are then backpropagated into the 3D representation space to optimize it [17, 25, 30, 38, 43, 54]. However, applying this approach to our problem is not straightforward.

In our case, the 3D scene consists of an input object mesh and a skinned human mesh. The goal is to pose the human, *i.e.*, to optimize the articulation parameters of the skinned human mesh. In theory, we can differentially render the two meshes together and apply SDS. However, this does not work well in practice (Sec. 4.4) due to the *local* nature of the image-space SDS gradients. Intuitively, the image-space gradients are *local* edits, indicating where in the image something should be added, removed, or slightly modified, which makes it challenging to directly derive the

*structural* changes required for the skeleton. Moreover, gradient-based optimization is prone to local optima. Updating the skeleton to its optimal location may necessitate human poses less aligned with the textual prompt during intermediate steps. Consequently, directly optimizing the articulation parameters using SDS often leads to convergence in sub-optimal poses.

To overcome this challenge, we draw inspiration from [16], which translates between an implicit pixel representation and an explicit skeleton-based human pose representation in 2D, and introduce a dual implicit-explicit 3D representation of the skinned human mesh. The implicit component is represented by a neural radiance field (NeRF) [32], while the explicit component comprises the input skinned mesh along with its articulation parameters, *i.e.*, bone rotations. To convert the implicit representation into the explicit one, we employ a regressor that predicts the bone rotations of the skinned mesh from multi-view renderings of the NeRF. Conversely, to revert to the implicit representation, we initialize the NeRF using the posed mesh. The implicit representation (*i.e.*, NeRF) then facilitates distillation from pre-trained diffusion models. However, it introduces a significant challenge in maintaining the identity of the articulated character during optimization. To address this issue, we periodically transition between the implicit and explicit representations. Specifically, after a number of optimization steps on the NeRF, we convert it to the explicit representation. By doing so, we ensure that the character’s identity is preserved, as we simply substitute the bone rotations of the skinned human with those predicted by the regressor, without changing the body shape or appearance. Subsequently, we re-initialize the NeRF using the posed human mesh and continue the optimization process.

We also make an important observation: while multi-view text-to-image diffusion models such as [43] are better at generating 3D assets than their single-view counterparts, we find that they tend to have an insufficient understanding of human-object interactions, potentially due to their extended fine-tuning on synthetic renderings. Therefore, we propose a simple and effective technique to combine the guidance from a multi-view text-to-image diffusion model with that from a state-of-the-art single-view model [19], leveraging their respective strengths.

To summarize, our contributions are:

1. We introduce a novel *zero-shot* approach dubbed DreamHOI for open-world synthesis of human-object interactions with a given human mesh, object mesh, and textual description of the interaction between them.
2. We propose a dual implicit-explicit representation of a skinned human mesh that allows us to harness the power of large text-to-image diffusion models to optimize its articulation parameters.
3. We demonstrate the effectiveness of our approach with

extensive qualitative and quantitative experiments.

## 2. Related Work

**3D Generation.** Early attempts on 3D generation [17, 24, 45, 48, 55, 60–62, 65] formulate the task as single/few-view reconstruction, and propose learning frameworks that regress the 3D model with various representations from the conditioning image(s). However, these models are category-specific and do *not* allow open-domain generation. Recent works have explored open-domain 3D generation using pre-trained image/video generators. Such generators, powered by diffusion models [15, 19, 40], are expected to acquire an implicit 3D understanding through their pre-training on Internet-scale images/videos. DreamFusion [38] and follow-ups [25, 30, 50, 54, 57] distill 2D diffusion models to extract 3D models from text or image prompts. To improve the generation process, authors have proposed to fine-tune pre-trained generators on datasets of multi-view images [11, 21, 27, 31, 43, 53, 70]. However, since these datasets consist mostly of scenes with single objects, the fine-tuned models often struggle with generating compositions and interactions between objects. In this work, we build upon both single-view and multi-view diffusion models to enable the generation of human-object interactions (HOIs) in 3D.

**Compositional Generation.** Prior work [26] has shown that existing image generators tend to struggle with composing multiple concepts into a single image. While efforts have been made to enhance compositionality in image generative models [9, 26], such techniques are not easily portable to distillation and hence 3D generation. Recent works instead extend the distillation approach to compositional 3D generation. Po and Wetzstein [37] introduce locally conditioned distillation, masking the 2D diffusion model’s prediction according to user-specified bounding boxes to generate composition. CG3D [52] optimizes individual objects and their relative rotation and translation via score distillation sampling (SDS, [38]), subject to additional physical constraints. Similarly, SceneWiz3D [66] proposes a hybrid representation to achieve scene generation, optimizing individual objects via SDS and a scene configuration. GALA3D [72] leverages a large language model to extract scene layout, which is then used to further refine the generated scene through layout-conditioned diffusion. ComboVerse [5] approaches compositional image-to-3D generation by leveraging a pre-trained single-view 3D reconstructor to generate the individual objects, followed by spatial relationship optimization using SDS. Concurrent to our work is InterFusion [6], which synthesizes human-object interactions based on textual prompts. This method selects an anchor pose from a codebook of poses in response

to a textual prompt. Subsequently, it synthesizes one neural radiance field (NeRF) representing a human body and another for an object, both designed to fit the anchor pose. However, these two NeRFs are optimized from scratch, offering no control over the identities of the human and the object. In contrast, our human-object interaction (HOI) generation pipeline takes existing 3D subjects (*i.e.*, the object mesh and the human mesh) as inputs and synthesizes the interaction with appropriate pose deformation. Our approach not only affords greater user control but is also more practical as it can be easily integrated into existing pipelines for the production of virtual 3D environments.

### **Data-Driven Human-object Interaction Synthesis.**

Since digital humans play an essential role in numerous applications ranging from AR/VR to gaming and movie production, prior works have considered *training* generative models of static human pose or dynamic human motion conditioned on an object or environment using curated HOI datasets, enabling HOI synthesis at test time. Earlier works [12, 68, 69] use a combination of affordance prediction and semantic heuristics to synthesize a static human pose. More recent works [8, 36, 59] learn human motion diffusion models [51] to generate a dynamic motion clip, conditioned on the encoded object or scene. The key difference between these training-based methods and our work is that our method extracts the HOI solely from a pre-trained diffusion model and does *not* require bespoke training data for this task. Consequently, our method can be easily extended to other categories by switching to a different deformation model, *e.g.*, using SMAL [73] for quadruped animals [73], whose motion data is not available at scale.

**Human Deformation Models.** To accurately reason about object pose and motion in images and videos, authors have developed deformation models capable of capturing the shape deformation space using low-dimensional latent codes [24, 28, 73, 74] or conditioning inputs [22, 23, 47]. The most popular model for human bodies is SMPL [28], which decomposes the 3D body into shape-dependent deformations, statistically learned from a large-scale 3D body scan dataset, and pose-dependent deformations, represented by joint rotations which are then used to deform body vertices with linear blend skinning [4, 29]. Numerous follow-ups have sought to improve this model. For instance, MANO [41], or SMPL+H, extends SMPL to also fit human hands. SMPL-X [35] goes one step beyond, facilitating the 3D model to capture both fully articulated hands and an expressive face. STAR [34] leverages a more compact shape representation and a larger human scan dataset than SMPL, which mitigates over-fitting and helps the model generalize better to new bodies. Meanwhile, orthogonal ef-

forts [2, 44, 56, 58, 64] have focused on training pose estimators which regress shape parameters of these body models from in-the-wild images and videos. In this work, we explore the possibilities of using these human body models and their estimators to allow high-quality HOI generation.

### 3. Method

Given an object mesh  $M_{\text{Obj}}$  (e.g., a 3D ball) and an intended interaction  $r$  as a textual prompt (e.g., “sit”), our goal is to automatically optimize the pose parameters  $\xi$  of a skinned human model  $M_\xi$ . The pose should realistically reflect the interaction  $r$  with the object (i.e., the character  $M_\xi$  is sitting on the ball  $M_{\text{Obj}}$ ) (Fig. 1). To achieve this, we propose a dual implicit-explicit representation of the skinned human model (Sec. 3.2), which enables us to leverage the image-space gradients from text-to-image diffusion models to optimize the 3D human’s pose parameters. We compose the human-object scene using the dual representation (Sec. 3.3) and apply a mixture of SDS guidance from both multi-view and single-view diffusion models (Sec. 3.4). We introduce additional regularizers to ensure the appropriate size of the human model and the consistency of the rendered scene (Sec. 3.5). Finally, we integrate all components and devise an iterative optimization procedure for the pose parameters  $\xi$  (Sec. 3.6). See Fig. 2 for an overview.

#### 3.1. Preliminaries

**Neural Radiance Fields (NeRFs).** A neural radiance field (NeRF [32]) represents a 3D scene as a function parametrized by  $\theta$ :

$$f_\theta : \boldsymbol{\mu} \mapsto (\mathbf{c}, \tau).$$

This function maps a 3D location  $\boldsymbol{\mu} = (x, y, z)$  to a color  $\mathbf{c} = (r, g, b)$  and a volume density  $\tau \geq 0$ <sup>1</sup>. NeRFs can be rendered differentially by aggregating the color of corresponding locations, and hence optimized via gradient descent using posed multi-view images.

Specifically, for each pixel  $\mathbf{u}$ , we sample  $N$  points  $\{\boldsymbol{\mu}_i\}_{i=1}^N$  along the ray cast from the camera  $\boldsymbol{\mu}_c$  to the pixel  $\mathbf{u}$ . We then obtain  $(\mathbf{c}_i, \tau_i) = f_\theta(\boldsymbol{\mu}_i)$ , sorted by distance to the camera  $d_i = \|\boldsymbol{\mu}_i - \boldsymbol{\mu}_c\|$ . Finally, we obtain the color  $\mathbf{c}$  of the pixel by alpha-compositing along the ray, as described in [32]:

$$\mathbf{c} = \sum_i w_i \mathbf{c}_i, \quad w_i = \alpha_i \prod_{j<i} (1 - \alpha_j), \quad (1)$$

$$\alpha_i = 1 - \exp(-\tau_i \|\boldsymbol{\mu}_i - \boldsymbol{\mu}_{i+1}\|). \quad (2)$$

**Score Distillation Sampling (SDS).** Score Distillation Sampling (SDS [38]) leverages pre-trained text-to-image

<sup>1</sup>In this work, for simplicity, the RGB color is *not* view-dependent.

diffusion models such as Stable Diffusion [40] to align a NeRF with a textual prompt. Assuming the learned denoiser  $\hat{\epsilon}(\mathbf{x}_t; y, t)$  which predicts the sampled noise  $\epsilon$  given the noisy 2D image  $\mathbf{x}_t$ , noise level  $t$  and textual prompt  $y$ , SDS computes the gradient:

$$\nabla_\theta \mathcal{L}_{\text{SDS}}(\mathbf{x}) = \mathbb{E}_{t,\epsilon} \left[ w(t) (\hat{\epsilon}(\mathbf{x}_t; y, t) - \epsilon) \frac{\partial \mathbf{x}}{\partial \theta} \right] \quad (3)$$

with respect to the NeRF parameters  $\theta$ . Here, the clean image  $\mathbf{x}$  is obtained via NeRF rendering, and  $w(t)$  is a weighting function. DreamFusion [38] achieves text-to-3D generation by iteratively updating  $\theta$  using Eq. (3) from a randomly initialized NeRF.

**Skinned 3D models.** A skinned 3D mesh  $M_\xi = (V \in \mathbb{R}^{|V| \times 3}, E)$  is parametrized by  $B$  bone rotations  $\xi = \{\xi_b\}_{b=1}^B \in SO(3)$ . To deform  $|V|$  vertices driven by  $B$  bones, we define rest-pose joint locations  $\mathbf{J} \in \mathbb{R}^{B \times 3}$ , a kinematic tree structure  $\pi$  defining a parent  $\pi(b)$  for each bone  $b$ , and skinning weights  $W \in \mathbb{R}^{|V| \times B}$ . The vertices are then posed by the linear blend skinning equation:

$$\hat{V}_i^{\text{posed}}(\xi) = \left( \sum_{b=1}^B W_{ib} G_b(\xi) G_b(\xi^*)^{-1} \right) \hat{V}_i,$$

$$G_{\text{root}} = g_{\text{root}}, \quad G_b = g_b \circ G_{\pi(b)}, \quad g_b(\xi) = \begin{bmatrix} R_{\xi_b} & \mathbf{J}_b \\ 0 & 1 \end{bmatrix},$$

where  $\xi^*$  denotes the bone rotations at the rest pose<sup>2</sup>. In practice,  $\mathbf{J}$  and  $W$  are either hand-crafted or learned. In this work, we adopt SMPL+H [41] as the skinned model for human body.

#### 3.2. Implicit-Explicit Skinned Mesh Representation

Given a skinned human mesh  $M_\xi$ , where  $\xi$  denotes bone rotations, our goal is to optimize these parameters so that the resulting posed mesh conforms to the desired interaction with a given object. As demonstrated in Sec. 4.4 and discussed in Appendix B, optimizing the pose parameters  $\xi$  directly by backpropagating image-space SDS gradients is not feasible due to the uninformative nature of these gradients. This issue is particularly pronounced for complex poses, such as those of humans. To address this challenge, we propose a dual implicit-explicit representation of a skinned mesh. The implicit representation consists of a NeRF and serves as a proxy to the skinned mesh, which can be efficiently optimized using image-space SDS gradients. In the following, we detail the process of translating from the implicit NeRF representation to the posed mesh and back.

<sup>2</sup>The  $\hat{\cdot}$  indicates  $\hat{V}_i$  and  $\hat{V}_i^{\text{posed}}$  are in homogeneous coordinates.

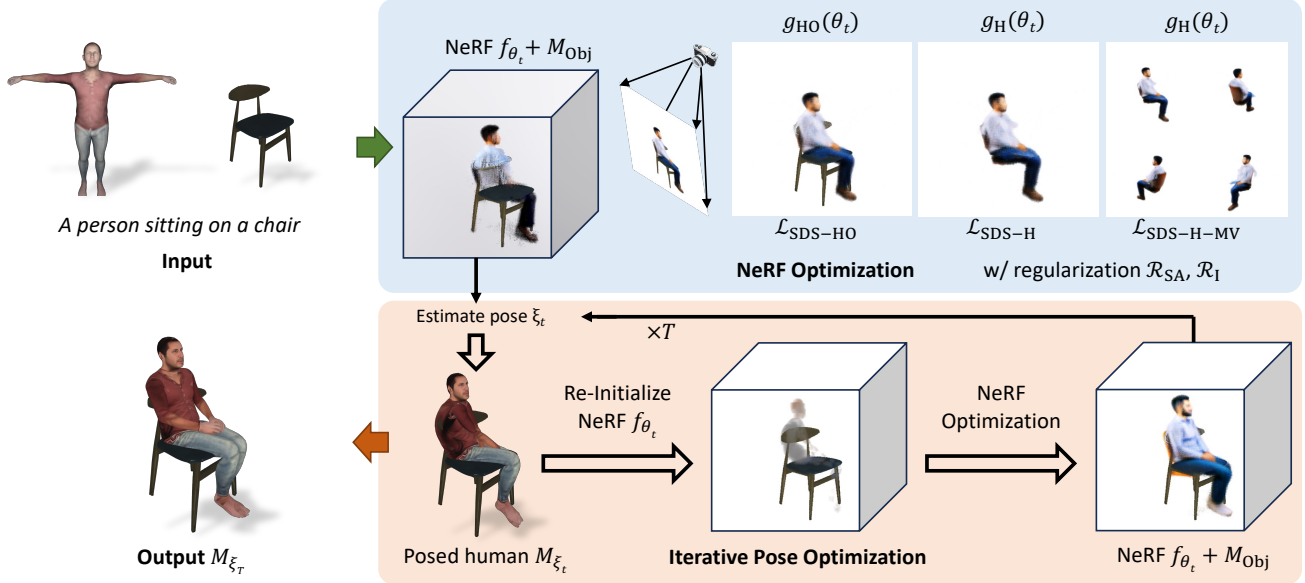


Figure 2. **Overview of DreamHOI.** Our method takes a human identity (in the form of a skinned body mesh) and an object mesh  $M_{\text{Obj}}$  (e.g., a 3D chair), together with their intended interaction (as a textual prompt, e.g., “sit”), as input. It first fits a NeRF  $f_{\theta_0}$  for the human using a mixture of diffusion guidance and regularizers, and then estimates its pose  $\xi_0$ . The posed human mesh  $M_{\xi_t}$  is used to re-initialize and further optimize the NeRF  $f_{\theta_t}$ , for iterations  $t \leq T$ . The final output is the posed human  $M_{\xi_T}$  at the last iteration. See Sec. 3.

**Translation from NeRF to Posed Mesh.** Given a NeRF  $f_{\theta}$  that represents a human, we render a set of images  $x_i$  of  $f_{\theta}$  from multiple camera viewpoints  $c_i$ . Using these multiview images  $x_i$ , we employ a pose estimator [71] to estimate the bone rotations  $\xi$  of the skinned mesh  $M_{\xi}$ . We note that our framework can, in principle, be applied to other categories, provided that a pose estimator for the skinned mesh is available, such as in the case of dogs [42, 73].

**Translation from a Posed Mesh to NeRF.** Given a posed mesh  $M_{\xi}$ , we construct a NeRF  $f_{\theta}$  by fitting  $f_{\theta}$  on 1k randomly sampled points  $\mu \sim \mathcal{N}(0, I_3)$ . Recall that  $f_{\theta}(\mu) = (c, \tau)$  maps a 3D location  $\mu$  to its color  $c$  and volume density  $\tau \geq 0$ . Specifically, the target color at each sampled point  $\mu$  is the color of the mesh vertex of  $M_{\xi}$  that is closest to  $\mu$ ; the target density is  $\infty$  if  $\mu$  is inside  $M_{\xi}$  and 0 if outside. We optimize  $\theta$  in a supervised manner for 10k iterations, taking about a minute.

### 3.3. Human-Object Scene Representation

The scene consists of the object mesh  $M_{\text{Obj}}$  and the dual implicit-explicit skinned mesh representation of the human which can take the form of either a skinned mesh  $M_{\xi}$  or a NeRF  $f_{\theta}$ , with the NeRF being used during optimization. To render the human-object scene  $x_{\text{HO}} = g_{\text{HO}}(\theta)$  consisting of the mesh  $M_{\text{Obj}}$  and the NeRF  $f_{\theta}$  we extend the standard volumetric rendering described in Section 3.1, which we denote as  $g_{\text{HO}}$ , to accommodate joint rendering with the

mesh.

Recall that to render a pixel  $u$ , we sample  $N$  points  $\{\mu_i\}_{i=1}^N$  along the ray and obtain  $P = \{(c_i, \tau_i)\}_{i=1}^N$  by evaluating  $f_{\theta}$ . In the case where the ray intersects the object mesh  $M_{\text{Obj}}$ , to account for the effect of  $M_{\text{Obj}}$  on pixel  $u$ , we insert an additional value  $(c, \tau = \infty)$  into  $P$  at the (first) intersection point  $\mu$ , where  $c$  is the color of  $M_{\text{Obj}}$  at  $\mu$ . The density value  $\infty$  reflects the opaque nature of  $M_{\text{Obj}}$ . We then apply the same alpha compositing method as in Eq. (2).

To enforce disentanglement between the object and the human, we render the NeRF  $f_{\theta}$  separately to obtain a human-only view  $x_{\text{H}} = g_{\text{H}}(\theta)$ . We then apply a guidance on  $x_{\text{H}}$  to ensure it contains a *single complete* human body, preventing the NeRF from generating extra limbs or additional people that might be occluded by the object in the full scene rendering  $x_{\text{HO}}$  (see Sec. 4.5 for its effect).

### 3.4. Guidance Mixture

Multi-view text-to-image diffusion models are superior in generating 3D assets [11, 21, 31, 43, 53, 70] when compared to their single-view counterparts [30, 38, 54]. Yet due to extensive fine-tuning on object-centric renderings, these multi-view models often struggle to capture the intended interaction  $r$  between humans and objects. Conversely, state-of-the-art single-view diffusion models [10, 19] exhibit a more fine-grained understanding of the conditioning prompts, including the textual description of the interaction. However, relying solely on them leads to 3D inconsistencies, including the multi-face Janus problem, where a per-

son is depicted with multiple faces to accommodate different viewing angles. To address this, we integrate the SDS guidance from both a multi-view model, MVDream [43], and a larger single-view model, DeepFloyd IF [19].

Specifically, we apply the SDS guidance using DeepFloyd IF on the scene rendering  $x_{\text{HO}}$  using the prompt  $y_{\text{HO}} = \text{“a photo of a person } \{r\} \text{ a } \{M_{\text{Obj}}\}, \text{ high detail, photography”}$ , where  $\{r\}$  is the intended interaction (e.g., “sitting on”) and  $\{M_{\text{Obj}}\}$  is the category of the object mesh (e.g., “ball”). We denote this loss as  $\mathcal{L}_{\text{SDS-HO}}$ . For the human-only rendering  $x_{\text{H}}$ , we employ a blend of SDS guidance using both DeepFloyd IF and MVDream, with the prompt  $y_{\text{HO}} = \text{“a photo of a person, high detail, photography”}$ . These losses are denoted as  $\mathcal{L}_{\text{SDS-H}}$  and  $\mathcal{L}_{\text{SDS-H-MV}}$ , respectively. We show that incorporating the MVDream guidance significantly enhances generation quality in Sec. 4.5. Note that we avoid conditioning on the interaction or the object for  $\mathcal{L}_{\text{SDS-H}}$  and  $\mathcal{L}_{\text{SDS-H-MV}}$  to ensure that the NeRF represents only a *single* person, as introducing additional objects can reduce the accuracy of the pose estimator when converting to the explicit mesh representation.

### 3.5. Regularizers

To further encourage realistic HOI generations, we introduce two regularizers. We define the *sparsity above threshold regularizer*:

$$\mathcal{R}_{\text{SA}} := \text{softplus}(\bar{x}_{\text{H}} - \eta), \quad (4)$$

where  $\bar{x}_{\text{H}}$  is the average opacity of the human-only rendering  $x_{\text{H}}$ . This regularizer controls the “size” of the human to ensure it occupies no more than  $\eta := 20\%$  of the camera field of view. The purpose of  $\mathcal{R}_{\text{SA}}$  is to compel the NeRF  $f_{\theta}$  to contain only the human; without  $\mathcal{R}_{\text{SA}}$ , we observe that the NeRF tends to expand and cover the object mesh  $M_{\text{Obj}}$  (see 4.5). To ensure robustness, we adjust the camera distance based on its randomly sampled focal length, ensuring that the 3D unit cube consistently occupies the same area in the renderings regardless of the focal length.

The *intersection regularizer*  $\mathcal{R}_{\text{I}}$  computes the average density  $\tau$  (as predicted by  $f_{\theta}$ ) of all ray points  $\mu$  *inside* the object mesh  $M_{\text{Obj}}$ . Informally, it measures the volume of intersection between the human NeRF and the object.  $\mathcal{R}_{\text{I}}$  discourages the model from generating body parts or other objects inside  $M_{\text{Obj}}$ , which would be invisible in  $x_{\text{HO}}$ .

### 3.6. Iterative Pose Optimization

Our objective function is a weighted sum of SDS guidances and regularizers:

$$\begin{aligned} \mathcal{L}(\theta) = & \lambda_{\text{SDS-HO}} \mathcal{L}_{\text{SDS-HO}}(g_{\text{HO}}(\theta)) + \lambda_{\text{SDS-H}} \mathcal{L}_{\text{SDS-H}}(g_{\text{H}}(\theta)) \\ & + \lambda_{\text{SDS-H-MV}} \mathcal{L}_{\text{SDS-H-MV}}(g_{\text{H}}(\theta)) \\ & + \lambda_{\text{SA}} \mathcal{R}_{\text{SA}}(g_{\text{H}}(\theta)) + \lambda_{\text{I}} \mathcal{R}_{\text{I}}(\theta). \end{aligned} \quad (5)$$

We initialize the NeRF  $f_{\theta_0}$  with a density bias centered at the origin. Throughout the optimization process, the implicit NeRF is periodically converted (every 10k optimization steps) into its explicit form as a posed human mesh  $M_{\xi_t}$  for identity grounding. The NeRF is then re-initialized from  $M_{\xi_t}$  and further refined by minimizing the loss  $\mathcal{L}(\theta)$ .

This iterative conversion between implicit and explicit representations is repeated for  $T$  times throughout the optimization. After the first NeRF re-initialization, we disregard the human-only SDS terms  $\mathcal{L}_{\text{SDS-H}}$  and  $\mathcal{L}_{\text{SDS-H-MV}}$  from the loss  $\mathcal{L}(\theta)$  since the human is already well-positioned, significantly reducing computational overhead. We also decrease the noise level and render at a higher resolution (see Sec. 4.2). The final output is the posed human mesh  $M_{\xi_T}$ .

## 4. Experiments

We conduct experiments on the subject-driven generation of 3D human-object interactions task. We use the SMPL+H model [28, 41] as the underlying rigged human model. We evaluate our proposed pipeline against several baselines and ablation settings.

### 4.1. Dataset

We select meshes suitable for human interaction from Sketchfab [46], a website that hosts CC-licensed 3D models uploaded by the community. We exclude meshes uploaded *before* the training cutoff date of DeepFloyd IF and MVDream to eliminate data contamination. We then use GPT-4 [33] to generate a set of prompts for the chosen meshes corresponding to a variety of human-object interactions, and filter out prompts that are too similar or impossible (e.g., “juggling balls” as there is only a single ball mesh), yielding a set of 12 prompts for comparisons and ablation studies. We manually position and scale  $M_{\text{Obj}}$  for each prompt such that it is natural to generate a human interacting with  $M_{\text{Obj}}$ . We source a set of skinned human models as input to our pipeline from SMPLitex [3].

### 4.2. Implementation Details

We estimate the SMPL+H [41] pose parameters  $\xi$  from NeRF  $f_{\theta}$  by using rendered images  $x_i$  from different camera angles  $c_i$  as in SMPLify-X [35]. This is done by optimizing the parameters  $\xi$  such that the resulting keypoints, when projected to 2D space by  $c_i$ , match with keypoints predicted from  $x_i$  using OpenPose [2, 44, 58]. Our implementation is adapted from an extension [71] of SMPLify-X that can train on multiple images.

We set  $\lambda_{\text{SDS-HO}} = 0.9$ ,  $\lambda_{\text{SDS-H}} = 0.05$ , and  $\lambda_{\text{SDS-H-MV}} = 0.05$ . For the first stage (i.e., before any implicit-explicit conversion and NeRF re-initialization), we set  $\lambda_{\text{I}} = 1$  and  $\lambda_{\text{SA}} = 10000$  to enforce a strict volume of the NeRF. After the NeRF re-initialization,  $\lambda_{\text{SA}}$  is decreased to 1000 for more flexibility. The SDS noise level (i.e., timestep)

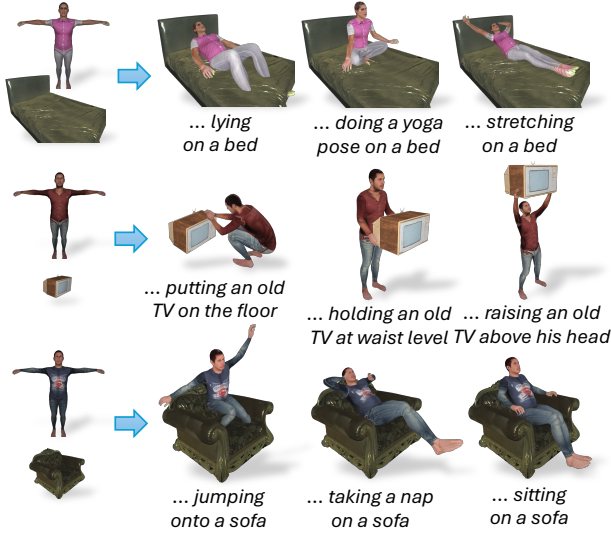


Figure 3. **Additional Results.** We demonstrate DreamHOI’s ability to control the pose based on different textual conditions.

is sampled from  $\mathcal{U}(0.02, 0.98)$  during the first stage and  $\mathcal{U}(0.02, 0.70)$  subsequently. These hyper-parameters apply to all 12 prompts and we do *not* tune them individually for each prompt.

During the first stage, we render  $64 \times 64$  images for the former half of the training and  $256 \times 256$  images for the latter half, following MVDream [43]. After NeRF re-initialization, the resolution is increased to  $512 \times 512$  to achieve finer quality, as in Magic3D [25]. We do not use shading for NeRF [38] as it is costly to compute.

### 4.3. Results

We present qualitative results in Fig. 1 and Fig. 3. We find DreamHOI can adapt to a wide variety of human-object interactions, including those that are complex or rare (e.g., “doing push-ups with hands on a ball”). It also automatically recognizes differences in object geometry (e.g. a sports motorcycle vs. a cruiser, Fig. 1) without mentioning the difference in the prompt (“riding a motorcycle”) and adapts the human pose accordingly. Given the same object, it can also respond to differences in the prompts (“lying” vs. “stretching” on a bed, Fig. 3). The open-world nature of these results contrasts with traditional methods for human-object interaction that only work on fixed categories of objects and a limited number of interactions as input.

### 4.4. Comparisons

We compare DreamHOI with baseline open-world approaches. In particular, we consider (1) using DreamFusion SDS loss [38] to fit a NeRF, and (2) using DreamFusion to fit SMPL+H [28, 41] pose parameters directly, as the baselines. Both baselines use DeepFloyd IF [19] as the guidance

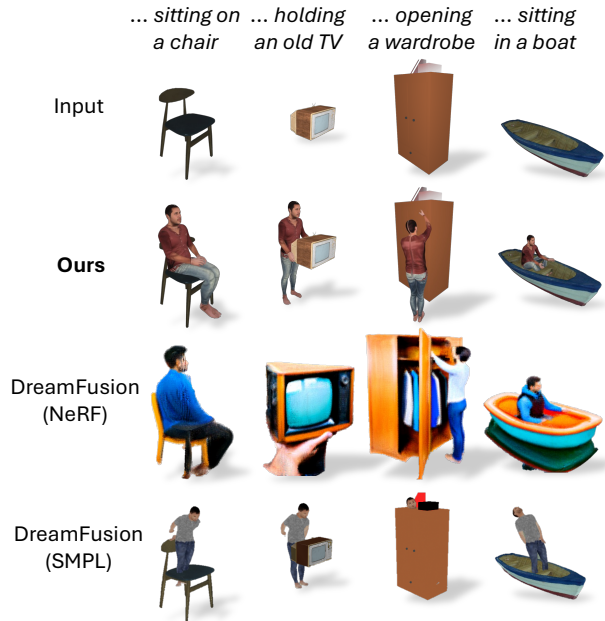


Figure 4. **Comparison** with baselines: (**third row**) using DreamFusion to optimize a NeRF with the given object mesh inserted; (**last row**) using DreamFusion to optimize the pose parameters of the skinned human mesh directly. See Sec. 4.4 for discussions.

model for SDS.

We show qualitative comparison in Fig. 4. We first note that DreamFusion (NeRF) does *not* achieve our goal of *subject-driven* human-object interaction generation, as it *cannot* preserve the identity and structure of the human mesh. It has limited success on some prompts (first and last columns) and also fails on others (second and third columns). We find that DreamFusion frequently generates a NeRF that is overly large and overtakes the object mesh, and its outputs are not view-consistent (e.g., the three-faced TV). DreamFusion (SMPL), on the other hand, completely fails to produce correct poses. This occurs because the image-space gradients provided by SDS *cannot* be effectively backpropagated into meaningful updates in the bone rotation space. For example, when adjusting the position of an arm, SDS does not provide a gradient that gradually shifts the arm. Instead, it deletes the arm and re-generates it in a new position, which may be far from the original. As a result, no gradients from the re-generated arm propagate through differentiable mesh rendering to the bone rotation parameters.

To compare quantitatively, we render our final scenes, which consist of the object mesh and the posed human mesh, from multiple viewpoints. We then measure the average CLIP similarity [13, 39] between our prompts and the renderings of the composed scene. Our method consistently outperforms the baselines as shown in Tab. 1.

Table 1. **Quantitative Comparisons and Ablations.** We compute the mean and standard deviation of CLIP similarities between the renderings of generated HOIs and the corresponding text prompts.

Method	CLIP Score $\uparrow$
DreamFusion (NeRF)	0.2915 $\pm$ 0.0182
DreamFusion (SMPL)	0.2511 $\pm$ 0.0247
<b>DreamHOI (Ours)</b>	<b>0.2932</b> $\pm$ 0.0239
Remove $\mathcal{R}_{SA}, \mathcal{R}_I$	0.2879 $\pm$ 0.0315
Remove $\mathcal{L}_{SDS-H}, \mathcal{L}_{SDS-H-MV}$	0.2869 $\pm$ 0.0306

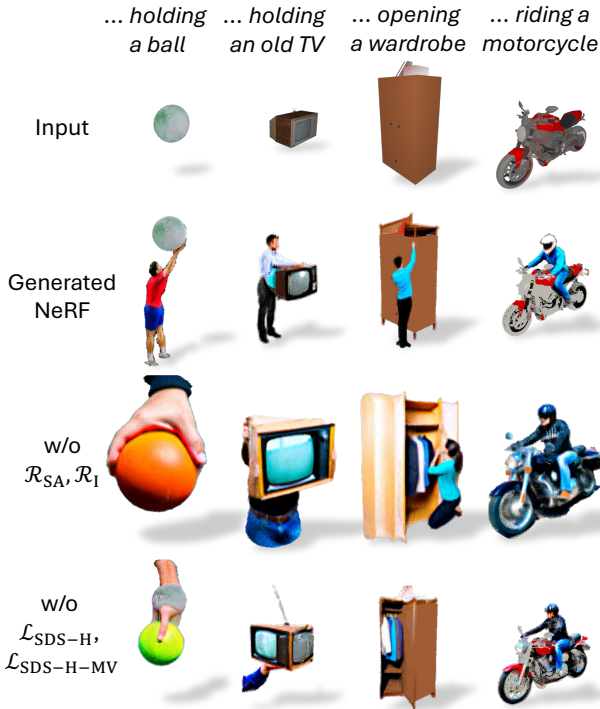


Figure 5. **Ablations.** We visualize output of first-round NeRF optimization (*i.e.*,  $f_{\theta_0}$ ). (**Row 2 vs. row 3**) Our regularizers (Sec. 3.5) effectively prevent the NeRF from encroaching upon and overriding the mesh. (**Row 2 vs. row 4**) The human-only SDS motivates a complete human while MVDream enhances view consistency (note the TV and wardrobe in row 4 have multiple faces).

#### 4.5. Ablations

To demonstrate the effectiveness of our techniques, we test our pipeline in the following ablation settings: (1) without regularizers  $\mathcal{R}_{SA}, \mathcal{R}_I$ ; and (2) without the human-only rendering  $g_H$  (*i.e.*, removing  $\mathcal{L}_{SDS-H}, \mathcal{L}_{SDS-H-MV}$ ). We present qualitative results in Fig. 5 and quantitative results in Tab. 1.

The regularizers  $\mathcal{R}_{SA}, \mathcal{R}_I$  effectively limit the size of the generated NeRF of the human; without them, the NeRF expands to include both the human and the object, covering and overriding the existing object mesh  $M_{Obj}$ , resulting in wrong human and object positions (Fig. 5 third row). On



Figure 6. **Effect of Iterative Pose Optimization.** When the initial human pose is sub-optimal, by translating the explicit skinned mesh back to the implicit NeRF representation and optimizing the NeRF further, the pose improves significantly. Note the movement of the feet toward the deck and the hands toward the handle.

the other hand, SDS on the human-only rendering  $g_H$  ensures the presence of a complete human body, as evidenced by the absence of the person in the wardrobe example or the presence of only an arm in the ball and TV examples (Fig. 5 last row). Furthermore, the MVDream [43] guidance enforces directional consistency, thereby eliminating the multi-face Janus problem, as seen from the three-faced TV and wardrobe in the last row. MVDream takes camera positions as input, and since most objects in its training data (Objaverse [7] renderings) face the forward direction, this also gives a canonical “front” direction for generating the HOI. See Sup. Mat. for details.

In Fig. 6, we visualize the effect of iterative pose optimization (Sec. 3.6) on a sub-optimal initial NeRF: the dual implicit-explicit optimization process improves the body pose. The improvement is attributed to the capability of SDS to effectively manipulate shapes represented as NeRFs, together with the re-initialization of NeRF through the corresponding posed mesh, which grounds the NeRF.

#### 4.6. Limitations

Similar to other 3D generation methods reliant on 2D guidance, the quality of our results is constrained by the ability of image diffusion models to accurately follow prompts. Future improvements in the underlying guidance models could lead to significant advancements in our approach.

DreamHOI depends on existing pose estimators (*i.e.*, OpenPose [2, 44, 58]) to regress the pose from the implicit NeRF. Consequently, our results depend on their ability to predict keypoints in the NeRF renderings. It is also limited by the training data and targets (*i.e.*, humans) of these estimators. We leave it as future work to develop an automatic method that can find correspondences between two articulated shapes to perform this alignment. This would enable us to extend DreamHOI to other categories where pose estimators are not readily available, while also circumventing the challenges associated with non-human data collection.

#### 5. Conclusions

We have proposed DreamHOI, an approach for posing rigged human models to interact naturally with given 3D



objects conditioned on textual prompts. At the core of DreamHOI is a dual implicit-explicit representation that enables the use of SDS to optimize the explicit pose parameters of skinned meshes, supplemented by a guidance mixture technique and novel regularization terms. We experiment on a diverse set of prompts and demonstrate our method achieves superior generation quality compared to various baseline approaches. Ablation studies verify the importance of the components that contributed to this improvement. Our work has the potential to simplify the creation of virtual environments populated by realistically interacting humans, which can be used in many applications, such as film and game production.

**Acknowledgments.** The authors would like to thank Lorenza Prospero for her helpful feedback on the manuscript.

## References

- [1] Bharat Lal Bhatnagar, Xianghui Xie, Ilya Petrov, Cristian Sminchisescu, Christian Theobalt, and Gerard Pons-Moll. Behave: Dataset and method for tracking human object interactions. In *CVPR*, 2022. 2
- [2] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In *CVPR*, 2017. 4, 6, 8
- [3] Dan Casas and Marc Comino-Trinidad. SMPLitex: A Generative Model and Dataset for 3D Human Texture Estimation from Single Image. In *BMVC*, 2023. 6
- [4] John E Chadwick, David R Haumann, and Richard E Parent. Layered construction for deformable animated characters. *ACM SIGGRAPH*, 1989. 3
- [5] Yongwei Chen, Tengfei Wang, Tong Wu, Xingang Pan, Kui Jia, and Ziwei Liu. Comboverse: Compositional 3d assets creation using spatially-aware diffusion guidance. In *ECCV*, 2024. 3
- [6] Sisi Dai, Wenhao Li, Haowen Sun, Haibin Huang, Chongyang Ma, Hui Huang, Kai Xu, and Ruizhen Hu. Interfusion: Text-driven generation of 3d human-object interaction. In *ECCV*, 2024. 3
- [7] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of annotated 3d objects. In *CVPR*, 2023. 8
- [8] Christian Diller and Angela Dai. Cg-hoi: Contact-guided 3d human-object interaction generation. In *CVPR*, 2024. 2, 3
- [9] Yilun Du, Shuang Li, and Igor Mordatch. Compositional visual generation with energy based models. *NeurIPS*, 2020. 3
- [10] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, Dustin Podell, Tim Dockhorn, Zion English, Kyle Lacey, Alex Goodwin, Yan-nik Marek, and Robin Rombach. Scaling rectified flow transformers for high-resolution image synthesis, 2024. 5
- [11] Ruiqi Gao, Aleksander Holynski, Philipp Henzler, Arthur Brussee, Ricardo Martin-Brualla, Pratul Srinivasan, Jonathan T Barron, and Ben Poole. Cat3d: Create anything in 3d with multi-view diffusion models. *arXiv preprint arXiv:2405.10314*, 2024. 3, 5
- [12] Mohamed Hassan, Partha Ghosh, Joachim Tesch, Dimitrios Tzionas, and Michael J Black. Populating 3d scenes by learning human-scene interaction. In *CVPR*, 2021. 3
- [13] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning, 2022. 7
- [14] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. 3
- [15] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, 2020. 3
- [16] Tomas Jakab, Ankush Gupta, Hakan Bilen, and Andrea Vedaldi. Self-supervised learning of interpretable keypoints from unlabelled videos. In *CVPR*, 2020. 2
- [17] Tomas Jakab, Ruining Li, Shangzhe Wu, Christian Rupprecht, and Andrea Vedaldi. Farm3D: Learning articulated 3d animals by distilling 2d diffusion. In *3DV*, 2024. 2, 3
- [18] Nan Jiang, Tengyu Liu, Zhexuan Cao, Jieming Cui, Zhiyuan Zhang, Yixin Chen, He Wang, Yixin Zhu, and Siyuan Huang. Full-body articulated human-object interaction. In *ICCV*, 2023. 2
- [19] DeepFloyd Lab. If. <https://github.com/deep-floyd/IF>, 2023. 2, 3, 5, 6, 7
- [20] Jiaman Li, Jiajun Wu, and C Karen Liu. Object motion guided human motion synthesis. *ACM TOG*, 2023. 2
- [21] Jiahao Li, Hao Tan, Kai Zhang, Zexiang Xu, Fajun Luan, Yinghao Xu, Yicong Hong, Kalyan Sunkavalli, Greg Shakhnarovich, and Sai Bi. Instant3D: Fast text-to-3D with sparse-view generation and large reconstruction model. In *ICLR*, 2024. 3, 5
- [22] Ruining Li, Chuanxia Zheng, Christian Rupprecht, and Andrea Vedaldi. Dragapart: Learning a part-level motion prior for articulated objects. In *ECCV*, 2024. 3
- [23] Ruining Li, Chuanxia Zheng, Christian Rupprecht, and Andrea Vedaldi. Puppet-master: Scaling interactive video generation as a motion prior for part-level dynamics. *arXiv preprint arXiv:2408.04631*, 2024. 3
- [24] Zizhang Li, Dor Litvak, Ruining Li, Yunzhi Zhang, Tomas Jakab, Christian Rupprecht, Shangzhe Wu, Andrea Vedaldi, and Jiajun Wu. Learning the 3d fauna of the web. In *CVPR*, 2024. 3
- [25] Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiao-hui Zeng, Xun Huang, Karsten Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. Magic3d: High-resolution text-to-3d content creation. In *CVPR*, 2023. 2, 3, 7
- [26] Nan Liu, Shuang Li, Yilun Du, Antonio Torralba, and Joshua B Tenenbaum. Compositional visual generation with composable diffusion models. In *ECCV*, 2022. 3
- [27] Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick. Zero-1-to-3: Zero-shot one image to 3d object. In *ICCV*, 2023. 3
- [28] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. SMPL: a skinned multi-person linear model. In *ACM TOG*, 2015. 3, 6, 7, 2

- [29] Nadia Magnenat-Thalmann, E Primeau, and Daniel Thalmann. Abstract muscle action procedures for human face animation. *TVC*, 1988. 3
- [30] Luke Melas-Kyriazi, Iro Laina, Christian Rupprecht, and Andrea Vedaldi. Realfusion: 360deg reconstruction of any object from a single image. In *CVPR*, pages 8446–8455, 2023. 2, 3, 5
- [31] Luke Melas-Kyriazi, Iro Laina, Christian Rupprecht, Natalia Neverova, Andrea Vedaldi, Oran Gafni, and Filippos Kokkinos. Im-3d: Iterative multiview diffusion and reconstruction for high-quality 3d generation. In *ICLR*, 2024. 3, 5
- [32] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020. 2, 4
- [33] OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 6
- [34] Ahmed A A Osman, Timo Bolkart, and Michael J. Black. STAR: A sparse trained articulated human body regressor. In *ECCV*, 2020. 3
- [35] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed A. A. Osman, Dimitrios Tzionas, and Michael J. Black. Expressive body capture: 3D hands, face, and body from a single image. In *CVPR*, 2019. 3, 6, 1
- [36] Xiaogang Peng, Yiming Xie, Zizhao Wu, Varun Jampani, Deqing Sun, and Huaizu Jiang. Hoi-diff: Text-driven synthesis of 3d human-object interactions using diffusion models. *arXiv preprint arXiv:2312.06553*, 2023. 2, 3
- [37] Ryan Po and Gordon Wetzstein. Compositional 3d scene generation using locally conditioned diffusion. In *3DV*, 2024. 3
- [38] Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. DreamFusion: Text-to-3d using 2d diffusion. In *ICLR*, 2023. 2, 3, 4, 5, 7
- [39] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021. 7
- [40] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 3, 4
- [41] Javier Romero, Dimitrios Tzionas, and Michael J. Black. Embodied hands: Modeling and capturing hands and bodies together. In *ACM TOG*, 2017. 3, 4, 6, 7
- [42] Nadine Rüegg, Shashank Tripathi, Konrad Schindler, Michael J. Black, and Silvia Zuffi. Bite: Beyond priors for improved three-d dog pose estimation. In *CVPR*, pages 8867–8876, 2023. 5
- [43] Yichun Shi, Peng Wang, Jianglong Ye, Mai Long, Kejie Li, and Xiao Yang. MVDream: Multi-view diffusion for 3D generation. In *ICLR*, 2024. 2, 3, 5, 6, 7, 8
- [44] Tomas Simon, Hanbyul Joo, Iain Matthews, and Yaser Sheikh. Hand keypoint detection in single images using multi-view bootstrapping. In *CVPR*, 2017. 4, 6, 8
- [45] Vincent Sitzmann, Justus Thies, Felix Heide, Matthias Nießner, Gordon Wetzstein, and Michael Zollhofer. Deep-voxels: Learning persistent 3d feature embeddings. In *CVPR*, 2019. 3
- [46] Sketchfab. Sketchfab. <https://sketchfab.com/>, 2024. 6
- [47] Robert W Sumner, Johannes Schmid, and Mark Pauly. Embedded deformation for shape manipulation. In *ACM SIGGRAPH*, 2007. 3
- [48] Stanislaw Szymanowicz, Christian Rupprecht, and Andrea Vedaldi. Splatter image: Ultra-fast single-view 3d reconstruction. In *CVPR*, 2024. 3
- [49] Omid Taheri, Nima Ghorbani, Michael J Black, and Dimitrios Tzionas. Grab: A dataset of whole-body human grasping of objects. In *ECCV*, 2020. 2
- [50] Jiayang Tang, Jiawei Ren, Hang Zhou, Ziwei Liu, and Gang Zeng. Dreamgaussian: Generative gaussian splatting for efficient 3d content creation. In *ICLR*, 2024. 3
- [51] Guy Tevet, Sigal Raab, Brian Gordon, Yoni Shafir, Daniel Cohen-or, and Amit Haim Bermano. Human motion diffusion model. In *ICLR*, 2022. 3
- [52] Alexander Vilesov, Pradyumna Chari, and Achuta Kadambi. Cg3d: Compositional generation for text-to-3d via gaussian splatting. *arXiv preprint arXiv:2311.17907*, 2023. 3
- [53] Vikram Voleti, Chun-Han Yao, Mark Boss, Adam Letts, David Pankratz, Dmitry Tochilkin, Christian Laforte, Robin Rombach, and Varun Jampani. Sv3d: Novel multi-view synthesis and 3d generation from a single image using latent video diffusion. *arXiv preprint arXiv:2403.12008*, 2024. 3, 5
- [54] Haochen Wang, Xiaodan Du, Jiahao Li, Raymond A Yeh, and Greg Shakhnarovich. Score jacobian chaining: Lifting pretrained 2d diffusion models for 3d generation. In *CVPR*, 2023. 2, 3, 5
- [55] Nanyang Wang, Yinda Zhang, Zhuwen Li, Yanwei Fu, Wei Liu, and Yu-Gang Jiang. Pixel2mesh: Generating 3d mesh models from single rgb images. In *ECCV*, 2018. 3
- [56] Yufu Wang, Ziyun Wang, Lingjie Liu, and Kostas Daniilidis. Tram: Global trajectory and motion of 3d humans from in-the-wild videos. *arXiv preprint arXiv:2403.17346*, 2024. 4
- [57] Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. Prolificdreamer: High-fidelity and diverse text-to-3d generation with variational score distillation. In *NeurIPS*, 2023. 3
- [58] Shih-En Wei, Varun Ramakrishna, Takeo Kanade, and Yaser Sheikh. Convolutional pose machines. In *CVPR*, 2016. 4, 6, 8
- [59] Qianyang Wu, Ye Shi, Xiaoshui Huang, Jingyi Yu, Lan Xu, and Jingya Wang. Thor: Text to human-object interaction diffusion via relation intervention. *arXiv preprint arXiv:2403.11208*, 2024. 2, 3
- [60] Shangzhe Wu, Christian Rupprecht, and Andrea Vedaldi. Unsupervised learning of probably symmetric deformable 3d objects from images in the wild. In *CVPR*, 2020. 3
- [61] Shangzhe Wu, Tomas Jakab, Christian Rupprecht, and Andrea Vedaldi. DOVE: Learning deformable 3d objects by watching videos. *IJCV*, 2023.
- [62] Shangzhe Wu, Ruining Li, Tomas Jakab, Christian Rupprecht, and Andrea Vedaldi. Magicpony: Learning articulated 3d animals in the wild. In *CVPR*, 2023. 3

- [63] Xiang Xu, Hanbyul Joo, Greg Mori, and Manolis Savva. D3d-hoi: Dynamic 3d human-object interactions from videos. *arXiv preprint arXiv:2108.08420*, 2021. 2
- [64] Vickie Ye, Georgios Pavlakos, Jitendra Malik, and Angjoo Kanazawa. Decoupling human and camera motion from videos in the wild. In *CVPR*, 2023. 4
- [65] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *CVPR*, 2021. 3
- [66] Qihang Zhang, Chaoyang Wang, Aliaksandr Siarohin, Peiye Zhuang, Yinghao Xu, Ceyuan Yang, Dahua Lin, Bolei Zhou, Sergey Tulyakov, and Hsin-Ying Lee. Scenewiz3d: Towards text-guided 3d scene composition. *arXiv preprint arXiv:2312.08885*, 2023. 3
- [67] Xiaohan Zhang, Bharat Lal Bhatnagar, Sebastian Starke, Vladimir Guzov, and Gerard Pons-Moll. Couch: Towards controllable human-chair interactions. In *ECCV*, 2022. 2
- [68] Yan Zhang, Mohamed Hassan, Heiko Neumann, Michael J Black, and Siyu Tang. Generating 3d people in scenes without people. In *CVPR*, 2020. 3
- [69] Kaifeng Zhao, Shaofei Wang, Yan Zhang, Thabo Beeler, and Siyu Tang. Compositional human-scene interaction synthesis with semantic control. In *ECCV*, 2022. 3
- [70] Chuanxia Zheng and Andrea Vedaldi. Free3d: Consistent novel view synthesis without 3d representation. In *CVPR*, 2024. 3, 5
- [71] Zerong Zheng, Tao Yu, Yebin Liu, and Qionghai Dai. Pamir: Parametric model-conditioned implicit representation for image-based human reconstruction. *PAMI*, 2021. 5, 6
- [72] Xiaoyu Zhou, Xingjian Ran, Yajiao Xiong, Jinlin He, Zhiwei Lin, Yongtao Wang, Deqing Sun, and Ming-Hsuan Yang. Gala3d: Towards text-to-3d complex scene generation via layout-guided generative gaussian splatting. In *ICML*, 2024. 3
- [73] Silvia Zuffi, Angjoo Kanazawa, David W. Jacobs, and Michael J. Black. 3D menagerie: Modeling the 3D shape and pose of animals. In *CVPR*, 2017. 3, 5
- [74] Silvia Zuffi, Ylva Mellbin, Ci Li, Markus Hoeschle, Hedvig Kjellström, Senya Polikovsky, Elin Hernlund, and Michael J Black. Varen: Very accurate and realistic equine network. In *CVPR*, 2024. 3

# DreamHOI: Subject-Driven Generation of 3D Human-Object Interactions with Diffusion Priors

## Supplementary Material

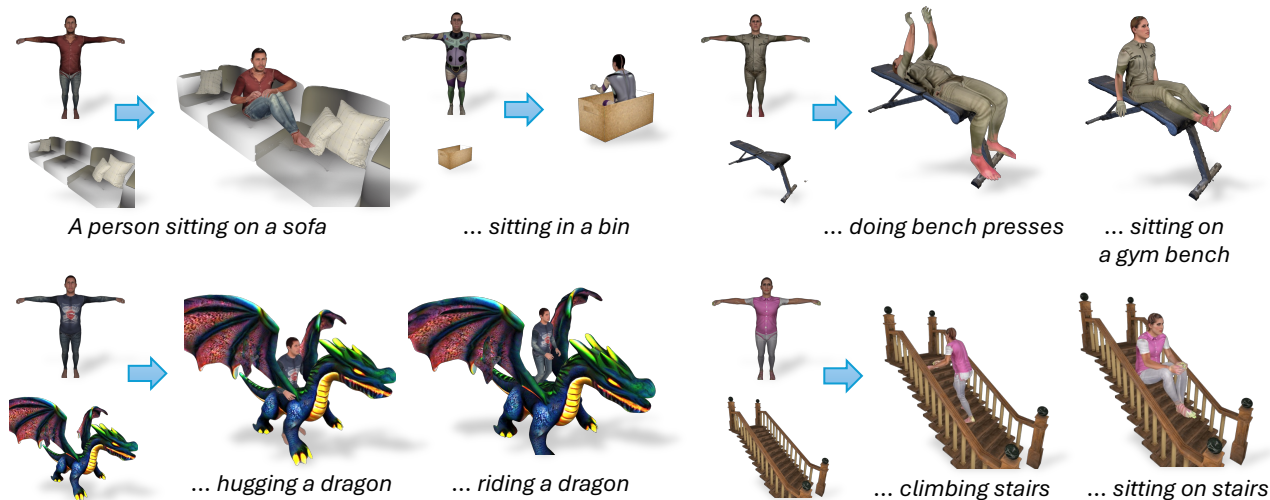


Figure 7. **Additional Results.** DreamHOI is able to generate HOIs for diverse objects and corresponding prompts.

### A. Additional Results

**Additional qualitative results.** Additional results on a variety of prompts and objects can be found in Fig. 7. DreamHOI is capable of realistically deforming the human pose to interact with these objects faithfully to the corresponding textual prompts.

**Failure cases.** We show some cases where DreamHOI failed in Fig. 8. From a manual inspection, in most cases this was due to the underlying diffusion model not understanding the semantic composition, because it was too complex, vague, or exotic (respectively, in Fig. 8). In other cases this was due to the pose prediction (SMPLify-X [35]) not working properly. Therefore we believe an improvement to either would improve DreamHOI’s ability to generate realistic HOIs.

**Additional comparisons.** In our baseline comparisons (Sec. 4.4), we considered comparing to the case where we generate a NeRF by DreamFusion using DeepFloyd IF guidance. We showed that, although mostly related to the prompt, the outputs often had problems like not view-consistent or being too large. We now compare against the same baseline but with IF replaced with MVDream guidance, in Fig. 9. Note that MVDream guidance is able to produce very detailed and view-consistent NeRFs, but fails to

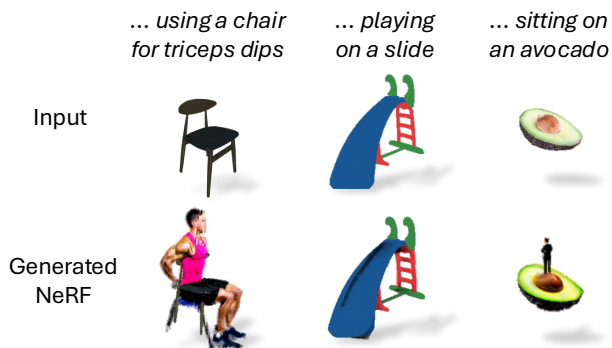


Figure 8. **Failure Cases.** We visualize failed outputs from the first-round NeRF optimization (*i.e.*,  $f_{\theta_0}$ ). Our pipeline is unlikely to recover if the NeRF from the initial round fails to capture the approximate spatial relationship between the human and the object.

understand the basic compositional relations in the prompt (*e.g.*, “sit”). This provides motivation to use DeepFloyd IF as our base model for better textual understanding.

### B. Discussions

**Failure analysis of direct optimization.** The most straightforward way to solve the HOI generation task we



Figure 9. **Additional Comparison.** Following Fig. 4, we additionally compare DreamHOI to a DreamFusion baseline where we use only MVDream [43] for the diffusion guidance without DeepFloyd IF [19].



Figure 10. **Visualization** of SDS gradient. Left: rendered  $x_{\text{HOI}}$  ( $M_\xi$  with  $M_{\text{Obj}}$ ); Middle: gradient  $\nabla_x \mathcal{L}_{\text{SDS-HOI}}$ ; Right: norm of the gradient.

propose is to only use explicit SMPL (or any other body model) pose parameters  $\xi$  instead of an implicit NeRF  $f_\theta$  as the object for optimization. In this solution, we render the resulting SMPL mesh  $M_\xi$  with an object mesh  $M_{\text{Obj}}$ , and use SDS to directly optimize  $\xi$ . This avoids the problem of translation between explicit and implicit forms and makes optimization much faster (*e.g.* SMPL only has 69 pose parameters [28]).

However, in our extensive tests, this method does not work at all (even if we initialize from a near-optimal pose,  $\xi$  regresses to a nonsensical pose as in Fig. 4; additional changes like making vertex colors and global position learnable do not help either). This was observed in Sec. 4.4.

To illustrate the reason, we monitor the SDS loss and gradient during the optimization of  $\xi$  in Fig. 10, for the prompt “a person sitting on a chair”. In the middle and right panels, the guidance tries to add legs to the position where legs would be expected, on the seat and to its front. However, there is no way for this gradient to add legs to be propagated to  $\xi$ , because  $\xi$  can only receive gradients on pixels that the rendered  $M_\xi$  occupies. In other words, the tendency of diffusion models to add and delete limbs globally, instead of gradually moving a limb, means that SDS is not suitable for optimizing  $\xi$  directly. This necessitates the dual implicit-explicit optimization we propose.

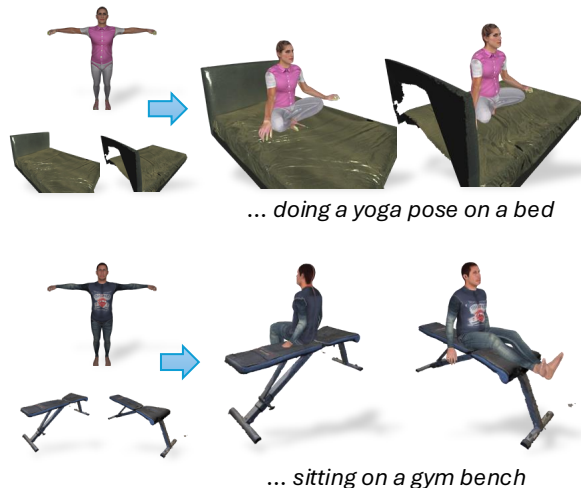


Figure 11. **Front Direction Control.** MVDream [43] guidance enables us to give an implicit “front” direction: the generated humans consistently face the  $\pm x$  direction regardless of the orientation of the object.

**MVDream front direction.** We claimed in Sec. 4.5 that MVDream takes camera positions as input and based on bias in its training data, MVDream gives a prior of a “front” direction (in  $+x$  direction) for generating the HOI. We demonstrate this in Fig. 11. This gives us the ability to control the forward face of the entire HOI (typically the direction the person is facing, or its opposite) with respect to the object by rotating the mesh  $M_{\text{Obj}}$  in the generation.

## C. Technical Details

**Optimization.** We follow MVDream [43] and optimize the initial NeRF in 2 stages. In each stage, we optimize the NeRF with AdamW optimizer (learning rate and weight decay are both set to 0.01) for 5000 steps. We render  $64 \times 64$  and  $256 \times 256$  images and use batch size 8 and 4 respectively in two stages. After NeRF re-initialization, MVDream guidance is no longer used, and we increase the rendering resolution to  $512 \times 512$ , reduce the batch size to 1, and decrease the learning rate to 0.001.

The field of view  $f$  of the camera in each optimization step is sampled uniformly at random within  $[15^\circ, 60^\circ]$ , and the camera distance to origin is set to  $D / \tan(f/2)$  where the denominator is such that a unit volume in the 3D space corresponds roughly to a fixed area in the 2D space, as in Sec. 3.5 for  $\mathcal{R}_{\text{SA}}$  to work properly, and  $D \sim [0.8, 1.0]$  is a perturbation. The elevation angle is sampled uniformly from  $[0^\circ, 30^\circ]$ . Although not done in this work, we recommend lowering it to  $[-30^\circ, 30^\circ]$  if parts of the human may be below the object for better supervision. For rendering views  $x_i$  of the NeRF for pose estimation (Sec. 4.2), we use

an array of cameras with distance 3 to the origin, elevated at  $40^\circ$ .

The NeRF representing the human is constrained in a ball of radius 1. We initialize it to be at the origin. The number of parameters of the NeRF MLP ( $f_\theta$ ) is about 12.6 million.

The background color is learned, with a lower learning rate 0.001, and is replaced during training with a random color with probability 0.5 (increased to 1 after re-initialization) in training for augmentation.

The NeRF renderer uses one ray per 2D pixel and 512 samples per ray.

**Guidance.** For SDS, we use classifier-free guidance [14] with guidance weight set to  $\omega = 50$ . We include the *negative* prompt “missing limbs, missing legs, missing arms” during optimization.