

# HOIAnimator: Generating Text-prompt Human-object Animations using Novel Perceptive Diffusion Models

Anonymous CVPR submission

Paper ID 1669

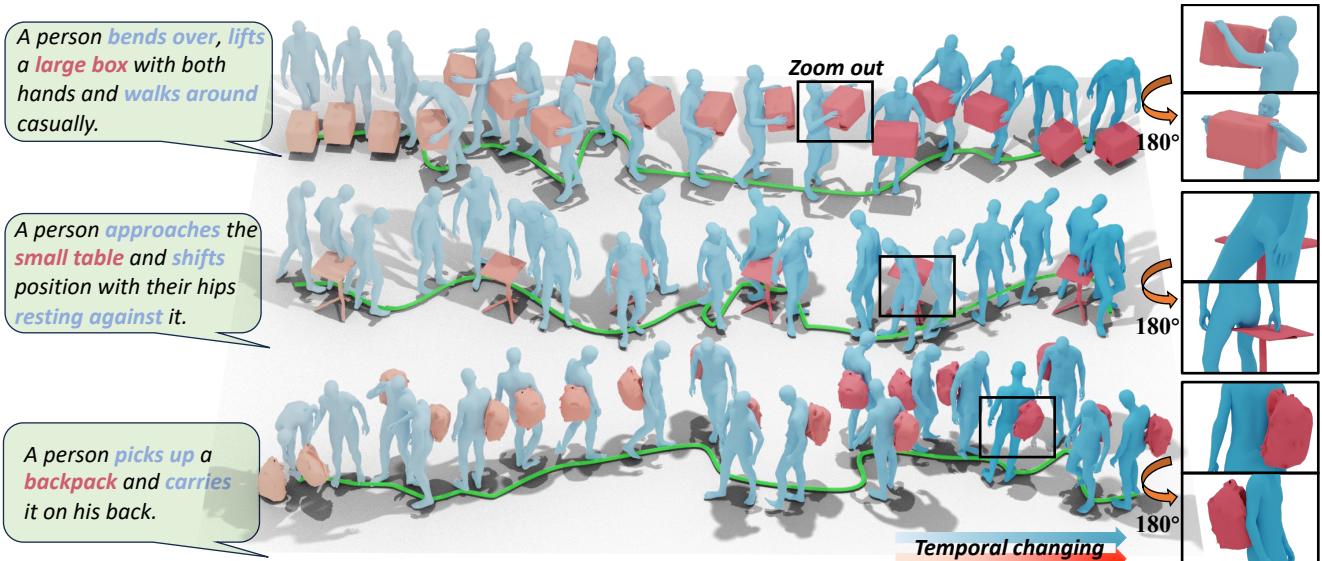


Figure 1. Our HOIAnimator excels in turning text descriptions into realistic animations of human-object interactions. It's adept at depicting a variety of actions, such as bending, lifting boxes, and picking up bags, with believable contact between the human and the objects.

## Abstract

To date, the quest to rapidly and effectively produce human-object interaction (HOI) animations directly from textual descriptions stands at the forefront of computer vision research. The underlying challenge demands both a discriminating interpretation of language and a comprehensive physics-centric model supporting real-world dynamics. To ameliorate, this paper advocates HOIAnimator, a novel and interactive diffusion model with perception ability and also ingeniously crafted to revolutionize the animation of complex interactions from linguistic narratives. The effectiveness of our model is anchored in two ground-breaking innovations: (1) Our Perceptive Diffusion Models (PDM) brings together two types of models: one focused on human movements and the other on objects. This combination allows for animations where humans and objects move in concert with each other, making the overall motion more realistic. Additionally, we propose a Perceptive Message Passing (PMP) mechanism to enhance the communication bridging the two models, ensuring that the animations are smooth and unified; (2) We devise an Interaction Contact

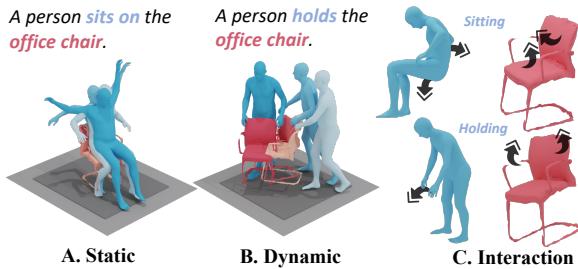
Field (ICF), a sophisticated model that implicitly captures the essence of HOIs. Beyond mere predictive contact points, the ICF assesses the proximity of human and object to their respective environment, informed by a probabilistic distribution of interactions learned throughout the denoising phase. Our comprehensive evaluation showcases HOIAnimator's superior ability to produce dynamic, context-aware animations that surpass existing benchmarks in text-driven animation synthesis. We will open the source codes upon the paper's acceptance.

## 1. Introduction and Motivation

In the dynamic landscape of AI-guided creation (AIGC), 3D animation has emerged as a crucial and challenging domain. This challenge is epitomized in the realm of human-object interactions (HOIs), where the goal is to generate realistic animations from textual descriptions. The intricacy lies in translating written language into visual narratives that accurately capture the nuanced dynamics between humans and objects. Achieving this requires a sophisticated under-

021  
022  
023  
024  
025  
026  
027  
028  
029  
030

031  
032  
033  
034  
035  
036  
037  
038  
039



**Figure 2. Navigating the complexity of HOIAnimator.** (A): the ‘static’ interaction is depicted with a stationary office chair, showcasing a human sitting on the chair. (B): the ‘dynamic’ interaction portrays both the human and the object in motion, exemplified by the act of holding an object. (C): Arrows denote forces and trajectories involved in HOI.

standing of linguistic cues and a deep knowledge of physical interaction principles. The intersection of language and physics makes 3D animation creation both challenging and innovative.

The field of animation creation, driven by advancements in natural language processing and generative modeling as demonstrated by Li et al.[21] and Zhang et al.[52], faces a formidable challenge: **effectively mapping the low-dimensional latent spaces of textual descriptions to the high-dimensional, complex spaces of human and object motions**. While recent endeavors like those by previous effort [7, 26, 36] showcase the potential of converting text prompts into visual content, the intricate task of accurately rendering realistic dynamics from concise text remains largely unmastered. This issue is particularly evident when the model interprets a simple text-based action like ‘holding’ a bag but cannot adequately replicate the diverse, context-dependent ways this action might manifest, such as the various methods of carrying a bag.

Building upon the foundation laid by innovative projects like InterGen [22], Scene Diffuser [16], Narrator [46], and InterDiff [45], the field has advanced significantly in generating animations through interactions with scenarios or human figures. Yet, a formidable challenge persists: **accurately modeling and animating the complex forces and reciprocal influences between humans and objects**. As illustrated in Fig. 2, the relationship between a human and an office chair, for example, is not merely a static or one-dimensional interaction. It is a dynamic interplay, influenced by various factors such as the intent behind the interaction, the trajectory of movement, and the contextual use of the object, all of which are dictated by the narrative text.

Our HOIAnimator introduces dual perceptive diffusion models (PDM) to simplify the first challenge, adeptly capturing spatial dynamics between humans and objects, creating animations consistent with the textual prompts. HOIAnimator utilizes two diffusion models: a human-centric model for human movements and an object-centric model for object dynamics. The models communicate with each

other through a novel Perceptive Message Passing (PMP) mechanism. The PMP adaptively learns the weight and bias of object clues embedded into human motion flow. This collaborative approach ensures the active engagement of both entities in the animation, leading to more complete and accurate representations of the narrative.

To tackle the challenge of accurately representing complex forces in HOI, we introduce a novel concept of Interaction Contact Field (ICF), a model that learns the patterns of contact as described in text prompts through a diffusion model. This diffusion model is skilled at interpreting textual cues and translating them into ICF, effectively capturing the spatial dynamics of HOIs. The ICF offers a comprehensive perspective on interaction probabilities, advancing beyond traditional collision detection methods. By incorporating a probabilistic approach that considers object affordance, human intent, and ergonomics, the ICF is able to predict potential points of interaction. This leads to animations that are both dynamic and adaptive, more accurately mirroring the complexities of real-world object manipulation.

To summarize, our salient contributions are listed as follows.

- We propose a **brand new HOIAnimator**, a framework that utilizes dual Perceptive Diffusion Models, human-centric model and object-centric model, to accurately render human-object interactions in animations. The core of PMD is a cutting-edge PMP mechanism designed to enable seamless and effective communication between human and object-centric models, ensuring lifelike and engaging animations.
- We present a **simple yet powerful ICF** to proactively identify and assess potential contact points between entities. The key is learning the distribution of interaction probability between human and object, and mapping text-based interaction cues into spatial dynamics. Our method goes beyond basic collision detection, using a probabilistic field informed by object characteristics, human intentions, and ergonomics. This leads to animations that are dynamically responsive and closely mimic real-world object interactions.
- We conduct extensive experiments in both public and wild datasets. The results demonstrate the superior ability of our HOIAnimator to generate human-object interaction animations. Our dataset and project will be public.

## 2. Related Work

### 2.1. Text-to-animation Generation

Animation generation commonly employed neural network models such as the Variational AutoEncoder (VAE)[12, 28, 42], Vector Quantized-Variational AutoEncoder (VQ-VAE)[13, 29, 53], and Generative Adversarial Networks (GAN) [24, 43, 47, 48] to acquire representations and pat-

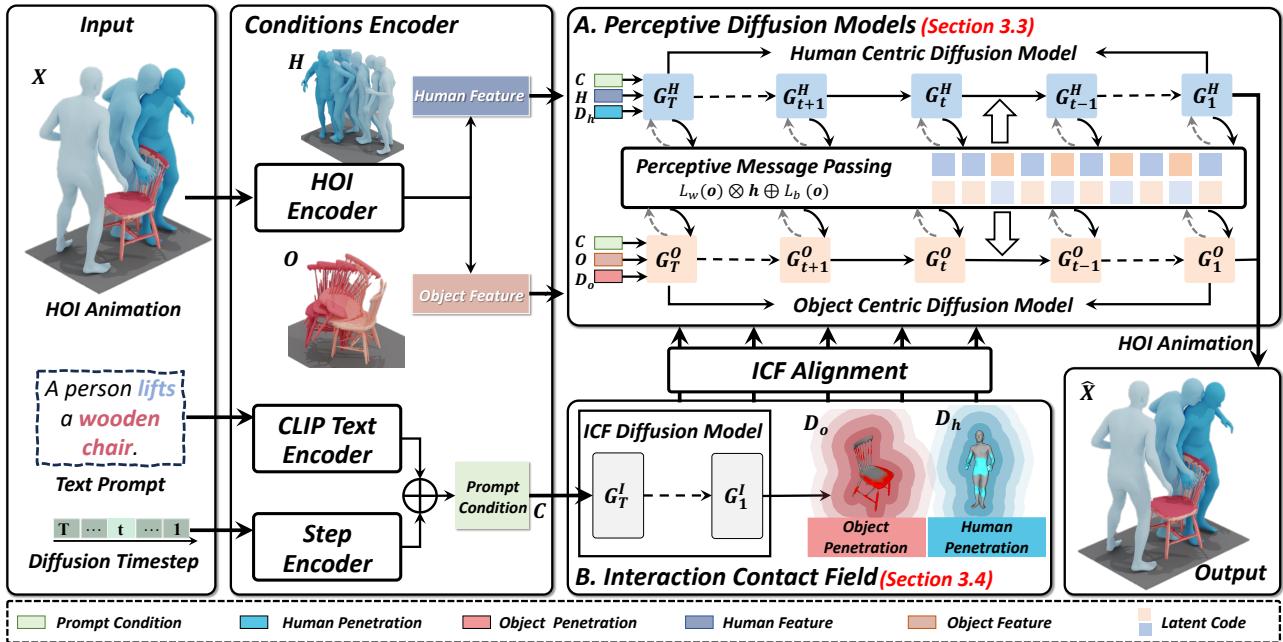


Figure 3. **Method overview**. We propose the HOIAnimator with two key parts: (1) Perceptive Diffusion Models (PDM). This part combines the movements of both people and objects in the animation, making sure they move together in a realistic way. (2) Interaction Contact Field (ICF). The ICF provides the clues that humans and objects interact and contact each other (Training phase of HOIAnimator).

terns necessary for generating animations from text descriptions. These models were trained on textual descriptions and generated representations and patterns for animations. However, the recent introduction of diffusion models [15, 23, 33, 51] greatly improved the ability to reason about text and represent animation [1, 3, 6, 41, 52]. For instance, MDM [37] introduced a transformer-based generative model that was adapted for the many-to-many nature of the domain. Building on MDM, SinMDM [30] learned the internal motifs of a single motion sequence with arbitrary topology and synthesized motions of arbitrary length that were faithful to them. Furthermore, PriorMDM [32] proposed using a pre-trained diffusion-based model as a generative prior to fine-tuning for few-shot and zero-shot settings. Inspired by the two-person generation, we propose a bidirectional diffusion model to generate HOI animations. This model utilizes the diffusion process for both inference and generation, enabling it to effectively handle the intricate relationships and uncertainties associated with humans and objects.

## 2.2. Human-object Dynamic Interaction

Current research prominently focuses on unraveling the intricacies of human-object interactions. Recent studies [10, 14, 39, 49] explored the detailed modeling of whole-body interactions. However, most works [2, 5, 25, 35, 44, 50] focused on human-object relations within static environments, where objects are treated as passive. In contrast, recent works [16, 20, 45] integrated objects and scenes as dynamic components in motion prediction models. In

multi-human interactions, some works [22, 32] introduced diffusion-based approaches for generating text-driven interaction motions for two people. The complex nature of interactions between humans and objects, significantly different from multi-human interactions, is further complicated by the disparity in data features. Addressing this, several works [8, 34, 40, 54] attempted to unravel these complexities. Utilizing the expanding array of 3D datasets that capture human interaction [4, 9, 17, 18, 25, 25, 46], our work introduces a novel text-prompt paradigm for streamlining the generation of HOI animations.

## 3. New Methodology

### 3.1. Overview

Our HOIAnimator aims to achieve end-to-end conversion from text description to 3D HOI animations. As illustrated in Fig. 3, the pipeline of HOIAnimator begins with a novel representation for HOI animation, described in Section 3.2, which underpins our text-prompt-based animation generator, aimed at minimizing inconsistencies between humans and objects. Subsequently, to synchronize text-driven prompts with corresponding dynamic visual representations, we introduce the PDM, a novel interactive diffusion mechanism specifically described in Section 3.3. Finally, to optimize the details between the surfaces of humans and objects, we introduce the ICF, engineered to evaluate the probability of contact within the interaction spaces afforded by objects in Section 3.4.

159  
160  
161  
162  
163  
164  
165  
166  
167  
168  
169

170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185

186 

### 3.2. HOIAnimator Definition and Preliminaries

187 Our HOIAnimator is specifically designed to handle the animation with the positions of 3D coordinates for both humans and objects. We define the **HOI Animation Definition**  
 188 in rigid mathematics. Meanwhile, we introduce **Diffusion Model for HOI Animation Generation** for generating HOI animation. We utilize diffusion models' generative power, employing a stochastic diffusion process for  
 189 dynamic, precise HOI animations.

190 **HOI Animation Definition.** Generating the spatial co-  
 191 ordinates of humans and objects and providing pose information  
 192 for both are essential for creating consistent animations. Inconsistencies often occur when various methods are  
 193 used to represent these elements. For instance, human bodies are commonly represented using SMPL-H [27], while  
 194 objects are typically depicted through translation and rotation. To resolve this, we propose a unified approach involving four key parameters: the human shape parameter  
 195 ( $\beta \in \mathbb{R}^{10}$ ), the human pose parameter ( $\theta \in \mathbb{R}^{159}$ ), the object's translation parameter ( $\tau \in \mathbb{R}^3$ ), and the object's rotation parameter ( $\gamma \in \mathbb{R}^3$ ). By integrating these parameters,  
 196 we form the HOI animation, denoted as  $x_{1:i} = \{\beta, \theta, \tau, \gamma\}$ , where  $x_i \in \mathbb{R}^{175}$  represents the pose state in frame  $i$ .  
 197  $i \in [0, N]$  and  $N$  is the maximum animation length. The HOI animation effectively captures the dynamics of human-object interactions in animation.

198 **Diffusion Model for HOI Generation.** Drawing inspiration  
 199 from previous works [11, 19, 31], we opt for the diffusion  
 200 model [16] to generate HOI animations. Similarly  
 201 to the text-driven motion generation task, our training set  
 202 for text-driven HOI animations consists of pairs  $(x_i, text_i)$ ,  
 203 where  $text_i$  is the textual description of the HOI animation  
 204 ( $x_i$ ). During inference, given a textual description of the  
 205 animation, we can generate an animation that matches the  
 206 description. We build our text-driven HOI animation pipeline  
 207 based on diffusion models. This diffusion can be modeled  
 208 as a Markov noising process ( $\{x_{1:i}^t\}_{t=0}^T$ ), gradually adding  
 209 Gaussian noise to the ground truth ( $x_{1:i}^0$ ) until it eventually  
 210 becomes pure Gaussian noise ( $x_{1:i}^T$ ):

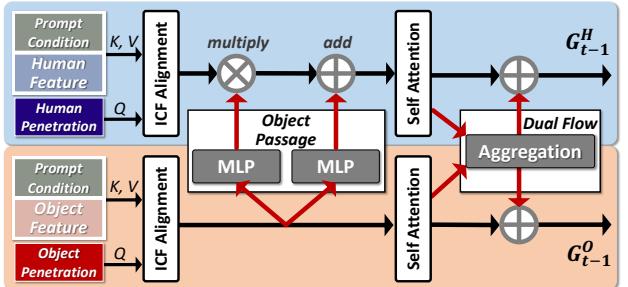
$$q(x_{1:i}^t | x_{1:i}^{t-1}) = \mathcal{N}(\sqrt{1 - \alpha_t} x_{1:i}^{t-1}, \alpha_t \mathbf{I}), \quad (1)$$

212 where  $t$  denotes diffusion step,  $t \in [1, T]$ ,  $\alpha_t \in [0, 1]$  are  
 213 fixed set of values, generated by formula [15]. Thus, at a  
 214 sufficiently large step  $T$ ,  $\alpha_t$  approaches 1, at which point  
 215 it can be approximated as a Gaussian distribution  $x_{1:i}^T \sim$   
 216  $\mathcal{N}(0, \mathbf{I})$ .

217 

### 3.3. Perceptive Diffusion Models

218 Specific correlations between humans and objects are es-  
 219 sential to address the challenge of significant differences in  
 220 pose parameters between humans and objects. Extended  
 221 from the single diffusion model [37], our approach employs



222 **Figure 4. Perceptive Message Passing.** Between object and human centric diffusion models, we use object passage and dual flow to adjust the features of humans and objects dynamically.

223 the PDM for more nuanced processing as depicted in Fig. 3-  
 224 A. PDM consists of two specialized components: (1) The  
 225 human centric diffusion model and object centric diffusion  
 226 model. (2) PMP exchanges clues for the two separate diffu-  
 227 sion models as depicted in Fig. 4.

228 **Human and Object Centric Diffusion Models.** In  
 229 PDM, the human centric diffusion model adapts to re-  
 230 fine human motion, emphasizing human movement dynamics.  
 231 In contrast, the object centric diffusion model fo-  
 232 cuses on optimizing object trajectories, ensuring precision  
 233 in object motion. We duplicate the original animation se-  
 234 quence ( $x_{1:i} = \{\beta, \theta, \tau, \gamma\}$ ), resulting in two identical  
 235 copies. Each copy is then specialized: one for the object se-  
 236 quences ( $x_{obj} = \{\tau, \gamma\}$ ) and the other for human sequences  
 237 ( $x_{hum} = \{\beta, \theta\}$ ). We distinctively handle the feature of the  
 238 human sequence (**H**) and the feature of the object sequence (**O**). This separation allows for tailored processing of each  
 239 sequence type. Further, as elaborated in Equation 1, we  
 240 develop two diffusion models: the object centric diffusion,  
 241 denoted as  $G^O$ , and the human centric diffusion, denoted as  
 242  $G^H$ . We get the final output ( $\hat{x}_{obj}, \hat{x}_{hum}$ ) as:

$$\begin{aligned} \hat{x}_{obj} &= G^O(x_{obj}, E_{text}(text) + E_{step}(t)), \\ \hat{x}_{hum} &= G^H(x_{hum}, E_{text}(text) + E_{step}(t)), \end{aligned} \quad (2)$$

243 where  $E_{step}$  is diffusion step encoder.  $E_{text}$  is text encoder.  
 244 The  $G^O, G^H$  predict the final clean animation in each sam-  
 245 pling step. We further break down this objective into dis-  
 246 tinct components: rotation and translation losses, applicable  
 247 to both humans and objects. We then focus on optimizing  
 248 the diffusion model for humans and objects separately, ad-  
 249 dressing each aspect as follows:

$$\begin{aligned} \mathcal{L} &= \mathcal{L}_{human} + \mathcal{L}_{obj} \\ &= \mathbb{E}_{t \sim [1:T]} [\|x_{hum} - G^H\|_2 + \|x_{obj} - G^O\|_2], \end{aligned} \quad (3)$$

250 where  $\mathbb{E}_{t \sim [1:T]}$  denotes the average loss for all time steps  $T$ .  
 251 Hence, it is possible to systematically eliminate noise from  
 252 both the human and object sequences collectively, yielding  
 253 a coherent animated sequence aligned with the provided text  
 254 condition.

**Perceptive Message Passing.** PMP facilitates efficient information exchange between these two components. This exchange is enhanced by two mechanisms: object passage, which deals with the movement of objects passage, and dual flow, which intertwines the processing of human and object motions. Together, these elements of PDM ensure a comprehensive and accurate representation of motion dynamics, avoiding the oversimplifications common in single diffusion models.

In the PMP, the first step is utilizing the object passage method. This is pivotal for improving the incorporation of object-specific details into human motion. Object passage processes object latent code ( $\mathbf{o}$ ) to modify human latent code ( $\mathbf{h}$ ) dynamically ( $\mathbf{o}$  and  $\mathbf{h}$  have a detailed description in Section 3.4). This dynamic adjustment aligns human-centered data effectively. The adaptive mechanism is instrumental in enhancing the model's proficiency in capturing the complex relationship between human kinetics and object dynamics, a key factor in preserving the realism of the generated HOI animation. The object passage ( $F_{obj}(\mathbf{h}|\varphi, \phi)$ ) can be written as:

$$\begin{aligned}\mathbf{h}' &= F_{obj}(\mathbf{h}|\varphi, \phi) \\ &= \varphi \cdot \mathbf{h} + \phi \\ &= L_w(\mathbf{o}) \cdot \mathbf{h} + L_b(\mathbf{o}),\end{aligned}\quad (4)$$

where  $L_w$  and  $L_b$  are two fully connected networks.  $\varphi$  and  $\phi$  denote dynamically adjusted weights and biases. The input processed by the converter, which is responsible for refining the human position sequence, is influenced by its intrinsic characteristics and dynamic interactions with object movement. This approach ensures that the generated object positional sequences harmonize with contemporaneous human motion, resulting in a faithful portrayal of human-object interactions.

The second step in PMP involves implementing the dual flow approach. This method boosts bidirectional communication between the human and object centric diffusion modules. Before the decoding phase of these modules, we integrate the human latent features ( $\mathbf{h}'$ ) and object latent features ( $\mathbf{o}'$ ) using an aggregation module ( $F_{dual}$ ). These integrated features are then added to the original features as residuals, enhancing the overall process. This integration can be articulated as:

$$\langle \hat{\mathbf{h}}, \hat{\mathbf{o}} \rangle = F_{dual}(\langle \mathbf{o}', \mathbf{h}' \rangle, \langle \mathbf{h}', \mathbf{o}' \rangle) \oplus \langle \mathbf{o}', \mathbf{h}' \rangle, \quad (5)$$

where  $\langle , \rangle$  operations follow the specified order, with the first element interacting with prior elements and the latter with the following ones.  $\oplus$  is the element-wise add, which is to learn the residual value.  $F_{dual}(a, b)$  denotes the concatenation of the two feature vectors  $a$  and  $b$ . First, we perform self-attention on  $a, b$  to obtain features. After transformation by the aggregation module, the concatenated features are truncated to match the dimensional of  $a$ . Last,  $\hat{\mathbf{h}}$

and  $\hat{\mathbf{o}}$  serve as feature inputs to the latent decoder, which then reconstructs them into  $\hat{x}_{1:i}^0$ .

### 3.4. Interaction Contact Field

Our primary objective is to create realistic and interactive HOI animations. Previous models [22, 32, 45] often struggle to capture this information about the contact between humans and objects. To address this challenge, we introduce a revolutionary method: ICF. The ICF is meticulously designed to calculate the probability of contact between human bodies and specific object regions crucial for interaction, as depicted in Fig. 3-B. By focusing on these contact probabilities, the ICF facilitates the generation of realistic HOI animations. This notably improves the realism of interactions. Further strengthening this innovation is a sophisticated ICF embedding scheme tailored for both granular (ICF) and comprehensive (HOI animation) latent spaces. This scheme ensures exceptional precision in capturing and visualizing the intricacies of HOIs.

**ICF Prediction.** Simply calculating the SDF for humans and objects yields only basic positional data. However, during interactions between humans and objects, more intricate details such as contact and penetration are crucial. To address these complex interactions, we propose calculating the ICF via the contact area. Specifically, we assess the contact and penetration states within the length ( $S$ ) of HOI animation. For humans, we represent the vertices as  $v_h \in \mathbb{R}^{S*V_h*3}$ , where  $V_h$  is the count of human vertices. Similarly, for objects, the vertices are denoted as  $v_o \in \mathbb{R}^{S*V_o*3}$  representing the object vertex count. Then, we randomly sample  $N$  points from the object vertices ( $V_o$ ), and for each point, we calculate the nearest distance to the human contact information ( $C_h$ ), assigning a symbol to represent object penetration ( $D_o$ ). The calculation of human penetration ( $D_h$ ), follows a similar approach. Therefore,  $D_{<h,o>}$  can be defined as follows:

$$\begin{aligned}D_{<h,o>} &= F_{ICF}\left(\underbrace{C_{<h,o>}, \text{sample}(v_{<o,h>}, N)}_{\uparrow}\right), \\ C_h[j] &= \| v_h[j] - v_o[i] \|_2, j = 1, \dots, V_h \\ C_o[i] &= \| v_o[i] - v_h[j] \|_2, j = 1, \dots, V_o\end{aligned}\quad (6)$$

where  $v_h[i] \in \mathbb{R}^3, v_h[j] \in \mathbb{R}^3$  are  $j$ -th and  $k$ -th vertex on humans and objects, respectively.  $F_{ICF}()$  is the function that directly computes the signed distance field.  $\text{sample}()$  is a sequence of point clouds.  $D_h, D_o \in \mathbb{R}^N$  provide us with information on the spatial relationship between the object mesh and the human mesh. Then we pre-train the ICF diffusion model  $G^I(text)$  similar to Equation 7, and the corresponding interaction contact field can be generated through text description. We predict the ICF process as:

$$\langle \hat{\mathbf{D}}_h, \hat{\mathbf{D}}_o \rangle = G^I(\mathbf{D}_{<h,o>}, E_{text}(text) + E_{step}(t)), \quad (7)$$

366 where *text* remains consistent with the text of HOI animation.  
 367 In this way, we can get the ICF based on textual cues,  
 368 which helps the PDM learn the probability of contact.

369 **ICF Alignment for HOI.** We incorporate the ICF em-  
 370 bedding scheme to effectively align the interaction contact  
 371 field ( $\hat{\mathbf{D}}_h, \hat{\mathbf{D}}_o$ ) with the HOI animations, as shown in Fig. 4.  
 372 Using an object latent code ( $\mathbf{o}$ ) as an example, we merge  
 373 the object feature ( $\mathbf{O}$ ) with the condition ( $\mathbf{c}$ ) to form a com-  
 374 bined feature map  $\mathbf{L} = \{\mathbf{c}, \mathbf{O}\}$ . Following this, we em-  
 375 ploy cross-attention (*Attn*) to calculate the desired atten-  
 376 tion weights, which are crucial for integrating the human  
 377 and object features in the HOI animation. The ICF Embed-  
 378 ding process can be formulated as:

$$\mathbf{o} = \text{Attn}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}, \quad (8)$$

$$\mathbf{Q} = \mathbf{W}^Q \hat{\mathbf{D}}_o, \mathbf{Q} = \mathbf{W}^K \mathbf{L}, \mathbf{V} = \mathbf{W}^V \mathbf{L},$$

380 where  $\mathbf{W}^K, \mathbf{W}^V \in \mathbb{R}^{d_s} \times \mathbb{R}^{d_k}$  and  $\mathbf{W}^Q \in \mathbb{R}^{d_q} \times \mathbb{R}^{d_k}$  are  
 381 trainable weights.  $d_s, d_q$  and  $d_k$  are the channel numbers of  
 382 the corresponding weights. ICF alignment guides the gen-  
 383 eration of HOI animation, ensuring excellent accuracy in  
 384 generating the complexity of HOI animations.

## 385 4. Experiments

386 This section presents our HOIAnimator’s implementa-  
 387 tion details and experimental results, comparing it with previous  
 388 state-of-the-art methods. In addition, it includes an ablation  
 389 study and a user study. More results are provided in the  
 390 supplementary material.

### 391 4.1. Implementation Details

392 **Datasets.** Our dataset is compiled from publicly available  
 393 Behave [4] and InterCap [17] datasets, which consist of motion  
 394 captured from 3D human interactions. These datasets  
 395 offer diverse human interaction actions, featuring common  
 396 objects such as tables, backpacks, and chairs and typical inter-  
 397 active actions such as carrying, sitting on, and playing  
 398 with. However, these datasets are limited because they are  
 399 labeled in a multi-class format and lack detailed textual  
 400 descriptions. To enhance our dataset for this application, we  
 401 undertake three normalization steps. First, we standardize  
 402 the frame rate of all motions to a consistent 30 FPS. For ani-  
 403 mations that are longer than 6 seconds, we randomly trim  
 404 them to a maximum of 10 seconds. Finally, these anima-  
 405 tions are aligned with our HOI animation template, ensur-  
 406 ing uniformity and coherence in the dataset. Following this,  
 407 we describe the actions in complete sentences and annotate  
 408 them using the SpaCy <sup>1</sup>. The final step in our data prepara-  
 409 tion process involves manual post-processing, where we  
 410 meticulously filter out any anomalies in the textual descrip-  
 411 tions, ensuring the dataset’s quality and relevance to the  
 412 HOI animation.

<sup>1</sup><https://spacy.io/models>

413 **Evaluation Metrics.** We follow the performance mea-  
 414 sures [12] for quantitative evaluations, including Frechet In-  
 415ception Distance (FID), R Precision, Diversity, and Multi-  
 416Modal Distance (MM Dist). Additionally, our evaluation  
 417 is broadened to include an examination of the Vertex dis-  
 418tance and Penetration score [20] of the generated objects.  
 419 (1) FID measures the quality of HOI animation generation  
 420 by contrasting features of real and synthetically generated  
 421 HOI animation. (2) R precision quantifies the alignment be-  
 422tween generated HOI animations and their textual descrip-  
 423tions, ranking actual text within the top 1, 2, or 3 positions.  
 424 (3) Diversity evaluates the range and depth of the HOI ani-  
 425mation produced. (4) MM Dist calculates the average Eu-  
 426clidean distance between motion features and correspond-  
 427ing textual descriptions. (5) Vertex distance evaluates gen-  
 428eration quality by comparing distances between vertices in  
 429 real and generated objects. (6) Penetration score assesses  
 430realism based on the human-object interaction proximity in  
 431the animations.

432 **Parameters.** For the HOI Encoder, we utilize a 2-layer  
 433 linear architecture with a latent dimensional of 1024. Re-  
 434 garding the ICF, as well as object and human centric dif-  
 435fusion modes, we use a 4-layer transformer with a latent  
 436 dimension of 512. For the variance settings in 3 diffusion  
 437 models, we preset the variance value to increase linearly  
 438 from 0.0001 to 0.02 within  $T = 1000$  noise steps. The  
 439 text encoder incorporates a frozen CLIP ViT-B/32 model  
 440 complemented by two additional transformer encoder lay-  
 441 ers. The Adam optimization algorithm is employed to train  
 442 the model, with a learning rate set at 0.0002. Training is  
 443 executed on 4 NVIDIA 3090Ti GPUs, with a batch size of  
 444 64 per GPU. The model is trained over 250,000 steps.

### 445 4.2. Quantitative Evaluation

446 **Baselines.** In this study, we propose a novel approach to  
 447 the generation of HOI animations with prompt text and  
 448 compare it with several state-of-the-art methods, including  
 449 MoitonCLIP [36], T2M [12], MDM [16], PriorMDM [32],  
 450 MLD [7], InterGen [22].

451 To the best of our knowledge, rare existing work has  
 452 explored text-driven 3D HOI animation generation. To  
 453 thoroughly evaluate the effectiveness of our HOIAnimator,  
 454 we conduct a comprehensive comparison with the above-  
 455 mentioned state-of-the-art. Our method takes textual de-  
 456 scriptions as input and produces HOI animations. For Pri-  
 457 orMDM and InterGen, we retained the core structure of  
 458 communication diffusion but tailored the parts involving  
 459 interactions between two humans to match the format of  
 460 our dataset. For MLD, we employed a VAE to encode our  
 461 dataset into a latent code representation. Subsequently, we  
 462 utilized a diffusion model to generate the latent code, which  
 463 was then decoded to produce the final HOI animations.

464 **Quantitative Results and Analysis.** Tab. 1 shows our

Methods	R Precision↑			Vertex Distance↓	FID↓	MM Dist↓	Diversity→	Penetration↑
	Top 1	Top 2	Top 3					
Real motions	0.508 $\pm$ 0.004	0.725 $\pm$ 0.005	0.821 $\pm$ 0.006	—	0.012 $\pm$ 0.002	6.754 $\pm$ 0.005	9.534 $\pm$ 0.065	—
MoitonCLIP [36]	0.322 $\pm$ 0.006	0.493 $\pm$ 0.005	0.614 $\pm$ 0.005	0.979 $\pm$ 0.110	1.389 $\pm$ 0.049	10.424 $\pm$ 0.009	8.192 $\pm$ 0.075	0.529 $\pm$ 0.003
T2M [12]	0.384 $\pm$ 0.005	0.582 $\pm$ 0.006	0.673 $\pm$ 0.005	0.813 $\pm$ 0.003	0.944 $\pm$ 0.042	8.492 $\pm$ 0.011	8.724 $\pm$ 0.132	0.561 $\pm$ 0.006
MDM [16]	0.363 $\pm$ 0.007	0.573 $\pm$ 0.006	0.692 $\pm$ 0.006	0.783 $\pm$ 0.021	0.859 $\pm$ 0.080	9.382 $\pm$ 0.017	<b>9.537</b> $\pm$ 0.043	0.568 $\pm$ 0.003
MLD [13]	0.448 $\pm$ 0.007	0.628 $\pm$ 0.006	0.701 $\pm$ 0.006	0.711 $\pm$ 0.005	0.859 $\pm$ 0.080	8.382 $\pm$ 0.017	8.543 $\pm$ 0.132	0.578 $\pm$ 0.003
PriorMDM [32]	0.461 $\pm$ 0.006	0.636 $\pm$ 0.005	0.727 $\pm$ 0.035	0.683 $\pm$ 0.073	0.853 $\pm$ 0.028	8.776 $\pm$ 0.012	9.213 $\pm$ 0.042	0.601 $\pm$ 0.002
InterGen [7]	0.491 $\pm$ 0.005	0.652 $\pm$ 0.005	0.734 $\pm$ 0.005	0.523 $\pm$ 0.005	0.717 $\pm$ 0.055	7.932 $\pm$ 0.021	9.344 $\pm$ 0.023	0.613 $\pm$ 0.001
Ours	<b>0.526</b> $\pm$ 0.006	<b>0.719</b> $\pm$ 0.006	<b>0.781</b> $\pm$ 0.005	<b>0.118</b> $\pm$ 0.063	<b>0.623</b> $\pm$ 0.063	<b>7.521</b> $\pm$ 0.014	9.526 $\pm$ 0.029	<b>0.643</b> $\pm$ 0.001

Table 1. **Quantitative evaluation on BEHAVE [4].** To ensure a fair comparison, we conducted 20 experiments.  $x^{\pm y}$  denotes that  $x$  represents the average value of the metric, while  $y$  corresponds to the confidence interval 95% around this mean. ‘↑’ (‘↓’, ‘→’) indicates that the values are better if the metric is larger (smaller, closer); The **bold fonts** denote best performers. The results show that the HOI animations synthesized by our model outperform other baselines in terms of semantic matching.

quantitative comparison results with 6 baselines on BEHAVE. Our HOIAnimator marks a notable advancement over InterGen, as evidenced by measurable improvements across several key metrics. Firstly, it demonstrates enhanced precision (Top-3), boosting the score from 0.734 to 0.781. Furthermore, there is a significant enhancement in the vertices distance metric, with a reduction in the score from 0.523 to 0.118, reflecting a more accurate representation. In addition, the HOIAnimator has achieved greater fidelity in generated animations, evidenced by a decrease in the FID score from 0.717 to 0.623 and a 3% improvement in the penetration score. These advancements collectively signify a substantial improvement in the performance and quality of our HOIAnimator.

Methods	Precision↑	FID ↓	Penetration↑
Real motions	0.821 $\pm$ 0.005	0.012 $\pm$ 0.002	—
w/o ICF	0.756 $\pm$ 0.004	0.714 $\pm$ 0.027	0.615 $\pm$ 0.003
w/o PMP	0.723 $\pm$ 0.006	0.789 $\pm$ 0.042	0.593 $\pm$ 0.002
w/o PDM	0.696 $\pm$ 0.008	0.823 $\pm$ 0.034	0.564 $\pm$ 0.003
Ours	<b>0.781</b> $\pm$ 0.005	<b>0.623</b> $\pm$ 0.063	<b>0.643</b> $\pm$ 0.001

Table 2. **Ablation study.** We show precision (Top-3), FID, and penetration. Our configuration can achieve the best results.

### 4.3. Ablation Study

In this section, we examine the roles of three crucial components in our method: PDM, PMP, and ICF. We present the comparative results in Tab. 2 and Fig. 5.

**PDM.** We evaluate the effect of PDM through an ablation study. Specifically, we benchmark our HOIAnimator (‘Ours’) against a variant devoid of the PDM (‘w/o PDM’), which employs a single diffusion model. Our results are visually represented in Fig. 5, where we demonstrate that our proposed method is more effective in accurately capturing the spatial relationships between humans and objects. Furthermore, when assessed in terms of precision (Top-3), Frechet Inception Distance (FID), and penetration score, our configuration outperforms the ‘w/o PDM’ model, indicating its superior performance.

**PMP.** To evaluate the effectiveness of our proposed PMP. This section evaluates our model in both the absence (‘w/o

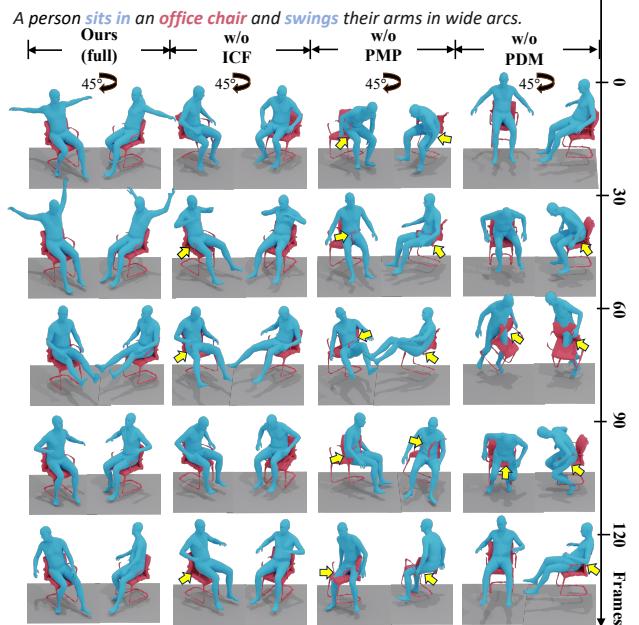


Figure 5. **Ablation study.** Based on textual descriptions, our model generates predicted HOI animation. Simultaneously, we apply a 45° rotation to the right side of each result and use yellow arrows to point out interaction errors to facilitate the comparison of their quality.

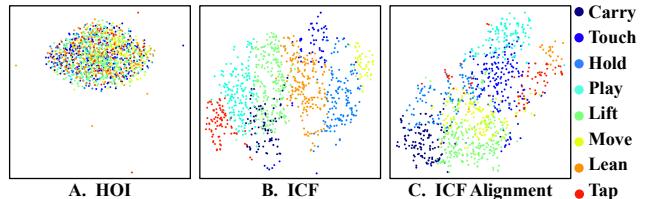


Figure 6. **Tsne [38] of HOI animation.** The feature is by HOI encoder (A), ICF (B), and HOI animation after ICF alignment (C).

PMP) and our approach. Fig. 5 clearly illustrates that while the spatial arrangement between individuals and objects appears normal, there is an absence of interactive dynamics.

**ICF.** We remove the ICF (‘w/o ICF’). As depicted in Fig. 5, the exclusion of the components results in generally acceptable HOI animations; nevertheless, these animations are marred by specific inaccuracies, such as unrealistic penetrations. This observation highlights the enhanced posi-

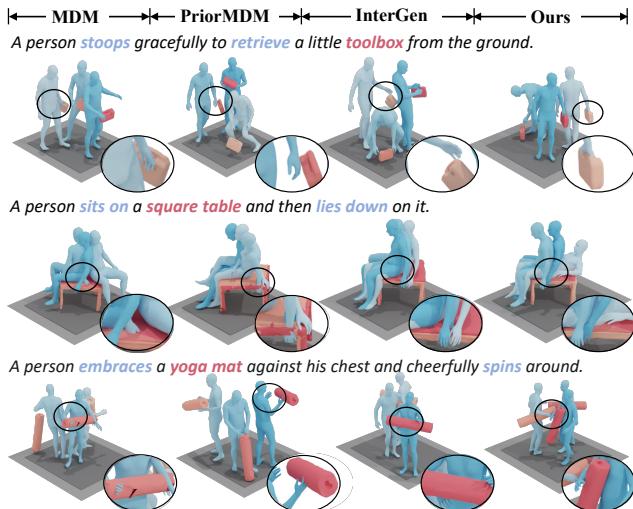


Figure 7. **Qualitative evaluation.** We present zoomed-in details highlighted within black boxes. For any specified text description, only our HOIAnimator is capable of accurately depicting the spatial relationships and the dynamic interactions involved.

504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
tional accuracy afforded by our proposed configuration. At the same time, it can be seen from Fig. 6 that ICF has great guidance on HOI animation.

#### 4.4. Qualitative Evaluation

To illustrate the effectiveness of HOIAnimator, we provide a qualitative comparison between previous works [13, 16, 32] and HOIAnimator. As shown in Fig. 7, HOIAnimator stands out as the only method capable of effectively translating textual descriptions that encompass the positional relationships and interactive dynamics between humans and objects. In comparison, MDM struggles with precise spatial positioning of individuals and objects. Although Prior MDM is adept at classifying human-object interaction (HOI) actions, it lacks in detailing the nuances of interaction dynamics. On the other hand, InterGen effectively understands these dynamics, yet it does not consistently execute interactions accurately. Through evaluated examples, HOIAnimator consistently demonstrates its capability in both structuring and generating these complex interactions effectively.

#### 4.5. User Study

To further evaluate the quality of our generated HOI animations, we conduct a user study to evaluate the quality of HOI animations. The study involves 40 participants of various backgrounds, including 22 students, 3 sales staff, 6 software engineers, 2 teachers, 3 managers, and 4 individuals of other professions. In the study, we randomly select 9 motion labels and 20 object labels and combine them to create 10 meaningful descriptions of HOI animations. Based on these coherent HOI descriptions, we generate synthetic animations and shuffle them using scoring methods for pre-

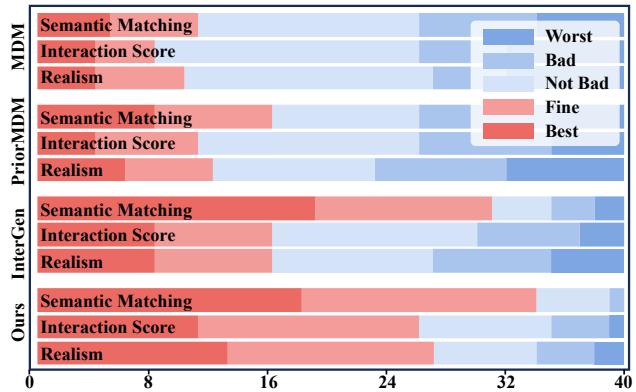


Figure 8. **User study.** The color bars in the figure indicate the percentage of the scores. The X-axis represents the number of participants. The results show that our method outperforms other baselines in terms of semantic matching, interaction score, and realism.

sentation. After that, we ask users to rate the synthetic animations on three aspects: (1) Semantic Matching: The generated animations match the semantics of the given text descriptions. (2) Interaction Score: The quality of the poses and interactions between humans and objects in the animation. (3) Realism: The level of realism in the motion of the characters. As shown in Fig. 8, the results demonstrate that our method surpasses other baselines in terms of semantic matching, interaction score, and realism.

#### 4.6. Limitation and Discussion

Although HOIAnimator generates realistic HOI animation, it still has some limitations. First, the HOIAnimator is less adept at depicting complex sequences of interactions. Second, it is limited to scenes involving multiple objects interacting simultaneously. Furthermore, our method does not support nonrigid objects animations (e.g., water). Including the deformation prior into a current framework for high-quality deformable generation is promising, but it requires much more diverse training data.

### 5. Conclusion and Future Work

In this paper, we propose the HOIAnimator to convert text instructions into detailed animations of human and object interactions. Our key innovation is the PDM, which could closely align human and object movements with their corresponding text descriptions. Additionally, we have developed an ICF. This field actively influences animation, ensuring it mirrors the precise and diverse nature of interactions observed in the real world. The results demonstrate that HOIAnimator excels at creating dynamic and context-aware animations. In future work, we will improve HOIAnimator to better handle complex, sequential actions and interactions involving multiple objects.

567 **References**

- [1] Hyemin Ahn, Esteve Valls Mascaro, and Dongheui Lee. Can we use diffusion probabilistic models for 3d motion prediction? *arXiv preprint arXiv:2302.14503*, 2023. 3
- [2] Samaneh Azadi, Thomas Hayes, Akbar Shah, Guan Pang, Devi Parikh, and Sonal Gupta. Text-conditional contextualized avatars for zero-shot personalization. *arXiv preprint arXiv:2304.07410*, 2023. 3
- [3] German Barquero, Sergio Escalera, and Cristina Palmero. Belfusion: Latent diffusion for behavior-driven human motion prediction. In *International Conference on Computer Vision*, pages 2317–2327, 2023. 3
- [4] Bharat Lal Bhatnagar, Xianghui Xie, Ilya A Petrov, Cristian Sminchisescu, Christian Theobalt, and Gerard Pons-Moll. Behave: Dataset and method for tracking human object interactions. In *Computer Vision and Pattern Recognition*, pages 15935–15946, 2022. 3, 6, 7
- [5] Zhe Cao, Hang Gao, Karttikeya Mangalam, Qi-Zhi Cai, Minh Vo, and Jitendra Malik. Long-term human motion prediction with scene context. In *European Conference on Computer Vision*, pages 387–404, 2020. 3
- [6] Ling-Hao Chen, Jiawei Zhang, Yewen Li, Yiren Pang, Xiaobo Xia, and Tongliang Liu. Humanmac: Masked motion completion for human motion prediction. In *International Conference on Computer Vision*, pages 9544–9555, 2023. 3
- [7] Xin Chen, Biao Jiang, Wen Liu, Zilong Huang, Bin Fu, Tao Chen, and Gang Yu. Executing your commands via motion diffusion in latent space. In *Computer Vision and Pattern Recognition*, pages 18000–18010, 2023. 2, 6, 7
- [8] Yixin Chen, Sai Kumar Dwivedi, Michael J Black, and Dimitrios Tzionas. Detecting human-object contact in images. In *Computer Vision and Pattern Recognition*, pages 17100–17110, 2023. 3
- [9] Zicong Fan, Omid Taheri, Dimitrios Tzionas, Muhammed Kocabas, Manuel Kaufmann, Michael J Black, and Otmar Hilliges. Arctic: A dataset for dexterous bimanual hand-object manipulation. In *Computer Vision and Pattern Recognition*, pages 12943–12954, 2023. 3
- [10] Anindita Ghosh, Rishabh Dabral, Vladislav Golyanik, Christian Theobalt, and Philipp Slusallek. Imos: Intent-driven full-body motion synthesis for human-object interactions. In *Computer Graphics Forum*, pages 1–12, 2023. 3
- [11] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in Neural Information Processing Systems*, 27, 2014. 4
- [12] Chuan Guo, Shihao Zou, Xinxin Zuo, Sen Wang, Wei Ji, Xingyu Li, and Li Cheng. Generating diverse and natural 3d human motions from text. In *Computer Vision and Pattern Recognition*, pages 5152–5161, 2022. 2, 6, 7
- [13] Chuan Guo, Xinxin Zuo, Sen Wang, and Li Cheng. Tm2t: Stochastic and tokenized modeling for the reciprocal generation of 3d human motions and texts. In *European Conference on Computer Vision*, pages 580–597, 2022. 2, 7, 8
- [14] Sanjay Haresh, Xiaohao Sun, Hanxiao Jiang, Angel X Chang, and Manolis Savva. Articulated 3d human-object interactions from rgb videos: An empirical analysis of approaches and challenges. In *International Conference on 3D Vision*, pages 312–321, 2022. 3
- [15] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020. 3, 4
- [16] Siyuan Huang, Zan Wang, Puha Li, Baoxiong Jia, Tengyu Liu, Yixin Zhu, Wei Liang, and Song-Chun Zhu. Diffusion-based generation, optimization, and planning in 3d scenes. In *Computer Vision and Pattern Recognition*, pages 16750–16761, 2023. 2, 3, 4, 6, 7, 8
- [17] Yinghao Huang, Omid Taheri, Michael J Black, and Dimitrios Tzionas. Intercap: Joint markerless 3d tracking of humans and objects in interaction. In *DAGM German Conference on Pattern Recognition*, pages 281–299, 2022. 3, 6
- [18] Nan Jiang, Tengyu Liu, Zhexuan Cao, Jieming Cui, Yixin Chen, He Wang, Yixin Zhu, and Siyuan Huang. Chairs: Towards full-body articulated human-object interaction. *arXiv preprint arXiv:2212.10621*, 2022. 3
- [19] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 4
- [20] Nilesh Kulkarni, Davis Rempe, Kyle Genova, Abhijit Kundu, Justin Johnson, David Fouhey, and Leonidas Guibas. Nifty: Neural object interaction fields for guided human motion synthesis. *arXiv preprint arXiv:2307.07511*, 2023. 3, 6
- [21] Manling Li, Ruochen Xu, Shuohang Wang, Luowei Zhou, Xudong Lin, Chenguang Zhu, Michael Zeng, Heng Ji, and Shih-Fu Chang. Clip-event: Connecting text and images with event structures. In *Computer Vision and Pattern Recognition*, pages 16420–16429, 2022. 2
- [22] Han Liang, Wenqian Zhang, Wenxuan Li, Jingyi Yu, and Lan Xu. InterGen: Diffusion-based multi-human motion generation under complex interactions. *arXiv preprint arXiv:2304.05684*, 2023. 2, 3, 5, 6
- [23] Nan Liu, Shuang Li, Yilun Du, Antonio Torralba, and Joshua B Tenenbaum. Compositional visual generation with composable diffusion models. In *European Conference on Computer Vision*, pages 423–439, 2022. 3
- [24] Xingang Pan, Ayush Tewari, Thomas Leimkühler, Lingjie Liu, Abhimitra Meka, and Christian Theobalt. Drag your gan: Interactive point-based manipulation on the generative image manifold. *arXiv preprint arXiv:2305.10973*, 2023. 2
- [25] Jeesung Park, Jin-Woo Park, and Jong-Seok Lee. Viplo: Vision transformer based pose-conditioned self-loop graph for human-object interaction detection. In *Computer Vision and Pattern Recognition*, pages 17152–17162, 2023. 3
- [26] Minho Park, Jooyeon Yun, Seunghwan Choi, and Jaegul Choo. Learning to generate semantic layouts for higher text-image correspondence in text-to-image synthesis. In *International Conference on Computer Vision*, pages 7591–7600, 2023. 2
- [27] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed AA Osman, Dimitrios Tzionas, and Michael J Black. Expressive body capture: 3d hands, face, and body from a single image. In *Computer Vision and Pattern Recognition*, pages 10975–10985, 2019. 4

- 680 [28] Mathis Petrovich, Michael J Black, and G  l Varol. Action- 737  
681 conditioned 3d human motion synthesis with transformer 738  
682 vae. In *International Conference on Computer Vision*, pages 739  
683 10985–10995, 2021. 2 740  
684 [29] Mathis Petrovich, Michael J Black, and G  l Varol. Temos: 741  
685 Generating diverse human motions from textual descriptions. 742  
686 In *European Conference on Computer Vision*, pages 480– 743  
687 497, 2022. 2 744  
688 [30] Sigal Raab, Inbal Leibovitch, Guy Tevet, Moab Arar, Amit H 745  
689 Bermano, and Daniel Cohen-Or. Single motion diffusion. 746  
690 *arXiv preprint arXiv:2302.05905*, 2023. 3 747  
691 [31] Robin Rombach, Andreas Blattmann, Dominik Lorenz, 748  
692 Patrick Esser, and Bj  rn Ommer. High-resolution image 749  
693 synthesis with latent diffusion models. In *Computer Vision and 750  
694 Pattern Recognition*, pages 10684–10695, 2022. 4 751  
695 [32] Yonatan Shafir, Guy Tevet, Roy Kapon, and Amit H 752  
696 Bermano. Human motion diffusion as a generative prior. 753  
697 *arXiv preprint arXiv:2303.01418*, 2023. 3, 5, 6, 7, 8 754  
698 [33] Jiaming Song, Chenlin Meng, and Stefano Ermon. 755  
699 Denoising diffusion implicit models. *arXiv preprint 756  
700 arXiv:2010.02502*, 2020. 3 757  
701 [34] Wenfeng Song, Xinyu Zhang, Yuting Guo, Shuai Li, Aimin 758  
702 Hao, and Hong Qin. Automatic Generation of 3D Scene 759  
703 Animation Based on Dynamic Knowledge Graphs and Contextual 760  
704 Encoding. *International Journal of Computer Vision*, 761  
705 131(11):2816–2844, 2023. 3 762  
706 [35] Purva Tendulkar, D  dac Sur  s, and Carl Vondrick. Flex: 763  
707 Full-body grasping without full-body grasps. In *Computer Vision 764  
708 and Pattern Recognition*, pages 21179–21189, 2023. 3 765  
709 [36] Guy Tevet, Brian Gordon, Amir Hertz, Amit H Bermano, 766  
710 and Daniel Cohen-Or. Motionclip: Exposing human motion 767  
711 generation to clip space. In *European Conference on 768  
712 Computer Vision*, pages 358–374, 2022. 2, 6, 7 769  
713 [37] Guy Tevet, Sigal Raab, Brian Gordon, Yoni Shafir, Daniel 770  
714 Cohen-or, and Amit Haim Bermano. Human motion diffusion 771  
715 model. In *International Conference on Learning Representations*, 772  
716 2023. 3, 4 773  
717 [38] Laurens Van der Maaten and Geoffrey Hinton. Visualizing 774  
718 data using t-sne. *Journal of machine learning research*, 9 775  
719 (11), 2008. 7 776  
720 [39] Jingbo Wang, Yu Rong, Jingyuan Liu, Sijie Yan, Dahua Lin, 777  
721 and Bo Dai. Towards diverse and natural scene-aware 3d 778  
722 human motion synthesis. In *Computer Vision and Pattern 779  
723 Recognition*, pages 20460–20469, 2022. 3 780  
724 [40] Xi Wang, Gen Li, Yen-Ling Kuo, Muhammed Kocabas, 781  
725 Emre Aksan, and Otmar Hilliges. Reconstructing action- 782  
726 conditioned human-object interactions using commonsense 783  
727 knowledge priors. In *International Conference on 3D Vision*, 784  
728 pages 353–362, 2022. 3 785  
729 [41] Dong Wei, Xiaoning Sun, Huaijiang Sun, Bin Li, Shengxiang 786  
730 Hu, Weiqing Li, and Jianfeng Lu. Understanding text- 787  
731 driven motion synthesis with keyframe collaboration via 788  
732 diffusion models. *arXiv preprint arXiv:2305.13773*, 2023. 3 789  
733 [42] Aming Wu and Cheng Deng. Discriminating known from 790  
734 unknown objects via structure-enhanced recurrent variational 791  
735 autoencoder. In *Computer Vision and Pattern Recognition*, 792  
736 pages 23956–23965, 2023. 2