# HOI-Diff: Text-Driven Synthesis of 3D Human-Object Interactions using Diffusion Models

Xiaogang Peng[1,2*], Yiming Xie[1*], Zizhao Wu[2], Varun Jampani[3], Deqing Sun[4], and Huaizu Jiang[1]

[1] Northeastern University, USA
[2] Hangzhou Dianzi University, China
[3] Stability AI
[4] Google Research
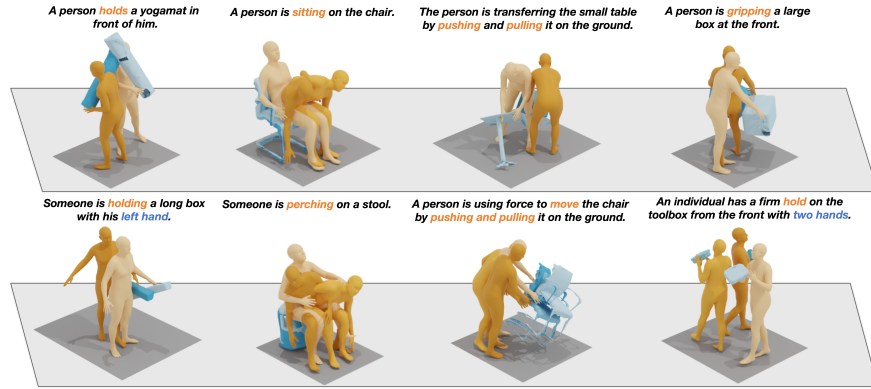https://neu-vi.github.io/HOI-Diff/

**Fig. 1: HOI-Diff generates realistic motions for 3D human-object interactions given a text prompt and object geometry.** Please see the sup. mat. for video results. *Darker color indicates later frames in the sequence. Best viewed in color.*

**Abstract.** We address the problem of generating realistic 3D human-object interactions (HOIs) driven by textual prompts. To this end, we take a modular design and decompose the complex task into simpler subtasks. We first develop a dual-branch diffusion model (HOI-DM) to generate both human and object motions conditioned on the input text, and encourage coherent motions by a cross-attention communication module between the human and object motion generation branches. We also develop an affordance prediction diffusion model (APDM) to predict the contacting area between the human and object during the interactions driven by the textual prompt. The APDM is independent of the results by the HOI-DM and thus can correct potential errors by the latter. Moreover, it stochastically generates the contacting points to diversify the generated motions. Finally, we incorporate the estimated contacting points into the classifier-guidance to achieve accurate and close contact

---

* Equal contribution.

between humans and objects. To train and evaluate our approach, we annotate BEHAVE dataset with text descriptions. Experimental results on BEHAVE and OMOMO demonstrate that our approach produces realistic HOIs with various interactions and different types of objects.

## 1   Introduction

Text-driven synthesis of 3D human-object interactions (HOIs) aims to generate motions for both the human and object that form coherent and semantically meaningful interactions. It enables virtual humans to naturally interact with objects, which has a wide range of applications in AR/VR, video games, and filmmaking, etc.

The generation of natural and physically plausible 3D HOIs involves humans interacting with *dynamic* objects in *various* ways according to the text prompts, thereby posing several challenges. First, the variability of object shapes makes it particularly challenging to generate semantically meaningful contact between the human and object to avoid floating objects. Second, the generated HOIs should be faithful to the input text prompts as there are many plausible interactions between human and the same object (*e.g.*, a person carries a chair, sits on a chair, pushes or pulls a chair). Text-driven 3D HOI synthesis with a diverse set of interactions is still under-explored. Third, the development and evaluation of 3D HOI synthesis models requires a high-quality human motion dataset with various HOIs and textual descriptions, but existing datasets lack either HOIs [11,34] or textual descriptions [3].

On one hand, recent methods [12,17,23,33,42,44,52,58] can synthesize realistic human motions for HOIs for *static* objects only. They usually synthesize the motion in the last mile of interaction, *i.e.* the motion between the given starting human pose and the final interaction pose, and overlook the movement of the objects when the human is interacting with them. On the other hand, existing methods for motion generation with dynamic objects do not adequately reflect real-world complexity. For instance, they focus on grasping small objects [10], provide the object motion as conditioning [25], predict deterministic interactions between the human and the same object without the diversity [37,54], consider only a small set of interactions (*e.g.* sit/lift [23], sit/lie down [12], sit [17,33,58], grasp [44,52]), or investigate a single type of object (*e.g.* chair [17,58]).

In this paper, we go beyond and propose **HOI-Diff** for 3D HOIs synthesis involving humans interacting with different types of objects in diverse ways, which are both physically plausible and semantically faithful to the textual prompt, as shown in Fig. 1. Our key insight is to decompose 3D HOIs synthesis into three modules: (a) **coarse 3D HOIs generation** that extends the human motion diffusion model [45] to a dual-branch diffusion model (HOI-DM) to generate both human and object motions conditioning on the input text prompt. To encourage coherent motions, we develop a cross-attention communication module, exchanging information between the human and object motion generation models; (b) **affordance prediction diffusion model** (APDM) that estimates

the contacting points between the human and object during the interactions driven by the textual prompt. Our APDM does not rely on the results of the HOI-DM and thus can recover from its potential errors. Moreover, it stochastically generates the contacting points to diversity the generated motions; and (c) **affordance-guided interaction correction** that incorporates the estimated contacting information and employs the classifier-guidance to achieve accurate and close contact between humans and objects, significantly alleviating the cases of floating objects. Compared with designing a monolithic model, our proposed HOI-Diff disentangles motion generation for humans and objects and estimation of their contacting points, which are later integrated to form coherent and diverse HOIs, reducing the complexity and burden for each of the three modules.

To our best knowledge, HOI-Diff is one of the first text-driven approachs capable of generating realistic motions for various HOIs and different types of objects during the interactions. While our approach pioneers this capability, we note the existence of concurrent works such as [9, 24, 49], which also explore similar directions. For both training and evaluation purposes, we annotate each video sequence in BEHAVE dataset [3] with text descriptions, which mitigates the issue of severe data scarcity for text-driven 3D HOIs generation. In addition, we evaluate our approach on the OMOMO dataset [25], which focuses on the manipulation of two hands. Extensive experiments validate the effectiveness and design choices of our approach, particularly for dynamic objects, thereby enabling a set of new applications in human motion generation.

## 2   Related Work

**Human Motion Generation with Diffusion Models.** The denoising diffusion models have been widely used 2D image generations [36,39,40] and achieved impressive results. Recent work [1, 2, 4–6, 18, 38, 41, 43, 45, 46, 51, 53, 55–57, 59] apply the diffusion model in the task of human motion generation. While these methods have successfully generated human motion, they usually generate isolated motions in the free space without considering the objects the human is interacting with. Our method is primarily focused on motion generation with human-object interactions.

**Scene- and Object-Aware Human Motion Generation.** Recent works condition motion synthesis on scene geometry [15, 48, 50, 60]. This facilitates the understanding of human-scene interactions. However, the motion fidelity is compromised due to the lack of paired full scene-motion data. Other approaches [12, 17, 23, 33, 42, 58] instead focus on the interactions with the objects and can produce realistic motions. However, they focus on interacting with static objects with limited interactions. OMOMO [25] proposes a conditional diffusion framework that can generate full-body motion from the object motion. The object motion is needed as input in OMOMO, whereas our method can jointly synthesize human motion and object motion. IMoS [10] synthesizes the full-body human along with the 3D object motions from textual inputs. Although

they can generate realistic interactions, IMoS only focuses on grasping small objects with hands. InterDiff [54] predicts whole-body interactions with dynamic objects. However, it focuses on the motion prediction task, where the future motion needs to be forecasted based on past motion information and object information. Note that the interaction type is deterministic, given the object and past motion information. Different from this, we tackle the motion synthesis task, where the interaction with the same object can be controlled by the text prompt. Recently, there has been a surge of interest in the text-driven synthesis of 3D human-object interactions for dynamic objects, resulting in the development of concurrent works [9, 24, 49].

**Multi-Human Motion Generation.** Recent works [26, 41, 41] incorporate human-to-human interactions into the motion diffusion process. They usually also adopt cross-attention to exchange information from different human motions. Different from them, our goal is to generate coarse HOIs instead of doing so in a single run, which alleviates the burden of the model.

**Affordance Estimation.** The affordance estimation on 3D point cloud is studied in [7, 16, 20–22, 29, 30]. [20] formulates this task as the point cloud segmentation task from the extracted geometric features. [21] voxelizes the point cloud and creates affordance maps with interactive manipulation. [16,22] focuses on the object grasping affordances and some specific applications. [29] learns the affordance for object-object interactions. [30] addressed the task of open-vocabulary affordance learning. Overall affordance learning is a very challenging task even the shape of point cloud is given. Instead of predicting the point-wise contact labels, in our method, we propose to make some hypotheses to simplify human-object interactions, making it more tractable without significantly compromising accuracy.

## 3     Method

### 3.1     Overview

**Motion Representations.** We denote a 3D HOI sequence as $\boldsymbol{x} = \{\boldsymbol{x}^h, \boldsymbol{x}^o\}$. It consists of human motion sequence $\boldsymbol{x}^h \in \mathbb{R}^{L \times D^h}$ and object motion sequence $\boldsymbol{x}^o \in \mathbb{R}^{L \times D^o}$, where $L$ denotes the length of the sequence. For $\boldsymbol{x}^h$, we adopt the redundant representation widely used in human motion generation [11] with $D^h = 263$, which include pelvis velocity, local joint positions, velocities and rotations of other joints in the pelvis space, and binary foot-ground contact labels. For the object motion sequence $\boldsymbol{x}^o$, we assume the object geometry is given as an input, and thus we only need to estimate its 6DoF poses in the generation, $i.e.$, $D^o = 6$. We represent each object instance as a point cloud of 512 points $\boldsymbol{p} \in \mathbb{R}^{512 \times 3}$.

**Diffusion Model for 3D HOI Generation.** Given a prompt $\boldsymbol{c} = (\boldsymbol{d}, \boldsymbol{p})$, consisting of a textual description $\boldsymbol{d}$ and the object instance's point cloud $\boldsymbol{p}$, a
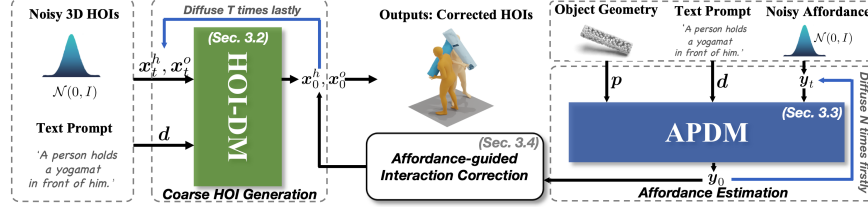
**Fig. 2: Overview of HOI-Diff for 3D HOIs generation using diffusion models.**
Our key insight is to decompose the generation task into three modules: (a) coarse 3D
HOI generation using a dual-branch diffusion model (HOI-DM), (b) affordance predic-
tion diffusion model (APDM) to estimate the contacting points of humans and objects,
and (c) afforance-guided interaction correction, which incorporates the estimated con-
tacting information and employs the classifier-guidance to achieve accurate and close
contact between humans and objects to form coherent HOIs.

diffusion model $p_\theta(\boldsymbol{x}_{t-1}|\boldsymbol{x}_t, \boldsymbol{c})$[5] learns the reverse diffusion process to generate
clean data from a Gaussian noise $\boldsymbol{x}_t$ with $T$ consecutive denoising steps

$$p_\theta(\boldsymbol{x}_{t-1}|\boldsymbol{x}_t, \boldsymbol{c}) := \mathcal{N}(\boldsymbol{x}_{t-1}, \mu_\theta(\boldsymbol{x}_t, t, \boldsymbol{c}), (1 - \alpha_t)\mathbf{I}), \quad (1)$$

where $t$ is the denoising step. Following [45], instead of predicting the noise $\epsilon_t$ in
each diffusion step, our diffusion model $M_\theta$ with parameters $\theta$ predicts the final
clean motion $\boldsymbol{x}_0 = M_\theta(\boldsymbol{x}_t, t, \boldsymbol{c})$. We sample $\mathbf{x}_{t-1} \sim \mathcal{N}(\boldsymbol{\mu}_t, \Sigma_t)$ and compute the
mean as in [31]

$$\boldsymbol{\mu}_t = \frac{\sqrt{\alpha_{t-1}}\beta_t}{1 - \alpha_t}\boldsymbol{x}_0 + \frac{\sqrt{1 - \beta_t}(1 - \alpha_{t-1})}{1 - \alpha_t}\boldsymbol{x}_t, \quad (2)$$

where $\alpha_t = \prod_{s=1}^{t}(1-\beta_s)$ and $\beta_t \in (0, 1)$ are the variance schedule. $\Sigma_t = \frac{1-\alpha_{t-1}}{1-\alpha_t}\beta_t$
[14] is a variance scheduler of choice. Similar to $\boldsymbol{x}_t$, $\boldsymbol{\mu}_t$ consists of $\boldsymbol{\mu}_t^h$ and $\boldsymbol{\mu}_t^o$,
corresponding to human and object motion, respectively.

Simply adopting the diffusion model described in Eq.(1) would impose a huge
burden on the model, which requires joint generation of human and object mo-
tion and more critically, enforcement of their intricate interactions to follow the
input textual description. In this paper, we propose **HOI-Diff** for 3D HOIs gen-
eration, disentangling motion generation for humans and objects and estimation
of their contacting points. They are later integrated to form coherent and diverse
HOIs, which reduces the complexity and burden for each of the three modules,
leading to better generation performance as evidenced by our experiments.

Fig. 2 shows the overview of our proposed approach. We introduce a dual-
branch Human-Object Interaction Diffusion Model (HOI-DM), which can pro-
duce diverse yet consistent motions, capturing the intricate interplay and mu-
tual interactions between humans and objects (Sec. 3.2). To ensure physically

---

[5] We use superscripts $h$ and $o$ to denote human and object sequence, respectively.
Without a superscript, it means the 3D HOI sequence, containing both $\boldsymbol{x}^h$ and $\boldsymbol{x}^o$.
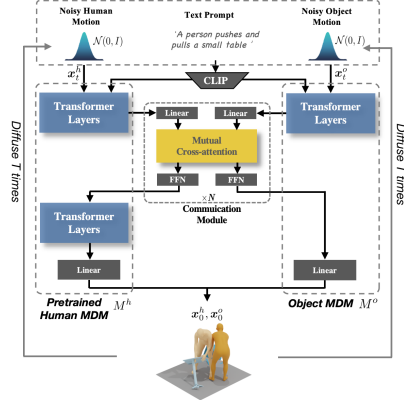Subscript is used for the diffusion denoising step.

**Fig. 3: Illustration of HOI-DM architecture for coarse 3D HOIs generation.** It has two branches designed for generating human and object motions individually. A mutual cross-attention is introduced to allow information exchange between two branches to generate coherent motions. The human motion model $M^h$ finetunes a pretrained human MDM [45].
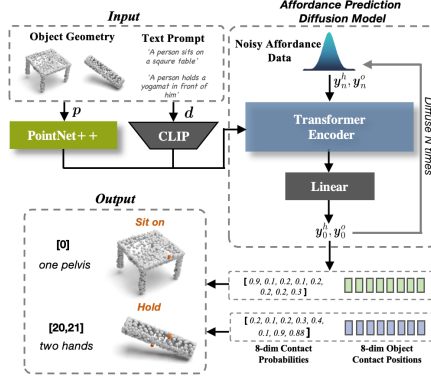
**Fig. 4: Illustration of APDM architecture for affordance estimation.** Affordance information of human contact labels, object contact positions, and binary object states are represented together as a noise variable, which is fed into the Transformer encoder to generate clean estimation. The object point cloud and textual prompt are taken as conditional input.

plausible contact between humans and objects, we propose a novel affordance prediction diffusion model (APDM) (Sec. 3.3), whose output will be used as classifier guidance (Sec. 3.4) to correct the interactions at each diffusion step of human/object motion generation.

### 3.2 Coarse 3D HOIs Generation

First, we introduce a dual-branch Human-Object Interaction Diffusion Model (HOI-DM) to generate human and object motions that are roughly coherent. As shown in Fig. 3, it consists of two Transformer models [47], human motion diffusion model (MDM) $M^h$ and object MDM $M^o$, which work similar to [45]. Specifically, at the diffusion step $t$, they take the text description and noisy motions $\boldsymbol{x}_t^h$ and $\boldsymbol{x}_t^o$ as input and predict clean human and object motions $\boldsymbol{x}_0^h$ and $\boldsymbol{x}_0^o$, respectively.

To enhance the learning of interactions of the human and object when generating their motion, we introduce a Communication Module ($CM$) designed for exchanging feature representations between the human MDM $M^h$ and the object MDM $M^o$. $CM$ is a Transformer block that receives the intermediate feature $\boldsymbol{f}^h, \boldsymbol{f}^o$ from both $M^h$ and $M^o$. It then processes these inputs to generate refined updates based on the cross attention mechanism [47]. The updated feature representations $\tilde{\boldsymbol{f}}_h$ and $\tilde{\boldsymbol{f}}_o$ of the human and object are then conditioned on each other, which are then fed into the subsequent layers of their respective

branches to estimate clean human and object motion $\boldsymbol{x}_0^h$ and $\boldsymbol{x}_0^o$, respectively. The $CM$ is inserted at the 4th transformer layer for human MDM and the last layer for object MDM, which was empirically found to work better.

Given the limited data availability for 3D HOI generation, during training, the human motion model $M^h$ finetunes a pretrained human MDM [45]. This fine-tuning is critical to ensure the smoothness of the generated human motions. We ablate this design choice in Sec. 4.3. Object MDM is trained from scratch. We modify the input and output linear layers to take in the object motion which has a different dimension from the human motion. More methodology details of HOI-DM can be found in the supplementary material.

### 3.3   Affordance Estimation

Due to the complexity of the interactions between a human and object, HOI-DM alone usually fails to produce physically plausible results, leading to floating objects or penetrations. To improve the generation of intricate interactions, the problem that needs to be solved is to *identify where the contacting areas are* between the human and object. InterDiff [54] defines the contacting area based on the distance measurement between the surface of human and object. This approach, however, heavily relies on the quality of the generated human and object motions and cannot recover from errors in the coarse 3D HOI results. In addition, the contact area is diverse even with the same object and interaction type, *e.g.*, "sit" can happen on either side of a table. To this end, we introduce an Affordance Prediction Diffusion Model (APDM) for affordance estimation. As illustrated in Fig. 4, the input includes a text description $\boldsymbol{d}$ and the object point cloud $\boldsymbol{p}$. Our APDM doesn't rely on the results of the HOI-DM and thus can recover from the potential errors in HOI-DM. In addition, it stochastically generates the contacting points to ensure the diversity of the generated motions.

Affordance estimation in 3D point clouds itself is a notably challenging problem [7, 16, 20–22, 29, 30], especially in the context of 3D HOI generation involving textual prompt. In this paper, we consider eight primary body joints – the `pelvis, neck, feet, shoulders`, and `hands` – as the interacting parts in HOI scenarios. It can effectively model common interactions such as grasping an object with both hands, sitting actions involving the pelvis and back, or lifting with a single hand. We use binary contact labels to determine which joints are in contact with the object. Subsequently, we predict eight corresponding contact points on the object surface, identified as the points closest to the selected body joints. Note that the binary contact label estimation for different body joints are independent, allowing us to handle complex HOIs.

Specifically, at each diffusion time step $n$ of APDM[6], the noisy data consists of human contact labels representing the contact status for the eight primary body joints, denoted as $\boldsymbol{y}_n^h \in \{0,1\}^8$, and the eight corresponding contact points on the object surface, denoted as $\boldsymbol{y}_n^o \in \mathbb{R}^{8 \times 3}$. The model is designed to predict both

---

[6] We note that APDM and HOI-DM work independently. We thus use two symbols to denote the different diffusion time steps to avoid confusion.

contact probabilities and contact positions. Subsequently, dynamic selection of contacting body joints is performed by considering predicted probabilities over a specific threshold $\tau$ (set to be 0.6). The corresponding contact points on the object are then determined based on the selected joints. APDM works similar to the diffusion denoising process described in Eq.(1). Besides, we utilize a large language model (ChatGPT) to determine whether the object state $\boldsymbol{y}_0^s \in \{0,1\}$ should be set to static ($\boldsymbol{y}_0^s = 1$) based on the textual description, which can help us better process static objects when synthesizing 3D HOIs, as discussed in the following section. All the clean affordance data is grouped as $\boldsymbol{y}_0 = (\boldsymbol{y}_0^h, \boldsymbol{y}_0^o, \boldsymbol{y}_0^s)$.

### 3.4  Affordance-guided Interaction Correction

With the estimated affordance, we can better align human and object motions to form coherent interactions. To this end, we propose to use the classifier guidance [8] to achieve accurate and close contact between humans and objects, significantly alleviating the cases of floating objects.

---

**Algorithm 1** Affordance-guided Interaction Correction

**Require:** Input $\boldsymbol{c} = (\boldsymbol{d}, \boldsymbol{p})$ consisting of a textual description $\boldsymbol{d}$ and object point cloud $\boldsymbol{p}$, HOI-Diff model $\mathbf{M}_\theta$, objective function $G(\boldsymbol{\mu}_t^h, \boldsymbol{\mu}_t^o, \boldsymbol{y}_0)$, and estimated affordance $\boldsymbol{y}_0 = (\boldsymbol{y}_0^h, \boldsymbol{y}_0^o, \boldsymbol{y}_0^s)$.
1:  $\boldsymbol{x}_T^h, \boldsymbol{x}_T^o \leftarrow$ sample from $\mathcal{N}(\mathbf{0}, \mathbf{I})$
2:  $K = 1$
3:  **for all** $t$ from $T$ to 1 **do**
4:      $\boldsymbol{x}_0^h, \boldsymbol{x}_0^o \leftarrow M_\theta(\boldsymbol{x}_t^h, \boldsymbol{x}_t^o, t, \boldsymbol{c})$ # Get $\mu_t^h, \mu_t^o$ according to Eq.(2) with $\Sigma_t$
5:      **if** $t = 1$ **then**
6:          $K = 100$
7:      **end if**
8:      **for all** $k$ from $K$ to 1 **do** # Separately perturb
9:          $\boldsymbol{\mu}_t^h \leftarrow \boldsymbol{\mu}_t^h - \tau_1 \Sigma_t \nabla_{\boldsymbol{\mu}_t^h} G(\boldsymbol{\mu}_t^h, \boldsymbol{\mu}_t^o, \boldsymbol{y}_0), \quad \boldsymbol{\mu}_t^o \leftarrow \boldsymbol{\mu}_t^o - \tau_2 \Sigma_t \nabla_{\mu_t^o} G(\boldsymbol{\mu}_t^h, \boldsymbol{\mu}_t^o, \boldsymbol{y}_0)$
10:     **end for**
11:     $\mathbf{x}_{t-1}^h \sim \mathcal{N}(\boldsymbol{\mu}_t^h, \Sigma_t), \quad \mathbf{x}_{t-1}^o \sim \mathcal{N}(\boldsymbol{\mu}_t^o, \Sigma_t)$
12: **end for**
13: **return** $\boldsymbol{x}_0^h, \boldsymbol{x}_0^o$

---

Specifically, in a nutshell, we define an analytic function $G(\boldsymbol{\mu}_t^h, \boldsymbol{\mu}_t^o, \boldsymbol{y}_0)$ that assesses how closely the generated human joints and object's 6DoF pose align with a desired objective. In our case, it enforces the contact positions of human and object to be close to each other and their motions are smooth temporally. Based on the gradient of $G(\boldsymbol{\mu}_t^h, \boldsymbol{\mu}_t^o, \boldsymbol{y}_0)$, we can perturb the generated human and object motion at each diffusion step $t$ as in [19, 23, 53],

$$\boldsymbol{\mu}_t^h = \boldsymbol{\mu}_t^h - \tau_1 \Sigma_t \nabla_{\boldsymbol{\mu}_t^h} G(\boldsymbol{\mu}_t^h, \boldsymbol{\mu}_t^o, \boldsymbol{y}_0), \tag{3}$$

$$\boldsymbol{\mu}_t^o = \boldsymbol{\mu}_t^o - \tau_2 \Sigma_t \nabla_{\boldsymbol{\mu}_t^o} G(\boldsymbol{\mu}_t^h, \boldsymbol{\mu}_t^o, \boldsymbol{y}_0). \tag{4}$$

Here $\tau_1$ and $\tau_2$ are different strengths to control the guidance for human and object motion, respectively. Due to the sparseness of object motion features, we assign a larger value to $\tau_2$ compared to $\tau_1$. This applies greater strength to perturb object motion, facilitating feasible corrections for contacting joints. During the denoising stage, to eliminate diffusion models' bias that can suppress the guidance signal, we iteratively perturb $K$ times in the last denoising step, as illustrated in Algorithm 1.

How can we define the objective function $G(\boldsymbol{\mu}_t^h, \boldsymbol{\mu}_t^o, \boldsymbol{y}_0)$? We consider three terms here. First, in the generated 3D HOIs, the human and object should be close to each other on the contacting points. We therefore minimize the distance between human contact joints and object contact points

$$G_{con} = \sum_{i \in \{1,2,\ldots,8\}} \left\| R\big(\boldsymbol{\mu}_t^h(i)\big) - V\big(\boldsymbol{y}_t^o(i)\big) \right\|^2, \tag{5}$$

where $\boldsymbol{\mu}_t^h(i)$ and $\boldsymbol{y}_t^o(i)$ denote the $i$-th available contacting joint indexed by $\boldsymbol{y}_0^h$ and $i$-th object contact point, respectively. $R(\cdot)$ converts the human joint's local positions to global absolute locations, and $V(\cdot)$ obtains the object's contact point sequence from the predicted mean of object pose $\boldsymbol{\mu}_t^o$.

Second, the generated motion of dynamic objects typically follows human movement. However, we observe that when the human interacts with a static object, such as sitting on a chair, the object appears slightly moved. To address this, we immobilize the object's movement in the generated samples if the state is static ($\boldsymbol{y}_0^s = 1$), ensuring that proper contact is established between the human and the static object. The objective is defined as

$$G_{sta} = \boldsymbol{y}_0^s \cdot \sum_{l=1}^{L} \left\| \boldsymbol{\mu}_t^o(l) - \bar{\boldsymbol{\mu}}_t^o \right\|^2, \tag{6}$$

where $\boldsymbol{\mu}_t^o(l)$ denotes the object's 6DoF pose in the $l$-th frame. $\bar{\boldsymbol{\mu}}_t^o = \frac{1}{L} \sum_l \boldsymbol{\mu}_t^o(l)$, which is the average of predicted means of the object's pose.

Third, we define a smoothness term $G_{smo}(\mu)$ for the object motion to mitigate motion jittering during contact. Due to the space limit, we explain it in the supplementary material.

Finally, we combine all these goal functions to as the final objective

$$G = G_{cont} + \alpha G_{sta} + \beta G_{smo}, \tag{7}$$

where $\alpha$ and $\beta$ are weights for balance.

## 4    Experiments

### 4.1    Setup

**Dataset.** Since the data designed for studying text-driven 3D HOIs generation is severely scarce, we manually label interaction types, interacting subjects, and contact body parts on top of the BEHAVE dataset [3]. We then use ChatGPT to assist in generating three text descriptions for each HOI sequence, increasing the diversity of the data. Specifically, BEHAVE encompasses the interactions of 8 subjects with 20 different objects. It provides the human SMPL-H representation [27], the object mesh, as well as its 6DoF pose information in each HOI sequence. To ensure consistency in our approach, we follow the processing method used in HumanML3D [11] to extract representations for 22 body joints.

All the models are trained to generate $L = 196$ frames in our experiments. In the end, we have 1451 3D HOI sequences along with textual descriptions to train and evaluate our proposed approach. We follow the official train/test split on BEHAVE. We provide more details of the dataset, our annotation process, and annotated textual examples in the supplementary material.

In addition, we evaluate our approach on OMOMO dataset [25]. OMOMO focuses on full-body manipulation with hands. It consists of human-object interaction motion for 15 objects in daily life, with a total duration of approximately 10 hours. It provides text descriptions for each interaction motion. We utilize their object split strategy for both training and evaluation, ensuring the objects between the training and testing sets are different. Additionally, we preprocess human and object motion, similar to our way for the BEHAVE dataset. More details of the OMOMO dataset are in the supplementary material.

**Evaluation metrics.** We first assess different models for human motion generation using standard metrics as introduced by [11], namely *Fréchet Inception Distance (FID)*, *R-Precision*, and *Diversity*. *FID* quantifies the discrepancy between the distributions of actual and generated motions via a pretrained motion encoder. *R-Precision* gauges the relevance between generated motions and their corresponding text prompts. *Diversity* evaluates the range of variation in the generated motions. Additionally, we compute the *Foot Skating Ratio* to measure the proportion of frames exhibiting foot skid over a threshold (2.5 cm) during ground contact (foot height < 5 cm).

To evaluate the effectiveness of HOIs generation, we report the *Contact Distance* metric, which quantitatively measures the proximity between the ground-truth human contact joints and the object contact points. Ideally, we should develop similar metrics, *e.g.*, *FID*, to evaluate the *stochastic* HOI generation. However, due to the limited data available in BEHAVE [3], training a motion encoder would produce biased evaluation results. To mitigate this issue, we resort to user studies to quantify the effectiveness of different models. Details will be introduced later.

### 4.2   Comparisons with Existing Methods

**Baselines.** Our work introduces a novel 3D HOIs generation task not addressed by existing text-to-motion methods, which focus exclusively on human motion generation without accounting for human-object interactions. To compare with existing works, we mainly focus on evaluating human motion generation. We then design different variants of our models for comparing 3D HOIs generation. Specifically, we adopt the prominent text-to-motion methods MDM [45] and PriorMDM* [41] with the following settings. (a) MDM$^{finetuned}$: In this setup, we fine-tune the original MDM model [45] on the BEHAVE dataset [3]. (b) MDM*: This variant involves adapting the input and output layers' dimensions of the MDM model [45] to accommodate the input of 3D HOI sequences. This adjustment allows for the simultaneous learning of both human and object motions within a singular, integrated model. (c) PriorMDM* [41]: We adapt the Com-MDM architecture proposed in [41], originally designed for two-person motion

generation, to suit our needs for HOIs synthesis by modifying one of its two branches for object motion generation. (d) InterDiff [54]: While InterDiff is not designed for text-driven synthesis of 3D HOI, we added text conditioning to InterDiff as the baseline. More details in the supplementary material.

| | FID ↓ | R-precision ↑ (Top-3) | Diversity → | Contact Distance ↓ | Foot Skate Ratio ↓ | User Study ↑ (preferred) |
|---|---|---|---|---|---|---|
| Real | 0.096 | 0.324 | 7.244 | - | - | - |
| MDM$^{finetuned}$ [45] | 3.832 | 0.198 | 7.976 | - | - | - |
| MDM* [45] | 3.169 | 0.206 | 7.606 | 0.448 | 0.184 | 0.87% |
| PriorMDM* [41] | 4.352 | 0.215 | 8.186 | 0.416 | 0.270 | 16.53% |
| InterDiff [54] | 7.192 | 0.170 | 7.686 | 0.506 | 0.218 | 0.00% |
| Ours (Full) | **1.501** | **0.235** | **7.391** | **0.347** | **0.182** | **82.60%** |

Table 1: **Quantitative results on the BEHAVE [3] dataset.** We compare our method with baselines adapted from existing models. MDM$^{finetuned}$: fine-tune the original MDM model on the BEHAVE dataset. MDM*: adapting the input and output layers' dimensions of the MDM model [45] to accommodate both human and object motions. PriorMDM*: We adapt the ComMDM architecture proposed in [41]. The right arrow → means closer to real data is better.

| | FID ↓ | R-precision ↑ (Top-3) | Diversity → | Contact Distance ↓ | Foot Skate Ratio ↓ |
|---|---|---|---|---|---|
| Real | 1.539 | 0.193 | 4.833 | - | - |
| MDM* | 9.155 | 0.137 | 5.353 | 0.768 | 0.191 |
| PriorMDM | 11.192 | 0.128 | 7.404 | 0.523 | 0.344 |
| InterDiff | 10.567 | 0.110 | 6.493 | 0.906 | 0.239 |
| Ours | **8.661** | **0.168** | **4.780** | **0.326** | **0.141** |

Table 2: **Quantitative results on OMOMO dataset.** The right arrow → means closer to real data is better.

**Quantitative Results.** Tab. 1 reports the quantitative results on BEHAVE dataset [3]. Compared with the baseline methods, our full method achieves the best performance. Specifically, it achieves state-of-the-art results in both *FID*, *R-precision*, and *Diversity*, underscoring its ability to generate high-quality human motions in the context of coherently interacting with objects. The best *Contact Distance* also suggests that our approach can generate physically plausible HOIs, capturing the intricate interplay interactions between humans and objects.

We also show the results of our *User Study* in Tab. 1 to illustrate user preferences between our method and other baselines. In our user study, we sample 24 sequences from each method and 23 participants are asked to choose their most preferred generation results from these samples. This user study is designed to directly compare user preferences across different methods. The result indicates a
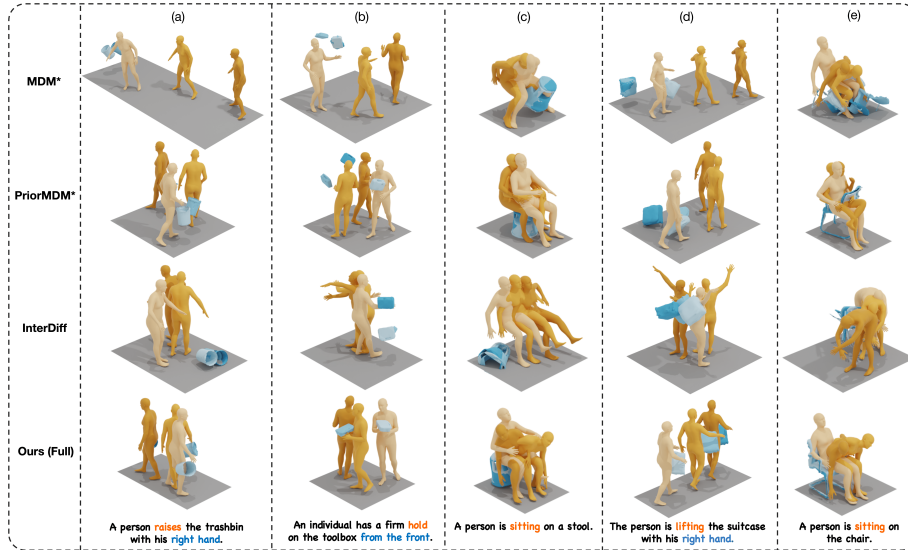
**Fig. 5: Qualitative comparisons of our approach and baseline methods on BEHAVE dataset.** The bottom row, showcasing our method, demonstrates the generation of realistic 3D HOIs with plausible contacts, particularly evident in columns 2 and 4. This contrasts with the baselines, which fail to achieve a similar level of realism and contact plausibility in the interactions. As an additional visual aid, the mesh color gradually darkens over time to represent progression. (Best viewed in color.)

strong preference for our method: it is favored over the baselines in a substantial majority of 82.60% of the cases.

Tab. 2 presents the quantitative results on the OMOMO dataset. Our method consistently outperforms other baselines by a considerable margin across all metrics. Notably, due to the distinctiveness of objects in the training and testing sets, the results in Tab. 2 indicate the effectiveness of our approach in *generalizing to unseen objects*, proving superior performance compared to other models.

**Qualitative Results.** We showcase qualitative comparisons, rendered with SMPL [27] shapes, between our approach and the baseline methods in Fig. 5. It is observed that the generated HOI motion by other baselines lacks smoothness and realism, where the object may float in the air (*e.g.*, the toolbox in Fig. 5 (b)). Furthermore, these baseline methods struggle to accurately capture the spatial relationships between humans and objects (*e.g.*, the chair in Fig. 5 (e)). In stark contrast, our approach excels in creating visually appealing and realistic HOIs. Notably, it adeptly reflects the intricate details outlined in text descriptions, capturing both the nature of the interactive actions and the specific body parts involved (*e.g.*, raising the trash bin with the right hand in Fig. 5 (a)).

| | FID $\downarrow$ | R-precision $\uparrow$ (Top-3) | Diversity $\rightarrow$ | Contact Distance $\downarrow$ | Foot Skate Ratio $\downarrow$ | User Study $\uparrow$ (preferred) |
|---|---|---|---|---|---|---|
| *w/o Interaction Correction* | | | | | | |
| Ours w/o CM | 2.825 | 0.201 | 7.850 | 0.524 | 0.265 | 0.00% |
| Ours w/o pretrain | 2.414 | 0.226 | 7.915 | 0.402 | **0.158** | 3.90% |
| Ours$^{global}$ | 11.057 | 0.187 | 7.545 | 0.375 | 0.274 | 0.00% |
| Ours | 1.808 | 0.217 | 7.718 | 0.416 | 0.205 | 11.69% |
| *w/ Interaction Correction* | | | | | | |
| Ours w/o $M^o$ & CM | 3.206 | 0.205 | 7.909 | 0.365 | 0.310 | 0.00% |
| Ours $^{Joint}$ | 3.412 | 0.195 | 7.878 | 0.421 | 0.342 | - |
| Ours w/o $G_{con}$ | 1.983 | 0.205 | 7.826 | 0.417 | 0.196 | - |
| Ours w/o $G_{sta}$ | 1.636 | 0.213 | 7.785 | 0.367 | **0.181** | - |
| Ours w/o $G_{smo}$ | 1.665 | 0.205 | 7.840 | 0.370 | 0.182 | - |
| Ours (Full) | **1.501** | **0.235** | **7.391** | **0.347** | 0.182 | **84.41%** |

**Table 3: Ablation studies of our model's variants on the BEHAVE dataset.** The right arrow $\rightarrow$ means closer to real data is better.
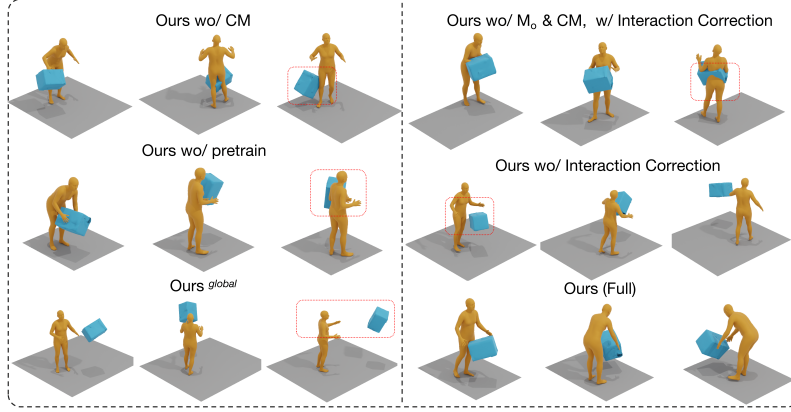


**Fig. 6: Visual results of different variants of our model in ablation studies on the BEHAVE dataset [3].**

## 4.3 Ablation Studies

We conduct extensive ablation studies in Tab. 3 and Fig. 6 to validate the effectiveness of different components . Tab. 3 categorizes the results into two parts depending on whether to use the affordance-guided interaction correction : *w/o Interaction Correction* and *w/ Interaction Correction*.

**w/o Interaction Correction.** In this part, we remove the affordance-guided interaction correction in the inference stage and compare each ablated variant with our method. When the Communication Module (CM) in HOI-DM is disabled, all the results drop substantially, especially in *Contact Distance*. "Ours w/o pretrain" demonstrates the effectiveness of fine-tuning the pretrained MDM [45] for human motion generation. Furthermore, as illustrated in "Ours$^{global}$", although we could also represent human joints in the same global space as the object 6DoF pose, we empirically find that the hybrid combination of the human joints

articulated in a local coordinate system and the object pose defined in a global space works better.

**w/ Interaction Correction.** A straightforward approach to generating 3D HOIs can be generating only human motion first and subsequently binding the object to the human body. To validate this design, we eliminate $M^o$ and $CM$ in HOI-DM from our training process, focusing solely on learning human motion. During inference, we apply the affordance-guided interaction correction to align the object motion with the human. However, as shown in the Tab. 3, it affects hu-

|            | AP (%) ↑ | L2 Dist ↓ |
|------------|----------|-----------|
| Ours $^{Joint}$ | 53.67 | 0.384 |
| Ours $^{APDM}$  | **78.54** | **0.272** |

**Table 4:** APDM evaluation. The reported metrics include Average Precision (AP) for predicted human contact probabilities and L2 Distance (Dist) error for predicted object contact points.

man motion performance, yielding unsatisfactory results. We also evaluate the effectiveness of the three analytic functions used in our interaction correction process. Notably, the function $G_{con}$, which incorporates the estimated affordance information to ensure close contact between humans and objects, proves to be crucial for generating realistic HOIs. The results clearly demonstrate the significance of affordance estimation in enhancing the authenticity of the interactions.

In addition to our primary quantitative evaluations, we also conduct *User Study* for ablation analysis. The results from this study indicate a clear user preference for our full method. This preference is observed consistently across various variants, underscoring the effectiveness of the complete set of design choices in our approach.

**Joint prediction of HOIs and affordance.** We attempt to generate HOIs and affordance jointly within the same diffusion model, without any specific design. However, as indicated in the $6^{th}$ row of Tab. 3 (Ours$^{Joint}$), the performance of the joint prediction method is notably inferior compared to our modular design.

**Performance of APDM.** To assess the efficacy of the proposed APDM, we present the results in the Tab. 4. APDM demonstrates superiority over the joint prediction method for both metrics. Additional visual results can be found in the supplementary material.

### 4.4   Limitations

The existing datasets for 3D HOIs are limited in terms of action and motion diversity, posing a challenge for synthesizing long-term interactions in our task. Furthermore, the effectiveness of our model's interaction correction component is contingent on the precision of affordance estimation. Despite simplifying this task, achieving accurate affordance estimation remains a significant challenge, impacting the overall performance of our model. A promising direction for future research involves integrating a sophisticated affordance model pre-trained on an

extensive 3D object dataset, along with text prompts. Such an advancement could significantly enhance the realism and accuracy of human-object contact in our model, leading to more natural and precise HOIs synthesis.

## 5    Conclusion

In summary, we presented a novel approach HOI-Diff to generate realistic 3D HOIs driven by textual prompts. By employing a modular design, we effectively decompose the complex task of HOI synthesis into simpler sub-tasks, enhancing the coherence and realism of the generated motions. Our HOI-Diff model successfully generates coarse dynamic human and object motions, while the affordance prediction diffusion model adds precision in predicting contact areas. The integration of estimated affordance data into classifier-guidance further ensures accurate human-object interactions. The promising experimental results on our annotated BEHAVE dataset demonstrate the efficacy of our approach in producing diverse and realistic HOIs.

# HOI-Diff: Text-Driven Synthesis of 3D Human-Object Interactions using Diffusion Models

## A    Additional Details of Methodology

In Sec. 3 of our main paper, we presented the foundational design of each key component in our HOI-Diff pipeline. Here, we delve into an elaborate explanation of model architecture, learning objectives and additional details associated with each crucial component.

**HOI diffusion model (HOI-DM).** The Communication Module (CM) in HOI-DM is based on cross attention mechanism. Formally,

$$\tilde{\boldsymbol{f}}^h = \mathrm{MLP}(\mathrm{Attn}(\boldsymbol{f}^h\mathbf{W}_Q, \boldsymbol{f}^o\mathbf{W}_K, \boldsymbol{f}^o\mathbf{W}_V)), \tag{1}$$

$$\tilde{\boldsymbol{f}}^o = \mathrm{MLP}(\mathrm{Attn}(\boldsymbol{f}^o\mathbf{W}_Q, \boldsymbol{f}^h\mathbf{W}_K, \boldsymbol{f}^h\mathbf{W}_V)), \tag{2}$$

where $\mathrm{MLP}(\cdot)$ denotes fully-connected layers, $\mathrm{Attn}(\cdot)$ is the attention block [47], and $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V$ are learned projection matrices for query, key, and value, respectively.

The training objective of this full model is based on reconstruction loss

$$\mathcal{L}_{hoi} = \mathbb{E}_{t\sim[1,T]}\|M_\theta(\boldsymbol{x}_t, t, \boldsymbol{c}) - \boldsymbol{x}_0\|_2^2, \tag{3}$$

where $\boldsymbol{x}_0$ is the ground truth of the HOI sequence.

**Affordance prediction diffusion model (APDM).** The affordance prediction diffusion model comprises eight Transformer layers for the encoder with a PointNet++ [35] to encode the object's point clouds. The training objective of this diffusion model is also based on reconstruction loss

$$\mathcal{L}_{aff} = \mathbb{E}_{t\sim[1,T]}\|A_\theta(\boldsymbol{y}_t, t, \boldsymbol{p}, \boldsymbol{d}) - \boldsymbol{y}_0\|_2^2, \tag{4}$$

where $\boldsymbol{y}_0$ is the ground-truth affordance data. $\boldsymbol{p}$ and $\boldsymbol{d}$ denote object point cloud and text description (prompt), respectively. $A_\theta$ represents the affordance prediction diffusion model.

**Affordance-guided interaction correction.** During the inference stage, it's found that the predicted object contact positions may occasionally be inaccurately positioned, residing either inside or outside the object. To rectify this, we implement post-processing steps that replace these predicted contact points, denoted as $\boldsymbol{y}_0^o$, with their nearest neighbors from the object's point clouds. This

adjustment aims to enhance the accuracy of the updated contact points, aligning them more closely with their actual positions on the object's surface. However, employing these updated contact points directly for contact constraints, particularly in the absence of detailed human shape information, introduces a new challenge. It can potentially lead to penetration issues within the contact area while reconstructing the human mesh in the final stage. To mitigate contact penetration, we adopt a method that recalculates points at a specified distance outward, perpendicular to the normal, originating from the object's contact points.

As for smoothness term, we formulate it as

$$G_{smo} = \sum_{l=1}^{L-1} \left\| \boldsymbol{x}_0^o(l+1) - \boldsymbol{x}_0^o(l) \right\|^2 , \tag{5}$$

where $\boldsymbol{x}_0^o(l)$ is the predicted 6DoF pose of the object in the $l$-th frame.

## B    Additional Details of Baselines

– $MDM^{finetuned}$: We finetune MDM [45] on BEHAVE dataset without considering the object motion.
– MDM*: We extend the original feature dimensions of the input and output processing in MDM [45] from $D^h$ to $D^h + D^o$, enabling support for HOIs sequences. The model is trained from scratch on BEHAVE dataset [3].
– PriorMDM*: The proposed approach for dual-person motion generation employs paired fixed MDMs [45] per individual to ensure uniformity within generated human motion distributions. This design leverages a singular Com-MDM to coordinate between the two branches of fixed MDM instances, streamlining training and maintaining consistency across generated motions. Given that both branches are based on MDM that pretrained on human motion datasets, direct utilization of them for human-object interactions in our task is infeasible. We maintain one branch dedicated to humans, leveraging pre-trained weights, while adapting the input and output processing of another branch specifically for generating object motion. Following this, we fine-tune the human MDM branch while initiating the learning of object motion from scratch within the object branch. Eventually, we integrate Com-MDM to facilitate communication and coordination between these distinct branches handling human and object interactions.
– InterDiff: InterDiff [54] is originally designed for a prediction task rather than text-driven HOIs generation. To tailor it to our task, we replace its Transformer encoder with a CLIP encoder and modify its feature dimensions of the input and output layers.

To ensure fair comparisons, all the above baselines as well as our own models are all trained on BEHAVE for 20k steps.
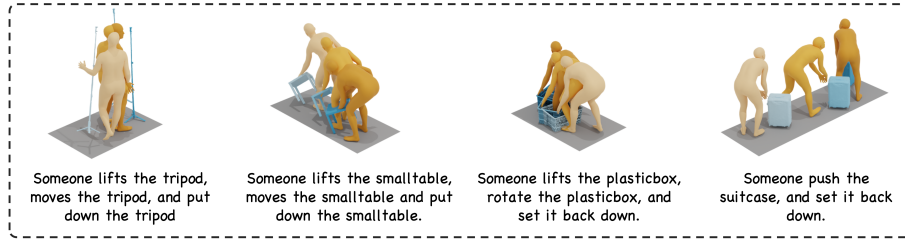
**Fig. 1:** Additional qualitative evaluation on OMOMO dataset. Given object geometry and text description, our method can generate high-quality human-object interactions even for the unseen objects.

## C   Additional Details of Evaluation Metrics

For detailed information regarding metrics employed in human motion generation, including *FID*, *R-Precision*, and *Diversity*, we refer readers to [11, 45] for comprehensive understanding.

Expanding on the concept of *Contact Distance*, we utilize the *chamfer distance* metric to quantify the closeness between human body joints and the object surface. This computation leverages ground-truth affordance data that includes human contact labels and object contact points,

$$ContactDistance = \frac{1}{L} \sum_{l}^{L} CD(\hat{\boldsymbol{x}}_l^h, \hat{\boldsymbol{p}}_l), \tag{6}$$

where $\hat{\boldsymbol{x}}_l^h$ represents two human contact joints at the $l$-th frame, indexed according to ground-truth contact labels. Additionally, $\hat{\boldsymbol{p}}_l$ denotes two object contact points derived from the object motion $\boldsymbol{x}_l^o$ at frame $l$, also indexed based on ground-truth information. $CD$ denotes the *chamfer distance*.

## D   Implementation Details

Both our HOI-MDM and APDM are built on the Transformer [47] architecture. Similar to MDM [45], we employ the CLIP model to encode text prompts, adhering to a classifier-free generation process. Our models are trained using PyTorch [32] on 1 NVIDIA A5000 GPU. We set control strength of guidance as $\tau_1 = 1$, $\tau_2 = 100$, and $\Sigma_t = min(\Sigma_t, 0.01)$.

**Model architecture.** Both the HOI-DM and APDM architectures of HOI-Diff are based on Transformers with 4 attention heads, a latent dimension of 512, a dropout of 0.1, a feed-forward size of 1024, and the GeLU activation [13]. The number of learned parameters for each model is stated in Tab. 1.

**Training hyperparameters.** Our training setting involves 20k iterations for the HOI-MDM and 10k iterations for the APDM model. These iterations utilize a batch size of 32 and employ the AdamW optimizer [28] with a learning rate set
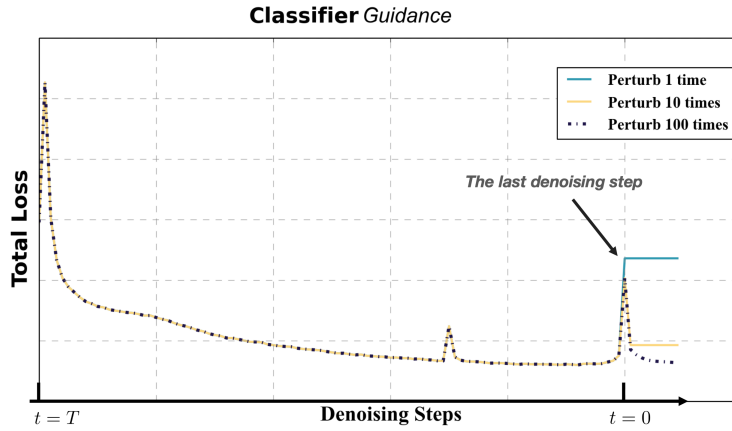
**Fig. 2: Effect of different number of perturbations in the final denoising step.**

| Model | HOI-DM | APDM |
|---|---|---|
| Parameters ($\cdot 10^6$) | 8.82 | 38.92 |

**Table 1: Model Parameters.** The number of learned parameters of our two core architectures.

| Method | MDM* | PriorMDM* | Ours (Full) |
|---|---|---|---|
| Time (s) | 32.3 | 38.6 | 118.0 |
| Component | APDM | HOI-DM | Interaction Correction |
| Time (s) | 24.2 | 46.4 | 47.4 |

**Table 2: Inference Time.** We report the inference time for baselines, our full method, and its key components.

at $10^{-4}$. We use $T$=1000 and $N$=500 diffusion steps in HOI-DM and APDM, respectively.

## E    Inference Time

In Tab. 2, we provide the inference times for both baselines and our full method, including its key components. All measurements were conducted using an NVIDIA A5000 GPU. Training an additional model for affordance information and using classifier guidance for interaction correction do contribute to increased inference costs. However, despite the longer inference time, our complete method notably enhances the accuracy of 3D HOIs generation.

## F    Additional Qualitative Results

In this section, we present additional qualitative results showcasing the model's performance evaluated on the OMOMO dataset, and the effectiveness of APDM.

**Qualitative results on OMOMO dataset.** We present additional qualitative results on the OMOMO dataset, rendered with SMPL [27] shapes. It is evident
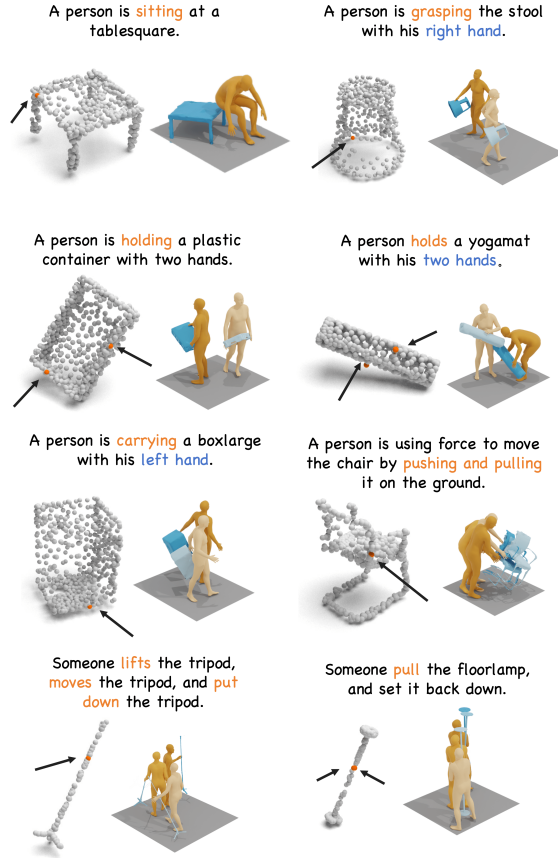
**Fig. 3:** Visual results of estimated contact points. Our APDM, trained on the BEHAVE dataset, can accurately estimating contact positions for objects based on textual descriptions. Furthermore, it showcases the capability to generalize to unseen objects in the OMOMO dataset, as demonstrated in the last row.

that our method can generalizes effectively to unseen objects and produce realistic 3D human-object interactions.

**Qualitative results of APDM.** To verify the accuracy of estimated contact points on object surface, we provide additional visual results in Fig. 3. It can be seen that our method can predict realistic and practical contact points based on text descriptions.

## G  Additional Ablation Studies

**Different perturbing times in classifier guidance.** As discussed in Sec. 3.4, in the later stage of classifier guidance, diffusion models tend to strongly attenu-
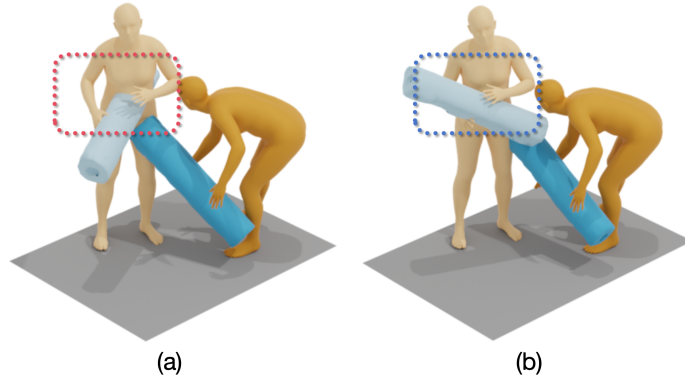
**Fig. 4: Effect of different control strengths for classifier guidance.** (a) We use equal strengths of $\tau_1 = 1, \tau_2 = 1$ to perturb the predicted mean of human motion and object motion, respectively. (b) We use different strengths of $\tau_1 = 1, \tau_2 = 100$ for the perturbation. We can see that different strengths work better.

ate the introduced signals. Therefore, we iteratively perturb the predicted mean of motion for $K$ times at the final denoising step. In Fig. 2, we present the ablation results, illustrating the impact of different numbers of perturbations. Notably, we observe that employing 100 perturbations leads to re-convergence and yields the desired results.

**Different guidance strength.** As detailed in Sec. 3.4, we employ distinct control strengths for classifier guidance, considering the varying feature densities in predicted human and object motion. Rather than employing equal control strengths, we opt to assign a higher control strength to object motion, allowing it to closely align with human contact joints, as illustrated in Fig. 4.

## H    Annotation for BEHAVE Dataset

**Text Annotating Process.** Initially, we manually annotate the interaction types and the specific human body parts involved, delineating actions like "lift" associated with the "left hand" or "hold" involving "two hands". Subsequently, to generate complete sentences, we leverage the capabilities of GPT-3.5 to assist in formulating the entirety of the description.

**Examples of Annotated Textual Descriptions.** In Tab. 3, we showcase a selection of our annotated textual descriptions for the BEHAVE dataset [3].

**Affordance Data.** Our affordance data includes 8-dimensional human contact labels and object contact points. We employ *chamfer distance* to measure the distance between all human body joints and object surface points. Following a predefined distance threshold $\gamma = 0.12$, we identify the 8 contact points on the object surface corresponding to the 8 primary human body joints. Subsequently,

| Object | Textual Descriptions |
|--------|----------------------|
| *backpack* | A person is carrying the backpack in front. |
| | The person is raising a backpack with his right hand. |
| | The person at the front presently has control over the backpack. |
| *chairwood* (*wooden chair*) | A person is using the chairwood for sitting. |
| | The person is propelling the chairwood on the ground. |
| | Someone is hoisting a chairwood by his left hand. |
| *tablesquare* (*square table*) | A person is lifting the tablesquare, utilizing his left hand. |
| | Someone is clutching onto a tablesquare from the front. |
| | An individual is moving the tablesquare back and forth. |
| *boxlong* (*long box*) | A person is gripping the boxlong from the front. |
| | A person is raising the boxlong using his left hand. |
| | Someone hoists the boxlong with his left hand. |
| *toolbox* | Someone is grasping the toolbox upfront. |
| | The person has a firm hold on the toolbox with his right hand. |
| | A person is gripping the toolbox with his left hand. |
| *yogaball* | A person is shifting a yogaball back and forth on the floor using his hands. |
| | The person is occupying a yogaball. |
| | A person is employing an yogaball to engage in an upper body game. |

**Table 3: Examples of our annotated textual descriptions for the BEHAVE dataset generated by GPT-3.5.**

we derive the human contact labels by encoding the indexes of contact joints into an 8-dimensional vector represented by binary values.

# I  Additional Details of OMOMO Dataset

The OMOMO dataset comprises data captured for a total of 15 objects. Adhering to their official split strategy depicted in [25](Figure 5), we allocate 10 objects for training and 5 objects for testing. This split allows us to further evaluate the model's generalization ability to new objects. Notably, the OMOMO dataset itself provides text annotation, and we use GPT-3.5 to add subjects to it and embellish it appropriately. For affordance data, we preprocess it the same way we handle BEHAVE.

## J   Supplementary Video

Beyond the qualitative results presented in the main paper, our supplementary materials offer comprehensive demos that provide an in-depth visualization of our task, further showcasing the effectiveness of our approach.

In these demonstrations, we highlight the better performance of our method, HOI-Diff, in producing diverse and realistic 3D HOIs while maintaining adherence to physical validity. Notably, the visualizations show that HOI-Diff consistently generates smooth, vivid interactions, accurately capturing human-object contacts.

Additionally, we present the visual ablation results and emphasize the significance and effectiveness of our affordance-guided interaction correction, underscoring its substantial impact on improving the overall performance and quality of the generated 3D HOIs.

## References

1. Ahn, H., Mascaro, E.V., Lee, D.: Can we use diffusion probabilistic models for 3d motion prediction? arXiv (2023) 3
2. Barquero, G., Escalera, S., Palmero, C.: Belfusion: Latent diffusion for behavior-driven human motion prediction. In: ICCV (2023) 3
3. Bhatnagar, B.L., Xie, X., Petrov, I., Sminchisescu, C., Theobalt, C., Pons-Moll, G.: Behave: Dataset and method for tracking human object interactions. In: CVPR (2022) 2, 3, 9, 10, 11, 13, 6
4. Chen, L.H., Zhang, J., Li, Y., Pang, Y., Xia, X., Liu, T.: Humanmac: Masked motion completion for human motion prediction. arXiv (2023) 3
5. Chen, X., Jiang, B., Liu, W., Huang, Z., Fu, B., Chen, T., Yu, G.: Executing your commands via motion diffusion in latent space. In: CVPR (2023) 3
6. Dabral, R., Mughal, M.H., Golyanik, V., Theobalt, C.: Mofusion: A framework for denoising-diffusion-based motion synthesis. In: CVPR (2023) 3
7. Deng, S., Xu, X., Wu, C., Chen, K., Jia, K.: 3d affordancenet: A benchmark for visual object affordance understanding. In: CVPR (2021) 4, 7
8. Dhariwal, P., Nichol, A.: Diffusion models beat gans on image synthesis. In: NeurIPS (2021) 8
9. Diller, C., Dai, A.: Cg-hoi: Contact-guided 3d human-object interaction generation (2024) 3, 4
10. Ghosh, A., Dabral, R., Golyanik, V., Theobalt, C., Slusallek, P.: Imos: Intent-driven full-body motion synthesis for human-object interactions. In: CGF (2023) 2, 3
11. Guo, C., Zou, S., Zuo, X., Wang, S., Ji, W., Li, X., Cheng, L.: Generating diverse and natural 3d human motions from text. In: CVPR (2022) 2, 4, 9, 10, 3
12. Hassan, M., Ceylan, D., Villegas, R., Saito, J., Yang, J., Zhou, Y., Black, M.: Stochastic scene-aware motion prediction. In: ICCV (2021) 2, 3
13. Hendrycks, D., Gimpel, K.: Gaussian error linear units (gelus). arXiv (2016) 3
14. Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models (2020) 5
15. Huang, S., Wang, Z., Li, P., Jia, B., Liu, T., Zhu, Y., Liang, W., Zhu, S.C.: Diffusion-based generation, optimization, and planning in 3d scenes. In: CVPR (2023) 3

16. Iriondo, A., Lazkano, E., Ansuategi, A.: Affordance-based grasping point detection using graph convolutional networks for industrial bin-picking applications. Sensors (2021) 4, 7
17. Jiang, N., Liu, T., Cao, Z., Cui, J., Chen, Y., Wang, H., Zhu, Y., Huang, S.: Chairs: Towards full-body articulated human-object interaction. arXiv (2022) 2, 3
18. Karunratanakul, K., Preechakul, K., Suwajanakorn, S., Tang, S.: Gmd: Controllable human motion synthesis via guided diffusion models. In: ICCV (2023) 3
19. Karunratanakul, K., Preechakul, K., Suwajanakorn, S., Tang, S.: Guided motion diffusion for controllable human motion synthesis. In: ICCV (2023) 8
20. Kim, D.I., Sukhatme, G.S.: Semantic labeling of 3d point clouds with object affordance for robot manipulation. In: ICRA (2014) 4, 7
21. Kim, D.I., Sukhatme, G.S.: Interactive affordance map building for a robotic task. In: IROS (2015) 4, 7
22. Kokic, M., Stork, J.A., Haustein, J.A., Kragic, D.: Affordance detection for task-specific grasping using deep learning. In: International Conference on Humanoid Robotics (Humanoids) (2017) 4, 7
23. Kulkarni, N., Rempe, D., Genova, K., Kundu, A., Johnson, J., Fouhey, D., Guibas, L.: Nifty: Neural object interaction fields for guided human motion synthesis. arXiv (2023) 2, 3, 8
24. Li, J., Clegg, A., Mottaghi, R., Wu, J., Puig, X., Liu, C.K.: Controllable human-object interaction synthesis (2023) 3, 4
25. Li, J., Wu, J., Liu, C.K.: Object motion guided human motion synthesis. TOG (2023) 2, 3, 10, 7
26. Liang, H., Zhang, W., Li, W., Yu, J., Xu, L.: Intergen: Diffusion-based multi-human motion generation under complex interactions. arXiv (2023) 4
27. Loper, M., Mahmood, N., Romero, J., Pons-Moll, G., Black, M.J.: SMPL: A skinned multi-person linear model. ACM Trans. Graphics (Proc. SIGGRAPH Asia) (2015) 9, 12, 4
28. Loshchilov, I., Hutter, F.: Decoupled weight decay regularization. In: ICLR (2017) 3
29. Mo, K., Qin, Y., Xiang, F., Su, H., Guibas, L.: O2o-afford: Annotation-free large-scale object-object affordance learning. In: CoRL (2022) 4, 7
30. Ngyen, T., Vu, M.N., Vuong, A., Nguyen, D., Vo, T., Le, N., Nguyen, A.: Open-vocabulary affordance detection in 3d point clouds. arXiv (2023) 4, 7
31. Nichol, A.Q., Dhariwal, P.: Improved denoising diffusion probabilistic models. In: ICML (2021) 5
32. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al.: Pytorch: An imperative style, high-performance deep learning library. NeurIPS (2019) 3
33. Pi, H., Peng, S., Yang, M., Zhou, X., Bao, H.: Hierarchical generation of human-object interactions with diffusion probabilistic models. In: ICCV (2023) 2, 3
34. Plappert, M., Mandery, C., Asfour, T.: The kit motion-language dataset. Big Data (2016) 2
35. Qi, C.R., Yi, L., Su, H., Guibas, L.J.: Pointnet++: Deep hierarchical feature learning on point sets in a metric space. NeurIPS (2017) 1
36. Ramesh, A., Pavlov, M., Goh, G., Gray, S., Voss, C., Radford, A., Chen, M., Sutskever, I.: Zero-shot text-to-image generation. In: ICML (2021) 3
37. Razali, H., Demiris, Y.: Action-conditioned generation of bimanual object manipulation sequences. In: AAAI (2023) 2

38. Rempe, D., Luo, Z., Peng, X.B., Yuan, Y., Kitani, K., Kreis, K., Fidler, S., Litany, O.: Trace and pace: Controllable pedestrian animation via guided trajectory diffusion. In: CVPR (2023) 3

39. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: CVPR (2022) 3

40. Saharia, C., Chan, W., Saxena, S., Li, L., Whang, J., Denton, E.L., Ghasemipour, K., Gontijo Lopes, R., Karagol Ayan, B., Salimans, T., et al.: Photorealistic text-to-image diffusion models with deep language understanding. In: NeurIPS (2022) 3

41. Shafir, Y., Tevet, G., Kapon, R., Bermano, A.H.: Human motion diffusion as a generative prior. arXiv (2023) 3, 4, 10, 11

42. Starke, S., Zhang, H., Komura, T., Saito, J.: Neural state machine for character-scene interactions. TOG (2019) 2, 3

43. Sun, J., Chowdhary, G.: Towards globally consistent stochastic human motion prediction via motion diffusion. arXiv (2023) 3

44. Taheri, O., Choutas, V., Black, M.J., Tzionas, D.: GOAL: Generating 4D whole-body motion for hand-object grasping. In: CVPR (2022) 2

45. Tevet, G., Raab, S., Gordon, B., Shafir, Y., Cohen-or, D., Bermano, A.H.: Human motion diffusion model. In: ICLR (2023) 2, 3, 5, 6, 7, 10, 11, 13

46. Tian, S., Zheng, M., Liang, X.: Transfusion: A practical and effective transformer-based diffusion model for 3d human motion prediction. arXiv (2023) 3

47. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. NeurIPS (2017) 6, 1, 3

48. Wang, J., Rong, Y., Liu, J., Yan, S., Lin, D., Dai, B.: Towards diverse and natural scene-aware 3d human motion synthesis. In: CVPR (2022) 3

49. Wang, Y., Lin, J., Zeng, A., Luo, Z., Zhang, J., Zhang, L.: Physhoi: Physics-based imitation of dynamic human-object interaction. arXiv preprint arXiv:2312.04393 (2023) 3, 4

50. Wang, Z., Chen, Y., Liu, T., Zhu, Y., Liang, W., Huang, S.: Humanise: Language-conditioned human motion generation in 3d scenes. NeurIPS (2022) 3

51. Wei, D., Sun, X., Sun, H., Li, B., Hu, S., Li, W., Lu, J.: Understanding text-driven motion synthesis with keyframe collaboration via diffusion models. arXiv (2023) 3

52. Wu, Y., Wang, J., Zhang, Y., Zhang, S., Hilliges, O., Yu, F., Tang, S.: Saga: Stochastic whole-body grasping with contact. In: ECCV (2022) 2

53. Xie, Y., Jampani, V., Zhong, L., Sun, D., Jiang, H.: Omnicontrol: Control any joint at any time for human motion generation. arXiv (2023) 3, 8

54. Xu, S., Li, Z., Wang, Y.X., Gui, L.Y.: Interdiff: Generating 3d human-object interactions with physics-informed diffusion. In: ICCV (2023) 2, 4, 7, 11

55. Zhang, J., Zhang, Y., Cun, X., Huang, S., Zhang, Y., Zhao, H., Lu, H., Shen, X.: T2m-gpt: Generating human motion from textual descriptions with discrete representations. In: CVPR (2023) 3

56. Zhang, M., Cai, Z., Pan, L., Hong, F., Guo, X., Yang, L., Liu, Z.: Motiondiffuse: Text-driven human motion generation with diffusion model. arXiv (2022) 3

57. Zhang, M., Guo, X., Pan, L., Cai, Z., Hong, F., Li, H., Yang, L., Liu, Z.: Remodiffuse: Retrieval-augmented motion diffusion model. arXiv (2023) 3

58. Zhang, X., Bhatnagar, B.L., Starke, S., Guzov, V., Pons-Moll, G.: Couch: Towards controllable human-chair interactions. In: ECCV (2022) 2, 3

59. Zhang, Z., Liu, R., Aberman, K., Hanocka, R.: Tedi: Temporally-entangled diffusion for long-term motion synthesis. arXiv (2023) 3

60. Zhao, K., Zhang, Y., Wang, S., Beeler, T., Tang, S.: Synthesizing diverse human motions in 3d indoor scenes. arXiv (2023) 3