

THOR: Text to Human-Object Interaction Diffusion via Relation Intervention

Qianyang Wu¹, Ye Shi¹, Xiaoshui Huang², Jingyi Yu¹, Lan Xu¹, and Jingya Wang¹

¹ ShanghaiTech University
² Shanghai AI Laboratory

Abstract. This paper addresses new methodologies to deal with the challenging task of generating dynamic Human-Object Interactions from textual descriptions (Text2HOI). While most existing works assume interactions with limited body parts or static objects, our task involves addressing the variation in human motion, the diversity of object shapes, and the semantic vagueness of object motion simultaneously. To tackle this, we propose a novel Text-guided Human-Object Interaction diffusion model with Relation Intervention (THOR). THOR is a cohesive diffusion model equipped with a relation intervention mechanism. In each diffusion step, we initiate text-guided human and object motion and then leverage human-object relations to intervene in object motion. This intervention enhances the spatial-temporal relations between humans and objects, with human-centric interaction representation providing additional guidance for synthesizing consistent motion from text. To achieve more reasonable and realistic results, interaction losses is introduced at different levels of motion granularity. Moreover, we construct Text-BEHAVE, a Text2HOI dataset that seamlessly integrates textual descriptions with the currently largest publicly available 3D HOI dataset. Both quantitative and qualitative experiments demonstrate the effectiveness of our proposed model.

1 Introduction

The synthesis of human-object interactions (HOI) is pivotal for various applications, including VR/AR, embodied AI, computer animation, and robotics. Just like humans possess a strong imaginative capability, allowing us to envision interactions with a given object. This imaginative capacity proves beneficial for subsequent planning and execution abilities. Recent advancements in 3D human-object motion capture and generation models have made HOI synthesis increasingly achievable. Given that text serves as a natural and interactive modality for controlling the generation, the imperative need for user-friendly text guided HOI generation method is underscored.

Recently, there have been notable efforts in generating 3D human motions from textual descriptions, showcasing a plausible transition from text to motion [11, 58, 73, 77]. Generating 3D Human-Object Interaction from text (Text2HOI)

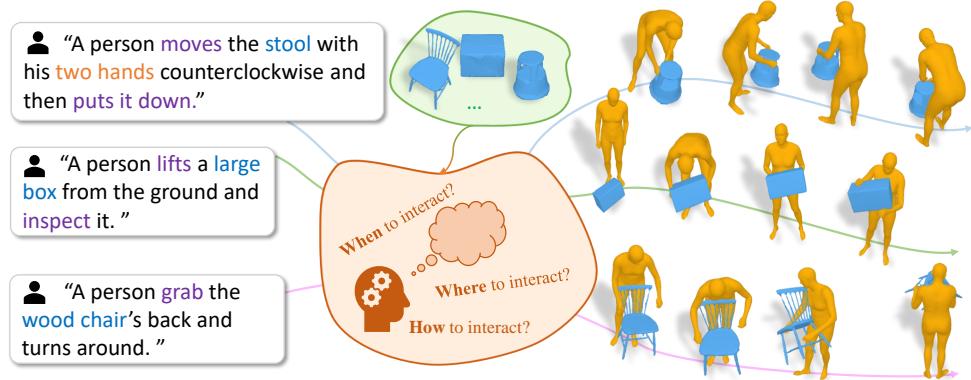


Fig. 1: A novel task of generating 3D Human-Object Interaction based on a Text prompt (Text2HOI), reflecting When, Where, and How humans interact with the object.

poses a significant challenge due to the inherent ambiguity in object motion, particularly when considering textual conditions. The ambiguity in object motion introduces complexities, making it challenging to precisely understand and combine the nuanced interactions between humans and objects in three-dimensional space.

Prior works on generating human-object interactions are either limited with static objects [5, 19, 27, 45, 79], or only generate the human motion of upper body [14, 54, 61]. Generating whole-body interactions with dynamic objects is extremely challenging due to intricate spatial dynamics, the need for accurate trajectory prediction of dynamic objects, and the difficulty in coordinating timing and synchronization.

Taking advantage of the recent advancements in generative models with diffusion [21, 48, 58], a straightforward strategy for Text2HOI involves extending existing diffusion models to generate human-object interactions. However, these models encounter challenges in accurately capturing the precise relations between human and object motions. The difficulty arises when explicit contact regions from historical motions [68] or predefined criteria [31] are absent, potentially leading to inconsistencies between the object motion and human motion directly generated from text. This discrepancy may stem from the inherent relative structure of joints and their parent joints with fixed bone length, constraining the motion space for humans. Since object dynamics require understanding not only where the interaction occurs but also when and how to interact, textual descriptions fall short in providing such precise contextual information, contributing to the vagueness of the object motion.

To tackle these challenges, we introduce THOR (**T**ext-conditioned **H**uman-**O**bject interaction diffusion with **R**elation intervention). THOR is a cohesive diffusion model that integrates interaction and intervention mechanisms in a single end-to-end framework. As objects play a passive role in interactions, they

exhibit a specific spatial distribution determined by their spatial relations with the active human participants. Drawing inspiration from this, we model human-centric relations to refine the object motion. Specifically, we design the Text2HOI Diffusion with a specially designed Human object relation intervention mechanism that models the human object kinematic relations from the rotation and translation perspectives, respectively. At inference, the intervention still remains in each denoising step.

To further improve the generation, supervision is introduced on intervened motions, kinematic relations, and geometric distance in a hierarchical manner. To enhance the precise awareness of interactions between human and object, the objective of our model is to denoising not only their motion, but also reasonable kinematic relations between them. While kinematic relations reflect the relative motion pattern between human and object, geometric distance gives point-wise relation for human object interaction. Considering the variation of object shape, it provides fine-grained description of their interactions. These two additional supervision on kinematic relation and geometric distance encourage the generation model to further capture the interaction pattern implicitly.

In summary, our contributions are as follows:

- We propose THOR, a diffusion model specifically tailored for Text2HOI that integrates human-object interactions and intervention mechanisms in a single end-to-end framework. Notably, the intervention mechanism is designed to refine implausible interactions sampled from textual prompts by leveraging the human-object kinematic relations to intervene in the object motion.
- We introduce the supervision on human and object kinematic relation and geometric distance to capture multi-level interactions. Through two special objective functions, our model embeds the human object relations into the diffusion process, facilitating to generate diverse and plausible human object interactions.
- We construct a Text2HOI dataset, named Text-BEHAVE that integrates textual descriptions to the currently largest publicly 3D HOI dataset. Both quantitative and qualitative evaluations demonstrate the capability of our model to handle the complexity of this task and produce meaningful and coherent interactions from textual prompts.

2 Related Work

Text-to-Human Motion Generation The realm of human motion synthesis has witnessed extensive exploration, with recent endeavors incorporating auxiliary modalities such as audio [1, 28, 32, 83] and action [42, 43, 67] to enhance generation performance and semantic richness. Text-conditioned human motion generation has become a focal point in recent research, with approaches leveraging various techniques including Variational Autoencoders (VAEs) [2, 9, 16, 44], Vector Quantized VAEs (VQ-VAEs) [17, 37, 75, 81], and Diffusion models [11, 26, 50, 58, 73, 77]. As a powerful pre-trained text-to-image model, CLIP [47]

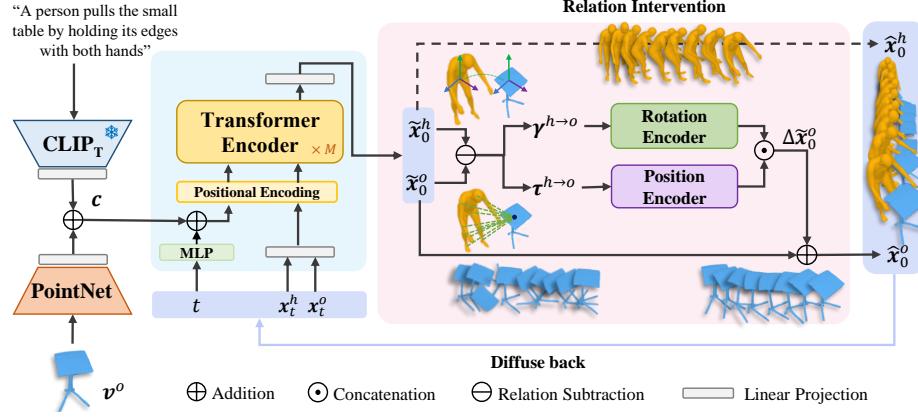


Fig. 2: The overview of our model THOR, designed to address a novel task of generating human-object interactions from textual descriptions (Text2HOI). The key innovation lies in leveraging the human-object spatial relations to further intervene the object motion and then diffuse back to benefit the whole Text2HOI diffusion framework.

has been widely adopted for encoding textual descriptions in numerous models [11, 12, 57, 77]. Beyond single-motion generation, efforts have been directed towards long-term human motion synthesis [2, 3, 8, 15, 69] and fine-grained motion control on trajectories or joints [52, 59, 65]. Multi-human interactions under text guidance have also been explored, with models like [33] employing cooperative denoising and [56] modeling asymmetric human interactions based on different roles. However, human-only motion usually lack the contextual information as in the real world, this paper we will tackle with the similar text guided generation but involving interactions with objects.

Human-Object Interaction In addition to the synthesis of human-only motion, several studies dive into the intricate interactions between humans and objects. Reconstruct human and object interactions from multi-view cameras [6, 23, 76] and monocular images [24, 63, 64, 74] has been both investigated. And the generation of human-object interactions has emerged as a prominent focus in recent research. Several approaches leverage human motion priors to generate object positions based on human trajectories and poses [18, 41, 70–72], providing a human-centric perspective on human-object interactions. In contrast, other works attempts to generate natural human motions in 3D indoor scenes [22, 80] or interactions with seating furniture [10, 25, 27, 45, 78]. These endeavors are further extended by introducing semantic descriptions [19, 34, 60, 62, 79].

Besides, generating interactions with dynamic objects have been also been investigated, particularly for upper limb interactions with prehensile objects [14, 49, 54, 55, 61, 82], or specific skills involving limited object categories [4, 20, 35, 38, 53, 66]. Recently, InterDiff [68] employs a diffusion model to forecast future interactions from historical observations and OMOMO [31] builds a two-stage

diffusion model to synthesising human motion from moving objects. Also, there are some concurrent works in the mean time. Overall, most existing approaches either focus on static scenes and limited interactions, or lack the text descriptions which is an import interface for generation, leaving a gap in addressing dynamic human-object interactions and the crucial role of textual guidance. Concurrent works for Text2HOI like CG-HOI [13] and CHOIS [29], explicitly model human object contact in the diffusion process. HOI-Diff [40] generates HOI in a dual branch diffusion model with affordance guided correction. Without explicit modeling of contact points or affordance, our model interprets interactions based on the kinematic relationship and geometric distance between humans and object surfaces.

3 Method

From a textual description, our goal is to generate plausible human-object interaction. Different from text-to-human motion generation, the object’s motion is intricately linked to its shape. Therefore the object shape is also taken as the condition for generation. This problem equals to that, given a text prompt and a 3D object model, generate human-object interactions that satisfy the textual description. This paper combines these two conditional signals and generate human-object interactions in a single diffusion framework.

We first encode two conditions and generate a primitive result. And this imperfect interaction provides relatively better human motion, which becomes additional guidance to facilitating network to learn object motion. Sec. 3.1 will describe the whole diffusion framework that how to encode two conditions and produce primitive interactions. Sec. 3.2 explains the details of the novel intervention mechanism and multi-level supervision on interactions, which benefits the whole Text2HOI diffusion framework to refine the primitive generation results by intervening the object motion based on the constructed human-object kinematic relations.

3.1 Text2HOI Diffusion Framework

Motion representation Human and object motion are expected to be in consistent representation. We keep the representation in BEHAVE [6], where human motion follows the representation in SMPL models [36, 39, 51] and object motion is a series of 6D poses along time dimension. Therefore, human and object are treated as two instances that composed of translations and rotations. In specific, human motion and object motion are denoted as

$$\mathbf{x}^h = [\mathbf{q}, \mathbf{j}], \quad \mathbf{x}^o = [\mathbf{r}, \mathbf{o}], \quad (1)$$

where $\mathbf{q} \in \mathbb{R}^{N \times J \times D_{rot}}$, $\mathbf{j} \in \mathbb{R}^{N \times J \times 3}$ refer to joint rotations and joint positions of the human, $\mathbf{r} \in \mathbb{R}^{N \times D_{rot}}$ and $\mathbf{o} \in \mathbb{R}^{N \times 3}$ are rotations and center translations of object. Here, N is the number of frames, J is the number of human joints, and

D_{rot} refer to the dimensions of rotation. The complete motion is the combination as $\mathbf{x} \triangleq [\mathbf{x}^h, \mathbf{x}^o]$. Similar to [16, 30], human and object motion are canonicalized with the same root orientation and root position in the first frame.

Conditional diffusion Unlike human-only text-to-motion generation, Text2HOI needs to encode both the text prompts and object shape. CLIP [47] is a powerful pre-trained model that has been verified in many motion generation models [11, 33, 58], thus we encode text prompts with the frozen text encoder of CLIP. Similarly, the object shape, noted as \mathbf{v}_o is embedded by a lightweight shape encoder PointNet [46]. After respective linear projection, the output text embedding \mathbf{c}_{text} and shape embedding \mathbf{c}_{shape} are added together to obtain the final condition \mathbf{c} .

The complete conditional HOI diffusion framework comprises forward process and reverse process [21]. The forward process is articulated as T Markov steps, producing noisy interactions \mathbf{x}_t from real interactions \mathbf{x}_0 :

$$q(\mathbf{x}_{1:T}|\mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1}), \quad (2)$$

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_{t-1}}\mathbf{x}_{t-1}, (1 - \alpha_{t-1})\mathbf{I}), \quad (3)$$

where $t = 1, 2, \dots, T$, $\alpha_t \in (0, 1)$ and $\bar{\alpha}_t = \prod_{i=0}^t \alpha_i$.

The reverse process aims to reconstruct realistic interactions \mathbf{x}_0 from noisy \mathbf{x}_T and condition \mathbf{c} composed of text and shape. Following [30, 58], we predict denoised interactions \mathbf{x}_0 with network \mathcal{G}_θ in each step and noise it back to \mathbf{x}_{t-1} iteratively following the Markov chain. Denoting $\hat{\mathbf{x}}_0 = \mathcal{G}_\theta(\mathbf{x}_t, t, \mathbf{c})$, the reverse process can be formulated as

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta, \sigma^2 \mathbf{I}), \quad (4)$$

$$\boldsymbol{\mu}_\theta = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\mathbf{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\hat{\mathbf{x}}_0}{1 - \bar{\alpha}_t}, \quad (5)$$

where σ^2 is a fixed variance.

Classifier-free guidance Previous diffusion models [11, 33, 48, 58] achieve classifier-free guidance by the linear combination of unconditional generation and conditional generation. However, in the context of Text2HOI generation, object shape has to serve as an essential condition modality. This is mainly because, even within the same class, objects with different shapes possess distinguishable motion spaces. This distinction may arise from variations in size or geometric structure, impacting how objects move or interact. Consequently, the classifier-free guidance is reformulated as:

$$\mathcal{G}_\theta(\mathbf{x}_t, \mathbf{c}, t) = (1 - s)\mathcal{G}_\theta(\mathbf{x}_t, \mathbf{c}_{shape}, t) + s\mathcal{G}_\theta(\mathbf{x}_t, \mathbf{c}, t), \quad (6)$$

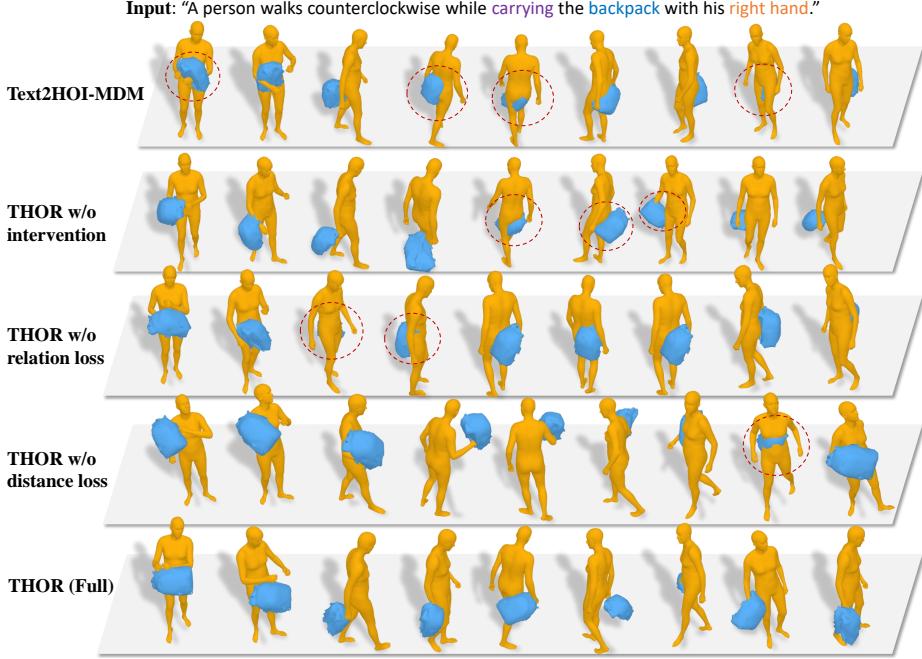


Fig. 3: Qualitative comparison for Text2HOI generation. Artifacts are highlighted in red circle. Our model, THOR, can generate realistic and plausible human-object interaction in response to the textual guidance. Through the intervention mechanism and two interaction losses, it corrects the drifting object motion and alleviate the implausible human-object spatial relations.

where s is the guidance scale. With this modification, THOR maintains a consistent awareness of the object’s shape, striking a balance between fidelity and diversity guided by the provided text prompt.

In each diffusion step t , the noisy interactions \mathbf{x}_t is linearly projected, and concatenated with the condition \mathbf{c} . We use a transformer encoder with M layers to predict the denoised $\tilde{\mathbf{x}}_0$ with positional encoding. This builds a general Text2HOI generation model if directly go to next step. However, such a simple diffusion model tends to learn the entire distribution of input data, while lacking the nuanced awareness regarding interactions. Our model, THOR, addresses this limitation by introducing human object relation intervention as in the next subsection.

3.2 Diffusion with Human-Object Relation Intervention

As highlighted in Sec. 1, text is a weak condition on object motion, which can lead to deficiencies in the overall interaction generation if object motion is directly sampled from it. Given object shape and coarse human motion, human ourselves can infer and envision plausible object motion in response. Acknowledging this,

the human motion can serve as auxiliary guidance to enhance semantic meanings of interactions learned from text, providing a more intuitive and contextually rich perspective in generation.

Modeling Human-object Kinematic Relations To bridge the gap between human motion and object motion in interaction, we introduce a novel intervention network designed to rectify implausible object motion. Rather than predicting object motion directly from human motion, which might neglect the initial learned interaction, our approach leverages their spatial relations to predict residual correction terms on object motion. These relations encapsulates the object's rotations and translations relative to each human joint from a human-centric perspective.

$$\gamma^{h \rightarrow o} = \mathbf{q} \ominus \mathbf{r}, \quad \tau^{h \rightarrow o} = \mathbf{j} \ominus \mathbf{o}. \quad (7)$$

Here, $\gamma^{h \rightarrow o} \in \mathbb{R}^{N \times J \times 3}$ represents the rotation relations, $\tau^{h \rightarrow o} \in \mathbb{R}^{N \times J \times 3}$ represents translation relations, and \ominus denotes human joint-wise subtraction. At the t -th step of diffusion model, the primitive motion $\tilde{\mathbf{x}}_0^h$ and $\tilde{\mathbf{x}}_0^o$ would be reformulated as $\tilde{\gamma}^{h \rightarrow o}$ and $\tilde{\tau}^{h \rightarrow o}$ similarly.

This unique human-centric relation representation enhances the model's understanding of the relationships between the human body and objects. By concentrating on object motion directly associated with human body poses, our approach not only improves the model's contextual awareness but also lays the groundwork for generating realistic, consistent, and semantically meaningful human-object interactions.

Relation Intervention Network Recognizing that rotations and translations are distinct transformations, we intervene separately to guarantee that the output residual intervention preserve a scale aligned with the input relation. After jointly modeling of human and object motion with two transformation representation, separating them can enhance the awareness of distinction between these two motion components. The intervention network is composed of two lightweight transformer encoders, named **Rotation encoder** and **Position encoder** respectively. Rotation relations and translation relations are individually input to these encoders and then we predict two residual terms $\Delta\tilde{\mathbf{r}}$ and $\Delta\tilde{\mathbf{o}}$. Consequently, we obtain a perturbation term on object motion, $\Delta\tilde{\mathbf{x}}_0^o = [\Delta\tilde{\mathbf{r}}, \Delta\tilde{\mathbf{o}}]$. Through the attention mechanism in transformer, it aggregates intervention from each human joint to improve the object motion. Finally, this $\Delta\tilde{\mathbf{x}}_0^o$ is added back to the primitive object motion $\tilde{\mathbf{x}}_0^o$ to obtain the final estimated $\hat{\mathbf{x}}_0^o$.

Benefited from the cyclical nature of the diffusion model, early rectification of object motion not only enhances the object motion but also contributes to the refinement of human motion throughout the diffusion process. This seamless integration eliminates the need for additional fusion operations, emphasizing the efficiency and coherence of our approach. We take both human motion and object motion as input to model their joint distribution. The simple objective [21, 58]

for our diffusion framework is

$$\mathcal{L}_{simple} = \mathbb{E}_{t \sim [1, T]} \|\hat{\mathbf{x}}_0 - \mathbf{x}_0\|. \quad (8)$$

Interaction Losses We present two interaction losses to supervise the kinematic relation and geometric distance between human and object to further enhance the generation ability. The first relation loss aims to supervise the kinematic relations to encourage interaction network to generate reasonable relations. Then the model is guided to generate interactions that not only adhere to textual descriptions but also align with reliable spatial configurations of human and object.

This enhances the overall realism of the generated human-object interactions. The relation loss is formulated as:

$$\mathcal{L}_{rel} = \mathbb{E}_{t \sim [1, T]} (\|\hat{\boldsymbol{\gamma}}_0^{h \rightarrow o} - \boldsymbol{\gamma}_0^{h \rightarrow o}\| + \|\hat{\boldsymbol{\tau}}_0^{h \rightarrow o} - \boldsymbol{\tau}_0^{h \rightarrow o}\|). \quad (9)$$

However, relation loss do not explicitly take the object shape into consideration, which do not modeling the geometric transformation relationship between object pose and points on its surface, and this may cause severe penetrations. Instead, the distance field between human joints and object surface can provide finer-grained geometry cues for HOI. The explicit distance restricts the object spatial distribution around human body by considering its shape. Therefore, we reconstruct their surface and calculate the signed distance between human joints and object points. Specifically, the mutual signed distance is represented as $\hat{\mathbf{d}}_0^{h \rightarrow o} \in \mathbb{R}^{N \times M \times 3}$ and $\hat{\mathbf{d}}_0^{o \rightarrow h} \in \mathbb{R}^{N \times M \times 3}$, And the distance loss can be written as:

$$\mathcal{L}_{dist} = \mathbb{E}_{t \sim [1, t']} (\|\hat{\mathbf{d}}_0^{h \rightarrow o} - \mathbf{d}_0^{h \rightarrow o}\| + \|\hat{\mathbf{d}}_0^{o \rightarrow h} - \mathbf{d}_0^{o \rightarrow h}\|). \quad (10)$$

Similar to [33, 68, 73], this is applied in only few diffusion steps under a threshold t' .

With the simple objective mentioned before, our approach supervise the interactions at multiple levels, which enhances the model's capacity to discern and understand interactions in a progressive way, providing a more comprehensive supervision of the underlying dynamics. The total training loss of our diffusion is summarized as:

$$\mathcal{L} = \mathcal{L}_{simple} + \lambda_{rel} \mathcal{L}_{rel} + \lambda_{dist} \mathcal{L}_{dist} + \lambda_{vel} \mathcal{L}_{vel}, \quad (11)$$

where $\mathcal{L}_{vel} = \mathbb{E}_{t \sim [1, T]} \|\hat{\mathbf{x}}_{0,1:N} - \hat{\mathbf{x}}_{0,0:N-1}\|$ represents velocity regularization.

4 Our Text-BEHAVE Dataset

We have enriched the BEHAVE dataset, currently the largest publicly accessible dataset for 3D Human-Object Interaction [6]. This enrichment involves the manual annotation of textual descriptions for each sequence. Specifically, each

sequence is partitioned into multiple clips based on the annotations, ensuring semantic consistency throughout every clip.

Besides, we discard clips that are meaningless or too short and randomly crop the clips with duration longer than 10s. This process results in a text annotated human-object interaction dataset, including total of 2377 interaction clips ranging from 2 to 10s. The total interaction length amounts 440, 840 frames at 30 fps. It contains 18 objects belonging to 12 categories (e.g. box, table, chair, backpack), and various common interactions (e.g. ‘lifting’, ‘sitting’, ‘dragging’) with these objects. The average length of textual descriptions is 19.7, which includes consecutive interactions with conjunction. As for interaction representation, human motion follows the format in SMPL-H [51] and object is represented in 6D pose including rotation and translation. In total, there are 2144 clips for training and 233 clips for testing, respectively. For brevity, we refer to this augmented dataset as **Text-BEHAVE**.

5 Experiments

Datasets and Metrics Our model are trained and tested on the Text-BEHAVE. As in [16], we evaluate the generation ability of our model with following metrics. *FID* evaluates the difference of distribution between the generated motion and real motion. *Diversity* evaluates the dissimilarity between generated motions, while *Modality* measures it within the same condition. *R precision* assesses the consistency between text and the generated motion in a retrieval way, *MM Dist* measures the feature distance between conditions and generated motions. We use contrastive loss to train a text encoder and motion encoder on the extended Text-BEHAVE dataset. Additionally, we also adopt the motion reconstruction metrics from [6, 31, 76]. We sample 20 times from input text and the best results are selected, as in [31, 68]. Furthermore, user study is introduced to evaluate the visual perceptual performance.

Baselines We compare our model with a VAE model whose encoder and decoder are from [44], and a motion diffusion model from [58]. They are modified to suit Text2HOI task, noted as ‘Text2HOI-VAE’ and ‘Text2HOI-MDM’ respectively.

Implementation details We use Adam optimizer and learning rate 10^{-4} . The batch size is set to 64 and the number of diffusion steps T is 1000. For the text encoder, we choose the CLIP pre-trained text encoder in the version of ‘ViT-B/32’. The text embedding is randomly masked by a ratio 0.1 for classifier-free guidance. The primitive generation uses a transformer encoder of 8 layers and intervention network comprises of 2 transformer layers. All the experiments were implemented on a single NVIDIA GeForce RTX 4090 with 24 GB of memory, and our model takes about 20 hours to converge.

5.1 Quantitative Results

As shown in Table 1, our approach outperforms the baselines across Text2HOI generation metrics. Through the intervention mechanism, our model can enhance the text guidance by the preliminary generated human motion, which further improve the performance on the retrieval-based metrics and multi-modality metrics. Additionally, THOR exhibits a significantly lower *FID* than Text2HOI-MDM, indicating a notable improvement in generating high-quality motions comparable to real ones. Our method also excels in terms of *diversity*, showcasing a more varied range of generated motions. This is crucial as it suggests that our approach can provide a broader range of human motions, contributing to more varied and realistic results. At last, the *MModality* score highlights the richness of multi-modality in our generated motions.

Table 1: Quantitative comparisons on the Text-BEHAVE test set for Text2HOI generation. All the evaluations run 20 times. \pm indicates the 95% confidence interval. **Bold** indicates best result.

Methods	R Precision↑			FID ↓	MM Dist↓	Diversity→	MModality ↑
	Top 1	Top 2	Top 3				
Real motions	0.994 \pm .003	0.998 \pm .001	1.000 \pm .000	0.031 \pm .002	3.377 \pm .004	22.832 \pm .074	—
Text2HOI-VAE	0.095 \pm .008	0.142 \pm .011	0.178 \pm .016	4.274 \pm .032	8.074 \pm .018	17.121 \pm .134	1.612 \pm .123
Text2HOI-MDM	0.186 \pm .012	0.286 \pm .010	0.348 \pm .007	2.680 \pm .028	7.065 \pm .014	22.376 \pm .088	2.261 \pm .079
THOR w/o intervention	0.194 \pm .007	0.306 \pm .011	0.387 \pm .015	1.909 \pm .015	7.008 \pm .018	23.015 \pm .141	2.247 \pm .109
THOR w/o dist. loss	0.198 \pm .007	0.314 \pm .008	0.390 \pm .015	1.936 \pm .018	6.933 \pm .010	23.025 \pm .086	2.393 \pm .053
THOR w/o rel. loss	0.218 \pm .005	0.338 \pm .008	0.420\pm.008	1.874\pm.022	6.887 \pm .011	22.475 \pm .077	2.093 \pm .070
THOR (Full)	0.250\pm.009	0.362\pm.008	0.411 \pm .001	1.983 \pm .026	6.874\pm.014	22.701\pm.081	2.575\pm.071

5.2 Qualitative Results

To demonstrate the visual performance of THOR, we provide qualitative results for evaluation. Text2HOI-VAE is excluded since it generate implausible results. As shown in Figure 3, we give an example text prompt and the generated interactions. To demonstrate the effectiveness of our model, we highlight the artifacts in other alternatives. In the absence of intervention, models may produce the drifting object or interactions that is misaligned with textual descriptions. And without relation loss and distance loss, generated interactions may encounter penetration or fail to establish contact with the human body. However, our model demonstrates the capability to generate consistent and realistic human-object interactions, ensuring the object motion aligns with the human motion through the intervention mechanism and interaction losses.

The user study is also conducted, serving as a complementary evaluation. We carefully selected 30 samples, encompassing 10 object classes, each generated

from 3 text prompts for every object. Subsequently, we invited 27 users to assess the generation quality and express their preferences between two choices. The results, depicted in Figure 5, demonstrate that our model outperforms others and is comparable to the ground truth in certain samples.

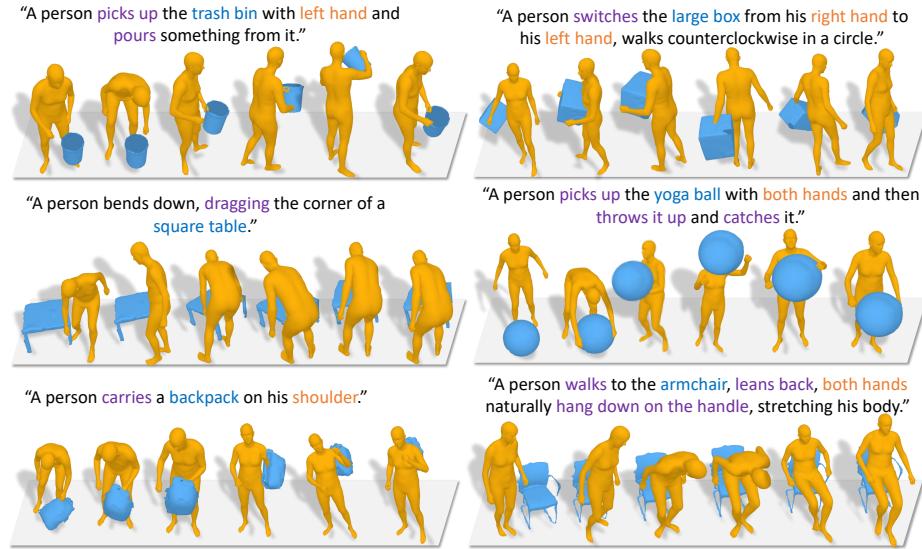


Fig. 4: Qualitative results of our model for Text2HOI generation. Our model can generate human-object interactions aligned with the text description involving static and dynamic objects with diverse categories and shapes.

If guidance scale $s = 0$ in sampling, the model turns to generate interactions solely conditioned on a object model. As shown in Figure 6, without textual guidance, our model can generate human-object interactions satisfying different object shape. This also verifies the effectiveness of the modified HOI classifier-free guidance, which strikes a balance for the fidelity matching the textual description and diversity with respect to the object geometry.

5.3 Ablation study

To better understand the contribution of each component in our proposed method, THOR, we conducted ablation studies by removing specific elements from the model architecture. The quantitative results are summarized in Table 1, and we discuss them as below: Our THOR model achieves the best overall performance. The absence of intervention mechanism causes a distinct decrease in *R Precision*, *MM Dist*, and *MModality*, confirming that the intervention from human motion guidance ensures a more precise alignment with textual guidance. The comparable *FID* and *Diversity* scores indicate that the model retains its ability

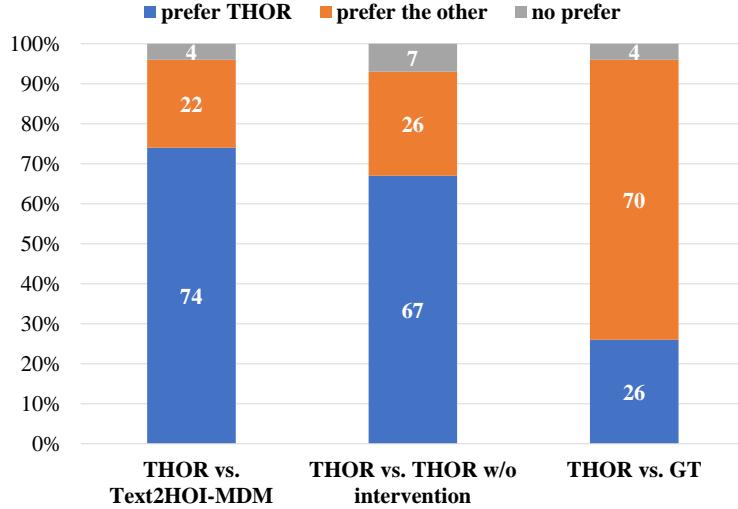


Fig. 5: User study on Text-BEHAVIE test set, where ‘w/o’ indicates ‘without’ and ‘GT’ indicates ground truth.

to generate diverse and plausible motions. The distance loss results in a minor decrease in *R Precision* and *Diversity*, indicating that this loss component contributes to the accuracy and variety of generated motions. Surprisingly, the *FID* score improves, suggesting that without distance loss, the model generates motions that are closer to real motions but potentially sacrifices diversity. Excluding the relation loss slightly reduces accuracy but leads to improved visual fidelity, suggesting that relation loss may introduce certain variations in generated motions. Applying these two interaction losses provides additional guidance to the model. By supervising the relations in the diffusion steps, the model is encouraged to generate reasonable spatial relations from the outset, contributing to the overall coherence of the generated interactions. This confirms that the combination of all proposed components leads to a balanced generation of fidelity, diversity guided by textual input.

Qualitative ablation studies are illustrated in Figure 3, highlighting the proficiency of our complete model in generating interactions with superior visual performance. The intervention mechanism and the relation loss addresses issues such as drifting object motion and meaningless trajectories, while the distance loss effectively prevents penetration. The ablation studies reveal that each component in THOR contributes to specific aspects of motion generation. The intervention mechanism and relation loss ensures precise and diverse motion synthesis, while distance loss improves the visual fidelity. THOR achieves a synergistic effect, striking a balance between accuracy and diversity in text-guided interaction generation.

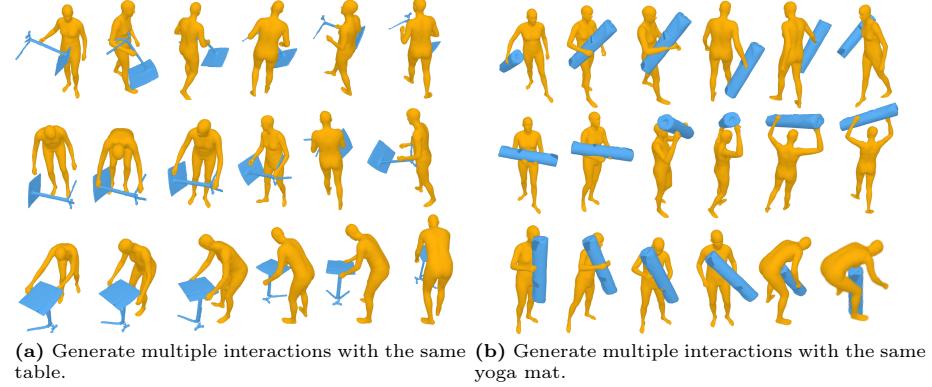


Fig. 6: Sampling interactions from only object shape. It shows that our model can generate plausible interactions in response to a given object shape.

6 Conclusion

In summary, we proposed a novel diffusion model to deal with Text2HOI, addressing the challenge of generating full-body interactions with dynamic objects guided by textual prompts. This task involves the simultaneous generation of consistent human and object motion, posing a significant difficulty. To address this, we present THOR, a diffusion framework seamlessly integrating human-object interaction and intervention mechanisms within a unified end-to-end framework. We interpret interactions based on the kinematic relation and geometric distance between humans and object surfaces, introducing two interaction losses to generate plausible and realistic results. Additionally, we contribute the Text-BEHAVE dataset, integrating textual descriptions into the largest publicly available 3D HOI dataset. Both quantitative and qualitative experiments are conducted to validate the effectiveness of THOR in generating coherent and realistic human-object interactions.

Limitations and future works Though our model generates realistic results, issues like penetration and floating still exists in some generated results, especially for objects with complex shape and unseen objects. The existing HOI datasets fall short in both scale and quality compared to human motion datasets, highlighting the clear need for more extensive and high-quality HOI data. The expansion of these datasets would be highly valuable to advance research in human-object interactions. Additionally, future works could benefit from integrating dexterous hand motion to generate more comprehensive human-object interactions. The exploration of long-term generation and fine-granularity control remains an area for further investigation.

References

1. Alexanderson, S., Nagy, R., Beskow, J., Henter, G.E.: Listen, denoise, action! audio-driven motion synthesis with diffusion models. *ACM Transactions on Graphics (TOG)* **42**(4), 1–20 (2023)
2. Athanasiou, N., Petrovich, M., Black, M.J., Varol, G.: Teach: Temporal action composition for 3D humans. In: 2022 International Conference on 3D Vision (3DV). pp. 414–423. IEEE (2022)
3. Athanasiou, N., Petrovich, M., Black, M.J., Varol, G.: Sinc: Spatial composition of 3D human motions for simultaneous action generation. arXiv preprint arXiv:2304.10417 (2023)
4. Bae, J., Won, J., Lim, D., Min, C.H., Kim, Y.M.: Pmp: Learning to physically interact with environments using part-wise motion priors. arXiv preprint arXiv:2305.03249 (2023)
5. Bao, F., Nie, S., Xue, K., Li, C., Pu, S., Wang, Y., Yue, G., Cao, Y., Su, H., Zhu, J.: One transformer fits all distributions in multi-modal diffusion at scale. In: International Conference on Machine Learning. pp. 1692–1717. PMLR (2023)
6. Bhatnagar, B.L., Xie, X., Petrov, I.A., Sminchisescu, C., Theobalt, C., Pons-Moll, G.: Behave: Dataset and method for tracking human object interactions. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 15935–15946 (2022)
7. Bhattacharyya, A., Schiele, B., Fritz, M.: Accurate and diverse sampling of sequences based on a “best of many” sample objective. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 8485–8493 (2018)
8. Cao, Z., Gao, H., Mangalam, K., Cai, Q.Z., Vo, M., Malik, J.: Long-term human motion prediction with scene context. In: Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16. pp. 387–404. Springer (2020)
9. Cervantes, P., Sekikawa, Y., Sato, I., Shinoda, K.: Implicit neural representations for variable length human motion generation. In: European Conference on Computer Vision. pp. 356–372. Springer (2022)
10. Chao, Y.W., Yang, J., Chen, W., Deng, J.: Learning to sit: Synthesizing human-chair interactions via hierarchical control. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 35, pp. 5887–5895 (2021)
11. Chen, X., Jiang, B., Liu, W., Huang, Z., Fu, B., Chen, T., Yu, G.: Executing your commands via motion diffusion in latent space. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 18000–18010 (2023)
12. Dabral, R., Mughal, M.H., Golyanik, V., Theobalt, C.: Mofusion: A framework for denoising-diffusion-based motion synthesis. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 9760–9770 (June 2023)
13. Diller, C., Dai, A.: Cg-hoi: Contact-guided 3d human-object interaction generation. arXiv preprint arXiv:2311.16097 (2023)
14. Ghosh, A., Dabral, R., Golyanik, V., Theobalt, C., Slusallek, P.: IMoS: Intent-driven full-body motion synthesis for human-object interactions. arXiv preprint arXiv:2212.07555 (2022)
15. Gong, K., Lian, D., Chang, H., Guo, C., Jiang, Z., Zuo, X., Mi, M.B., Wang, X.: Tm2d: Bimodality driven 3D dance generation via music-text integration. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 9942–9952 (2023)

16. Guo, C., Zou, S., Zuo, X., Wang, S., Ji, W., Li, X., Cheng, L.: Generating diverse and natural 3D human motions from text. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5152–5161 (2022)
17. Guo, C., Zuo, X., Wang, S., Cheng, L.: Tm2t: Stochastic and tokenized modeling for the reciprocal generation of 3D human motions and texts. In: European Conference on Computer Vision. pp. 580–597. Springer (2022)
18. Han, S., Joo, H.: CHORUS: Learning canonicalized 3D human-object spatial relations from unbounded synthesized images. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 15835–15846 (2023)
19. Hassan, M., Ghosh, P., Tesch, J., Tzionas, D., Black, M.J.: Populating 3D scenes by learning human-scene interaction. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 14708–14718 (2021)
20. Hassan, M., Guo, Y., Wang, T., Black, M., Fidler, S., Peng, X.B.: Synthesizing physical character-scene interactions. arXiv preprint arXiv:2302.00883 (2023)
21. Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models. Advances in neural information processing systems **33**, 6840–6851 (2020)
22. 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: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 16750–16761 (2023)
23. Huang, Y., Taheri, O., Black, M.J., Tzionas, D.: InterCap: Joint markerless 3D tracking of humans and objects in interaction. In: DAGM German Conference on Pattern Recognition. pp. 281–299. Springer (2022)
24. Huo, C., Shi, Y., Ma, Y., Xu, L., Yu, J., Wang, J.: Stackflow: Monocular human-object reconstruction by stacked normalizing flow with offset. In: Elkind, E. (ed.) Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence, IJCAI-23. pp. 902–910. International Joint Conferences on Artificial Intelligence Organization (8 2023). <https://doi.org/10.24963/ijcai.2023/100>, <https://doi.org/10.24963/ijcai.2023/100>, main Track
25. Jiang, N., Liu, T., Cao, Z., Cui, J., Zhang, Z., Chen, Y., Wang, H., Zhu, Y., Huang, S.: Full-body articulated human-object interaction. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 9365–9376 (2023)
26. Kim, J., Kim, J., Choi, S.: Flame: Free-form language-based motion synthesis & editing. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 37, pp. 8255–8263 (2023)
27. 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 preprint arXiv:2307.07511 (2023)
28. Li, B., Zhao, Y., Zhelun, S., Sheng, L.: Danceformer: Music conditioned 3D dance generation with parametric motion transformer. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 36, pp. 1272–1279 (2022)
29. Li, J., Clegg, A., Mottaghi, R., Wu, J., Puig, X., Liu, C.K.: Controllable human-object interaction synthesis. arXiv preprint arXiv:2312.03913 (2023)
30. Li, J., Liu, K., Wu, J.: Ego-body pose estimation via ego-head pose estimation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 17142–17151 (2023)
31. Li, J., Wu, J., Liu, C.K.: Object motion guided human motion synthesis. arXiv preprint arXiv:2309.16237 (2023)
32. Li, R., Yang, S., Ross, D.A., Kanazawa, A.: Ai choreographer: Music conditioned 3D dance generation with aist++. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 13401–13412 (2021)

33. Liang, H., Zhang, W., Li, W., Yu, J., Xu, L.: Intergen: Diffusion-based multi-human motion generation under complex interactions. arXiv preprint arXiv:2304.05684 (2023)
34. Lim, D., Jeong, C., Kim, Y.M.: MAMMOS: Mapping multiple human motion with scene understanding and natural interactions. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 4278–4287 (2023)
35. Liu, L., Hodgins, J.: Learning basketball dribbling skills using trajectory optimization and deep reinforcement learning. ACM Transactions on Graphics (TOG) **37**(4), 1–14 (2018)
36. Loper, M., Mahmood, N., Romero, J., Pons-Moll, G., Black, M.J.: Smpl: a skinned multi-person linear model. ACM Transactions on Graphics (TOG) **34**(6), 1–16 (2015)
37. Lucas, T., Baradel, F., Weinzaepfel, P., Rogez, G.: Posegpt: Quantization-based 3D human motion generation and forecasting. In: European Conference on Computer Vision. pp. 417–435. Springer (2022)
38. Merel, J., Tunyasuvunakool, S., Ahuja, A., Tassa, Y., Hasenclever, L., Pham, V., Erez, T., Wayne, G., Heess, N.: Catch & carry: reusable neural controllers for vision-guided whole-body tasks. ACM Transactions on Graphics (TOG) **39**(4), 39–1 (2020)
39. Pavlakos, G., Choutas, V., Ghorbani, N., Bolkart, T., Osman, A.A., Tzionas, D., Black, M.J.: Expressive body capture: 3d hands, face, and body from a single image. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 10975–10985 (2019)
40. Peng, X., Xie, Y., Wu, Z., Jampani, V., Sun, D., Jiang, H.: Hoi-diff: Text-driven synthesis of 3d human-object interactions using diffusion models. arXiv preprint arXiv:2312.06553 (2023)
41. Petrov, I.A., Marin, R., Chibane, J., Pons-Moll, G.: Object pop-up: Can we infer 3D objects and their poses from human interactions alone? In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 4726–4736 (2023)
42. Petrovich, M., Black, M.J., Varol, G.: Action-conditioned 3D human motion synthesis with transformer VAE. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 10985–10995 (2021)
43. Petrovich, M., Black, M.J., Varol, G.: Action-conditioned 3D human motion synthesis with transformer vae. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). pp. 10985–10995 (October 2021)
44. Petrovich, M., Black, M.J., Varol, G.: Temos: Generating diverse human motions from textual descriptions. In: European Conference on Computer Vision. pp. 480–497. Springer (2022)
45. Pi, H., Peng, S., Yang, M., Zhou, X., Bao, H.: Hierarchical generation of human-object interactions with diffusion probabilistic models. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 15061–15073 (2023)
46. Qi, C.R., Su, H., Mo, K., Guibas, L.J.: Pointnet: Deep learning on point sets for 3D classification and segmentation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 652–660 (2017)
47. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: International conference on machine learning. pp. 8748–8763. PMLR (2021)

48. Ramesh, A., Dhariwal, P., Nichol, A., Chu, C., Chen, M.: Hierarchical text-conditional image generation with clip latents. arXiv preprint arXiv:2204.06125 **1**(2), 3 (2022)
49. Razali, H., Demiris, Y.: Action-conditioned generation of bimanual object manipulation sequences. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 37, pp. 2146–2154 (2023)
50. Ren, Z., Pan, Z., Zhou, X., Kang, L.: Diffusion motion: Generate text-guided 3D human motion by diffusion model. In: ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). pp. 1–5. IEEE (2023)
51. Romero, J., Tzionas, D., Black, M.J.: Embodied hands: Modeling and capturing hands and bodies together. ACM Transactions on Graphics, (Proc. SIGGRAPH Asia) **36**(6) (Nov 2017)
52. Shafir, Y., Tevet, G., Kapon, R., Bermano, A.H.: Human motion diffusion as a generative prior. arXiv preprint arXiv:2303.01418 (2023)
53. Starke, S., Zhang, H., Komura, T., Saito, J.: Neural state machine for character-scene interactions. ACM Trans. Graph. **38**(6), 209–1 (2019)
54. Taheri, O., Choutas, V., Black, M.J., Tzionas, D.: GOAL: Generating 4D whole-body motion for hand-object grasping. In: Conference on Computer Vision and Pattern Recognition (CVPR) (2022), <https://goal.is.tue.mpg.de>
55. Taheri, O., Zhou, Y., Tzionas, D., Zhou, Y., Ceylan, D., Pirk, S., Black, M.J.: GRIP: Generating interaction poses using latent consistency and spatial cues. arXiv preprint arXiv:2308.11617 (2023)
56. Tanaka, M., Fujiwara, K.: Role-aware interaction generation from textual description. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). pp. 15999–16009 (October 2023)
57. Tevet, G., Gordon, B., Hertz, A., Bermano, A.H., Cohen-Or, D.: Motionclip: Exposing human motion generation to clip space. In: European Conference on Computer Vision. pp. 358–374. Springer (2022)
58. Tevet, G., Raab, S., Gordon, B., Shafir, Y., Cohen-or, D., Bermano, A.H.: Human motion diffusion model. In: The Eleventh International Conference on Learning Representations (2023), <https://openreview.net/forum?id=SJ1kSy02jwu>
59. Wang, Y., Leng, Z., Li, F.W.B., Wu, S.C., Liang, X.: Fg-T2M: Fine-grained text-driven human motion generation via diffusion model. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). pp. 22035–22044 (October 2023)
60. Wang, Z., Chen, Y., Liu, T., Zhu, Y., Liang, W., Huang, S.: Humanise: Language-conditioned human motion generation in 3D scenes. Advances in Neural Information Processing Systems **35**, 14959–14971 (2022)
61. Wu, Y., Wang, J., Zhang, Y., Zhang, S., Hilliges, O., Yu, F., Tang, S.: Saga: Stochastic whole-body grasping with contact. In: European Conference on Computer Vision. pp. 257–274. Springer (2022)
62. Xiao, Z., Wang, T., Wang, J., Cao, J., Zhang, W., Dai, B., Lin, D., Pang, J.: Unified human-scene interaction via prompted chain-of-contacts. arXiv preprint arXiv:2309.07918 (2023)
63. Xie, X., Bhatnagar, B.L., Pons-Moll, G.: Chore: Contact, human and object reconstruction from a single rgb image. In: European Conference on Computer Vision. pp. 125–145. Springer (2022)
64. Xie, X., Bhatnagar, B.L., Pons-Moll, G.: Visibility aware human-object interaction tracking from single rgb camera. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 4757–4768 (2023)

65. Xie, Y., Jampani, V., Zhong, L., Sun, D., Jiang, H.: Omnicontrol: Control any joint at any time for human motion generation. arXiv preprint arXiv:2310.08580 (2023)
66. Xie, Z., Tseng, J., Starke, S., van de Panne, M., Liu, C.K.: Hierarchical planning and control for box loco-manipulation. arXiv preprint arXiv:2306.09532 (2023)
67. Xu, L., Song, Z., Wang, D., Su, J., Fang, Z., Ding, C., Gan, W., Yan, Y., Jin, X., Yang, X., Zeng, W., Wu, W.: ActFormer: A gan-based transformer towards general action-conditioned 3D human motion generation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). pp. 2228–2238 (October 2023)
68. Xu, S., Li, Z., Wang, Y.X., Gui, L.Y.: Interdiff: Generating 3D human-object interactions with physics-informed diffusion. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 14928–14940 (2023)
69. Yang, Z., Su, B., Wen, J.R.: Synthesizing long-term human motions with diffusion models via coherent sampling. arXiv preprint arXiv:2308.01850 (2023)
70. Ye, S., Wang, Y., Li, J., Park, D., Liu, C.K., Xu, H., Wu, J.: Scene synthesis from human motion. In: SIGGRAPH Asia 2022 Conference Papers. SA '22, Association for Computing Machinery, New York, NY, USA (2022). <https://doi.org/10.1145/3550469.3555426>
71. Yi, H., Huang, C.H.P., Tripathi, S., Hering, L., Thies, J., Black, M.J.: Mime: Human-aware 3D scene generation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12965–12976 (2023)
72. Yi, H., Huang, C.H.P., Tzionas, D., Kocabas, M., Hassan, M., Tang, S., Thies, J., Black, M.J.: Human-aware object placement for visual environment reconstruction. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 3959–3970 (June 2022)
73. Yuan, Y., Song, J., Iqbal, U., Vahdat, A., Kautz, J.: Physdiff: Physics-guided human motion diffusion model. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 16010–16021 (2023)
74. Zhang, J.Y., Pepose, S., Joo, H., Ramanan, D., Malik, J., Kanazawa, A.: Perceiving 3D human-object spatial arrangements from a single image in the wild. In: European Conference on Computer Vision (ECCV) (2020)
75. 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. arXiv preprint arXiv:2301.06052 (2023)
76. Zhang, J., Luo, H., Yang, H., Xu, X., Wu, Q., Shi, Y., Yu, J., Xu, L., Wang, J.: NeuralDome: A neural modeling pipeline on multi-view human-object interactions. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 8834–8845 (2023)
77. Zhang, M., Guo, X., Pan, L., Cai, Z., Hong, F., Li, H., Yang, L., Liu, Z.: ReMoDiffuse: Retrieval-augmented motion diffusion model. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). pp. 364–373 (October 2023)
78. Zhang, X., Bhatnagar, B.L., Starke, S., Guzov, V., Pons-Moll, G.: Couch: Towards controllable human-chair interactions. In: European Conference on Computer Vision. pp. 518–535. Springer (2022)
79. Zhao, K., Wang, S., Zhang, Y., Beeler, T., Tang, S.: Compositional human-scene interaction synthesis with semantic control. In: European Conference on Computer Vision. pp. 311–327. Springer (2022)
80. Zhao, K., Zhang, Y., Wang, S., Beeler, T., Tang, S.: Synthesizing diverse human motions in 3D indoor scenes. In: International conference on computer vision (ICCV) (2023)

81. Zhong, C., Hu, L., Zhang, Z., Xia, S.: Attt2m: Text-driven human motion generation with multi-perspective attention mechanism. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 509–519 (2023)
82. Zhou, K., Bhatnagar, B.L., Lenssen, J.E., Pons-Moll, G.: Toch: Spatio-temporal object-to-hand correspondence for motion refinement. In: European Conference on Computer Vision. pp. 1–19. Springer (2022)
83. Zhuang, W., Wang, C., Chai, J., Wang, Y., Shao, M., Xia, S.: Music2dance: Dancenet for music-driven dance generation. ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM) **18**(2), 1–21 (2022)

Supplementary materials includes additional details of our Text-BEHAVE dataset, experiment results and evaluation model. Sec. A showcases the object models in our dataset and describes the details of annotation. Sec. C provides more experiment results on out-of-the-dataset generation, guidance scale analysis and object conditioned generation. Sec. D gives training details of the evaluation model. Sec. B adds other implementation details and Sec. E shows some failure cases.

A Data analysis

A.1 Object Categories

Text-BEHAVE comprises 18 object models, for which the original dataset furnishes comprehensive motion data captured at 30 fps. As depicted in Figure A1, these models span diverse categories, sizes, and geometries.

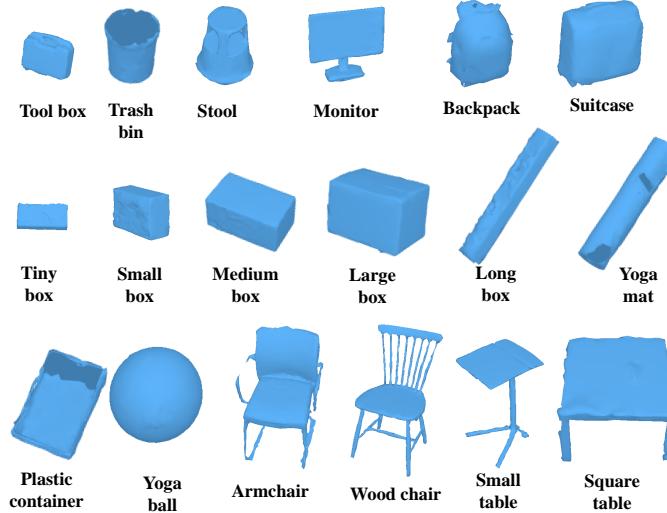
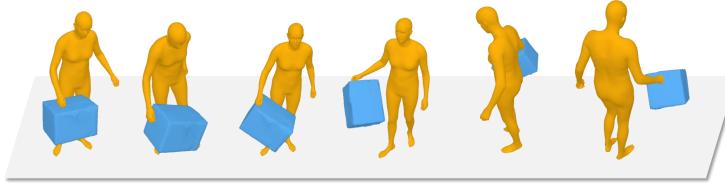


Fig. A1: Object models in Text-BEHAVE.

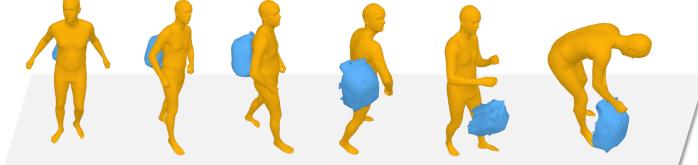
A.2 Details of Text Annotation

We provide a instruction for the Crowd-sourcing Platform to describe the interactions in the videos of BEHAVE [6], and then manually replace the object name. Since origin BEHAVE dataset includes amounts of non-interactive and redundant actions in a long motion sequence, we ask the annotators to divide the long sequence and mark out the meaningless segments. A detailed instruction and textual example are as the follows,

Annotation Requirement:



1. A man walks back and forth with the large box **held** in **his right hand**.
2. A person strolls back and forth, **grasping** a large **box** in **his right hand**.
3. A person **carries** a large box with **his right hand**, walking clockwise.



1. A person **takes** the **backpack** off his **right shoulder** and **puts it on the ground**.
2. A person turns around, **removes** the **backpack** from his **right shoulder** and **sets it down**.
3. A person **takes** the **backpack** off his **right shoulder** and **puts it down**.

Fig. A2: Examples in Text-BEHAVIE dataset. Each interactions segment is annotated with 3 textual descriptions.

- Divide the entire sequence into multiple interaction segments, ensuring each segment has complete semantics, and keep the duration within 3 to 10 seconds.
- Each interaction description should include **character's action**, **object name**, and try to include **body part** in contact with the object **changes in the object's position** and **changes in the person's position**.
- For long videos, provide both a summary description of the entire video and segmented descriptions.

Example Segments:

1. ‘A person bends down to pick up a backpack from the ground in front’.
2. ‘Holding the bottom of the backpack with both hands, he walks counter-clockwise’.
3. ‘He lifts the backpack onto his left shoulder, with his left hand supporting the bottom, walking clockwise’.
4. ‘He shifts the backpack from his left shoulder to the right shoulder, starting to walk counterclockwise’.
5. ‘Using both hands, he raises the backpack above his head’.
6. ‘With both hands, he moves the backpack from above his head to the front of his chest’.
7. ‘He walks, holding the backpack with both hands, pressed against his abdomen’.

Each interaction segment is annotated with 3 textual descriptions, as shown in Figure A2.

B Additional Implementation Details

Transformer layers in early generation are with hidden size 512, feed forward size 1024, dropout 0.1 and ‘gelu’ activation. Transformer layers in Position encoder and Rotation encoder are with hidden size 156, feed forward size 128, dropout 0.1 and ‘gelu’ activation. Loss weights are $\lambda_{rel} = 0.01$, $\lambda_{dist} = 0.01$ and $\lambda_{vel} = 0.0003$ respectively. Guidance scale s is set to 2.5 when sampling the interactions. The total training epochs are 200.

C More Results

C.1 Effectiveness of Guidance Scale

In Table A1, we test different guidance scales and show the effectiveness of large guidance scale.

Table A1: Quantitative comparisons of different guidance scales on the Text-BEHAVE test set for Text2HOI generation. All the evaluations run 20 times. \pm indicates the 95% confidence interval. **Bold** indicates best result.

Guidance scale	R Precision↑			FID ↓	MM Dist↓	Diversity→	MModality ↑
	Top 1	Top 2	Top 3				
Real motions	0.994 \pm .003	0.998 \pm .001	1.000 \pm .000	0.031 \pm .002	3.377 \pm .004	22.832 \pm .074	—
0.0	0.094 \pm .007	0.183 \pm .014	0.232 \pm .011	2.742 \pm .033	7.428 \pm .018	22.271 \pm .114	1.612 \pm .123
0.5	0.181 \pm .010	0.277 \pm .014	0.353 \pm .015	2.493 \pm .060	7.057 \pm .030	22.533 \pm .110	3.781 \pm .205
1.0	0.225 \pm .008	0.328 \pm .009	0.404 \pm .009	2.133 \pm .021	6.907 \pm .012	22.676 \pm .088	3.375 \pm .093
1.5	0.224 \pm .009	0.328 \pm .009	0.399 \pm .010	2.033 \pm .023	7.008 \pm .018	22.705 \pm .141	2.717 \pm .109
2.0	0.249 \pm .009	0.358 \pm .010	0.409 \pm .003	2.048 \pm .021	6.907 \pm .014	22.687 \pm .078	2.535 \pm .109
2.5	0.250 \pm .009	0.362 \pm .008	0.411 \pm .001	1.989 \pm .020	6.874 \pm .014	22.701 \pm .081	2.575 \pm .071
3.0	0.224 \pm .009	0.334 \pm .011	0.408 \pm .010	2.066 \pm .016	6.905 \pm .010	22.764 \pm .058	2.492 \pm .010
4.0	0.231 \pm .008	0.348 \pm .008	0.426 \pm .007	2.081 \pm .018	6.883 \pm .011	22.655 \pm .080	2.513 \pm .063

C.2 Out-of-the-dataset Generation

We additionally annotate the InterCap dataset [23]. However, due to its scale and the transition cost from SMPL-X [39] to SMPL-H [51], we opt to selectively include its object models for supplementary testing, particularly in scenarios involving out-of-the-dataset generation. Notably, our model demonstrates proficiency in generating interactions with objects that are beyond the scope of the training data, as in Figure A3.

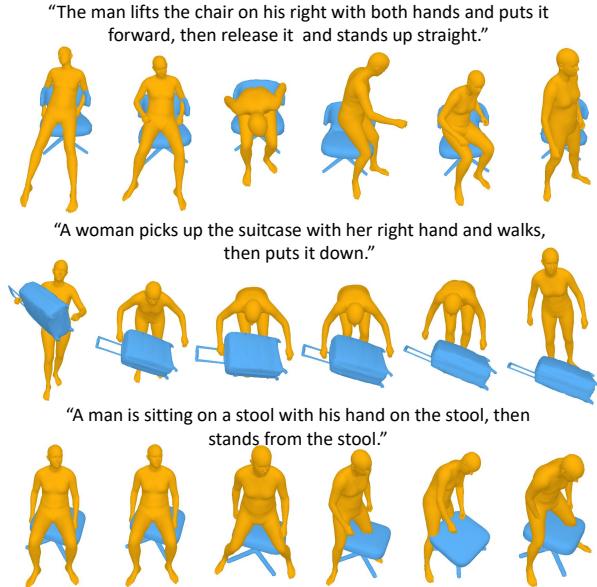


Fig. A3: We also test our model on the InterCap [23]. Even if the object shape is not in Text-BEHAVIE, our model demonstrates the ability of generalization with respect to the object shape.

C.3 Motion Reconstruction

Motion reconstruction metrics using the Best-of-many [7] method only provides a marginal reference for possible best result over the test set [31, 68], as in Tab. A2. Through 20 times repeated generation under the same text, THOR can generate interactions mostly close to the ground truth, especially for the object motion. It also turns out the human motion can also be refined during the diffusion process with intervention and intervention losses. *MPJPE* refers to mean per-joint position error, $v2v_h$ and $v2v_o$ refers to vertex position error of human and object [76], *Rot. Err.* refers to the Frobenius norm of rotation matrix difference and *Tr. Err.* is the translation error of object [31].

D Evaluation Details

Our evaluation model follows [33, 57], which is composed of a text feature extractor and human-object motion feature extractor. The text feature extractor comprises token embedding from CLIP [47] and a transformer encoder of 8 layers with 4 heads. The motion feature extractor is the same as in our model THOR, including a transformer encoder of 8 layers with 4 heads. Their hidden size is 512 and feed forward size is 1024. The latent dimension for text and motion are both 512. It is trained on the entire Text-BEHAVIE dataset with contrastive

Table A2: Best-of-many Motion reconstruction metrics. All the evaluations run 20 times and the best results are selected.

Methods	Human		Object		
	MPJPE↓ v2v _h ↓	Rot. Err. ↓	Tr. Err.↓ v2v _o ↓		
Text-MDM	60.60	56.30	1.80	53.82	60.03
THOR w/o intervention	59.71	55.61	1.76	51.80	58.40
THOR w/o dist. loss	59.50	55.18	1.77	51.74	57.90
THOR w/o rel. loss	59.85	55.56	1.80	51.79	58.09
THOR (Full)	59.69	55.44	1.73	51.25	57.71

loss [33, 47, 57]:

$$\begin{aligned} \mathcal{L}_{text} &= \text{CrossEntropy}(\mathbf{z}_{text}, \mathbf{z}_{cls}) \\ \mathcal{L}_{motion} &= \text{CrossEntropy}(\mathbf{z}_{motion}, \mathbf{z}_{cls}) \\ \mathcal{L} &= \frac{1}{2}(\mathcal{L}_{text} + \mathcal{L}_{motion}) \end{aligned} \quad (12)$$

where \mathbf{z}_{text} and \mathbf{z}_{motion} are extracted text and motion features, and \mathbf{z}_{cls} are labels following [47]. The evaluation model is trained with batch size of 64 and learning rate $1e^{-4}$ with weight decay $1e^{-4}$.

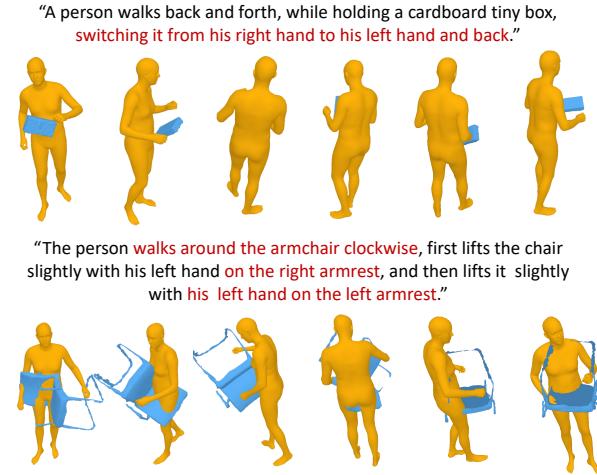


Fig. A4: Examples of failure cases. It mainly arises from discrepancies between complex text prompts and penetrations posed by objects with intricate geometries.

E Failure Cases

Text2HOI is a challenging task, and there are failed generation results as well. Some failure examples are shown in Figure A4. Floating movements of objects and penetrations between humans and objects may occur, particularly when dealing with objects that are excessively small or possess an exceedingly complex structure. Besides, there are failures dealing with the long-term textual descriptions, leading to missing interactions or transitions.