

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

Supplementary Material

Paper ID 1669

001 1. Overview

002 In this supplementary material, we commence by delineating
 003 the detailed HOIAnimator used for computing the metrics
 004 (Section 2). Next, we provide an in-depth presentation
 005 of the user study details (Section 3). Subsequently, we
 006 delve into additional experiments (Section 4), and provide
 007 an expanded set of results (Section 5). Finally, we provide a
 008 comprehensive illustration of the network architecture (Section 6).

010 2. Evaluation Criteria

011 In this section, we detail the text and Human-Object Inter-
 012 action (HOI) animation feature extractors, which are cru-
 013 cial components of our metric computation framework. Ad-
 014 ditionally, our evaluation criteria for the HOIAnimator en-
 015 compasses six key metrics: the Fréchet Inception Distance
 016 (FID), Diversity, MultiModality Diversity (MM-Dist), R
 017 Precision, and Penetration. The comprehensive testing pro-
 018 tocol is meticulously designed to rigorously evaluate the ef-
 019 ficacy and robustness of our approach in generating HOI
 020 animations across diverse scenarios.

021 **Text and HOI Animation Feature Extractor.** Following
 022 the approach outlined in T2M [2], our HOIAnimator
 023 involves utilizing a text extractor to convert raw text into a
 024 semantic feature vector, denoted as (s). Concurrently, HOI
 025 animations are processed through an HOI animation extrac-
 026 tor, resulting in another feature vector (m). In this process,
 027 we aim to minimize the distance between matched pairs of
 028 text and HOI animation feature vectors, ensuring feature
 029 vectors' close correspondence.

030 **FID.** FID calculates the distance between real samples
 031 and generated samples in latent space. Following the an-
 032 imation generation work [4], we define $FID(x, \hat{x}) = \|$
 $\mu_x - \mu_{\hat{x}}\|_2^2 + Tr(\Sigma_p + \Sigma_{\hat{x}} - 2(\Sigma_p \Sigma_{\hat{x}})^{0.5})$. Here, x and
 \hat{x} represent the real HOI animations and the generated HOI
 033 animations. FID is an objective metric calculating the dis-
 034 tance between features extracted from real and generated
 035 motion sequences, which reflects the generation quality.

036 **Diversity.** Diversity evaluates the variability of the gen-
 037 erated HOI animations across a range of descriptions. To
 038 measure this, we randomly sample two subsets of equal size
 (S_d) , from the entire collection of motions generated from
 039 various descriptions. We extract respective sets of HOI an-
 040 imation feature vectors $\{m_1, \dots, m_{S_d}\}$ and $\{m'_1, \dots, m'_{S_d}\}$

041 . The diversity of the set of motions is defined as,
 $Diversity = \frac{1}{S_d} \sum_{t=1}^{S_d} \|m_i - m'_i\|$, where $S_d = 300$ is
 042 used in experiments.

043 **MM-Dist.** Distinct from the concept of diver-
 044 sity, MM-Dist quantifies the extent to which the gen-
 045 erated motions vary within each individual text descrip-
 046 tion. It is measured across a dataset containing moti-
 047 ons paired with C distinct descriptions. For c -th de-
 048 scription, we randomly sample two subsets with the same
 049 size S_m and then extract two subsets of feature vectors
 $\{m_{c,1}, \dots, m_{c,S_m}\}$ and $\{m'_{c,1}, \dots, m'_{c,S_m}\}$. The multimodal-
 050 ity of the motion set is formalized as, $MM - Dis =$
 $\frac{1}{C*S} \sum_{c=1}^C \sum_{t=1}^{S_m} \|m_{c,i} - m'_{c,i}\|$, where $S_m = 10$ is used
 051 in experiments.

052 **R Precision.** The R precision metric evaluates the simi-
 053 larity between the textual description and the generated mo-
 054 tion sequence. It represents the likelihood that the actual
 055 text ranks within the top k positions after sorting. In this
 056 study, we set k to be 1, 2, and 3.

057 **Vertex distance.** Vertex distance evaluates generation
 058 quality by comparing distances between vertices in real and
 059 generated objects.

060 **Penetration.** Similar to prior work [3], the penetration
 061 score evaluates whether the human gets close to the object
 062 during the interaction. We define the approach phase as the
 063 initial motion frames from a sequence N_A . Then the pen-
 064 etration distance for a trajectory is $\frac{1}{N_A} \sum_v \sum_i^{N_A} \text{sdf}(v) \cdot$
 $\mathbb{1}_{\text{sdf}(v)>0}$, where $\mathbb{1}$ is indicator function, sdf_i is the signed
 065 distance function of the human in the i^{th} frame and v is one
 066 of 2K points on the object's surface. We report the percent-
 067 age of trajectories with penetration distance $\leq 2cm$, ignor-
 068 ing trajectories with distance to object $> 2cm$, since tra-
 069 jectories that do not approach the object will trivially avoid
 070 penetration.

071 3. Details of User Study

072 In this section, we provide more details of our user study.
 073 We use the WenJuanxing [1] website to design and collect
 074 our questionnaires. We show our designed user interface,
 075 where users should rate each HOI animation as shown in
 076 Fig. 1. We invite 40 participants from varied backgrounds,
 077 comprising 22 students, 3 salespeople, 6 software engi-
 078 neers, 2 teachers, 3 managers, and 4 individuals from other
 079 professional fields. Among all participants, 65% are male,
 080 and a significant majority of 79% fell within the age range
 081 082 083 084 085 086 087



Text-driven human and object interaction animation quality assessment questionnaire

This is a quality assessment of the text-driven generation of interactive animations between people and objects. We will give 7 groups of animations, each with 4 animations related to text. We hope you can give a score based on the following criteria.

Semantic Matching: The generated animations match the semantics of the given text descriptions. [0-5]

Interaction Score: The quality of the poses and interactions between humans and objects in the animation. [0-5]

Realism: The level of realism in the motion of the characters. [0-5]

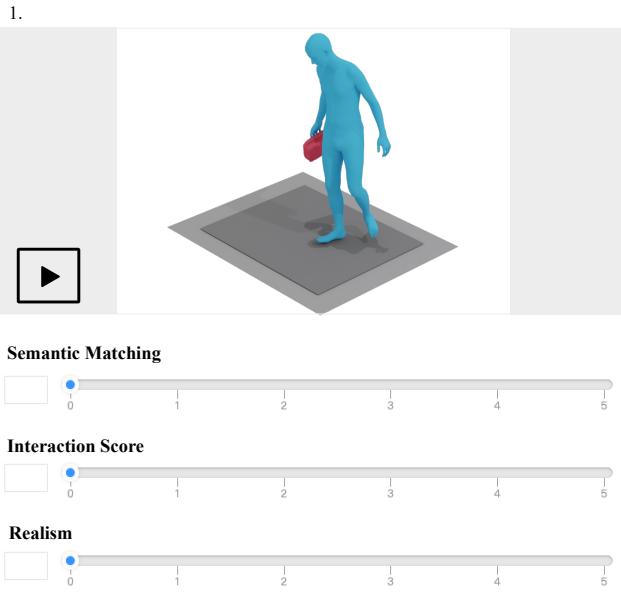


Figure 1. **User interface in our user study.** We ask users to rate the synthetic animations on three aspects: Semantic Matching, Interaction Score, and Realism.

088 of 18 to 25 years.

4. Additional Experiments

4.1. PMP Modules Analysis

091 In this section, we explore the impact of the Object Passage and Dual Flow modules within our network's architecture. To assess modules' contributions to PMP, we conducted a series of comparative experiments, which included removing or altering the positions of modules. The results, presented in Table 1, indicate that placing the Object Passage module at the beginning of the network significantly improves performance, increasing precision from 0.699 to 0.771. In contrast, positioning the Dual Flow module towards the end of the network is more beneficial, as demonstrated by an increase in precision from 0.768 to 0.772,

compared to its initial placement. Our method ensures text congruence and enhances the overall quality of the generated HOI animations.

Methods	Precision↑	FID ↓	Penetration↑
Real motions	0.821 ± 0.005	0.012 ± 0.002	—
— DF	0.772 ± 0.007	0.631 ± 0.027	0.623 ± 0.003
DF —	0.768 ± 0.003	0.632 ± 0.027	0.622 ± 0.003
— OP	0.699 ± 0.004	0.634 ± 0.047	0.621 ± 0.002
OP —	0.771 ± 0.004	0.632 ± 0.032	0.627 ± 0.005
DF OP	0.778 ± 0.007	0.625 ± 0.052	0.638 ± 0.003
Ours	0.781 ± 0.005	0.623 ± 0.063	0.643 ± 0.001

Table 1. **PMP module analysis.** We adjust the position of Object Passage (OP) and Dual Flow (DF) or remove a certain part. Our configuration can achieve the best results. Here, ‘|’ indicates the order. ‘—’ means do not use any operation.

4.2. Importance of Object Motion Guided Human Motion

In this section, we demonstrate the symbiotic guidance relationship between humans and objects within HOIAnimator. We dynamically influence the features of objects by interchanging and comparing them with the motion features of humans in the Object Passage. As shown in Table 2, network performance is significantly enhanced when object features guide human motion features. Specifically, precision increases from 0.772 to 0.781, the FID improves from 0.689 to 0.623, and the Penetration metric rises from 0.631 to 0.643. The results evidence suggests that object motion features are more effective in guiding human motion features. We attribute results to the more sensitive nature of object location features. Conversely, using human features rich in information content to guide object features tends to be counterproductive. We produce high-quality HOI animations by steering human motion features using the nuanced aspects of object motion features.

Methods	Precision↑	FID ↓	Penetration↑
Real motions	0.821 ± 0.005	0.012 ± 0.002	—
Human Object	0.772 ± 0.005	0.689 ± 0.027	0.631 ± 0.002
Object Human	0.781 ± 0.005	0.623 ± 0.063	0.643 ± 0.001

Table 2. **Object passage.** We adjust the order of human and object centric diffusion model passage.

5. More Results

In this section, we showcase an expanded range of our experimental results. The results encompass three key areas: firstly, the generation of diverse HOI animations from tex-

102
103
104

105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123

124
125
126
127

128 tual descriptions; secondly, the interaction of various motions
129 with the single object label; and thirdly, the exchange
130 of diverse objects with the same motion label. We demon-
131 strate the versatility and adaptability of our approach in cre-
132 ating varied and dynamic HOI animations.

133 **Text to diverse HOI animations.** Our HOIAnimator
134 can generate diverse object interaction animations from
135 simple text descriptions. When the text prompt is “A person
136 moves a stool with his hands”, HOIAnimator can generate
137 diverse animations of the stool moving in diverse directions,
138 as shown in Fig 2.

139 **Interaction of diverse objects with the same motion.**
140 We demonstrate HOIAnimator’s capability to generate a
141 range of animations when the motion label prompt is “Hold-
142 ing”. We show how it adeptly depicts holding diverse ob-
143 jects, such as a yoga ball, a backpack, and a box, as shown
144 in Fig. 3. Our HOIAnimator can generate diverse object
145 interaction motions based on the same motion label.

146 **Interaction of diverse motions with the same object.**
147 We demonstrate how our HOIAnimator, using only the sim-
148 ple prompt “wooden chair”, can generate a range of anima-
149 tions that interact with the wooden chair. The HOI anima-
150 tions encompass diverse actions like moving, lifting, and
151 sitting, as depicted in Fig. 4. The results showcase the
152 HOIAnimator’s capability to create diverse interactive mo-
153 tions from a singular object label.

154 6. Network Architecture

155 In order to enable our method to be successfully repro-
156 duced, we elaborate our HOIAnimator network structure in
157 Table 3.

158 References

- 159 [1] Wenjuanxing. <https://www.wjx.cn/>. 1
- 160 [2] Chuan Guo, Shihao Zou, Xinxin Zuo, Sen Wang, Wei Ji,
161 Xingyu Li, and Li Cheng. Generating diverse and natural 3d
162 human motions from text. In *Computer Vision and Pattern
163 Recognition*, pages 5152–5161, 2022. 1
- 164 [3] Nilesh Kulkarni, Davis Rempe, Kyle Genova, Abhijit Kundu,
165 Justin Johnson, David Fouhey, and Leonidas Guibas. Nifty:
166 Neural object interaction fields for guided human motion syn-
167 thesis. *arXiv preprint arXiv:2307.07511*, 2023. 1
- 168 [4] Mathis Petrovich, Michael J Black, and Gü̈l Varol. Action-
169 conditioned 3d human motion synthesis with transformer
170 vae. In *International Conference on Computer Vision*, pages
171 10985–10995, 2021. 1

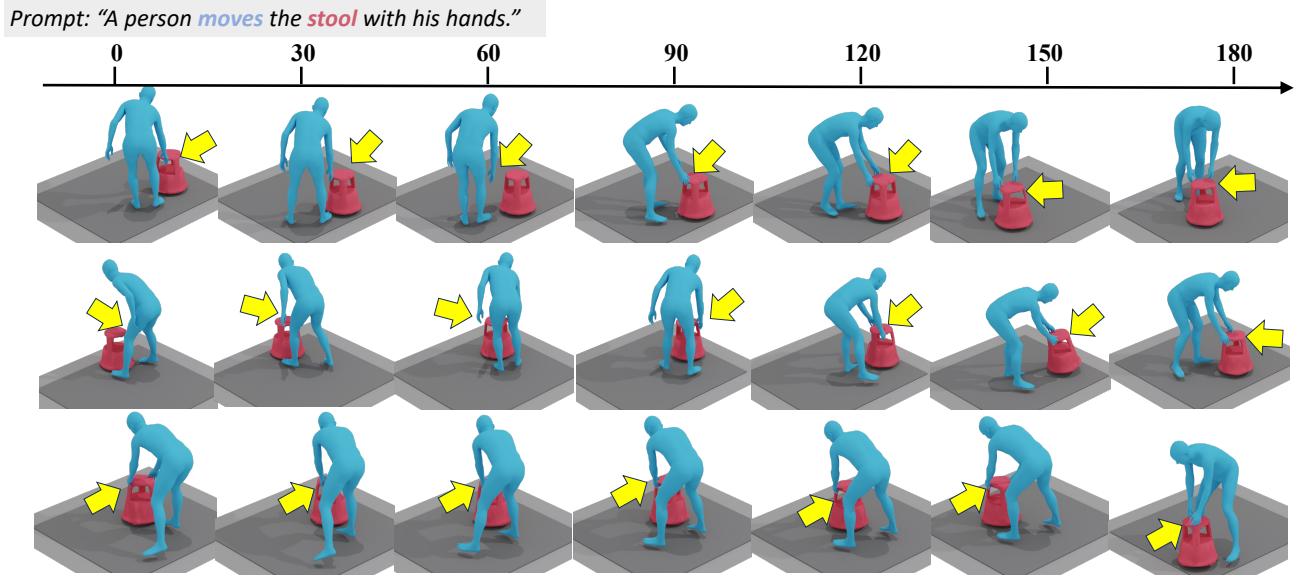


Figure 2. **Text to diverse HOI animations.** HOIAnimator can produce highly consistent HOI animations that align seamlessly with the given text, featuring rational interactions with high diversity. The yellow arrows indicate interaction areas.

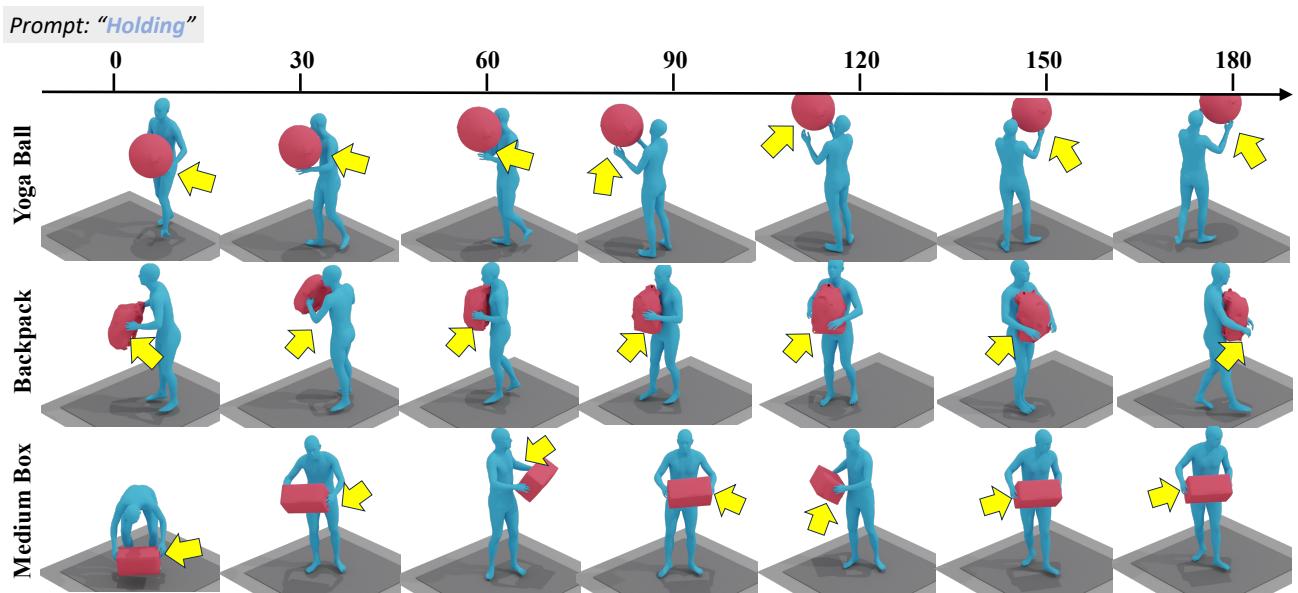


Figure 3. **Interaction of various objects with the same motion.** HOIAnimator can generate diverse object interaction motions based on the same motion label.

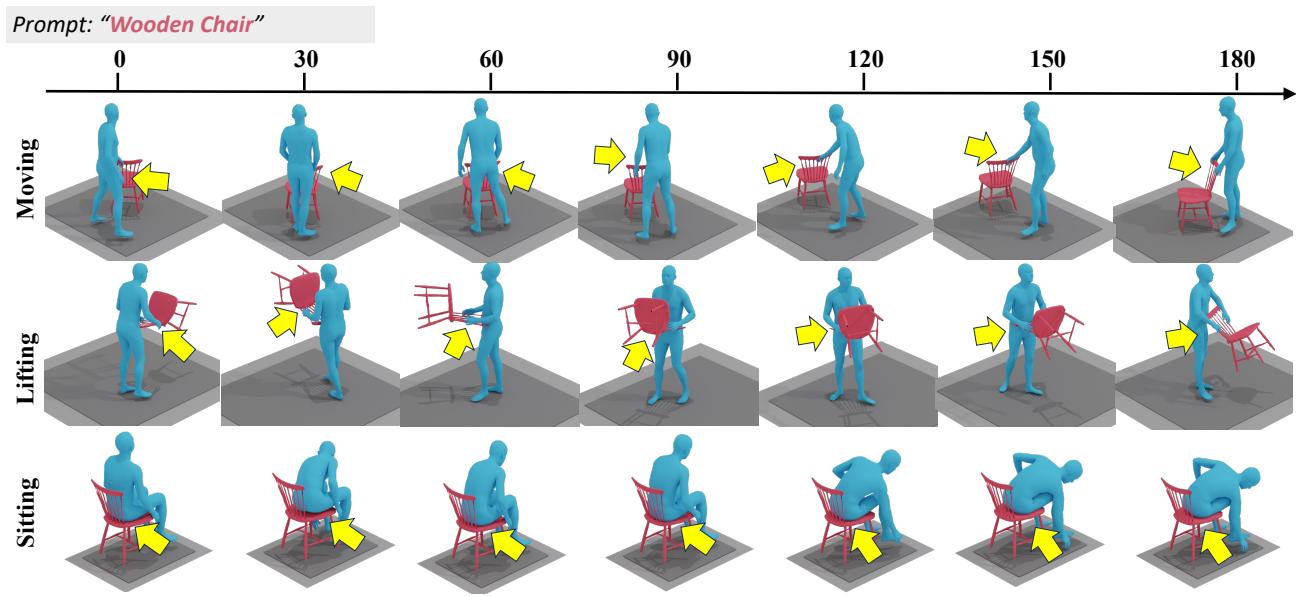


Figure 4. **Interaction of diverse motions with the same object.** HOIAnimator can generate diverse interactive motions based on the same object label.

Text Encoder	Frozen CLIP ViT-B/32 TransformerEncoderLayer(d_model=256, num_heads=4, dim_feedforward=1024) \times 2
Times Encoder	Linear(in_features=1000, out_features=512) Mish()
HOI Encoder	Linear(in_features=165, out_features=1024)
Object Passage	Linear(in_features=1024, out_features=1024) Linear(in_features=1024, out_features=1024)
Dual Flow	TransformerEncoderLayer(d_model=256, num_heads=4, dim_feedforward=2048) \times 2 Linear(in_features=2048, out_features=1024)
Latent Encoder	Linear(in_features=512, out_features=512) Linear(in_features=1024, out_features=1024, bias=True) (Vertice) LeakyReLU(negative_slope=0.2, inplace=True) Linear(in_features=7, out_features=1024, bias=True) (Emotion) LeakyReLU(negative_slope=0.2, inplace=True) Conv1d(in_channels=1024, out_channels=1024, kernel_size=5, stride=1, padding=2, padding_mode='replicate') LeakyReLU(negative_slope=0.2, inplace=True) InstanceNorm1d(num_features=1024) Linear(in_features=1024, out_features=1024, bias=True) Transformer(in_size=1024, hidden_size=1024, num_hidden_layers=6, num_attention_heads=8, intermediate_size=1536) Linear(in_features=1024, out_features=1024, bias=True)
TransformerEncoder	Linear(in_features=512, out_features=165)

Table 3. **Architecture of our method.** We provide detailed network architecture of our key components, including HOI Encoder, Object Passage, and Dual Flow.