036

037

038

039

044

045

046

047

048

049

050

051

052

053

054

055

056

057

058

059

060

061

062

063

064

065

066

067

068

069

070

071

HOIGen-1M: A Large-scale Dataset for Human-Object Interaction Video Generation

Supplementary Material

001 1. Overview

002 In this supplementary file, we provide more implementation details about the dataset curation and the proposed caption 003 method in Section 2. We also present the human evaluation 004 on massive generated videos to verify the proposed metrics' 005 006 alignment with human perception in Section 3. Additionally, 007 in Section 4, we discuss the limitations and societal impact of 008 our dataset. Finally, we offer more visual results in Section 5 to support the conclusion and insights in the main paper. 009

2. Implementation Details

011 2.1. Dataset Curation

In this section, we present more details of all process in the
dataset construction. Table 1 presents our data processing
pipeline, explaining how we identify one million high-quality
HOI videos from an initial collection of eight million.

016Data Collection. Initially, we select three public datasets017designed for HOI perception: BEHAVE [1], InterCap [5],018and HOI4D [12], which have high resolution, and each clip019showcases at least one HOI instance. However, these datasets020only contain 22K valid videos, which is far from sufficient021to support the training of large video generation models.

To collect more HOI videos, we expand our dataset se-022 lection from HOI perception data to general videos. We 023 incorporate five large-scale datasets for T2V generation: 024 Panda-70M [2], ViSR [11], and Mixkit, Pixabay, and Pexels 025 026 introduced in OpenSoraPlan [8]. By this means, we obtain about eight million raw videos. However, these videos are 027 of varying quality and most of them contain non-HOI con-028 029 tent. Therefore, to filter high-quality HOI videos from such a vast corpus, we adopt a data processing pipeline that is 030 031 introduced as follows.

032Metadata. We initially analyze video metadata, including033duration, resolution, and frames per second (FPS). Videos034are filtered to retain those exceeding one second in duration,035with a resolution \geq 720p and a frame rate \geq 20 FPS.

Optical Character Recognition. The DBNet++ [9] model is employed for optical character recognition to exclude video samples with over one text, which are considered as noises for the training of video generation models.

Aesthetics Score. Video aesthetic scores are a key factor
 in determining the visual quality of videos. To ensure high quality videos, we utilize the Laion Aesthetic Predictor ¹ to
 assess the visual appeal of videos. Only those videos that

PipelineToolOptical Character Recognition
Aesthetics ScoreDBNet++ [9]
LAION Aesthetics Predictor
UniMatch [16]Human Body Detection
Cartoon Style Detection
LLM EvaluationRTMPose Model [6]
ResNet-50 Binary Classifier [4]
PLLaVA [17] and Qwen2.5 [14]

Volunteers

Table 1. The pipeline of constructing HOIGen-1M.

achieve high aesthetic scores are retained.

Human Verification

Motion Score. To further ensure video quality, we enhance video quality by computing smooth motion. Specially, we employ UniMatch [16] to calculate the optical flow scores of the videos. Videos with excessively high or low optical flow scores are excluded, while those with moderate scores are retained.

Human Body Detection. After the previous processing steps, we have obtained high-quality video data, next we introduce how to filter the data for human-object interaction. The goal of this step is to exclude videos that lack human body and show bodies that are too small. We first evenly sampled five frames from each video and then employed RTMPose model [6] to identify the whole body 2D keypoints (133 keypoints) of human figures. Furthermore, we utilize the detected results to classified the whole dataset into three HOI categories: close-up hand interactions, single-person, and multi-person.

Cartoon Style Detection. Despite previous filtering, we observed some cartoon-style video clips remaining in the dataset. These clips, being fictional, often depict interactions that defy real-world physical laws. To exclude such videos, we curated a dataset comprising 100 cartoon-style videos and 125 realistic videos, segmented them into frames, and employed these to train a ResNet-50 [4] as a binary classifier. We then extracted two frames from each video and discarded any video where both frames were recognized as cartoon-style.

LLM Evaluation. To rapidly identify videos containing 072 HOI from a vast dataset, we initially employ PLLaVA [17] to 073 annotate videos and produce detailed captions. Subsequently, 074 we retain captions that include terms indicative of human 075 presence to filter out the videos without human. Then, we 076 ask the Qwen2.5 [14] whether there is interaction in these 077 captions using a predefined question and require the LLM 078 to explain the reason behind its answer. The advantage of 079

¹https://github.com/christophschuhmann/improved-aesthetic-predictor

determining whether there is HOI from captions is that captions are an essential part of the generation process, and text
tokens are fewer than video tokens, thereby reducing the
overall computational cost.

Human Verification. Ultimately, to ensure the high quality of the dataset, we enlisted seven volunteers to manually
check the remaining videos that showcases HOI. Specially,
volunteers were asked to scan videos to remove these featuring inconspicuous human-object interactions and minimal
object visibility.

090 2.2. Video Captioning with MoME

We design the Mixture-of-Multimodal-Experts (MoME) to
eliminate the hallucination by individual MLLM. In summary, MoME first adopts two captions and one decision
expert to detect the hallucination. Then, an additional set of
decision experts and caption experts is introduced to eliminate these hallucinations.

097The complete prompt for three caption experts (e.g.,098PLLaVa [17], Qwen2-VL [15]), and LLaVA [10] is shown099as followings:

Please describe the content of this video in as much detail as possible, including the people, objects, colors, scenery, environment, and camera movements within the video. Focus on describing the interaction between people and objects in the video, such as what actions people take, changes in the position or shape of objects, etc. Please describe the content of the video and the changes that occur, in chronological order. The description should be useful for AI to re-generate the video. The description should not be less than six sentences.

100 101 We adopt Llama3.1 [3] to act as a decision expert, where the prompt for determining hallucinations is as follows:

You are an AI assistant to check whether the given descriptions are about the same video. You will be provided with a description $\langle A \rangle$ and another one $\langle B \rangle$.

Please pay attention to the objects, attributes, and other relationships mentioned in these descriptions. Finally, make a judgment on whether these two descriptions are from the same video. Please ignore some of the very detailed differences. Please answer Yes or No first, and then explain the reason.

For example:

No, The main object in description A is an apple, but the object in description B is a pear.

Yes, The main objects in these two descriptions are consistent, both depicting scenes where an apple is placed on a table. If the hallucination exists, the prompt used for the decision expert to eliminate the hallucination is as follows: 103

You are an AI assistant and need to choose a description of a video. You will be given a description $\langle A \rangle$, a description $\langle B \rangle$, a description $\langle C \rangle$ and a difference $\langle D \rangle$.

Descriptions A, B, and C are detailed descriptions of the same video. Difference D explains the differences between the objects, properties, and other relationships mentioned in description A and descrip tion B. Description C gives a more accurate description of the objects in the video. Your task is to choose from Description A and Description B the one that more closely matches the description of the object and property in Description C.

You only need to answer "A" or "B". If descriptions A, B, and C are all completely different (ignoring very detailed points), answer "None".

If there is no hallucination, the prompt used for decision expert to choose a better caption is as follows:

You are an AI assistant that needs to pick the description of one video. You will be given a description $\langle A \rangle$ and a description $\langle B \rangle$.

Descriptions A and B are detailed descriptions of the same video. Your task is to select a description from description A and description B. The requirement is for a more detailed and accurate description of the video, e.g., a more detailed and specific description of people, objects, scenes, colors, environments, state changes, human actions, and human-object interactions.

You only need to answer "A" or "B".

3. Human Evaluation on Generated Videos

To validate the proposed evaluation metrics' alignment with human perception, we conduct human preference annotations on massive generated videos.

Data Preparation. Given a prompt p_i , and six models 110 to be evaluated {A,B,C,A',B',C'}, including three open-111 sourced T2V models and corresponding models fine-tuned 112 on our dataset, we employ each model to produce a video, 113 constructing a set of videos $S_i = \{V_{i,A}, V_{i,B}, V_{i,C}, V_{i,A'}, V_{i,A'}\}$ 114 $V_{i,B'}, V_{i,C'}$. For every prompt p_i , we create pairs of videos 115 in pairwise combinations, resulting in three pairs: $(V_{i,A},$ 116 $V_{i,A'}$), $(V_{i,B}, V_{i,B'})$, $(V_{i,C}, V_{i,C'})$, and ask human annotators 117 to judge their preferred video for each pair. Within the our 118 proposed framework, a prompt suite of N prompts produces 119 N×3 pairwise video comparisons. The sequence of videos in 120 each pair is shuffled to guarantee unbiased labeling. 121

104 105

106

107

108

109

183

196

Methods	Coarse	Fine	Win
	HOIScore	HOIScore	Ratio
OpenSora [19]	31.86%	91.86%	30.32%
Fine-tuned	35.38%	94.91%	69.68%
CogVideoX-2B [18]	31.34%	92.06%	35.23%
Fine-tuned	39.13%	93.82%	64.77%
CogVideoX-5B [18]	32.84%	94.25%	32.86%
Fine-tuned	44.04%	96.04%	67.14%

Table 2. The result of human preference annotation.

123

122 Human Annotation Rules. Specifically, human annotators are instructed to focus solely on the particular evaluation 124 dimension of interest and choose the video they prefer. For 125 instance, regarding the CoarseHOIScore, the question is "Does the video exhibit the people, objects, and correspond-126 ing interaction in the prompts?". Annotators are asked to 127 concentrate only on whether the generated video's content 128 129 aligns with the interaction, disregarding other quality factors 130 such as background consistency. For each metric, we meticulously prepare guidelines and train the human evaluators 131 to comprehend the definition of the metric. For example, 132 CoarseHOIScore pay more attention on the rough exist of 133 the interaction while FineHOIScore concentrates on the de-134 tails of interaction, such as the contact distance between 135 the humans and objects. We also conduct two rounds of 136 post-labeling checks to ensure the annotation quality. 137

Human Annotation Results. After performing a large-138 139 scale human annotation, we calculate the win ratio between the original model and the fine-tuend models. In pairwise 140 comparisons, if a model's video is chosen as superior, the 141 model receives a score of 1, while the other model receives a 142 score of 0. In the event of a tie, both models earn a score of 143 144 0.5. The win ratio for each model is determined by dividing the total score by the total number of pairwise comparisons 145 in which it participated. We present the result of human 146 preference annotations (e.g., win ratio) and in Table 2. We 147 can see that all fine-tuned models achieve a higher win ratio, 148 which is consistent with achieving a higher HOI score. The 149 above observation validate that our evaluation metrics can 150 faithfully reflect human perception. 151

4. Discussion 152

4.1. Data Availability Statement 153

154 We are committed to maintaining transparency and compliance in our data collection and sharing methods. According 155 156 to this principle, we have adhered to the following rules:

157 Public Data Sources: The data utilized in our work is pub-

licly available. We do not rely on any exclusive or confiden-158 tial data sources. 159 Data Distribution Policy: Our policy on data distribution 160 builds upon the precedent set by previous works like Kinetics, 161 InternVid, and others. Rather than offering the original raw 162 data, we only provide the YouTube video IDs necessary for 163 downloading the relevant content. 164 Usage Rights: The data released by us is solely for academic 165 and research purposes. This agreement does not authorize 166 any commercial use. 167 Compliance with YouTube Policies: Our methods for col-168 lecting and disseminating data are fully compliant with 169 YouTube's privacy policies. We ensure that no personal data 170 or privacy rights are compromised during the process. 171 Data Licensing: We employ the protocol of CC BY 4.0. 172 4.2. The Potential of HOIGen-1M 173

While HOIGen-1M is designed for video generation, it still 174 has great potential in HOI perception. For example, HOIGen-175 1M contains over 15,000 objects and more than 7,000 interac-176 tion action types, which significantly extends the categories 177 of existing HOI detection datasets. Consequently, it is a 178 golden testbed for HOI perception, such as open-vocabulary 179 HOI detection. Besides, This dataset also has the potential to 180 be utilized for training robotic arms to interact with objects, 181 facilitating the development of embodied AI. 182

4.3. Ethical Issues

There are three main ethical issues of this paper: 1) privacy, 184 2) data bias, and 3) human annotation bias. For the first 185 issue, the dataset will be distributed only for research pur-186 poses and we do not plan to provide the original raw data. 187 To eliminate data bias, we select the videos from multiple 188 sources, including YouTube, movies, sports competitions, 189 public area surveillance, etc. To reduce human annotation 190 bias, we first recruit annotators from diverse backgrounds, 191 including academia and industry, to ensure the diversity of 192 the annotation team. Moreover, we provide a clear annota-193 tion guide and conduct regular calibration tests to ensure 194 consistent understanding of the annotation standards. 195

5. Visual Results

In this section, more visual results are provided in order 197 to verify the effectiveness of HOIGen-1M. We show the 198 results of CogVideoX [18] and its fine-tuned version on 199 HOIGen-1M, as well as two commercial models: Gen3 [13] 200 and Kling1.5 [7]. For video comparisons of more models 201 please watch our supplementary video. 202

Dashed boxes of different colors in the below figures 203 indicate that the content does not match the corresponding 204 cue, and solid boxes indicate that the content matches. 205



CoarseHOIScore: 100% FineHOIScore: 98.25%

Figure 1. A person is opening a microwave door. Our generated video is the only one that generates the smooth interaction, which follows the physical law of the real world. Details in the prompt, such as the hand with a ring and the steaming soup, are also shown in full.

Figure 2. A man lifts and places a box in the trunk of the car. Our method generates the whole actions, including lifting and placing in, as well as detailed scenarios, such as a blanket and basket inside the trunk.

Figure 3. A person is running and flying a kite. Our approach successfully generates not only the interaction between the person and the kite, but also the wide shot as well as the action of the person raising his arms high.

Figure 4. A man is carrying a suitcase and walking on the street. Our method successfully follows the scene transitions in the prompt: "close-up of his hand gripping the suitcase handle", and eliminates the unreasonable interactions of the CogvideoX-5B.

Figure 5. A complex interaction: kickflip with a skateboard. Our generated video maintains the best consistency in skateboarding. In addition, compared to CogvideoX-2B, our video also enhances the details of human and environments.

Figure 6. A woman is twisting the orange off the stem and placing the orange into a woven basket. Our method performs well on twisting, placing, and shot switching.

Figure 7. A woman is smelling Pizza. Other methods do not generate the action correctly, generating eating pizza or showing pizza to the camera.

Figure 8. A young athlete is kicking a bright soccer. The video we generated not only correctly shows the full action of the kick, but also shows the precise scene transitions from the young athlete to the close-up of the soccer ball.

220

221

222

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

263

264

References 206

- 207 [1] Bharat Lal Bhatnagar, Xianghui Xie, Ilya A. Petrov, Cristian 208 Sminchisescu, Christian Theobalt, and Gerard Pons-Moll. 209 BEHAVE: dataset and method for tracking human object 210 interactions. In CVPR, pages 15914–15925. IEEE, 2022. 1
- 211 [2] Tsai-Shien Chen, Aliaksandr Siarohin, Willi Menapace, Eka-212 terina Deyneka, Hsiang-wei Chao, Byung Eun Jeon, Yuwei 213 Fang, Hsin-Ying Lee, Jian Ren, Ming-Hsuan Yang, and 214 Sergey Tulyakov. Panda-70m: Captioning 70m videos with multiple cross-modality teachers. In CVPR, pages 13320-215 216 13331. IEEE, 2024. 1
- [3] Abhimanyu Dubey, Abhinav Jauhri, and Abhinav Pandey et 217 218 al. The llama 3 herd of models. CoRR, abs/2407.21783, 2024. 219
 - [4] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In CVPR, pages 770-778. IEEE, 2016. 1
- 223 [5] Yinghao Huang, Omid Taheri, Michael J. Black, and Dim-224 itrios Tzionas. Intercap: Joint markerless 3d tracking of hu-225 mans and objects in interaction from multi-view RGB-D images. IJCV, 132(7):2551-2566, 2024. 1
 - [6] Tao Jiang, Peng Lu, Li Zhang, Ningsheng Ma, Rui Han, Chengqi Lyu, Yining Li, and Kai Chen. Rtmpose: Realtime multi-person pose estimation based on mmpose. CoRR, abs/2303.07399, 2023. 1
 - [7] Kwai. Keling. https://klingai.kuaishou.com, 2024.3
 - [8] PKU-Yuan Lab and Tuzhan AI etc. Open-sora-plan. https://github.com/PKU-YuanGroup/Open-Sora-Plan, 2024. 1
 - [9] Minghui Liao, Zhisheng Zou, Zhaoyi Wan, Cong Yao, and Xiang Bai. Real-time scene text detection with differentiable binarization and adaptive scale fusion. TPAMI, 45(1):919-931, 2023. 1
 - [10] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In NeurIPS, 2024. 2
- 242 [11] Xinchen Liu, Wu Liu, Meng Zhang, Jingwen Chen, Lianli 243 Gao, Chenggang Yan, and Tao Mei. Social relation recognition from videos via multi-scale spatial-temporal reasoning. 244 In CVPR, pages 3566-3574. IEEE, 2019. 1 245
- 246 [12] Yunze Liu, Yun Liu, Che Jiang, Kangbo Lyu, Weikang Wan, Hao Shen, Bogiang Liang, Zhoujie Fu, He Wang, and Li Yi. 247 248 HOI4D: A 4d egocentric dataset for category-level human-249 object interaction. In CVPR, pages 20981-20990. IEEE, 2022. 250
- 251 [13] Runway. Gen-3 alpha. https://runwayml.com, 2024. 252
- 253 [14] Qwen Team. Qwen2.5: A party of foundation models, 2024. 254
- 255 [15] Peng Wang, Shuai Bai, and Sinan Tan et al. Qwen2-vl: En-256 hancing vision-language model's perception of the world at 257 any resolution. CoRR, abs/2409.12191, 2024. 2
- 258 [16] Haofei Xu, Jing Zhang, Jianfei Cai, Hamid Rezatofighi, Fisher 259 Yu, Dacheng Tao, and Andreas Geiger. Unifying flow, stereo 260 and depth estimation. TPAMI, 45(11):13941-13958, 2023. 1
- 261 [17] Lin Xu, Yilin Zhao, Daquan Zhou, Zhijie Lin, See-Kiong 262 Ng, and Jiashi Feng. Pllava : Parameter-free llava extension

from images to videos for video dense captioning. CoRR, abs/2404.16994, 2024, 1, 2

- [18] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu 265 Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaohan 266 Zhang, Guanyu Feng, Da Yin, Xiaotao Gu, Yuxuan Zhang, 267 Weihan Wang, Yean Cheng, Ting Liu, Bin Xu, Yuxiao Dong, 268 and Jie Tang. Cogvideox: Text-to-video diffusion models 269 with an expert transformer. CoRR, abs/2408.06072, 2024. 3 270
- [19] Zangwei Zheng, Xiangyu Peng, Tianji Yang, Chenhui Shen, 271 Shenggui Li, Hongxin Liu, Yukun Zhou, Tianyi Li, and Yang 272 You. Open-sora: Democratizing efficient video production 273 for all. https://github.com/hpcaitech/Open-274 Sora, 2024. 3 275