# InterGen: Diffusion-based Multi-human Motion Generation under Complex Interactions

**Han Liang**[1] · **Wenqian Zhang**[1] · **Wenxuan Li**[1] · **Jingyi Yu**[1] · **Lan Xu**[1]

**Abstract** We have recently seen tremendous progress in diffusion advances for generating realistic human motions. Yet, they largely disregard the rich multi-human interactions. In this paper, we present InterGen, an effective diffusion-based approach that incorporates human-to-human interactions into the motion diffusion process, which enables layman users to customize high-quality two-person interaction motions, with only text guidance. We first contribute a multimodal dataset, named InterHuman. It consists of about 107M frames for diverse two-person interactions, with accurate skeletal motions and 16,756 natural language descriptions. For the algorithm side, we carefully tailor the motion diffusion model to our two-person interaction setting. To handle the symmetry of human identities during interactions, we propose two cooperative transformer-based denoisers that explicitly share weights, with a mutual attention mechanism to further connect the two denoising processes. Then, we propose a novel representation for motion input in our interaction diffusion model, which explicitly formulates the global relations between the two performers in the world frame. We further introduce two novel regularization terms to encode spatial relations, equipped with a corresponding damping scheme during the training of our interaction diffusion model. Extensive experiments validate the effectiveness and generalizability of InterGen. Notably, it can generate more diverse and compelling two-person motions than previous methods and enables various downstream applications for human interactions.

**Keywords** Motion synthesis · Multimodal generation · Diffusion model · Text-driven generation

## 1 Introduction

Digital human motions should reflect how we humans interact and communicate with each other, in order to depict the diverse cultures and societies that make up our physical world. A successful motion creation tool will hence allow users to customize realistic human motions under interactions. The produced motions also need to match specific themes, e.g., from as complicated as the movie script, or as simple as textual descriptions by novice users. Such human motion generation serves as a core computer vision problem, with various applications in VR/AR, games, or films.

Recent years have witnessed impressive progress in human motion generation under various user-specified conditioning, such as action categories (Guo et al., 2020; Petrovich et al., 2021), music pieces (Li et al., 2022, 2021), speeches (Habibie et al., 2022; Ao et al., 2022), or natural text prompts (Petrovich et al., 2022; Tevet et al., 2022b). The key idea is to learn a conditional generative model for the complex multimodal distribution of human motions, equipped with powerful neural techniques, from variational autoencoders (VAEs) (Kingma and Welling, 2013), generative adversarial networks (GANs) (Goodfellow et al., 2020),

Han Liang
E-mail: lianghan@shanghaitech.edu.cn

Wenqian Zhang
E-mail: zhangwq2022@shanghaitech.edu.cn

Wenxuan Li
E-mail: liwx2@shanghaitech.edu.cn

Jingyi Yu
E-mail: yujingyi@shanghaitech.edu.cn

Lan Xu
E-mail: xulan1@shanghaitech.edu.cn
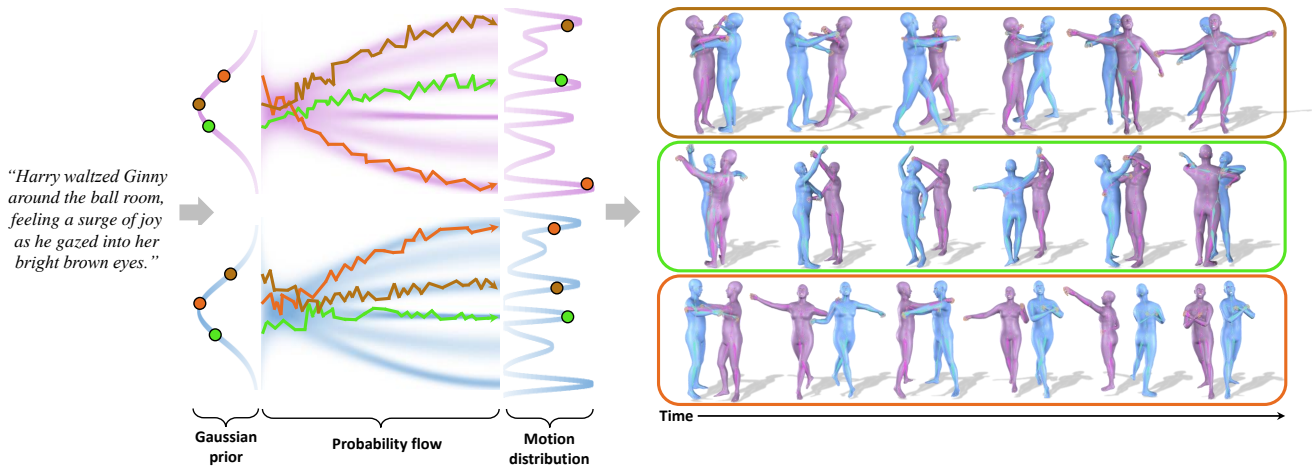[1] ShanghaiTech University, Shanghai, China

**Fig. 1** InterGen is the most advanced model for generating human interactive motion at present. **(Left)** We take text prompt as input, and a CLIP text encoder is used to process text into CLIP embedding. **(Middle)** We randomly sample noised motion, then denoise them in the standard diffusion denoising process. **(Right)** The final output denoised motion with high quality and diversity demonstrates the superiority of our approach.

normalization flows (Rezende and Mohamed, 2015), to the latest Diffusion Models (Ho et al., 2020; Song et al., 2020b). Only recently, the elaborate large-scale language models (Devlin et al., 2018; Brown et al., 2020) and diffusion methods (Ho et al., 2020; You et al., 2020) have quickly found their way into this field, due to their natural and convenient input controls and the strong ability to model complex distributions. These prompt-guided motion diffusion methods (Yuan et al., 2022; Zhang et al., 2022; Tevet et al., 2022b; Chen et al., 2022a) significantly democratize the accessible and high-quality generation of human motions for novices.

However, most of the above motion diffusion models are tailored for single-person setting, and hence overlook one essential aspect of human motions – the rich human-to-human interactions. The challenges are manifold. First, existing datasets Punnakkal et al. (2021); Liu et al. (2019) fail to simultaneously provide accurate captured results and natural prompt labels for diverse human interactions. The former usually relies on expensive dome-like capturing devices while the latter requires tedious and costly manual labeling. As a result, existing methods focus on generating the kinematic structure of a single human body, without exploring the diverse and complex spatial relationships between various human identities during interactions. The recent concurrent work Shafir et al. (2023) fine-tunes the single-person motion generator MDM Tevet et al. (2022b) into two-person scenario. Yet, it still suffers from unnatural interactions, inherently due to the limited interaction patterns in the single-person training datasets. In a nutshell, the lack of both multimodal datasets and corresponding explicit modeling schemes constitute barriers in two-person text-guided motion generation.

In this paper, we tackle the above challenges and present *InterGen* – an effective diffusion-based approach that enables layman users to customize high-quality two-person interaction motions with only text guidance, as illustrated in Fig. 1. Specifically, we first contribute a novel multimodal human motion dataset, named *InterHuman*, covering a wide range of two-person interactions, from daily ones like hugging, to professional motions, i.e., boxing or Latin dancing. In our InterHuman Dataset, we record dense video sequences with 76 RGB cameras, resulting in about 107 million video frames for 6,022 motion sequences that last for 6.56 hours. We then recover the ground-truth human skeletal motions under interactions from such rich RGB image modalities using off-the-shelf motion capture approach (He et al., 2021). Besides, we provide natural language labels for the captured motions, with 16,756 unique descriptions composed of 5,656 distinct words. Note that our InterHuman dataset is the first of its kind to open up the research direction for prompt-guided two-person motion generation under interaction setting. Its multi-modality also brings huge potential for future direction like multi-modal human interaction and behavior analysis.

Based on our InterHuman dataset, the key idea of our InterGen approach is to carefully bridge the general diffusion pipeline (Song et al., 2020b) into the human motion domain under two-person interactions. Specifically, we observe the symmetrical fact that exchanging the identities of performers during interactions does not change the semantics of motions. Thus, in our interaction diffusion model, we introduce two cooperative

transformer-style denoisers to correspondingly generate the motions of two performers. These denoisers explicitly share weights, with the aid of a novel mutual attention mechanism to further connects the two denoising processes at different feature levels. Such design encourages the two denoisers to perform the same operations and yield the same motion capacity, so as to avoid severe mode collapse when generating interaction motions, e.g., one person can dance professionally while the other cannot.

We further observe the widely adopted canonical representation for single-person motion (Guo et al., 2022; Tevet et al., 2022b; Zhang et al., 2022) discards the precise spatial relations in our interaction scenarios. Yet, naively adding the relative translation and rotation into the representation will lead to motion drifting during generation. We hence propose a non-canonical motion representation for our interaction diffusion model, where the relations between two people are explicitly encoded by global positions in the same world frame, facilitating the networks to learn relative relations. Besides, to generate more realistic two-person motions, we introduce two novel regularization terms to model the spatial relations during human-to-human interactions, including a masked joint distance map (DM) loss and relative orientation (RO) loss. The former DM loss encodes the spatial interference between two people with implicit physical constraints, while the latter RO loss encodes the orientation information since we humans pay more attention to our frontal orientation while interacting. We further adapt a damping scheme for these two losses during training, especially when the sampled timestamp of the diffusion process reaches specific thresholds, achieving a more diverse generation. Finally, we perform extensive experiments to demonstrate that our approach can generate more compelling two-person motions than previous methods, and showcase its various downstream applications for human interactions, i.e., trajectory control, interactive motion inbetweening, and person-to-person generation.

To summarize, our main contributions include:

- We contribute a new human interaction dataset with rich text/motion modalities, and present a novel diffusion-based approach to generate realistic two-person motions from only prompt inputs.
- In our interaction diffusion model, we introduce cooperative denoisers with novel weights-sharing and a mutual attention mechanism to significantly improve the generation quality.
- We propose an effective motion representation, as well as two additional regularization losses with a damping schedule to model the complex spatial relation under human-to-human interactions.

## 2 Related Work

### 2.1 Human Motion Generation

The field of Human Motion Generation is greatly facilitated by the integration of extensive multimodal data inputs, including text (Petrovich et al., 2022; Tevet et al., 2022b,a; Yuan et al., 2022; Chen et al., 2022a; Shafir et al., 2023; Guo et al., 2022; Kim et al., 2022), action (Guo et al., 2020; Petrovich et al., 2021), incomplete motion sequences (Duan et al., 2021; Harvey et al., 2020), control signals (Starke et al., 2022; Peng et al., 2021; Starke et al., 2019), music (Li et al., 2021, 2022; Lee et al., 2019), speech (Habibie et al., 2022; Ao et al., 2022), scene (Wang et al., 2021, 2022) and images (Rempe et al., 2021; Chen et al., 2022b).

Currently, there exists a variety of works focused on action label-based human motion generation (Guo et al., 2020; Petrovich et al., 2021; Song et al., 2022). As the emergence of large language models (LLMs) (Brown et al., 2020; Touvron et al., 2023) paves the way for LLM-based multimodal (OpenAI, 2023) models, the potential of text-based multimodal generation becomes increasingly apparent. However, action-based human motion generation is fundamentally a class-based generation approach, which can not enable high-level natural language and multimodal control. Consequently, it is crucial to explore human motion generation based on natural language input. Early approaches such as Ahn et al. (2018) employ a sequence-to-sequence model to generate upper body motion. Subsequently, Ahuja and Morency (2019) and Ghosh et al. (2021) concentrate on developing a unified language and pose representation to enable autoregressive synthesis of human movements. Additionally, Petrovich et al. (2022) and Guo et al. (2022) both adopt Variational Autoencoder (VAE) (Kingma and Welling, 2013) based architectures for motion generation.Multimodal pre-trained models (Radford et al., 2021) facilitate more seamless integration between textual and motion spaces (Tevet et al., 2022a).

Recent advancements in Diffusion Models (Ho et al., 2020; Song et al., 2020b) have significantly propelled text-driven motion generation. Kim et al. (2022) introduces a groundbreaking transformer-based architecture that effectively manages motion data, which is crucial for handling variable-length motions and attending to free-form text. Chen et al. (2022a) presents a Motion Latent-based Diffusion model (MLD) that generates plausible human motion sequences based on various conditional inputs, such as action classes or textual descriptors. Tevet et al. (2022b) proposes the Motion Diffusion Model (MDM), a meticulously adapted

classifier-free diffusion-based generative model for the human motion domain. Building on this work, Yuan et al. (2022) integrates physical constraints into the diffusion process to generate physically plausible human motions.

The aforementioned motion generation techniques primarily focus on single-person motion generation, lacking the capability to generate and model dual or multi-person human motions. Recent work by (Shafir et al., 2023) suggests that the gap in data availability for motion generation can be addressed using a pretrained diffusion-based model (Tevet et al., 2022b) as a generative prior and demonstrates the prior's effectiveness for fine-tuning in a few-shot manner. However, this approach is limited to generating interactions observed during training, resulting in constrained generalization and suboptimal generation quality capabilities.

### 2.2 Human Motion Capture

Motion capture techniques have been well-developed in the last decade. Marker-based techniques, such as Vicon (2019) and those presented by Vlasic et al. (2007), have been successful in capturing high-quality human motions for professional applications. However, these methods are not suitable for daily use due to their costly and laborious setup. To overcome this limitation, markerless motion capture methods have been developed (Bregler and Malik, 1998; De Aguiar et al., 2008; Theobalt et al., 2010). Advances in parametric human models (Anguelov et al., 2005; Loper et al., 2015; Pavlakos et al., 2019; Osman et al., 2020) have led to data-driven approaches for estimating 3D human pose and shape using optimization (Huang et al., 2017; Lassner et al., 2017; Bogo et al., 2016; Kolotouros et al., 2019) or direct regression (Kanazawa et al., 2019; Kocabas et al., 2020; Zanfir et al., 2021) of human model parameters. Template-based methods, utilizing specific template meshes as priors, have been proposed for both multi-view (Gall et al., 2010; Stoll et al., 2011; Liu et al., 2013; Robertini et al., 2016; Pavlakos et al., 2017; Simon et al., 2017; Xu et al., 2018a) and monocular (Xu et al., 2018b; Habermann et al., 2019; Xu et al., 2020; Habermann et al., 2020) setups. Another research line is inertial measurement units (IMUs). Commercially available systems, such as Xsens MVN (Movella, 2022), have employed large numbers of sensors, but the intrusive nature of these systems has prompted research into sparser sensor setups. Von Marcard et al. (2017) presented a pioneering exploration called SIP, which employs only six IMUs. However, the traditional optimization framework used in SIP hampers real-time application. Data-driven approaches (Huang et al., 2018;



**Fig. 2** Our motion capture studio **(Top)** and our collected InterHuman dataset illustration **(Bottom)**. The system comprises 76 calibrated multi-view cameras. InterHuman covers a wide range of two- person interactions.

Yi et al., 2021, 2022) utilizing sparse sensors have shown significant improvements in accuracy and efficiency, but substantial drift remains an issue for challenging motions. Previous sensor-aided solutions have combined IMUs with videos (Gilbert et al., 2019; Henschel et al., 2020; Malleson et al., 2019, 2017; Liang et al., 2022), RGB-D cameras (Helten et al., 2013; Zheng et al., 2018), optical markers (Andrews et al., 2016), or even LiDAR (Zhao et al., 2022) to address the scene-occlusion problem and effectively correct drift. These methods achieve highly accurate capture of human motions given various modalities of signals, however, the high-level control of synthesizing captured motion and even non-seen motions with more modalities of input remains challenges.

### 2.3 Human Motion Dataset

In recent times, the accessibility of extensive motion datasets has played a crucial role in propelling motion generation research forward. Action label datasets such as BABEL (Punnakkal et al., 2021) and NTU RGB+D 120 (Liu et al., 2019), although having labeled annotations, differ significantly from natural language when composed of verb or verb-object structures, and furthermore, the action labels only enable direct classification of actions, making it impossible to support text-based human motion generation. Datasets such as KIT (Plappert et al., 2016), and HumanML3D (Guo et al., 2022), which feature text annotations, have been especially valuable for the progression of text-driven

Two martial-art masters bow to each other and get ready to attack.

Two people lift something together and square up to fight.

Two people raise their hands in front of chests and raise the both fists, do martial-art squats.

Two people stand still and bend forward.

Two people stand face by face and bow to each other.

Two people stand looking at each other and lower their upper body.

One runs to the other and shakes the fists right and left, the other tries to catch opponent's arm.

One steps forward with a right fist followed by a left, and the other raises the left arm in defense.

One does a right uppercut next to a left uppercut, the other lowers the raised left arm.

Two people get closer and raises their right hand to shake hands.

Two people walk with right arms lifted and bump hands with stopped walking.

Two people walk towards each other and lift right fists then align their hands.

One punches left jab, the other throws a hard left that dodged by another man.

One is raising his left arm, the other extends their right arm to catch another.

Two were bumping fists and one raises left arm, the other waves right hard and retracted the body.

Two people hug and separate.

Two persons stand close, one raises arms to grab the other and the other one pushes away.

One person falls into the other's chest, the other lifts arms to another's back.

**Fig. 3** InterHuman dataset consists of diverse human professional and daily interactions with diverse natural language annotations from different annotators. The figure showcases two examples of our dataset, martial arts, and social manners, with thorough descriptions from different perspectives.

**Table 1 Dataset comparisons.** We compare our InterHuman dataset with existing human motion datasets. **Motions** refers to the total number of motion clips. **Vocab.** shows the number of distinct words used in the annotations, while **Descriptions** summarizes the total number of textual descriptions.

| Dataset | Natural Language | Interactive | Motions | Vocab. | Descriptions | Duration |
|---|---|---|---|---|---|---|
| KIT (Plappert et al., 2016) | ✓ | - | 3911 | 1623 | 6278 | 11.23h |
| HumanML3D (Guo et al., 2022) | ✓ | - | 14616 | 5371 | 44970 | 28.59h |
| BABEL (Punnakkal et al., 2021) | - | - | 13220 | - | - | 43.5h |
| NTU RGB+D 120 (Liu et al., 2019) | - | ✓ | 739 | - | - | 0.47h |
| UMPM (Van der Aa et al., 2011) | - | ✓ | 36 | - | - | 2.22h |
| You2Me (Ng et al., 2020) | - | ✓ | 42 | - | - | 1.4h |
| **InterHuman(Ours)** | ✓ | ✓ | **6022** | **5656** | **16756** | **6.56h** |

motion generation research (Chen et al., 2022a; Tevet et al., 2022b; Petrovich et al., 2022). However, these datasets only consist of single-person motions and annotations, making it difficult to apply and generalize to the generation of interactive motions involving two or more individuals.

Various multi-person motion datasets have been developed, including 3DPW (Von Marcard et al., 2018), You2Me (Ng et al., 2020), and UMPM (Van der Aa et al., 2011). However, while these datasets contain two-person and multi-person motion data, they are limited in size and annotations. In particular, there is a lack of textual or other modal annotations in these datasets. Efforts to annotate existing datasets with text, as demonstrated in the annotation of 3DPW (Von Marcard et al., 2018) in ComMDM (Shafir et al., 2023), establish a foundation for future advancements in text-guided multi-person motion generation. However, since it only contains 27 two-person motion sequences, the issue of limited availability of two-person interaction datasets still persists. Our proposed InterHuman

dataset is currently the largest interaction-language dataset, addressing the lack of suitable datasets in text-based human-to-human interactive motion generation research.

## 3 InterHuman Dataset

InterHuman[1] is a comprehensive, large-scale 3D human interactive motion dataset encompassing a diverse range of 3D motions of two interactive people, each accompanied by natural language annotations. To the best of our knowledge, it is the most extensive 3D human-to-human interaction dataset available. Unlike some previous datasets that focus only on single-person motion or particular actions, such as dancing (Li et al.,

---

[1] Our dataset consists of dense multi-view and motion information of various identities. Due to privacy concerns of performers included in our dataset, we cannot directly share it publicly. But it is available on reasonable request to the corresponding author at `https://github.com/tr3e/InterGen` (for research purposes only).

2021). Our dataset comprises various interaction motions, broadly classified into two categories: daily motion, which encompasses everyday routines involving two people (e.g., passing objects, greeting, communicating, etc.), and professional motions, which include typical human-to-human interactions (e.g., Taekwondo, Latin dance, boxing, etc.).

**Data collection.** We utilized a system with 76 calibrated Z-CAM (Z-cam, 2022) RGB cameras to capture human interactions, the Fig. 2 shows our system and datasets. In daily motion sessions, participants were given a guiding script, which they interpreted and performed. In professional motion sessions, participants were asked to demonstrate as many interactions as possible with their partners based on their past experiences, ensuring the dataset's diversity. Subsequently, we employed a data processing pipeline as same as Li et al. (2021) to obtain ground truth SMPL (Loper et al., 2015) parameters from the multi-view videos. We conduct the textual annotation process via Amazon Mechanical Turk (AMT). We ask the annotators to split a video into several clips up to 10 seconds and collect 3 text descriptions for each clip from distinct annotators. The annotation examples are shown in Fig. 3, where the detailed interactions between the two people are described from different perspectives.

**Dataset comparison.** Our InterHuman dataset was captured using a system with 76 cameras, and all motions have been meticulously annotated, consisting of 6,022 motions derived from various categories of human actions, labeled with 16,756 unique descriptions composed of 5,656 distinct words, with a total duration of 6.56 hours, which makes it the largest and most diverse known scripted dataset of human-to-human interactions. For specific durations of each motion category, please refer to Fig. 4. Tab. 1 provides a comprehensive comparison of our dataset with several existing human datasets from various perspectives, highlighting that our dataset is currently the most suitable for tasks involving human interactions.

## 4 InterGen Approach

Our goal is to generate diverse and high-quality human interaction motions conditioning on text prompts, within a diffusion-based framework. To this end, as illustrated in Fig. 5, our approach consists of three key technical designs. The first is an effective motion representation (Sec. 4.1) that preserves the spatial relations of interacting people in the common world frame. Then, we adopt a novel denoising architecture that involves two cooperative networks (Sec. 4.2) with
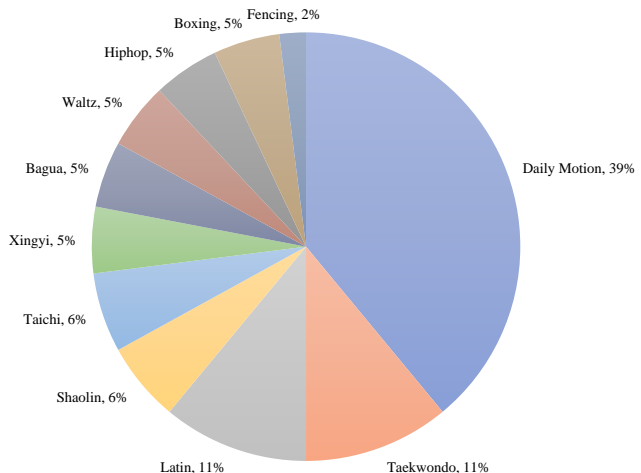


**Fig. 4** InterHuman dataset covers a wide range of two- person interactions, from the daily ones like hugging, handshake, and argument to the professional motions ranging from dance to martial arts.

sharing-weights and mutual attention to connect the two branches at hierarchical feature levels, so as to balance the motion capacity of two interacting performers. To train such interaction diffusion model, we further propose additional regularization terms (Sec. 4.3), consisting of a joint distance map (DM) loss and a relative orientation (RO) loss to enforce the networks to depend on each others, especially under continuous and close interactions. In addition, we propose a novel loss damping schedule during training to improvive the generation results.

### 4.1 Human Interaction Representation

We model a two-person interaction $\mathbf{x}$ as a collection of two single-person motion sequences $\mathbf{x}_h$, i.e., $\mathbf{x} = \{\mathbf{x}_a, \mathbf{x}_b\}$, where $\mathbf{x}_h = \{x^i\}_{i=1}^L$ is a fixed-framerate sequence of motion states $x^i$. Thus these two sequences are naturally synchronized. The core problem is to encode the spatial relationships between them.

**Canonical representation.** HumanML3D (Guo et al., 2022) proposed a human motion representation for single-person scenarios that incorporate ground contact information and motion features. This representation is over-parameterized, expressive, and neural network friendly, and has been adopted by several recent works. However, this representation cannot be directly applied to multi-person scenarios because it canonicalizes joint positions and velocities to the root frame, which loses global spatial information. ComMDM (Shafir et al., 2023) tries to mitigate this issue by predicting the initial relative rotation and translation between two people. In contrast, we extend this
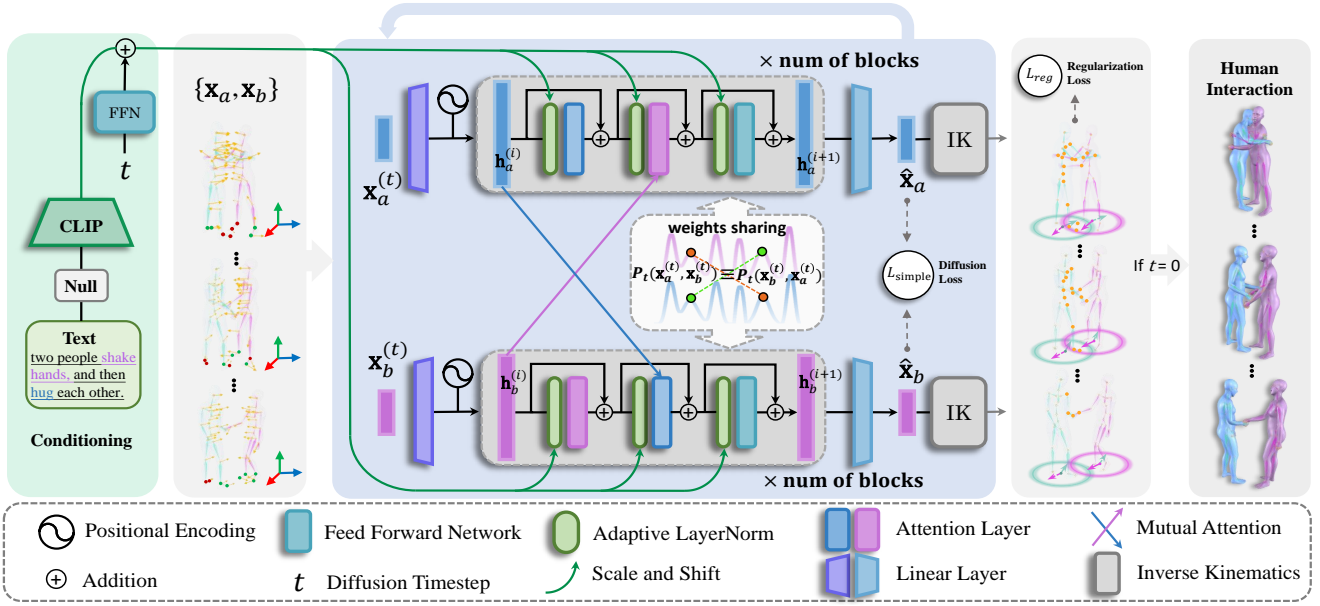
**Fig. 5 The overview of our InterGen.** we contribute three primary technical designs. First, we propose an efficient two-person interaction motion representation. Second, we introduce two cooperative transformer-style weights-sharing networks with mutual attention to interactively perform denoising. Lastly, we introduce an effective loss function that significantly improves the quality of two-person interaction generation.

representation by introducing global relative rotation $\mathbf{r}^h \in \mathbb{R}^2$ along Y-axis and translation $\mathbf{t}^h \in \mathbb{R}^2$ on XZ-plane over other humans.

$$x^i = [\dot{r}^a, \dot{r}^x, \dot{r}^z, r^y, \mathbf{j}_l^p, \mathbf{j}_l^v, \mathbf{j}^r, \mathbf{c}^f, \mathbf{r}^h, \mathbf{t}^h], \tag{1}$$

where the $i$-th motion state $x^i$ is defined as a collection of root angular velocity $\dot{r}^a \in \mathbb{R}$ along the Y-axis, root linear velocities $\dot{r}^x, \dot{r}^z \in \mathbb{R}$ on the XZ-plane, root height $r^y \in \mathbb{R}$, local joint positions $\mathbf{j}_l^p \in \mathbb{R}^{3N_j}$, local velocities $\mathbf{j}_l^v \in \mathbb{R}^{3N_j}$, rotations $\mathbf{j}^r \in \mathbb{R}^{6N_j}$ in root space, and binary foot-ground contact features $\mathbf{c}^f \in \mathbb{R}^4$ by thresholding the heel and toe joint velocities, where $N_j$ denotes the joint number.

**Non-canonical representation.** However, the canonical representation of motion suffers from drifts, as its global rotation and translation are obtained by accumulating local angular and linear velocities, which leads to error accumulation and produces unbounded exponential drifts (Von Marcard et al., 2017) over time. This absolute trajectory error can be ignored for tasks such as single-person short-sequence motion synthesis, which focuses on the plausibility of local motion and does not care much about the absolute trajectory. Whereas it is fatal for tasks such as multi-person motion synthesis, which require person-to-person precise spatial relationships.

Thus we propose a non-canonicalization representation for multi-person interaction motion. Instead of transforming joint positions and velocities to the root

frame, we keep them in the world frame, which allows us to directly access the global translation and rotation from the root position and inverse kinematics (IK), respectively, to avoid drift effectively. The representation is formulated as:

$$x^i = [\mathbf{j}_g^p, \mathbf{j}_g^v, \mathbf{j}^r, \mathbf{c}^f], \tag{2}$$

where the $i$-th motion state $x^i$ is defined as a collection of global joint positions $\mathbf{j}_g^p \in \mathbb{R}^{3N_j}$, velocities $\mathbf{j}_g^v \in \mathbb{R}^{3N_j}$ in the world frame, 6D representation of local rotations $\mathbf{j}^r \in \mathbb{R}^{6N_j}$ in the root frame, and binary foot-ground contact features $\mathbf{c}^f \in \mathbb{R}^4$.

### 4.2 Human Interaction Diffusion

**Diffusion models.** We sometimes omit the explicit dependence on condition c for simplicity of notation. Note that we can always train diffusion models with some condition $c$; even for the unconditional case, we can condition the model on a universal null token $\emptyset$.

Let $p_0(\mathbf{x})$ denote the human interactive motion data distribution, diffusion models Ho et al. (2020) injects time-dependent $i.i.d$ Gaussian noise to the samples from $p_0$, giving a diffusion process $\{p_t(\mathbf{x})\}_{t=0}^T$ with a continuous variable $t \in [0, T]$. Then a generative model can be obtained by reversing the process, starting from samples $\mathbf{x}^{(T)} \sim p_T$ that is a standard Gaussian distribution, and then solving the following reverse-time
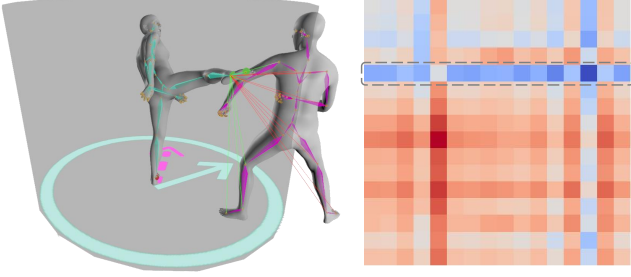
**Fig. 6** (**Left**) visualize our proposed interactive losses, where the relative orientation loss is the angular separation between the frontal orientations of the two people. And the partial joint distance map of the heel joint is truncated with the region of the cylinder, which is shown in (**Right**), where the highlighted row encodes the spatial relations between the heel and the other person.

SDE from $t = T$ to $t = 0$:

$$dx = [\mathbf{f}(\mathbf{x}, t) - \sigma_t^2 \nabla_\mathbf{x} log p_t(\mathbf{x})]dt + \sigma_t d\mathbf{w}, \tag{3}$$

where $\mathbf{f}(\cdot, t)$ is a deterministic drift function, $\sigma_t$ is the diffusion coefficient that is increasing over time to control the noise level, $dt$ is infinitesimal negative timestep, $d\mathbf{w}$ is infinitesimal noise, and $\nabla_\mathbf{x} log p_t(\mathbf{x})$ is the score function that is the only intractable term. Note that it can be obtained from the expectation of $\mathbf{x}$ given $\mathbf{x}^{(t)}$:

$$\nabla_{\mathbf{x}^{(t)}} log p_t(\mathbf{x}^{(t)}) = (\mathbb{E}[\mathbf{x}|\mathbf{x}^{(t)}] - \mathbf{x}^{(t)})/\sigma_t^2. \tag{4}$$

When we drop the noise term at Eq. (3), an ordinary differential equation (ODE) is obtained, which is a corresponding deterministic process sharing the same marginal distribution $\{p_t(\mathbf{x})\}_{t=0}^T$ and is referred to as probability flow ODE, which can accelerate sampling process by performing a linear interpolation between $\mathbf{x}^{(t)}$ and the predicted $\mathbb{E}[\mathbf{x}|\mathbf{x}^{(t)}]$ (Song et al., 2020a).

**Interaction Diffusion.** Our approach is based on a fundamental assumption, commutative property, which means that two-person interactions $\{\mathbf{x}_a, \mathbf{x}_b\}$ and $\{\mathbf{x}_b, \mathbf{x}_a\}$ are equivalent, i.e., the order of every single motion does not change the semantics of the interaction itself. In other words, the distribution of interaction data satisfies the following property:

$$p(\mathbf{x}_a, \mathbf{x}_b) \equiv p(\mathbf{x}_b, \mathbf{x}_a). \tag{5}$$

Under this assumption, the two people share the same single-person motion marginal distribution. Since noise $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ is independent of data distribution, we have:

$$p_t(\mathbf{x}_a^{(t)}, \mathbf{x}_b^{(t)}) \equiv p_t(\mathbf{x}_b^{(t)}, \mathbf{x}_a^{(t)}). \tag{6}$$

Thus based on the above conclusion the score function can be reformulated as:

$$\nabla_{\mathbf{x}^{(t)}} log p_t(\mathbf{x}^{(t)})$$
$$= [\nabla_{\mathbf{x}_a^{(t)}} log p_t(\mathbf{x}_a^{(t)}, \mathbf{x}_b^{(t)}), \nabla_{\mathbf{x}_b^{(t)}} log p_t(\mathbf{x}_a^{(t)}, \mathbf{x}_b^{(t)})]$$
$$= [\nabla_{\mathbf{x}_a^{(t)}} log p_t(\mathbf{x}_a^{(t)}, \mathbf{x}_b^{(t)}), \nabla_{\mathbf{x}_b^{(t)}} log p_t(\mathbf{x}_b^{(t)}, \mathbf{x}_a^{(t)})], \tag{7}$$

where the two parts for $\mathbf{x}_a$ and $\mathbf{x}_b$ are the same function $\nabla_a log p_t(a, b)$, which can be approximated by employing the same network with the following denoising autoencoder objective:

$$\mathcal{L}_{simple} = \mathbb{E}_{\mathbf{x}, t, \epsilon}[\lambda_t ||\mathbf{x}_a - D_\theta(\mathbf{x}_a + \sigma_t \epsilon_a, \mathbf{x}_b + \sigma_t \epsilon_b, t, c)||_2^2$$
$$+ \lambda_t ||\mathbf{x}_b - D_\theta(\mathbf{x}_b + \sigma_t \epsilon_b, \mathbf{x}_a + \sigma_t \epsilon_a, t, c)||_2^2], \tag{8}$$

where $D_\theta$ is the denoisers sharing the common network weights, whose input consists of its own noisy motion to denoise, the cooperator's noisy motion, the time $t$, and the condition $c$, noise $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, and $\lambda_t$ is the loss weighting factor.

**Cooperative denoisers.** Based on the aforementioned conclusions, we adopt interactively cooperative transformer-style networks sharing the common weights to model $D_\theta$, as demonstrated in Fig. 5. The networks are fed their own noisy motions, $\mathbf{x}_a^{(t)}$ and $\mathbf{x}_b^{(t)}$, as inputs for denoising and subsequently output the corresponding denoised versions of the motions, $\mathbf{x}_a$ and $\mathbf{x}_b$. This prediction process is conditioned on the diffusion timestep $t$, control condition $c$, and the hidden states $\mathbf{h}^{(i)}$ of the counterpart network.

Specifically, the noisy motion is first embedded into a common latent space and positionally encoded into an internal representation often referred to as the hidden states $\mathbf{h}^{(0)}$. Then, it is processed by $N$ attention-based blocks to obtain denoised hidden states $\mathbf{h}^{(N)}$. Finally, a common inverse embedding layer is applied to output the denoised motion.

Each block consists of two multi-head attention layers ($Attn$) followed by one feed-forward network ($FF$). The first attention layer is a self-attention layer, which embeds the current hidden states $\mathbf{h}^{(i)}$ into a context vector $\mathbf{c}^{(i)}$. The computation of $\mathbf{c}_a^{(i)}$ part is formulated as the following:

$$\mathbf{c}_a^{(i)} = Attn(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = softmax(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{C}})\mathbf{V},$$
$$\mathbf{Q} = \mathbf{h}_a^{(i)}\mathbf{W}_s^Q, \mathbf{K} = \mathbf{h}_a^{(i)}\mathbf{W}_s^K, \mathbf{V} = \mathbf{h}_a^{(i)}\mathbf{W}_s^V, \tag{9}$$

where $C$ is the number of channels in the attention layer and $\mathbf{W}_s$ are trainable weights, and $\mathbf{c}_b^{(i)}$ is calculated with shared weights $\mathbf{W}_s$ in the same way.

The second attention layer is a mutual attention layer, where the key $\mathbf{K}$ value $\mathbf{V}$ pair is provided by the

hidden states $\mathbf{h}^{(i)}$ of the counterpart block. Here the computation process of the next hidden states $\mathbf{h}^{(i+1)}$ is formulated as:

$$\mathbf{h}_a^{(i+1)} = FF(Attn(\mathbf{Q}_a, \mathbf{K}_b, \mathbf{V}_b)),$$
$$\mathbf{h}_b^{(i+1)} = FF(Attn(\mathbf{Q}_b, \mathbf{K}_a, \mathbf{V}_a)),$$
$$\mathbf{Q}_a = \mathbf{c}_a^{(i)}\mathbf{W}_m^Q, \mathbf{K}_a = \mathbf{h}_a^{(i)}\mathbf{W}_m^K, \mathbf{V}_a = \mathbf{h}_a^{(i)}\mathbf{W}_m^V,$$
$$\mathbf{Q}_b = \mathbf{c}_b^{(i)}\mathbf{W}_m^Q, \mathbf{K}_b = \mathbf{h}_b^{(i)}\mathbf{W}_m^K, \mathbf{V}_b = \mathbf{h}_a^{(i)}\mathbf{W}_m^V, \quad (10)$$

where $\mathbf{W}_m$ are trainable weights shared by the two branches.

In addition, the adaptive layer normalization is employed before all attention layers and the feed-forward network to condition on the control condition $c$ and timestep $t$.

### 4.3 Additional Regularization Losses

**Geometric losses.** We adopt the common geometric losses in the field of human motion, such as foot contact loss $\mathcal{L}_{foot}$ and joint velocity loss $\mathcal{L}_{vel}$, to regularize the generative models and enforce physical plausibility and coherence for each single-person motion. For more details, we refer the reader to MDM (Tevet et al., 2022b). In addition, for our non-canonical representation, we introduce bone length loss $\mathcal{L}_{BL}$ to constrain the global joint positions of each person to satisfy skeleton consistency, which implicitly encodes the human body's kinematic structure. We formulate the bone length loss as follows:

$$\mathcal{L}_{BL} = ||B(\hat{\mathbf{x}}_a) - B(\mathbf{x}_a)||_2^2 + ||B(\hat{\mathbf{x}}_b) - B(\mathbf{x}_b)||_2^2, \quad (11)$$

where $B$ represents the bone lengths in a pre-defined human body kinematic tree, derived from the global joint positions in $\mathbf{x}$.

**Interactive losses.** To handle the complexity of spatial relations in multi-person interactions, we further introduce interactive losses, comprising masked joint distance map (DM) loss and relative orientation (RO) loss, as illustrated in Fig. 6. The DM loss measures the $N_j \times N_j$ joint distance map of two people and matches it with the ground truth, where $N_j$ is the number of joints per person. Thus we design the DM loss as follows:

$$\mathcal{L}_{DM}$$
$$= ||(M(\hat{\mathbf{x}}_a, \hat{\mathbf{x}}_b) - M(\mathbf{x}_a, \mathbf{x}_b)) \odot I(M_{xz}(\mathbf{x}_a, \mathbf{x}_b) < \bar{M})||_2^2, \quad (12)$$

where $M$ denotes the joint distance map of two people, obtained from the global joint positions in their motions, $I(\cdot)$ is the indicator function that masks the loss by applying a 2D distance threshold on the XZ-plane,

which activates this loss only when the horizontal distance between the two people is small enough, $M_{xz}$ represents the distance map projected onto the XZ-plane, $\bar{M}$ is the distance threshold, and $\odot$ indicates Hadamard product.

The RO loss estimates the relative orientation of two people and aligns it with the ground truth. The RO loss is formulated as:

$$\mathcal{L}_{RO} = ||O(IK(\hat{\mathbf{x}}_a), IK(\hat{\mathbf{x}}_b)) - O(IK(\mathbf{x}_a), IK(\mathbf{x}_b))||_2^2, \quad (13)$$

where $IK(\cdot)$ represents the inverse kinematics process, which outputs the joint rotations, and $O$ indicates the 2D relative orientation between the two people around the Y-axis obtained from rotations. These losses constrain the positional spatial relations and the relative frontal orientation of the two people to be consistent with the nature of human interactions.

Here, we summarize our additional regularization loss as follows:

$$\mathcal{L}_{reg} = \lambda_{vel}\mathcal{L}_{vel} + \lambda_{foot}\mathcal{L}_{foot} + \lambda_{BL}\mathcal{L}_{BL}$$
$$+ \lambda_{DM}\mathcal{L}_{DM} + \lambda_{RO}\mathcal{L}_{RO}. \quad (14)$$

**Regularization loss schedule.** Physdiff (Yuan et al., 2022) informed that the network-predicted denoised motion is increasingly implausible as the diffusion timestep $t$ is larger (noise level is higher). Since denoisers estimate the expectation $\mathbb{E}[\mathbf{x}|\mathbf{x}^{(t)}]$, i.e., the mean motion $\mathbf{x}$ given $\mathbf{x}^{(t)}$, and empirically they tend to output average poses with some root translations and rotations when the noise level is high. This results in not only severe physical implausible motions but also unrealistic interactions between two people in our two-person scenario. If we apply the above regularization losses when $t$ is large, the network output will become the minimum mean squared error (MMSE) estimation of biased losses, deviating from $\mathbb{E}[\mathbf{x}|\mathbf{x}^{(t)}]$.

Inspired by that, we devise a novel diffusion training scheme. We truncate diffusion timesteps with a threshold $\bar{t}$ and only apply regularization loss to the network when the sampled timestep $t$ is below the threshold. Thus the total loss is formulated as:

$$\mathcal{L} = \mathcal{L}_{simple} + \lambda_{reg}\mathbb{E}_t[I(t \leq \bar{t}) \cdot \mathcal{L}_{reg}^{(t)}], \quad (15)$$

where $I(t \leq \bar{t})$ is an indicator function, which drops the regularization term when $t > \bar{t}$.

### 4.4 Implementation Details

We implement our InterGen with $N = 8$ blocks whose latent dimension is set to 1024 and each attention layer
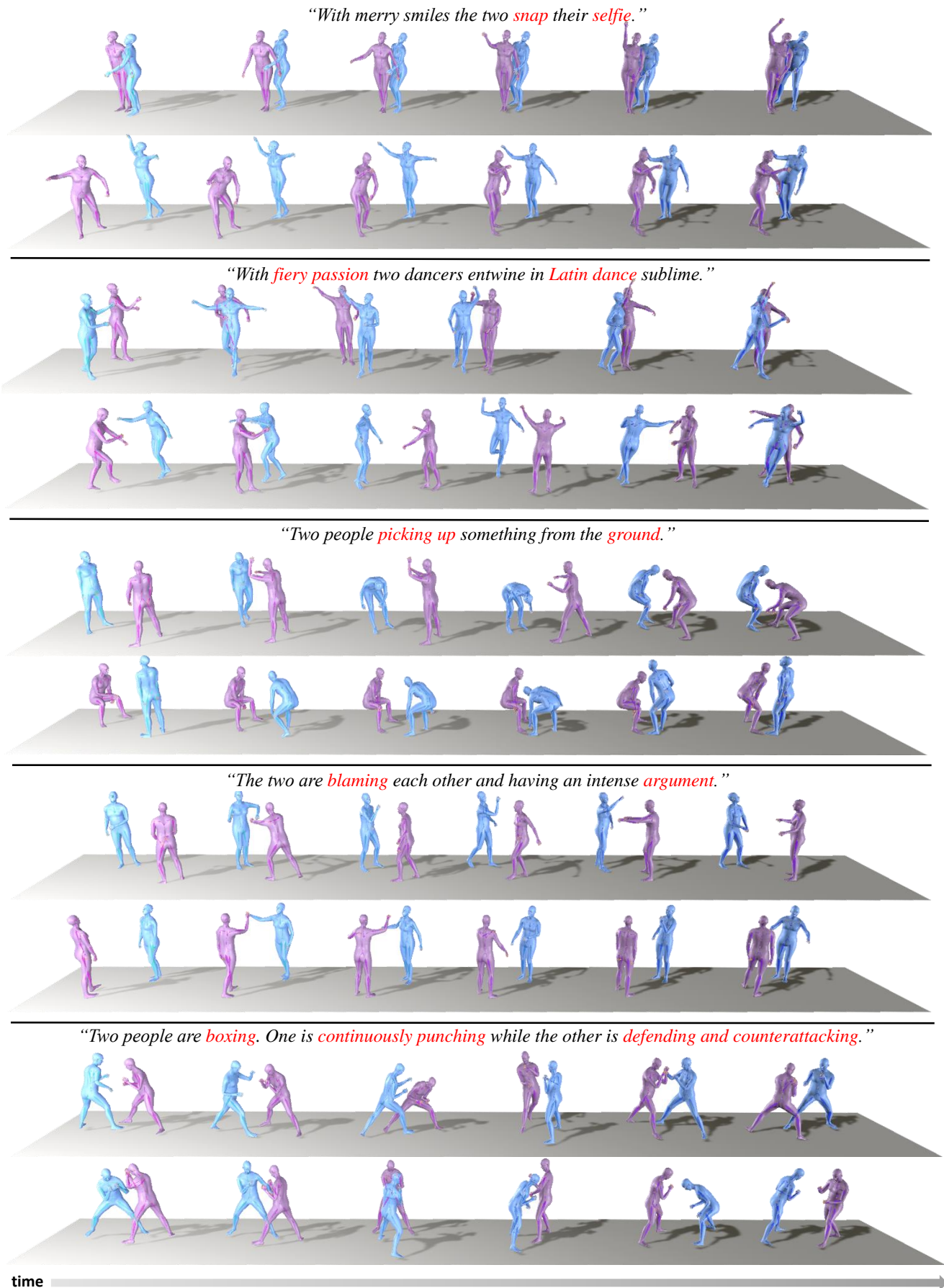
*"With merry smiles the two snap their selfie."*

*"With fiery passion two dancers entwine in Latin dance sublime."*

*"Two people picking up something from the ground."*

*"The two are blaming each other and having an intense argument."*

*"Two people are boxing. One is continuously punching while the other is defending and counterattacking."*

**time**

**Fig. 7 Qualitative results** generated by our InterGen model. We showcase two different samples per text prompt, which demonstrate the high quality and diversity of our interaction generation.

**Table 2 Quantitative comparisons on the InterHuman test set.** We run all the evaluations 20 times except MModality runs 5 times. ± indicates the 95% confidence interval. **Bold** indicates best result.

| Methods | R Precision↑ | | | FID ↓ | MM Dist↓ | Diversity→ | MModality ↑ |
|---|---|---|---|---|---|---|---|
| | Top 1 | Top 2 | Top 3 | | | | |
| Real | $0.452^{\pm.008}$ | $0.610^{\pm.009}$ | $0.701^{\pm.008}$ | $0.273^{\pm.007}$ | $3.755^{\pm.008}$ | $7.948^{\pm.064}$ | - |
| TEMOS | $0.224^{\pm.010}$ | $0.316^{\pm.013}$ | $0.450^{\pm.018}$ | $17.375^{\pm.043}$ | $6.342^{\pm.015}$ | $6.939^{\pm.071}$ | $0.535^{\pm.014}$ |
| T2M | $0.238^{\pm.012}$ | $0.325^{\pm.010}$ | $0.464^{\pm.014}$ | $13.769^{\pm.072}$ | $5.731^{\pm.013}$ | $7.046^{\pm.022}$ | $1.387^{\pm.076}$ |
| MDM | $0.153^{\pm.012}$ | $0.260^{\pm.009}$ | $0.339^{\pm.012}$ | $9.167^{\pm.056}$ | $7.125^{\pm.018}$ | $\mathbf{7.602}^{\pm.045}$ | $\mathbf{2.355}^{\pm.080}$ |
| ComMDM* | $0.067^{\pm.013}$ | $0.125^{\pm.018}$ | $0.184^{\pm.015}$ | $38.643^{\pm.098}$ | $14.211^{\pm.013}$ | $3.520^{\pm.058}$ | $0.217^{\pm.018}$ |
| ComMDM | $0.223^{\pm.009}$ | $0.334^{\pm.008}$ | $0.466^{\pm.010}$ | $7.069^{\pm.054}$ | $6.212^{\pm.021}$ | $7.244^{\pm.038}$ | $1.822^{\pm.052}$ |
| InterGen (Ours) | $\mathbf{0.371}^{\pm.010}$ | $\mathbf{0.515}^{\pm.012}$ | $\mathbf{0.624}^{\pm.010}$ | $\mathbf{5.918}^{\pm.079}$ | $\mathbf{5.108}^{\pm.014}$ | $7.387^{\pm.029}$ | $2.141^{\pm.063}$ |

consists of 8 heads, as same as the re-implemented MDM and ComMDM. We employ a frozen *CLIP-ViT-L-14* model as the text encoder. The number of diffusion timesteps is set to 1,000 during training and we apply the DDIM (Song et al., 2020a) sampling strategy with 50 timesteps and $\eta = 0$. We adopt the cosine noise level schedule (Nichol and Dhariwal, 2021) and classifier-free guidance (Ho and Salimans, 2022) where the 10% random CLIP embeddings are set to zero during training and the guidance coefficient is set to 3.5 during sampling. All the models are trained with AdamW optimizer using a fixed learning rate of $10^{-4}$. For Eqn. (14), we set $\lambda_{vel} = 10$, $\lambda_{foot} = 10$, $\lambda_{BL} = 5$, $\lambda_{DM} = 10$, and $\lambda_{RO} = 5$, and for Eqn. (15), we set $\lambda_{reg} = 1$ in all the experiments. We train our diffusion denoisers with minibatch size 64 for 2,000 epochs on two Nvidia Tesla V100 GPUs.

## 5 Experimental Results

Here, we demonstrate the capability of our approach in a variety of scenarios. We first introduce the evaluation dataset and metrics and then showcase a gallery of our generation results for two-person interactions in Fig. 7. We then provide the comparison with previous methods as well as the evaluation of our technical components, both qualitatively and quantitatively, followed by the analysis of various downstream applications.

**Evaluation dataset.** The existing available human motion datasets lack sufficient categories of human interactions and corresponding text descriptions. We hence contribute a new dataset, InterHuman, to evaluate our approach. It provides accurate skeletal human motions and rich natural language descriptions, covering diverse two-person interaction scenarios (see Sec. 3 for more details). We also augment our data following the HumanML3D (Guo et al., 2022), which involves mirroring all motions and replacing relevant keywords in the descriptions, and swapping the order of two peo-

ple in all interactions. We then split the data into training, validation, and test sets using the same protocol.

**Evaluation metrics.** We adopt the same evaluation metrics as Heusel et al. (2017), namely Frechet Inception Distance (FID) measuring the latent distance between generated dataset and real dataset, interaction recognition accuracy (R Precision) measuring the text-motion matching, diversity (Diversity) measuring latent variance, multimodality (MModality) measuring diversity within the same text, and multimodality distance (MM Dist) measuring the distance between motions and texts. These metrics jointly evaluate the quality and diversity of the generated interactions. To calculate these metrics, we train an interaction feature extractor and a text feature extractor using contrastive loss following Radford et al. (2021), which encourages matched text-interaction pairs to have geometrically close feature vectors.

### 5.1 Comparison

We compare our InterGen with various representative text-to-motion methods in two-person interactive scenarios. Specifically, we apply single-person methods VAE-based TEMOS (Petrovich et al., 2022) and T2M (Guo et al., 2022), diffusion-based MDM (Tevet et al., 2022b), and recent two-person method ComMDM (Shafir et al., 2023). To thoroughly evaluate our method and conduct fair comparisons, we retrain the above methods with the same InterHuman training set and test on the test set. Note that for extending the above single-person motion synthesis models to handle two-person interaction, we modify the input and output dimensions of their networks to accommodate our non-canonical representation of two-person interaction. For fair comparisons, we report the results of ComMDM pre-trained and fine-tuned in the original few-shot setting with 10 training samples (with *) and retrained on the same InterHuman training set (without *) with the
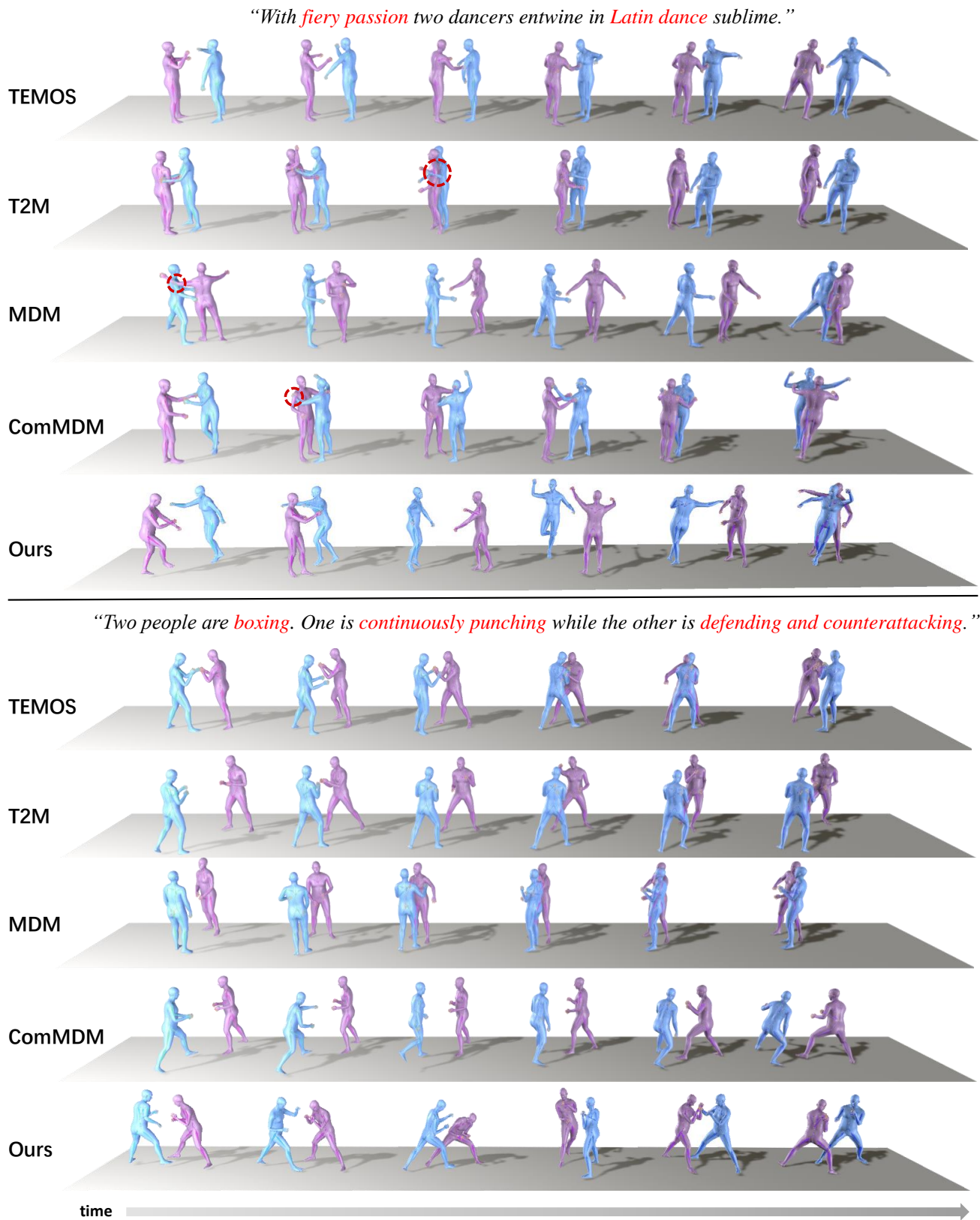
**Fig. 8** Qualitative comparison with SOTA techniques. The inputs to the model are listed at the top and middle, while the outputs of different models  (Petrovich et al., 2022; Guo et al., 2022; Tevet et al., 2022b; Shafir et al., 2023) are listed below. Intersecting portions of the motions are highlighted with red dashed circles.

**Table 3** Quantitative evaluation of our key designs. The performance drop-off highlights our technical contributions.

| Methods | R Precision↑ | | | FID ↓ | MM Dist↓ | Diversity→ | MModality ↑ |
|---|---|---|---|---|---|---|---|
| | Top 1 | Top 2 | Top 3 | | | | |
| Real | $0.452^{\pm.008}$ | $0.610^{\pm.009}$ | $0.701^{\pm.008}$ | $0.273^{\pm.007}$ | $3.755^{\pm.017}$ | $7.948^{\pm.008}$ | - |
| canonical rep | $0.134^{\pm.011}$ | $0.233^{\pm.013}$ | $0.315^{\pm.013}$ | $11.322^{\pm.055}$ | $7.865^{\pm.015}$ | $\mathbf{8.534}^{\pm.028}$ | $\mathbf{2.963}^{\pm.068}$ |
| concat-conditioning | $0.303^{\pm.010}$ | $0.451^{\pm.010}$ | $0.575^{\pm.012}$ | $6.273^{\pm.080}$ | $5.579^{\pm.012}$ | $7.123^{\pm.035}$ | $1.790^{\pm.052}$ |
| w/o weights sharing | $0.153^{\pm.015}$ | $0.257^{\pm.013}$ | $0.337^{\pm.012}$ | $8.059^{\pm.077}$ | $7.242^{\pm.010}$ | $7.508^{\pm.026}$ | $2.288^{\pm.076}$ |
| w/o DM loss | $0.293^{\pm.012}$ | $0.437^{\pm.010}$ | $0.533^{\pm.016}$ | $6.653^{\pm.064}$ | $5.934^{\pm.011}$ | $7.328^{\pm.033}$ | $2.181^{\pm.055}$ |
| w/o RO loss | $0.310^{\pm.013}$ | $0.466^{\pm.009}$ | $0.587^{\pm.010}$ | $6.311^{\pm.052}$ | $5.515^{\pm.012}$ | $7.309^{\pm.050}$ | $2.019^{\pm.053}$ |
| InterGen (Ours) | $\mathbf{0.371}^{\pm.010}$ | $\mathbf{0.515}^{\pm.012}$ | $\mathbf{0.624}^{\pm.010}$ | $\mathbf{5.918}^{\pm.079}$ | $\mathbf{5.108}^{\pm.014}$ | $7.387^{\pm.029}$ | $2.141^{\pm.063}$ |

**Table 4** Quantitative evaluation of our regularization loss schedule training scheme. The strategy of different treatments for different noise levels improves the performance significantly.

| RegLoss schedule | R Precision↑ | | | FID ↓ | MM Dist↓ | Diversity→ | MModality ↑ |
|---|---|---|---|---|---|---|---|
| | Top 1 | Top 2 | Top 3 | | | | |
| Real motion | $0.452^{\pm.008}$ | $0.610^{\pm.009}$ | $0.701^{\pm.008}$ | $0.273^{\pm.066}$ | $3.755^{\pm.015}$ | $7.948^{\pm.008}$ | - |
| None | $0.201^{\pm.010}$ | $0.315^{\pm.014}$ | $0.406^{\pm.009}$ | $7.862^{\pm.074}$ | $6.919^{\pm.011}$ | $7.301^{\pm.044}$ | $\mathbf{2.198}^{\pm.042}$ |
| $t \leq 0.1T$ | $0.285^{\pm.013}$ | $0.443^{\pm.010}$ | $0.544^{\pm.014}$ | $6.556^{\pm.065}$ | $5.822^{\pm.013}$ | $7.315^{\pm.032}$ | $2.156^{\pm.055}$ |
| $t \leq 0.2T$ | $0.310^{\pm.013}$ | $0.464^{\pm.010}$ | $0.561^{\pm.014}$ | $6.178^{\pm.086}$ | $5.443^{\pm.012}$ | $7.342^{\pm.037}$ | $2.122^{\pm.078}$ |
| $t \leq 0.3T$ | $0.338^{\pm.010}$ | $0.483^{\pm.014}$ | $0.582^{\pm.011}$ | $6.034^{\pm.069}$ | $5.461^{\pm.009}$ | $7.323^{\pm.040}$ | $2.114^{\pm.039}$ |
| $t \leq 0.4T$ | $0.353^{\pm.010}$ | $0.496^{\pm.014}$ | $0.598^{\pm.012}$ | $5.945^{\pm.078}$ | $5.263^{\pm.013}$ | $7.303^{\pm.035}$ | $2.095^{\pm.061}$ |
| $t \leq 0.5T$ | $0.362^{\pm.011}$ | $0.506^{\pm.013}$ | $0.610^{\pm.014}$ | $5.938^{\pm.052}$ | $5.175^{\pm.016}$ | $7.318^{\pm.041}$ | $2.129^{\pm.046}$ |
| $t \leq 0.6T$ | $0.367^{\pm.013}$ | $0.513^{\pm.009}$ | $0.619^{\pm.010}$ | $\mathbf{5.887}^{\pm.064}$ | $5.209^{\pm.012}$ | $7.356^{\pm.022}$ | $2.133^{\pm.067}$ |
| $t \leq 0.7T$ | $\mathbf{0.371}^{\pm.010}$ | $\mathbf{0.515}^{\pm.012}$ | $\mathbf{0.624}^{\pm.010}$ | $5.918^{\pm.079}$ | $\mathbf{5.108}^{\pm.014}$ | $\mathbf{7.387}^{\pm.029}$ | $2.141^{\pm.063}$ |
| $t \leq 0.8T$ | $0.299^{\pm.012}$ | $0.453^{\pm.011}$ | $0.555^{\pm.014}$ | $6.521^{\pm.070}$ | $5.745^{\pm.012}$ | $7.346^{\pm.024}$ | $2.177^{\pm.062}$ |
| $t \leq 0.9T$ | $0.278^{\pm.010}$ | $0.432^{\pm.009}$ | $0.532^{\pm.013}$ | $6.664^{\pm.062}$ | $6.042^{\pm.012}$ | $7.314^{\pm.039}$ | $2.103^{\pm.072}$ |
| $t \leq T$ | $0.232^{\pm.012}$ | $0.365^{\pm.011}$ | $0.468^{\pm.010}$ | $7.037^{\pm.053}$ | $6.620^{\pm.010}$ | $7.282^{\pm.028}$ | $2.135^{\pm.080}$ |

same data representation. Note that the source code and training data of ComMDM are not publicly available yet and we re-implement it with the same setting on our dataset.

Tab. 2 summarizes the quantitative comparison results. Specifically, our approach outperforms other baselines in terms of FID, R precision, and MM Dist. Note that these metrics numbers are calculated using the whole test sets, which indicates that our InterGen achieves more compelling interaction motion generation with more accurate text/motion matching. The corresponding representative qualitative results are provided in Fig. 8, which demonstrates that our approach ensures the diversity and plausibility of the generated two-person interaction motions. Note that our approach can generate more natural interaction states, more diverse motions, and more accurate alignment of the global relative orientations and translations between the two performers.

## 5.2 Evaluations

Here, we provide ablation analysis for the key technical designs in our InterGen approach.

**Interaction motion representation.** We first compare our non-canonical representation with the canonical representation, as shown in *canonical rep* row in Tab. 3. We replace the motion representation using canonical representation and retrain our model in the same setting. The R precision and FID drop off significantly, which indicates the lower quality of motions and text prompt matching. The qualitative result is shown in the first row of Fig. 9. As time flows, the cumulative error of the trajectory becomes increasingly larger, which leads to two people performing their own motions without any spatial relation to each other, making interactions become unrealistic. This demonstrates the superiority of our two-person non-canonical representation in the common frame, which explicitly preserves the spatial relations of two people in the same space, which facilitates the model to learn.
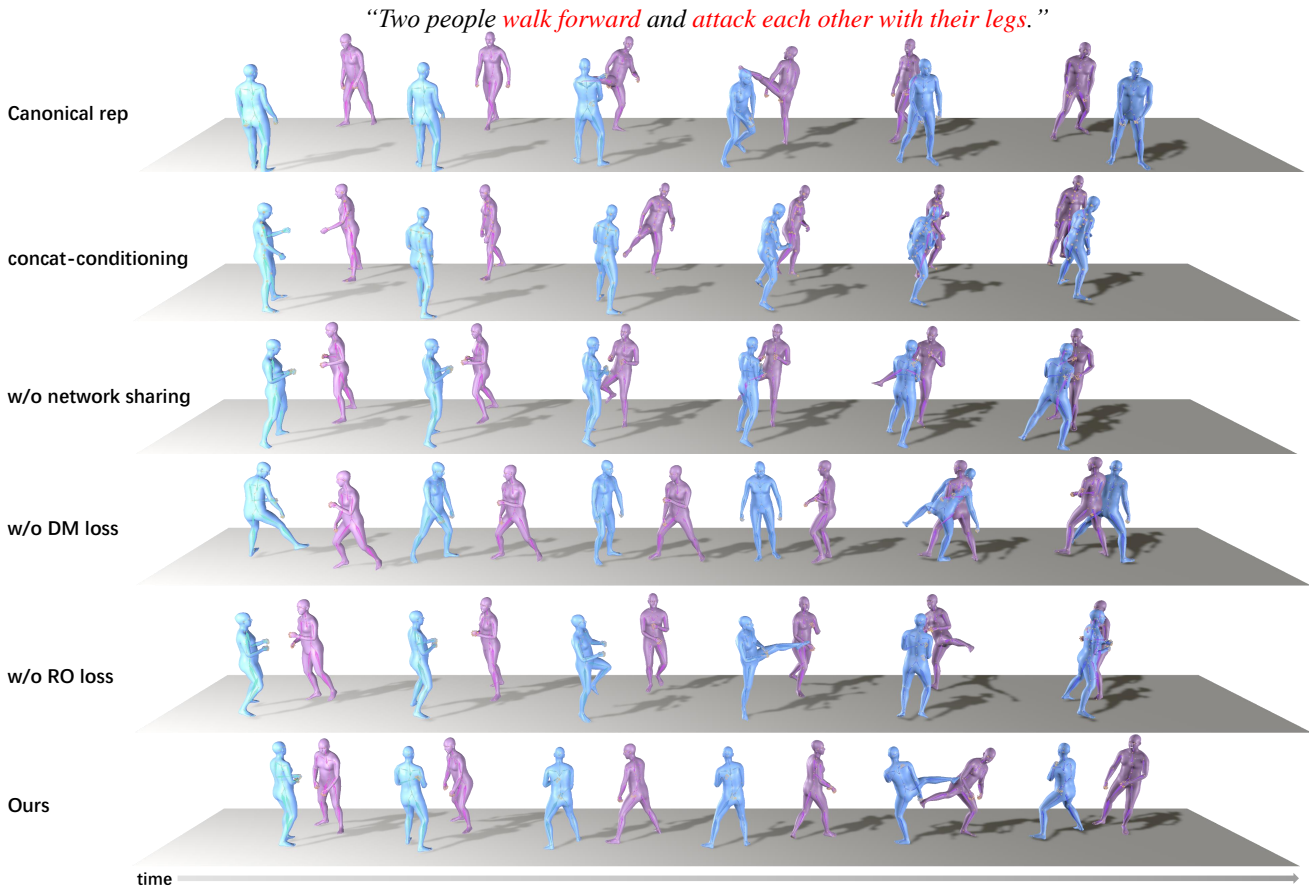
**Fig. 9** Qualitative results of ablation study. The top of the figure displays text prompts, while the lower illustrates the results of different ablation experiments and our best result. For quantitative comparisons of experimental results, please refer to Tab. 3.

**Cooperative networks.** Next, we evaluate the effectiveness of our cooperative networks. We first replace the mutual attention conditioning mechanism with concatenation, where each network concatenates its own noisy motion with the counterpart noisy motion as the input and drop the mutual attention layers in its blocks and retrain the model. The quantitative result is shown in the $concat-conditioning$ row in Tab. 3, which shows a decline in performance. The qualitative result is shown in the second row of Fig. 9, where the plausibility of the interaction and alignment of the motions are significantly reduced compared to ours.

We then ablate the weights sharing mechanism and train two networks without explicitly enforcing them to be the same. The quantitative result is presented in the $w/o\_weights\_sharing$ row in Tab. 3, where R precision and FID are significantly worse than ours, as also illustrated in the third row of Fig. 9, where two people do not show similar activation and motion capacity during the interaction. These indicate that our cooperative networks with mutual attention and weights sharing mechanisms can effectively handle not only the

complexity of interaction but also the balance of motion capacity between two people.

**Interaction losses.** We also apply ablation experiments on our two interactive losses and regularization loss schedule. As shown in the $w/o\_DM\_loss$ row of Tab. 3, we drop the DM loss and retrain our model, the interaction generation performance drops off, as shown in the fourth row of Fig. 9 qualitatively, the two people perform a weird interaction that they pass through each other's body since the absence of distance constraints.

And we drop the RO loss and apply the same process, as shown in the $w/o\_RO\_loss$ row of Tab. 3, getting a worse performance too. The qualitative result in the fifth row of Fig. 9 demonstrates that the two people perform an unrealistic interaction where the relative orientation between them does not match the motions they perform. These qualitative and quantitative results demonstrate the effectiveness and necessity of our interactive losses.

Tab. 4 demonstrates the effect of the choice of diffusion timesteps to apply regularization loss. We choose the model with the best performance occurring on $t \leq$
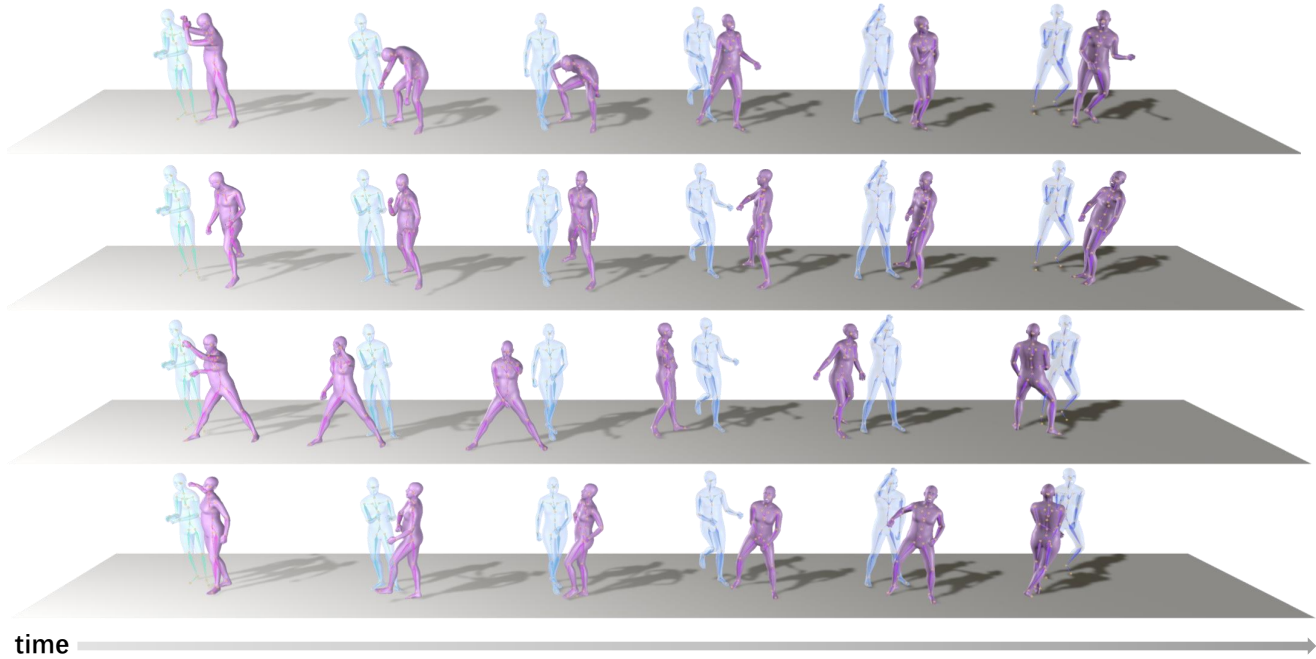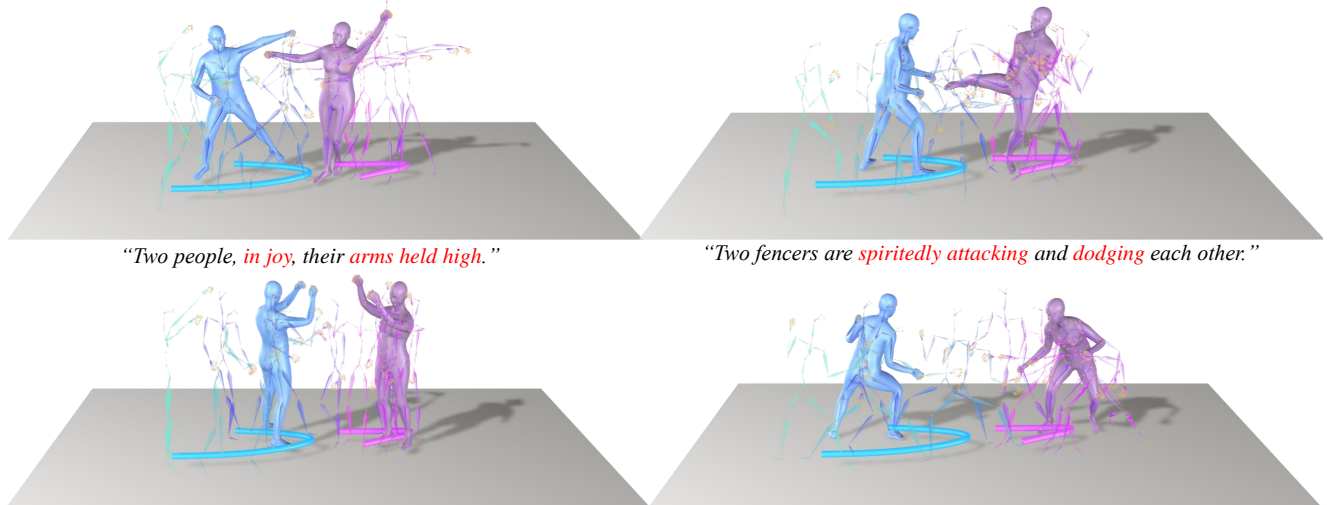
**Fig. 10 Person-to-person generation.** The above five motions are generated with the premise of freezing the motion of one person (represented by a semi-transparent SMPL model) while generating the motion of the other person (represented by an opaque SMPL model).



*"Two dancers, Latin steps they traced, arms softly swaying, interlaced."*

*"Two taekwondo masters train in offense and defense."*

*"Two people, in joy, their arms held high."*

*"Two fencers are spiritedly attacking and dodging each other."*

**Fig. 11 Trajectory control.** The curve beneath the peoples' feet represents their motion trajectories. The skeleton representation displays the position and motion of each frame over time, with the SMPL model (Loper et al., 2015) indicating a specific frame of motion. The text input for each motion is provided above the respective persons.

$0.7T$ with a cosine diffusion noise level schedule (Nichol and Dhariwal, 2021) as our final model, which outperforms the naively applying to all timesteps ($t \leq T$) and not applying such a regularization loss ($None$). This indicates the effectiveness of our loss schedule training scheme.

### 5.3 Application

Once trained, our InterGen model can be easily customized and employed in various human interaction-related tasks. Here, we showcase a series of downstream applications using our InterGen, i.e., human-to-human motion generation, additional trajectory control on top
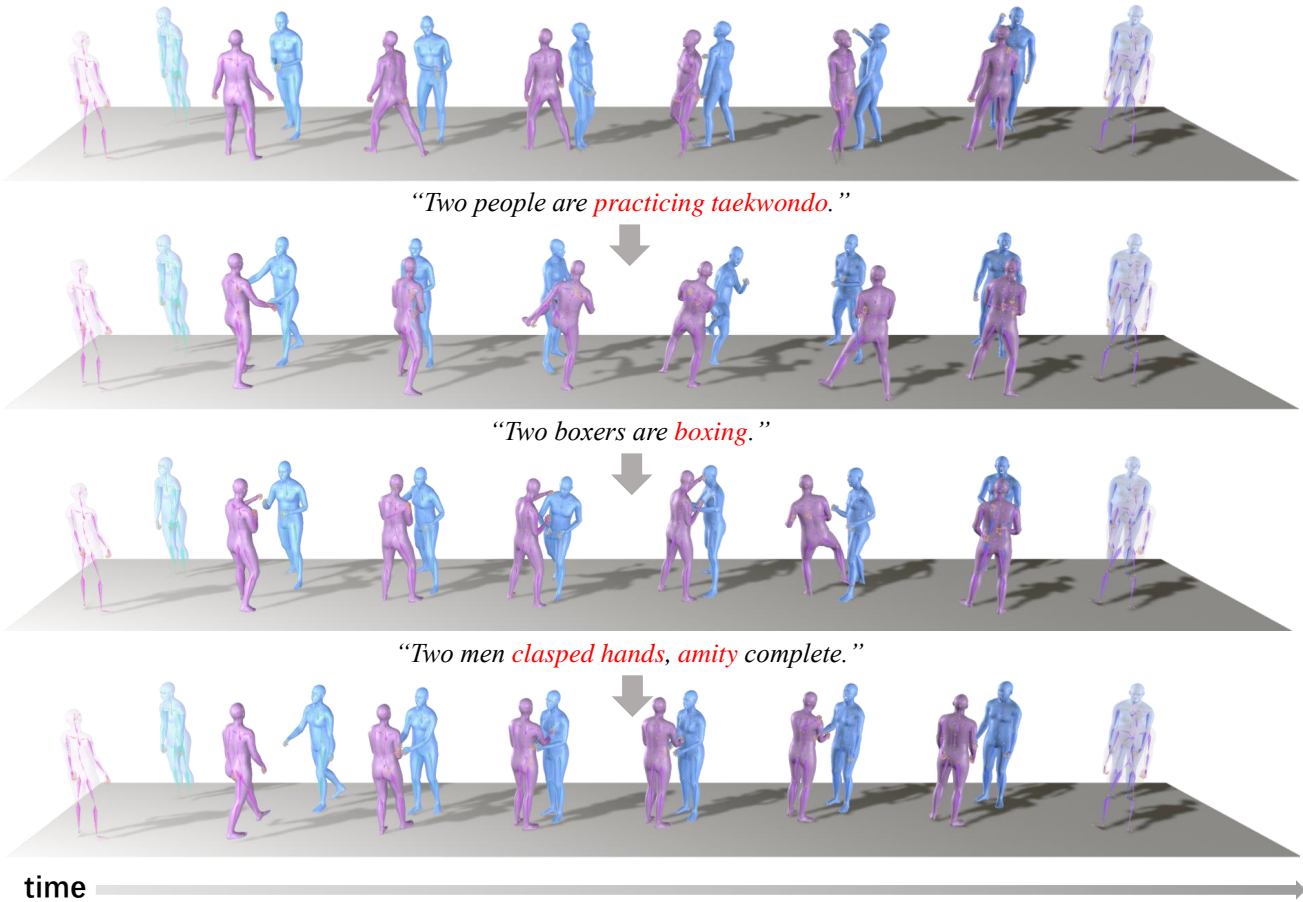
*"Two people are practicing taekwondo."*

*"Two boxers are boxing."*

*"Two men clasped hands, amity complete."*

time

**Fig. 12 Interaction inbetweening**. The first row depicts the outcome of motion freezing the starting and ending motions, without the use of any text prompts. Subsequently, the output generated in the second row is used as the input for text prompts, resulting in the inbetweening results presented in the next row. By iterating this process, the motion sequence is gradually transformed into a daily motion, as shown in the final row.

of text-prompt, and generating diverse motion inbetweening for interaction scenes.

**Person-to-person generation.** Here we fine-tune the original InterGen to enable the generation of human-to-human motions by simply taking a single-person motion sequence as input. Our method works by masking the noise applied to the motions which we wish to freeze during the forward pass of the diffusion process. The frozen motion serves as ground truth that propagates into the model and hence the model learns to rely on these motions when attempting to reconstruct the counterpart motions. This workflow is inspired by the fine-tuning strategy from MDM (Tevet et al., 2022b) and please refer to the original paper for implementation details. Differently, we freeze one person's motion instead of setting the trajectory to zero and utilize two inputs $\{\mathbf{x}_a, \mathbf{x}_b\}$ rather than one input $\mathbf{x}_0$. The visualization of this person-to-person generation process is shown in Fig. 10. The translucent person represents the frozen motion, i.e., the ground-truth motion at each iteration, while the opaque person represents the con-

ditionally generated motions. It can be observed that our diffusion process produces diverse motion results.

**Trajectory control.** We further demonstrate combining additional trajectory controls with the text inputs for more controllable human-to-human interaction generation. Specifically, it requires pre-inputting the trajectory during the text processing stage to mask the global transform of motion during the diffusion process. Such a strategy allows us to apply the trajectory control injected in the input stage to the diffusion at each time step, thus easily obtaining the interaction motions of two people that conform to the trajectory. As shown in Fig. 11, the trajectory of the human-to-human interaction and the differently stylized interaction actions under this trajectory are clearly depicted. The skeleton in the figure represents the historical trajectory of the motion, and the two visualized SMPL models (Loper et al., 2015) represent the end of the motion. Under the constraint of the same trajectory, our approach can generate diverse and naturally conformed motions.

**Interaction inbetweening.** Similar to person-to-person generation, we can also fine-tune the InterGen to freeze the beginning and ending segments of the motion sequence, so as to enforce InterGen to generate diverse Interaction motions that fill the frozen segments in between for motion inbetweening applications. Specifically, we implement time freezing of human-to-human motion in the pipeline and freeze motions instead of the training step during sampling. We then replace the part of in-between motion sequence inputs at each timestep when sampling. As shown in Fig. 12, the translucent double models represent the motion that has been frozen at the initial and end time steps, and the opaque double models represent the generated motions. Note that our approach can generate vivid motion inbetweening results that are also incorporated with the text-prompt inputs, even under various interaction scenarios.

## 5.4 Limitation and Discussion

Our InterGen has achieved compelling results for generating realistic two-person interaction motions from text prompts, yet still yielding some limitations. First, our approach only considers the interactions between two people, which may limit the applicability of our approach to more complex scenarios (e.g., crowd dynamics simulation). It's promising to model more complex group interactions such as team sports or group dances, yet requiring much more diverse training datasets. Besides, our approach only generates motion sequences based on given text prompts, without further considering the feedback from the user. This may limit the user's creativity over the generated motions, especially when the text prompts are vague or ambiguous. It's an interesting direction to generate motions that satisfy the user's specific preferences or expectations, such as the style, speed, or intensity of the motions. Our generated motion is also limited to a fixed largest length, which cannot support extremely long-sequence generation. It hence hinders the diversity and coherence of the generated motions, especially when the text prompts are complex. More advanced strategies like involving multiple sub-motions or transitions to span a longer period of time will help, where our approach can serve as a solid cornerstone.

## 6 Conclusion

We gave presented InterGen, a diffusion-based approach to conveniently generate two-person motion under diverse interactions, from only text-prompt con-trols. Specifically, we contribute a novel multimodal dataset with rich motion results and natural language descriptions, covering a wide range of interaction scenarios. Then, in our interaction diffusion model, our cooperative denoisers with sharing weights and a mutual attention mechanism can effectively model the symmetric fact of human identities during interactions. Our non-canonical motion representation also effectively models the global relations between performers for the interaction setting. Our regularization design with a specific damping scheme further encodes the spatial relations to generate more diverse and reasonable interactions. Extensive experimental results have demonstrated the effectiveness of InterGen for the generation of compelling two-person motions and a series of downstream interaction applications. We believe our approach and multimodal dataset can serve as a solid step towards text-guided generation and understanding of human-to-human interactions, with numerous potential applications for entertainment, gaming, and immersive experience in VR/AR.

## References

Van der Aa N, Luo X, Giezeman GJ, Tan RT, Veltkamp RC (2011) Umpm benchmark: A multi-person dataset with synchronized video and motion capture data for evaluation of articulated human motion and interaction. In: 2011 IEEE international conference on computer vision workshops (ICCV Workshops), IEEE, pp 1264–1269

Ahn H, Ha T, Choi Y, Yoo H, Oh S (2018) Text2action: Generative adversarial synthesis from language to action. In: 2018 IEEE International Conference on Robotics and Automation (ICRA), IEEE, pp 5915–5920

Ahuja C, Morency LP (2019) Language2pose: Natural language grounded pose forecasting. In: 2019 International Conference on 3D Vision (3DV), IEEE, pp 719–728

Andrews S, Huerta I, Komura T, Sigal L, Mitchell K (2016) Real-time physics-based motion capture with sparse sensors. In: Proceedings of the 13th European conference on visual media production (CVMP 2016), pp 1–10

Anguelov D, Srinivasan P, Koller D, Thrun S, Rodgers J, Davis J (2005) Scape: shape completion and animation of people. In: ACM SIGGRAPH 2005 Papers, pp 408–416

Ao T, Gao Q, Lou Y, Chen B, Liu L (2022) Rhythmic gesticulator: Rhythm-aware co-speech gesture synthesis with hierarchical neural embeddings. ACM Transactions on Graphics (TOG) 41(6):1–19

Bogo F, Kanazawa A, Lassner C, Gehler P, Romero J, Black MJ (2016) Keep it smpl: Automatic estimation of 3d human pose and shape from a single image. In: Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part V 14, Springer, pp 561–578

Bregler C, Malik J (1998) Tracking people with twists and exponential maps. In: Proceedings. 1998 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (Cat. No. 98CB36231), IEEE, pp 8–15

Brown T, Mann B, Ryder N, Subbiah M, Kaplan JD, Dhariwal P, Neelakantan A, Shyam P, Sastry G, Askell A, et al.

(2020) Language models are few-shot learners. Advances in neural information processing systems 33:1877–1901

Chen X, Jiang B, Liu W, Huang Z, Fu B, Chen T, Yu J, Yu G (2022a) Executing your commands via motion diffusion in latent space. arXiv preprint arXiv:221204048

Chen X, Su Z, Yang L, Cheng P, Xu L, Fu B, Yu G (2022b) Learning variational motion prior for video-based motion capture. arXiv preprint arXiv:221015134

De Aguiar E, Stoll C, Theobalt C, Ahmed N, Seidel HP, Thrun S (2008) Performance capture from sparse multi-view video. In: ACM SIGGRAPH 2008 papers, pp 1–10

Devlin J, Chang MW, Lee K, Toutanova K (2018) Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:181004805

Duan Y, Shi T, Zou Z, Lin Y, Qian Z, Zhang B, Yuan Y (2021) Single-shot motion completion with transformer. arXiv preprint arXiv:210300776

Gall J, Rosenhahn B, Brox T, Seidel HP (2010) Optimization and filtering for human motion capture. International Journal of Computer Vision (IJCV) 87(1–2):75–92

Ghosh A, Cheema N, Oguz C, Theobalt C, Slusallek P (2021) Synthesis of compositional animations from textual descriptions. In: Proceedings of the IEEE/CVF international conference on computer vision, pp 1396–1406

Gilbert A, Trumble M, Malleson C, Hilton A, Collomosse J (2019) Fusing visual and inertial sensors with semantics for 3d human pose estimation. International Journal of Computer Vision 127:381–397

Goodfellow I, Pouget-Abadie J, Mirza M, Xu B, Warde-Farley D, Ozair S, Courville A, Bengio Y (2020) Generative adversarial networks. Communications of the ACM 63(11):139–144

Guo C, Zuo X, Wang S, Zou S, Sun Q, Deng A, Gong M, Cheng L (2020) Action2motion: Conditioned generation of 3d human motions. In: Proceedings of the 28th ACM International Conference on Multimedia, pp 2021–2029

Guo C, Zou S, Zuo X, Wang S, Ji W, Li X, Cheng L (2022) 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

Habermann M, Xu W, Zollhöfer M, Pons-Moll G, Theobalt C (2019) Livecap: Real-time human performance capture from monocular video. ACM Transactions on Graphics (TOG) 38(2):14:1–14:17

Habermann M, Xu W, Zollhofer M, Pons-Moll G, Theobalt C (2020) Deepcap: Monocular human performance capture using weak supervision. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)

Habibie I, Elgharib M, Sarkar K, Abdullah A, Nyatsanga S, Neff M, Theobalt C (2022) A motion matching-based framework for controllable gesture synthesis from speech. In: ACM SIGGRAPH 2022 Conference Proceedings, pp 1–9

Harvey FG, Yurick M, Nowrouzezahrai D, Pal C (2020) Robust motion in-betweening. ACM Transactions on Graphics (TOG) 39(4):60–1

He Y, Pang A, Chen X, Liang H, Wu M, Ma Y, Xu L (2021) Challencap: Monocular 3d capture of challenging human performances using multi-modal references. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp 11400–11411

Helten T, Muller M, Seidel HP, Theobalt C (2013) Real-time body tracking with one depth camera and inertial sensors. In: Proceedings of the IEEE international conference on computer vision, pp 1105–1112

Henschel R, Von Marcard T, Rosenhahn B (2020) Accurate long-term multiple people tracking using video and body-worn imus. IEEE Transactions on Image Processing 29:8476–8489

Heusel M, Ramsauer H, Unterthiner T, Nessler B, Hochreiter S (2017) Gans trained by a two time-scale update rule converge to a local nash equilibrium. Advances in neural information processing systems 30

Ho J, Salimans T (2022) Classifier-free diffusion guidance. arXiv preprint arXiv:220712598

Ho J, Jain A, Abbeel P (2020) Denoising diffusion probabilistic models. Advances in Neural Information Processing Systems 33:6840–6851

Huang Y, Bogo F, Lassner C, Kanazawa A, Gehler PV, Romero J, Akhter I, Black MJ (2017) Towards accurate marker-less human shape and pose estimation over time. In: 2017 international conference on 3D vision (3DV), IEEE, pp 421–430

Huang Y, Kaufmann M, Aksan E, Black MJ, Hilliges O, Pons-Moll G (2018) Deep inertial poser: Learning to reconstruct human pose from sparse inertial measurements in real time. ACM Transactions on Graphics (TOG) 37(6):1–15

Kanazawa A, Zhang JY, Felsen P, Malik J (2019) Learning 3d human dynamics from video. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp 5614–5623

Kim J, Kim J, Choi S (2022) Flame: Free-form language-based motion synthesis & editing. arXiv preprint arXiv:220900349

Kingma DP, Welling M (2013) Auto-encoding variational bayes. arXiv preprint arXiv:13126114

Kocabas M, Athanasiou N, Black MJ (2020) Vibe: Video inference for human body pose and shape estimation. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp 5253–5263

Kolotouros N, Pavlakos G, Black MJ, Daniilidis K (2019) Learning to reconstruct 3d human pose and shape via model-fitting in the loop. In: Proceedings of the IEEE/CVF international conference on computer vision, pp 2252–2261

Lassner C, Romero J, Kiefel M, Bogo F, Black MJ, Gehler PV (2017) Unite the people: Closing the loop between 3d and 2d human representations. In: Proceedings of the IEEE conference on computer vision and pattern recognition, pp 6050–6059

Lee HY, Yang X, Liu MY, Wang TC, Lu YD, Yang MH, Kautz J (2019) Dancing to music. Advances in neural information processing systems 32

Li B, Zhao Y, Zhelun S, Sheng L (2022) Danceformer: Music conditioned 3d dance generation with parametric motion transformer. In: Proceedings of the AAAI Conference on Artificial Intelligence, vol 36, pp 1272–1279

Li R, Yang S, Ross DA, Kanazawa A (2021) Ai choreographer: Music conditioned 3d dance generation with aist++. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp 13401–13412

Liang H, He Y, Zhao C, Li M, Wang J, Yu J, Xu L (2022) Hybridcap: Inertia-aid monocular capture of challenging human motions. arXiv preprint arXiv:220309287

Liu J, Shahroudy A, Perez M, Wang G, Duan LY, Kot AC (2019) Ntu rgb+ d 120: A large-scale benchmark for 3d human activity understanding. IEEE transactions on pattern analysis and machine intelligence 42(10):2684–2701

Liu Y, Gall J, Stoll C, Dai Q, Seidel HP, Theobalt C (2013) Markerless motion capture of multiple characters using multiview image segmentation. IEEE transactions on pat-

tern analysis and machine intelligence 35(11):2720–2735

Loper M, Mahmood N, Romero J, Pons-Moll G, Black MJ (2015) Smpl: A skinned multi-person linear model. ACM transactions on graphics (TOG) 34(6):1–16

Malleson C, Gilbert A, Trumble M, Collomosse J, Hilton A, Volino M (2017) Real-time full-body motion capture from video and imus. In: 2017 international conference on 3D vision (3DV), IEEE, pp 449–457

Malleson C, Collomosse J, Hilton A (2019) Real-time multi-person motion capture from multi-view video and imus. International Journal of Computer Vision pp 1–18

Movella (2022) Movella xsens products. `https://www.movella.com/products/xsens`, accessed: 2023-03-26

Ng E, Xiang D, Joo H, Grauman K (2020) You2me: Inferring body pose in egocentric video via first and second person interactions. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp 9890–9900

Nichol AQ, Dhariwal P (2021) Improved denoising diffusion probabilistic models. In: International Conference on Machine Learning, PMLR, pp 8162–8171

OpenAI (2023) Gpt-4 technical report. `2303.08774`

Osman AA, Bolkart T, Black MJ (2020) Star: Sparse trained articulated human body regressor. In: Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part VI 16, Springer, pp 598–613

Pavlakos G, Zhou X, Derpanis KG, Daniilidis K (2017) Harvesting multiple views for marker-less 3d human pose annotations. In: Computer Vision and Pattern Recognition (CVPR)

Pavlakos G, Choutas V, Ghorbani N, Bolkart T, Osman AA, Tzionas D, Black MJ (2019) 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

Peng XB, Ma Z, Abbeel P, Levine S, Kanazawa A (2021) Amp: Adversarial motion priors for stylized physics-based character control. ACM Transactions on Graphics (TOG) 40(4):1–20

Petrovich M, Black MJ, Varol G (2021) Action-conditioned 3d human motion synthesis with transformer vae. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp 10985–10995

Petrovich M, Black MJ, Varol G (2022) Temos: Generating diverse human motions from textual descriptions. In: Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXII, Springer, pp 480–497

Plappert M, Mandery C, Asfour T (2016) The kit motion-language dataset. Big data 4(4):236–252

Punnakkal AR, Chandrasekaran A, Athanasiou N, Quiros-Ramirez A, Black MJ (2021) Babel: bodies, action and behavior with english labels. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp 722–731

Radford A, Kim JW, Hallacy C, Ramesh A, Goh G, Agarwal S, Sastry G, Askell A, Mishkin P, Clark J, et al. (2021) Learning transferable visual models from natural language supervision. In: International conference on machine learning, PMLR, pp 8748–8763

Rempe D, Birdal T, Hertzmann A, Yang J, Sridhar S, Guibas LJ (2021) Humor: 3d human motion model for robust pose estimation. In: Proceedings of the IEEE/CVF international conference on computer vision, pp 11488–11499

Rezende D, Mohamed S (2015) Variational inference with normalizing flows. In: International conference on machine learning, PMLR, pp 1530–1538

Robertini N, Casas D, Rhodin H, Seidel HP, Theobalt C (2016) Model-based outdoor performance capture. In: 2016 Fourth International Conference on 3D Vision (3DV), IEEE, pp 166–175

Shafir Y, Tevet G, Kapon R, Bermano AH (2023) Human motion diffusion as a generative prior. arXiv preprint arXiv:230301418

Simon T, Joo H, Matthews I, Sheikh Y (2017) Hand keypoint detection in single images using multiview bootstrapping. In: Computer Vision and Pattern Recognition (CVPR)

Song J, Meng C, Ermon S (2020a) Denoising diffusion implicit models. arXiv preprint arXiv:201002502

Song Y, Sohl-Dickstein J, Kingma DP, Kumar A, Ermon S, Poole B (2020b) Score-based generative modeling through stochastic differential equations. arXiv preprint arXiv:201113456

Song Z, Wang D, Jiang N, Fang Z, Ding C, Gan W, Wu W (2022) Actformer: A gan transformer framework towards general action-conditioned 3d human motion generation. arXiv preprint arXiv:220307706

Starke S, Zhang H, Komura T, Saito J (2019) Neural state machine for character-scene interactions. ACM Trans Graph 38(6):209–1

Starke S, Mason I, Komura T (2022) Deepphase: Periodic autoencoders for learning motion phase manifolds. ACM Transactions on Graphics (TOG) 41(4):1–13

Stoll C, Hasler N, Gall J, Seidel HP, Theobalt C (2011) Fast articulated motion tracking using a sums of Gaussians body model. In: International Conference on Computer Vision (ICCV)

Tevet G, Gordon B, Hertz A, Bermano AH, Cohen-Or D (2022a) Motionclip: Exposing human motion generation to clip space. In: Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXII, Springer, pp 358–374

Tevet G, Raab S, Gordon B, Shafir Y, Cohen-Or D, Bermano AH (2022b) Human motion diffusion model. arXiv preprint arXiv:220914916

Theobalt C, de Aguiar E, Stoll C, Seidel HP, Thrun S (2010) Performance capture from multi-view video. In: Image and Geometry Processing for 3-D Cinematography, Springer, pp 127–149

Touvron H, Lavril T, Izacard G, Martinet X, Lachaux MA, Lacroix T, Rozière B, Goyal N, Hambro E, Azhar F, et al. (2023) Llama: Open and efficient foundation language models. arXiv preprint arXiv:230213971

Vicon (2019) Vicon Motion Systems. `https://www.vicon.com/`

Vlasic D, Adelsberger R, Vannucci G, Barnwell J, Gross M, Matusik W, Popović J (2007) Practical motion capture in everyday surroundings. ACM transactions on graphics (TOG) 26(3):35–es

Von Marcard T, Rosenhahn B, Black MJ, Pons-Moll G (2017) Sparse inertial poser: Automatic 3d human pose estimation from sparse imus. In: Computer Graphics Forum, Wiley Online Library, vol 36, pp 349–360

Von Marcard T, Henschel R, Black MJ, Rosenhahn B, Pons-Moll G (2018) Recovering accurate 3d human pose in the wild using imus and a moving camera. In: Proceedings of the European conference on computer vision (ECCV), pp 601–617

Wang J, Yan S, Dai B, Lin D (2021) Scene-aware generative network for human motion synthesis. In: Proceedings of the

IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp 12206–12215

Wang Z, Chen Y, Liu T, Zhu Y, Liang W, Huang S (2022) Humanise: Language-conditioned human motion generation in 3d scenes. arXiv preprint arXiv:221009729

Xu L, Liu Y, Cheng W, Guo K, Zhou G, Dai Q, Fang L (2018a) Flycap: Markerless motion capture using multiple autonomous flying cameras. IEEE Transactions on Visualization and Computer Graphics 24(8):2284–2297

Xu L, Xu W, Golyanik V, Habermann M, Fang L, Theobalt C (2020) Eventcap: Monocular 3d capture of high-speed human motions using an event camera. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp 4968–4978

Xu W, Chatterjee A, Zollhöfer M, Rhodin H, Mehta D, Seidel HP, Theobalt C (2018b) Monoperfcap: Human performance capture from monocular video. ACM Transactions on Graphics (TOG) 37(2):27:1–27:15

Yi X, Zhou Y, Xu F (2021) Transpose: Real-time 3d human translation and pose estimation with six inertial sensors. ACM Transactions on Graphics (TOG) 40(4):1–13

Yi X, Zhou Y, Habermann M, Shimada S, Golyanik V, Theobalt C, Xu F (2022) Physical inertial poser (pip): Physics-aware real-time human motion tracking from sparse inertial sensors. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)

You J, Leskovec J, He K, Xie S (2020) Graph structure of neural networks. In: International Conference on Machine Learning, PMLR, pp 10881–10891

Yuan Y, Song J, Iqbal U, Vahdat A, Kautz J (2022) Physdiff: Physics-guided human motion diffusion model. arXiv preprint arXiv:221202500

Z-cam (2022) Z CAM Cinema Camera. `https://www.z-cam.com`, accessed :2023-03-26

Zanfir A, Bazavan EG, Zanfir M, Freeman WT, Sukthankar R, Sminchisescu C (2021) Neural descent for visual 3d human pose and shape. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp 14484–14493

Zhang M, Cai Z, Pan L, Hong F, Guo X, Yang L, Liu Z (2022) Motiondiffuse: Text-driven human motion generation with diffusion model. arXiv preprint arXiv:220815001

Zhao C, Ren Y, He Y, Cong P, Liang H, Yu J, Xu L, Ma Y (2022) Lidar-aid inertial poser: Large-scale human motion capture by sparse inertial and lidar sensors. arXiv preprint arXiv:220515410

Zheng Z, Yu T, Li H, Guo K, Dai Q, Fang L, Liu Y (2018) Hybridfusion: Real-time performance capture using a single depth sensor and sparse imus. In: Proceedings of the European Conference on Computer Vision (ECCV), pp 384–400