

MoFusion: A Framework for Denoising-Diffusion-based Motion Synthesis

Rishabh Dabral¹ Muhammad Hamza Mughal^{1,2} Vladislav Golyanik¹ Christian Theobalt¹

¹Max Planck Institute for Informatics, SIC

²Saarland University

Figure 1. Our MoFusion approach synthesises long sequences of human motions in 3D from textual and audio inputs (e.g., by providing music samples). Our model has significantly improved generalisability and realism, and can be conditioned on modalities like text and audio. The resulting dance movements match the rhythm of the conditioning music, even if it is outside the training distribution.

Abstract

[mpg.de/projects/MoFusion/](https://www.mpg.de/projects/MoFusion/)

1. Introduction

Conventional methods for human motion synthesis have either been deterministic or have had to struggle with the trade-off between motion diversity vs motion quality. In response to these limitations, we introduce MoFusion, i.e., a new denoising-diffusion-based framework for high-quality unconditional human motion synthesis that can synthesise long, temporally plausible, and semantically accurate motions based on a range of conditioning contexts (such as music and text). We also present ways to introduce well-known kinematic losses for motion plausibility within the diffusion framework through our scheduled weighting strategy. The learned latent space can be used for several interactive motion-editing applications like in-betweening, seeding, and text-based editing, thus, providing crucial abilities for virtual-character animation and robotics. Through comprehensive quantitative evaluations and a perceptual user study, we demonstrate the effectiveness of MoFusion compared to the state of the art on established benchmarks in the literature. We urge the reader to watch our supplementary video <https://vcai.mpi-inf.de/>.

3D human motion synthesis is an important generative computer vision problem that often arises in robotics, virtual character animation and video games and movie production (e.g., for crowd dynamics simulation). It saw impressive progress over the last years; several works recently tackled it with reinforcement learning [41, 60, 64], deep generative models [2, 42, 43, 50] or using deterministic approaches [12, 29, 35]. Despite the progress, multiple open challenges remain, such as improving motion variability, enabling higher motion realism and enhancing synthesis fidelity under user-specified conditioning. Under conditioning, we understand in the context of conditioning the model outputs according to a control signal (e.g., “walking counter-clockwise”). The key goal of conditional human motion synthesis is to generate motions that semantically agree with the conditioning signal. To facilitate the same, the recent state-of-the-art approaches have widely adopted generative techniques like conditional variational auto-encoders (CVAE) [17, 32,

[42, 43], normalizing flows [2, 3], as well as GANs [19, 30]. Naturally, each of them has strengths and limitations. GAN-based synthesis methods suffer from mode-collapse, thus resulting in insufficient diversity of synthesis, especially for less common input conditioning. On the other hand, methods using CVAEs and normalizing flows typically have to deal with the trade-off between synthesis quality and the richness of the latent space (i.e., diversity) [3, 50].

The seminal works of Sohl-Dickstein et al. [20] and Ho et al. [21] recently demonstrated the ability of Denoising Diffusion Probabilistic Models (DDPM) to learn the underlying data distribution while also allowing for diverse sampling. Recent works [37, 49, 52] exhibited remarkable capabilities in the conditional synthesis of images and audio with high-frequency details while also allowing interactive applications like editing and inpainting. However, it has remained unclear how DDPM could be trained for such a problem with the temporal component as human motion synthesis.

Motivated by the recent advances in diffusion models, we propose MoFusion, i.e., a new approach for human motion synthesis with DDPM. This paper shows that diffusion models are highly effective for this task; see Fig. 1 for an overview. Our proposal includes a lightweight 1D U-Net network for reverse diffusion to reduce the rather long inference times. Furthermore, we demonstrate how domain-inspired kinematic losses can be introduced to diffusion framework during training, thanks to our time-varying weight schedule, which is our primary contribution. The result is a new versatile framework for human motion synthesis that produces diverse, temporally and kinematically plausible, and semantically accurate results.

We analyse DDPM for motion synthesis on two relevant sub-tasks: music-conditioned choreography generation and text-conditioned motion synthesis. While most existing choreography generation methods produce repetitive (loopy) motions, and text-to-motion synthesis methods struggle with left-right disambiguation, directional awareness and kinematic implausibility, we show that MoFusion barely suffers from these limitations. Finally, formulating motion synthesis in a diffusion framework also affords us the ability to perform interactive editing of the synthesised motion. To that end, we discuss the applications of a pre-trained MoFusion, like motion forecasting and in-betweening (both are important applications for virtual character animation). We show improvements in both the sub-tasks through quantitative evaluations on AIST++ [29] and HumanML3D [15] datasets as well as a user study. In summary, our core technical contributions are as follows:

- The first method for conditional 3D human motion synthesis using denoising diffusion models. Thanks to the proposed time-varying weight schedule, we incorporate several kinematic losses that make the syn-

thesised outputs temporally plausible and semantically accurate with the conditioning signal.

- Model conditioning on various signals, i.e., music and text, which is reflected in our framework's architecture. For a music-to-choreography generation, our results generalise well to new music and do not suffer from degenerate repetitiveness.

2. Related Works

We discuss the relevant literature from two vantage points, i.e., prior methods for human motion synthesis and literature on diffusion models.

2.1. Conditional Human Motion Synthesis

Traditionally, the problem of Human Motion Synthesis has been approached either by statistical modelling [6, 11] or sequence modelling techniques [10, 35]. Both approaches employed an initial seed sequence corresponding to a starting pose or past motion, which helps guide future motion prediction. However, synthesising motion sequences from scratch proves to be a harder task, where synthesis is guided by a conditioning mechanism.

A common approach in conditioned human motion synthesis is to guide the motion generation by using of class descriptions corresponding to actions [17, 42]. These approaches typically employ generative models like conditional VAEs [26] and learn a latent representation for motion based on action conditioning. Among such methods, Action2Motion [17] uses a frame-level motion representation with temporal VAEs, while ACTOR [42] improves results using a sequence-level motion representation with transformer-VAEs to synthesise motions based on action input. However, action conditioning does not provide a rich description of the target motion.

Text-Conditioned Motion Synthesis: The methods discussed above were followed by text-conditioned motion synthesis, developed on textually-annotated motion datasets like KIT [44], BABEL [46] and HumanML3D [15]. Such methods typically learn a shared latent space upon which both text and motion signals are projected [12, 65]. Lin et al. [31] use an LSTM encoder and a GRU decoder to predict future pose sequences. Ahuja et al. [1] and Ghosh et al. [12] focus on creating a joint language and pose representation to synthesise the motions autoregressively. TEMOS [43] builds upon the ideas by Ghosh et al. [12] and ACTOR [42] by using a transformer-VAE-based generative model with conditioning from a pre-trained language model. Finally, Guo et al. [15] use a temporal VAE to synthesise motions by extracting text-based features and then auto-regressively generating motion sequences.

Dance-Conditioned Motion Synthesis: Besides text, audio has also been applied to guide human motion synthe-

sis. Speech is used to learn gesture animations to mimic face, hand and body movements while speaking [2, 18]. Similarly, dance music has also been used extensively to synthesise motions. Various works [13, 24, 51, 59, 63] for music-conditioned motion synthesis tackle this problem by predicting motion from audio without seed motion. However, they converge to a mean pose, as dance typically consists of repetitive poses. Lai et al. [28] address this problem by providing an initial pose and the audio as input to a transformer-based architecture. DanceNet [73] proposes an autoregressive generative model, while Dance Revolution [22] uses a curriculum learning approach and a seq2seq architecture to synthesise dance motion. AI Choreographer [29] also approaches this problem by providing seed motion along with music to a cross-modal transformer for future dance motion prediction. Zhou et al. [71] enhance dance motions with music-to-dance alignment, and Aristidou et al. [4] enforce a global structure of the dance theme over the motion synthesis pipeline. The recent Bailando method [60] achieves impressive results in music-to-dance generation through a two-stage generation process. Their method learns to encode dance features into a codebook using a VQ-VAE [68] and then employs GPT [48] to predict a future pose code sequence given input music and starting seed pose. Finally, the pose code sequence is converted into a dance sequence via the learned codebook and CNN decoder. Unlike Bailando [60], our motion generation does not require multiple stages during inference.

Most of the existing methods depend on seed motion as input and usually produce repetitive dance choreography, and we differ from previous music-conditioned choreography generation methods by producing non-repetitive choreographies while also not requiring any seed motion sequence. Moreover, earlier methods [28, 29, 60] use hand-crafted music features (such as beats, chroma and onset strength) along with MFCC representation of audio signals for predicting music-aligned dance sequences. In contrast, our method learns to predict dance sequences on raw Mel spectrograms without auxiliary features like beats.

2.2. Diffusion Models

Diffusion models [61] have shown great promise in terms of generative modelling by showing outstanding results in synthesis applications ranging from image generation [21, 49, 52, 54], speech synthesis [27, 45], to point-cloud generation [33]. The seminal work of Sohl-Dickstein et al. [61] gradually diffuses Gaussian noise into a training sample and trains a neural network to reverse-diffuse the noise. Ho et al. [20] apply the same modelling technique in DDPM to achieve high-quality image synthesis, and Song et al. [62] improve the efficiency of the generative process by introducing faster sampling in the reverse process.

These models have been applied for various computer

Figure 2. An illustration of our diffusion for 3D motion synthesis. During forward diffusion, we iteratively add Gaussian noise $(M^{(t)} | M^{(t-1)}) = N(M^{(t)} | (1 - \beta_t)M^{(t-1)}; \beta_t I)$ to initial motion at $t=0$. A neural network $f(\cdot; \theta)$ is trained to denoise the noisy motion $M^{(t)}$ at time t based on the conditioning signal

vision tasks like text-to-image generation. Paradigms like classifier guidance [9] and classifier-free guidance [21] for the diffusion process have been introduced to improve image synthesis quality. CLIP-based guidance strategies are also used by GLIDE [37]. Ramesh et al. [49] also utilise text-image embeddings by CLIP and a diffusion decoder to achieve high-quality image synthesis. Other than text-to-image generation tasks, diffusion models have also been popular in other vision applications [7]. Besides image generation, diffusion models have also been applied to synthesise audio. Grad-TTS [45] and DiffWave [27] apply the diffusion paradigm to text-to-speech synthesis. Work by Luo et al. [33] also uses diffusion models for 3D point cloud generation tasks.

We note the presence of three concurrent works (published on arXiv at the time of submission) that are similar to our approach [25, 66, 70]. However, all three methods differ in their network design and loss functions. While FLAME [25], Tevet et al. [66] and Zhang et al. [70] use a transformer network, we instead choose a 1D U-Net with cross-modal transformers to learn the denoising function. We also train our network differently using a time-varying weighting schedule on the kinematic losses. Finally, all concurrent works use diffusion models to synthesise motions conditioned on text and action. On the other hand, we focus not only on text-driven motion synthesis but also on dance choreography generation using raw music

3. Method

Given a conditioning signal $c \in \mathbb{R}^{k \times d}$, our goal is to synthesise human motion $M^{(0)} = f(m_1; m_2; \dots; m_N; g)$. The pose at each time step t is parameterised as $a_t \in \mathbb{R}^{3J}$, which includes the root-relative 3D coordinates of each of the J joints and the camera-relative translation of the root joint. This representation is flexible and one could, if desired, train for joint angles instead (see supplementary materials). The conditioning signal c could either be an audio

Figure 3. Illustration of the 1D U-Net architecture with cross-modal transformer blocks with multi-head attention (bottom right). The network's input is a noisy motion sample at time step t and the output is an estimate of the noise ϵ . Additionally, it can be conditioned on either music or text prompts. In both cases, we learn a projection function to map the conditioning features to 1D U-Net features.

clip or a text prompt. It is represented as a d -dimensional embedding of Mel spectrogram features (for audio) or word tokens (for text).

In the following, we first discuss the basics of denoising diffusion models (Sec. 3.1). Next, we discuss how our kinematic losses can be incorporated within the diffusion framework (Sec. 3.2). Finally, the neural architecture design and the modifications required for task-specific conditioning are introduced (Sec. 3.3).

3.1. Diffusion for Motion Synthesis

The motion generation task is formulated as a reverse diffusion process that requires sampling a random noise vector, $z \in \mathbb{R}^N \times \mathbb{R}^J$, from a noise distribution to generate a meaningful motion sequence (see Fig. 2). While training, the forward diffusion process requires successively corrupting motion sequence $M^{(0)}$ by adding Gaussian noise to a motion sequence for T timesteps in a Markovian fashion. This results in the conversion of a meaningful motion sequence $M^{(0)}$ in the training set into a noise distribution $M^{(T)}$:

$$q(M^{(1:T)} | M^{(0)}) = \prod_{t=1}^T q(M^{(t)} | M^{(t-1)}); \quad (1)$$

where $q(M^{(t)} | M^{(t-1)}) = \mathcal{N}(M^{(t)} | (1 - \alpha_t)M^{(t-1)}; \alpha_t I)$ is the Markov diffusion kernel that adds Gaussian noise to the motion at time step t , and α_t is a hyperparameter that controls the rate of diffusion. In practice, there exists a re-parameterisation trick that allows closed-form sampling at any timestep t :

$$M^{(t)} = \sqrt{1 - \alpha_t} M^{(0)} + \sqrt{\alpha_t} \epsilon; \quad (2)$$

wherein ϵ is the random noise matrix and $\epsilon = \sum_{s=0}^{Q_t} \epsilon_s$ (1). With sufficiently large T , one can assume $\epsilon \sim z$.

To generate a motion sequence from a random noise matrix z , we need to iteratively reverse-diffuse for T timesteps. The reverse-diffusion is formulated as [61]:

$$p(M^{(0:T)}) = p(M^{(T)}) \prod_{t=1}^T p(M^{(t-1)} | M^{(t)}); \quad (3)$$

The reverse transition probability $p(M^{(t-1)} | M^{(t)})$ is approximated using a neural network that learns the function $f(M^{(t-1)} | M^{(t)}; t)$. While several variations of $f(\cdot; \cdot)$ exist, we follow [21] and train the network to predict the original noise ϵ . For the conditional synthesis setting, the network is additionally subjected to the conditioning signal as $f(M^{(t-1)} | M^{(t)}; t; c)$.

3.2. Training Objectives

We now discuss how kinematic loss terms inspired by domain knowledge can be introduced within the diffusion framework. The overall loss for training MoFusion is a weighted sum of two broad loss types:

$$L_t = L_{da} + \sum_k \lambda_k L_k; \quad (4)$$

The primary data term L_{da} , is the commonly-used L_2 distance between the noise ϵ used for forward diffusion (2) and the estimated ϵ ($M^{(t)}; t; c$).

While L_{da} is strong enough to approximate the underlying data distribution, the synthesised motions are not guaranteed to be physically and anatomically plausible. Consequently, it allows for artefacts like motion jitter, illegal skeletons and foot-sliding. Fortunately, human motion capture literature consists of several kinematic and physical constraints that can be used to regularize the synthesised motion [8, 50, 58, 72]. These kinematic loss functions are

well established in the motion synthesis literature and have nature of the network allows us to train the network with been consistently used to avoid synthesis artefacts. However, motions of various lengths. Fig. 3 illustrates the schema ever, since the denoising network is trained to estimate of the network. The network consists of three downsam- the noise, it is not straightforward to apply such con- pling blocks that first successively reduce the feature length, straints. One workaround is to apply the losses to the final motion, from N to $bN=8c$ before being upsampled using cor- reverse-diffused motion, which can be estimated using the corresponding upsampling blocks. Each 1D residual block re-parameterisation trick:

$$\hat{M}^{(0)} = \frac{1}{p} \frac{1}{t} M^{(t)} \quad (5)$$

However, naïvely using $\hat{M}^{(0)}$ to approximate the reverse-diffusion outputs leads to unstable training because the generated motion is extremely noisy when close to T .

Therefore, we introduce a time-varying weight schedule for L_k by varying the schedule as $p_k^{(t)} = \frac{1}{t}$. This ensures that the motions at T receive an exponentially lower weight compared to 0. Within $L_k = L_s + \alpha L_a + \beta L_m$, we include three loss terms: First, we use the skeleton-consistency loss, that ensures that the bone lengths in the synthesised motion remain consistent across time. To achieve this, we minimize the temporal variance of the bone lengths:

$$L_s = \frac{1}{n} \sum_{i=1}^n (l_i - \bar{l})^2; \quad (6)$$

where \bar{l} is the vector of mean bone lengths. Secondly, we use an anatomical constraint L_a , that penalizes left/right asymmetry of the bone lengths: $L_a = \sum_{j_1, j_2} |BL(j_1) - BL(j_2)|$, where $BL(\cdot)$ computes the bone-lengths between the input joints and d provides the index of the corresponding symmetrically opposite joint. When using joint angle representation instead of joint positions, it is possible to use joint-angle limit regularisations as in [23] instead of bone length constraints. Finally, we again add ground-truth supervision on motion synthesis, this time with:

$$L_m = \|\hat{M}^{(0)} - M^{(0)}\|_2; \quad (7)$$

It is worth noting that these kinematic loss terms are not exhaustive and there exist several other loss terms that can attend to different aspects of motion synthesis. For example, it is possible to add the foot-sliding loss of [56], or physics-based constraints of [50, 57, 58]. Through our formulation, we demonstrate how such losses can be incorporated within the diffusion framework

3.3. The MoFusion Architecture

Drawing inspiration from successful 1D-Convnet architectures for motion synthesis [40] and pose estimation [39], we use a 1D U-Net [53] to approximate $f(\cdot; \cdot)$. This is also consistent with several state-of-the-art diffusion-based im- age generation methods [49, 52, 54] that use a U-Net archi- tecture for the denoising network. The fully-convolutional

is followed by a cross-modal transformer that incorporates the conditioning context, into the network. The time- embedding is generated by passing the sinusoidal time em- bedding through a two-layer MLP. For incorporating the context, we treat the intermediate residual motion features, $x \in \mathbb{R}^{n \times d}$, to get the query vector while using the condi- tioning signal, $c \in \mathbb{R}^{m \times d}$, to compute the key and value vectors. Specifically, we first estimate

$$Q = W_q x; K = W_k c; \text{ and } V = W_v c; \quad (8)$$

where W_q, W_k ; and W_v are the Query, Key and Value ma- trices, respectively. As in standard cross-attention [69], the relevance scores are first computed with the softmax, and then used to weigh the values

$$\text{Attention}(Q; K; V) = \text{softmax} \left(\frac{QK^T}{d} \right) V; \quad (9)$$

In the case of unconditional generation, the formulation switches to self-attention by also getting the key and value vectors K, V from x . We now discuss the task-dependent processing of the conditioning input.

Music-to-Dance Synthesis: For conditioning the network to music signals, we choose to represent them using the Mel spectrogram representation [45, 55]. This is unlike sev- eral existing music-to-dance synthesis methods [29, 60] that use MFCC features along with music-specific features like beats or tempograms. Thus, we leave it up to the context- embedding layer to learn an appropriate projection to the feature space of U-Net. In theory, this also allows our method to be trained on other audio (not necessarily mu-

sic) conditioning such as speech. To extract the Mel spectrograms, we re-sample audio signals to 16kHz and convert them to log-Mel spectrograms with $k=80$ Mel bands by using hop-length of 12 and the minimum and maximum frequencies of 0 and 8 kHz, respectively. As a result, we obtain a conditioning signal $c \in \mathbb{R}^{(m \times k)}$, where $m=32$ for one second of the audio signal. We use a linear layer to project the input Mel spectrogram onto the context embedding

Text-to-Motion Synthesis: Text-conditioned Diffusion Models [49, 52] have recently shown impressive genera- tion capabilities. For synthesizing motion from textual de- scriptions, we use the pre-trained CLIP [47] token embed- dings. We first retrieve the tokenised embedding for each

word in the input prompt. Next, these token embeddings are position-encoded and subjected to CLIP's transformer. Finally, we project the token embeddings using an MLP that maps the transformer embedding onto

4. Experiments

We next evaluate the proposed MoFusion framework in two scenarios, i.e., conditioned by audio and text. We first discuss music-to-choreography generation details (Sec. 4.1), followed by text-conditioned motion generation (Sec. 4.2) and, finally, show applications like seed-motion forecasting, editing and inbetweening in Sec. 4.4.

4.1. Music-to-Dance Synthesis

Datasets: We train MoFusion for music-conditioned dance synthesis on the AIST++ Dataset [29]. The dataset contains 1408 unique dance motion sequences with lengths ranging from 7:4 to 48:0 seconds. There are ten different dance motion genres with multiple dance choreographies for each genre, which provides a rich diversity in terms of types of dance motions. The data has been annotated using multi-view capture and we use the provided 3D motion sequences as the target motion and their corresponding music as our conditioning input. More importantly, we use the dataset split based on music choreography, which ensures that the validation/test set contains unheard music and, correspondingly, choreography *vis-a-vis* the training set.

Evaluation Metrics: We perform the quantitative evaluation for music-conditioned synthesis by using Frechet Inception Distance (FID) score, Diversity (Div), Beat Alignment Score (BAS) and Multi-Modality. The FID score is evaluated following the method used in Siyabal. [60]. We measure and compare FID using a kinetic feature extractor [38], which includes hand-crafted features regarding velocity and acceleration in its feature representation. We use implementation by the fairmotion toolbox [14] to measure FID. To measure the diversity of generated motions, the diversity metric (Div) computes the average pairwise Euclidean distance of the kinetic features of the motions synthesised from audios in the test set. We also measure Beat Alignment Score (BAS) [29], which expresses the similarity between the kinematic and music beats. Here, kinematic beats refer to the local minima of the kinetic velocity of a motion sequence showing beats as the “stopping points” during the motion. Moreover, music beats in the audio signal are extracted using Librosa toolbox [36]. The score is defined as the mean distance between every kinematic beat and its nearest music beat:

$$BAS = \frac{1}{|B^m|} \sum_{b^d \in B^d} \exp \left(-\frac{\min_{b^m \in B^m} \|b^d - b^m\|_2}{2} \right); \quad (10)$$

where b^d represents a kinematic beat with B^d being a set

Method	Quality		Diversity	
	BAS [†]	FID #	Div [†]	M.-Modality [†]
Ground Truth	0.237	17.10	8.19	n/a
Li et al. [28]	0.160	86.43	6.85	n/a
DanceNet [73]	0.143	69.13	2.86	n/a
Dance Revolution [22]	0.195	73.42	3.52	n/a
AI Choreographer [29]	0.221	35.35	5.94	n/a
Bailando [60]	0.233	28.16	7.83	n/a
MoFusion (Ours)	0.253	50.31	9.09	11.38

Table 1. Comparison of our method with the previous methods. We achieve state-of-the-art performance on beat alignment score as well as Diversity. [†]: Unlike Bailando [60], we do not explicitly train our method using BAS as a reward or a loss function.

of all kinematic beats and B^m represents a music beat with B^m being a set of all music beats. We follow [29, 60] and keep $k = 3$ in our experiments. Finally, we also measure multi-modality for our approach by calculating the average Euclidean distance between the kinetic features of 50 generated motion sequences for the same music input. This expresses the multi-modality of the dance generation.

Quantitative Results: The quantitative results are summarised in Table 1. Our method improves upon the diversity scores of the state of the art and achieves a multi-modality score of 11.38. These results confirm the variability claims of DDM for motion synthesis. In contrast, state-of-the-art methods like Bailando and AI Choreographer are deterministic and produce similar outputs given the same input music. Therefore, measuring multi-modality is not applicable to them. In addition, we observe a better beat alignment score than Bailando [60] and the ground truth, showing that MoFusion learns better motion alignment with beats. It is also worth noting that Bailando explicitly uses BAS in its reward formulation, whereas we do not. Finally, we observe subpar performance on FID compared to [29, 60]. Upon visual inspection, we notice that both Bailando and AI-Choreographer produce repetitive, loopy dance motions which are very similar to the ground truth (and the training set). On the other hand, our diffusion-based model seldom produces repetitive or loopy motions and, therefore, differs significantly from the hand-crafted kinetic-feature profile used to compute the FID.

Analysis: Fig. 4 depicts the cross-modal attention weights of the audio signal against the generated motion. Interestingly, we observe that the transformer learns to associate high attention with the occurrence of beats in the music. Here, the beats are not provided as input features, and beat alignment is automatically learnt from the Mel spectrogram by the network. This is in contrast to methods that either explicitly use music-specific hand-crafted features [29] or train the network with a beat alignment loss [60]. Upon qualitative inspection, we also notice that, unlike other methods, our synthesised choreography rarely

Methods	Diversity	Multi-Modality	R-Precision
Real Motions	9.503	n/a	0.797
Language2Pose [1]	7.676	n/a	0.486
Text2Gesture [5]	6.409	n/a	0.345
MoCoGAN [67]	0.462	0.019	0.106
Dance2Music [28]	0.725	0.043	0.097
Guo et al. [15]	9.188	2.090	0.740
MoFusion (Ours)	8.82	2.521	0.492

Table 2. Comparison of our method with the previous state of the art on HumanML3D.

Figure 4. Visualisation of cross-modal attention weights at different levels of the U-Net. Notice the alignment of the attention weights to the specific beats in the audio. Also, while the shallower levels (top) have scattered attention, the attention heads at the bottleneck layer (bottom) degenerate to specific audio sections corresponding to the music beats.

repeats (see supplementary video). MoFusion manages to avoid this phenomenon since we do not require a seed motion input that can bias the network towards loopy motion.

4.2. Text-to-Motion Synthesis

Datasets: For the sub-task of text-to-motion synthesis, we train our method on HumanML3D [15] dataset. It consists of 28k text-annotated motion sequences from AMASS dataset [34]. Each sequence in the dataset is on average 7:1s long and has been annotated 3-4 times, thus providing a rich corpus of textual annotation for motion data. We also use the BABEL dataset [46] for qualitative evaluation that contains shorter phrase-level motion annotations.

Evaluation Metrics: Similar to dance synthesis, we evaluate the synthesised motions using the conventionally used evaluation metrics on HumanML3D dataset: Average Pairwise Euclidean Distance (Diversity) and Multi-Modality. The multi-modality metric evaluates the per-prompt diversity claims of the method by sampling the method multiple times for the same text input and computing the average pairwise Euclidean distance of the synthesised motions; a higher Euclidean distance signifies higher variations. Similarly, the diversity metric computes the average pairwise Euclidean distance between random pairs in the dataset, irrespective of the input prompt. Finally, the R-Precision score measures the classification accuracy of the synthesised motions on a pre-trained classifier [15]. However, our network represents two motions, the users are asked to answer the following motion using joint positions, whereas the classifier network requires an over-parameterised representation of motion involving joint positions, 6D joint angles, local velocities and root translation. To make our method compatible for eval-

uation, we derive the remaining inputs based on the joint positions using inverse kinematics on the estimated joints. We provide qualitative results in Fig. 1 of music-to-dance and text-to-motion synthesis results as well as in the supplementary video.

Quantitative Results: Table 2 illustrates the performance of our method on the HumanML3D dataset. Similar to the case of Music-to-Dance synthesis, our method achieves state-of-the-art results in terms of synthesis variety. This is exhibited by our performance on the multi-modality metric (2:52 vs 2:09). Further, our diversity score of 8.82 is similar to the ground truth 9.5 and second only to Guo et al.'s [15] 9:18. Recall that to calculate the R-Precision (and FID) defined in T2M, we needed to perform IK on the synthesized joint positions (1740). Naturally, this is suboptimal and prohibitively prone to jittery motion (esp. around the head joint) and manifests itself in worse R-Precision score of 0:492 compared to 0:74 of the state of the art method. We now discuss the perceptual evaluation of our results through a user study.

4.3. User Study

It is worth noting that all the evaluation metrics discussed above are imperfect performance indicators. Thus, Diversity and Multi-modality can be fooled by an untrained network that produces random, but meaningless, motion every time it is sampled. Likewise, the FID metric used in AI Choreographer [29] uses hand-crafted features and incorrectly rewards overfitting. We, therefore, conduct a user study wherein we invite participants to perceptually evaluate the quality of our synthesis. To that end, we randomly sample audio (or text) queries from the test set and present each participant with two options to choose from. One of the two options shows our synthesis, and the other option can come from either the ground truth or from other state-of-the-art methods; Bailando [60] for audio and T2M [15] and MotionDiffuse [70] for text. After having seen the two questions: "Which motion best justifies the music/text prompt?" and "Which motion looks more realistic?" This way, we evaluate the methods on their Semantic accuracy as well as Realism. Fig. 6 informs the results of the user

Figure 5. Examples of diverse motion generation for a given text prompt. Notice the variations in terms of the direction of movement as well as the difference in stances. More results, especially for choreography synthesis, can be found in the supplementary video.

motion corresponding to this snapshot.

Motion Inbetweening: In a manner similar to seed-conditioned synthesis, we perform motion inbetweening by fixing a set of keyframes in the motion sequence and reverse-diffusing the remaining frames. This application is of significant utility for virtual character animation as it provides an easy way to interpolate the keyframes.

Figure 6. The results of the user study are based on metrics of Realism and Semantics. As explained by Sec. 4.3, realism measures how realistic is the motion shown as a prompt and semantics measures how well a motion corresponds to music/text. Each bar indicates the user preference for motion generated by MoFusion compared to another motion.

study. We achieve better semantic accuracy than T2M and MotionDiffuse. It is interesting to note that our synthesis was considered more realistic than the ground-truth choreography on 51.4% of occasions. We also do well on semantics (52.3%), primarily because the ground-truth choreography consists of several basic motions in which not much dancing takes place. More comparisons can be found in our video.

4.4. Interactive Motion Editing

Seed-Conditioned Motion Synthesis: In this setting, the goal is to forecast future motion frames based on a user-provided seed sequence of a few frames. For our analysis, we consider a seed sequence of 40 frames (2s) and synthesise the future $t = 160$ frames. To achieve this, we first construct noise vectors by forward-diffusing the seed frames to produce $\epsilon^{(t)} \sim \mathcal{R}^{(S+N) \times 3J}$ for each time step t wherein the remaining frames are populated with random noise. Then, at each denoising step, we use a mask to ensure that the seed frames are not denoised. Thus, the resulting motion looks at a snapshot of the diffused seed sequence at every denoising step and generates a faithful

5. Discussion and Conclusion

Discussion: Through our analysis, we highlighted the ability of Denoising-Diffusion Probabilistic Models for conditional motion synthesis. A less-discussed aspect of MoFusion is its ability to avoid convergence to mean pose, especially since the motion is synthesised in a non-autoregressive manner. Thanks to a large latent space, it also avoids motion flicker artefacts that quantised codebook-based methods [16] are prone to. Finally, two aspects of our model that could be improved in future are 1) the inference time and 2) comparably restricted vocabulary for textual conditioning. At the same time, we foresee that MoFusion will benefit in future from fundamental advances in diffusion models and more richly annotated datasets.

Concluding Remarks: All in all, we introduced the first approach for 3D human motion synthesis based on diffusion models. The proposed MoFusion method accepts audio or textual conditioning signals and produces temporally-coherent human motion sequences that are longer, more diverse and more expressive compared to the outputs of previous approaches. Our claims are supported by thorough experiments and a user study. Moreover, MoFusion has direct applications in computer graphics, such as virtual character animation and crowd simulation. We interpret the obtained results as an encouraging step forward in cross-modal generative synthesis in computer vision.

Acknowledgements: This work was supported by the ERC Consolidator Grant 4DRPLY(770784).

References

- [1] Chaitanya Ahuja and Louis-Philippe Morency. Language2pose: Natural language grounded pose forecasting. 3DV, 2019. 2, 7
- [2] Simon Alexanderson, Gustav Eje Henter, Taras Kucherenko, and Jonas Beskow. Style-Controllable Speech-Driven Gesture Synthesis Using Normalising Flows. Computer Graphics Forum 2020. 1, 2, 3
- [3] Sadegh Aliakbarian, Pashmina Cameron, Federica Bogo, Andrew Fitzgibbon, and Tom Cashman. Flag: Flow-based 3d avatar generation from sparse observations. CVPR, 2022. 2
- [4] Andreas Aristidou, Anastasios Yiannakidis, Kr Aberman, Daniel Cohen-Or, Ariel Shamir, and Yiorgos Chrysanthou. Rhythm is a dancer: Music-driven motion synthesis with global structure. IEEE Transactions on Visualization and Computer Graphics, 2022. 3
- [5] Uttaran Bhattacharya, Nicholas Rewkowski, Abhishek Banerjee, Pooja Guhan, Aniket Bera, and Dinesh Manocha. Text2gestures: A transformer-based network for generating emotive body gestures for virtual agents. 7
- [6] Richard Bowden. Learning statistical models of human motion. In Proceedings of CVPR - IEEE Workshop on Human Modeling, Analysis and Synthesis, 2000. 2
- [7] Florinel-Alin Croitoru, Vlad Hondru, Radu Tudor Ionescu, and Mubarak Shah. Diffusion models in vision: A survey. arXiv, 2022. 3
- [8] Rishabh Dabral, Anurag Mundhada, Uday Kusupati, Safeer Afaq, Abhishek Sharma, and Arjun Jain. Learning 3d human pose from structure and motion. ECCV, 2018. 4
- [9] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. NeurIPS 2021. 3
- [10] Katerina Fragkiadaki, Sergey Levine, Panna Felsen, and Jitendra Malik. Recurrent network models for human dynamics. In ICCV, 2015. 2
- [11] Aphrodite Galata, Neil Johnson, and David Hogg. Learning variable-length markov models of behavior. Computer Vision and Image Understanding, 2001. 2
- [12] Anindita Ghosh, Noshaba Cheema, Cennet Oguz, Christian Theobalt, and P. Slusallek. Synthesis of compositional animations from textual descriptions. ICCV, 2021. 1, 2
- [13] Shiry Ginossar, Amir Bar, Gefen Kohavi, Caroline Chan, Andrew Owens, and Jitendra Malik. Learning individual styles of conversational gesture. CVPR 2019. 3
- [14] Deepak Gopinath and Jungdam Won. fairmotion - tools to load, process and visualize motion capture data, 2020. 6
- [15] Chuan Guo, Shihao Zou, Xinxin Zuo, Sen Wang, Wei Ji, Xingyu Li, and Li Cheng. Generating diverse and natural 3d human motions from text. ICCVPR 2022. 2, 7
- [16] Chuan Guo, Xinxin Zuo, Sen Wang, and Li Cheng. Tm2t: Stochastic and tokenized modeling for the reciprocal generation of 3d human motions and texts. ECCV, 2022. 8
- [17] Chuan Guo, Xinxin Zuo, Sen Wang, Shihao Zou, Qingyao Sun, Annan Deng, Minglun Gong, and Li Cheng. Action2motion: Conditioned generation of 3d human motions. In ACM MM, 2020. 1, 2
- [18] Ikhsanul Habibie, Weipeng Xu, Dushyant Mehta, Lingjie Liu, Hans-Peter Seidel, Gerard Pons-Moll, Mohamed Elgharib, and Christian Theobalt. Learning speech-driven 3d conversational gestures from video. IVA, 3
- [19] Alejandro Hernandez, Jurgen Gall, and Francesc Moreno-Noguer. Human motion prediction via spatio-temporal inpainting. In ICCV. 2
- [20] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In NeurIPS 2020. 2, 3
- [21] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications, 2021. 2, 3, 4
- [22] Ruozhi Huang, Huang Hu, Wei Wu, Kei Sawada, Mi Zhang, and Daxin Jiang. Dance revolution: Long-term dance generation with music via curriculum learning. ICLR, 2021. 3, 6
- [23] Angjoo Kanazawa, Michael J. Black, David W. Jacobs, and Jitendra Malik. End-to-end recovery of human shape and pose. In CVPR 2018. 5
- [24] Hsuan-Kai Kao and Li Su. Temporally guided music-to-body-movement generation. ACM MM, 2020. 3
- [25] Jihoon Kim, Jiseob Kim, and Sungjoon Choi. Flame: Free-form language-based motion synthesis editing. 2023. 3
- [26] Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes. In ICLR, 2014. 2
- [27] Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. Diffwave: A versatile diffusion model for audio synthesis. ICLR, 2021. 3
- [28] Jiaman Li, Yihang Yin, Hang Chu, Yi Zhou, Tingwu Wang, Sanja Fidler, and Hao Li. Learning to generate diverse dance motions with transformer. ArXiv, abs/2008.08171, 2020. 3, 6, 7
- [29] Ruilong Li, Sha Yang, David A. Ross, and Angjoo Kanazawa. Ai choreographer: Music conditioned 3d dance generation with aist++. In ICCV, 2021. 1, 2, 3, 5, 6, 7
- [30] Xiaodan Liang, Lisa Lee, Wei Dai, and Eric P. Xing. Dual motion gan for future-flow embedded video prediction. In ICCV, 2017. 2
- [31] Angela S. Lin, Lemeng Wu, and Qixing Huang Raymond J. Mooney Rodolfo Corona, Kevin Tai. Generating animated videos of human activities from natural language descriptions. In Proceedings of the Visually Grounded Interaction and Language Workshop at NeurIPS, 2018. 2
- [32] Hung Yu Ling, Fabio Zinno, George Cheng, and Michiel van de Panne. Character controllers using motion vectors. ACM TOG, 2020. 1
- [33] Shitong Luo and Wei Hu. Diffusion probabilistic models for 3d point cloud generation. ICCVPR 2021. 3
- [34] Naureen Mahmood, Nima Ghorbani, Nikolaus F. Troje, Gerard Pons-Moll, and Michael J. Black. AMASS: Archive of motion capture as surface shapes. ICCV, 2019. 7
- [35] Julieta Martinez, Michael J. Black, and Javier Romero. On human motion prediction using recurrent neural networks. In CVPR 2017. 1, 2
- [36] Brian McFee, Colin Raffel, Dawen Liang, Daniel PW Ellis, Matt McVicar, Eric Battenberg, and Oriol Nieto. librosa: Audio and music signal analysis in python. Proceedings of the 14th python in science conference, 2016. 6

- [37] Alexander Quinn Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. GLIDE: towards photorealistic image generation and editing with text-guided diffusion models. In *ICML*, 2022. 2, 3
- [38] Kensuke Onuma, Christos Faloutsos, and Jessica K. Hodgins. FMDistance: A Fast and Effective Distance Function for Motion Capture Data. In *Eurographics 2008 - Short Papers*, 2008. 6
- [39] Dario Pavullo, Christoph Feichtenhofer, David Grangier, and Michael Auli. 3d human pose estimation in video with temporal convolutions and semi-supervised training. *CVPR*, 2019. 5
- [40] Dario Pavullo, David Grangier, and Michael Auli. Quaternet: A quaternion-based recurrent model for human motion. In *BMVC*, 2018. 5
- [41] Xue Bin Peng, Ze Ma, Pieter Abbeel, Sergey Levine, and Angjoo Kanazawa. Amp: Adversarial motion priors for stylized physics-based character control. *ACM TOG*, 2021. 1
- [42] Mathis Petrovich, Michael J. Black, and Michael J. Black. Action-conditioned 3D human motion synthesis with transformer VAE. In *ICCV*, 2021. 1, 2
- [43] Mathis Petrovich, Michael J. Black, and Michael J. Black. TEMOS: Generating diverse human motions from textual descriptions. In *ECCV*, 2022. 1, 2
- [44] Matthias Plappert, Christian Mandery, and Tamim Asfour. The KIT motion-language dataset. *Big Data*, 2016. 2
- [45] Vadim Popov, Ivan Vovk, Vladimir Gogoryan, Tasnima Sadekova, and Mikhail A. Kudinov. Grad-tts: A diffusion probabilistic model for text-to-speech. 2021. 3, 5
- [46] Abhinanda R. Punnakal, Arjun Chandrasekaran, Nikos Athanasiou, Alejandra Quiros-Ramirez, and Michael J. Black. BABEL: Bodies, action and behavior with english labels. In *CVPR*, 2021. 2, 7
- [47] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. 2021. 5
- [48] Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. 2019. 3
- [49] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv*, 2022. 2, 3, 5
- [50] Davis Rempe, Tolga Birdal, Aaron Hertzmann, Jimei Yang, Srinath Sridhar, and Leonidas J. Guibas. Humor: 3d human motion model for robust pose estimation. In *ICCV*, 2021. 1, 2, 4, 5
- [51] Xuanchi Ren, Haoran Li, Zijian Huang, and Qifeng Chen. Self-supervised dance video synthesis conditioned on music. In *ACM MM*, 2020. 3
- [52] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2021. 2, 3, 5
- [53] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, 2015. 5
- [54] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Rapha Gontijo Lopes, Tim Salimans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv*, 2022. 3, 5
- [55] Jonathan Shen, Ruoming Pang, Ron J. Weiss, Mike Schuster, Navdeep Jaitly, Zongheng Yang, Zhifeng Chen, Yu Zhang, Yuxuan Wang, Rj Skerrv-Ryan, Rif A. Saurous, Yannis Agiomvrgiannakis, and Yonghui Wu. Natural tts synthesis by conditioning wavenet on mel spectrogram predictions. In *ICASSP*, 2018. 5
- [56] Mingyi Shi, K r Aberman, Andreas Aristidou, Taku Komura, Dani Lischinski, Daniel Cohen-Or, and Baoquan Chen. Motionet: 3d human motion reconstruction from monocular video with skeleton consistency. *ACM TOG*, 2020. 5
- [57] Soshi Shimada, Vladislav Golyanik, Weipeng Xu, Patrick Pérez, and Christian Theobalt. Neural monocular 3d human motion capture with physical awareness. *ACM TOG*, 2021. 5
- [58] Soshi Shimada, Vladislav Golyanik, Weipeng Xu, and Christian Theobalt. Physcap: Physically plausible monocular 3d motion capture in real time. *ACM TOG*, 2020. 4, 5
- [59] Eli Shlizerman, Lucio M Dery, Hayden Schoen, and Ira Kemelmacher-Shlizerman. Audio to body dynamics. In *CVPR*, 2017. 3
- [60] Li Siyao, Weijiang Yu, Tianpei Gu, Chunze Lin, Quan Wang, Chen Qian, Chen Change Loy, and Ziwei Liu. Bailando: 3d dance generation via actor-critic gpt with choreographic memory. In *CVPR*, 2022. 1, 3, 5, 6, 7
- [61] Jascha Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. 2015. 3, 4
- [62] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *ICLR*, 2022. 3
- [63] Sebastian Starke, Ian Mason, and Taku Komura. Deepphase: Periodic autoencoders for learning motion phase manifolds. *ACM TOG*, 2022. 3
- [64] Sebastian Starke, He Zhang, Taku Komura, and Jun Saito. Neural state machine for character-scene interaction. *ACM TOG*, 2019. 1
- [65] Guy Tevet, Brian Gordon, Amir Hertz, Amit H Bermano, and Daniel Cohen-Or. Motionclip: Exposing human motion generation to clip space. In *ECCV*, 2022. 2
- [66] Guy Tevet, Sigal Raab, Brian Gordon, Yonatan Shalrit, Amit H Bermano, and Daniel Cohen-Or. Human motion diffusion model. In *ICLR*, 2022. 3
- [67] Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, and Jan Kautz. MoCoGAN: Decomposing motion and content for video generation. In *CVPR*, 2018. 7
- [68] Aaron van den Oord, Oriol Vinyals, and koray kavukcuoglu. Neural discrete representation learning. In *NeurIPS*, 2017. 3
- [69] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017. 5

- [70] Mingyuan Zhang, Zhongang Cai, Liang Pan, Fangzhou Hong, Xinying Guo, Lei Yang, and Ziwei Liu. Motiondiffuse: Text-driven human motion generation with diffusion model. arXiv preprint arXiv:2208.15001, 2022. 3, 7
- [71] Qiu Zhou, Manyi Li, Qiong Zeng, Andreas Aristidou, Xiaojing Zhang, Lin Chen, and Changhe Tu. Let's all dance: Enhancing amateur dance motion. Computational Visual Media, 2022. 3
- [72] Xingyi Zhou, Qixing Huang, Xiao Sun, Xiangyang Xue, and Yichen Wei. Towards 3d human pose estimation in the wild: A weakly-supervised approach. ICCV, 2017. 4
- [73] Wenlin Zhuang, Congyi Wang, Jinxiang Chai, Yangang Wang, Ming Shao, and Siyu Xia. Music2dance: Dancenet for music-driven dance generation. ACM Trans. Multimedia Comput. Commun. App, 2022. 3, 6