

# Single Motion Diffusion

Sigal Raab\*, Inbal Leibovitch\*, Guy Tevet, Moab Arar, Amit H. Bermano, and Daniel Cohen-Or

Tel-Aviv University, Israel

sigalraab@tauex.tau.ac.il, {inball1, guytevet}@mail.tau.ac.il

Synthesizing realistic animations of humans, animals, and even imaginary creatures, has long been a goal for artists and computer graphics professionals. Compared to the imaging domain, which is rich with large available datasets, the number of data instances for the motion domain is limited, particularly for the animation of animals and exotic creatures (e.g., dragons), which have unique skeletons and motion patterns. In this work, we present a Single Motion Diffusion Model, dubbed SinMDM, a model designed to learn the internal motifs of a single motion sequence with arbitrary topology and synthesize motions of arbitrary length that are faithful to them. We harness the power of diffusion models and present a denoising network explicitly designed for the task of learning from a single input motion. SinMDM is designed to be a lightweight architecture, which avoids overfitting by using a shallow network with local attention layers that narrow the receptive field and encourage motion diversity. SinMDM can be applied in various contexts, including spatial and temporal in-betweening, motion expansion, style transfer, and crowd animation. Our results show that SinMDM outperforms existing methods both in quality and time-space efficiency. Moreover, while current approaches require additional training for different applications, our work facilitates these applications at inference time. Our code and trained models are available at <https://sinmdm.github.io/SinMDM-page>.

## 1 INTRODUCTION

3D character animation is a long pursued task in computer graphics with many applications, from the big screen to virtual reality headsets. It is notoriously known as a time-consuming task done by expert artists. In recent years, neural models have offered faster and less expensive tools for modeling motion [Holden et al. 2016; Petrovich et al. 2022; Raab et al. 2022]. In particular, the very recent adaptation of diffusion models into the motion domain provides unprecedented results in both quality and diversity [Kim et al. 2022; Tevet et al. 2022b].

These data-driven methods typically require large amounts of data for training. However, motion data is quite scarce and, moreover, for a non-human skeleton, it is barely existent. The few available datasets contain humanoids only, whose topology and bone ratio are fixed. Animators often customize a skeleton per character (human, animal, or magical creature), for which common data-driven techniques are irrelevant.

In this work, we present a Single Motion Diffusion Model, dubbed SinMDM, that trains on a *single motion* input sequence. Our model enables modeling motions of arbitrary skeletal topology, which often have no more than one animation sequence to learn from. SinMDM can synthesize a variety of variable-length motion sequences that retain the core motion elements of the input and can handle complex and non-trivial skeletons. For example, our model can generate a diverse clan based on one flying dragon or one hopping ostrich.

Learning from a single instance has been explored for the imaging domain [Shaham et al. 2019; Shocher et al. 2019] and for the motion domain [Li et al. 2022], using the GAN architecture [Goodfellow

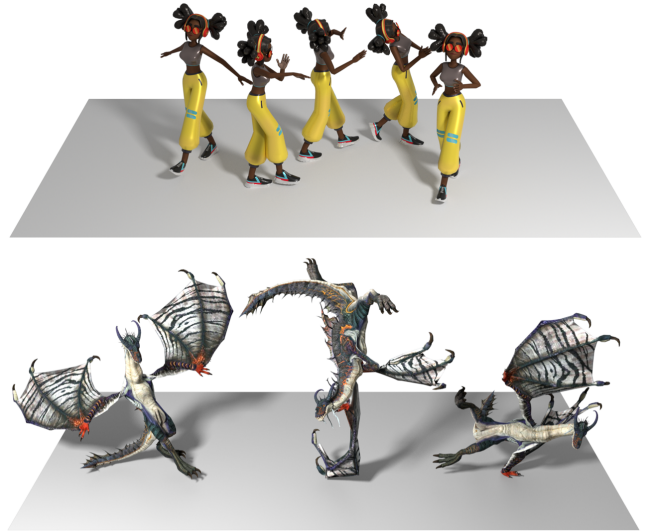


Fig. 1. SinMDM learns the internal motion motifs from a single motion sequence with arbitrary topology and synthesizes motions that are faithful to the learned core motifs of the input sequence. Top: a girl exercising while walking. Bottom: a breakdancing dragon. Left to right: breakdance uprock, breakdance freeze, and breakdance flair.

et al. 2014]. Indeed until recently, GANs have been the dominant approach for generative models. We find diffusion models [Ho et al. 2020] more suitable for single input learning, as the descriptive ability attained by gradual denoising yields a lightweight model that, compared to prior art, is simpler in architecture and more efficient in terms of the number of parameters and training time. Furthermore, we demonstrate that diffusion models can be effectively utilized with limited data, challenging the notion that they solely rely on large datasets.

To learn local motion sequences, the receptive field must be small enough, analogously to the use of patch-based discriminators [Isola et al. 2017; Li and Wand 2016] in GAN-based techniques. The use of a narrow receptive field (Fig. 2) promotes diversity and reduces overfitting. We show the importance of narrow receptive fields in our ablation studies.

Most motion diffusion models use transformers. However, vanilla transformers are not suitable for learning a single sequence, as their receptive field encompasses the entire motion. A similar challenge exists in the UNet architecture, which is common for image diffusion models. Its depth, combined with global attention layers, induces a receptive field that covers the whole motion.

SinMDM leverages the concept of narrow receptive fields and introduces a motion architecture specifically designed with this concept. It combines a shallow UNet [Ronneberger et al. 2015] model adapted for motion with a QnA [Arar et al. 2022] local attention

\*equal contribution.

mechanism, instead of global attention. As a result, SinMDM outperforms prior art both quantitatively and qualitatively, and demonstrates high efficiency with shorter training time and less memory consumption. Imputed to its lightweight architecture, SinMDM can be trained on a single mid-range GPU.

We present many use cases of SinMDM. While prior works require designated training per application, ours are applied at inference time, with no need to re-train. Moreover, applications that require different dedicated algorithms in prior art, are here grouped together as special cases of the same technique, significantly simplifying their use. One of the applications we present is *Motion Composition*, where a given motion sequence is composed jointly with a synthesized one, either temporally or spatially. Special cases of motion composition include in-betweening and motion expansion. Another application that we present is *Harmonization*, along with its special case, style transfer. Here, a reference motion is modified to match the learned motion motifs. It should be emphasized that implementing style transfer using a denoising model is a non-trivial task, and enabling it through motion harmonization is unique. We further present two more applications: *long sequence generation* and *crowd animation*.

In our presented work, we suggest two comprehensive benchmarks for single-motion evaluation. The first is built upon the artistically crafted MIXAMO [2021] dataset, utilizing metrics that do not require an additional feature-extracting model. The second is based on the HumanML3D [2022] dataset and enables metrics that use latent features, such as single-FID. We show that our model outperforms current works on both benchmarks.

## 2 RELATED WORK

### 2.1 Single-Instance Learning

The goal of single-instance generation is to learn an unconditional generative model from a single instance and generate diverse samples with similar content by capturing the internal statistics of patches. The type of instance depends on the input domain. Most single-instance learning research has been focused on the domain of imaging. The first works on this topic are SinGAN [Shaham et al. 2019] and InGAN [Shocher et al. 2019]. SinGAN uses a patch-based discriminator [Isola et al. 2017; Li and Wand 2016] and an image pyramid to generate diverse results hierarchically. InGAN [Shocher et al. 2019], uses a conditional GAN to solve the same problem using geometry transformation. More recent approaches include ExSinGAN [Zhang et al. 2021b], which trains multiple modular GANs to model the distribution of structure, semantics, and texture, and ConSinGAN [Hinz et al. 2021], which trains several stages sequentially and improves SinGAN. Many works in the imaging domain follow and improve the aforementioned pioneering works [Asano et al. 2020; Chen et al. 2021; Granot et al. 2022; Lin et al. 2020; Sun and Liu 2020; Sushko et al. 2021; Yoo and Chen 2021; Zhang et al. 2022b; Zheng et al. 2021].

Several works have been introduced in other domains, such as for the shapes domain [Son et al. 2022] and for the 3D scenes domain [Son et al. 2022]. In the motion domain, the only work that learns a single motion is Ganimator [Li et al. 2022]. Ganimator follows SinGAN, hence it uses a GAN architecture, with a patch-based discriminator and a temporal pyramid.

The vast majority of single-instance learning works use a GAN architecture [Goodfellow et al. 2014]. Until recently, GANs have been the dominant approach for generative models. However, we are currently seeing a trend towards using diffusion models [Ho et al. 2020; Song et al. 2020a] as an alternative to GANs.

A number of concurrent works in the imaging domain use diffusion models to learn from single images. Similar to our approach, Wang et al. [2022] and Nikankin et al. [2022] drop the image pyramid structure and use a UNet [Ronneberger et al. 2015] with limited depth. A different work [Kulikov et al. 2022] constructs a multi-scale diffusion process from down-sampled versions of the training image, as well as their blurry versions.

Ganimator [Li et al. 2022] is our immediate comparison reference, as it is the only single-motion learning work. Sec. 6 and our supplementary video show that SinMDM outperforms it quantitatively and qualitatively. In addition, Ganimator uses a complex architecture that combines a temporal hierarchy of motions with a skeletal hierarchy of joints. Our model uses neither hierarchies, which makes it simple in architecture and efficient in time and space, while achieving better results.

### 2.2 Diffusion Models

Diffusion models use a stochastic diffusion process, as modeled in thermodynamics [Sohl-Dickstein et al. 2015; Song and Ermon 2020], to generate samples from a data distribution. These models are adapted for image generative applications. Dhariwal and Nichol [2021] introduce the concept of classifier-guided diffusion for conditional generation, which is later adapted in the GLIDE [Nichol et al. 2022] model. Ho and Salimans [2022] propose the Classifier-Free Guidance approach, which can trade-off between fidelity and diversity in the generated samples. This approach has been demonstrated to achieve better results compared to other methods, as shown by Nichol et al. [2022].

Local editing of images may be viewed as an inpainting problem, in which a portion of the image is held constant while the model denoises the remaining part [Saharia et al. 2022; Song et al. 2020b]. In our work, we adapt this technique for motion composition of specific body parts or temporal intervals.

In the motion domain, several very recent works [Kim et al. 2022; Tevet et al. 2022b; Zhang et al. 2022a] introduce diffusion-based synthesis, where the most prominent one is MDM [Tevet et al. 2022b]. MDM utilizes a lightweight network, uses a transformer rather than the common UNet and predicts motion rather than noise. Its general design has already been used for various motion applications [Shafir et al. 2023; Tseng et al. 2022; Yuan et al. 2022]. Like MDM, SinMDM presents a lightweight architecture and predicts motion rather than noise. However, unlike MDM, our work uses a QnA-based UNet rather than a transformer, as the receptive field of a transformer is the full motion, inducing over-fitting.

### 2.3 Motion Synthesis Models

In recent years, we witness prosperity in the domain of motion synthesis using neural networks [Holden et al. 2016, 2015]. Most of these models focus on specific motion-related tasks, conditioned on some limiting factors, which can be high-level guidance such

as action [Cervantes et al. 2022; Guo et al. 2020; Petrovich et al. 2021; Tevet et al. 2022b] or text [Ahuja and Morency 2019; Bhat-tacharya et al. 2021; Guo et al. 2022; Petrovich et al. 2022; Tevet et al. 2022a,b; Zhang et al. 2021c], can be parts of a motion such as motion prefix [Aksan et al. 2019; Barsoum et al. 2018; Habibie et al. 2017; Hernandez et al. 2019; Yuan and Kitani 2020; Zhang et al. 2021a] or in-betweening [Duan et al. 2021; Harvey and Pal 2018; Harvey et al. 2020; Kaufmann et al. 2020], motion retargeting or style transfer [Aberman et al. 2020a,b, 2019; Holden et al. 2017; Villegas et al. 2018], and even music [Aristidou et al. 2022; Lee et al. 2018; Li et al. 2021; Sun et al. 2020]. Fewer models are fully unconditioned [Holden et al. 2016; Raab et al. 2022; Starke et al. 2022] and they learn the motion manifold from the input data in an unsupervised manner.

The architecture of motion synthesis models can be roughly divided into recurrent [Fragkiadaki et al. 2015; Ghorbani et al. 2020; Habibie et al. 2017; Zhou et al. 2018], auto encoder based [Guo et al. 2020; Jang and Lee 2020; Maheshwari et al. 2022; Petrovich et al. 2021], GAN based [Degardin et al. 2022; Wang et al. 2020; Yan et al. 2019; Yu et al. 2020], normalizing flows based [Henter et al. 2020], and more recently, neural field based [He et al. 2022] and diffusion based [Kim et al. 2022; Tevet et al. 2022b; Zhang et al. 2022a]. Our work belongs to the latter category.

### 3 PRELIMINARY

In this work, we present SinMDM, a novel framework that learns the internal motion motifs of a *single motion* of arbitrary topology, and generates a variety of synthesized motions that retain the core motion elements of the input sequence.

At the crux of our approach lays a denoising diffusion probabilistic model (DDPM) [Ho et al. 2020]. We consider diffusion models to be more appropriate for single input learning compared to previous methods and suggest a lightweight model, efficient in time and space and simple in architecture. This is achieved through the gradual denoising process, which enhances the model’s descriptive capability. Our generative network is a UNet [Ronneberger et al. 2015] whose attention layers are replaced by the recently introduced QnA layers [Arar et al. 2022].

In the rest of this section, we briefly recap DDPM and describe our motion representation. In the following section, we describe our method and focus on our design choices. Next, we describe various applications enabled by SinMDM (Sec. 5), detail the experiments conducted to validate our approach (Sec. 6), and summarise with conclusions (Sec. 7). The readers are encouraged to watch the supplementary video in order to get a full impression of our results.

#### 3.1 Denoising Diffusion Probabilistic Models (DDPM)

DDPMs [2020] have become the de-facto leading generative networks technique. While they have primarily dominated the imaging domain [Dhariwal and Nichol 2021], recent works have successfully applied this approach in the motion domain [Tevet et al. 2022b; Zhang et al. 2022a]. Denoising networks learn to convert unstructured noise to samples from a given distribution, through an iterative process of progressively removing small amounts of Gaussian noise.

Given an input motion sequence  $x_0$ , we apply a Markov noising process of  $T$  steps,  $\{x_t\}_{t=0}^T$ , such that

$$q(x_t|x_{t-1}) = \mathcal{N}(\sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)I), \quad (1)$$

where  $\alpha_t \in (0, 1)$  are constant hyper-parameters. When  $\alpha_t$  is small enough, we can approximate  $x_T \sim \mathcal{N}(0, I)$ .

We apply unconditional motion synthesis that models  $x_0$  as the reversed diffusion process of gradually cleaning  $x_T$ , using a generative network  $p_\theta$ . Following Tevet et al. [2022b] we choose to predict the input motion, denoted  $\hat{x}_0$  [Ramesh et al. 2022] rather than predicting  $\epsilon_t$ , hence

$$\hat{x}_0 = p_\theta(x_t, t). \quad (2)$$

We apply the widespread diffusion loss, via

$$\mathcal{L}_{\text{simple}} = \mathbb{E}_{t \sim [1, T]} \|x_0 - p_\theta(x_t, t)\|_2^2. \quad (3)$$

Synthesis at inference time is applied through a series of iterations, starting with pure noise  $x_T$ . In each iteration, a clean version of the current sample  $x_t$  is predicted using a generator  $p_\theta$ . This predicted clean sample  $\hat{x}_0$  is then "re-noised" to create the next sample  $x_{t-1}$ , with the process being repeated until  $t = 0$  is reached.

#### 3.2 Motion representation

A motion sequence is represented by its dynamic and static features, **D** and **S**, respectively. The former differ at each temporal frame (e.g., joint rotation angles), while the latter is temporally fixed (e.g., bone lengths). **D** and **S** can be combined into global 3D pose sequences using *forward kinematics* (FK). In our research, we focus on synthesizing the *dynamic features*, leaving the static features intact. That is, we predict dynamics for a fixed skeleton topology with fixed bone lengths. For simplicity, we use the term *motion synthesis* for the generation of dynamic features only.

Let  $N$  denote the number of frames in a motion sequence, and  $F$  denote the length of the features describing a single frame. In the HumanML3D dataset, a frame is redundantly represented with the root position and joint positions, angles, velocities, and foot contact [Guo et al. 2022]. For the other datasets used in this work, a frame is represented by joint angles, root positions, and foot contact labels. We represent the dynamic features of a motion by a tensor  $\mathbf{D} \in \mathbb{R}^{N \times F}$ . Naturally, the convolution for this representation is 1D, convolving on the temporal dimension (of size  $N$ ) and holding  $F$  features. Let  $J$  denote the number of skeletal joints, and let  $Q$  denote the number of rotational features, where rotational features may be Euler angles, quaternions, 6D rotations, etc. Let  $C$  denote the number of joints that are prone to contact the ground. Clearly, a human, a spider, and a snake possess different values of  $C$ .

When using the HumanML3D [Guo et al. 2022] dataset, we adhere to its representation, in which a single pose is defined by

$$p = (\dot{r}^a, \dot{r}^x, \dot{r}^z, r^y, j^p, j^v, j^r, c^f) \in \mathbb{R}^F,$$

where  $\dot{r}^a \in \mathbb{R}$  is the root angular velocity along the Y-axis.  $\dot{r}^x, \dot{r}^z \in \mathbb{R}$  are root linear velocities on the XZ-plane, and  $r^y \in \mathbb{R}$  is the root height.  $j^p \in \mathbb{R}^{3J}$ ,  $j^v \in \mathbb{R}^{3J}$  and  $j^r \in \mathbb{R}^{6J}$  are the local joint positions, velocities, and rotations with respect to the root, and  $c^f \in \mathbb{R}^4$  are binary features denoting the foot contact labels for four foot joints (two for each leg).

When using data from other datasets, we adhere to the representation used by Ganimator [2022], so we can conduct a fair comparison with their results. Their representation consists of a 3D root location, a rotation angle for each joint, and foot contact labels. Altogether, for a general representation  $D \in \mathbb{R}^{N \times F}$ , we have got  $F = 3 + JQ + C$ .

The rotations in both representations are defined in the coordinate frame of their parent in the kinematic chain, and are represented by the 6D rotation features ( $Q = 6$ ) Zhou et al. [2019], which yields the best result in many works [Petrovich et al. 2021; Qin et al. 2022].

A growing number of works use foot contact labels [Gordon et al. 2022; Raab et al. 2022] to mitigate common foot sliding artifacts. Let  $C$  denote the set of joints that contact the ground in the subject whose motion is being learned such that  $C = |C|$ . The foot contact labels are represented by  $L \in \{0, 1\}^{N \times C}$ .

When a dataset provides foot contact label information [Guo et al. 2022], we use it as is. When a dataset does not provide them, we calculate it as done by Li et al. [2022], via

$$\forall j \in C, n \in [1, N] : \quad L^{nj} = \mathbf{1}[\|\Delta_n \text{FK}([D, S])^{nj}\|_2 < \epsilon], \quad (4)$$

where  $\Delta_n \text{FK}([D, S])^{nj}$  denotes the velocity of joint  $j$  in frame  $n$  retrieved by a forward kinematics operator, and  $\mathbf{1}[\cdot]$  is an indicator function.

## 4 GENERATIVE NETWORK

Our goal is to construct a model that can generate a variety of synthesized motions that retain the core motion motifs of a single learned input sequence. More formally, we would like to construct the generative network  $p_\theta$  (Eq. 2) that synthesizes a motion  $\hat{x}_0$  from a noised motion  $x_t$ .

While traditional single-instance techniques use a pyramid of down-sampled instances (images or motions) and learn in a coarse-to-fine fashion, our model introduces a simple architecture that requires no pyramids.

Our key insight is that internal motifs are learned more effectively with a limited receptive field (Fig. 2 Left). We design SinMDM, a novel generative architecture, accordingly. Our model is a QnA-based degenerated UNet (Fig. 2 Right). The UNet architecture [Ronneberger et al. 2015] is frequently used by diffusion models in the imaging domain [Nichol et al. 2022]. However, training a UNet on a single input leads to significant overfitting due to its large receptive field, resulting in synthesized sequences that closely resemble the input.

Our first design choice in mitigating this issue is to decrease the depth of the UNet and thereby limit the receptive field width. However, this step alone is not enough, since standard UNets employ global attention layers, resulting in a receptive field that encompasses the entire sequence. A possible solution would be using local attention in non-overlapping windows, like in ViT [Dosovitskiy et al. 2021]. Nonetheless, non-interleaving windows tend to limit the cross-window interaction, compromising the model’s performance. Our solution is to use QnA [Arar et al. 2022], a state-of-the-art shift-invariant local attention layer, that aggregates the input locally in an overlapping manner, much like convolutions, but with the expressive power of attention. The key idea behind QnA is to introduce learned queries, shared by all windows, allowing fast and efficient

implementation. In particular, QnA enables local attention with a temporally narrow receptive field. Our QnA-based UNet is the first to be used in the motion domain, where we plug QnA layers instead of the global attention layers of a vanilla UNet. QnA is substantially more efficient than global attention in terms of space and time, and our model benefits from this advantage as a byproduct. A detailed description of the QnA layers is available in Appendix B.

In Sec. 6, we validate these design choices. We show the effectiveness of a narrow receptive field, and justify the usage of QnA layers and the choice of a UNet rather than a transformer.

In Appendix A, we provide a comprehensive list of the hyperparameters that can be used to reconstruct our results.

## 5 APPLICATIONS

Single-motion learning using diffusion models enables various applications. All our applications are applied at inference time, with no need to re-train. This is in contrast to the only current single-motion synthesis work, Ganimator [2022], which requires specialized training for most of its applications. To meet paper length limits and given the variety of potential applications, we illustrate six selected ones. Note that applications that require different dedicated algorithms in prior art, are grouped together here as special cases of the same technique, significantly simplifying their use.

In the following, we show *Motion Composition* (Sec. 5.1), where a given motion sequence is composed jointly with a synthesized one, either temporally or spatially. Special cases of motion composition include *in-betweening*, *motion expansion*, *trajectory control*, and *joints control*. With our *Motion Harmonization* (Sec. 5.2), a reference input motion is altered to align with the learned motion motifs. We illustrate one important special case, *style transfer*. Lastly, we show how straightforward use (Sec. 5.3) of SinMDM enables *one shot long motion generation* and *crowd animation*. The applications presented here are also demonstrated in our supplementary video.

### 5.1 Motion Composition

Given a reference motion sequence  $y$ , and a region of interest (ROI) mask  $m$ , our goal is to synthesize a new motion  $\hat{x}_0$ , such that the regions of interest  $\hat{x}_0 \odot m$  are synthesized from random noise, while the complementary area remains as close as possible to the given motion  $y$ , i.e.,  $y \odot (1 - m) \approx \hat{x}_0 \odot (1 - m)$ , where  $\odot$  is element-wise multiplication. The model should output a coherent motion sequence, where the transition between given and synthesized parts is seamless. Moreover, the reference motion can be an arbitrary one, on which our model has *not* been trained.

When using a binary mask [Avrahami et al. 2022], as the reference motion  $y$  deviates from the motion the model was trained on, the blending between the given and synthesized parts becomes less smooth. To mitigate this issue, we change the ROI mask such that the borders between the given and the synthesized motion segments are linearly interpolated, as depicted in Fig. 3.

We fix the motion segments that need to remain unchanged and sample the parts that need to be filled in. Each step of the iterative inference process (described in Sec. 3.1) is slightly changed, such that parts of  $y$  are assigned into  $\hat{x}_0$  according to the indices of the mask. That is,  $\hat{x}_0 \odot (1 - m) \leftarrow y \odot (1 - m)$ .

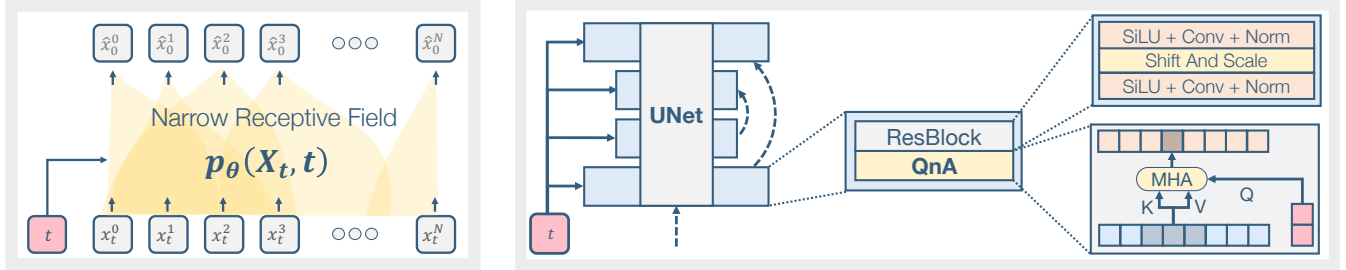


Fig. 2. Left: To allow training on a single motion, our denoising network is designed such that its overall receptive field covers only a portion of the input sequence. This effectively allows the network to simultaneously learn from multiple local temporal motion segments. Our denoiser predicts the input sequence from a noisy one.  $x_t^0 \dots x_t^N$  is a motion of  $N$  frames at diffusion step  $t$ . Right: Our network is a shallow UNet, enhanced with a QnA local attention layer.

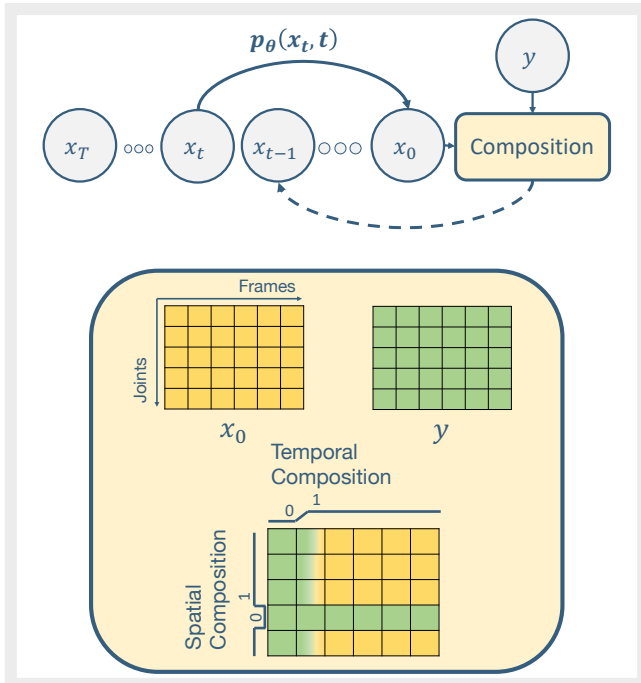


Fig. 3. Motion composition. Parts from a reference motion  $y$ , are composed with the synthesized motion  $\hat{x}_0$ , according to a composition map.

*Temporal composition – use cases: in-betweening, motion expansion.* Temporal composition is the action of filling in selected frame sequences. *In-betweening* [Harvey et al. 2020], depicted in Fig. 4, is a special case of temporal composition, where the filled-in part is at the temporal interior of the sequence, and the reference  $y$  is from the same distribution as the learned motion. Another special case of temporal composition is *motion expansion*, the motion domain’s equivalent of image outpainting [Lin et al. 2021; Teterwak et al. 2019; Yu et al. 2019], where the model generates content that resides beyond the edges of a reference motion sequence. In the case of motion expansion, the ROI mask is zeroed in the center frames, and assigned ones in the outer regions. See Fig. 5.

*Spatial composition – use cases: trajectory control, joints control.* Motion composition can be applied spatially, by assigning selected joint indices to the ROI mask. In Fig. 6 we illustrate control over the

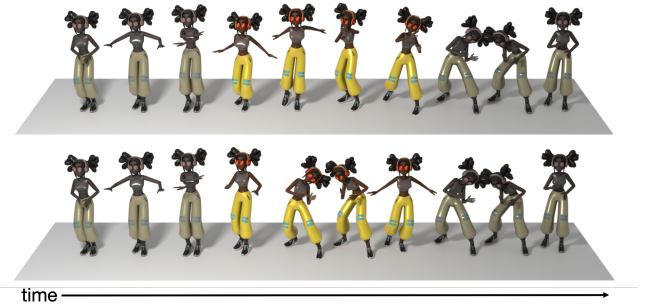


Fig. 4. Temporal composition – In-betweening. Both top and bottom show results for the same input, introducing diverse outputs. The beginning and the end of the motion are given by the reference sequence and can be distinguished according to their faded tone. Observe that the beginning and the end are identical in both sequences. The center of each motion is synthesized.

upper body, where the motion of the upper body is determined by a reference motion and assigned to the target motion. The model synthesizes the rest of the joints yielding a motion with the given sequence in the upper body, and with the learned motifs in the lower body. A composition can be both spatial and temporal, and all it takes is an ROI mask where several frame sequences are zeroed, *i.e.*, taken from the reference motion, and in the complementary part, several joints are zeroed (see Fig. 3).

## 5.2 Motion Harmonization

Given a synthesized motion sequence  $x_0$ , we would like to integrate a portion of an unseen motion,  $y$ , into it. The portion of  $y$  can be either temporal, *i.e.*, several frames, or spatial, *i.e.*, several joints, or both. As visualized in Fig. 7, SinMDM overrides a window in  $x_0$  with the desired portion of  $y$  and denotes the outcome  $y_0$ . Next,  $y_0$  is harmonized such that it matches the core motion elements learned by our model, using a linear low-pass filter  $\phi_N$  as suggested by Choi et al. [2021]. Let  $x'_{t-1}$  and  $y_{t-1}$  denote the noised version of motions  $p_\theta(x_t, t)$  and  $y_0$ , respectively. The high-frequency details of  $x'_{t-1}$  are added to the low-frequency of  $y_{t-1}$  via

$$x_{t-1} = \phi_N(y_{t-1}) + x'_{t-1} - \phi_N(x'_{t-1}). \quad (5)$$

Note the difference between harmonization and motion composition: Both assign parts of an unseen sequence  $y$  into a synthesized

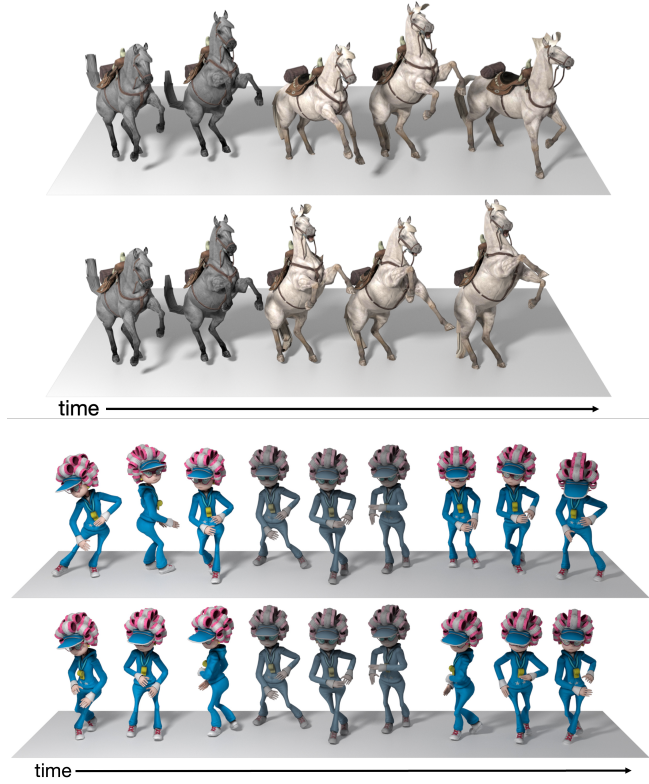


Fig. 5. Temporal composition – motion expansion. Pairs of motions exhibit diverse synthesis from a single input. The motion part provided by the reference sequence is identifiable by its faded color. Note that the parts given as input are identical in both sequences, while the synthesized parts differ. Top: synthesize a suffix given a temporal prefix. Bottom: synthesize a prefix and a suffix, given the middle part.

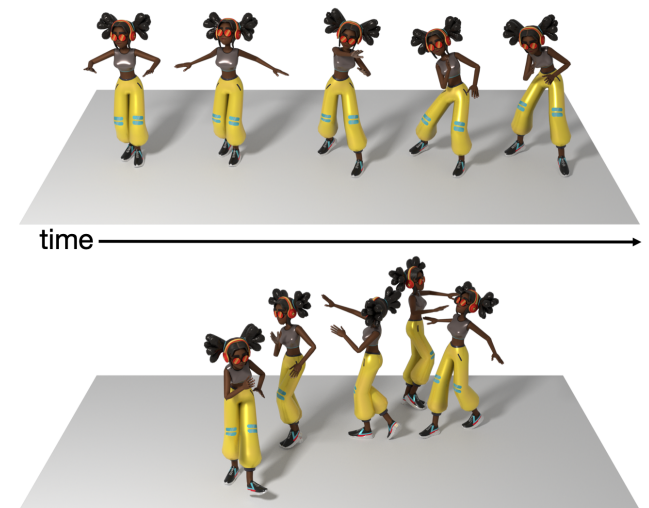


Fig. 6. Spatial composition. Top: reference motion, unseen by the network. Bottom: composed motion. The referenced sequence is *warm-up*, and the learned one is *walk in circle*. In the composed result, the top body part performs a warm-up activity, and the bottom body part walks in a curved path.

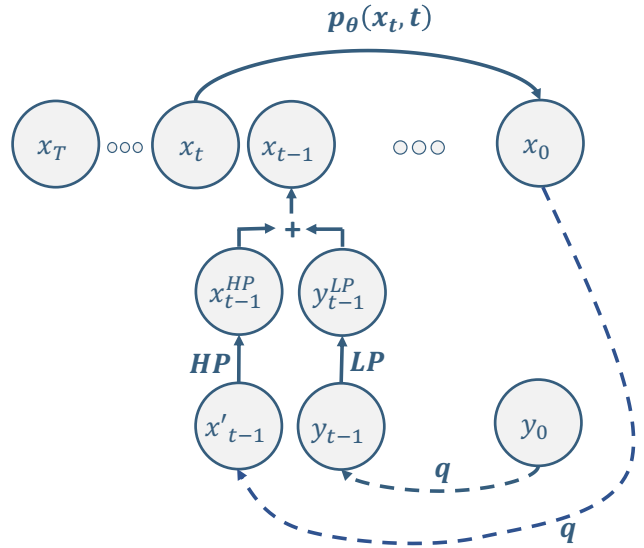


Fig. 7. Motion Harmonization. In order to inject guidance from the input motion  $y_0$  during synthesis, we follow Choi *et al.* [2021] and add its low frequencies  $y_{t-1}^{LP}$  at each denoising step  $t$ .

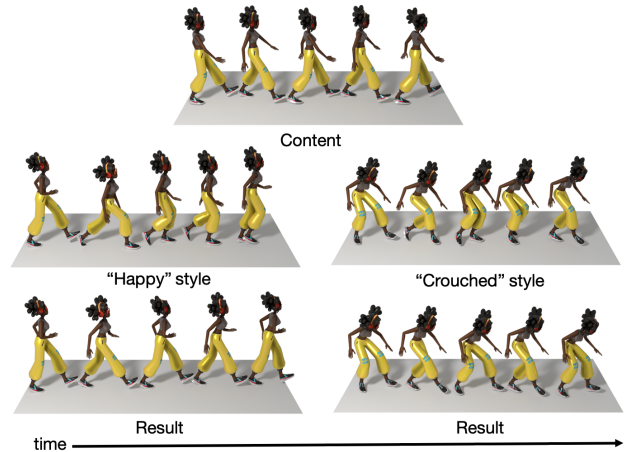


Fig. 8. Style transfer is a special case of the harmonization application, where a reference motion is adjusted such that it matches the learned motion motifs. The style motion is learned by the model, and the content motion is unseen by it. Top: one content, unseen by the network, is applied to both styles. Left: a "happy" style, learned by the network, and below it the harmonized result. Right: a "crouched" style, learned by the network, and below it the harmonized result. Note that the character in both results is using the exact step rhythm and size as the character in the content motion.

motion  $x_0$ . However, harmonization changes the assigned part such that it matches the learned distribution, while composition aims to keep it unchanged.

*Style transfer.* We implement style transfer as a special case of harmonization, where instead of using a portion from  $y$ , we use all of it. That is, we fully override  $x_0$ . We use a style motion  $x$  learned

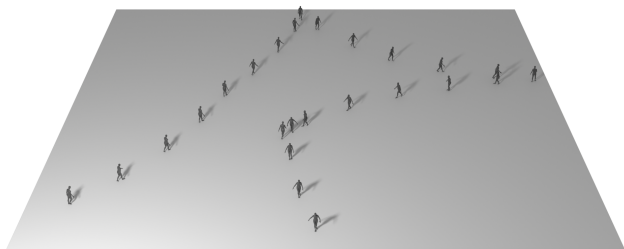


Fig. 9. Long motion. The learned sequence is a 10 seconds motion, depicting a person walking back, then turning and walking back again. The synthesized motion is a 60 seconds sequence, depicting a person walking back, and occasionally turning and walking back again.

by the model, and a content motion  $y$ , unseen by the model. Once applying harmonization, the result possesses the content of  $y$  and the style of  $x$ , as depicted in Fig. 8.

### 5.3 Straight-forward Applications

In this section we present applications that may require special techniques in existing works, but require no special technique when conducted using our model.

*Long motion sequences.* Our model can synthesize variable-length motions, even very long ones, with no additional training. Imputed to its small receptive field, the model can hallucinate a sequence as long as requested. An example of a one-minute animation is introduced in Fig. 9.

*Crowd animation.* Although trained on a single sequence, during inference SinMDM can generate a crowd performing a variety of similar motions, each sampled from a different Gaussian noise  $x_T \sim \mathcal{N}(0, I)$ , as illustrated in Fig. 10.

## 6 EXPERIMENTS

Our experiments are held on motion data from the HumanML3D [2022], Mixamo [2021], and Truebones Zoo [2022] datasets, and on an artist-created animation, using an NVIDIA GeForce RTX 2080 Ti GPU.

### 6.1 Benchmarks

We test our framework on two benchmarks. One consists of data from the HumanML3D dataset, and the other from the Mixamo dataset. These two datasets are different in many aspects. The data in HumanML3D fits the SMPL [Loper et al. 2015] topology, and its users normally use SMPL’s mean body definition. In contrast, Mixamo provides 70 characters, each possessing their unique bone lengths and some possessing unique topologies. In addition, the motions in the Mixamo dataset are more diverse and more dynamic.

### 6.2 Metrics

For each benchmark, we use a different set of metrics. For the Mixamo benchmark, we use the metrics introduced in Ganimator [Li et al. 2022]. Ganimator is our immediate comparison reference, as it

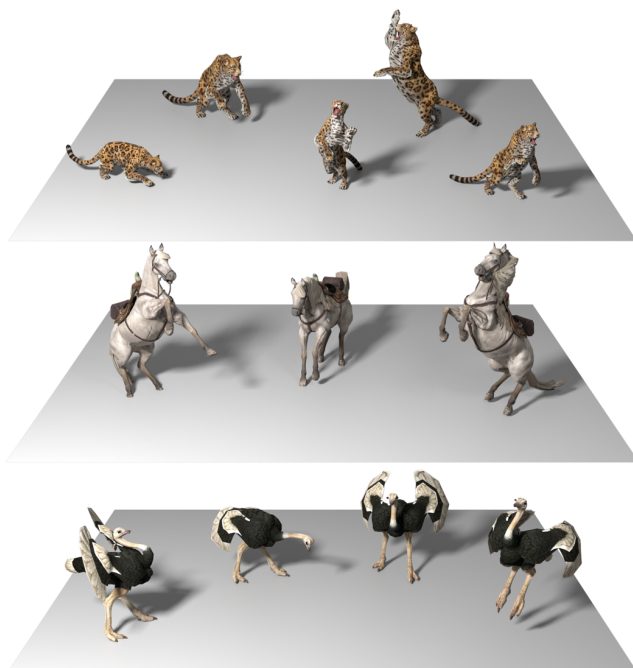


Fig. 10. Crowd animation. Groups of jaguars, horses, and ostriches. In each group, no motion is like the other, and yet they are all learned from a single motion sequence.

is the only single-motion learning work. For fairness, we compare with it using its own metrics. However, these metrics are based on the values of motion features (e.g., rotation angles) while the usage of deep features is the current best practice [Zhang et al. 2018]. Given HumanML3D’s capability for deep feature calculation, we utilize it to present our results specifically on these features.

A *good* score varies depending on the metric, being a *high* value when higher is better and a *low* value when lower is preferred. Note that attaining a good score on some metrics, but a bad score on others, is insufficient: Good diversity scores with bad fidelity indicate deviation from the input motion, while good fidelity scores with bad diversity suggest overfitting.

An ideal outcome is a combination of good values for all metrics. For models with mixed scores, a better-scoring model is the one whose scores are more balanced. To this end, we follow established literature [Chinchor 1992; Rijsbergen 1979] and suggest the Harmonic Mean metric, which is widely used in Machine Learning for this purpose [Taha and Hanbury 2015]. We compute it as follows: first, we normalize the scores for each metric. Normalization is between zero to the metric’s maximum value. If the maximum is not known, we select the 90% percentile of the computed scores. For metrics where lower is better, we subtract the score from the maximum value. Note that a negative value is therefore valid. We compute the Harmonic Mean via

$$HM = E / \left( \frac{1}{s_1} + \dots + \frac{1}{s_E} \right), \quad (6)$$

where  $E$  is the number of metrics in a table and  $s_i$  is the normalized score of metric  $i$ .

Table 1. Results on the Mixamo benchmark, comparing our work with state-of-the-art Ganimator. SinMDM leads in all metrics but one. In particular, it demonstrates a significant advantage in the Harmonic Mean metric.

	Coverage ↑	Global Div. ↑	Local Div. ↑	Inter Div. ↑	Intra Div. Diff. ↓	#Param. (M) ↓	#Iter. (K) ↓	Iter. Time (s) ↓	Tot. Time (h) ↓	Harmon. Mean ↑
Ganimator	94.3	1.24	1.17	0.09	0.13	21.7	60 (15×4)	0.36	6.0	-0.22
SinMDM (Ours)	94.3	<b>1.42</b>	1.00	<b>0.13</b>	<b>0.03</b>	<b>5.26</b>	60	<b>0.09</b>	<b>1.5</b>	<b>0.85</b>

*Metrics on the Mixamo benchmark.* We use the Mixamo dataset to compare SinMDM and Ganimator. For a fair comparison, we use the metrics suggested by them. However, their metrics do not measure the difference between synthesized motions (inter-diversity) nor the difference between sub-motions within one motion (intra-diversity), thus we add metrics to measure these missing qualities. This group of metrics is applied to the motion itself, and not to deep features.

The metrics in Ganimator consist of (a) *coverage*, which is the rate of temporal windows in the input motion  $x_0$  that are reproduced by the synthesized one, (b) *global diversity*, measuring the distance between  $tess(\hat{x}_0)$  and  $x_0$ , where  $tess(\cdot)$  is a tessellation that minimizes the L2 distance to the input sequence, and (c) *local diversity*, which is the average distance between windows in the synthesized motion  $\hat{x}_0$  and their nearest neighbors in the input one.

The aforementioned metrics are measured relative to the input motion sequence. We add two metrics that are not related to the input motion. The first is (d) *inter diversity*, the diversity between synthesized motions. We define *intra diversity* to be the diversity between sub-windows internal to a motion and define (e) *intra diversity diff*, which is the difference between the intra diversity of the synthesized motions and that of the input motion.

In addition, we measure time-space efficiency values: (f) the number of network parameters, (g) the number of required iterations, (h) the time required for each iteration, and (d) the total running time, which is a multiplication of the last two. For metrics (a)-(d), a higher score is better. For metrics (e)-(h), a lower score is better.

*Metrics on the HumanML3D benchmark.* We use this benchmark to measure metrics that are applied on deep features, obtained with a motion encoder by Guo et al. [2022]. The computed metrics are (a) *SiFID* [Shaham et al. 2019], which measures the distance between the distribution of sub-windows in the learned motion and a synthesized one, (b) *inter diversity*, which is the LPIPS distance [Zhang et al. 2018] between various motions synthesized out of one input, and (c) *intra diversity diff*, which is the difference between the intra diversity of the synthesized motions and that of the input motion, where intra diversity is the LPIPS distance between sub-windows in one synthesized motion. For metrics (a) and (c) lower is better, and for metric (b) higher is better.

### 6.3 Quantitative Results

In table 1 we compare SinMDM with Ganimator. The table shows that our work outperforms Ganimator in all metrics except one. Specifically, SinMDM exhibits a notable advantage in the Harmonic Mean metric, which effectively captures the collective strength of all scores.

Table 2. Results on the Gangnam-style motion. We mark the table leaders for the bottom part only, as MotionTexture and acRNN exhibit either overfit or divergence. In the lower part, where both models achieve high scores in all metrics, our model takes the lead, despite the fact that the subject motion has been selectively chosen by the authors of Ganimator.

	Coverage ↑	Global Diversity ↑	Local Diversity ↑	Harmonic Mean ↑
MotionTexture [2002]	84.6	1.03	1.04	0.32
MotionTexture (Single)	100	0.21	0.33	0.09
acRNN [2018]	11.6	5.63	6.69	0.30
Ganimator [2022]	97.2	1.29	<b>1.19</b>	0.38
SinMDM (Ours)	<b>98.1</b>	<b>1.55</b>	1.12	<b>0.39</b>

Table 3. Comparison with motion diffusion model MDM. MDM achieves good SiFID and intra-diversity but exhibits poor inter-diversity, indicating overfitting. MDM on crops attains good inter-diversity but bad scores for the other metrics, indicating deviation from the input. Our model attains good scores in all metrics, demonstrating that a balance of good scores across all metrics is more important than excelling in only a select few.

	SiFID ↓	Inter Diversity ↑	Intra Div. Diff. ↓	Harmonic Mean ↑
MDM [2022b]	0.01	0.03	0.14	0.14
MDM on crops	13.94	1.64	1.83	-1.01
SinMDM (Ours)	1.87	0.73	0.40	<b>0.82</b>

All the metrics are computed separately on each benchmark motion and then averaged. The metrics that measure time were computed on benchmark motion number 9 only.

The authors of Ganimator published a quantitative comparison solely on one selected motion, namely the Gangnam-style dancing sequence. We align with their study and measure our results on this motion as well, as shown in Tab. 2. In this table we compare with two other, non-single motion works, MotionTexture [Li et al. 2002] and acRNN [Zhou et al. 2018]. Note that for this specific motion, selected by the Ganimator authors, our results lead the table in two metrics and are comparable in the third.

To evaluate our performance vs. another motion diffusion model, we compare it with two variations of the MDM [Tevet et al. 2022b] framework. The first is a vanilla MDM, trained on a single-motion. The second is a variation of MDM in which we extract short crops out of the single-motion sequence and train an MDM on them. Note that the second variation holds a narrow receptive field.

The comparison is conducted on the HumanML3D dataset with metrics based on deep features. The results are shown in Tab. 3. As mentioned above, attaining a high score in one metric only, indicates either overfit or divergence from the input motion. MDM yields complete overfit, thus its SiFID and intra-diversity scores are perfect (indicating similarity to the input motion), but its inter-diversity scores are low. The overfit of MDM is caused by the global attention it uses. On the other hand, the quantitative results for the second MDM variation indicate divergence from the input motion motifs. These quantitative results are supported by the qualitative results in our supplementary video.

Finally, we perform a user study in which users are requested to judge which model is better in terms of diversity, fidelity, and quality. In the study, we compare our model vs. Ganimator and MDM trained on crops. Each pair of models is compared over 8 different motions,



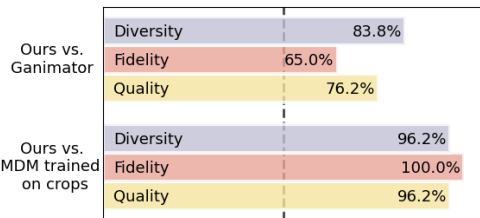


Fig. 11. User study. Users vote that our model performs better than state-of-the-art Ganimator and MDM trained on crops. The dashed line marks 50%.

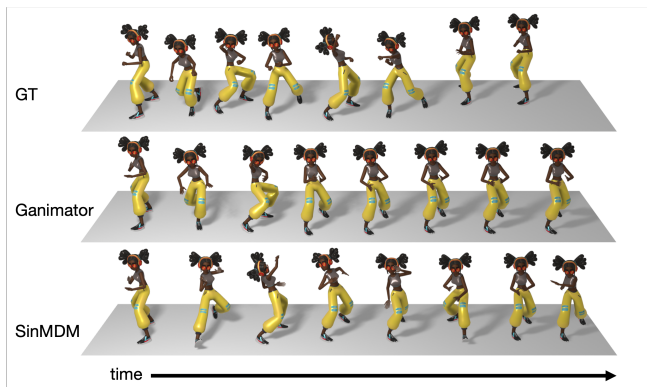


Fig. 12. Qualitative comparison, based on the motion *punch to elbow* from the Mixamo dataset. Observe the mode collapse in Ganimator’s synthesized motion, where over half of the motion is frozen.

Table 4. Ablation results on the HumanML3D benchmarks. Our selected architecture is framed in gray. Rows 1,2: comparing receptive field widths. Rows 3,4: vanilla attention vs. none (vs. QnA in row 2). Row 5: QnA-based transformer (vs. QnA-based UNet in row 2). Row 1 indicates overfit and row 6 indicates divergence. Rows 3,4 present good scores but not as good as ours. Abbr.: *r.f.* → receptive field, *d* → depth, *atn.* → attention.

	SiFID ↓	Inter Diversity ↑	Intra Div. Dist. ↓	Harmonic Mean ↑
UNet w/ QnA				
wide r.f. ( <i>d</i> =3)	0.69	0.20	0.34	0.54
narrow r.f. ( <i>d</i> =1)	1.87	0.73	0.40	<b>0.82</b>
UNet ( <i>d</i> =1)				
w/ vanilla atn.	1.88	0.72	0.45	0.79
w/o atn. w/o QnA	2.03	0.72	0.43	0.79
Transformer w/ QnA	5.99	1.74	0.57	0.56

and each such comparison is judged by 10 distinct users. The results (Fig. 11) show that our model is significantly preferred by the users. Screenshots from our user study are provided in Appendix C.

#### 6.4 Qualitative results

Our supplementary video reflects the quality of our results. It presents multiple synthesized motions, as well as a comparison to other works. In addition, Fig. 12 and 13 depict SinMDM vs. current work. In contrast to other works that exhibit mode collapse, overfitting, or produce jittery motion, SinMDM demonstrates none of these issues.

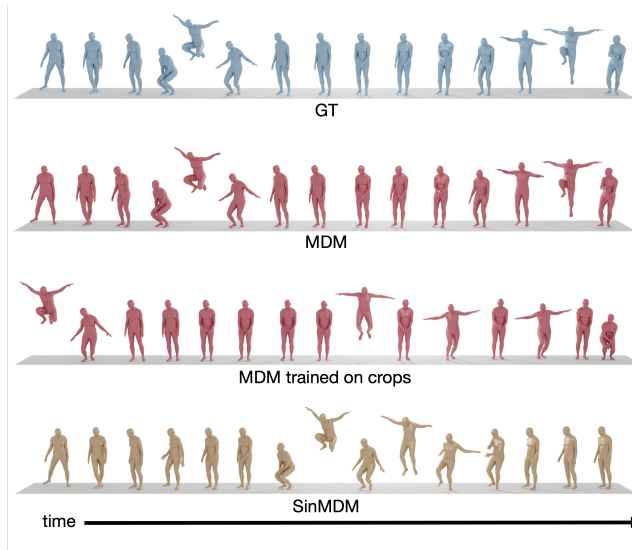


Fig. 13. Qualitative comparison on the HumanML3D dataset. MDM exhibits overfit, and MDM trained on crops exhibits jittery motion, e.g., when transferring from standing to jumping without bending the knees before and after.

#### 6.5 Ablation

We examine several architectural variations on the HumanML3D benchmark and present the results in Tab. 4. We start by confirming that a narrow receptive field produces plausible results while a wider one induces overfit (rows 1,2). In order to do so, we examine a fixed architecture (QnA-based UNet) with two receptive field widths. We control the width by tweaking the depth of the UNet. Indeed we observe that the model with the wide receptive field overfits (replicates) the input motion, as its inter-diversity is bad while its SinFid and intra-diversity are good.

Recall that the UNet architecture used by many diffusion models [Ho et al. 2020; Nichol et al. 2022] holds global attention layers in it. In the next experiment (rows 3,4), we confirm that replacing UNet’s global attention with a local one (QnA) is a good choice. We fix the network’s depth and examine alternatives to QnA. One alternative is the usage of vanilla attention, and the other is to use no attention whatsoever. Both alternatives show plausible metric results, and yet, our QnA-based UNet (row 2) performs better. Plausible results with vanilla attention (row 3) are noteworthy, considering its global receptive field. This can be attributed to the absence of temporal positional embedding, enabling the generative model to identify motion patterns across various temporal regions.

Finally, as many motion diffusion models favor a transformer over a UNet [Kim et al. 2022; Tevet et al. 2022b], we measure the scores for a QnA-based transformer (row 5). To refrain from overfitting, we apply QnA layers within the transformer as we do with the UNet. In addition, to promote diversity and permit the rearrangement of motion components, we employ relative temporal positional embeddings [Press et al. 2021; Shaw et al. 2018; Su et al. 2021] instead of the existing global ones. However, the QnA-based transformer attains a bad SiFid score, indicating poor fidelity to the input motion.

Note that due to the mixed scores (that indicate either overfit or divergence), the usage of the Harmonic Mean metric is essential as it allows for the assessment of the combined strength of all scores.

## 7 CONCLUSIONS

We have explored the use of diffusion models on single motion sequence synthesis and designed a motion denoising transformer with a narrow receptive field. Training on single motions is particularly useful in motion domains, where the number of data instances is scarce. Particularly, for animals and imaginary creatures, which have unique skeletons and motion motifs. The motion of such creatures cannot be captured easily nor learned from the human motion data available.

Our experiments on several datasets demonstrate that our lightweight diffusion-based method significantly outperforms current work both in quality and time-space performance. Moreover, our approach allows the synthesis of particularly long motions, and enables a variety of motion manipulation tasks, including spatial and temporal in-betweening, motion expansion, harmonization, style transfer, and crowd animation.

The innate limitation of our method, common to all models (in all domains) that learn a single instance, is the limited ability to synthesize out-of-distribution. However, the main limitation of our diffusion-based approach is the relatively long inference time. This is due to the iterative nature of diffusion models.

Finally, our work shows the competence of diffusion models to learn from limited data, which contradicts their reputation for requiring large amounts of data. Nevertheless, in the future, we would like to address the single input limitations, by possibly learning from available motion data of creatures with rather compatible skeletons.

## 8 ACKNOWLEDGMENTS

We are grateful to Panayiotis Charalambous, Andreas Aristidou and Brian Gordon for reviewing earlier versions of the manuscript. This research was supported in part by the Israel Science Foundation (grants no. 2492/20 and 3441/21), Len Blavatnik and the Blavatnik family foundation, and the Tel Aviv University Innovation Laboratories (TILabs).

## REFERENCES

- Kfir Aberman, Peizhuo Li, Dani Lischinski, Olga Sorkine-Hornung, Daniel Cohen-Or, and Baoquan Chen. 2020a. Skeleton-aware networks for deep motion retargeting. *ACM Transactions on Graphics (TOG)* 39, 4 (2020), 62–1.
- Kfir Aberman, Yijia Weng, Dani Lischinski, Daniel Cohen-Or, and Baoquan Chen. 2020b. Unpaired motion style transfer from video to animation. *ACM Transactions on Graphics (TOG)* 39, 4 (2020), 64–1.
- Kfir Aberman, Rundui Wu, Dani Lischinski, Baoquan Chen, and Daniel Cohen-Or. 2019. Learning Character-Agnostic Motion for Motion Retargeting in 2D. *ACM Transactions on Graphics (TOG)* 38, 4 (2019), 75.
- Adobe Systems Inc. 2021. Mixamo. <https://www.mixamo.com> Accessed: 2021-12-25.
- Chaitanya Ahuja and Louis-Philippe Morency. 2019. Language2pose: Natural language grounded pose forecasting. In *2019 International Conference on 3D Vision (3DV)*. IEEE, IEEE Computer Society, Washington, DC, USA, 719–728.
- Emre Aksan, Manuel Kaufmann, and Otmar Hilliges. 2019. Structured prediction helps 3d human motion modelling. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. IEEE Computer Society, Washington, DC, USA, 7144–7153.
- Moab Arar, Ariel Shamir, and Amit H Bermano. 2022. Learned Queries for Efficient Local Attention. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. IEEE Computer Society, Washington, DC, USA, 10841–10852.
- A Aristidou, A Yiannakidis, K Aberman, D Cohen-Or, A Shamir, and Y Chrysanthou. 2022. Rhythm is a Dancer: Music-Driven Motion Synthesis with Global Structure. *IEEE Transactions on Visualization and Computer Graphics* 1 (2022).
- Yuki Markus Asano, Christian Rupprecht, and Andrea Vedaldi. 2020. A critical analysis of self-supervision, or what we can learn from a single image. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, OpenReview.net. <https://openreview.net/forum?id=B1esx6EYvr>
- Omri Avrahami, Dani Lischinski, and Ohad Fried. 2022. Blended diffusion for text-driven editing of natural images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. IEEE Computer Society, Washington, DC, USA, 18208–18218.
- Emad Barsoum, John Kender, and Zicheng Liu. 2018. Hp-gan: Probabilistic 3d human motion prediction via gan. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*. IEEE Computer Society, Washington, DC, USA, 1418–1427.
- Uttaran Bhattacharya, Nicholas Rewkowski, Abhishek Banerjee, Pooja Guhan, Aniket Bera, and Dinesh Manocha. 2021. Text2gestures: A transformer-based network for generating emotive body gestures for virtual agents. In *2021 IEEE Virtual Reality and 3D User Interfaces (VR)*. IEEE, IEEE Computer Society, Washington, DC, USA, 1–10.
- Pablo Cervantes, Yusuke Sekikawa, Ikuro Sato, and Koichi Shinoda. 2022. Implicit Neural Representations for Variable Length Human Motion Generation.
- Jinshu Chen, Qihui Xu, Qi Kang, and MengChu Zhou. 2021. Mogan: Morphologic-structure-aware generative learning from a single image.
- N Chinchor. 1992. MUC-4 evaluation metrics in Proc. of the Fourth Message Understanding Conference 22–29.
- Jooyoung Choi, Sungwon Kim, Yonghyun Jeong, Youngjune Gwon, and Sungroh Yoon. 2021. Ilvr: Conditioning method for denoising diffusion probabilistic models. In *2021 IEEE. In CVF International Conference on Computer Vision (ICCV)*. IEEE Computer Society, Washington, DC, USA, 14347–14356.
- Bruno Degardin, João Neves, Vasco Lopes, João Brito, Ehsan Yaghoubi, and Hugo Proença. 2022. Generative Adversarial Graph Convolutional Networks for Human Action Synthesis. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. IEEE Computer Society, Los Alamitos, CA, USA, 1150–1159.
- Prafulla Dhariwal and Alexander Nichol. 2021. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems* 34 (2021), 8780–8794.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2021. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, OpenReview.net. <https://openreview.net/forum?id=YicbFdNTTy>
- Yinglin Duan, Tianyang Shi, Zhengxia Zou, Yanan Lin, Zhehui Qian, Bohan Zhang, and Yi Yuan. 2021. Single-shot motion completion with transformer.
- Katerina Fragkiadaki, Sergey Levine, Panna Felsen, and Jitendra Malik. 2015. Recurrent network models for human dynamics. In *Proceedings of the IEEE international conference on computer vision*. IEEE Computer Society, Washington, DC, USA, 4346–4354.
- Saeed Ghorbani, Calden Wloka, Ali Etemad, Marcus A. Brubaker, and Nikolaus F. Troje. 2020. Probabilistic Character Motion Synthesis using a Hierarchical Deep Latent Variable Model. *Computer Graphics Forum* 1 (2020). <https://doi.org/10.1111/cgf.14116>
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial nets. *Advances in neural information processing systems* 27 (2014).
- Brian Gordon, Sigal Raab, Guy Azov, Raja Giryes, and Daniel Cohen-Or. 2022. FLEX: Extrinsic Parameters-free Multi-view 3D Human Motion Reconstruction. In *European Conference on Computer Vision*. Springer, Springer International Publishing, Berlin/Heidelberg, Germany, 176–196.
- Niv Granot, Ben Feinstein, Assaf Shocher, Shai Bagon, and Michal Irani. 2022. Drop the gan: In defense of patches nearest neighbors as single image generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. IEEE Computer Society, Washington, DC, USA, 13460–13469.
- Chuan Guo, Shihao Zou, Xinxin Zuo, Sen Wang, Wei Ji, Xingyu Li, and Li Cheng. 2022. Generating Diverse and Natural 3D Human Motions From Text. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. IEEE Computer Society, Washington, DC, USA, 5152–5161.
- Chuan Guo, Xinxin Zuo, Sen Wang, Shihao Zou, Qingyao Sun, Annan Deng, Minglun Gong, and Li Cheng. 2020. Action2motion: Conditioned generation of 3d human motions. In *Proceedings of the 28th ACM International Conference on Multimedia*. ACM, New York, NY, USA, 2021–2029.
- Ikhsanul Habibie, Daniel Holden, Jonathan Schwarz, Joe Yearsley, and Taku Komura. 2017. A recurrent variational autoencoder for human motion synthesis. In *2017 British Machine Vision Conference*. BMVC, UK.
- Félix G Harvey and Christopher Pal. 2018. Recurrent transition networks for character locomotion. In *SIGGRAPH Asia 2018 Technical Briefs*. ACM, New York, NY, USA, 1–4.

- Félic G Harvey, Mike Yurick, Derek Nowrouzezahrai, and Christopher Pal. 2020. Robust motion in-betweening. *ACM Transactions on Graphics (TOG)* 39, 4 (2020), 60–1.
- Chengnan He, Jun Saito, James Zachary, Holly Rushmeier, and Yi Zhou. 2022. NeMF: Neural Motion Fields for Kinematic Animation. In *NeurIPS*.
- Gustav Eje Henter, Simon Alexanderson, and Jonas Beskow. 2020. Moglow: Probabilistic and controllable motion synthesis using normalising flows. *ACM Transactions on Graphics (TOG)* 39, 6 (2020), 1–14.
- Alejandro Hernandez, Jurgen Gall, and Francesc Moreno-Noguer. 2019. Human motion prediction via spatio-temporal inpainting. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. IEEE Computer Society, Washington, DC, USA, 7134–7143.
- Tobias Hinz, Matthew Fisher, Oliver Wang, and Stefan Wermter. 2021. Improved techniques for training single-image gans. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. IEEE Computer Society, Washington, DC, USA, 1300–1309.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems* 33 (2020), 6840–6851.
- Jonathan Ho and Tim Salimans. 2022. Classifier-free diffusion guidance.
- Daniel Holden, Ikhsanul Habibie, Ikuo Kusajima, and Taku Komura. 2017. Fast neural style transfer for motion data. *IEEE computer graphics and applications* 37, 4 (2017), 42–49.
- Daniel Holden, Jun Saito, and Taku Komura. 2016. A deep learning framework for character motion synthesis and editing. *ACM Transactions on Graphics (TOG)* 35, 4 (2016), 1–11.
- Daniel Holden, Jun Saito, Taku Komura, and Thomas Joyce. 2015. Learning motion manifolds with convolutional autoencoders. In *SIGGRAPH Asia 2015 technical briefs*. ACM, New York, NY, USA, 1–4.
- Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. 2017. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. IEEE Computer Society, Washington, DC, USA, 1125–1134.
- Deok-Kyeong Jang and Sung-Hee Lee. 2020. Constructing human motion manifold with sequential networks. In *Computer Graphics Forum*, Vol. 39. Wiley Online Library, The Eurographics Association and John Wiley and Sons Ltd., Hoboken, NJ, USA, 314–324.
- Manuel Kaufmann, Emre Aksan, Jie Song, Fabrizio Pece, Remo Ziegler, and Otmar Hilliges. 2020. Convolutional autoencoders for human motion infilling. In *2020 International Conference on 3D Vision (3DV)*. IEEE, IEEE Computer Society, Washington, DC, USA, 918–927.
- Jihoon Kim, Jiseob Kim, and Sungjoon Choi. 2022. FLAME: Free-form Language-based Motion Synthesis & Editing.
- Vladimir Kulikov, Shahar Yadin, Matan Kleiner, and Tomer Michaeli. 2022. SinDDM: A Single Image Denoising Diffusion Model.
- Juheon Lee, Seohyun Kim, and Kyogu Lee. 2018. Listen to dance: Music-driven choreography generation using autoregressive encoder-decoder network.
- Chuan Li and Michael Wand. 2016. Precomputed real-time texture synthesis with markovian generative adversarial networks. In *European conference on computer vision*. Springer, Springer International Publishing, Berlin/Heidelberg, Germany, 702–716.
- Peizhuo Li, Kfir Aberman, Zihan Zhang, Rana Hanocka, and Olga Sorkine-Hornung. 2022. GANimator: Neural Motion Synthesis from a Single Sequence. *ACM Transactions on Graphics (TOG)* 41, 4 (2022), 138.
- Ruilong Li, Shan Yang, David A Ross, and Angjoo Kanazawa. 2021. Learn to dance with aist++: Music conditioned 3d dance generation. , arXiv–2101 pages.
- Yan Li, Tianshu Wang, and Heung-Yeung Shum. 2002. Motion texture: a two-level statistical model for character motion synthesis. In *Proceedings of the 29th annual conference on Computer graphics and interactive techniques*. ACM, New York, NY, USA, 465–472.
- Chieh Hubert Lin, Hsin-Ying Lee, Yen-Chi Cheng, Sergey Tulyakov, and Ming-Hsuan Yang. 2021. InfinityGAN: Towards Infinite-Pixel Image Synthesis. In *International Conference on Learning Representations*. OpenReview.net, OpenReview.net.
- Jianxin Lin, Yingxue Pang, Yingce Xia, Zhibo Chen, and Jiebo Luo. 2020. TuiGAN: Learning versatile image-to-image translation with two unpaired images. In *European Conference on Computer Vision*. Springer, Springer International Publishing, Berlin/Heidelberg, Germany, 18–35.
- Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. 2015. SMPL: A skinned multi-person linear model. *ACM transactions on graphics (TOG)* 34, 6 (2015), 1–16.
- Shubh Maheshwari, Debtanu Gupta, and Ravi Kiran Sarvadevabhatla. 2022. MUGL: Large Scale Multi Person Conditional Action Generation with Locomotion. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. IEEE Computer Society, Los Alamitos, CA, USA, 257–265.
- Alexander Quinn Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. 2022. GLIDE: Towards Photorealistic Image Generation and Editing with Text-Guided Diffusion Models. In *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA (Proceedings of Machine Learning Research, Vol. 162)*, Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato (Eds.). PMLR, PMLR, 16784–16804. <https://proceedings.mlr.press/v162/nichol22a.html>
- Yaniv Nikankin, Niv Haim, and Michal Irani. 2022. SinFusion: Training Diffusion Models on a Single Image or Video.
- Niki Parmar, Prajit Ramachandran, Ashish Vaswani, Irwan Bello, Anselm Levskaya, and Jon Shlens. 2019. Stand-Alone Self-Attention in Vision Models. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett (Eds.). 68–80. <https://proceedings.neurips.cc/paper/2019/hash/3416a75f4cea9109507cac8e2faefc-Abstract.html>
- Mathis Petrovich, Michael J. Black, and Gül Varol. 2021. Action-Conditioned 3D Human Motion Synthesis with Transformer VAE. In *International Conference on Computer Vision (ICCV)*. IEEE Computer Society, Washington, DC, USA, 10985–10995.
- Mathis Petrovich, Michael J. Black, and Gül Varol. 2022. TEMOS: Generating diverse human motions from textual descriptions. In *European Conference on Computer Vision (ECCV)*. Springer International Publishing, Berlin/Heidelberg, Germany.
- Ofir Press, Noah A Smith, and Mike Lewis. 2021. Train short, test long: Attention with linear biases enables input length extrapolation.
- Jia Qin, Youyi Zheng, and Kun Zhou. 2022. Motion In-Betweening via Two-Stage Transformers. *ACM Transactions on Graphics (TOG)* 41, 6 (2022), 1–16.
- Sigal Raab, Inbal Leibovitch, Peizhuo Li, Kfir Aberman, Olga Sorkine-Hornung, and Daniel Cohen-Or. 2022. MoDi: Unconditional Motion Synthesis from Diverse Data.
- Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. 2022. Hierarchical text-conditional image generation with clip latents.
- CJ Rijsbergen. 1979. Information retrieval. online book <http://www.dcs.gla.ac.uk/Keith.Preface.html> (1979).
- Olaf Ronneberger, Philipp Fischer, and Thomas Brox. 2015. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*. Springer, Springer International Publishing, Berlin/Heidelberg, Germany, 234–241.
- Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi. 2022. Palette: Image-to-image diffusion models. In *ACM SIGGRAPH 2022 Conference Proceedings*. ACM, New York, NY, USA, 1–10.
- Yonatan Shafir, Guy Tevet, Roy Kapon, and Amit H Bermanto. 2023. Human Motion Diffusion as a Generative Prior.
- Tamar Rott Shaham, Tali Dekel, and Tomer Michaeli. 2019. SinGAN: Learning a generative model from a single natural image. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. IEEE Computer Society, Washington, DC, USA, 4570–4580.
- Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. 2018. Self-Attention with Relative Position Representations. In *Proceedings of NAACL-HLT*. 464–468.
- Assaf Shocher, Shai Bagon, Phillip Isola, and Michal Irani. 2019. Ingan: Capturing and retargeting the “dna” of a natural image. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. IEEE Computer Society, Washington, DC, USA, 4492–4501.
- Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. 2015. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*. PMLR, PMLR, 2256–2265.
- Minjung Son, Jeong Joon Park, Leonidas Guibas, and Gordon Wetzstein. 2022. SinGRAF: Learning a 3D Generative Radiance Field for a Single Scene.
- Jiaming Song, Chenlin Meng, and Stefano Ermon. 2020a. Denoising Diffusion Implicit Models. In *International Conference on Learning Representations*. OpenReview.net, OpenReview.net.
- Yang Song and Stefano Ermon. 2020. Improved techniques for training score-based generative models. *Advances in neural information processing systems* 33 (2020), 12438–12448.
- Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. 2020b. Score-Based Generative Modeling through Stochastic Differential Equations. In *International Conference on Learning Representations*. OpenReview.net, OpenReview.net.
- Sebastian Starke, Ian Mason, and Taku Komura. 2022. Deepphase: Periodic autoencoders for learning motion phase manifolds. *ACM Transactions on Graphics (TOG)* 41, 4 (2022), 1–13.
- Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. 2021. Reformer: Enhanced transformer with rotary position embedding.
- Guofei Sun, Yongkang Wong, Zhiyong Cheng, Mohan S Kankanhalli, Weidong Geng, and Xiangdong Li. 2020. DeepDance: music-to-dance motion choreography with adversarial learning. *IEEE Transactions on Multimedia* 23 (2020), 497–509.
- Wenyue Sun and Bao-Di Liu. 2020. ESinGAN: Enhanced single-image GAN using pixel attention mechanism for image super-resolution. In *2020 15th IEEE International Conference on Signal Processing (ICSP)*, Vol. 1. IEEE, IEEE Computer Society, Washington, DC, USA, 181–186.

- Vadim Sushko, Dan Zhang, Juergen Gall, and Anna Khoreva. 2021. Generating Novel Scene Compositions from Single Images and Videos.
- Abdel Aziz Taha and Allan Hanbury. 2015. Metrics for evaluating 3D medical image segmentation: analysis, selection, and tool. *BMC medical imaging* 15, 1 (2015), 1–28.
- Piotr Teterwak, Aaron Sarna, Dilip Krishnan, Aaron Maschinot, David Belanger, Ce Liu, and William T Freeman. 2019. Boundless: Generative adversarial networks for image extension. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. IEEE Computer Society, Washington, DC, USA, 10521–10530.
- Guy Tevet, Brian Gordon, Amir Hertz, Amit H Bermano, and Daniel Cohen-Or. 2022a. MotionCLIP: Exposing Human Motion Generation to CLIP Space.
- Guy Tevet, Sigal Raab, Brian Gordon, Yonatan Shafir, Daniel Cohen-Or, and Amit H Bermano. 2022b. Human motion diffusion model.
- Truebones Motions Animation Studios. 2022. Truebones. <https://truebones.gumroad.com/> Accessed: 2022-1-15.
- Jonathan Tseng, Rodrigo Castellon, and C Karen Liu. 2022. EDGE: Editable Dance Generation From Music.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).
- Ruben Villegas, Jimei Yang, Duygu Ceylan, and Honglak Lee. 2018. Neural kinematic networks for unsupervised motion retargetting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. IEEE Computer Society, Washington, DC, USA, 8639–8648.
- Weilun Wang, Jianmin Bao, Wengang Zhou, Dongdong Chen, Dong Chen, Lu Yuan, and Houqiang Li. 2022. SinDiffusion: Learning a Diffusion Model from a Single Natural Image.
- Zhenyi Wang, Ping Yu, Yang Zhao, Ruiyi Zhang, Yufan Zhou, Junsong Yuan, and Changyou Chen. 2020. Learning diverse stochastic human-action generators by learning smooth latent transitions. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 34. AAAI Press, Washington, DC, USA, 12281–12288.
- Sijie Yan, Zhizhong Li, Yuanjun Xiong, Huahan Yan, and Dahua Lin. 2019. Convolutional sequence generation for skeleton-based action synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. IEEE Computer Society, Washington, DC, USA, 4394–4402.
- Jihyeong Yoo and Qifeng Chen. 2021. SinIR: Efficient General Image Manipulation with Single Image Reconstruction. In *International Conference on Machine Learning*. PMLR, PMLR, 12040–12050.
- Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S Huang. 2019. Free-form image inpainting with gated convolution. In *Proceedings of the IEEE/CVF international conference on computer vision*. IEEE Computer Society, Washington, DC, USA, 4471–4480.
- Ping Yu, Yang Zhao, Chunyuan Li, Junsong Yuan, and Changyou Chen. 2020. Structure-aware human-action generation. In *European Conference on Computer Vision*. Springer, Springer International Publishing, Berlin/Heidelberg, Germany, 18–34.
- Ye Yuan and Kris Kitani. 2020. Dlow: Diversifying latent flows for diverse human motion prediction. In *European Conference on Computer Vision*. Springer, Springer International Publishing, Berlin/Heidelberg, Germany, 346–364.
- Ye Yuan, Jiaming Song, Umar Iqbal, Arash Vahdat, and Jan Kautz. 2022. PhysDiff: Physics-Guided Human Motion Diffusion Model.
- Jia-Qi Zhang, Xiang Xu, Zhi-Meng Shen, Ze-Huan Huang, Yang Zhao, Yan-Pei Cao, Pengfei Wan, and Miao Wang. 2021c. Write-An-Animation: High-level Text-based Animation Editing with Character-Scene Interaction. In *Computer Graphics Forum*, Vol. 40. Wiley Online Library, The Eurographics Association and John Wiley and Sons Ltd., Hoboken, NJ, USA, 217–228.
- Mingyuan Zhang, Zhongang Cai, Liang Pan, Fangzhou Hong, Xinying Guo, Lei Yang, and Ziwei Liu. 2022a. MotionDiffuse: Text-Driven Human Motion Generation with Diffusion Model.
- Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. 2018. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. IEEE Computer Society, Washington, DC, USA, 586–595.
- Yan Zhang, Michael J Black, and Siyu Tang. 2021a. We are more than our joints: Predicting how 3d bodies move. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. IEEE Computer Society, Washington, DC, USA, 3372–3382.
- ZiCheng Zhang, CongYing Han, and TianDe Guo. 2021b. ExSinGAN: Learning an Explainable Generative Model from a Single Image.
- Zicheng Zhang, Yinglu Liu, Congying Han, Hailin Shi, Tiande Guo, and Bowen Zhou. 2022b. PetsGAN: Rethinking Priors for Single Image Generation. In *Thirty-Sixth AAAI Conference on Artificial Intelligence*, AAAI. AAAI Press, Washington, DC, USA, 3408–3416. <https://ojs.aaai.org/index.php/AAAI/article/view/20251>
- Zilong Zheng, Jianwen Xie, and Ping Li. 2021. Patchwise generative convnet: Training energy-based models from a single natural image for internal learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. IEEE Computer Society, Washington, DC, USA, 2961–2970.
- Yi Zhou, Connelly Barnes, Jingwan Lu, Jimei Yang, and Hao Li. 2019. On the continuity of rotation representations in neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. IEEE Computer Society, Washington, DC, USA, 5745–5753.
- Yi Zhou, Zimo Li, Shuangjiu Xiao, Chong He, Zeng Huang, and Hao Li. 2018. Auto-Conditioned Recurrent Networks for Extended Complex Human Motion Synthesis. In *International Conference on Learning Representations*. OpenReview.net, OpenReview.net.

## Appendix

### A HYPERPARAMETERS AND TRAINING DETAILS

In Tab. 5 we detail the values of the hyperparameters that have been used to produce the results shown in this work. Our models have been trained on an NVIDIA GeForce RTX 2080 Ti GPU.

### B QNA RECAP

QnA layers [Arar *et al.* 2022] are a fundamental component in our suggested architecture. In this section, we provide an overview of its underlying implementation and illustrate it in Fig. 14. In particular, QnA is an efficient attention-based layer, which operates in a shift-invariant manner. For every  $k$ -size window, the output is calculated using the self-attention mechanism which is commonly used in the transformer architecture [Vaswani *et al.* 2017]. The self-attention is calculated by first projecting the input features into keys  $K = XW_K$ , values  $V = XW_V$ , and queries  $Q = XW_Q$  via three linear projection matrices  $W_K, W_V, W_Q \in \mathbb{R}^{D \times D}$ . Then, the output of the self-attention operation is defined by:

$$\begin{aligned} \mathbf{SA}(X) &= \mathbf{Attention}(Q, K) \cdot V \\ &= \mathbf{Softmax}\left(\frac{QK^T}{\sqrt{D}}\right) \cdot V. \end{aligned} \quad (7)$$

Instead of performing the pricey query-key operation, QnA detours from extracting the queries from the window itself and directly learns them for the whole-training data (see Fig. 14c). Learning the queries preserves the expressive capability of the self-attention mechanism and enables an efficient implementation that relies on simple and fast operations. In particular, a single query  $\tilde{q}$  is learned, and the attention is applied locally for every  $k$ -size window. Therefore, the output at entry  $z_i$  becomes:

$$z_i = \mathbf{Attention}(\tilde{q}, K_{\mathcal{N}_i}) \cdot V_{\mathcal{N}_i}, \quad (8)$$

where  $\mathcal{N}_i$  is the set  $k$ -neighbourhood of frame  $i$ .

QnA exhibits state-of-the-art accuracy-efficiency trade-off, as depicted in Fig. 15.

### C USER STUDY – SCREENSHOTS

Our user study displays several video clips on each screen, requesting the user to select the one that is more suitable to the examined attribute, which is either quality, fidelity, or diversity. Screenshots from a representative video for each attribute are shown in Fig. 16.

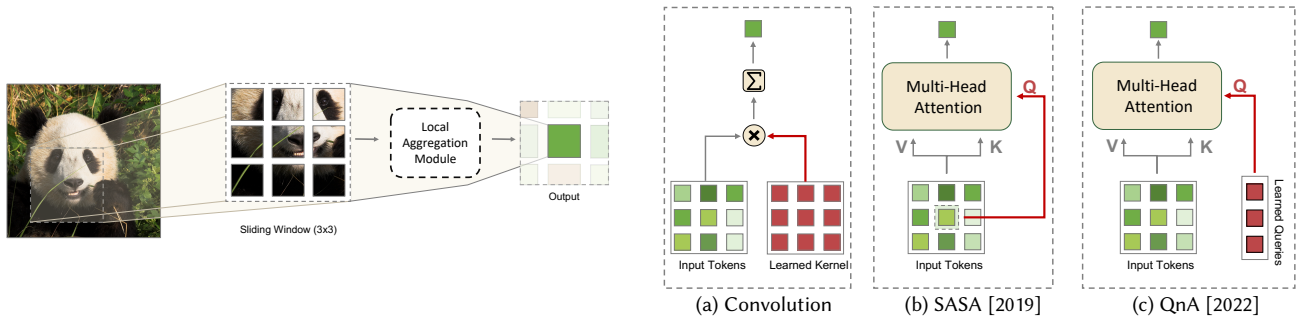


Fig. 14. QnA overview (extracted from the QnA paper). Left: Local layers may utilize various approaches to overlapping windows. (a) Convolutions apply aggregation by learning shared weighted filters. (b) SASA [2019] combines window tokens through self-attention. (c) QnA use shared learned queries across windows, maintaining the expressive power of attention while achieving linear space complexity.

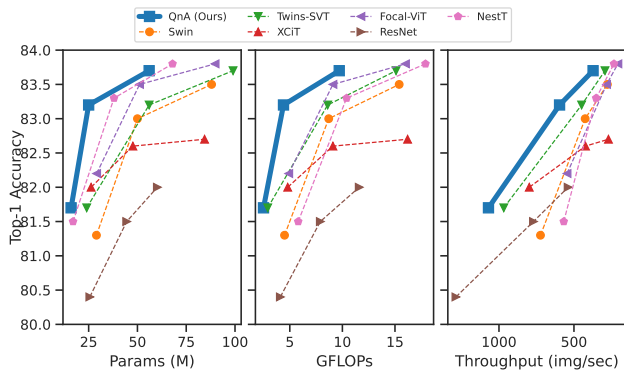
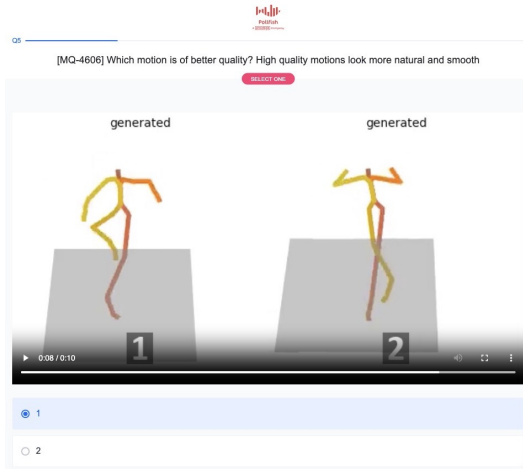
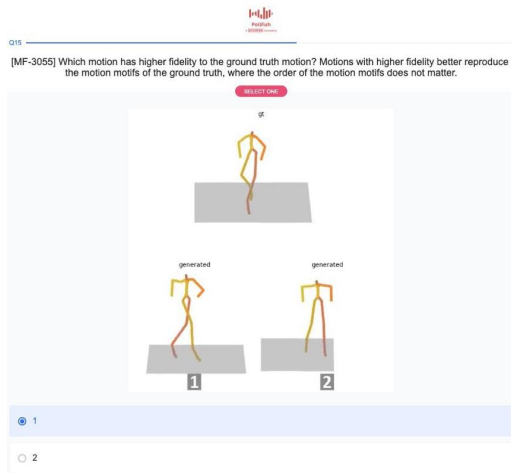


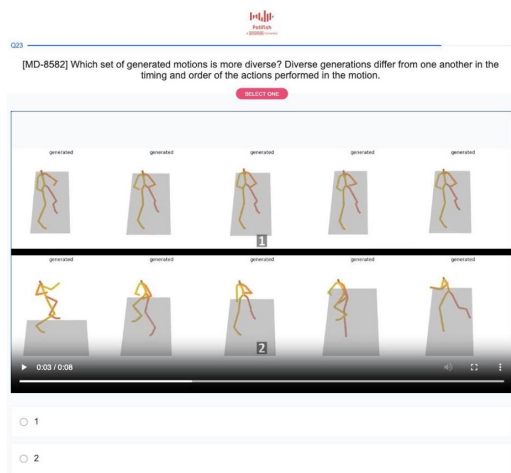
Fig. 15. QnA demonstrates better accuracy-efficiency trade-off compared to state-of-the-art baselines (extracted from the QnA paper).



(a) Quality.



(b) Fidelity.



(c) Diversity.

Fig. 16. Screenshots from our user study. Note that each human figure in the screenshot is played as a video.

Table 5. Our choice of hyperparameters, given with the same names as used in the code.

Name	Value
<u>UNet related</u>	
num_channels	256
channel_mult	1
num_res_blocks	1
kernel_size	3
use_scale_shift_norm	True
use_checkpoint	True
use_attention	True
use_qna	True
<u>QnA related</u>	
head_dim	32
num_heads	4
<u>Diffusion related</u>	
diffusion_steps	1000
noise_schedule	cosine
<u>Training related</u>	
batch_size	64
dropout	0.5
lr_method	ExponentialLR
lr_gamma	0.99998
num_steps	60000
padding_mode	zeros
warmup_steps	0
weight_decay	0