Implicit Neural Representations for Variable Length Human Motion Generation

Pablo Cervantes¹[0000-0002-5256-9317]</sup>, Yusuke Sekikawa²[0000-0003-1111-5949]</sup>, Ikuro Sato^{1,2}[0000-0001-5234-3177]</sup>, and Koichi Shinoda¹[0000-0003-1095-3203]</sup>

¹ Tokyo Institute of Technology ² Denso IT Laboratory Inc.

Abstract. We propose an action-conditional human motion generation method using variational implicit neural representations (INR). The variational formalism enables action-conditional distributions of INRs, from which one can easily sample representations to generate novel human motion sequences. Our method offers variable-length sequence generation by construction because a part of INR is optimized for a whole sequence of arbitrary length with temporal embeddings. In contrast, previous works reported difficulties with modeling variable-length sequences. We confirm that our method with a Transformer decoder outperforms all relevant methods on HumanAct12, NTU-RGBD, and UESTC datasets in terms of realism and diversity of generated motions. Surprisingly, even our method with an MLP decoder consistently outperforms the stateof-the-art Transformer-based auto-encoder. In particular, we show that variable-length motions generated by our method are better than fixedlength motions generated by the state-of-the-art method in terms of realism and diversity. Code at https://github.com/PACerv/ImplicitMotion.

Keywords: Motion Generation, Implicit Neural Representations

1 Introduction

Generative models of human motion serve as a basis for human motion prediction [2,3,12,6,15,5], human animation [37,36], and data augmentation for downstream recognition tasks [23,38,9,34]. There has been intensive research on generative models for realistic and diverse human motions [39,13,19] and in particular methods that can generate motions while controlling some semantic factors such as emotion [13], rhythm [19] or action class [31,10]. For tasks such as rare action recognition, data-efficient action-conditional motion generation has great potential, since it may provide data augmentation even for rare actions.

For motion generation, the quality of generations is evaluated by their realism and diversity. Models need the ability to sample novel and rich representations to generate high-quality motions. A suitable generative model yields distributions of representations in a latent space, where a simple distance measure corresponds to semantic similarity between motions so that interpolations provide novel and high-quality motions. A common generative modeling approach is Variational

Auto-Encoders (VAE) [17,11,10,31], which employ an encoder to infer a distribution from which representations of motions can be sampled and a decoder which reconstructs the data from the representation. The reconstruction loss provides strong supervision, while the variational approach results in a representation space, which allows sampling of novel data with high realism and diversity.

Since human motions naturally vary in length depending on persons or action, it is important to consider variable lengths in motion generation. For example, we would like the representations of quick (short) and slow (long) sitting motions to be different but closer to each other than the representation of a walking motion. In RNN-based VAEs [10], representations are updated each time-step; thus, it is not obvious how to sample a particular action such as quick sitting. Also, their recursive generation may accumulate error when generating long sequences. In contrast, ACTOR, a Transformer VAE [31], should conceptually provide timeindependent representations and generate variable-length motions without accumulating error. Nevertheless, [31] reports directly training with variable-length motions results in almost static motions, and accordingly ACTOR requires an additional fine-tuning scheme to enable variable-length motion generation. It remains unclear what causes such issues with the Transformer architecture.

A recently proposed generative modeling approach is Implicit Neural Representations (INR), which have been shown to be highly efficient in modeling complex data such as 3D scenes [24,26,27]. INRs are representations that encode information without an explicit encoder, but through an optimization procedure as shown in Fig. 1. INRs are usually constructed with respect to a decoder that takes a target coordinate and the representation of a target sample as input and returns the signal of the target sample at the target coordinate. Such representations are optimized individually for each sample with respect to the reconstruction loss at all coordinates. For a time-series, an INR is a time-independent, optimal representation that represents one whole sequence, regardless of the sequence length. Since human motions are naturally variable-length, INRs are a very promising modeling approach. However, to the best of our knowledge, there is no INR-based motion generation method that serves as a strong baseline.

To construct distributions from which one can sample a representation to generate a novel and high-quality motion, we propose variational INRs. Compared to VAEs which infer variational distributions with an encoder, our variational INR framework models each sequence by a distribution with optimized parameters, e.g., mean and covariance, in the representation space. We further decompose INRs into an action-wise part and sequence-wise part. The action-wise part, whose distribution parameters are optimized for all sequences within the same action class, provides generalized components of an action. The sequence-wise part, whose distribution parameters are optimized for an individual sequence, adds fine details of a specific sequence on top of the generalized components.

The average of the sequence distributions within an action class in the representation space serves as the action-specific generative model together with an appropriately trained decoder. In our method, we further split the averaged distribution into several distributions depending on different intervals of sequence



Fig. 1. Comparison between a Variational Auto-Encoder (VAE) baseline (top) and our variational implicit neural representation approach (bottom). In VAEs the encoder weights are optimized with respect to a full dataset and no guarantee of optimal representations for each individual sequence. In contrast, our sequence representations are directly optimized for each individual sequence and by construction offer variablelength sequence generation because a part of each INR is optimized for a whole sequence of arbitrary length. In this figure, we drop the temporal embeddings for simplicity.

lengths. This allows sampling of novel sequences with a target action and a target length. To parameterize the action and sequence-length conditional distribution, we employ Gaussian Mixture Models (GMM). Note that existing high-performing methods are unable to control sequence length. This often results in poor motion generation with sequences ending before the action completes.

However, when fitting a GMM with a high degree of freedom to the representation space we risk simply reproducing training samples. Previous evaluation metrics such as the Fréchet Inception Distance (FID) and Diversity are not sensitive to this problematic model behavior, because they assign a high value to generated motions with a similar distribution as the training set. In this regard, we propose a novel metric, the Mean Maximum Similarity (MMS), to measure such reproducing behavior. By using this metric we confirm that our and previous studies successfully generate motions distinct from the training sequences.

We find that our proposed approach outperforms the current SOTA for action-conditional motion generation, ACTOR [31], in terms of realism and diversity. By employing an identical decoder architecture as ACTOR [31], we conduct a fair comparison between our INR-approach and a VAE-approach and find that our INR approach improves motion generation. Furthermore, since Transformer models can be difficult and expensive to train and we also explore the use of an MLP decoder and find that even such a simple, lightweight (6x fewer parameters) model can reach the SOTA performance.

Our key contributions are summarized as follows:

- We propose a variational INR framework for motion generation, which gives time-independent, optimal representations for variable-length sequences distributed such that representations for novel motions can be easily sampled.
- To improve action-conditional motion generation, we propose INRs that are decomposed into action-wise and sequence-wise INRs. The action-wise INR generalizes to features across an action-class and helps generating realistic and novel motions for a target action class.

- We show in experiments that our method outperforms SOTA (ACTOR [31]) on the HumanAct12, NTU13 and UESTC datasets in term of realism and diversity, and confirm that it generates high-quality variable-length sequences. For example on HumanAct12 we generate sequences with lengths between 8 -470 time-steps and find that our motions generated with variable-length even outperform fixed-length motions generated by previous works (ACTOR[31], Action2Motion[10]) in terms of realism and diversity.

2 **Related Works**

In the following we review the context of our work first regarding human motion modeling and then regarding implicit neural representations.

Human Motion Modeling: The modeling of human motions is important for understanding and predicting human behavior. Most modern approaches regard a human motion as a time-series of either skeleton poses or full 3D body shapes [22,29] and previous works have proposed methods to estimate motions from videos, predicting future motions based on past motions and generating such motions conditional on signals such as emotion [13] and rhythm [19]. Our work is similar to [31,10], which generate motions conditional on the action class.

Previous works for motion generation are mostly based on Variational Auto-Encoders (VAE) [17,11,10,31], which employs an encoder to infer a variational distribution from which representations of motions can be sampled and a decoder which reconstructs data from a representation. This encoder is optimized with respect to a reconstruction loss for a whole dataset, without a guarantee that the representation for each individual sequence is optimal. The encoder may focus on the most common features in the dataset and become insensitive to rare features. In contrast, Implicit Neural Representations (INR), which optimize the representation of each sample directly, can be sensitive even to unique features.

A similarity of all sequence modeling approaches is the use of model architectures such as RNNs or Transformers. RNN are typically formulated as an auto-regressive model [2,10,13], which generates motions by recursively predicting the pose at time-step t based on the prediction of the pose at time-step t-1. This recursive nature of RNNs means that their sequence-representations are time-dependent and representations of variable-length sequences can not be easily compared. Furthermore, the recursive generation procedure accumulates error and may result in poor performance when generating long sequences.

Petrovich et al. [31] proposes ACTOR, a Transformer VAE, which yields a single fixed-length representation for a variable-length sequence through a Transformer encoder. This representation is decoded by a Transformer decoder, which receives the representation and the temporal embedding of the target time-steps as input and generates the target sequence in one forward-pass. Since such a Transformer VAE should conceptually handle variable-length sequences well, we choose this work as our main baseline. However, [31] reports that even for ACTOR a fine-tuning scheme is needed to enable good performance for variable-length sequences.

4

P. Cervantes et al.

5

Implicit Neural Representations: INR as proposed in [27,26,7] are encoder-free models which instead optimize their parameters to represent and fit a single sample. They have been popularized particularly in 3D modeling and have shown great performance on tasks such as inverse graphics [24,40], image synthesis [16,33] or scene generation [25,8]. While this work is, to the best of our knowledge, the first to explore implicit neural representations in the context of motion modeling, previous works have considered other time-series [20,26,35].

Previous works for data synthesis using INR use a GAN-like approach [16,1,33] for image synthesis. Such approaches don't optimize the INR, but sample representations from a predetermined distribution. These representations are then used by a generator to generate images, which can fool a discriminator. However, for our task the amount of training data required by GANs is problematic.

Another approach that takes inspiration from VAEs are variational INR [27,4]. Most similar to ours is [4], however, this work doesn't optimize the INRs but rather approximates the optimal INR using an empirical Bayes. Furthermore, they predict a variational distribution from the INR and then apply a regularizing loss to this intermediate representation instead of directly regularizing the INR. Instead we directly optimize the mean and variance of the variational distribution as persistent parameters per sample. The approach in [27] optimizes point-estimates of the representations of a sample and regularizes the distribution of these estimates, but doesn't sample stochastically from this distribution for each sample and samples stochastically from it during training.

With this work we would like to show that INR are not only powerful for high dimensional data such as dense 3D point clouds, but that their flexibility is also useful for other domains.

3 Methodology

In this section, we will first describe how to apply INRs to model human motions and decompose the INR into sequence-wise and action-wise representations (3.1). Then we will introduce the proposed variational INRs (3.2), before we discuss how we fit a conditional Gaussian Mixture Model (GMM) to the representation space and how we sample novel sequence-wise representations from it (3.3). Finally we will describe the Mean Maximum Similarity (MMS) as a measure to detect models that only reproduce training samples (3.4).

3.1 Implicit Neural Representations for Motion Modeling

We consider a human motion as a sequence of poses represented by a lowdimensional skeleton. Formally, we denote a skeleton pose of sequence *i* at time *t* as $x_t^i \in \mathbb{R}^{P \times B}$ where *P* is the number of joints and *B* is the dimensionality of the joint representation. The *i*-th motion (a sequence of poses) is denoted as $\mathbf{x}^i = \{x_t^i\}_{t=1}^{T^i}$ with the sequence length T^i .



Fig. 2. Overview of Implicit Motion Modeling. Each representation is composed of two components, the action representation α and the sequence representation β . Instead of inferring these representations from an encoder, we directly optimize the parameters of a posterior normal distribution for both the action representation $(\mu_{\alpha}, \Sigma_{\alpha})$ shared by all sequences with the same action class and sequence representation $(\mu_{\beta}, \Sigma_{\beta})$. The representation, together with a temporal embedding (PE) τ_t of time t is then input to an MLP, which predicts the pose at time t.

For each sequence i, we construct an Implicit Neural Representation (INR) c_i and a decoder D_{θ} (shared among all sequences) that predicts a pose \hat{x}_t^i of sequence i from the INR c_i and a temporal embedding τ_t of time t

$$\hat{x}_t^i = D_\theta(c^i, \tau_t). \tag{1}$$

Note that depending on the decoder architecture, the decoder may process all time-steps of a sequence independently (MLP) or multiple time-steps simultaneously (Transformer). We obtain an INR c_i , shared by all time steps $(t \in \{1, 2, ..., T_i\})$, by minimizing the reconstruction loss \mathcal{L}_{rec}^i . Thus, INRs can represent a sequence of any sequence-length T^i . Also, for a given INR, the decoder can interpolate between time-steps (e.g. t = 0.5) or extrapolate $(t > T_i)$.

To generalize INRs to all features of the same action class, we decompose the INR and introduce an action representation shared across all samples of the same action class. Formally, we divide each INR c^i into a sequence-wise representation $\beta^i \in \mathbb{R}^S$ with in $i \in \mathcal{M}$ with a set of motions \mathcal{M} and an actionwise representation $\alpha^z \in \mathbb{R}^A$ shared by all sequences with the same action label $z \in \mathcal{Z}$. Here \mathcal{Z} is the set of action classes (e.g. $\alpha^z \in {\alpha^{\text{sit}}, \alpha^{\text{walk}}, \alpha^{\text{run}} \dots})$ and S and A denote the size of each representation respectively.

3.2 Variational Implicit Neural Representations

Note that each INR c^i is optimized to reconstruct a single sample with an overparameterized decoder D_{θ} . This can make the distribution of INRs complex and result in a representation space where a simple distance measure doesn't correspond to semantic similarity. Accordingly, interpolations between representations in this space may not be meaningful. To avoid such a complex representation space, we introduce a variational approach as regularization [17,11]. We formulate each INR as a normal distribution, whose mean μ^i and covariance matrix Σ^i are optimized and from which we sample an instance with the re-parameterization trick during training. This makes the representation space smoother so that close representations are semantically similar. We summarize the sequence-wise and action-wise variational representations as

$$c^{i} \sim \mathcal{N}(\mu^{i}, \Sigma^{i}) \text{ with}$$

$$\mu^{i} = \operatorname{concat}(\mu^{z}_{\alpha}, \mu^{i}_{\beta}),$$

$$\Sigma^{i} = \begin{bmatrix} \Sigma^{z}_{\alpha} & 0^{A \times S} \\ 0^{S \times A} & \Sigma^{i}_{\beta} \end{bmatrix},$$
(2)

where concat denotes the concatenation operation.

Furthermore, by assuming a standard normal distribution as the prior of each INR c_i , we further encourage a simple and compact representation space. We then use the Kullback-Leibler (KL) Divergence \mathcal{L}_{KL}^i as a regularizing loss

$$\mathcal{L}_{\mathrm{KL}}^{i} = \mathcal{D}_{\mathrm{KL}}(\mathcal{N}(\mu^{i}, \Sigma^{i}) \| \mathcal{N}(0, I)).$$
(3)

The sequence wise training objective of our method is thus

$$\mathcal{L}^{i} = \mathcal{L}^{i}_{\rm rec} + \lambda \mathcal{L}^{i}_{\rm KL}, \qquad (4)$$

where \mathcal{L}_{rec}^{i} is the reconstruction term

$$\mathcal{L}_{\text{rec}}^{i} = -\mathbb{E}_{c^{i} \sim \mathcal{N}(\mu^{i}, \Sigma^{i})} \sum_{t=1}^{T^{i}} \log p(x_{t}^{i} | c^{i}, \theta),$$
(5)
$$\log p(x_{t}^{i} | c^{i}, \theta) \propto \|x_{t}^{i} - D_{\theta}(\tau_{t}, \alpha^{z}, \beta^{i})\|_{2} + \text{const.},$$

and $\mathcal{L}_{\mathrm{KL}}^{i}$ is the regularizing KL divergence moderated by a weight λ .

We define the optimization problem for the model parameters as:

$$\theta^{\star} = \underset{\theta}{\operatorname{argmin}} \sum_{z=1}^{Z} \underset{\mu_{\alpha}^{z}, \Sigma_{\alpha}^{z}, i \in \mathcal{M}^{z}}{\min} \underbrace{\underset{\mu_{\beta}^{i}, \Sigma_{\beta}^{i}}{\min}}_{\text{sequence-wise minimum}}, \quad (6)$$

where \mathcal{M}^z denotes a set of sequence indices within action class z. We optimize action-wise parameters $\mu_{\alpha}, \Sigma_{\alpha}$ for each action z:

$$(\mu_{\alpha}^{z*}, \Sigma_{\alpha}^{z*}) = \operatorname*{argmin}_{\mu_{\alpha}^{z}, \Sigma_{\alpha}^{z}} \sum_{i \in \mathcal{M}^{z}} \min_{\mu_{\beta}^{i}, \Sigma_{\beta}^{i}} \mathcal{L}^{i}.$$
 (7)

Likewise, for the sequence-wise parameters we define the optimization problem for each sequence i as:

$$(\mu_{\beta}^{i\star}, \Sigma_{\beta}^{i\star}) = \underset{\mu_{\beta}^{i}, \Sigma_{\beta}^{i}}{\operatorname{argmin}} \mathcal{L}^{i}.$$
(8)

3.3 Conditional GMM of Representation Space

To generate new sequences for a target action class, we need novel samples from the distribution of sequence-wise representations. In the distribution of sequence-wise representations obtained during training, semantic factors such as sequence-lengths and action classes may be entangled. Accordingly, the actionconditional distribution of sequence-wise representations may differ from the standard normal distribution. To control sequence-length and action class for motion generation, we fit a conditional Gaussian Mixture Model (GMM) to the sequence-wise representations β^i sampled 50 times from the variational distributions $\beta^i \sim \mathcal{N}(\mu_{\beta}^{i\pi}, \Sigma_{\beta}^{i\pi})$ for each training sequence.

We fit such a conditional GMM by first constructing subsets of sequence-wise representations that have the same action class z and a sequence-length within the range $[T, T + \Delta T]$. We choose the size of the sequence-length range ΔT to ensure a minimum number of samples in each subset and then fit an independent GMM to each subset of sequence-wise representations. The details for how we select such a set of sequence-length ranges are provided in Appendix A.2. Finally, we obtain the GMM of $p(\beta|z, [T, T + \Delta T])$.

To sample new sequence representations and generate corresponding novel motion sequences, we need to provide a target action class and sequence length. We sample a new sequence representation β^{new}

$$\beta^{\text{new}} \sim p(\beta|z, [T, T + \Delta T]),$$
(9)

by sampling from the GMM corresponding to the target action class and sequence length. With a new sequence representation we generate a new motion

$$\mathbf{x}^{\text{new}} = \{ D_{\theta^{\star}}(\alpha^{z\star}, \beta^{\text{new}}, \tau_t) \}_{t=t_0}^{T'}$$
(10)

with the target action code α^{z*} (obtained during the training stage) and the target sequence length $T' \in [T, T + \Delta T]$.

3.4 Mean Maximum Similarity

By increasing the number of components of the GMM, it can better fit the training distribution, which improves the realism of generated motions. However, we also risk fitting a GMM which only reproduces motions in the training set. Previous metrics such as the Fréchet Inception Distance (FID) or the Diversity compute the feature distribution of training and generated motions and compare these distributions. Generated motions that have a similar distribution (FID) or variance (Diversity) as real motions are considered high-quality. Generated motions identical to the training set would be considered best by such metrics.

To detect models that just reproduce training samples, we introduce the Mean Maximum Similarity (MMS) as a complementary metric. Similarly to previous metrics we extract the features from all training sample and generated motions. Then for each generated motion, we find the training sample with the smallest feature distance (most similar) to it. The mean distance over a large set of generated motions should be small for models that reproduce training samples and large for models that generate novel motions. Formally we denote the features of a motion as f and the sets of generated and training motion sequences \mathcal{M}_{gen} and $\mathcal{M}_{\text{train}}$ respectively, and compute the MMS as

$$\mathcal{D}_{MMS}(\mathcal{M}_{\text{gen}}, \mathcal{M}_{\text{train}}) = \frac{1}{|\mathcal{M}_{\text{gen}}|} \sum_{i \in \mathcal{M}_{\text{gen}}} \min_{j \in \mathcal{M}_{\text{train}}} (\|f_i - f_j\|_2).$$
(11)

We estimate the MMS of model that only reproduces motions as baseline by computing $\mathcal{D}_{MMS}(\mathcal{M}_{\text{train}}, \mathcal{M}_{\text{train}})$ of the set of training motions $\mathcal{M}_{\text{train}}$ against itself. A large gap between $\mathcal{D}_{MMS}(\mathcal{M}_{\text{gen}}, \mathcal{M}_{\text{train}})$ and $\mathcal{D}_{MMS}(\mathcal{M}_{\text{train}}, \mathcal{M}_{\text{train}})$ indicates novel generated motions distinct from the training set.

4 Experiments

To verify the quality of motions generated by variational INR we perform experiments with a Transformer and an MLP decoder. The Transformers is a powerful, but costly and difficult to train modeling tool, while the MLP is simple and comparatively light-weight. The comparison should highlight the efficiency of the variational INR framework independent of decoder architecture. In this section we will first explain the implementation details of our models (4.1) and the datasets for our experiments (4.2). Then we will describe how we quantify the realism, diversity and novelty of generated motions (4.3). Finally we will discuss the quantitative (4.4) and qualitative (4.5) results.

4.1 Implementation

Skeleton representation: We represent the human body as a kinematic tree defined by joint rotations, bone-lengths and the root joint. More specifically, we use the SMPL model [22] with pose parameters consisting of 23 joint rotations, 1 global rotation and 1 root trajectory. During training we only predict the pose parameters, which are independent of the body shape and can be used to animate any body at test time. We represent rotations with a 6D rotation parameterization as proposed by [41] which means the full body pose has 147 dimensions $(24 \times 6+3)$. We use a reconstruction loss composed of a loss on the pose parameters (joint rotations and root joint locations) as well as the vertices of the SMLP model since [31]'s findings suggest the best performance for this configuration. On the NTU13 dataset, where at the time of writing the SMPL data was no longer available we represent the pose with a 6D rotation parameterization, but use a reconstruction loss on the joint locations (through forward-kinematics) as proposed by [10] and find similarly high performance with our method.

Model Architecture: We implement our MLP decoder with ELU activations and 5 hidden layers (1000, 500, 500, 200, 100). The input are temporal embeddings with 256 dimensions and sequence-wise representations/action-wise representations, which are both 128 dimensions respectively, and the decoder outputs 147

dimensional pose parameters. This results in a network with 1,399,147 parameters. Due to the larger dataset size of UESTC we also implement a larger model (2000, 2000, 1000, 1000, 200, 100) with 8,265,147 parameters which is only used on UESTC. We also implement a Transformer-decoder (same as ACTOR [31]) with 8 layers, 4 attention heads, a dropout rate of 0.1 and a feedforward network of 1024 dimensions. With temporal embeddings with 256 dimensions and the same pose parameterization this results in a network with 8,465,299 parameters (6× more than the MLP model). More details are provided in Appendix A.1.

Note that the Transformer-decoder is sensitive to the initialization of the implicit representations. If the variance parameters are initialized with a high variance the Transformer-decoder may fail to converge, while the MLP decoder is not sensitive to this phenomenon. We explore this more in Appendix B.1.

The Transformer-based decoder has an identical structure to ACTOR [31] and thus allows us a direct comparison between an auto-encoder and an implicit framework. The MLP decoder is simpler to train than the Transformer decoder and doesn't rely on self-attention. The comparison of these decoders allows us to determine if the choice of decoder architecture is critical for good performance.

4.2 Datasets

To evaluate the quality of action-conditional human motion generation, we used the UESTC, NTU-RGBD and HumanAct12 dataset curated by [10].³

HumanAct12 [10]: This dataset is based on PHSPD [42] and consists of 1191 motion clips and 90099 frames in total. Action labels for 12 actions are provided with at least 47 and at most 218 samples per label. Sequence-lengths range from 8 to 470. We follow the procedure by [31] to align the poses to frontal view.

NTU13 [21]: The NTU-RGBD dataset originally contains pose annotations from a MS Kinect sensor and label annotations for 120 actions. [10] re-estimated the data of a subset of 13 action, which we denote NTU13, with a state-of-the-art pose estimation method [18] to reduce noise. In this refined subset each action label has between 286 - 309 samples. The refined poses have 18 body joints and the sequence lengths range from 20 - 201.⁴

UESTC [14]: This dataset with 40 action classes, 40 subjects and 25K samples is the largest dataset we perform experiments on and the only dataset with a train/test split. We use the SMPL sequences provided by [31] and apply the same pre-proprocessing, namely we rotate all sequences to frontal view. Using the same cross-subject testing protocol we have a training split with between 225 - 345 samples per action class and sequence lengths between 24 and 2891 time steps (on average 300 time steps).

³ We considered the CMU Mocap dataset, but manual inspection found the label annotations for some actions such as "Wash" and "Step" to be extremely noisy.

⁴ Due to the release agreement of NTU RGBD, this subset can no longer be distributed. We report results to provide a complete comparison to previous studies.

4.3 Evaluation Metrics

We use the same evaluation metrics as [10,31] (Fréchet Inception Distance (FID), action recognition accuracy, diversity and multimodality) to measure the realism and diversity of generated motions. Also, we measure the proposed **Mean Maximum Similarity** (Section 3.4) to detect models that reproduce training samples. We report a 95% confidence interval computed of 20 evaluations.

The features for these evaluation metrics are extracted from motions of a predetermined length (60 time-steps) by an RNN-based action recognition model (weights provided by [10]) for the NTU13 and HumanAct12 dataset and by an ST-GCN-based action recognition model (weights provided by [31]) for the UESTC dataset. However, since the real training data is variable-length, we follow [10]'s procedure for feature extraction during evaluation. This procedure adjust all sequences to a target length, by repeating the last pose of short sequences and by sampling random sub-sequence from longer sequences.

Such stochastic feature extraction means the MMS may not be zero, even for sets of identical motions. Thus we first compute a baseline MMS for identical real motions and then evaluate the MMS between real motions and generated motions. If the MMS for generated motions is larger than that of real motions only, we conclude that the generated motions are distinct from the real motions.

We find that there is a difference in the evaluation procedure of previous works in the sampling frequency of different action classes for generation. The approach by [10] generates motions uniformly for all action classes. On datasets with an action imbalance, this creates an inflated FID score. We follow [31]'s approach which generates motions according to the frequency of the action class in the training dataset, since this leads to more consistent results.

Our GMM samples novel representations conditional on the sequence-length. For model evaluation we sample sequence-lengths according to their distribution in the training dataset. We then sample corresponding representations and generate motions with the corresponding sequence-length. We perform the same feature extraction as for the variable-length real motions. More details can be found in the Appendix A.7

4.4 Quantitative Results

We compare our method to an RNN [10] and a Transformer [31] baseline and present some ablations for the proposed novel components of our model on HumanAct12 and NTU13 in Table 1 and UESTC in Table 2. Furthermore, we present a new state-of-the-art with the results for our Transformer-based and MLP-based models. We also investigate the contribution of variational INR by comparing them to a non-variational version and the contribution of the decomposed representations by comparing to a version with no action code. Note that by construction our motion generation procedure can generate high-quality motions for arbitrarily specified sequence lengths (as in Table 1) within the variation of training sequence lengths, whereas previous works reported a performance drop for variable-length generation.

Table 1. Comparison on HumanAct12 and NTU13 (The best in bold, the second best underlined). *Non-variational* uses action codes and *no action code* uses the variational approach. (\pm indicates 95% confidence interval, \rightarrow closer to real is better)

Matha J	HumanAct12				
Method	$\mathrm{FID}\downarrow$	Accuracy \uparrow	$\text{Diversity} \rightarrow$	Multimod. \rightarrow	
Real	$0.020^{\pm.010}$	$0.997^{\pm.001}$	$6.850^{\pm.050}$	$2.450^{\pm.040}$	
Action2Motion [10]	$0.338^{\pm.015}$	$0.917^{\pm.003}$	$6.879^{\pm.066}$	$2.511^{\pm.023}$	
ACTOR [31]	$0.12^{\pm.00}$	$0.955^{\pm.008}$	$6.84^{\pm.03}$	$2.53^{\pm.02}$	
INR (Transformer)	$0.088^{\pm.004}$	$0.973^{\pm.001}$	$6.881^{\pm.048}$	$2.569^{\pm.040}$	
INR (MLP)	$0.114^{\pm.001}$	$0.970^{\pm.001}$	$6.786^{\pm.057}$	$2.507^{\pm.034}$	
- (Non-variational)	$0.551^{\pm.005}$	$0.795^{\pm.002}$	$6.800^{\pm.046}$	$3.700^{\pm.032}$	
- (No action code)	$0.146^{\pm.003}$	$0.955^{\pm.001}$	$6.797^{\pm.066}$	$2.769^{\pm.045}$	
	NTU13				
Real	$0.031^{\pm.004}$	$0.999^{\pm.001}$	$7.108^{\pm.048}$	$2.194^{\pm.025}$	
Action2Motion [10]	$0.351^{\pm.011}$	$0.949^{\pm.001}$	$7.116^{\pm.037}$	$2.186^{\pm.033}$	
ACTOR [31]	$0.11^{\pm.00}$	$0.971^{\pm.002}$	$7.08^{\pm.04}$	$2.08^{\pm.01}$	
INR (Transformer)	$0.097^{\pm001}$	$0.977^{\pm001}$	$7.060^{\pm.040}$	$2.108^{\pm.025}$	
INR (MLP)	$0.113^{\pm.001}$	$0.976^{\pm.001}$	$7.070^{\pm.052}$	$2.070^{\pm.043}$	
- (Non-variational)	$0.646^{\pm.003}$	$0.849^{\pm.001}$	$6.905^{\pm.056}$	$3.244^{\pm.049}$	
- (No action code)	$0.202^{\pm.002}$	$0.912^{\pm.001}$	$7.025^{\pm.050}$	$2.648^{\pm.043}$	

The results show that our proposed method improves over both Action2Motion [10] and ACTOR[31] especially on the FID and accuracy metric. We show that our optimized INR outperforms methods with representations produced by an optimized encoder. This is most apparent when comparing our implicit Transformer model and ACTOR, since both models use the same decoder architecture.

Furthermore, we show high performance even with a simple MLP decoder. This shows that the self-attention mechanism is not necessary. We argue, that the common property of the Transformer-based models and our implicit MLP-based model, namely time-independent sequence-wise representations, are critical for motion generation performance. Such representations can avoid the error accumulation of current RNN-based models and represent variable-length sequences.

Comparing the non-variational and variational approach, we find that the realism and diversity is improved with the variational approach (see (Non-variational) in Table 1). This finding suggests that the variational approach strongly regularizes the latent space and improves sampling of new motions. However, both approaches for implicit representations are able to reach similar reconstruction performance for training samples and learn effective representations that allow high quality reconstruction.

Comparing the approach with an action-wise and sequence-wise representation to an approach only with a sequence-wise representation, we find a clear advantage from using a decomposed action representation. Even the approach with

Method	$\mathrm{FID}_{\mathrm{train}}\downarrow$	$\mathrm{FID}_{\mathrm{test}}\downarrow$	Accuracy \uparrow	$\text{Diversity} \rightarrow$	Multimod. \rightarrow
Real	$2.92^{\pm.26}$	$2.79^{\pm.29}$	$0.988^{\pm.001}$	$33.34^{\pm.320}$	$14.16^{\pm.06}$
ACTOR [31]	$20.49^{\pm 2.31}$	$23.43^{\pm 2.20}$	$0.911^{\pm.003}$	$31.96^{\pm.33}$	$14.52^{\pm.09}$
Ours (MLP)	$9.55^{\pm.06}$	$15.00^{\pm.09}$	$0.941^{\pm.001}$	$31.59^{\pm.19}$	$14.68^{\pm.07}$

Table 2. Baseline comparison with UESTC. (\pm indicates 95% confidence interval, \rightarrow closer to real is better

Table 3. Mean Maximum Similarity as a sanity check to detect overfitting.

Method	HumanAct12	NTU13	UESTC
Real	$0.329^{\pm.003}$	$0.209^{\pm.002}$	$4.925^{\pm.007}$
Action2Motion[10]	$0.945^{\pm.006}$	$0.667^{\pm.006}$	—
ACTOR[31]	$0.921^{\pm.001}$	$0.701^{\pm.001}$	$8.645^{\pm 0.008}$
INR (MLP)	$0.941^{\pm.005}$	$0.620^{\pm.001}$	$7.113^{\pm 0.006}$
INR (Transformer)	$0.778^{\pm.003}$	$0.570^{\pm.002}$	—

only sequence-wise representations performs comparable to the RNN-baseline (See (No action code) in Table 1) However, a decomposition representation is needed to outperform the Transformer baseline.

For further ablation studies we refer to Appendix B, where we investigate various modeling choices. Among others, we investigate the effect of the number of components in the GMM and show in Appendix B.2 high performance, on-par with Action2Motion[10], even when using GMMs with a single component.

Finally we perform a sanity check to see if any model is reproducing training samples by checking the proposed Mean Maximum Similarity between real and generated sequences defined by Eq. (11) and show the results in Table 3. We interpret the gap between this baseline of real motions and all other models as an indication that no model is just reproducing training samples. Note that while the INR (Transformer) model has a lower MMS than other models, the gap to the baseline is significant. This suggests that it generates motions that are more similar to training motions than other models yet distinct from them.

4.5 Qualitative Results

We manually inspect the quality of the generated motions and find that our methods consistently generates high-quality motions even over a long time range. In particular we observe, that the RNN-based Action2Motion[10] shows a slow-down effect when predicting long sequences, which neither our methods nor the Transformer-based ACTOR[31] suffer from.

Also, we observe that we reliably generate complete actions, due to our model's ability to model the sequence-length. In contrast, both Action2Motion



Fig. 3. Motions generated by our MLP model trained on HumanAct predicted with 40 time-steps (each third frame shown) for the actions *throw*, *sit* and *jump*.

as well as ACTOR tend to generate incomplete actions, particularly when generating short sequences. As shown in Fig. 3, even for short sequences and actions with a clear start and end our generated actions are complete. For more qualitative results, we refer to supplementary videos and Appendix C.

5 Limitations and Future Work

While our implicit sequence representations are parameter-efficient, the effort to train sequence-wise parameters scales linearly with the size of the training dataset. Furthermore, we observe a sensitivity to the ratio of parameter updates between the sequence-wise parameters and the decoder parameters. If the decoder parameters are updated significantly more often than sequence-wise parameters, our implicit models might perform poorly.

6 Conclusion

We present an MLP-based model for action-conditional human motion generation. The proposed approach improves over previous RNN-based and Transformerbased baselines by employing variational implicit neural representations. We argue that the likely reason for the success of our method is implicit neural representations, which are optimized representations that can represent full variablelength sequences and we supported this hypothesis experimentally by reaching state-of-the-art performance on commonly used metrics for motion generation. While the results of our work may improve technologies for animation and action recognition, there is the potential for malicious use such a deep fakes as well. For more detailed discussions about such potential negative societal impacts and personal data of human subjects we refer to Appendices A.5 and A.6.

Acknowledgements: This work is an outcome of a research project, Development of Quality Foundation for Machine-Learning Applications, supported by DENSO IT LAB Recognition and Learning Algorithm Collaborative Research Chair (Tokyo Tech.). It was also supported by JST CREST JPMJCR1687.

References

- Anokhin, I., Demochkin, K., Khakhulin, T., Sterkin, G., Lempitsky, V., Korzhenkov, D.: Image generators with conditionally-independent pixel synthesis. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 14278–14287 (2021) 5
- Barsoum, E., Kender, J., Liu, Z.: HP-GAN: Probabilistic 3D human motion prediction via GAN. In: Proceedings of the IEEE conference on computer vision and pattern recognition workshops (CVPRW). pp. 1418–1427 (2018) 1, 4
- Battan, N., Agrawal, Y., Rao, S.S., Goel, A., Sharma, A.: GlocalNet: Class-aware long-term human motion synthesis. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV). pp. 879–888 (2021) 1
- Bond-Taylor, S., Willcocks, C.G.: Gradient origin networks. In: 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021. OpenReview.net (2021) 5
- Butepage, J., Black, M.J., Kragic, D., Kjellstrom, H.: Deep representation learning for human motion prediction and classification. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 6158–6166 (2017) 1
- Chen, Y., Liu, C., Shi, B.E., Liu, M.: CoMoGCN: Coherent motion aware trajectory prediction with graph representation. In: British Machine Vision Conference (BMVC) (2020) 1
- Chen, Z., Zhang, H.: Learning implicit fields for generative shape modeling. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2019) 5
- DeVries, T., Bautista, M.A., Srivastava, N., Taylor, G.W., Susskind, J.M.: Unconstrained scene generation with locally conditioned radiance fields. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). pp. 14304–14313 (2021) 5
- Doersch, C., Zisserman, A.: Sim2real transfer learning for 3d human pose estimation: motion to the rescue. Advances in Neural Information Processing Systems (NeurIPS) 32 (2019) 1
- Guo, C., Zuo, X., Wang, S., Zou, S., Sun, Q., Deng, A., Gong, M., Cheng, L.: Action2Motion: Conditioned generation of 3D human motions. In: Proceedings of the 28th ACM International Conference on Multimedia (MM '20) (2020) 1, 2, 4, 9, 10, 11, 12, 13, 19, 21, 22, 23, 24
- Higgins, I., Matthey, L., Pal, A., Burgess, C., Glorot, X., Botvinick, M., Mohamed, S., Lerchner, A.: β-vae: Learning basic visual concepts with a constrained variational framework. In: 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings. OpenReview.net (2017) 2, 4, 6
- Honda, Y., Kawakami, R., Naemura, T.: Rnn-based motion prediction in competitive fencing considering interaction between players. In: British Machine Vision Conference (BMVC) (2020) 1
- Hou, Y., Yao, H., Sun, X., Li, H.: Soul dancer: Emotion-based human action generation. ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM) 15(3s), 1–19 (2020) 1, 4
- Ji, Y., Xu, F., Yang, Y., Shen, F., Shen, H.T., Zheng, W.S.: A large-scale RGB-D database for arbitrary-view human action recognition. In: Proceedings of the 26th ACM International Conference on Multimedia (MM '18) (2018) 10

- 16 P. Cervantes et al.
- Kanazawa, A., Zhang, J.Y., Felsen, P., Malik, J.: Learning 3D human dynamics from video. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2019) 1
- Karras, T., Aittala, M., Laine, S., Härkönen, E., Hellsten, J., Lehtinen, J., Aila, T.: Advances in Neural Information Processing Systems (NeurIPS) 34 (2021) 5
- Kingma, D.P., Welling, M.: Auto-encoding variational bayes. In: 2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings (2014) 2, 4, 6
- Kocabas, M., Athanasiou, N., Black, M.J.: Vibe: Video inference for human body pose and shape estimation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 5253–5263 (2020) 10
- Li, R., Yang, S., Ross, D.A., Kanazawa, A.: Ai choreographer: Music conditioned 3d dance generation with AIST++. In: Proc. of the IEEE International Conf. on Computer Vision (ICCV) (2021) 1, 4
- Li, Z., Niklaus, S., Snavely, N., Wang, O.: Neural scene flow fields for space-time view synthesis of dynamic scenes. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2021) 5
- Liu, J., Shahroudy, A., Perez, M., Wang, G., Duan, L.Y., Kot, A.C.: NTU RGB+D 120: A large-scale benchmark for 3d human activity understanding. IEEE transactions on pattern analysis and machine intelligence 42(10), 2684–2701 (2019) 10
- Loper, M., Mahmood, N., Romero, J., Pons-Moll, G., Black, M.J.: SMPL: A skinned multi-person linear model. ACM Trans. Graphics (Proc. SIGGRAPH Asia) 34(6), 248:1–248:16 (Oct 2015) 4, 9
- Meng, F., Liu, H., Liang, Y., Tu, J., Liu, M.: Sample fusion network: An endto-end data augmentation network for skeleton-based human action recognition. IEEE Transactions on Image Processing 28(11), 5281–5295 (2019) 1
- Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: NeRF: Representing scenes as neural radiance fields for view synthesis. In: European conference on computer vision (ECCV). pp. 405–421. Springer (2020) 2, 5
- Niemeyer, M., Geiger, A.: GIRAFFE Representing scenes as compositional generative neural feature fields. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2021) 5
- Niemeyer, M., Mescheder, L., Oechsle, M., Geiger, A.: Occupancy flow: 4D reconstruction by learning particle dynamics. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) (2019) 2, 5
- Park, J.J., Florence, P., Straub, J., Newcombe, R., Lovegrove, S.: DeepSDF Learning continuous signed distance functions for shape representation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 165–174 (2019) 2, 5
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al.: Pytorch: An imperative style, highperformance deep learning library. Advances in Neural Information Processing Systems (NeurIPS) **32** (2019) 20
- Pavlakos, G., Choutas, V., Ghorbani, N., Bolkart, T., Osman, A.A.A., Tzionas, D., Black, M.J.: Expressive body capture: 3D hands, face, and body from a single image. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2019) 4, 20
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A.,

Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, E.: Scikit-learn: Machine learning in Python. Journal of Machine Learning Research **12**, 2825–2830 (2011) 20

- Petrovich, M., Black, M.J., Varol, G.: Action-conditioned 3D human motion synthesis with Transformer VAE. In: International Conference on Computer Vision (ICCV) (2021) 1, 2, 3, 4, 9, 10, 11, 12, 13, 18, 19, 20, 21, 24
- Ravi, N., Reizenstein, J., Novotny, D., Gordon, T., Lo, W.Y., Johnson, J., Gkioxari, G.: Accelerating 3D deep learning with pytorch3d. arXiv:2007.08501 (2020) 20
- Schwarz, K., Liao, Y., Niemeyer, M., Geiger, A.: GRAF Generative radiance fields for 3d-aware image synthesis. In: Advances in Neural Information Processing Systems (NeurIPS) (2020) 5
- Sengupta, A., Budvytis, I., Cipolla, R.: Synthetic training for accurate 3d human pose and shape estimation in the wild. In: British Machine Vision Conference (BMVC) (September 2020) 1
- Sitzmann, V., Martel, J.N., Bergman, A.W., Lindell, D.B., Wetzstein, G.: Implicit neural representations with periodic activation functions. In: Advances in Neural Information Processing Systems (NeurIPS) (2020) 5
- Starke, S., Zhao, Y., Komura, T., Zaman, K.: Local motion phases for learning multi-contact character movements. ACM Trans. Graph. 39(4) (Jul 2020) 1
- Starke, S., Zhao, Y., Zinno, F., Komura, T.: Neural animation layering for synthesizing martial arts movements. ACM Trans. Graph. 40(4) (Jul 2021) 1
- Varol, G., Laptev, I., Schmid, C., Zisserman, A.: Synthetic humans for action recognition from unseen viewpoints. International Journal of Computer Vision 129(7), 2264–2287 (2021) 1
- Yan, S., Li, Z., Xiong, Y., Yan, H., Lin, D.: Convolutional sequence generation for skeleton-based action synthesis. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) (2019) 1
- Yariv, L., Kasten, Y., Moran, D., Galun, M., Atzmon, M., Ronen, B., Lipman, Y.: Multiview neural surface reconstruction by disentangling geometry and appearance. Advances in Neural Information Processing Systems (NeurIPS) 33 (2020)
- Zhou, Y., Barnes, C., Lu, J., Yang, J., Li, H.: On the continuity of rotation representations in neural networks. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 5745–5753 (2019) 9
- Zou, S., Zuo, X., Qian, Y., Wang, S., Xu, C., Gong, M., Cheng, L.: 3D human shape reconstruction from a polarization image. In: Proceedings of the European Conference on Computer Vision (ECCV) (2020) 10