

Weakly-Supervised Deep Recurrent Neural Networks for Basic Dance Step Generation

Nelson Yalta^{*}, Shinji Watanabe[†], Kazuhiro Nakadai[‡] and Tetsuya Ogata^{*}

^{*}*Department of Intermedia Art and Science, Waseda University, Tokyo, Japan*

[†]*Johns Hopkins University, Baltimore, USA*

[‡]*Honda Research Institute Japan, Saitama, Japan*

**nelson.yalta@ruri.waseda.jp*

Abstract—Synthesizing human’s movements such as dancing is a flourishing research field which has several applications in computer graphics. Recent studies have demonstrated the advantages of deep neural networks (DNNs) for achieving remarkable performance in motion and music tasks with little effort for feature pre-processing. However, applying DNNs for generating dance to a piece of music is nevertheless challenging, because of 1) DNNs need to generate large sequences while mapping the music input, 2) the DNN needs to constraint the motion beat to the music, and 3) DNNs require a considerable amount of hand-crafted data. In this study, we propose a weakly supervised deep recurrent method for real-time basic dance generation with audio power spectrum as input. The proposed model employs convolutional layers and a multilayered Long Short-Term memory (LSTM) to process the audio input. Then, another deep LSTM layer decodes the target dance sequence. Notably, this end-to-end approach has 1) an auto-conditioned decode configuration that reduces accumulation of feedback error of large dance sequence, 2) uses a contrastive cost function to regulate the mapping between the music and motion beat, and 3) trains with weak labels generated from the motion beat, reducing the amount of hand-crafted data. We evaluate the proposed network based on i) the similarities between generated and the baseline dancer motion with a cross entropy measure for large dance sequences, and ii) accurate timing between the music and motion beat with an F-measure. Experimental results revealed that, after training using a small dataset, the model generates basic dance steps with low cross entropy and maintains an F-measure score similar to that of a baseline dancer.

Index Terms—Deep recurrent networks; Contrastive loss; Dance generation

I. INTRODUCTION

Dancing is a performing art and expresses meaning or people’s emotion [1]. Dance is composed of sequential rhythmical motion units called basic movements (i.e., basic steps) [2]. Recently, methods to synthesize dance movements are actively investigated in various domains. Some research has led to the development of applications that go beyond merely generating dance motion for robots [3], [4], animated computer graphics, animated choreographies [5], and video games [2]. Besides, various motion and dance generation approaches have been proposed in recent years, including probabilistic models [5], Boltzmann machines, and artificial neural networks [6]. Recent findings demonstrate the remarkable performance of

applying deep learning in motion generation [7]–[9], music processing [10], [11] and several other tasks. An advantage of using deep learning is the relatively low effort for feature engineering [12]. Deep learning also enables an end-to-end mapping between the input and output features, reducing the requirement for intermediate hand-crafted annotations [10].

Motivated by the benefits of deep learning, in this study, we explore the application of deep learning models in the generation of basic dance steps. Though deep learning offers the benefits mentioned above, applying it to dance generation can be challenging because of the following three main issues: 1) deep learning models may not generate variable-length non-linear sequences [13] such as dance; 2) the given model may not be able to constrain the motion beat [1], [2], [14] to music beat, and 3) the performance of deep learning models is proportional to the number of training datasets; thus, they require large carefully-labeled datasets for good performance [12].

In this study, we propose a weakly-supervised deep recurrent model for generating basic dance synchronized to the music rhythm. Following the research carried out in [8] and [15], the model proposed in this paper employs multilayered Long Short-Term Memory (LSTM) layers and convolutional layers to encode the audio power spectrum. The convolutional layers reduce the frequency variation of the input audio and the LSTM layers model the time sequence features. Besides, we employ another deep LSTM layer with an **auto-conditioned** configuration [9] to decode the motion. This configuration enables the model to handle a more extended dance sequence with low noise accumulation, which is fed back into the network. To ensure alignment between the motion and music beat, we utilize a **contrastive cost function** [16] for music-motion regulation. The contrastive cost function is a measure of similarities between the given inputs, and it minimizes the distance of the input patterns if the inputs are similar; otherwise, the distance is maximized. Taking advantage of the sequential characteristic of the training, we utilize the Euclidean distance measured by the contrastive loss to track the music beat between two consecutive audio input frames. Then, the model is trained to maximize the distance when the beat is present, i.e., different; otherwise, it minimizes the distance. To constrain the motion beat to the music beat [1], we assume that a music beat may be present when a motion beat appears.

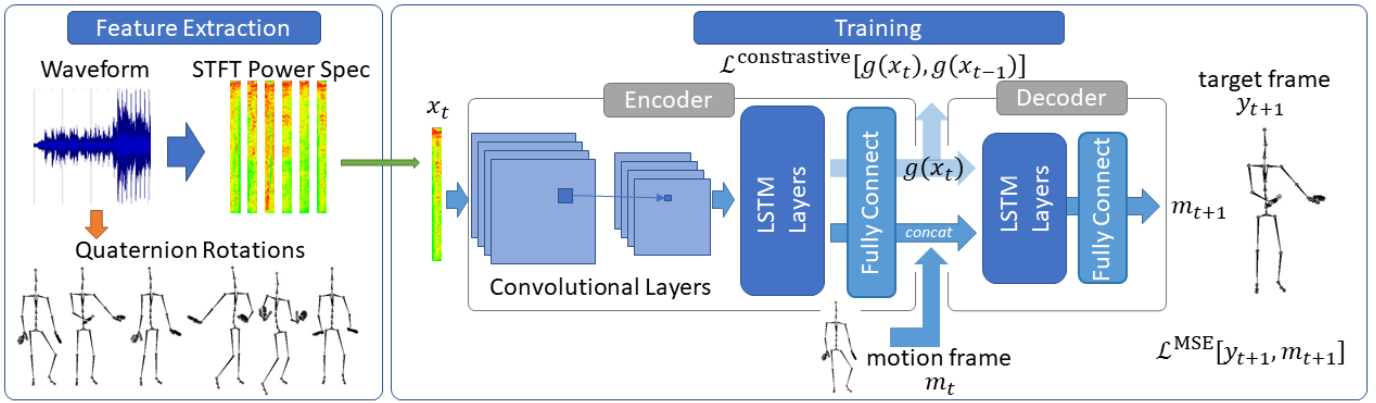


Fig. 1. Framework.

To avoid additional hand-crafted annotations and reduce the amount of training data, the contrastive cost function employs **weak labels** (see [12], [17], [18]) that are generated by motion direction. *Weak* labels may yield incomplete, inaccurate or incorrect data [19]. However, a contrastive cost function trained with weak labels supports training with a small number of samples and avoids pre-training; therefore, it reduces the need for high computational capabilities. The proposed model demonstrates improved music-motion regulation.

The primary contributions of this study are summarized as follows:

- We propose a deep recurrent neural network (DRNN) (Section III-A) with an auto-condition configuration (Section III-B) to generate long dance sequences using the audio spectrum as input.
- Using contrastive cost function, we explore the regulation of the alignment between music and motion beat; detecting the music beat between two consecutive audio input frames (Section III-C).
- Additionally, using motion direction, we explore the generation of *weak* labels for motion-music alignment. Section III-C shows that this configuration reduces the need for additional annotations or hand-crafted labeled data, while we describe the feature extraction and training setup in Section IV.
- We conduct evaluations on the motion beat accuracy and cross entropy of the generated dance relative to the trained music (Section V). Furthermore, we demonstrate that the proposed approach increases the precision of the motion beat along with the music beat; moreover, the approach models basic dance steps with lower cross entropy.
- Conclusions and suggestions for potential future enhancements of the proposed model are given in Section VI.

II. RELATED WORKS

Several studies have considered different approaches to handle variable-length sequences, such as dance. In [7], a factored conditional restricted Boltzmann machine and recurrent neural network (RNN) was employed to tackle the

non-linear characteristic of dance. The model maps non-linear characteristics between audio and motion features and generates a new dance sequence. A generative model was presented in [6] to generate a new dance sequence for a solo dancer. However, the model requires significant computational capabilities or large datasets; moreover, the generated dance is constrained by the trained data.

Dancing involves significant changes in motion that occur at regular intervals, i.e., a motion beat (see, e.g., [1], [2], [14]); and when dancing to music, the music and motion beat should be synchronized. In earlier studies, the music beat [8] and other musical features [5] were utilized to improve dance generation. However, the generated features require additional effort to obtain additional annotations.

In this paper, we perform dance generation with an audio spectrum as an input and motion-music alignment with weak labels.

III. PROPOSED FRAMEWORK

An overview of the proposed system is shown Fig. 1.

A. Deep Recurrent Neural Network

Mapping high-dimensional sequences such as motion is a challenging task for deep neural networks (DNN) [13] because such sequences are not constrained to a fixed size. Besides, to generate motion from music, the model under consideration must map highly non-linear representations between music and motion [7]. In time signal modeling [20], DRNNs implemented with LSTM layers displayed remarkable performance and stable training when deeper networks are employed. Furthermore, the application of stacked convolutional layers in a DRNN, referred to as CLDNN, has demonstrated promising results for speech recognition tasks [15]. In the CLDNN, the convolutional layers reduce the spectral variation of the input sound while the LSTM layers perform time signal modeling. To construct our model, we consider a DRNN with LSTM layers separated into two blocks [20]: one block plays the role of reducing the music input sequence (encoder), and the other is for the motion output sequence (decoder). This configuration can handle non-fixed dimensional signals, such

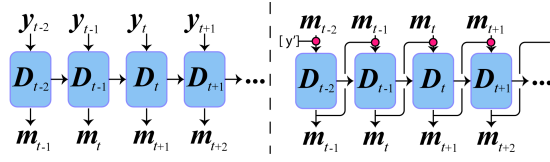


Fig. 2. Auto-conditioned decoder.

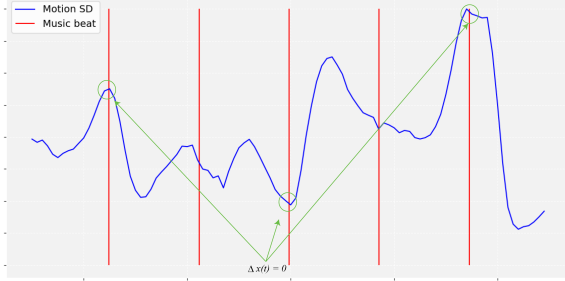


Fig. 3. Synchronization of motion and music beat.

as motion, and avoids performance degradation due to the long-term dependency of RNNs.

The input to the network is the power spectrum from the audio represented as $x_{1:n} = (x_t \in \mathbb{R}^b | t = 1, \dots, n)$ with n frames and b frequency bins, and the ground truth sequence is represented as $y_{1:n} = (y(t) \in \mathbb{R}^j | t = 1, \dots, n)$ with n frames and j joint axes. The following equations show the relations of the motion modeling:

$$g(x_t) = \text{LSTM}^{le}(x'_t); \quad (1)$$

$$m_{t+1} = \text{LSTM}^{ld}(g(x_t)), \quad (2)$$

where $g(x_t)$ is the output processed by the encoder with le layers, and x'_t is the output from the convolutional layers. The network output m_{t+1} is processed from the current input and previous states of the decoder with ld layers (Fig. 2 left). Then, m_{t+1} is a l_2 -norm model.

However, with time-series data (such as dance data), the model may freeze, or the output may diverge from the target due to accumulated feedback errors. To address these issues, the value of the output of the decoder is set to consider auto-regressive noise accumulation by including previously generated step in the motion generation process.

B. Auto-conditioned Decoder

A conventional method uses the ground truth of the given sequence as an input to train sequences with RNN models. During evaluations, the model accustomed to the ground truth in the training process may freeze or diverge from the target due to the accumulation of slight differences between the trained and a self-generated sequence.

By conditioning the network using its output during training, the auto-conditioned LSTM layer handles errors accumulated during sequence generation. Thus, the network can handle

large sequences from a single input, maintain accuracy, and mitigate error accumulation.

In the research carried out in [9] where an auto-conditioned LSTM layer for complex motion synthesis was employed, the conditioned LSTM layer was trained by shifting the input from the generated motion with the ground truth motion after fixed repetitive steps. In the method proposed in this paper, we only employ ground truth motion at the beginning of the training sequence (see Fig. 2 right). By modifying Eq. 2, the generated output is expressed as follows:

$$m_{t+1} = \text{LSTM}^{ld}([g(x_t), y'_t]), \quad (3)$$

where

$$y'_t = \begin{cases} y_t & \text{if } t = 0, \\ m_t & \text{otherwise.} \end{cases} \quad (4)$$

The motion error is calculated using a mean squared error (MSE) cost function expressed as follows:

$$\mathcal{L}^{\text{MSE}} = \frac{1}{k} \sum_{i=1}^k (y_{t+1} - m_{t+1})^2, \quad (5)$$

where k is the training batch size, $y(t+1)$ is the ground truth, and $m(t+1)$ is the generated motion.

We employ a zero vector as the input of the first step in our evaluations, followed by a self-generated output to generate the dance until the music stops.

C. Music-motion Alignment Regulation with Weak Labels

Motion beat is defined as significant changes in movement at regular intervals. An earlier study [1] revealed that motion beat frames occur when the direction of the movement changes; thus, a motion beat occurs when the speed drops to zero. Furthermore, harmony is a fundamental criterion when dancing to music; hence, the music and motion beat should be synchronized.

For basic dance steps, repetitions of dance steps are given by a repetitive music beat, where the direction of the movement changes drastically (see Fig. 3). To avoid the use of additional information, employing the previous definition we formalize *the extracted music features, which are found to be different compared to the previous frame (i.e., $g(x_t) \neq g(x_{t-1})$) when a beat occurs; otherwise, it may maintain a similar dimension, i.e., $g(x_t) \approx g(x_{t-1})$* (see Fig. 4). This procedure generates *weak labels*.

Regarding regulation, we employ a contrastive cost function [16] that maps a similarity metric to the given features.

To utilize the contrastive loss, we extract the standard deviation (SD) of the ground truth motion at each frame and compare it to that of the next frame. Then, we assign a label equal to 1 when the motion maintains its direction; otherwise, the label is 0.

At t :

$$dv_t = \text{SD}(y_t) \quad (6)$$

$$s_t = \text{sign}(dv_t - dv_{t-1}), \quad (7)$$

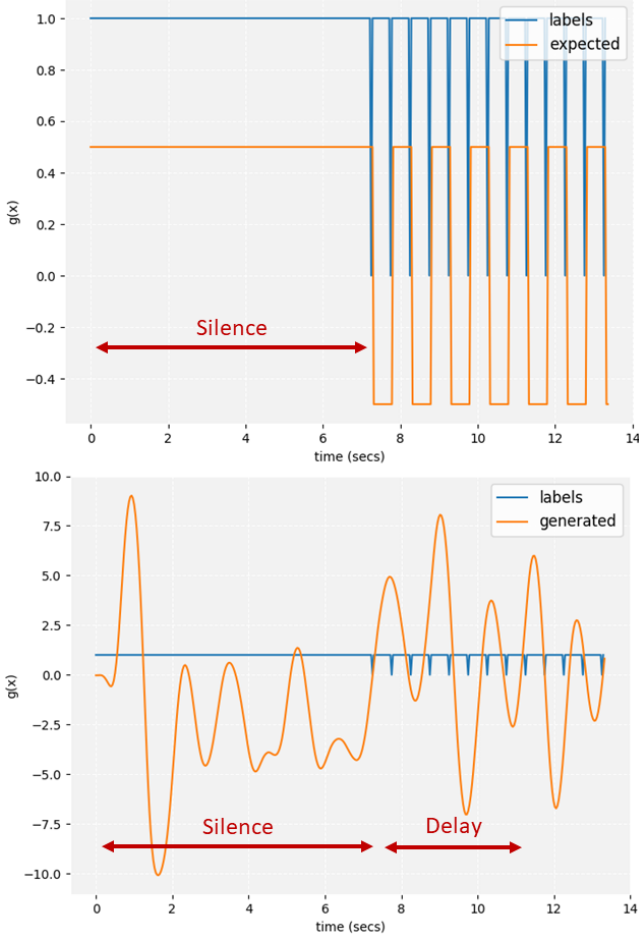


Fig. 4. Music-motion regulation. Top: Expected distance, Bottom: Generated distance.

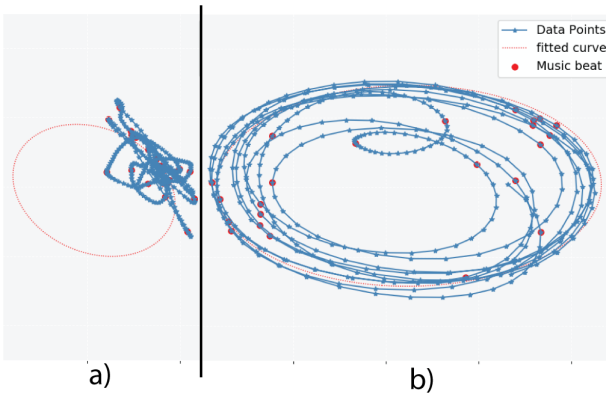


Fig. 5. Encoder PCAs: a) DRNN model without contrastive cost function, b) DRNN with contrastive cost function.

TABLE I
MODEL ARCHITECTURE.

Layer Name	Parameters
conv1	33×2 , 16 channels, stride 1
conv2	33×2 , 32 channels, stride 1
conv3	33×2 , 64 channels, stride 1
conv4	33×2 , 65 channels, stride 1
enc_lstm1	500 units
enc_lstm2	500 units
enc_lstm3	500 units
fc01	65-d fc, ELU
dec_lstm1	500 units
dec_lstm2	500 units
dec_lstm3	500 units
out	71-d fc, ELU

and at $t + 1$:

$$s_{t+1} = \text{sign}(dv_{t+1} - dv_t). \quad (8)$$

A label d is expressed as follows:

$$d = \begin{cases} 1 & \text{if } s_{t+1} = s_t, \\ 0 & \text{otherwise.} \end{cases} \quad (9)$$

The contrastive cost function at frame $t + 1$ is expressed as:

$$\mathcal{L}^{\text{contrastive}} = \frac{1}{2}(d_{t+1}g^2 + (1 - d_{t+1})\max(1 - g, 0)^2) \quad (10)$$

where $g = \|g(x_{t+1}) - g(x_t)\|^2$.

Finally, the cost function of the model is formulated as follows:

$$\mathcal{L}_y = \mathcal{L}^{\text{MSE}}[y_{t+1}, m_{t+1}] + \max(\mathcal{L}^{\text{contrastive}}[g_{t+1}, g_t], 0). \quad (11)$$

In this manner, we synchronize the music beat to the motion beat without requiring additional annotations or further information regarding the music beat. Figure 5 shows the behavior of features from the output encoder after being trained by the contrastive loss. The principal component analysis (PCA) features of the encoder output have a repetitive pattern; moreover, an experimental outcome revealed that they move in an elliptical shape and group the music beats at specific areas.

D. Model Description

The DRNN topology employed in our experiments is comprised of a CLDNN encoder and a deep recurrent decoder (see Table I). The CLDNN architecture follows a similar configuration as that considered in [15].

The input audio features are reduced by four convolutional layers, each followed by a batch normalization [21] and an exponential linear unit activation [22]. Then, three LSTM layers with 500 units each and a fully-connected layer with 65 dimensions complete the structure of the encoder.

The input to the decoder consists of the previous motion frame (71 dimensions) and the decoder output (65 dimensions) with a width of 136 dimensions. The decoder is also comprised of three LSTM layers with 500 units each and a fully-connected layer with 71 dimensions.

Regarding music-motion control, we add the contrastive cost function after calculating the next step and the MSE cost function.

IV. EXPERIMENTS

In this study, we conducted a series of experiments to 1) improve motion generation with weakly-supervised learning, and 2) examine the effect of music with different characteristics on the results.

A. Data

Because of a lack of available datasets that include dance synchronized to music, we prepared three datasets only. We restricted the data to small samples using different music genres with different rhythms.

Hip hop bounce: This dataset comprises two hip hop music tracks with a repetitive lateral bouncing step to match the rhythm. Each track is three minutes long on the average at 80 to 95 beats per second (bpm).

Salsa: This dataset comprises seven salsa tracks (four minutes long on average). All tracks include vocals and rhythms between 95 to 130 bpm. Furthermore, this dataset includes a lateral salsa dance step during instrumental moments and a front-back salsa step during vocal elements.

Mixed: This dataset comprises 13 music tracks with and without vocal elements (six genres: salsa, bachata, ballad, hip hop, rock, and bossa nova) with an average length of three minutes. Depending on the genre, each track includes up to two dance steps.

B. Audio Feature Extraction

We extracted the power features for each file that was sampled at 16 KHz as follows:

- First, we utilize white noise to contaminate the audio input at different signal-to-noise ratio (SNR). This allows to augment the number of training files and reduce possible overfitting. The SNR values for the training set were 0 and 10 dB.
- To synchronize the audio with the motion frame, we extracted a slice of 534 samples (33 milliseconds) of the corresponding position. This extracted slice was converted to H short-time Fourier transform (STFT) frames of 160 samples (10 milliseconds) with a shift of 80 samples (5 milliseconds).
- From the STFT frames, we used the power information, which was normalized between -0.9 and 0.9 on the W frequency bin axis.
- Finally, we stacked the H frames; thus, the input of the network was a $1 \times W \times H$ -dimensional file.

C. Motion Representation

For each audio track, we employed the manually selected rotations and root translation (Table II) captured by a single Kinect v2 device at a regular rate of 30 frames per second (see Fig 6). Then, the captured motion was post-processed and synchronized with the audio data using a motion beat

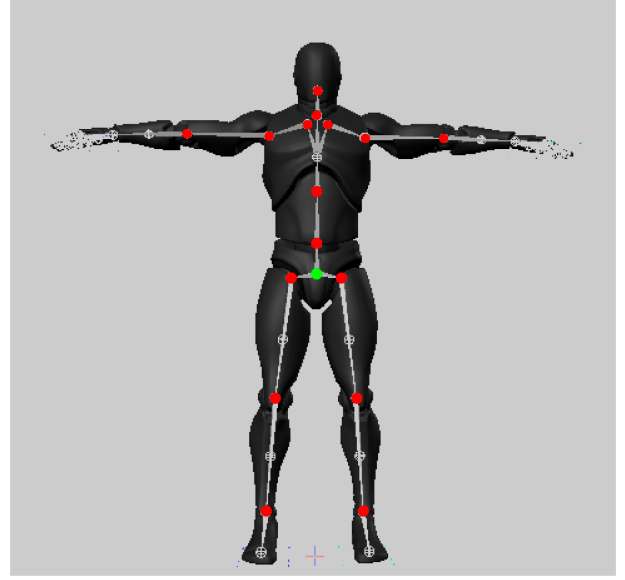


Fig. 6. Skeleton of the employed model.

TABLE II
MOTION FEATURES.

Joint	Type	Index
Root	Translation	y_t^0, y_t^1, y_t^2
(Root) Pelvis	Rotation	$y_t^3, y_t^4, y_t^5, y_t^6$
Head	Rotation	$y_t^7, y_t^8, y_t^9, y_t^{10}$
Neck	Rotation	$y_t^{11}, y_t^{12}, y_t^{13}, y_t^{14}$
Spine1	Rotation	$y_t^{15}, y_t^{16}, y_t^{17}, y_t^{18}$
Spine2	Rotation	$y_t^{19}, y_t^{20}, y_t^{21}, y_t^{22}$
Left Clavicle	Rotation	$y_t^{23}, y_t^{24}, y_t^{25}, y_t^{26}$
Left Shoulder	Rotation	$y_t^{27}, y_t^{28}, y_t^{29}, y_t^{30}$
Left Forearm	Rotation	$y_t^{31}, y_t^{32}, y_t^{33}, y_t^{34}$
Right Clavicle	Rotation	$y_t^{35}, y_t^{36}, y_t^{37}, y_t^{38}$
Right Shoulder	Rotation	$y_t^{39}, y_t^{40}, y_t^{41}, y_t^{42}$
Right Forearm	Rotation	$y_t^{43}, y_t^{44}, y_t^{45}, y_t^{46}$
Left Thigh	Rotation	$y_t^{47}, y_t^{48}, y_t^{49}, y_t^{50}$
Left Knee	Rotation	$y_t^{51}, y_t^{52}, y_t^{53}, y_t^{54}$
Left Foot	Rotation	$y_t^{55}, y_t^{56}, y_t^{57}, y_t^{58}$
Right Thigh	Rotation	$y_t^{59}, y_t^{60}, y_t^{61}, y_t^{62}$
Right Knee	Rotation	$y_t^{63}, y_t^{64}, y_t^{65}, y_t^{66}$
Right Foot	Rotation	$y_t^{67}, y_t^{68}, y_t^{69}, y_t^{70}$

algorithm introduced in [1]. The motion features are then post-processed from as follows:

- From the captured data in a hierarchical translation-rotation format, we processed the spatial information (i.e., translation) of the body as a three dimensional vector (x, y, z) in meters.
- Then, the joint rotation in degrees are converted into quaternions [23] $(q_t^x, q_t^y, q_t^z, q_t^w)$ using the tool Transforms3D¹ and concatenated them to the translation vector.
- To avoid the saturation of the activation functions, we normalized each vector component using the maximum

¹Python Library: <https://github.com/matthew-brett/transforms3d>

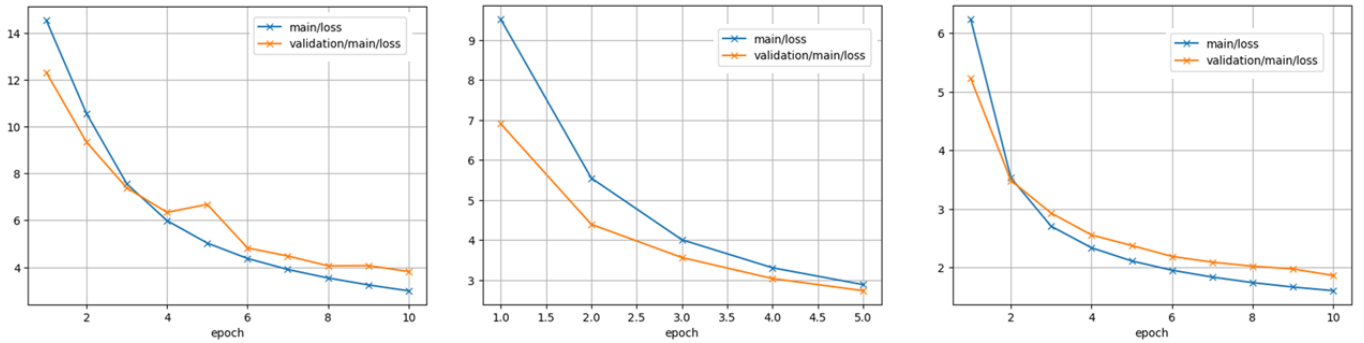


Fig. 7. Training loss. Left: Hip hop dataset, Center: Salsa dataset, Right: Mixed genres dataset.

value of each component in the range of -0.9 to 0.9 . The resultant vector (71 dimensions) was the target of the neural network.

Note that we did not apply any filtering or denoising method to maintain the noisy nature of the motion.

D. Training procedure

The models were trained for five to ten epochs using each dataset and the “chainer” framework [24] as an optimization tool. Each training epoch took an average of 60 minutes.

The models were trained using an NVIDIA GTX TITAN graphic processing unit. For the optimization, we employed the “Adam” solver [25] with a training mini-batch of 50 files and white noise added to the gradient. Each training batch employed sequences of 150 steps.

V. RESULTS

A. Training Loss

Figure 7 shows the training process of the models proposed in this study. The validation data is prepared by contaminating the audio with white noise at an SNR of 5 dB (a different SNR value from those employed for the training set). The models were observed to have fast convergence for all the datasets. However, the salsa dataset required five epochs only while the others required ten epochs each. We assume that the following two characteristics accounts for the short training of the salsa dataset: 1) The salsa dataset has more files than the hip hop bounce dataset which generate more iterations for a single epoch, and 2) the salsa dataset has fewer dance steps than the mixed dataset.

B. Metrics

In this study, a quantitative evaluation of the proposed models was performed using the f-score. Using the method employed in [1], we extract the motion beat² from each dance and obtain the f-score regarding the music beat. To demonstrate the benefits of adding a motion cost function, we compared the performance of the sequence-to-sequence (S2S) and music-motion-loss sequence-to-sequence (S2SMC)

²Video samples can be found online: <https://www.youtube.com/watch?v=VTy0yf-saDQ>

TABLE III
F-SCORE AND CROSS ENTROPY OF BOUNCE DATASET.

Method	F-Score							Entropy
	Hip Hop				Other genres			
	Clean*	White* noise	Claps	Crowd	Clean	Clean	White noise	
Madmom (Music beat)	89.18	-	-	-	-	80.13	-	-
Marsyas (Music beat)	54.11	-	-	-	-	48.89	-	-
Dancer (baseline)	62.31	-	-	-	-	54.49	-	-
S2S	55.98	46.33	50.58	54.11	32.95	35.84	34.86	1.98
S2S-MC	64.90	55.58	60.62	56.37	37.69	34.63	34.05	1.41

*Trained data for S2S and S2SMC

models with music beat retrieval frameworks. The models were evaluated under different conditions: clean trained music and white noise at a music input with a SNR of 20 dB, untrained noises (e.g., claps and crowd noise) at an SNR of 20 dB, clean untrained music tracks of the same genre and music tracks of different genres under clean conditions, and white noise at an SNR of 20 dB. Besides, we evaluated the cross entropy of the generated motion and the dancer for the hip hop dataset.

C. Music vs. Motion beat

We compared the results to two music beat retrieval systems; “madmom” and the music analysis, retrieval and synthesis for audio signals (MARSYAS) systems. “Madmom” [26] is a Python audio and music signal processing library that employs deep learning to process the music beat, while MARSYAS is an open-source framework that obtains a music beat using an agent-based tempo that processes a continuous audio signal [27]. The beat annotations of each system were compared to manually annotated ground truth beat data. Furthermore, we compared the ground truth beat data to the motion beat of the dancer for each dance. Here the average f-score was used as a baseline.

Table III compares the results of the music beat retrieval frameworks, the motion beat of the dancer (baseline), and the motion beat of the generated dances. It is evident from the evaluation results that the proposed models demonstrate better performance than MARSYAS, and S2SMC outperformed S2S

TABLE IV
F-SCORE OF SALSA DATASET.

Method	Clean*	White*	Claps	Crowd
Madmom	51.62	-	-	-
Marsyas	23.38	-	-	-
Dancer (baseline)	52.82	-	-	-
S2S	53.79	52.88	52.76	51.98
S2S-MC	53.96	53.09	53.61	52.48

*Trained data for S2S and S2SMC

TABLE V
F-SCORE OF MIXED GENRES.

Method	Bachata*	Ballad*	Bossa* Nova	Rock*	Hip Hop*	Salsa*
Dancer (baseline)	62.55	52.07	45.02	62.32	55.84	53.96
S2S	60.72	49.92	46.19	60.06	64.30	52.31
S2S-MC	56.63	48.48	40.91	64.87	63.85	53.71

*Trained data for S2S and S2SMC

in the evaluations that used clean and noisy data for training. It was also found that the addition of different noises to the music input does not critically affect the models when the music input is obtained from the training set. However, the performance is reduced to almost half when an untrained music track of the same or different genres was used as the input. Additionally, the proposed models did not outperform a model trained to only process music beat, i.e., MADMOM. We suppose that the music processing approach in our models accounts for this low performance. The proposed model requires three input frames of 33 milliseconds each to generate the motion of a given music beat, while the evaluation only allows a threshold of 70 milliseconds. Furthermore, the accuracy of the training model is reduced due to the delayed reaction of the baseline dancer to the music beat. S2SMC demonstrated lower cross entropy than S2S, signifying that the dance generated by S2SMC was similar to the trained dance.

Table IV shows the f-score for the music tracks trained using the salsa dataset. Both models show better performance than the dancer when tested under the same training conditions, and S2SMC shows better performance than S2S under all conditions. It is worth noting that the size of the dataset influences the performance of the models, besides we employed a larger dataset compared to the that in previous experiment.

Table V shows the results of the mixed genre dataset. As can be seen, different results were obtained for each model. The proposed methods do not outperform the baseline, whereas S2S outperformed S2SMC for most genres. The main reason for this difference in the results is the complexity of the dataset and the variety of dance steps relative to the number of music samples; thus, the model could not group the beat correctly.

VI. CONCLUSION

In this paper, we have presented an optimization technique for weakly-supervised deep recurrent neural networks for dance generation tasks. The proposed model was trained end-to-end and performed better than using only a mean squared cost function. We have demonstrated that the models can

generate a correlated motion pattern with a motion beat f-score similar to that of a dancer and lower cross entropy. Besides, the models could be used for real-time tasks because of the low forwarding time (approximately 12 ms). Furthermore, the models show low training time and can be trained from scratch.

The proposed model demonstrates reliable performance for motion generation, including music track inputs with a different type of noises. However, the motion pattern is affected by the diversity of the trained patterns and is constrained to the given dataset; this issue will be the focus of future research.

REFERENCES

- [1] C. Ho, W. T. Tsai, K. S. Lin, and H. H. Chen, "Extraction and alignment evaluation of motion beats for street dance," in *2013 IEEE International Conference on Acoustics, Speech and Signal Processing*, May 2013, pp. 2429–2433.
- [2] T. Kim, S. I. Park, and S. Y. Shin, "Rhythmic-motion synthesis based on motion-beat analysis," in *ACM SIGGRAPH 2003 Papers*, New York, NY, USA, 2003, SIGGRAPH '03, pp. 392–401, ACM.
- [3] J. L. Oliveira, G. Ince, K. Nakamura, K. Nakadai, H. G. Okuno, and *et al.*, "Beat tracking for interactive dancing robots," *International Journal of Humanoid Robotics*, vol. 12, no. 04, pp. 1550023, 2015.
- [4] J. J. Aucouturier, K. Ikeuchi, H. Hirukawa, S. Nakaoka, and T. Shiratori *et al.*, "Cheek to chip: Dancing robots and ai's future," *IEEE Intelligent Systems*, vol. 23, no. 2, pp. 74–84, March 2008.
- [5] S. Fukayama and M. Goto, "Music Content Driven Automated Choreography with Beat-wise Motion Connectivity Constraints," in *Proceedings of SMC*, 2015.
- [6] L. Crnkovic-Friis and L. Crnkovic-Friis, "Generative choreography using deep learning," *CoRR*, vol. abs/1605.06921, 2016.
- [7] O. Alemi and P. Pasquier, "Groovenet : Real-time music-driven dance movement generation using artificial neural networks," in *23rd ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2017.
- [8] T. Tang, J. Jia, and H. Mao, "Dance with melody: An lstm-autoencoder approach to music-oriented dance synthesis," in *2018 ACM Multimedia Conference on Multimedia Conference, MM 2018*, 2018, pp. 1598–1606.
- [9] Z. Li, Y. Zhou, S. Xiao, C. He, and H. Li, "Auto-conditioned LSTM network for extended complex human motion synthesis," *CoRR*, vol. abs/1707.05363, 2017.
- [10] F. Korzeniowski and G. Widmer, "End-to-end musical key estimation using a convolutional neural network," *CoRR*, vol. abs/1706.02921, 2017.
- [11] V. Kuleshov, S. Z. Enam, and S. Ermon, "Audio super resolution using neural networks," *CoRR*, vol. abs/1708.00853, 2017.
- [12] Y. Meng, J. Shen, C. Zhang, and J. Han, "Weakly-Supervised Hierarchical Text Classification," *arXiv e-prints*, p. arXiv:1812.11270, Dec. 2018.
- [13] Z. Gan, C. Li, R. Henaio, D. Carlson, and L. Carin, "Deep Temporal Sigmoid Belief Networks for Sequence Modeling," *Advances in Neural Information Processing Systems*, pp. 1–9, 2015.
- [14] W. T. Chu and S. Y. Tsai, "Rhythm of motion extraction and rhythm-based cross-media alignment for dance videos," *IEEE Transactions on Multimedia*, vol. 14, no. 1, pp. 129–141, Feb 2012.
- [15] T. N. Sainath, O. Vinyals, A. Senior, and H. Sak, "Convolutional, long short-term memory, fully connected deep neural networks," in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, April 2015, pp. 4580–4584.
- [16] S. Chopra, R. Hadsell, and Y. LeCun, "Learning a similarity metric discriminatively, with application to face verification," in *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, June 2005, vol. 1, pp. 539–546 vol. 1.
- [17] A. Kumar and B. Raj, "Audio event detection using weakly labeled data," in *Proceedings of the 24th ACM International Conference on Multimedia*, New York, NY, USA, 2016, MM '16, pp. 1038–1047, ACM.
- [18] M. I. Mandel and D. P. W. Ellis, "Multiple-instance learning for music information retrieval," in *ISMIR*, 2008.

- [19] F. Tian and X. Shen, "Image annotation with weak labels," in *Web-Age Information Management*, Jianyong Wang, Hui Xiong, Yoshiharu Ishikawa, Jianliang Xu, and Junfeng Zhou, Eds., Berlin, Heidelberg, 2013, pp. 375–380, Springer Berlin Heidelberg.
- [20] I. Sutskever, O. Vinyals, and Q. V. Le, "Sequence to sequence learning with neural networks," *CoRR*, vol. abs/1409.3215, 2014.
- [21] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *ICML*, 2015.
- [22] D.-A. Clevert, T. Unterthiner, and S. Hochreiter, "Fast and accurate deep network learning by exponential linear units (elus)," *CoRR*, vol. abs/1511.07289, 2015.
- [23] A. Shenitzer, B.A. Rosenfeld, and H. Grant, *A History of Non-Euclidean Geometry: Evolution of the Concept of a Geometric Space*, Studies in the History of Mathematics and Physical Sciences. Springer New York, 2012.
- [24] S. Tokui, K. Oono, S. Hido, and J. Clayton, "Chainer: a next-generation open source framework for deep learning," in *Proceedings of Workshop on Machine Learning Systems (LearningSys) in The Twenty-ninth Annual Conference on Neural Information Processing Systems (NIPS)*, 2015.
- [25] D. Kingma and J. Ba, "Adam: A Method for Stochastic Optimization," *arXiv:1412.6980 [cs]*, pp. 1–15, 2014.
- [26] S. Böck, F. Korzeniowski, J. Schlüter, F. Krebs, and G. Widmer, "madmom: a new Python Audio and Music Signal Processing Library," in *Proceedings of the 24th ACM International Conference on Multimedia*, Amsterdam, The Netherlands, 10 2016, pp. 1174–1178.
- [27] J. L. Oliveira, "Ibt: A real-time tempo and beat tracking system," in *Proceedings of International Society for Music Information Retrieval Conference (ISMIR)*, 2010.