

1 We thank the reviewers for the constructive comments and suggestions. Below we address the main
2 issues. We will fix all other minor issues in the revised paper.

3 **R1, R3: Dance unit.** We will elaborate on the dance units, which are not a simple division of
4 a dance sequence. We first extract kinematic beats and then align the extracted beat times to the
5 pre-defined time steps (i.e., 8, 16, 24 and 32 in a dance unit of length 32). Finally, we interpolate and
6 extrapolate the original pose sequences to obtain the temporally normalized dance units, which are
7 further encoded into the disentangled initial pose space and movement space.

8 **R2: Title.** Thanks for the suggestion. We will change the title in the revision.

9 **R2: Term usage.** (1) **Abstract:** We agree with the comment and will remove the terms “top-down”
10 and “bottom-up” in the abstract. (2) **Multimodality:** We defined the term “multimodality” in L272
11 and in the caption of Table 1. However, we agree that the usage of this term might be ambiguous. We
12 will replace it with more explicit expressions (e.g., “overall diversity” and “instance diversity”). (3)
13 **Style-consistency:** We will replace it with a more explicit term such as “music-dance consistency” or
14 “music-dance matching”.

15 **R2: Unclear observations.** The observations in L21 refer to “dance to music is a creative process
16 that is both innate and acquired”. Similarly, we leverage our prior knowledge (innate) to design the
17 framework and then use a data-driven learning process, which corresponds to the acquired properties.

18 **R2: What is assumed on the model.** To facilitate the training process, we assume that the latent
19 factors of the music-to-dance generation process can be disentangled to two components: beat and
20 style. In the decomposition, we learn how to dance according to the beat. In the composition, we
21 learn how to dance to the beat and style at the same time.

22 **R2: More details.** Our code and trained model will be released, where all implementation details
23 can be found. (1) **Number of poses:** We use 32 frames for each dance unit in our experiments (i.e.,
24 two seconds at 16 fps). (2) **Music style classifier:** We have evaluated several methods to extract
25 features: features from a pre-trained network SoundNet, from a music autoencoder, and from a
26 music classifier. We adopt the music classifier as our music style extractor due to the capability to
27 better separate different types of music. This classifier consists of multiple fully-connected layers
28 and takes as inputs the MFCC features extracted from music of various types. The training data are
29 the audios from videos used for training our dance model, totaling 360K clips and three categories.
30 (3) **Baselines:** LSTM baseline is basically the same as [28] (see Figure 5 in [28]). It takes a sequence
31 of extracted audio features as inputs and directly predicts a sequences of poses. Aud-MoCoGAN is
32 an extension of MoCoGAN [31]. The original MoCoGAN takes a sequence of random variables as
33 inputs. We use audio features (MFCC) and beats (a sequence of binary variables) as additional inputs.
34 (4) **Action classifier:** We utilize an RNN to take as inputs the pose sequences of arbitrary length, and
35 append multiple fully-connected layers on the last hidden state to perform classification. We use the
36 features from the last fully-connected layer to compute the FID. (5) **Subjective test:** We conduct the
37 subjective test online. The professional backgrounds (high-school students, doctors, professors) and
38 ages (from 15 to 61) of the subjects are diverse.

39 **R2: Quantitative and qualitative results.** We will reorganize the section order according to the
40 suggestion. Since the trained model will be released, it is easy to qualitatively evaluate the results.
41 For the quantitative results, only motion realism and style consistency are from the subjective test.
42 Other evaluation metrics including FID, beat hit-rate, beat coverage, diversity, and multimodality are
43 all measured quantitatively.

44 **R2, R3: More references.** The LSTM-based methods are similar to our first baseline and suffer
45 from diversity and multimodality. We will cite and discuss these papers in the revised manuscript.

46 **R3: Long sequence generation.** In our framework, we can generate up to five consecutive dance
47 units (about 10 seconds) in each step, and a sequence of arbitrary length can be seamlessly generated
48 in an iterative manner. In this work, we use the last pose of the previous step as the initial pose
49 of the next step, as shown in the bottom right of Figure 1 and Eq. (9). We will design methods to
50 better capture the long-term relationship among sub-sequences (e.g., through a hierarchical recurrent
51 network) in our future work.

52 **R3: Complicated dance.** In term of complicated poses, the hip-hop dance in our selected dance
53 styles is quite complicated. We will continue to collect and incorporate more dance styles including
54 pop dance and partner dance.