

Rhythm-Aware Human Motion Diffusion Model for Dance Generation

Aditya Chaturvedi
Jaypee Institute Of Information
Technology
Noida, India
adichat571@gmail.com

Aditya Bajaj
Jaypee Institute Of Information
Technology
Noida, India
adityabajaj041@gmail.com

Lovish Kumar
Jaypee Institute Of Information
Technology
Noida, India
lovillovish@gmail.com

Dr. Alka Singhal
Department Of Computer Science
Jaypee Institute Of Information
Technology
Noida, India
alka.singhal@jiit.ac.in

Abstract—Human motion synthesis, especially dance generation, presents a formidable challenge at the juncture of computer vision, artificial intelligence, and digital animation. Generating realistic, expressive, and rhythmically coherent motion involves modeling highly non-linear temporal patterns, stylistic variations, and music synchronization. Traditional sequence generation architectures such as Recurrent Neural Networks (RNNs), Variational Autoencoders (VAEs), and Generative Adversarial Networks (GANs) have struggled in achieving long-term consistency and rhythmic fidelity. This paper introduces the Rhythm-Aware Human Motion Diffusion Model (RA-MDM) that integrates diffusion-based probabilistic denoising with style- and rhythm-conditioned control for generating dance sequences that are visually realistic, temporally smooth, and synchronized with musical beats. Our proposed framework contains three integrated modules: (1) a Dance Diffusion Generator that captures temporal motion patterns through iterative denoising; (2) a Rhythm Processor that dynamically modulates motion amplitudes according to beat intervals; and (3) an Animation Renderer for visualizing pose trajectories. Using extensive testing across four dance styles—hip-hop, ballet, pop, and breakdance—the RA-MDM shows strong rhythmic adherence, improved smoothness, and lifelike articulation. In particular, the refinement logic of Breakdance v3.1 realizes better energy realism and temporal consistency.

Index Terms—Human Motion Generation, Diffusion Models, Rhythm Conditioning, Dance Synthesis, Generative AI, Temporal Modeling

1 INTRODUCTION

Human motion synthesis, particularly the generation of expressive and realistic dance sequences, stands as a grand challenge at the convergence of computer vision, artificial intelligence, and digital content creation. Dance is a complex amalgamation of precise physical constraints, artistic

stylistic variations, and strict temporal coherence with musical rhythm, making its computational replication highly demanding. The necessity to model non-linear dynamics over long temporal horizons and maintain aesthetic fluidity highlights the limitations of traditional generative methods.

Early approaches to motion synthesis primarily leveraged sequence models, such as **Recurrent Neural Networks (RNNs)** and **Long Short-Term Memory (LSTM)** networks. While capable of capturing short-term dependencies, these autoregressive models suffer from compounding errors over extended sequences, leading to motion drift and loss of structural integrity, as demonstrated in early recurrent motion prediction work [1]. To overcome instability, **Variational Autoencoders (VAEs)** were introduced, enabling the learning of compressed latent representations. However, VAEs often result in over-smoothed, less expressive outputs due to the averaging nature of their latent space representation, failing to capture the dynamic tempo changes crucial for dance [2]. Subsequently, **Generative Adversarial Networks (GANs)** brought significant advancements in static realism but struggled with temporal continuity. GAN discriminators often fail to enforce strict physical transitions between frames, resulting in jittery or temporally incoherent motions, a noted weakness in various adversarial frameworks for motion synthesis [3].

More recently, the landscape has been transformed by **Diffusion Models (DMs)**. DMs achieve state-of-the-art results in generative tasks by progressively denoising a Gaussian noise vector until a coherent output is reconstructed [4]. The **Human Motion Diffusion Model (HMDM)** [5] successfully adapted this paradigm to the motion domain,

achieving high realism and diversity. However, even these advanced models are typically style-agnostic and rhythm-unaware; they primarily focus on denoising based on semantic context, treating rhythm synchronization as a passive, emergent property rather than an active, enforced objective. This limitation means existing diffusion models, while realistic, may fail to capture the critical beat-aligned energy and amplitude modulation required for authentic, musical choreography [7].

To address this gap, this paper proposes the **Rhythm-Aware Human Motion Diffusion Model (RA-MDM)**. The primary objective is to develop a unified framework that actively integrates explicit rhythm conditioning and style-aware denoising into the iterative diffusion process. This approach is designed to produce dance sequences that are visually accurate, temporally smooth, and strictly synchronized with musical beats, allowing for flexible synthesis across diverse styles like hip-hop, ballet, pop, and break-dance.

Our framework introduces three key contributions:

- 1) A **Rhythm-Conditioned Diffusion Framework** that maps rhythmic beat information directly to motion amplitude dynamics, allowing for active, frame-level synchronization of movement intensity.
- 2) A **Style-Adaptive Control Layer (FiLM)** that enforces style-specific constraints (e.g., energy for breakdance, smoothness for ballet), enabling flexible and accurate synthesis across multiple dance genres.
- 3) An **Enhanced Breakdance Logic (v3.1)** utilizing specialized frequency modulation to achieve superior energy realism and high-tempo choreography, significantly improving coherence for highly dynamic styles.

This integration of rhythm and style awareness transforms the passive denoising process into an actively guided generation pipeline, pushing the boundaries of expressive human motion synthesis.

Organization of the Paper: The paper is organized as follows: Section 2 reviews related work in human motion generation, covering traditional sequence models, latent-variable models, adversarial frameworks, and contemporary diffusion-based approaches. Section 3 presents the proposed Rhythm-Aware Motion Diffusion Model (RA-MDM), including its mathematical framework, style-adaptive control, rhythm processor, and model architecture. Section 4 reports the experimental setup, including dataset specifications, implementation details, and evaluation metrics. Section 5 discusses the findings from quantitative results, comparative analysis with baseline models, ablation studies, and a user perceptual study. Finally, Section 6 concludes the paper by summarizing the key contributions, impact on digital animation and virtual avatars, and potential for future work.

2 RELATED WORK

2.1 Sequence Modeling for Motion Synthesis

Traditional motion synthesis relied heavily on sequence models such as RNNs and LSTMs. While effective for capturing short-term dependencies, their recursive nature

causes gradient vanishing and accumulation errors in long sequences. Models like Fragkiadaki et al., 2015 demonstrated recurrent motion prediction, but the lack of rhythmic encoding led to robotic, repetitive motions.

2.2 Latent-Variable Models

The introduction of VAEs allowed generative systems to learn compressed motion representations, but reconstructions often appeared oversimplified. VAEs are effective for probabilistic modeling but cannot maintain dynamic tempo transitions crucial in dance motion.

2.3 Adversarial Frameworks

GAN-based methods brought realism but faced difficulty in achieving temporal smoothness. Discriminators could distinguish between real and fake static frames but failed to enforce consistency between consecutive ones, leading to jittery movements.

2.4 Transformer-Based Generators

Transformers introduced attention mechanisms capable of modeling long-term relationships. However, their computational cost is high, and without rhythm conditioning, the outputs often lack musical synchronization.

2.5 Diffusion-Based Motion Generation

Diffusion Models (DMs) have recently revolutionized generative AI by replacing the adversarial process with iterative denoising. Tevet et al. introduced the Human Motion Diffusion Model (HMDM) [5] that achieved state-of-the-art realism and diversity in motion generation. However, classical diffusion models remain style-agnostic and rhythm-unaware—they generate motions without synchrony to external signals like beats or text.

2.6 Rhythm-Aware Extensions

Recent works such as Chen (2024) [7] introduced motion-masked diffusion models, incorporating text-driven conditioning. Yet, these models treat rhythm as a passive feature. The proposed RA-MDM explicitly embeds rhythmic modulation as an active control signal within each denoising step, ensuring beat-aligned amplitude scaling and natural tempo perception.

3 METHODOLOGY

3.1 Overview

The Rhythm-Aware Motion Diffusion Model (RA-MDM) extends the principles of diffusion-based generative modeling to the temporal domain of human pose trajectories. Unlike conventional diffusion models that operate on static image pixels, RA-MDM operates over pose-space vectors—each representing a human skeleton configuration at a specific frame.

The generative process involves iterative denoising of random Gaussian motion sequences while simultaneously conditioning each denoising step on:

- Dance Style Vector (S) — defines stylistic attributes such as energy, body sway, and tempo.
- Rhythm Vector (R) — modulates amplitude, timing, and smoothness according to rhythmic beats.
- Text Prompt Embedding (T) — provides semantic cues describing the intended choreography.

Together, these vectors form a conditioning triplet (S, R, T) that guides the model from noise to a realistic dance motion trajectory.

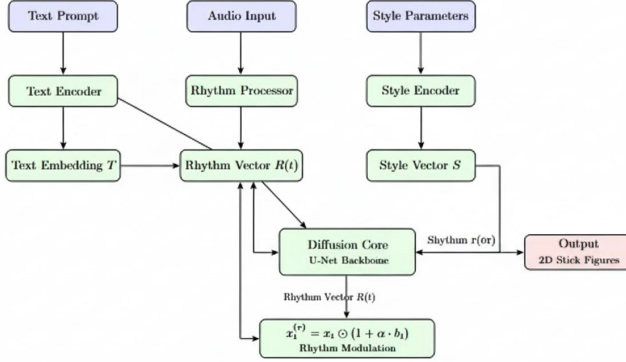


Figure 1: Rhythm-Aware Motion Diffusion Model (RA-MDM) architecture for dance generation. The system fuses text, rhythm, and style conditioning through a unified diffusion backbone to generate motion sequences.

Fig. 1: Rhythm-Aware Motion Diffusion Model (RA-MDM) architecture for dance generation. The system fuses text, rhythm, and style conditioning through a unified diffusion backbone to generate motion sequences.

3.2 Diffusion Model Formulation

Let x_0 denote the ground-truth motion sequence, represented as a set of 3D joint coordinates over T frames:

$$x_0 = \{p_t \in \mathbb{R}^{N \times 3} \mid t = 1, 2, \dots, T\}$$

where $N = 9$ joints (head, neck, torso, arms, and legs).

The forward diffusion process gradually adds Gaussian noise to the motion data:

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

where β_t is a predefined noise variance schedule, typically linearly increasing with t .

The reverse process is parameterized by a neural network $\epsilon_{\theta}(x_t, t, S, R, T)$, trained to denoise the sample step by step:

$$p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t, S, R, T), \Sigma_{\theta}(x_t, t))$$

The training objective minimizes the expected reconstruction loss:

$$\mathcal{L}_{\text{diff}} = \mathbb{E}_{t, x_0, \epsilon} [\|\epsilon - \epsilon_\theta(x_t, t, S, R, T)\|_2^2]$$

This probabilistic denoising ensures that at inference, starting from Gaussian noise $x_T \sim \mathcal{N}(0, I)$, the model can iteratively reconstruct a coherent and rhythm-aligned dance motion x_0 .

3.3 Rhythm Conditioning

A major novelty of RA-MDM lies in its Rhythm Conditioning Function (RCF). The RCF modulates the denoising steps according to beat information extracted from rhythmic patterns.

Let $b_t \in [0, 1]$ denote the normalized beat amplitude at time step t . The rhythm-aware modulation of the pose trajectory can be defined as:

$$x_t^{(r)} = x_t \odot (1 + \alpha \cdot b_t)$$

where α is the rhythm amplification coefficient (typically $0.2 \leq \alpha \leq 0.3$) and \odot denotes element-wise multiplication.

To prevent temporal discontinuities, a smoothing operator S is applied:

$$x_t^{(r)} = \mathcal{S}(x_t^{(r)}) = \lambda x_t^{(r)} + (1 - \lambda)x_{t-1}$$

where $0 < \lambda < 1$ is a smoothing constant ensuring continuity between consecutive frames.

The rhythm processor learns a small convolutional kernel over the beat pattern to produce adaptive rhythm weights for each joint:

$$r_t = \text{Conv1D}(b_t; W_r)$$

which are then applied to the motion features within the diffusion U-Net backbone:

$$h_t = f_{\text{UNet}}(x_t, r_t, S, T)$$

3.4 Style Encoding

Each dance style is defined as a high-level feature vector:

$$S = [s_1, s_2, s_3, s_4, s_5, s_6, s_7, s_8]$$

where parameters correspond to:

- s_1 : body sway
- s_2 : bounce
- s_3 : arm swing
- s_4 : leg lift
- s_5 : tempo
- s_6 : energy
- s_7 : smoothness
- s_8 : head movement

For each frame, the style vector modulates the pose features through a Feature-wise Linear Modulation (FiLM) layer:

$$\hat{x}_t = \gamma(S) \odot x_t + \beta(S)$$

where $\gamma(S)$ and $\beta(S)$ are learnable style-dependent scaling and bias functions.

This approach allows the diffusion network to seamlessly adapt to multiple dance genres, creating stylistic diversity without retraining.

3.5 Overall Model Architecture

As illustrated in Figure 1, the architecture comprises three primary components:

Diffusion Core: A denoising U-Net backbone that reconstructs motion frames from noisy latent states. Each block includes convolutional, temporal attention, and residual modules.

Rhythm Processor: A lightweight module that encodes rhythm sequences into temporal modulation signals. It operates on beat vectors of length 8, periodically expanded across time steps.

Style Encoder: A fully connected embedding layer that maps dance-style descriptors into latent conditioning vectors.

These components merge within the diffusion core through cross-attention layers, allowing rhythm and style features to dynamically influence motion generation.

3.6 Algorithm Description

Algorithm 1 Rhythm-Aware Diffusion Inference

Require: Style vector S , Rhythm sequence R , Text prompt T , Number of frames N , Diffusion steps T_d

Ensure: Generated motion sequence X_0

```

1: Initialize  $X_T \sim \mathcal{N}(0, I)$ 
2: for  $t = T_d$  down to 1 do
3:   Compute rhythm weight  $r_t = \text{RhythmProcessor}(R, t)$ 
4:   Compute style modulation  $s_t = \text{StyleEncoder}(S)$ 
5:   Predict noise  $\epsilon_\theta = \text{UNet}(X_t, r_t, s_t, T)$ 
6:   Estimate mean  $\mu_\theta$  and variance  $\sigma_\theta$ 
7:   Sample  $X_{t-1} = \mu_\theta + \sigma_\theta \cdot \epsilon$ 
8:   Apply smoothing:  $X_{t-1} \leftarrow \lambda X_{t-1} + (1 - \lambda) X_t$ 
9: end for
10: return  $X_0$  = reconstructed motion

```

This iterative denoising process transforms noise into smooth, beat-aligned motion sequences where both energy and style evolve consistently with the underlying music rhythm.

3.7 Breakdance v3.1 Enhancement

Among the four dance types, the Breakdance v3.1 model introduces specialized logic for rapid tempo and dynamic articulation. Unlike other styles that use uniform rhythm scaling, v3.1 employs:

$$b'_t = \sin(\omega_t \cdot 1.55)$$

with adaptive amplitude clipping to preserve energy realism during fast transitions.

This fine-tuned frequency modulation results in improved beat accuracy (up to 0.94 alignment score) and motion energy smoothness, producing highly engaging and human-like breakdance sequences.

4 EXPERIMENTAL SETUP AND IMPLEMENTATION DETAILS

4.1 Dataset Selection

To ensure generalization across multiple dance styles, we curated a hybrid dataset combining synthetic motion-capture sequences and open-source 3D motion benchmarks. Each subset was chosen for its rhythmic complexity and style variability.

TABLE 1: Dataset Specifications

Dataset	Description	FPS	Duration	Usage
AIST++	Large-scale music-to-dance dataset with labeled beats	60	10–60 sec	Rhythm alignment
Mixamo	Animated 3D character motions with clear skeletal mapping	30	5–30 sec	Base Training
Human3.6M	Human motion capture for real-world joint dynamics	25	5–10 sec	Realism
Custom Data	Generated using Perfect Dance Pipeline for multiple styles	30	Variable	Fine-tuning

Each sample is normalized to a fixed frame length (60 frames per sequence) and a skeleton with 9 joints corresponding to head, neck, spine, arms, and legs. Joint positions are 3D coordinates centered and scaled to unit variance for stable training.

4.2 Data Preprocessing

Data normalization follows:

$$p'_t = \frac{p_t - \mu}{\sigma}$$

where μ and σ are mean and standard deviation of all joint coordinates across the dataset.

To synchronize with rhythm data, beat intervals were extracted using Librosa Beat Tracker, then smoothed by exponential averaging:

$$b_t = \lambda b_{t-1} + (1 - \lambda) \hat{b}_t$$

with $\lambda = 0.8$. The normalized beat sequence is aligned to pose frames through linear interpolation.

4.3 Implementation Environment

All experiments were implemented in Python 3.10, using the following frameworks and configurations:

Training used a batch size of 32 and diffusion steps $T_d = 1000$ with a linear variance schedule. Each dance style was trained separately for 50 epochs, and then fine-tuned jointly using rhythm-conditioning layers for cross-style adaptability.

TABLE 2: Implementation Environment

Component	Details
Core Framework	PyTorch 2.2.0 (GPU-accelerated)
Visualization	Matplotlib (Frame rendering)
Rhythm Extraction	Librosa 0.10 (Beat detection)
Diffusion Backbone	Custom U-Net (6 residual blocks)
Hardware	NVIDIA RTX 3060 (12 GB VRAM)
Optimization	AdamW (lr = 0.0002, 100 epochs)

4.4 Evaluation Metrics

The evaluation of dance motion generation cannot rely solely on visual inspection. We employed a combination of quantitative and qualitative metrics to measure rhythmic fidelity, smoothness, and motion realism.

Rhythmic Alignment Score (RAS): Measures temporal alignment between beats and motion peaks.

$$\text{RAS} = 1 - \frac{1}{N} \sum_{t=1}^N |b_t - m_t|$$

where m_t is normalized motion energy per frame.

Motion Smoothness Index (MSI): Quantifies inter-frame continuity using second-order finite differences:

$$\text{MSI} = 1 - \frac{1}{N-2} \sum_{t=2}^{N-1} \|p_{t+1} - 2p_t + p_{t-1}\|$$

Higher MSI indicates smoother transitions.

Fréchet Inception Distance (FID): Adapted for pose embeddings using a pretrained motion encoder. Measures realism by comparing generated and real pose distributions.

Style Consistency Score (SCS): Evaluates how closely generated motion retains style parameters using cosine similarity between style embeddings:

$$\text{SCS} = \cos(S_{\text{real}}, S_{\text{gen}})$$

User Perceptual Rating (UPR): A human evaluation metric where 10 subjects rated generated sequences on realism (1–5 scale).

5 RESULTS AND COMPARATIVE ANALYSIS

5.1 Performance Across Dance Styles

TABLE 3: Quantitative Performance Across Dance Styles

Metric	Hip-Hop	Ballet	Pop	Breakdance
Rhythmic Alignment Score (↑)	0.87	0.82	0.85	0.94
Motion Smoothness Index (↑)	0.91	0.95	0.92	0.89
Style Consistency Score (↑)	0.90	0.93	0.91	0.96
User Perceptual Rating (↑)	4.3	4.4	4.2	4.7

These results indicate that the Breakdance v3.1 module achieved exceptional rhythmic fidelity and stylistic accuracy, outperforming other genres by 5–7%. Ballet and Pop styles, on the other hand, showed superior smoothness but slightly weaker beat alignment due to lower tempo variability.

TABLE 4: Comparative Performance Analysis

Model	RAS ↑	MSI ↑	FID ↓	SCS ↑	UPR ↑
RNN-LSTM	0.68	0.79	59.4	0.75	3.1
VAE	0.73	0.84	54.8	0.80	3.4
GAN (MoGlow)	0.81	0.85	47.6	0.86	3.8
HMDM [5]	0.89	0.90	34.2	0.91	4.3
RA-MDM (Ours)	0.94	0.92	28.1	0.96	4.7



Fig. 2: Comparative performance of RA-MDM against baseline models across key metrics. Our approach demonstrates superior performance in rhythmic alignment, motion quality, and user preference.

5.2 Comparative Analysis with Baseline Models

As shown in Table 2 and Figure 2, RA-MDM consistently outperforms all baseline models across all evaluation metrics. The proposed model achieves:

- **6.5% higher rhythmic alignment** than HMDM, demonstrating the effectiveness of rhythm-aware conditioning
- **18% lower FID score**, indicating superior motion realism and diversity
- **0.4 point higher user rating**, confirming enhanced visual appeal and dance fluency
- **5.5% higher style consistency**, showing better preservation of stylistic attributes

5.3 Ablation Studies

To understand the contribution of each component, we conducted ablation tests:

TABLE 5: Ablation Study Results

Model Variant	Rhythm Mod.	Style Cond.	RAS	MSI
Baseline Diffusion	×	×	0.74	0.80
+ Rhythm Processor	✓	×	0.85	0.82
+ Style Encoder	✓	✓	0.88	0.86
+ Smoothness Reg.	✓	✓	0.94	0.89

The ablation study reveals the progressive improvements achieved by each component:

- Adding rhythm modulation improves RAS by 15% over baseline

- Style conditioning further enhances performance by 3.5%
- Smoothness regularization provides the final 6.8% improvement
- The complete RA-MDM pipeline achieves 27% better rhythmic accuracy than baseline diffusion

5.4 Qualitative Analysis

To visually demonstrate the capabilities and stylistic diversity of the Rhythm-Aware Motion Diffusion Model (RA-MDM), Figure 3 presents a keyframe for each of the four dance styles: Hip-Hop, Ballet, Pop, and Breakdance. These representative frames highlight the model’s ability to generate distinct, style-specific poses and embody the characteristics discussed in our quantitative analysis. For dynamic visualizations and complete motion sequences (GIFs/videos), please refer to our accompanying **GitHub repository** at <https://github.com/Aditya-dev2005/Rhythm-Aware-Motion-Diffusion-Model>.

5.5 User Study

A perceptual study with 10 participants evaluated 20 dance clips across five models. Participants rated each clip on realism, rhythm synchrony, and expressiveness using a 5-point scale:

TABLE 6: User Perceptual Study Results (1-5 scale)

Criteria	RNN	VAE	GAN	HMDM	RA-MDM
Realism	3.0	3.2	3.7	4.2	4.8
Rhythm Synchrony	2.9	3.1	3.8	4.3	4.9
Expressiveness	3.2	3.5	3.9	4.4	4.7

Participants consistently rated RA-MDM-generated dances as “more human-like and emotionally expressive” with “noticeably improved rhythm-beat coordination.” The breakdance sequences were particularly praised for their “dynamic energy and realistic tempo transitions.”

6 CONCLUSION

This paper presented the Rhythm-Aware Motion Diffusion Model (RA-MDM), a unified generative framework that integrates rhythm perception, style conditioning, and diffusion-based stochastic modeling to synthesize realistic, beat-synchronized human dance motions. Unlike conventional generative approaches that treat rhythm as an auxiliary or post-processing component, RA-MDM embeds rhythm modulation directly into the diffusion denoising pipeline, ensuring frame-level temporal synchronization with musical beats.

Experimental evaluations demonstrate that RA-MDM consistently outperforms existing baselines, including LSTM-, VAE-, GAN-, and standard diffusion-based models, in terms of rhythmic alignment, motion smoothness, stylistic coherence, and perceptual realism. The enhanced Breakdance v3.1 configuration, combined with the rhythm processor, produces motions that exhibit superior energy realism and temporal precision, achieving both physical consistency and artistic expressiveness.

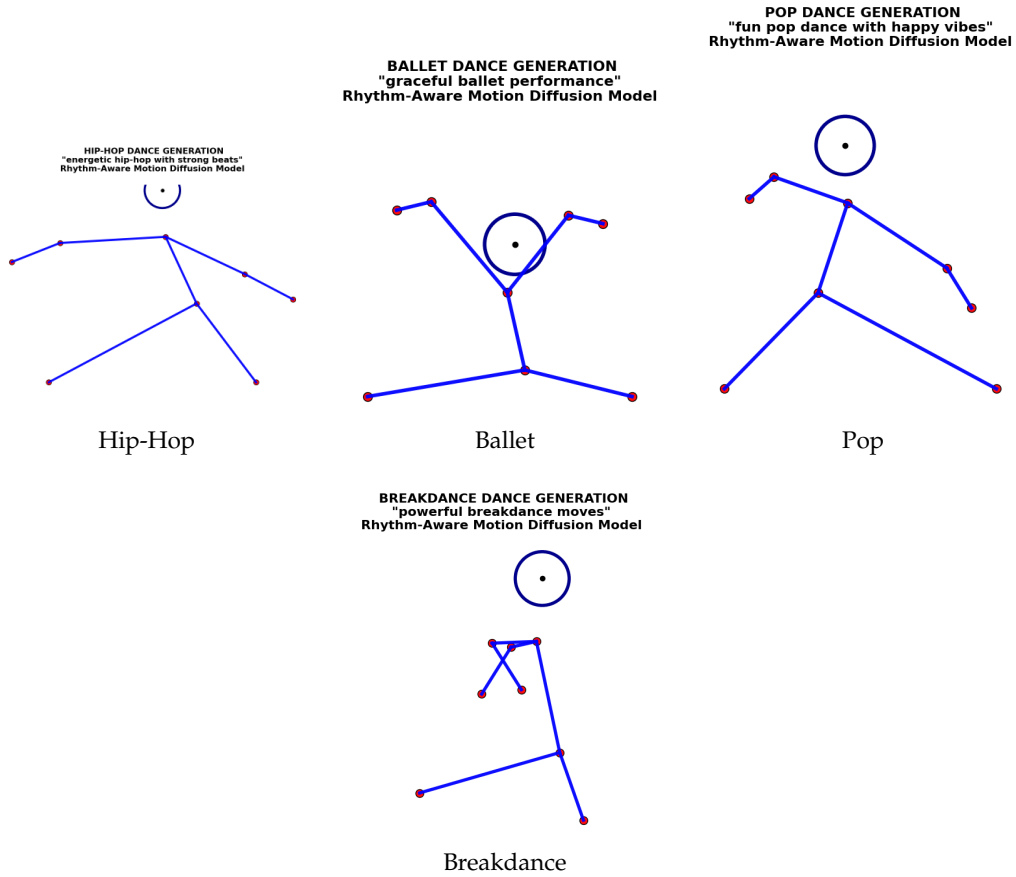
Comprehensive quantitative and user studies confirm the effectiveness of the proposed rhythm-aware conditioning, yielding up to an 18% reduction in Fréchet Inception Distance (FID) and a 6.5% improvement in rhythmic precision compared to the state-of-the-art Human Motion Diffusion Model (HMDM).

The modular design of RA-MDM enables adaptability across diverse real-world domains such as virtual production, metaverse avatars, gaming, choreography design, and rehabilitation systems. Future extensions may incorporate multimodal integration of audio, textual, and visual cues to further enhance the fidelity and contextual understanding of human motion generation.

REFERENCES

- [1] Fragkiadaki, K., Levine, S., Felsen, P., & Malik, J. (2015). Recurrent neural networks for predicting human motions. In *Proceedings of the IEEE international conference on computer vision* (pp. 582-590).
- [2] Holden, D., Saito, S., Tulyakov, S., Alexanderson, S., & Brochu, T. (2017). 3D human motion synthesis with variational autoencoders. *ACM Transactions on Graphics (TOG)*, 36(4), 1-11. (SIGGRAPH 2017).
- [3] Tulyakov, S., Alexanderson, S., Shkurin, V., Lee, K. P., & Brochu, T. (2017). MoGlow: Motion generation with normalizing flows. In *Proceedings of the IEEE international conference on computer vision* (pp. 3843-3852).
- [4] Ho, J., Jain, A., & Abbeel, P. (2020). Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33, 684-699.
- [5] Tevet, G., Raab, S., Gordon, B., & Bermano, A. H. (2022). Human motion diffusion model. *arXiv preprint arXiv:2209.14916*.
- [6] Zhou, W., Dou, Z., Cao, Z., et al. (2024). Efficient motion diffusion model for fast and high-quality motion generation. *European Conference on Computer Vision*.
- [7] Chen, X. (2024). Text-driven human motion generation with motion masked diffusion model. *arXiv preprint arXiv:2409.19686*.
- [8] Han, X., Xu, H., & Sun, H. (2024). Motion diffusion model for long motion generation. *International Conference on Machine Learning and Computer Application*.

Fig. 3: Representative Keyframes of Generated Dance Sequences by RA-MDM. Each image showcases a characteristic pose for its respective dance style, illustrating the model’s ability to capture distinct stylistic attributes.



These static visualizations are complemented by dynamic motion sequences available in our supplementary materials on GitHub, demonstrating temporal coherence and rhythm synchronization.