

MusicInfuser: Making Video Diffusion Listen and Dance

Susung Hong

Ira Kemelmacher-Shlizerman

Brian Curless

Steven M. Seitz

University of Washington

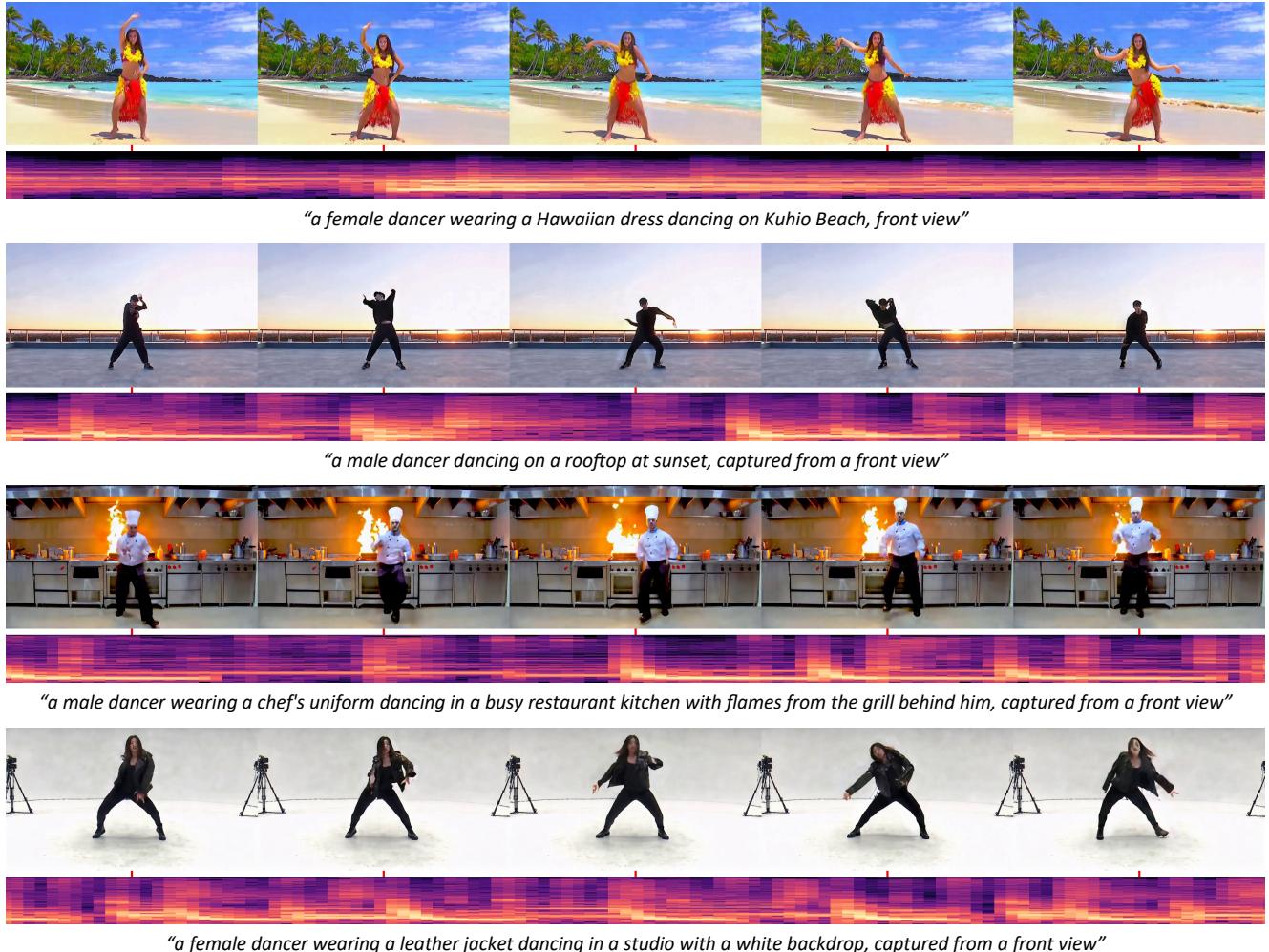


Figure 1. **MusicInfuser** adapts video diffusion models to music, making them listen and dance according to the music. This adaptation is done in a prior-preserving manner, enabling it to also accept style through the prompt while aligning the movement to the music. Please refer to the project page videos, as the movement appears slower due to the frame sampling rate.

Abstract

We introduce **MusicInfuser**, an approach for generating high-quality dance videos that are synchronized to a specified music track. Rather than attempting to design and train a new multimodal audio-video model, we show how

existing video diffusion models can be adapted to align with musical inputs by introducing lightweight music-video cross-attention and a low-rank adapter. Unlike prior work requiring motion capture data, our approach fine-tunes only on dance videos. **MusicInfuser** achieves high-quality music-driven video generation while preserving the flexi-

bility and generative capabilities of the underlying models. We introduce an evaluation framework using Video-LLMs to assess multiple dimensions of dance generation quality. The project page and code are available at <https://susunghong.github.io/MusicInfuser>.

1. Introduction

Today’s leading AI video generation tools (e.g., Sora, Gen, Veo) produce *silent* videos. While it’s possible to add music after the fact, it is difficult to generate motion that’s properly synchronized with a specified music track. Some research has begun to explore joint audio-video generation, e.g., [39], but this area is in its infancy and requires training larger and more complex joint audio-video models.

In this paper, we introduce an approach to adapt *pre-trained* text-to-video models to condition on music tracks. Our method, called *MusicInfuser*, generates output videos that are synchronized with the input music, with styles and appearances controllable via text prompts. We focus on the specific application of *dance*, i.e., generating realistic dancing figures that adjust and synchronize to the music.

The automatic generation of dance movements from music presents significant challenges due to the need to simultaneously consider multiple aspects such as style, beat, and the inherently multimodal nature of dance, where multiple valid movement sequences can follow from any given pose [30]. Principles of choreography [48] have inspired computational approaches to dance generation, leading researchers to explore methods ranging from graph-based approaches [12, 28] to modern deep neural networks [45, 50, 51] to synthesize dance movements. However, traditional methods for dance synthesis have relied on motion capture data [2], which is resource-intensive, or reconstructed motions [32], which often have floating and jitter issues [2].

We take a different approach by aligning music with videos generated by pre-trained Text-to-Video (T2V) models [44]. MusicInfuser does not require motion capture or motion reconstruction, but instead simply uses dance videos for training. Concretely, we propose adapter networks consisting of music-video cross-attention and a low-rank adapter, along with a novel training methodology and layer selection strategy for the cross-attention to maintain a balanced representation between multiple modalities with significant modality gaps.

Besides endowing the video diffusion model with the ability to listen to music, MusicInfuser preserves the rich knowledge in the text modality, enabling various forms of control and providing a flexible interface for the generation process. This means users can still leverage textual prompts to guide dance style, setting, and other aesthetic elements while maintaining synchronization with music (Fig. 1). Moreover, our framework can even general-

ize to make group choreography (Fig. 2) and longer dance videos with unseen music tracks (Fig. 3). To systematically evaluate our results, we developed an automatic evaluation framework using Video-LLMs [31] capable of processing video, audio, and language information simultaneously. This comprehensive approach allows us to assess multiple dimensions of dance quality, video quality, and prompt alignment within a single framework.

Our experiments demonstrate that MusicInfuser effectively bridges the gap between music and videos without requiring specialized motion data. By leveraging existing video diffusion models through targeted adaptation, we achieve high-quality and novel dance movements that respond naturally to musical rhythms and patterns, offering versatility in response to prompts and suggesting another direction for music-driven choreography synthesis.

2. Related Work

Music-to-Dance Generation Early music-to-dance generation research developed frameworks mapping music primitives to dance elements, using Hidden Markov Models [35]. Graph-based methods built movement transition graphs synchronized to musical beats [28], using quality rating functions and constraint-based dynamic programming [15]. Later, researchers incorporated Gaussian processes [16], Conditional Restricted Boltzmann Machines, Recurrent Neural Networks [1], and Convolutional Neural Networks [50, 51]. More recent approaches incorporate transformers [32, 50]. For instance, a Full-Attention Cross-Modal Transformer [32] predicts dance sequences based on seed motion and audio information. Traditional approaches typically generated movements synchronized with beats but lacking contextual meaning or showing excessive repetition [3], while showing limited choreographic diversity [4] and struggling with generalization. Recent dance generation advances have shifted toward diffusion-based approaches to address limitations of earlier methods [2, 29, 36, 37, 45]. Unlike these approaches that focus on motion skeleton synthesis from music, MusicInfuser directly synthesizes dance videos and choreography by uniquely adapting pre-trained text-to-video diffusion models to incorporate musical inputs while preserving their inherent knowledge of diverse dance styles and general human movements.

Controllable Approaches Dance generation systems have evolved to incorporate multiple input modalities for richer choreographic control [8, 18, 34], with text emerging as a powerful interface for its zero-shot capability and communicating choreographic ideas [34]. Transformer-based approaches using Vector Quantized-Variational Autoencoders create discrete motion tokens processable alongside text [41], while systems now process both text and mu-

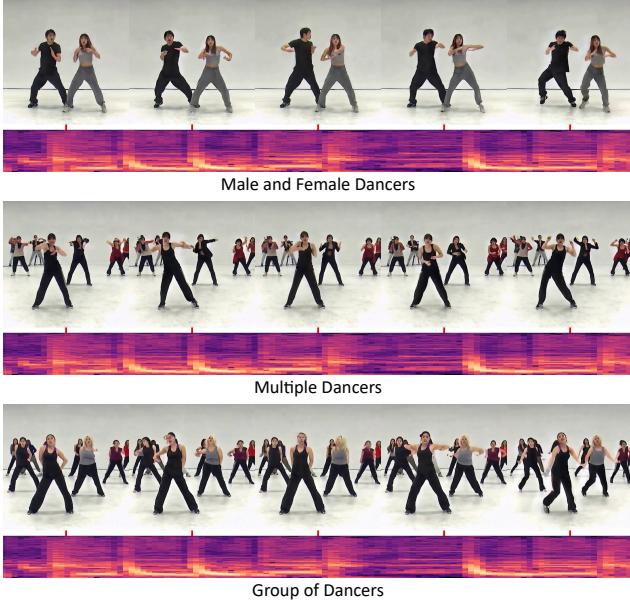


Figure 2. Due to the conservation of knowledge in video and text modalities, our model generalizes to generate group dance videos by modulating the prompt. To show this, the prompt is set to “*[DANCERS] dancing in a studio with a white backdrop, captured from a front view*,” where *[DANCERS]* denotes the description for each number of dancers.

sic inputs simultaneously [18]. Our MusicInfuser framework combines the flexibility of text-based interfaces with precise audio synchronization, allowing users to control stylistic and aesthetic elements of generated dance videos through prompts while ensuring movements remain aligned with musical features.

Audio-to-Video Generation Another domain that is adjacent to our method is audio-driven video generation. Pioneering this domain, Sound2Sight [9] introduced a deep variational encoder-decoder framework that predicts future frames by conditioning on both past frames and audio input. TATS [17] addressed audio-to-video generation challenges through a combination of time-agnostic VQGAN and time-sensitive transformer architectures. More recently, leveraging advances in diffusion models [21, 42], joint audio-video generation methods like MM-Diffusion [39] have been developed, enabling bidirectional generation where either modality can condition the other.

3. Preliminaries

Video Diffusion Models Diffusion models [6, 7, 21, 22, 42, 43, 49] represent a family of generative techniques that restore data via iterative denoising steps. The goal is to generate samples from a video distribution $p(\mathbf{x})$. To this end, we can define a convoluted distribution of $p(\mathbf{x})$ and a Gaussian

distribution with standard deviation σ , namely $p(\mathbf{x}, \sigma)$. In this paper, we follow [27] to construct a compact formulation of diffusion models. The denoiser D_θ is optimized with the following L2 objective:

$$\mathcal{L} = \mathbb{E}_{\mathbf{y} \sim p, \sigma \sim \Sigma_{\text{train}}, \mathbf{n} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})} \|D_\theta(\mathbf{y} + \mathbf{n}; \sigma) - \mathbf{y}\|_2^2, \quad (1)$$

where Σ_{train} denotes a noise distribution from which we sample noise during training, which is typically a uniform distribution. To sample with the denoiser D_θ , the ODE representing the change in the sample \mathbf{x} with the change in σ can be defined as:

$$\frac{d\mathbf{x}}{d\sigma} = -\frac{D_\theta(\mathbf{x}; \sigma) - \mathbf{x}}{\sigma}. \quad (2)$$

Text-Conditional Generation In a similar way, we can construct a conditional denoiser $D_\theta(\mathbf{x}|\mathbf{c}; \sigma)$ by training with a condition \mathbf{c} paired with each \mathbf{y} and replace the sampling process with a conditional denoiser. To boost generated content quality and alignment with prompts, classifier-free guidance (CFG) [20] has become widely used. Applying CFG, the modified ODE then becomes the linear combination form:

$$\frac{d\mathbf{x}}{d\sigma} = -\gamma_{\text{cfg}} \left[\frac{D_\theta(\mathbf{x}|\mathbf{c}; \sigma) - \mathbf{x}}{\sigma} \right] \quad (3)$$

$$+ (\gamma_{\text{cfg}} - 1) \left[\frac{D_\theta(\mathbf{x}; \sigma) - \mathbf{x}}{\sigma} \right]. \quad (4)$$

In this formulation, $D_\theta(\mathbf{x}; \sigma)$ shares the same parameters as $D_\theta(\mathbf{x}|\mathbf{c}; \sigma)$ but is trained by randomly omitting conditional information during training, and parameter γ_{cfg} denotes the guidance scale.

4. MusicInfuser

Text-to-Video models already know how to dance. These models, trained on diverse video datasets, have already internalized a wide range of human motion patterns including dance movements. The choreographic knowledge encoded in these models represents a valuable foundation that can be leveraged for music-driven video generation. Therefore, our goal is to adapt the models to construct a final denoiser, $D_\theta(\mathbf{x}|\mathbf{c}, \mathbf{a}; \sigma)$, that is also conditioned on a music condition \mathbf{a} . Rather than teaching these models to dance from scratch, our approach focuses on effectively aligning this pre-existing knowledge with audio input. In this section, we address methods that preserve the rich knowledge embedded in text-to-video models while establishing correlations between musical features and dance movements. The illustration of the overall architecture is displayed in Fig. 4.

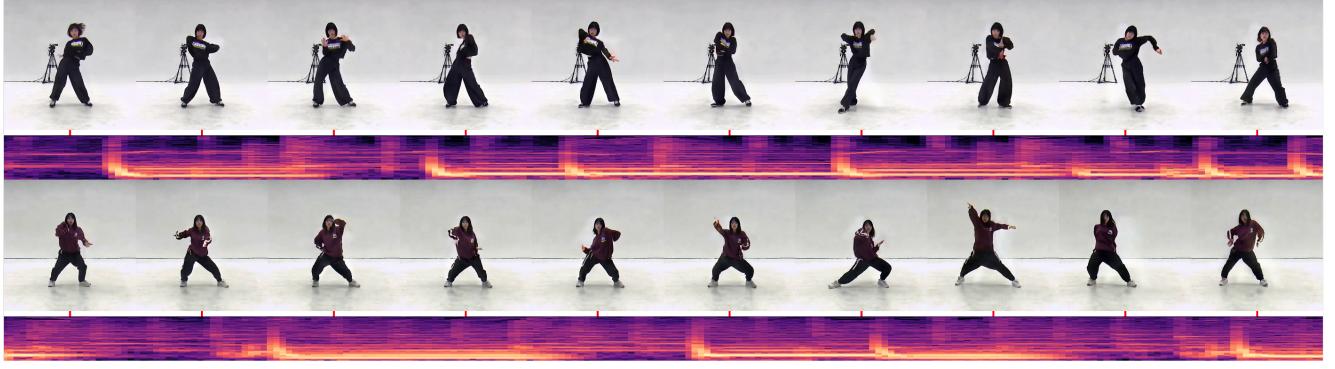


Figure 3. We generate longer dance videos that are 2 times longer than the videos used for training. For each row, we use synthetic in-the-wild music tracks with a keyword ‘‘K-pop,’’ a type of music not existing in AIST [46], and use a prompt ‘‘*a professional dancer dancing K-pop*’’ This shows our method is highly generalizable, even extending to longer videos with an unheard category of the music. The beat and style alignment can be more clearly observed in the videos on the project page.

4.1. Zero-Initialized cross-attention

Attention mechanisms [47] provide an effective mechanism for conditioning diffusion models [38, 40]. The cross and self attention are known to highly contribute to the structure of the generated images or videos [10, 13, 19, 23, 26, 40]. In our framework, we propose a Zero-Initialized cross-attention (ZICA) adapter to condition the musical information of the audio, while preserving the original capabilities of the base model.

Specifically, the audio signal is first processed through a specialized audio encoder [5] that extracts relevant temporal and spectral features. These encoded audio representations then pass through a learnable projector that maps them into the same embedding space as the video tokens. Then, the projected audio tokens interact with video tokens through cross-attention layers, allowing the model to establish correlations between audio patterns and corresponding visual choreography. Formally, if \mathbf{A} denotes the audio embeddings and \mathbf{V} denotes video tokens, the cross-attention operation can be expressed as:

$$\text{Attention}(\mathbf{V}, \mathbf{A}) = \text{softmax} \left(\frac{\mathbf{W}_Q \mathbf{V} \cdot (\mathbf{W}_K \mathbf{A})^T}{\sqrt{d}} \right) \mathbf{W}_V \mathbf{A}. \quad (5)$$

The complete cross-attention mechanism includes an output transformation matrix \mathbf{W}_O that projects the attention output back to the original feature space, with the full operation being:

$$\mathbf{Z} = \mathbf{V} + \mathbf{W}_O \text{Attention}(\mathbf{V}, \mathbf{A}). \quad (6)$$

By initializing \mathbf{W}_O to zero during the early stages of training, we effectively disable the influence of audio conditioning on the video features. This initialization strategy ensures that $\mathbf{Z} = \mathbf{V}$, meaning the output remains identical

to the input video features regardless of audio content. As training progresses, the parameters of \mathbf{W}_O gradually move away from zero, allowing the model to incrementally incorporate audio-conditioned information. This controlled integration preserves the rich choreographic knowledge of the base model while allowing for the gradual emergence of audio-visual correlations without disrupting the foundational video generation capabilities.

4.2. Low-Rank Adaptation with Higher Rank (HR-LoRA)

In addition to the cross-attention, we adapt the attention weights for diffusion transformer blocks. The adapter serves two key purposes: (1) to effectively integrate audio features into the text-video processing pipeline, and (2) to shift the domain toward our target application, synthesizing clear choreography. The ranks conventionally adopted for Low-Rank Adaptation (LoRA) [24] for visual models (typically 8 or 16) are optimized for image models and often lead to failure in capturing the full complexity of spatiotemporal information when directly applied to video models.

For adapting video tokens, a higher rank is required compared to image tokens, since video tokens contain temporal information. To effectively model motion adaptation separately from spatial adaptation, the optimal rank for the linear map should be increased compared to what’s needed for static images. For instance, adapting the full dimension of homography requires increasing the optimal rank by at least 8 (the degree of freedom for homography), while adapting general videos or complex human motion necessitates even higher ranks.

4.3. Beta-Uniform Scheduling

Current diffusion models, including those using LoRA fine-tuning, typically employ uniform scheduling for noise sampling during training. However, for adapter training for

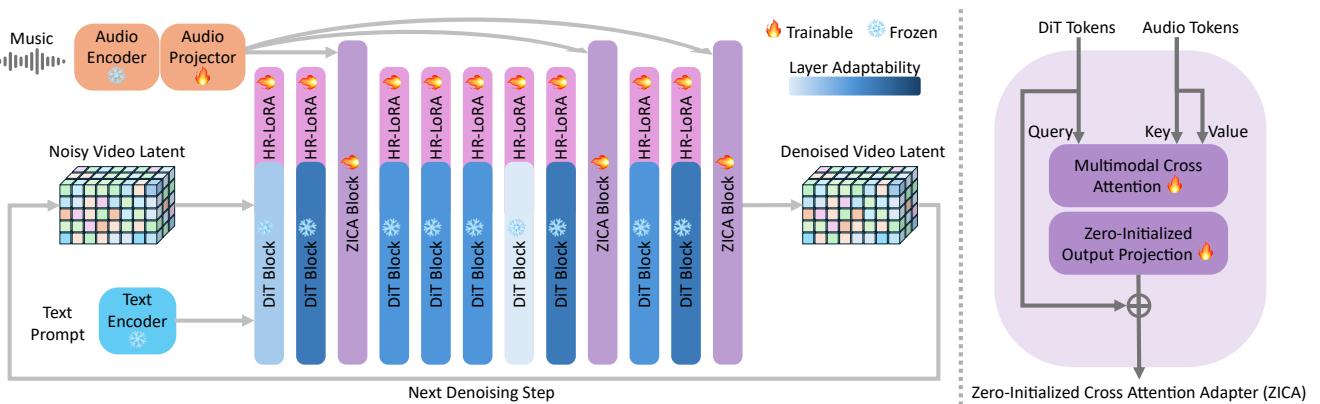


Figure 4. Overall architecture of MusicInfuser. Our framework adapts a pre-trained diffusion model with audio embedding using ZICA blocks (Sec. 4.1) and HR-LoRA blocks (Sec. 4.2). The placement of ZICA blocks is selected based on layer adaptability (Sec. 4.6).

dance video generation, our aim is to preserve the pre-trained models’ denoising capability of major components like coarse human motion at early stages and gradually learn the components over the course of the training process. To this end, we propose a Beta-Uniform scheduling strategy that makes the training noise distribution Σ_{train} evolve from a Beta distribution to a uniform distribution.

The Beta distribution with parameters $\alpha = 1$ and β is formally defined by the probability density function:

$$f(x; \alpha = 1, \beta) = \frac{(1-x)^{\beta-1}}{B(1, \beta)}, \quad 0 \leq x \leq 1 \quad (7)$$

where $B(\alpha, \beta)$ is the Beta function serving as a normalization constant. When $\beta > 1$, the distribution Beta($1, \beta$) concentrates probability mass near zero, which in our diffusion framework corresponds to sampling predominantly smaller noise scales. As β decays toward 1, the distribution gradually flattens, approaching Uniform($0, 1$), i.e., $\lim_{\beta \rightarrow 1} f(x; 1, \beta) = 1$, for all $0 \leq x \leq 1$.

This causes a smooth transition from focusing on high-frequency components at lower noise levels to broadly considering all frequencies. By first influencing the task-specific fine components of the dance and then the fundamental structure of dance movements, our approach preserves the rich knowledge of human motion and produces more coherent dance sequences. See the supplementary material for the details.

4.4. Utilizing In-the-Wild Videos

Training exclusively on highly structured datasets [32, 46] can lead to reduced generalizability and model degradation when confronted with diverse real-world scenarios. Therefore, we use a mixture of in-the-wild data with structured datasets. We gathered the data from YouTube playlists showing dance performances where we extracted video clips for training. These videos introduce valuable diversity in terms of camera angles, lighting conditions, performance

environments, and dance styles. The inclusion of in-the-wild data serves as regularization, preventing overfitting to specific dance patterns or environmental settings.

4.5. Obtaining and Diversifying Captions

We use caption templates for constrained setting datasets that provide consistent and structured textual descriptions. These templates contain placeholders for key attributes such as dance style, setting, and movement quality, which are populated based on the specific characteristics of each video. For in-the-wild videos, which lack standardized descriptions, we use VideoChat2 [31] for generating captions. VideoChat2 analyzes the visual content and generates detailed captions that capture the contextual information present in these diverse video samples.

Furthermore, we randomly replace a small portion of detailed captions with basic, simple captions. This way, the adapter network learns how to respond to the music without relying on the text. This effectively reduces the model’s dependence on textual cues and encourages it to develop stronger associations between musical features while still maintaining prompt adherence. See the supplementary material for the exact prompt templates we use for the diversification and replacement.

4.6. Selecting Layers Based on Adaptability

Applying cross-attention mechanisms to all layers of the model is computationally expensive and may introduce unnecessary complexity. However, exhaustively searching for an optimal combination of layers is intractable. To address this, we devise a novel selection criterion to selectively apply audio conditioning based on layer adaptability. Our method draws inspiration from Spatiotemporal Skip Guidance (STG) [26], which shares our motivation of identifying layers where modulation influences motion and structure without significantly causing predictions to deviate from the originally learned manifold. See the supplementary mate-

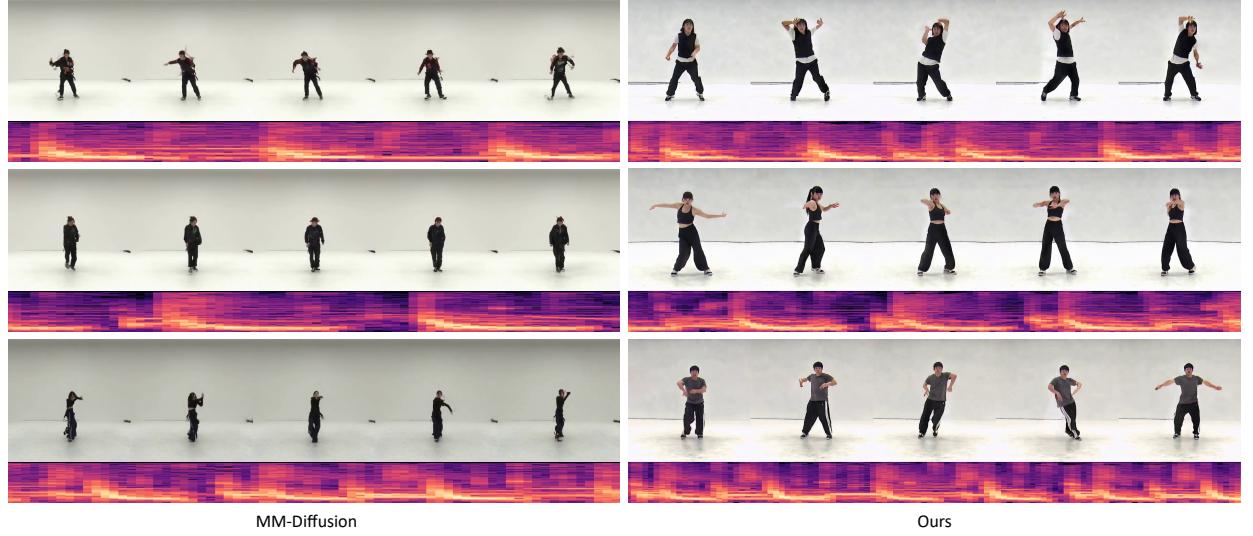


Figure 5. Comparison of audio-driven generation with MM-Diffusion [39]. Our method produces fewer artifacts (shown in the first and third rows), while generating more realistic dance videos with more natural movements (first row) and more dynamic motion (second and third rows). Note that we use the same music track for each row, and the spectrogram is stretched for MM-Diffusion since we generate longer videos. For our method, we use the fixed caption “a professional dancer dancing ...” across all music tracks.

Model	Modality	Style Alignment	Beat Alignment	Body Representation	Movement Realism	Choreography Complexity	Dance Quality Average
AIST Dataset (GT) [46]	A+V	7.46	8.95	7.53	8.67	7.45	8.01
MM-Diffusion [39]	A+V	7.16	8.56	5.52	7.05	7.53	7.16
Mochi [44]	T+V	7.20	8.34	7.47	7.68	7.82	7.70
MusicInfuser (Ours)	T+A+V	7.56	8.89	7.16	8.24	7.90	7.95

Table 1. Dance quality metrics comparing different models. A, V, and T denote audio, video, and text input modalities, respectively. For the models that have text input modality, we report an average of scores using a predefined benchmark of prompts.

rial for details on the selected layers.

5. Experiments

5.1. Implementation Details

Dataset The AIST dataset [46] includes 13,940 videos with 60 musical pieces, 10 dance genres, and 35 dancers. We utilize AIST as one of our primary data sources for training our audio-driven dance generation model. We extract 2,378 clips and strictly divide training and test sets with non-overlapping music tracks, following [32]. For each training instance, we randomly sample approximately 2.5-second clips from the full sequences. To enhance the quality and generalizability of our model, we supplement the structured AIST data with in-the-wild dance videos. Specifically, we extract 15,799 video clips from YouTube dance videos, encompassing a wide variety of dance styles, settings, and production qualities. These clips are mixed with those from AIST at a 1:1 ratio during training, creating a balanced dataset that leverages both the controlled environment of AIST and the rich diversity of real-world dance performances.

Model Details We train our model on a single NVIDIA A100 GPU for 4,000 steps with a learning rate of $1e-4$, which takes roughly 20 hours to complete. Our HR-LoRA uses rank 64, providing sufficient capacity to capture complex dance movements while maintaining parameter efficiency. For Beta-Uniform scheduling, we set initial $\beta = 3$ with exponential decay toward $\beta = 1$. We use Mochi [44] as our base model with a classifier-free guidance scale of $\gamma_{\text{cfg}} = 6.0$ during inference and employ Wav2Vec 2.0 [5] as the audio encoder.

Quantitative Metrics Evaluating generated content with plausible metrics presents a significant challenge. Inspired by VBench’s use of Visual-Language Models (VLMs) [31] for text-to-video assessment, we propose a novel metric based on Video-LLMs [14] to evaluate generated videos alongside their conditioning audio and prompts. Specifically, we leverage VideoLLaMA 2 [14] and formulate targeted queries to assess three key components: dance quality, video quality, and prompt alignment. For dance quality, we evaluate style alignment, beat alignment, body representation, movement realism, and choreography complexity. Video quality assessment includes imaging quality,



Top row: slowed down, middle row: original speed, bottom row: sped up

Figure 6. Speed control. The audio input is slowed down (the top row) or sped up (the bottom row) by 0.75 times and 1.25 times, respectively. This shows that speeding up generally results in more movements. Also see the change in the dynamicity, as speeding up the audio also increases the tone of the music.

Model	Modality	Imaging Quality	Aesthetic Quality	Overall Consistency	Video Quality Average
AIST Dataset (GT) [46]	A+V	9.76	8.17	9.77	9.23
MM-Diffusion [39]	A+V	8.94	6.52	8.38	7.94
Mochi [44]	T+V	9.46	7.90	8.98	8.78
MusicInfuser (Ours)	T+A+V	9.60	7.87	9.39	8.95

Table 2. Video quality metrics comparing different models. For the models that have text input modality, we report an average of scores using a predefined benchmark of prompts.

Model	Style Capture	Creative Interpretation	Overall Satisfaction	Prompt Align Average
Mochi [44]	7.98	9.04	9.55	8.86
MusicInfuser (Ours)	7.80	9.27	9.80	8.96
No in-the-Wild Data	6.80	8.69	8.40	7.96
Base Prompt 0%	7.45	8.85	9.43	8.58
Base Prompt 100%	7.33	9.06	9.36	8.58

Table 3. Prompt alignment metrics comparing different models.

aesthetic quality, and overall consistency. Prompt alignment evaluation covers style capture, creative interpretation, and overall satisfaction. In Table 1 and 2, we present results on the AIST test data to demonstrate the plausibility. In dance quality metrics, the AIST test data—which should serve as an upper bound for some metrics such as beat alignment, imaging quality, and movement realism—outperforms other models (including ours) in such metrics.

5.2. Experimental Results

Music- and Text-Driven Dance Video Generation
Fig. 1 showcases the model’s ability to combine textual control with musical synchronization. The generated videos successfully incorporate scene contexts (restaurant kitchen, beach at sunset) and dancer attributes (wearing a leather jacket, chef’s uniform) as specified in the prompts, while simultaneously aligning the choreographic style with the musical input. Fig. 2 demonstrates that our model can generalize to generate group dance videos by simply modifying the prompt to reference multiple dancers.

Music Responsiveness In Fig. 5, we show how MusicInfuser generates dance videos, including the movement and outfit of the dancer, based on the music condition, while keeping the prompt fixed. Additionally, we demonstrate the model’s responsiveness to musical features through exper-



Figure 7. Videos generated with three distinct in-the-wild music tracks created with SUNO AI. For each row, we use in-the-wild music tracks generated with a word “K-pop,” an unseen category.

iments with tempo modification. By accelerating the music track by 1.25 times or decelerating it by 0.75 times, the generated dance movements appropriately adjust pace while maintaining similar choreographic style, as shown in Fig. 6. Furthermore, acceleration and deceleration also result in changes in tone, which affect the dynamicity of the dance generated by our model. This shows that our model successfully captures the relationship between musical tempo and dance movement dynamicity, a critical aspect of dance-music synchronization.

Generalization to In-the-Wild Music and Longer Videos

To evaluate generalization beyond the AIST music distribution, we test our model with music tracks generated by SUNO AI, which represent styles not represented in the training data. Fig. 7 shows successful generation for these unseen music categories, confirming the model’s ability to map novel audio patterns to appropriate dance movements. Fig. 3 shows longer video generation results with the same setting but with twice as many frames as the videos we used for training. These experiments show flexible generalization of our method.

Model	Style Alignment	Beat Alignment	Body Representation	Movement Realism	Choreography Complexity	Dance Quality Average	Imaging Quality	Aesthetic Quality	Overall Consistency	Video Quality Average	Overall
Full	7.56	8.89	7.16	8.24	7.90	7.95	9.60	7.87	9.39	8.95	8.33
No ZICA Layer Selection	7.31	8.81	7.28	7.70	7.96	7.81	9.33	7.78	9.04	8.72	8.15
No Higher Rank	7.37	8.76	6.86	7.75	7.98	7.74	9.55	7.94	9.49	8.99	8.21
No LoRA	7.48	8.62	7.02	7.53	7.95	7.72	9.43	8.08	9.36	8.96	8.18
No Beta-Uniform Schedule	8.04	9.07	6.35	7.88	7.91	7.85	9.17	7.85	9.37	8.80	8.21
Feature Addition	7.62	8.90	6.78	7.97	7.88	7.83	9.44	7.88	9.31	8.88	8.22

Table 4. Ablation study. Feature addition denotes that we spatially expand the audio feature and add the feature to the corresponding frame.

Comparison with Prior Work Fig. 5 presents a side-by-side comparison with MM-Diffusion [39]. MusicInfuser produces more consistent human forms, fewer visual artifacts, and more fluid, realistic movements compared to both baselines. Notably, our method also demonstrates superior alignment between music beat patterns and dance movements, visible in the correspondence between spectrogram intensity and motion transitions.

Our approach overcomes key limitations of prior work: unlike MM-Diffusion, which generates shorter videos with limited style control, MusicInfuser produces longer sequences with both musical synchronization and prompt-based style control, while improving overall consistency of the video and having fewer artifacts. Compared to Mochi, our method adds music responsiveness while maintaining or improving video quality and prompt adherence.

Quantitative Evaluation Tables 1–2 present our quantitative results compared against several baselines [39, 44]. For dance quality (Table 1), our method outperforms previous approaches in style alignment, beat alignment, movement realism, and choreography complexity, while maintaining competitive scores across other metrics. Table 2 demonstrates our superiority in video quality metrics, particularly in imaging quality and overall consistency compared to MM-Diffusion [39] and Mochi [44]. In Table 3, MusicInfuser shows improved prompt alignment over the baseline Mochi model, with significant gains in creative interpretation and overall satisfaction.

Ablation Studies Our ablation study (Table 4) systematically evaluates the contribution of each component in our framework. The full model achieves the highest overall score, with the layer adaptability selection providing the most significant contribution. HR-LoRA contributes substantially to movement realism, while our Beta-Uniform scheduling improves body representation. The naive feature addition baseline, where instead of using the ZICA adapter we simply spatially expand the audio feature and add it to the corresponding frame, performs worse than our approach in most metrics, confirming the effectiveness of our ZICA strategy. Additionally, in Table 3, we show the trade-off between style capture and creative interpretation of the prompt depending on the base prompt ratio of 0% and 100%.

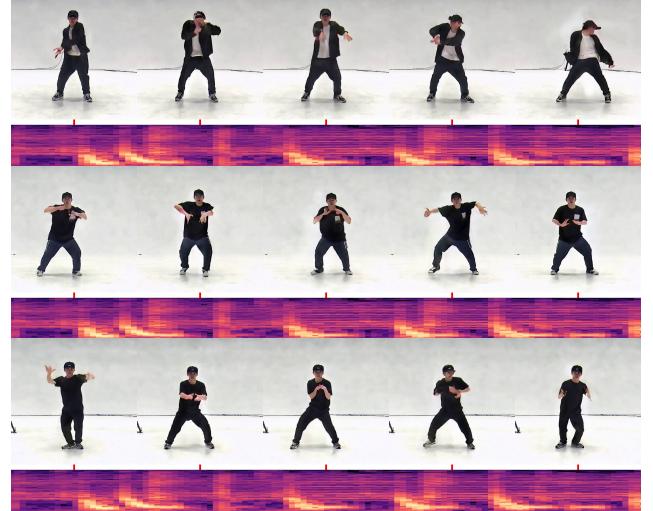


Figure 8. By changing the seed, our method can produce diverse results given the same music and text. The generated choreography of each dance is different from each other. We use the fixed prompt “a professional dancer dancing”

Diversity of Results By varying the random seed while keeping the prompt and music constant, our model generates diverse choreographies as shown in Fig. 8, demonstrating that it does not simply memorize specific dance routines for particular music tracks but instead learns a rich mapping between music and movement possibilities.

6. Conclusion

In this paper, we present MusicInfuser, a novel approach for generating dance videos synchronized with music by leveraging the rich choreographic knowledge embedded in pre-trained text-to-video diffusion models. Through our adaptation architecture and strategies, MusicInfuser successfully enables synchronized dance movements with musical inputs while preserving text-based control over style and scene elements. It achieves this without requiring expensive motion capture data, generalizes to novel music tracks, and supports the generation of diverse choreographies and group dance videos. By circumventing the prohibitive costs and restrictions of conventional approaches, MusicInfuser opens new possibilities for creative music-driven dance video synthesis.

Acknowledgments We thank Xiaojuan Wang and Jingwei Ma for their valuable feedback. This work was supported by the UW Reality Lab and Google.

References

- [1] Omid Alemi, Jules Fran oise, and Philippe Pasquier. Groovenet: Real-time music-driven dance movement generation using artificial neural networks. *networks*, 8(17):26, 2017. [2](#)
- [2] Simon Alexanderson, Rajmund Nagy, Jonas Beskow, and Gustav Eje Henter. Listen, denoise, action! audio-driven motion synthesis with diffusion models. *ACM Transactions on Graphics (TOG)*, 42(4):1–20, 2023. [2](#)
- [3] Andreas Aristidou, Anastasios Yiannakidis, Kfir Aberman, Daniel Cohen-Or, Ariel Shamir, and Yiorgos Chrysanthou. Rhythm is a dancer: Music-driven motion synthesis with global structure. *IEEE transactions on visualization and computer graphics*, 29(8):3519–3534, 2022. [2](#)
- [4] Ho Yin Au, Jie Chen, Junkun Jiang, and Yike Guo. Chorograph: Music-conditioned automatic dance choreography over a style and tempo consistent dynamic graph. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 3917–3925, 2022. [2](#)
- [5] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in neural information processing systems*, 33:12449–12460, 2020. [4, 6](#)
- [6] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023. [3](#)
- [7] Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 22563–22575, 2023. [3](#)
- [8] Caroline Chan, Shiry Ginosar, Tinghui Zhou, and Alexei A Efros. Everybody dance now. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5933–5942, 2019. [2](#)
- [9] Moitreya Chatterjee and Anoop Cherian. Sound2sight: Generating visual dynamics from sound and context. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXVII 16*, pages 701–719. Springer, 2020. [3](#)
- [10] Hila Chefer, Yuval Alaluf, Yael Vinker, Lior Wolf, and Daniel Cohen-Or. Attend-and-excite: Attention-based semantic guidance for text-to-image diffusion models. *ACM transactions on Graphics (TOG)*, 42(4):1–10, 2023. [4](#)
- [11] Hila Chefer, Uriel Singer, Amit Zohar, Yuval Kirstain, Adam Polyak, Yaniv Taigman, Lior Wolf, and Shelly Sheynin. Videojam: Joint appearance-motion representations for enhanced motion generation in video models. *arXiv preprint arXiv:2502.02492*, 2025. [14](#)
- [12] Kang Chen, Zhipeng Tan, Jin Lei, Song-Hai Zhang, Yuan-Chen Guo, Weidong Zhang, and Shi-Min Hu. Choreomaster: choreography-oriented music-driven dance synthesis. *ACM Transactions on Graphics (TOG)*, 40(4):1–13, 2021. [2](#)
- [13] Minghao Chen, Iro Laina, and Andrea Vedaldi. Training-free layout control with cross-attention guidance. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 5343–5353, 2024. [4](#)
- [14] Zesen Cheng, Sicong Leng, Hang Zhang, Yifei Xin, Xin Li, Guanzheng Chen, Yongxin Zhu, Wenqi Zhang, Ziyang Luo, Deli Zhao, et al. Videollama 2: Advancing spatial-temporal modeling and audio understanding in video-lmms. *arXiv preprint arXiv:2406.07476*, 2024. [6, 14](#)
- [15] Rukun Fan, Songhua Xu, and Weidong Geng. Example-based automatic music-driven conventional dance motion synthesis. *IEEE transactions on visualization and computer graphics*, 18(3):501–515, 2011. [2](#)
- [16] Satoru Fukayama and Masataka Goto. Automated choreography synthesis using a gaussian process leveraging consumer-generated dance motions. In *Proceedings of the 11th Conference on Advances in Computer Entertainment Technology*, pages 1–6, 2014. [2](#)
- [17] Songwei Ge, Thomas Hayes, Harry Yang, Xi Yin, Guan Pang, David Jacobs, Jia-Bin Huang, and Devi Parikh. Long video generation with time-agnostic vqgan and time-sensitive transformer. In *European Conference on Computer Vision*, pages 102–118. Springer, 2022. [3](#)
- [18] Kehong Gong, Dongze Lian, Heng Chang, Chuan Guo, Zihang Jiang, Xinxin Zuo, Michael Bi Mi, and Xinchao Wang. Tm2d: Bimodality driven 3d dance generation via music-text integration. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9942–9952, 2023. [2, 3](#)
- [19] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*, 2022. [4](#)
- [20] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. [3](#)
- [21] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. [3](#)
- [22] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *Advances in Neural Information Processing Systems*, 35:8633–8646, 2022. [3](#)
- [23] Susung Hong, Gyuseong Lee, Wooseok Jang, and Seungryong Kim. Improving sample quality of diffusion models using self-attention guidance. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7462–7471, 2023. [4](#)
- [24] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022. [4](#)
- [25] Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing Wu, Qingyang Jin,

- Nattapol Champaisit, et al. Vbench: Comprehensive benchmark suite for video generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21807–21818, 2024. 13
- [26] Junha Hyung, Kinam Kim, Susung Hong, Min-Jung Kim, and Jaegul Choo. Spatiotemporal skip guidance for enhanced video diffusion sampling. *arXiv preprint arXiv:2411.18664*, 2024. 4, 5, 13, 17
- [27] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. *Advances in neural information processing systems*, 35:26565–26577, 2022. 3
- [28] Tae-hoon Kim, Sang Il Park, and Sung Yong Shin. Rhythmic-motion synthesis based on motion-beat analysis. *ACM Transactions on Graphics (TOG)*, 22(3):392–401, 2003. 2
- [29] Nhat Le, Tuong Do, Khoa Do, Hien Nguyen, Erman Tjiputra, Quang D Tran, and Anh Nguyen. Controllable group choreography using contrastive diffusion. *ACM Transactions on Graphics (TOG)*, 42(6):1–14, 2023. 2
- [30] Hsin-Ying Lee, Xiaodong Yang, Ming-Yu Liu, Ting-Chun Wang, Yu-Ding Lu, Ming-Hsuan Yang, and Jan Kautz. Dancing to music. *Advances in neural information processing systems*, 32, 2019. 2
- [31] KunChang Li, Yinan He, Yi Wang, Yizhuo Li, Wenhui Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao. Videochat: Chat-centric video understanding. *arXiv preprint arXiv:2305.06355*, 2023. 2, 5, 6, 13
- [32] Ruilong Li, Shan Yang, David A Ross, and Angjoo Kanazawa. Ai choreographer: Music conditioned 3d dance generation with aist++. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 13401–13412, 2021. 2, 5, 6, 13
- [33] Gaojie Lin, Jianwen Jiang, Jiaqi Yang, Zerong Zheng, and Chao Liang. Omnihuman-1: Rethinking the scaling-up of one-stage conditioned human animation models. *arXiv preprint arXiv:2502.01061*, 2025. 14
- [34] Yimeng Liu and Misha Sra. Dancegen: Supporting choreography ideation and prototyping with generative ai. In *Proceedings of the 2024 ACM Designing Interactive Systems Conference*, pages 920–938, 2024. 2
- [35] Ferda Ofli, Engin Erzin, Yücel Yemez, and A Murat Tekalp. Learn2dance: Learning statistical music-to-dance mappings for choreography synthesis. *IEEE Transactions on Multimedia*, 14(3):747–759, 2011. 2
- [36] Qiaosong Qi, Le Zhuo, Aixi Zhang, Yue Liao, Fei Fang, Si Liu, and Shuicheng Yan. Diffdance: Cascaded human motion diffusion model for dance generation. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 1374–1382, 2023. 2
- [37] Liangdong Qiu, Chengxing Yu, Yanran Li, Zhao Wang, Haibin Huang, Chongyang Ma, Di Zhang, Pengfei Wan, and Xiaoguang Han. Vimo: Generating motions from casual videos. *arXiv preprint arXiv:2408.06614*, 2024. 2
- [38] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 4
- [39] Ludan Ruan, Yiyang Ma, Huan Yang, Huiguo He, Bei Liu, Jianlong Fu, Nicholas Jing Yuan, Qin Jin, and Baining Guo. Mm-diffusion: Learning multi-modal diffusion models for joint audio and video generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10219–10228, 2023. 2, 3, 6, 7, 8
- [40] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in neural information processing systems*, 35:36479–36494, 2022. 4
- [41] Li Siyao, Weijiang Yu, Tianpei Gu, Chunze Lin, Quan Wang, Chen Qian, Chen Change Loy, and Ziwei Liu. Bailando: 3d dance generation by actor-critic gpt with choreographic memory. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11050–11059, 2022. 2
- [42] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *ICLR*, 2021. 3
- [43] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *ICLR*, 2021. 3
- [44] Genmo Team. Mochi 1. <https://github.com/genmoai/models>, 2024. 2, 6, 7, 8, 12, 13, 14
- [45] Jonathan Tseng, Rodrigo Castellon, and Karen Liu. Edge: Editable dance generation from music. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 448–458, 2023. 2
- [46] Shuhei Tsuchida, Satoru Fukayama, Masahiro Hamasaki, and Masataka Goto. Aist dance video database: Multi-genre, multi-dancer, and multi-camera database for dance information processing. In *ISMIR*, page 6, 2019. 4, 5, 6, 7, 13
- [47] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 4
- [48] Zijie Ye, Haozhe Wu, Jia Jia, Yaohua Bu, Wei Chen, Fanbo Meng, and Yanfeng Wang. Choreonet: Towards music to dance synthesis with choreographic action unit. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 744–752, 2020. 2
- [49] Sihyun Yu, Kihyuk Sohn, Subin Kim, and Jinwoo Shin. Video probabilistic diffusion models in projected latent space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 18456–18466, 2023. 3
- [50] Mingao Zhang, Changhong Liu, Yong Chen, Zhenchun Lei, and Mingwen Wang. Music-to-dance generation with multiple conformer. In *Proceedings of the 2022 International Conference on Multimedia Retrieval*, pages 34–38, 2022. 2
- [51] Wenlin Zhuang, Congyi Wang, Jinxiang Chai, Yangang Wang, Ming Shao, and Siyu Xia. Music2dance: Dancenet for music-driven dance generation. *ACM Transactions on*

MusicInfuser: Making Video Diffusion Listen and Dance

Supplementary Material

A. Video Results

We present the flattened video results along the time axis and the corresponding spectrograms in the main paper. However, our frame sampling rate does not exceed the Nyquist frequency for the general musical beat, causing the movement to appear slower. Therefore, we encourage readers to visit the project page to view the video results.

B. Beta Distribution

Fig. 9 shows the evolution of graphs of Beta distributions from Beta(1, 3) to Beta(1, 1).

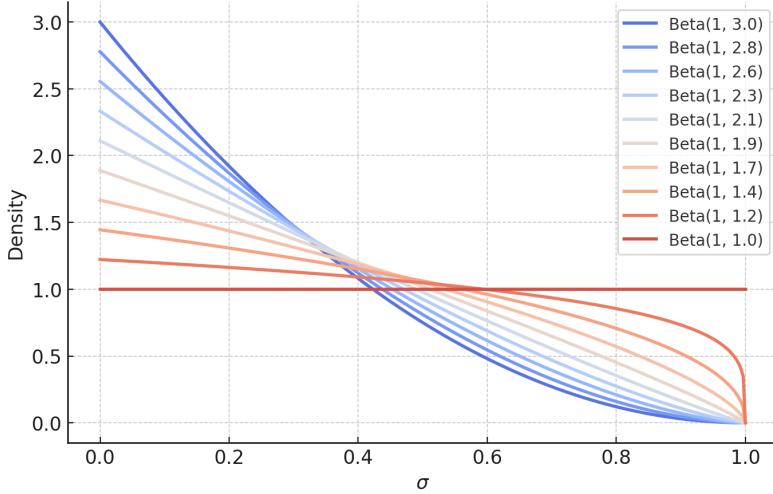


Figure 9. Beta distributions.

C. Difficulty Control

We demonstrate difficulty control of the choreography in Fig. 10, which is achieved using the same seed and music but with prompts of varying specificity. For basic dance, we use the general prompt “a professional dancing in a studio with a white backdrop.” For styled dance, we additionally specify the dance genre but use “basic dance setting,” and for advanced, we change it to “advanced dance setting.”

D. Limitations

Although our method adds listening capability to text-to-video models and improves dance generation, some properties such as style capture of the prompt and imaging quality are bounded by the capabilities of the models. Also, it inherits some problems from text-to-video models. Sometimes, fine parts such as fingers and faces fail to be generated properly, especially when our model synthesizes dance videos with fast movements. Additionally, our model is easily fooled by the silhouette of the dancers, which means under the same silhouette, they merge or change the positions of body parts, which is also a problem in the base model. We include some examples of the failure cases in Fig. 16.

E. Comparison with the Base Model

We show a comparison with Mochi [44] in Fig. 11. Note that Mochi is not able to hear the music.

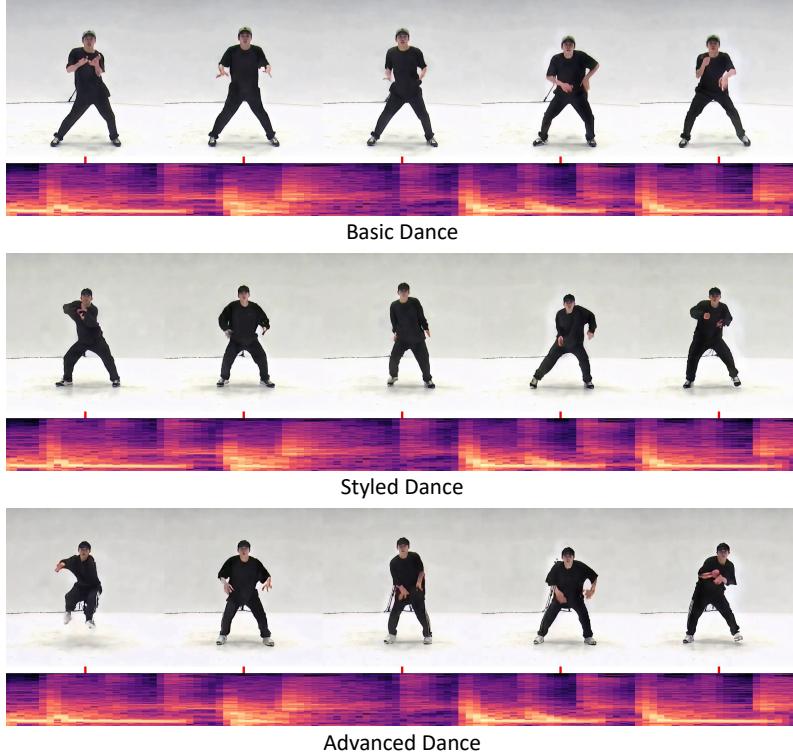


Figure 10. Changes in the complexity of choreography.

F. Ablation Study

We present qualitative results of our ablation study in Fig. 12 and Fig. 13. Our full model successfully generates consistent body shapes that align with the music, while preserving prior knowledge without introducing significant artifacts.

G. Layer Adaptability

The imaging and aesthetic quality of the base model [44] is presented in Fig. 14. This is analyzed with STG [26], an inference-time technique, and the score is calculated with VBench [25]. Based on the imaging quality, we choose the top 16 out of 48 layers in terms of image quality to select the layers.

H. Test Music Tracks

For evaluating our method, we use music tracks that are set aside from the training set [46], following AIST++ [32]. The full list of the test music codes is listed in Table 5.

I. In-the-Wild Music Tracks

We use in-the-wild music tracks to demonstrate the generalizability of our method. In the main paper, we use SUNO AI to generate the music tracks, which produces music for non-commercial use. In addition, in Fig. 15, we show the results generated with three distinct music tracks from TikTok videos.

J. Prompts

As mentioned in our main paper, we use a proper prompt format and base prompt for AIST [46]. The full list is shown in Table 6. Note that since we use VideoChat2 [31] to label YouTube videos, we have only the base prompt for the dataset. We also provide a predefined set of prompts in Table 7 that is used to generate samples for the evaluation, which ultimately



Figure 11. MusicInfuser infuses listening capability into the text-to-video model (Mochi [44]), while preserving the prompt adherence and improving overall consistency and realism (frames 2 and 5 of the top example, and frames 2–4 of the bottom example).

Test Music Code	Genre
mLH4	LA style Hip-hop
mKR2	Krump
mBR0	Break
mLO2	Lock
mJB5	Ballet Jazz
mWA0	Waack
mJS3	Street Jazz
mMH3	Middle Hip-hop
mHO5	House
mPO1	Pop

Table 5. List of test music codes with corresponding dance genres.

results in $10 \times 10 = 100$ videos for evaluation for each model configuration. The system prompts for VideoLLaMA [14] used for evaluation are in Table 8.

K. Concurrent Work

Several concurrent approaches have emerged alongside our research that address related challenges. Notable among these is VideoJAM [11], which enhances motion generation by jointly denoising both the motion maps and the video, an approach that is orthogonal to ours. Another line of research is OmniHuman-1 [33], which integrates audio and pose inputs into diffusion models. The application of OmniHuman-1 remains primarily confined to scenarios that do not require much creative movement, relies on a private model, and necessitates full fine-tuning procedures, which distinguishes it from our approach.



Figure 12. Ablation study. The prompt is set to “a male dancer dancing in an art gallery with some paintings, captured from a front view”. The seed and music are set the same across all methods.

L. More Results

We show more music-and-text-to-video generation examples in Fig. 17.

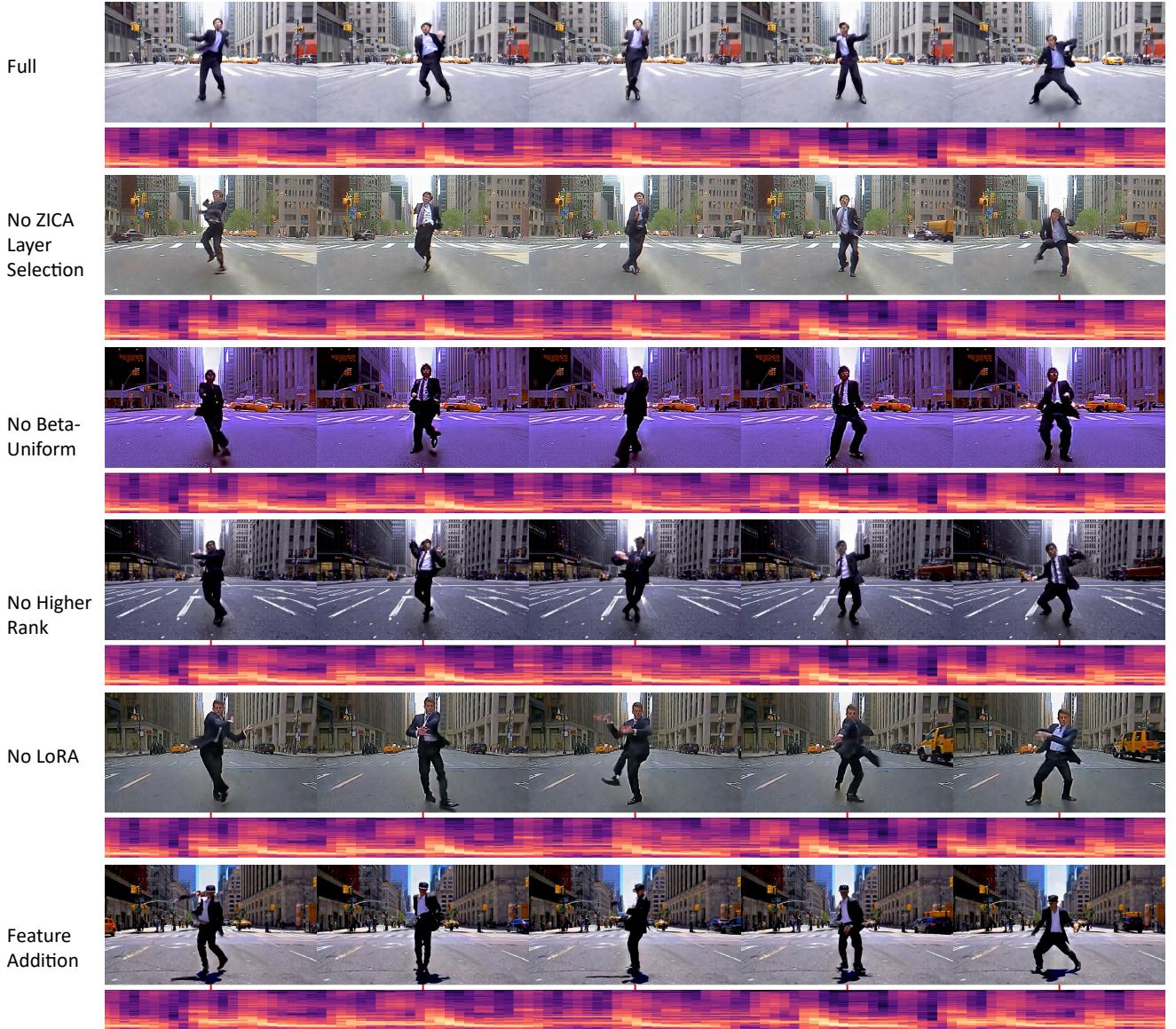


Figure 13. Ablation study. The prompt is set to “a male dancer wearing a suit dancing in the middle of a New York City, captured from a front view”. The seed and music are set the same across all methods.

Category	Dataset	Prompt Template
Prompt Format	AIST	{dancers_text} dancing {genre_name} in a {situation_name} setting in a studio with a white backdrop, captured from a {camera_view}
Prompt Format	AIST	a {camera_view} video of {dancers_text} performing {genre_name} choreography against a white background in a {situation_name} scene
Prompt Format	AIST	{dancers_text} executing {genre_name} movements in a minimalist studio space in a {situation_name} setting, shot from a {camera_view}
Prompt Format	AIST	a {genre_name} dance performance by {dancers_text} in a pristine white studio, {camera_view}, {situation_name}
Base Prompt	AIST	a professional dancer dancing in a studio with a white backdrop
Base Prompt	YouTube	a dance video

Table 6. Dance prompt templates categorized by type and dataset, including parameterized formats and simple base prompts.

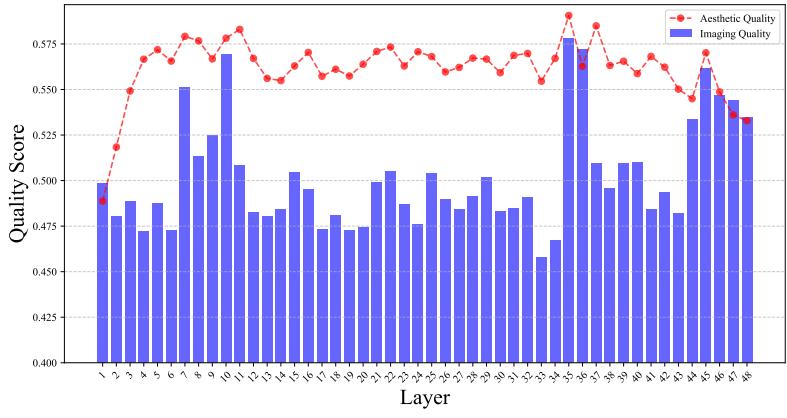


Figure 14. Layer adaptability graph from [26], showing imaging and aesthetic quality.

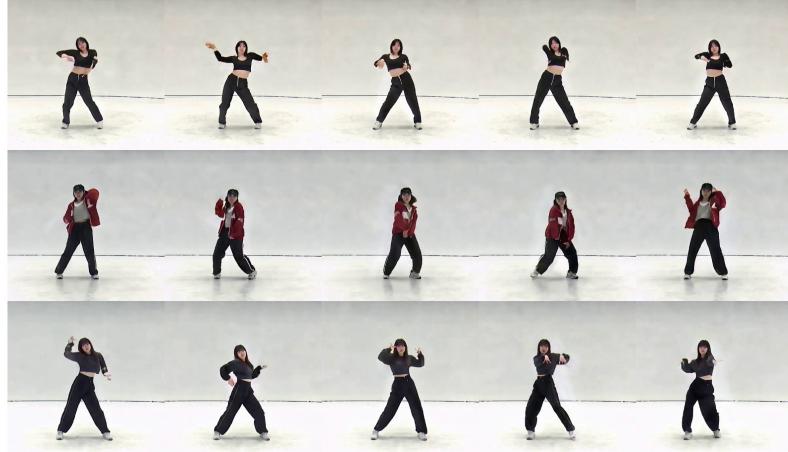


Figure 15. Videos generated with three distinct in-the-wild music tracks.



Figure 16. Failure cases. Our model inherits some issues from the base model, such as failing to generate fine details (e.g., fingers and faces) and being fooled by the silhouette of the dancers.

Prompts
a male dancer dancing on a rooftop at sunset, captured from a front view
a female dancer dancing in a subway station, captured from a front view
a male dancer dancing in an art gallery with some paintings, captured from a front view
a female dancer wearing a leather jacket dancing in a studio with a white backdrop, captured from a front view
a male dancer wearing a hoodie dancing in a studio with a white backdrop, captured from a front view
a female dancer wearing a denim vest dancing in a studio with a white backdrop, captured from a front view
a female dancer wearing a Hawaiian dress dancing on Waikiki Beach at sunset with Diamond Head in the background, captured from a front view
a male dancer wearing a suit dancing in the middle of a New York City, captured from a front view
a male dancer wearing a chef's uniform dancing in a busy restaurant kitchen with flames from the grill behind him, captured from a front view
a female dancer wearing a Renaissance gown dancing in a Venetian masquerade ball with ornate chandeliers overhead, captured from a front view

Table 7. Collection of dance scene prompts with various subjects, attire, and settings.

Metric	Prompt
Dance Quality	
Style Alignment	Rate the style alignment of the dance to music where: 0 means poor style alignment of the dance to music, 5 means moderate style alignment of the dance to music, and 10 means perfect style alignment of the dance to music. Output only the number.
Beat Alignment	Rate the beat alignment of the dance to music where: 0 means poor beat alignment of the dance to music, 5 means moderate beat alignment of the dance to music, and 10 means perfect beat alignment of the dance to music. Output only the number.
Body Representation	Rate the body representation of the dancer where: 0 means unrealistic/distorted proportions of the dancer, 5 means minor anatomical issues of the dancer, and 10 means anatomically perfect representation of the dancer. Output only the number.
Movement Realism	Rate the movement realism of the dancer where: 0 means poor movement realism of the dancer, 5 means moderate movement realism of the dancer, and 10 means perfect movement realism of the dancer. Output only the number.
Choreography Complexity	Rate the complexity of the choreography where: 0 means extremely basic choreography, 5 means intermediate choreography, and 10 means extremely complex/advanced choreography. Output only the number.
Video Quality	
Imaging Quality	Rate the imaging quality where: 0 means poor imaging quality, 5 means moderate imaging quality, and 10 means perfect imaging quality. Output only the number.
Aesthetic Quality	Rate the aesthetic quality where: 0 means poor aesthetic quality, 5 means moderate aesthetic quality, and 10 means perfect aesthetic quality. Output only the number.
Overall Consistency	Rate the overall consistency where: 0 means poor consistency, 5 means moderate consistency, and 10 means perfect consistency. Output only the number.
Prompt Alignment	
Style Capture	How well does the dance video capture the specific style mentioned in the prompt: '{prompt}'? Rate 0-10 where: 0 means completely missed the style, 5 means some elements of the style are present, and 10 means perfectly captures the style. Output only the number.
Creative Interpretation	Based on the prompt '{prompt}', rate the creativity in interpreting the prompt 0-10 where: 0 means generic/standard interpretation, 5 means moderate creativity, and 10 means highly creative and unique interpretation. Output only the number.
Overall Prompt Satisfaction	Rate the overall prompt satisfaction 0-10 where: 0 means the video fails to satisfy the prompt '{prompt}', 5 means it partially satisfies the prompt, and 10 means it fully satisfies all aspects of the prompt. Output only the number.

Table 8. System prompts for evaluation



Figure 17. More music-and-text-to-video generation results.