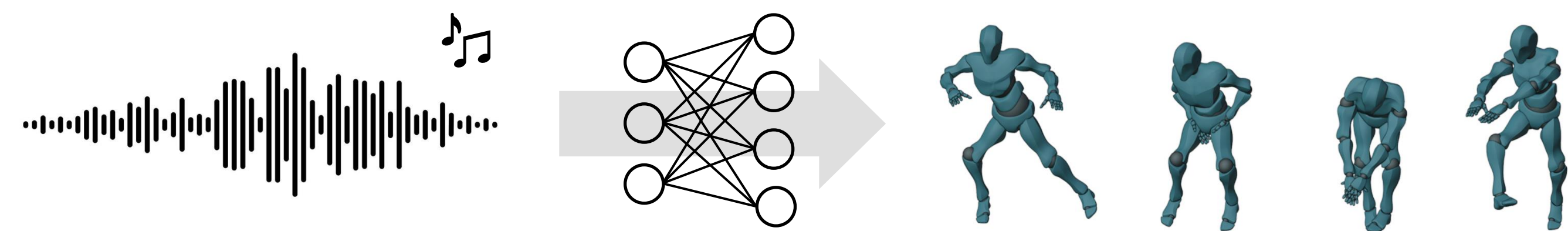


## Goal

### ✓ Music-driven 3D dance generation task



- Dance is a specific form of human motion along to **music**, which is expressive, symbolic and **rhythmic**
- Our goal is to generate plausible and diverse dance movements, conditioned by given music

## Motivation

### ✓ H1. **Mamba**<sup>[1]</sup> operates 3D dance data better than Transformer

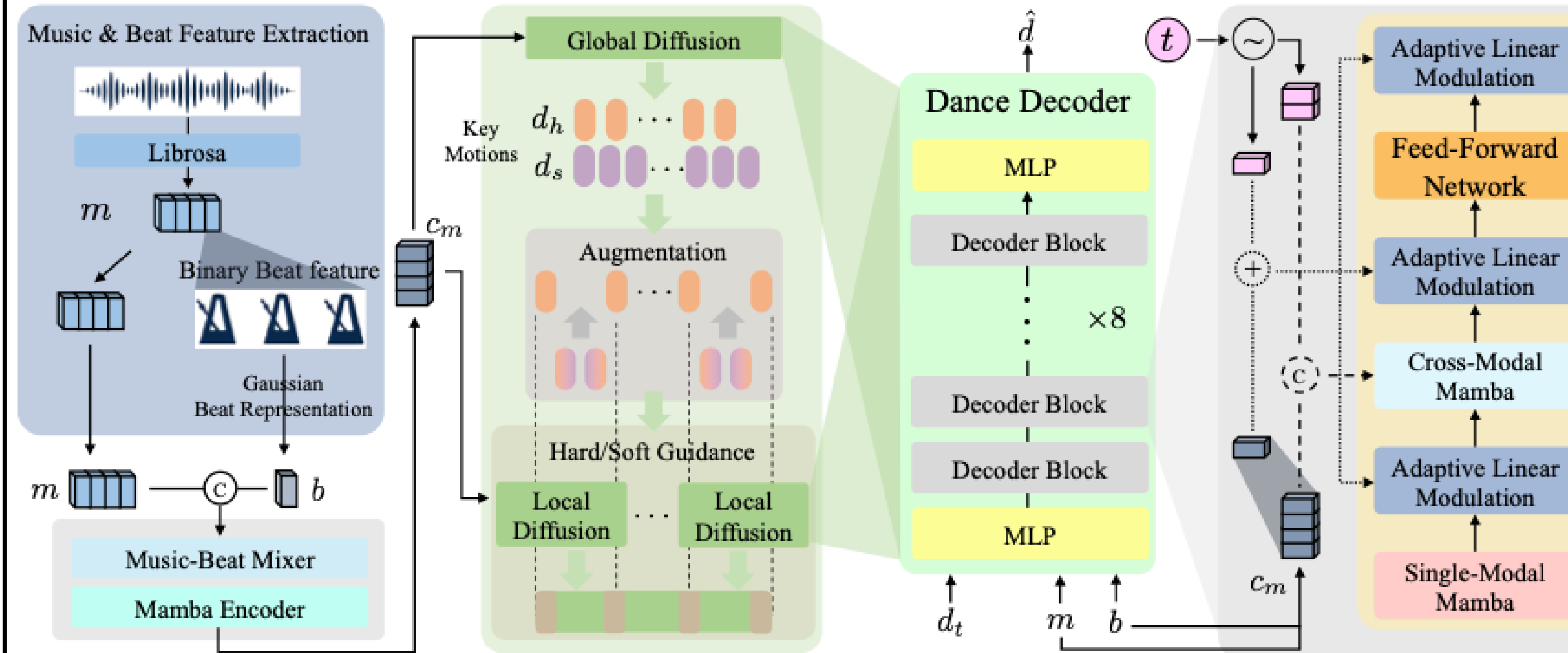
- Mamba's inductive bias helps modeling the autoregressive data
- Based on the previous Transformer-based two-stage diffusion framework<sup>[6]</sup>, we leverage Mamba for dance decoder

### ✓ H2. More **expressive beat representation** helps the generated dance better-aligned to the beat

- Properties of beats: (i) frames closer to beats carry stronger signals, (ii) this strength decays rapidly yet smoothly with temporal distance.
- We formulate previous simple beat representation NBD<sup>[3]</sup>, and use Gaussian decay function that can satisfy the properties

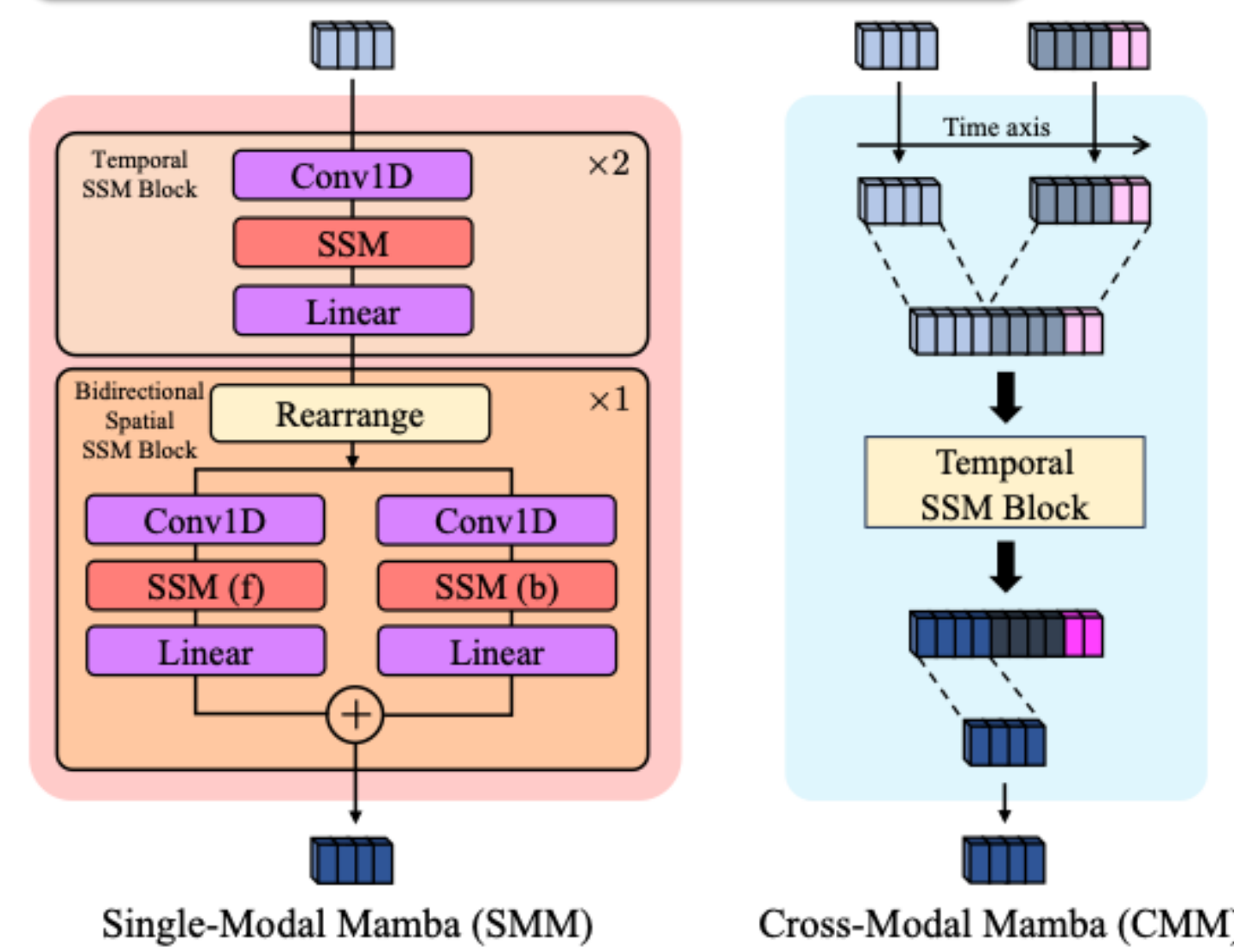
$$b(i) = \exp\left(-\frac{\text{NBD}(i)^2}{2(\alpha \cdot l(i))^2}\right)$$

## Method Overview



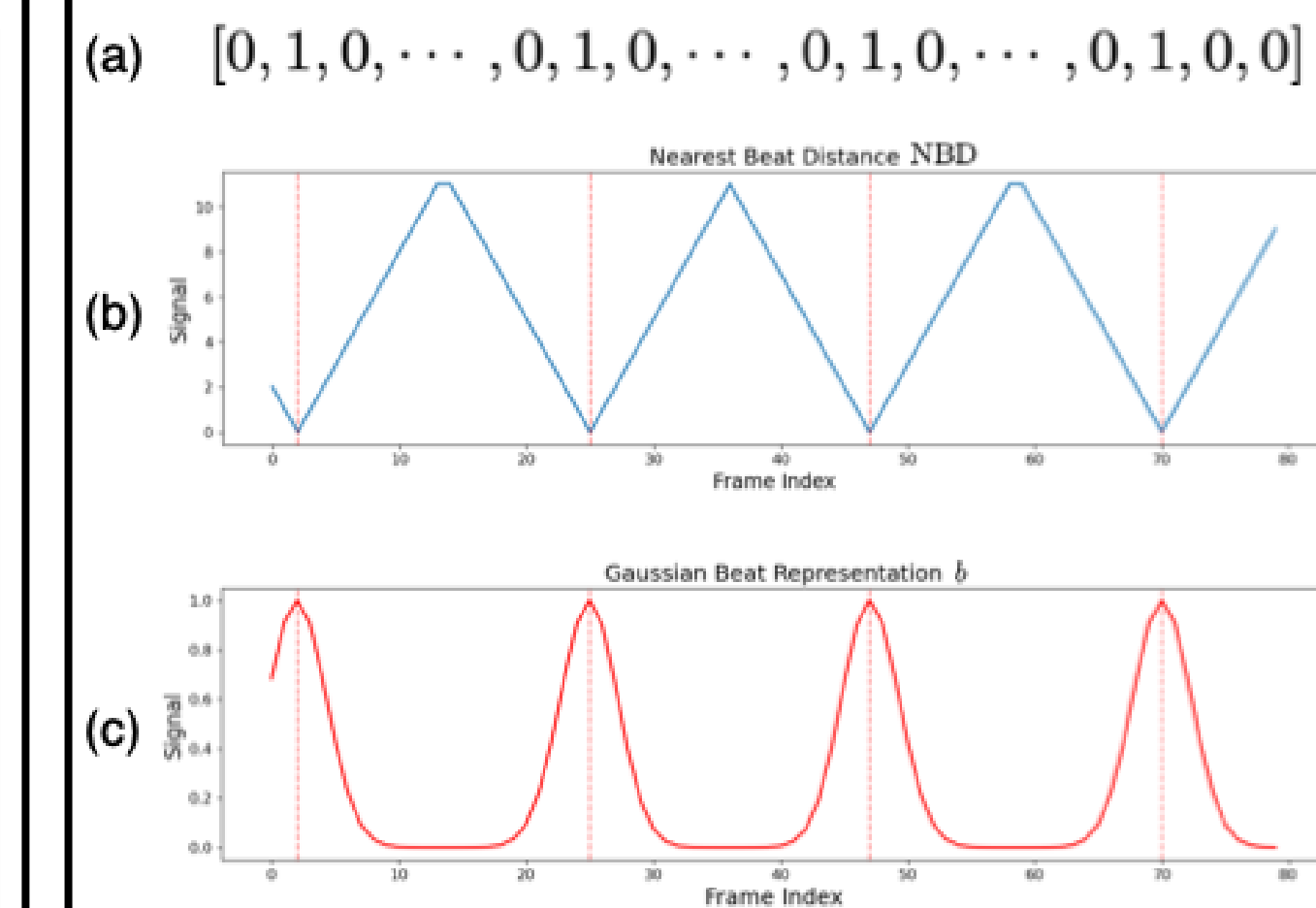
- Inspired by two-stage diffusion framework in Lodge<sup>[6]</sup>,
- Mamba-based Diffusion:  $d \approx \hat{d} = f_\theta(d_t, t, m, b)$
- Gaussian beat representation:  $m \in \mathbb{R}^{L \times 35} \rightarrow b_{\text{raw}} \in \mathbb{R}^{L \times 1} \rightarrow b \in \mathbb{R}^{L \times 1}$

## Mamba Modules



- **SMM** processes motion latent with SSM module of Mamba
- **CMM** processes concatenated sequence of motion latent, musical condition, and timestep embedding
- **AdaLM** is Adaptive Group Normalization (AdaGN<sup>[4]</sup>) for 1D sequences

## Gaussian Beat Representation



- (a) **1D raw beat**: 1-channel binary signal from music feature
- (b) **NBD**<sup>[3]</sup>: each entry denotes the temporal distance to the nearest beat frame
- (c) **Gaussian Beat Rep.**: Gaussian decay function provides explicit and interpretable rhythmic cues by locally stressing beat frames

## Experiments

### ✓ Quantitative Comparison

Dataset	Model	Fidelity			Beat	Diversity		Wins (↑)
		FID <sub>k</sub> (↓)	FID <sub>g</sub> (↓)	PFC (↓)	BAS (↑)	Div <sub>k</sub> (→)	Div <sub>g</sub> (→)	
FineDance <sup>[8]</sup>	GT	-	-	0.1852	-	10.9924	7.7424	-
	EDGE <sup>[2]</sup>	179.01 $\pm$ 3.10	1234.28 $\pm$ 236.95	0.3994 $\pm$ 0.0177	0.2261 $\pm$ 0.0013	10.34 $\pm$ 0.36	31.74 $\pm$ 2.76	3.5%
	POPDG <sup>[5]</sup>	190.23 $\pm$ 1.18	1479.14 $\pm$ 6.44	0.3765 $\pm$ 0.0124	0.2361 $\pm$ 0.0033	7.22 $\pm$ 0.14	14.53 $\pm$ 0.12	5.5%
	Lodge <sup>[6]</sup>	84.99 $\pm$ 2.07	64.57 $\pm$ 10.74	0.0585 $\pm$ 0.014	0.2410 $\pm$ 0.0063	7.98 $\pm$ 0.20	7.67 $\pm$ 0.68	35.0%
	<b>Ours</b>	<b>51.36<math>\pm</math>0.67</b>	<b>43.11<math>\pm</math>0.54</b>	<b>0.0119<math>\pm</math>0.0008</b>	<b>0.2441<math>\pm</math>0.0044</b>	6.38 $\pm$ 0.17	6.44 $\pm$ 0.88	<b>56.0%</b>
AIST++ <sup>[7]</sup>	GT	-	-	1.2544	-	9.61	7.78	-
	EDGE <sup>[2]</sup>	125.99 $\pm$ 128.69	28.72 $\pm$ 4.29	3.1883 $\pm$ 0.5318	0.2572 $\pm$ 0.0112	11.45 $\pm$ 3.25	4.91 $\pm$ 0.56	10.5%
	POPDG <sup>[5]</sup>	777.32 $\pm$ 711.65	60.08 $\pm$ 5.98	4.8615 $\pm$ 0.6010	0.2318 $\pm$ 0.0129	24.08 $\pm$ 7.40	7.87 $\pm$ 0.59	9.0%
	Lodge <sup>[6]</sup>	67.13 $\pm$ 2.79	28.93 $\pm$ 0.47	1.4087 $\pm$ 0.1296	0.2397 $\pm$ 0.0158	3.34 $\pm$ 0.32	3.54 $\pm$ 0.14	32.0%
	<b>Ours</b>	<b>65.86<math>\pm</math>3.11</b>	<b>26.58<math>\pm</math>1.02</b>	<b>1.0622<math>\pm</math>0.2343</b>	<b>0.2701<math>\pm</math>0.0116</b>	3.57 $\pm$ 0.37	4.98 $\pm$ 0.43	<b>48.5%</b>

### ✓ Qualitative Result (Details in VIDEO!!!)



## References

[1] Albert Gu and Tri Dao, "Mamba: Linear-time Sequence Modeling with Selective State Spaces." ArXiv 2023.  
 [2] Jonathan Tseng, et al., "EDGE: Editable Dance Generation from Music." CVPR 2023.  
 [3] Zikai Huang, et al., "Beat-It: Beat-Synchronized Multi-Condition 3D Dance Generation." ECCV 2024.  
 [4] Prafulla Dhariwal and Alexander Nichol, "Diffusion Models Beat GANs on Image Synthesis." NeurIPS 2021.

[5] Zhenye Luo, et al., "POPDG: Popular 3D Dance Generation with PopDanceSet." CVPR 2024.  
 [6] Ronghui Li, et al., "Lodge: A Coarse to Fine Diffusion Network for Long Dance Generation Guided by the Characteristic Dance Primitives." CVPR 2024.  
 [7] Ruilong Li, et al., "AI Choreographer: Music Conditioned 3D Dance Generation with AIST++." ICCV 2021.  
 [8] Ronghui Li, et al., "FineDance: A Fine-grained Choreography Dataset for 3D Full Body Dance Generation." ICCV 2023.