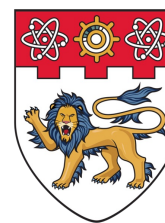


MotionMix: Weakly-Supervised Diffusion for Controllable Motion Generation

Nhat M. Hoang¹, Kehong Gong², Chuan Guo², Michael Bi Mi²

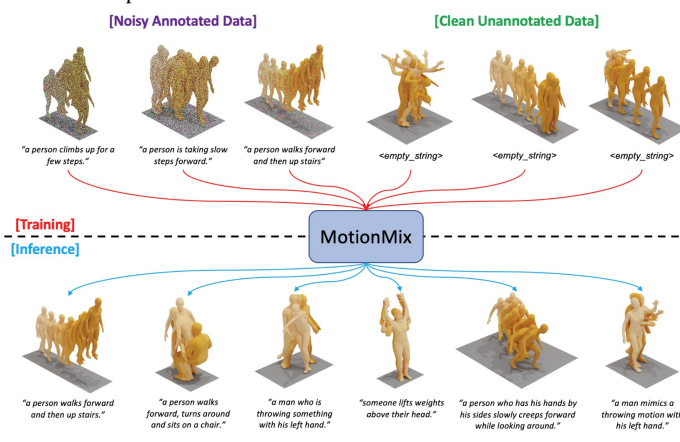
¹Nanyang Technological University, ²Huawei Technologies Co., Ltd.



Scan for motion videos

Overview

- Controllable generation of 3D human motions becomes an important topic as the world embraces digital transformation.
- Existing methods heavily rely on costly, annotated high-quality motion data.
- We propose MotionMix, a simple yet effective weakly-supervised approach for diffusion model to utilize both noisy and unannotated motion sequences.



Setup

Motion generation tasks and benchmarks:

- Text-to-Motion (T2M):** HumanML3D and KIT-ML datasets
- Action-to-Motion (A2M):** HumanAct12 and UESTC datasets
- Music-to-Dance (M2D):** AIST++ dataset

MotionMix was applied to *sota* diffusion model of different designs:

- MDM** for T2M, A2M: x0-parameterization, classifier-free guidance.
- MotionDiffuse** for T2M: epsilon-parameterization.
- EDGE** for M2D: x0-parameterization with classifier-free guidance.

Ablation on Denoising Pivot

→ More clean data does not lead to better performance while more annotated data

Method	R Precision (top 3)↑	FID↓	Multimodal Dist.↓	Diversity→	Multimodality↑
Real Motion	0.797±.002	0.002±.000	2.974±.008	9.503±.065	-
MDM	0.611±.007	0.544±.440	5.566±.027	9.559±.860	2.799±.072
<i>50% noisy, T₁=20, T₂=60</i>					
MDM (MotionMix) (T [*] =0)	0.598±.006	0.714±.045	5.503±.036	9.750±.123	3.044±.054
MDM (MotionMix) (T [*] =20)	0.601±.005	0.497±.048	5.562±.026	9.414±.092	2.935±.059
MDM (MotionMix) (T [*] =40)	0.604±.008	0.402±.032	5.524±.033	9.396±.094	2.747±.070
MDM (MotionMix) (T [*] =60)	0.632±.006	0.381±.042	5.325±.026	9.520±.090	2.718±.019
MDM (MotionMix) (T [*] =80)	0.594±.005	0.589±.059	5.670±.033	9.242±.086	2.602±.057

Ablation Studies

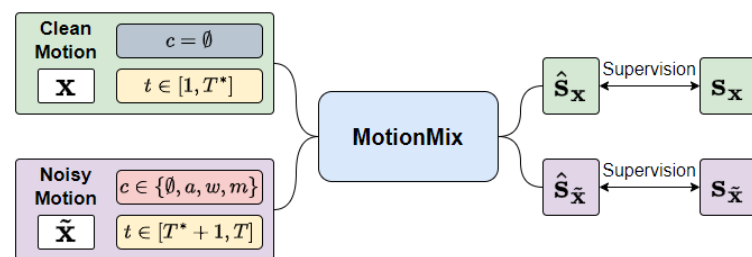
Ablation on Noisy Ratio

→ More clean data ≠ better, more annotated data (even noisy) = better

Method	R Precision (top 3)↑	FID↓	Multimodal Dist.↓	Diversity→	Multimodality↑
Real Motion	0.797±.002	0.002±.000	2.974±.008	9.503±.065	-
MDM	0.611±.007	0.544±.440	5.566±.027	9.559±.860	2.799±.072
<i>T₁=20, T₂=60, T[*]=60</i>					
MDM (MotionMix) (30% noisy)	0.601±.007	0.898±.045	5.581±.030	9.080±.092	2.856±.074
MDM (MotionMix) (50% noisy)	0.632±.006	0.381±.042	5.325±.026	9.520±.090	2.718±.019
MDM (MotionMix) (70% noisy)	0.615±.006	0.359±.030	5.545±.031	9.457±.098	2.867±.107

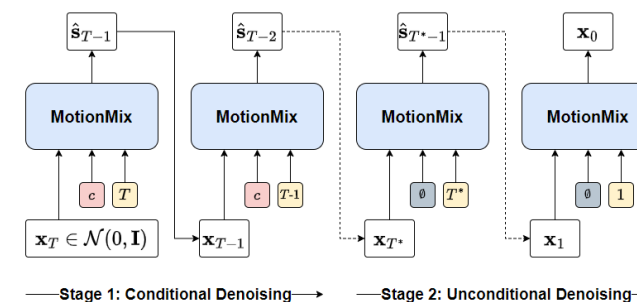
Methodology

(a) Training: Separate the diffusion steps into two distinct ranges regarding the data type. Clean samples are guided by an empty condition. Ground-truth s_x refers to both x-prediction and epsilon-prediction, two fundamental objective of diffusion models.



(b) Inference (two stages):

- Stage 1: generate rough motion guided by the conditional input.
- Stage 2: refine these rough motions to high-quality ones while the conditional input is masked.



Evaluation Results

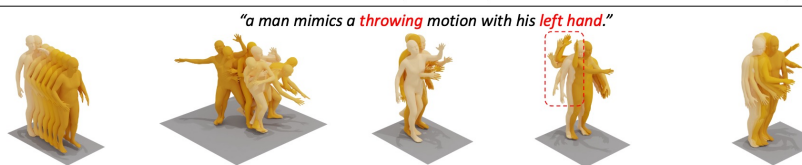
Text-to-Motion

Method	R Precision (top 3)↑	FID↓	Multimodal Dist.↓	Diversity→	Multimodality↑
ML3D					
Real Motion	0.797±.002	0.002±.000	2.974±.008	9.503±.065	-
MDM	0.611±.007	0.544±.440	5.566±.027	9.559±.860	2.799±.072
MDM (MotionMix)	0.632±.006 (↑3.4%)	0.381±.042 (↑30.0%)	5.325±.026 (↑4.3%)	9.520±.090 (↑69.6%)	2.718±.019 (↓2.9%)
KIT-ML					
Real Motion	0.779±.006	0.031±.004	2.788±.012	11.080±.097	-
MDM	0.396±.004	0.497±.021	9.191±.022	10.847±.109	1.907±.214
MDM (MotionMix)	0.404±.005 (↑2.0%)	0.322±.020 (↑35.2%)	9.068±.019 (↑1.3%)	10.781±.098 (↓28.3%)	1.946±.019 (↑2.0%)
MotionDiffuse					
Real Motion	0.739±.004	1.954±.062	2.958±.005	11.100±.143	0.730±.013
MotionDiffuse (MotionMix)	0.742±.005 (↑0.4%)	1.192±.073 (↑39.0%)	3.066±.018 (↓3.6%)	10.998±.072 (↓310%)	1.391±.111 (↑90.5%)

Action-to-Motion

Method	FID ↓	Accuracy ↑	Diversity →	Multimodality →
HumanAct12				
Real Motion	0.053±.003	0.995±.001	6.835±.045	2.604±.040
Action2Motion	0.338±.015	0.917±.001	6.850±.050	2.511±.023
ACTOR	0.130±.000	0.955±.008	6.840±.030	2.530±.020
MLD	0.077±.004	0.964±.002	6.831±.050	2.824±.038
UESTC				
Real Motion	0.100±.000	0.990±.000	6.860±.050	2.520±.010
MDM (MotionMix)	0.196±.007 (↓96%)	0.930±.003 (↓6.1%)	6.836±.062 (↓96%)	3.043±.054 (↓22.6%)

Qualitative samples from HumanML3D test set. Please view videos on our webpage.



Qualitative samples from HumanML3D test set. Please view videos on our webpage.

Music-to-Dance

Method	PFCC ↓	Beat Align. ↑	Dist _c →	Dist _g →
Real Motion	1.380	0.314	9.545	7.766
Bailando	1.754	0.23	10.58	7.72
FACT	2.2543	0.22	10.85	6.14
EDGE				
EDGE†	1.605±.224	0.224±.025	5.549±.783	4.831±.752
EDGE (MotionMix)	1.988±.120 (↓21.32%)	0.256±.013 (↑13.3%)	10.103±2.039 (↓95.0%)	6.595±.173 (↓15.1%)

Ablation on Noisy Range

→ MotionMix is robust on different levels of corrupted motions

Method	R Precision (top 3)↑	FID↓	Multimodal Dist.↓	Diversity→	Multimodality↑
Real Motion	0.797±.002	0.002±.000	2.974±.008	9.503±.065	-
MDM	0.611±.007	0.544±.440	5.566±.027	9.559±.860	2.799±.072
<i>50% noisy, T[*] = T₂</i>					
MDM (MotionMix) (T ₁ =20, T ₂ =40)	0.616±.006	0.451±.033	5.459±.027	9.585±.101	2.585±.076
MDM (MotionMix) (T ₁ =20, T ₂ =60)	0.632±.006	0.381±.042	5.325±.026	9.520±.090	2.718±.019
MDM (MotionMix) (T ₁ =20, T ₂ =80)	0.604±.004	0.614±.000	5.540±.024	9.554±.104	2.768±.098
<i>50% noisy, T[*] = T₁</i>					
MDM (MotionMix) (T ₁ =10, T ₂ =30)	0.592±.008	0.713±.048	5.634±.028	9.567±.109	2.783±.130
MDM (MotionMix) (T ₁ =10, T ₂ =40)	0.616±.006	0.451±.033	5.459±.027	9.585±.101	2.585±.076
MDM (MotionMix) (T ₁ =10, T ₂ =60)	0.598±.004	0.554±.076	5.600±.031	9.479±.100	2.815±.094
MDM (MotionMix) (T ₁ =10, T ₂ =80)	0.597±.008	0.437±.039	5.554±.033	9.452±.092	2.895±.079