



NeRM: Learning Neural Representations for High-Framerate Human Motion Synthesis

Dong Wei¹, Huaijiang Sun¹, Bin Li², Xiaoning Sun¹, Shengxiang Hu¹, Weiqing Li¹, Jianfeng Lu¹

1. School of Computer Science and Engineering, Nanjing University of Science and Technology

2. Tianjin AiForward Science and Technology Co., Ltd

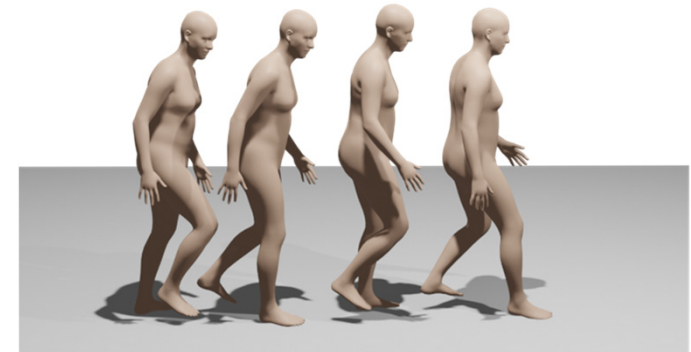
Problem Definition

Given an arbitrary condition c , such as:

- text description
- class label
- ...

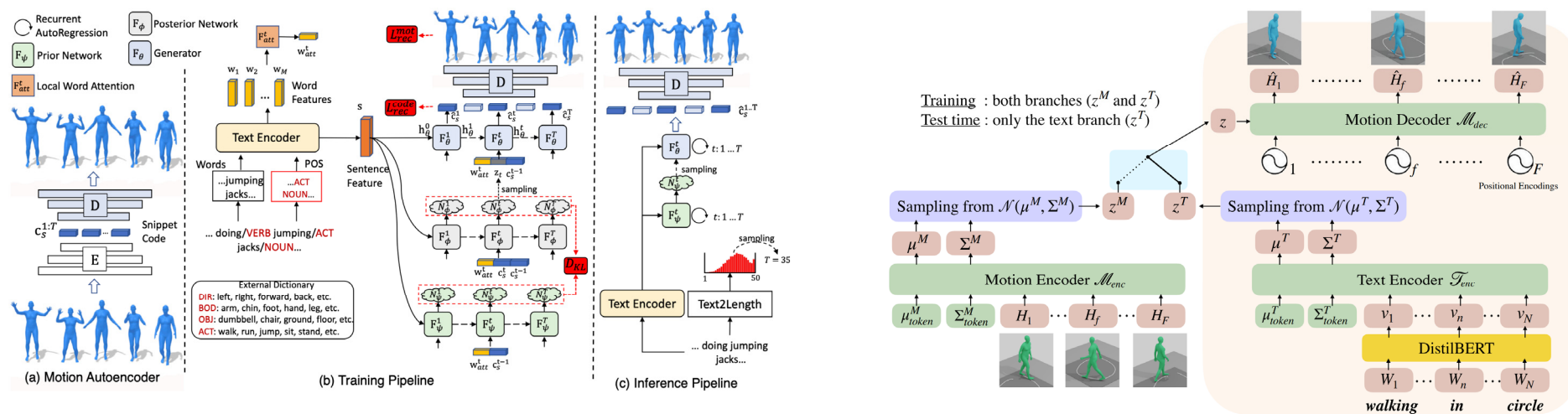
Goal: Generate plausible and diverse human motions.

Text: *“A person walks to the right slowly.”*



Related Work

T2M (Guo et al. 2022) & TEMOS (Petrovich et al. 2022): a conditional VAE model



Problem:

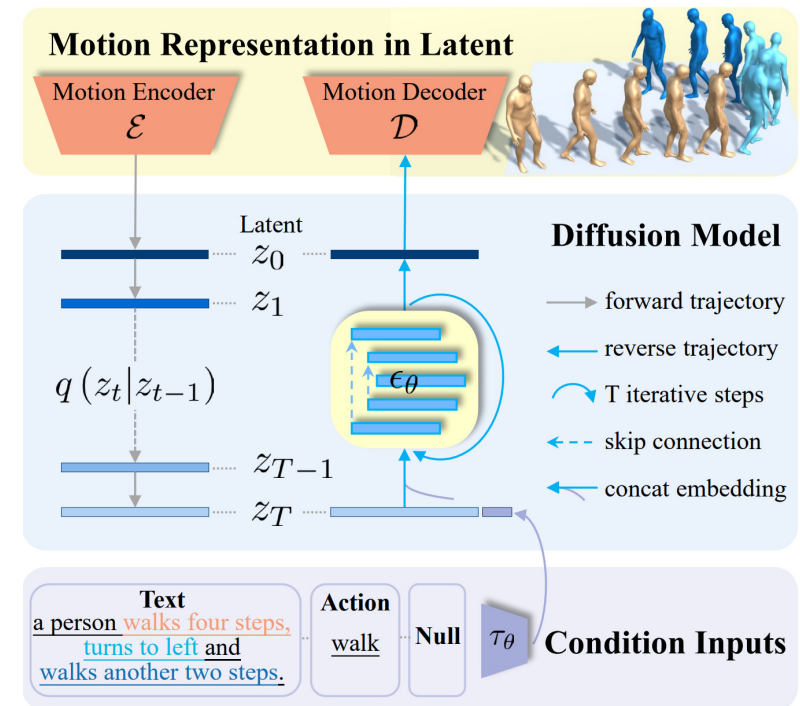
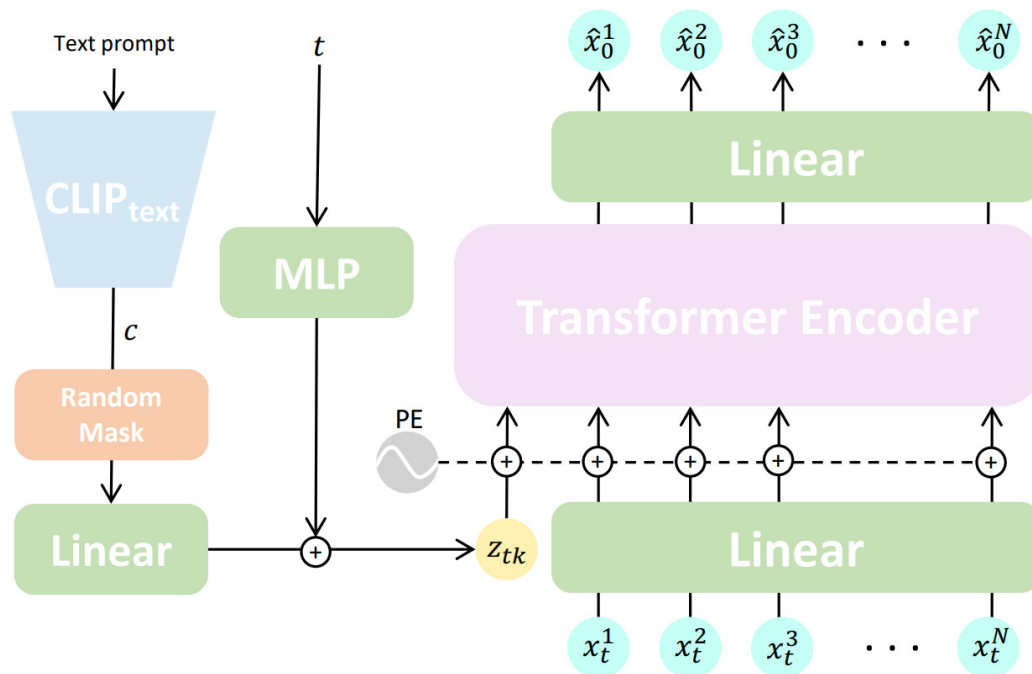
- VAE brings about posterior collapse with a powerful decoder \rightarrow not diverse enough
- Limited generative capability \rightarrow not plausible enough

Chuan Guo, Zhihao Zou, Xinxin Zuo, et al. Generating diverse and natural 3D human motions from text. CVPR, 2022.

Mathis Petrovich, Michael J. Black, Gul Varol. TEMOS: Generating diverse human motions from textual descriptions. ECCV, 2022.

Related Work

MDM (Tevet et al. 2023) & MLD (Chen et al. 2023): a diffusion model



Address the limitations of VAE models.

Guy Tevet, Sigal Raab, Brian Gordon, et al. Human motion diffusion model. ICLR, 2023.

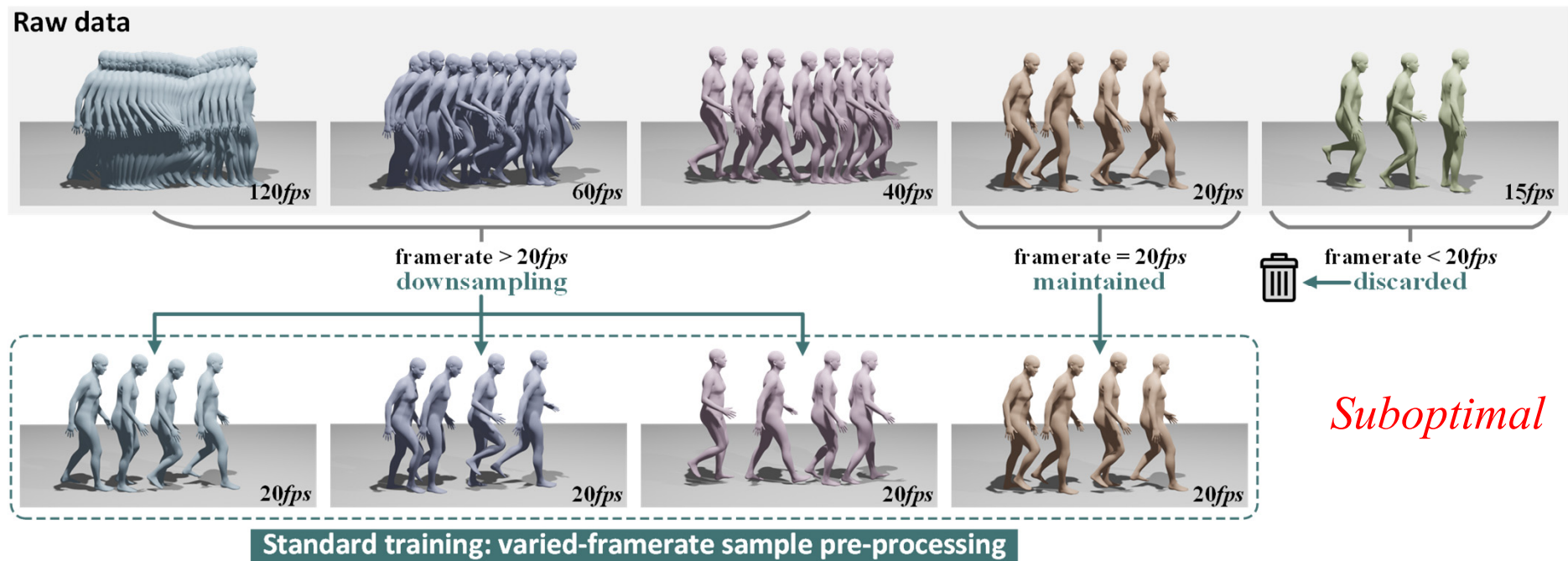
Xin Chen, Biao Jiang, Wen Liu, et al. Executing your commands via motion diffusion in latent space. CVPR, 2023.

Challenge

Existing methods cannot support high-framerate motion generation, due to:

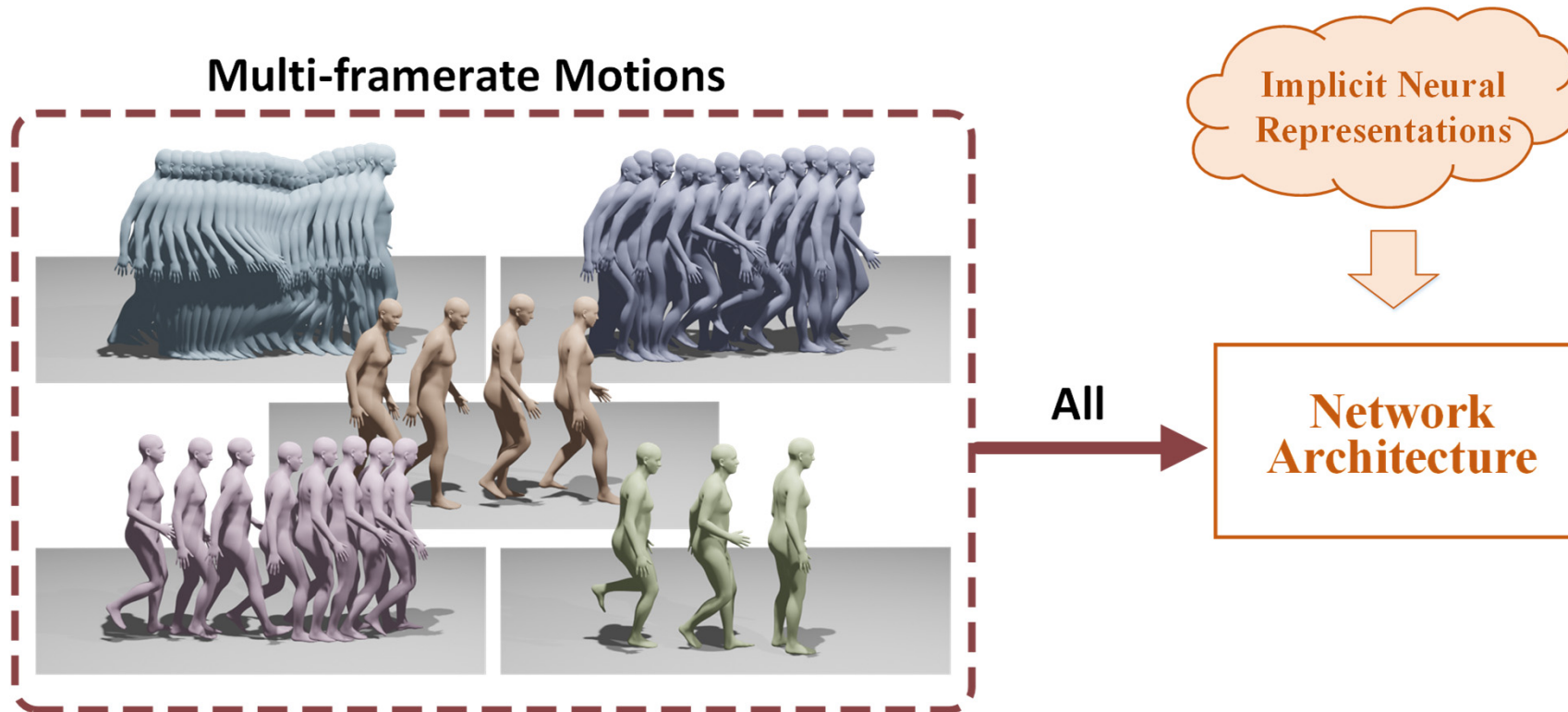
- Memory overloads
- Varied-size training dataset

Therefore, they employ a preprocessing step.

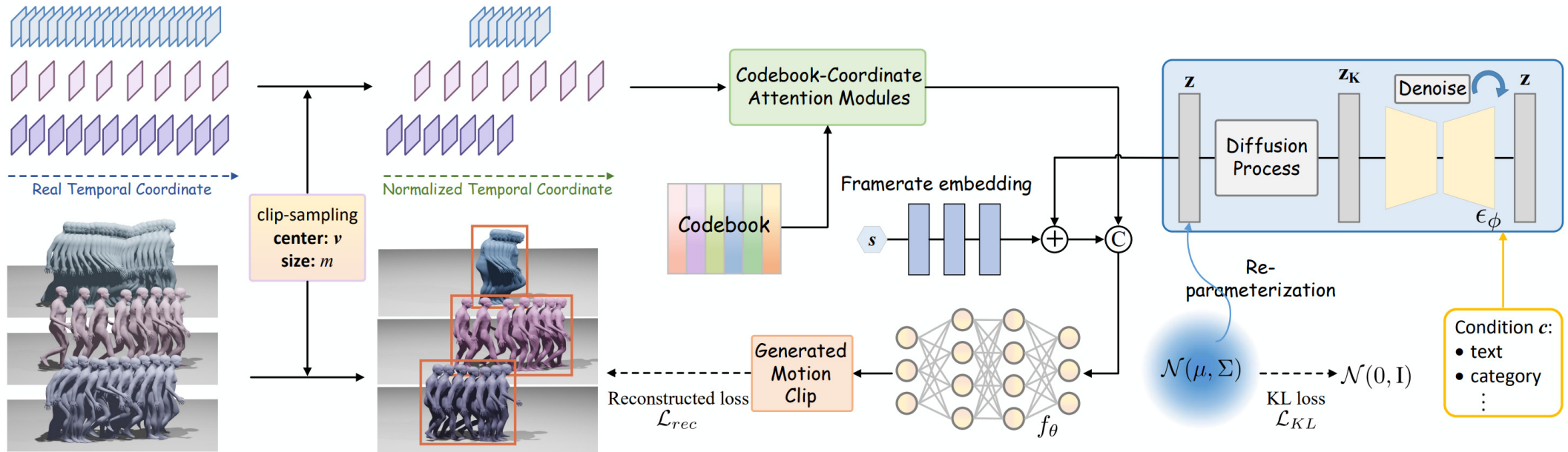


Key Insight

It is desired to embrace the natural diversity of motion framerates and process them at their **native** framerate.



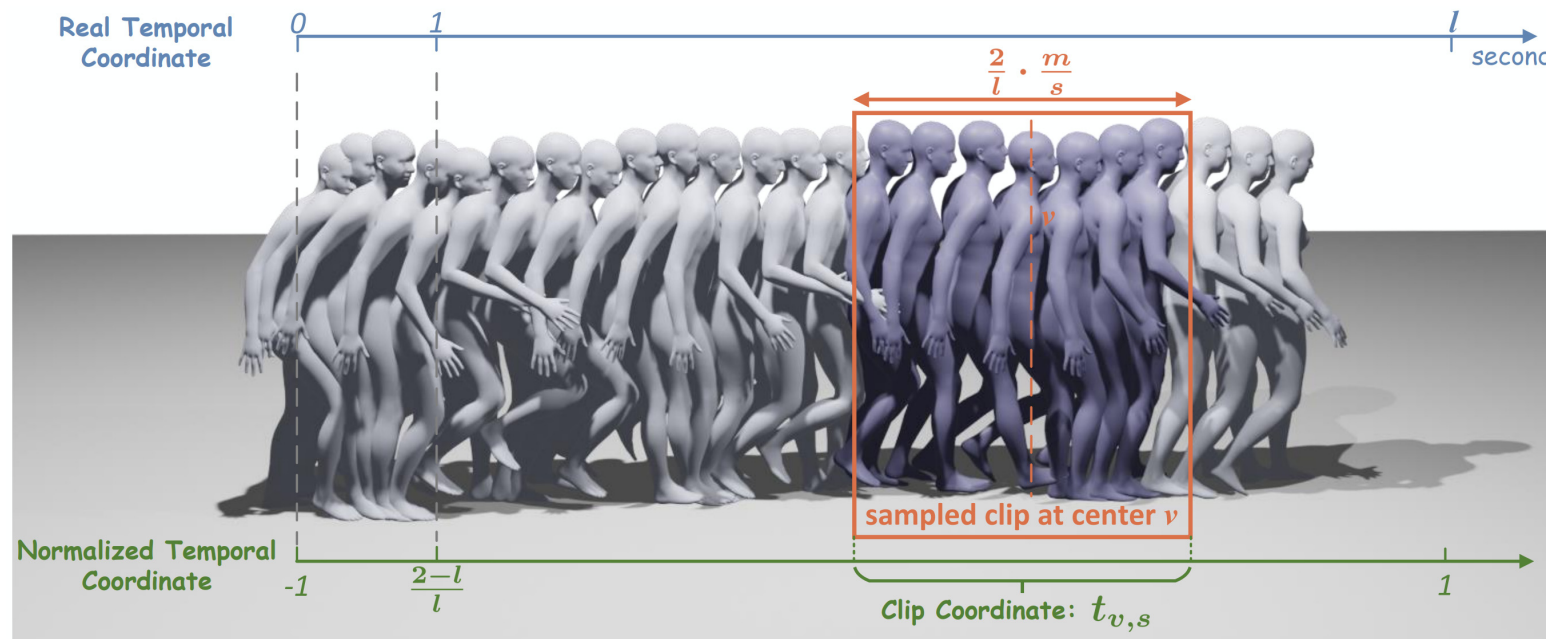
Overview



Two-stage pipeline of our NeRM. In the first stage (left), we sample motion clips at random framerates from full-size motion sequences for training. The second stage (right) takes the latent variables as input to our diffusion model and can be guided by various conditions.

Multi-framerate Training

To learn from multi-framerate motions in training sets, we design our approach by exploiting the temporal consistency in motions to generate motion *clips*.



Loss Function

The first stage: Variational implicit neural representation

$$\mathcal{L}^i = \mathcal{L}_{rec}^i + \lambda_{KL} \mathcal{L}_{KL}^i = \|\hat{x}_{clip}^i - x_{clip}^i\|^2 + \lambda_{KL} D_{KL}(\mathcal{N}(\mu_i, \Sigma_i) \| p(z))$$

The first stage: Model parameters optimization

$$\{(\mu_i^*, \Sigma_i^*)\}_{i=1,2,\dots,n} = \arg \min_{\mu_i, \Sigma_i} \mathcal{L}^i, \quad \theta^* = \arg \min_{\theta} \sum_{i=1}^n \min_{\mu_i, \Sigma_i} \mathcal{L}^i$$

The second stage: Motion diffusion in latent space

$$\min_{\phi} \mathbb{E}_{k, \mathbf{z} \sim Z, \epsilon \sim \mathcal{N}(0, \mathbf{I})} [\|\epsilon - \epsilon_{\phi}(\sqrt{\bar{\alpha}_k} \mathbf{z} + \sqrt{1 - \bar{\alpha}_k} \epsilon, k)\|^2]$$

The second stage: Classifier-free guidance

$$\epsilon_{\phi}(\mathbf{z}_k, k, c) = r \epsilon_{\phi}(\mathbf{z}_k, k, c) + (1 - r) \epsilon_{\theta}(\mathbf{z}_k, k, \emptyset)$$

Experiments

Text-to-Motion

Method	HumanML3D (Guo et al., 2022)					KIT (Plappert et al., 2016)				
	FID ↓	R-Precision (Top-3) ↑	Multimodal Dist ↓	Diversity →	MM ↑	FID ↓	R-Precision (Top-3) ↑	Multimodal Dist ↓	Diversity →	MM ↑
Real	0.002±.000	0.797±.002	2.974±.008	9.503±.065	-	0.031±.004	0.779±.006	2.788±.012	11.08±.097	-
JL2P (Ahuja & Morency, 2019)	11.02±.046	0.486±.002	5.296±.008	7.676±.058	-	6.545±.072	0.483±.005	5.147±.030	9.073±.100	-
Hier (Ghosh et al., 2021)	6.532±.024	0.552±.004	5.012±.018	8.332±.042	-	5.203±.107	0.531±.007	4.986±.027	9.563±.072	-
T2M (Guo et al., 2022)	1.067±.002	0.740±.003	3.340±.008	9.188±.002	2.090±.083	2.770±.109	0.693±.007	3.401±.008	<u>10.91</u> ±.119	1.482±.065
MoFusion (Dabral et al., 2023)	-	0.492	-	8.82	<u>2.521</u>	-	-	-	-	-
MDM (Tevet et al., 2023)	0.544±.044	0.611±.007	5.566±.027	<u>9.559</u> ±.086	2.799 ±.072	0.497±.021	0.396±.004	9.191±.022	10.85±.109	1.907±.214
PhysDiff (Yuan et al., 2023)	<u>0.433</u>	0.631	-	-	-	-	-	-	-	-
MLD (Chen et al., 2023)	0.473±.013	0.772±.002	3.196±.010	9.724±.082	2.413±.079	0.404 ±.027	<u>0.734</u> ±.007	<u>3.204</u> ±.027	10.80±.117	2.192 ±.071
NeRM (fixed-framerate train)	0.489±.013	<u>0.774</u> ±.003	<u>3.186</u> ±.015	9.692±.082	2.330±.075	0.522±.029	0.727±.005	3.843±.027	10.90±.094	<u>1.931</u> ±.133
NeRM (native-framerate train)	0.389 ±.011	0.779 ±.003	3.178 ±.016	9.547 ±.073	2.193±.081	0.472±.019	0.736 ±.007	3.189 ±.031	10.94 ±.114	1.785±.082

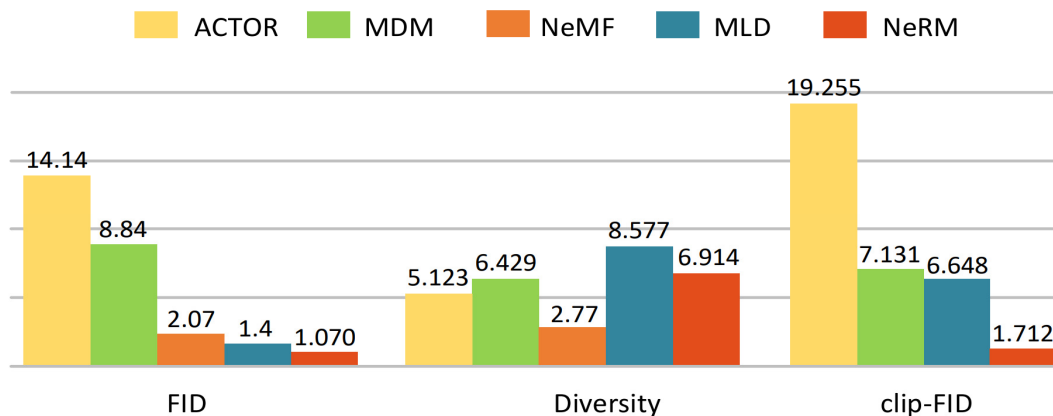
- SOTA FID
- Using original multi-framerate datasets can achieve much better performance.

Experiments

Action-to-Motion

Method	UESTC (Ji et al., 2018)					HumanAct12 (Guo et al., 2020)			
	FID _{train} ↓	FID _{test} ↓	Accuracy ↑	Diversity →	MM →	FID _{train} ↓	Accuracy ↑	Diversity →	MM →
Real	2.92 \pm .26	2.79 \pm .29	0.988 \pm .001	33.34 \pm .320	14.16 \pm .06	0.020 \pm .010	0.997 \pm .001	6.850 \pm .050	2.450 \pm .040
ACTOR (Petrovich et al., 2021)	20.5 \pm 2.3	23.43 \pm 2.20	0.911 \pm .003	31.96 \pm .33	14.52 \pm .09	0.120 \pm .000	0.955 \pm .008	6.840 \pm .030	2.530 \pm .020
MDM (Tevet et al., 2023)	<u>9.98</u> \pm 1.33	12.81 \pm 1.46	0.950 \pm .000	33.02 \pm .28	14.26 \pm .12	<u>0.100</u> \pm .000	0.990 \pm .000	6.680 \pm .050	2.520 \pm .010
MLD (Chen et al., 2023)	12.89 \pm .109	15.79 \pm .079	<u>0.954</u> \pm .001	<u>33.52</u> \pm .14	<u>13.57</u> \pm .06	0.077 \pm .004	0.964 \pm .002	<u>6.831</u> \pm .050	2.824 \pm .038
INR-MLP (Cervantes et al., 2022)	9.55 \pm .06	15.00 \pm .09	0.941 \pm .001	31.59 \pm .19	14.68 \pm .07	0.114 \pm .001	0.970 \pm .001	6.786 \pm .057	<u>2.507</u> \pm .034
NeRM (Ours)	11.75 \pm .31	<u>14.23</u> \pm .174	0.956 \pm .001	33.20 \pm .21	14.41 \pm .06	0.106 \pm .000	<u>0.977</u> \pm .001	6.866 \pm .032	2.492 \pm .048

Unconditional Generation



Experiments

Ablation Study

Effectiveness of time encoding

Simple	Codebook	Motion Reconstruction (MRE)					Motion Synthesis (clip-FID)				
		20	40	60	100	120	20	40	60	100	120
\times	\times	0.134	0.141	0.164	0.094	0.091	0.471	0.803	1.070	1.769	2.944
\checkmark	\times	0.053	0.049	0.039	0.057	0.043	0.397	0.519	0.701	0.142	1.717
\times	\checkmark	0.041	0.036	0.034	0.035	0.038	0.389	0.493	0.680	0.903	1.315

Effectiveness of variational INRs

Variational INRs	Motion Reconstruction (MRE)					Motion Synthesis (clip-FID)				
	20	40	60	100	120	20	40	60	100	120
\times	0.032	0.030	0.031	0.036	0.027	1.280	2.924	7.012	10.482	14.654
\checkmark	0.041	0.036	0.034	0.035	0.038	0.389	0.493	0.680	0.903	1.315

Experiments

Ablation Study

Effectiveness of different framerates

Method	20	40	60	100	120
T2M	1.067	2.831	6.442	9.182	11.264
MDM	0.544	1.882	4.605	5.966	8.401
MLD	0.473	1.465	3.816	5.138	7.878
NeRM	0.389	0.493	0.680	0.903	1.315

Effectiveness of time normalization

Time Normalization	MRE@20fps	FID	Diversity
✗	0.118	0.958	9.892
✓	0.041	0.389	9.547

Experiments

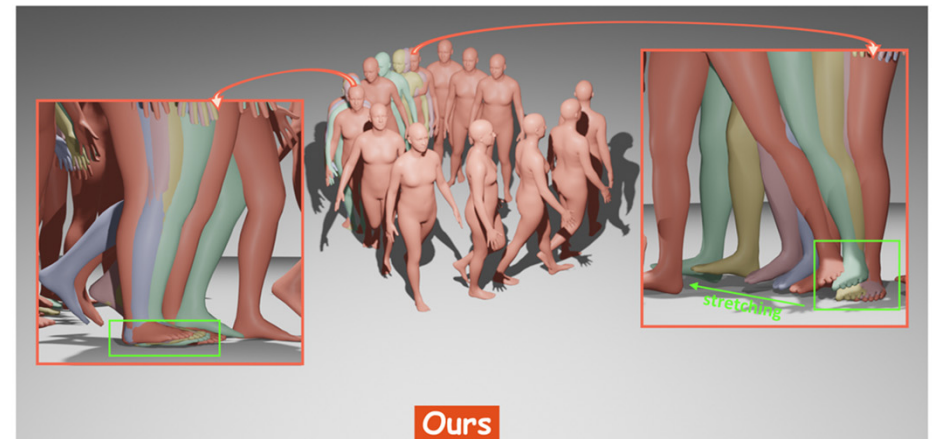
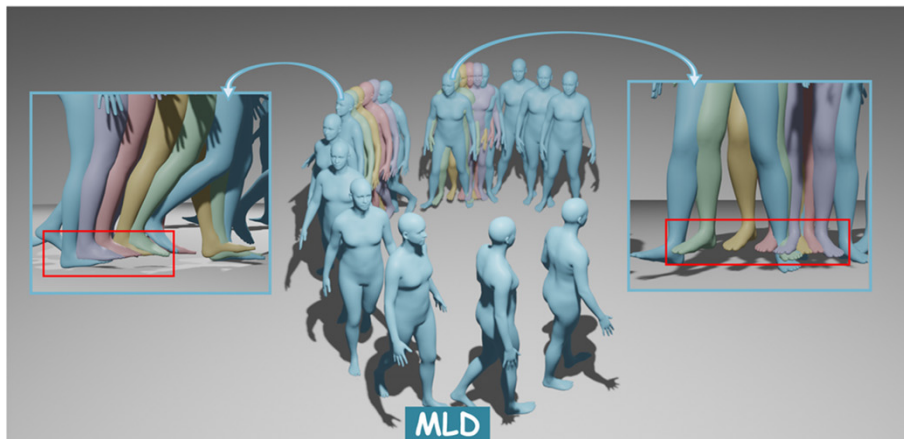
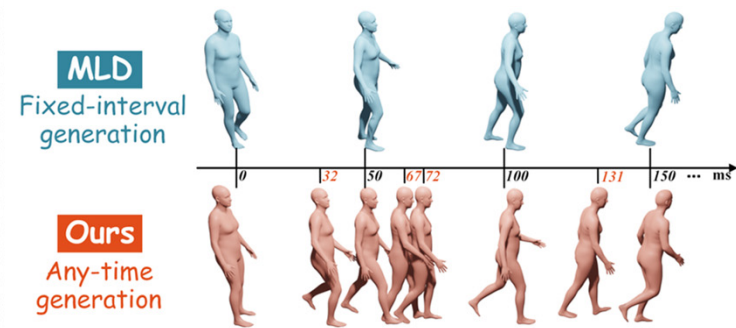
Visualization

(a) Conventional generation

Text: Person walks up then takes a large step to their left and then goes back onto the same path they were on.



(c) Generation at any time steps



Text: This person shuffles right then walks forward and circles back.

(b) High-framerate details

Thank You!