KINEMADIFF: TOWARDS DIFFUSION FOR COHERENT AND PHYSICALLY PLAUSIBLE HUMAN MOTION PRE-DICTION

Anonymous authors

000

001

002 003

004

006

008 009 010

011 012

013

014

015

016

018

019

021

023

024

025

026

027

028

029

031

032

034

035

037

040

041

042

043

044

046

047

048

051

052

Paper under double-blind review

ABSTRACT

Stochastic Human Motion Prediction (HMP) has become an essential task for the realm of computer vision, for its capacity to anticipate accurate and diverse future human trajectories. Current diffusion-based techniques typically enforce skeletal consistency by encoding structural priors into network architectures. Although effective in promoting plausible kinematics, this approach provides only indirect control over the generative process and often fails to guarantee strict physical constraint satisfaction. In this work, we propose a structure-aligned and joint-aware diffusion framework that enforces physical constraints by embedding skeletal topology and joint-specific dynamics directly into the diffusion process. Specifically, our framework consists of two key modules, the Joint-Adaptive Noise Generator and the Structure-Aligned Constraint Enforcer. The former component, Joint-Adaptive Noise Generator, infers joint-specific dynamics and injects heterogeneous, instance-aware noise per joint and sample to capture spatial variability and enhance motion diversity. The latter component, Structure-Aligned Constraint Enforcer, encodes skeletal topology by modeling joint connectivity and bone lengths from historical motions, and it constrains each denoising step to preserve anatomical consistency. Through their synergistic operation, these modules grant KinemaDiff direct control over physical realism and motion diversity, addressing the common limitations of indirect structural priors and uniform noise application. Extensive experiments on multiple benchmarks demonstrate the effectiveness of our method, attributable to tailoring the diffusion process through structural alignment and joint-adaptive noise modeling.

Introduction

Human Motion Prediction (HMP) (Barsoum et al., 2018) aims to forecast future human motion sequences based on past observations, which is crucial for applications like autonomous driving (Paden et al., 2016), assistive robotics (Gui et al., 2018), and virtual avatars. While early deterministic methods (Xu et al., 2023; Li et al., 2022; Ma et al., 2022) sought (a) Current: take the same noise distribution to predict a single most likely future, they fell short in capturing the inherent unpredictability of human actions. Therefore, accurately modeling the multimodal and physically plausible distribution of future motions emerges as a paramount yet daunting task. To address this challenge, stochastic methods have gained prominence, with denoising diffusion probabilistic models becoming the mainstream approach (Yuan & Kitani, 2020; Dang et al., 2022). These models demonstrate remarkable capabilities in generating diverse motion sequences by progressively refining random noise into coherent human poses, as Fig. 1.

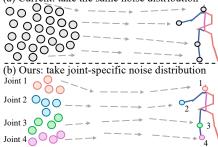


Figure 1: Comparisons about the noise taken between other baselines and ours.

Despite their generative power, these diffusion-based methods face two critical technical limitations in producing coherent and physically realistic human poses throughout the iterative diffusion process. On the one hand, a uniform noise schedule is typically applied across all human joints, failing to account for their heterogeneous motion patterns. Different joints exhibit vastly different degrees

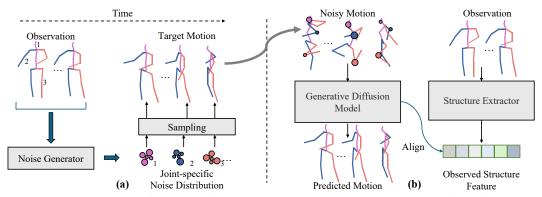


Figure 2: (a) Illustration of joint-adaptive noise generator. We propose a joint-adaptive noise determined by both the joint characteristics and the human motion observations. These noises are then added to the corresponding human motions to be predicted. (b) Representation of structure-aligned constraint enforcer, which identifies the human structural information from the historical motion and leverages the extracted structure to guide the motion generation during the diffusion process.

of freedom and dynamic behaviors. Applying identical noise profiles overlooks these unique kinematic properties, resulting in disordered or physically invalid predictions that compromise motion diversity and realism. On the other hand, prior methods tend to neglect the human skeleton's inherent anatomical structure. They often rely on implicitly learning the structural constraints (Chen et al., 2023; Sun & Chowdhary, 2024) or post-processing corrections (Wei et al., 2023), without integrating these constraints into the diffusion process, which leads to the generation of physically implausible poses with stretched or compressed bones, undermining motion realism, as in Fig. 2.

To address the aforementioned limitations, we present KinemaDiff, a novel kinematics-aware diffusion framework that fundamentally reshapes the denoising diffusion process, via explicitly embedding anatomical consistency and kinematic heterogeneity. Our framework consists of two core modules: the Joint-Adaptive Noise Generator and the Structure-Aligned Constraint Enforcer, as illustrated in Fig.2. The first component, the Joint-Adaptive Noise Generator, is responsible for capturing and injecting heterogeneous motion patterns. It learns and applies instance-specific, heterogeneous noise profiles to different joints, effectively adapting the noise characteristics based on their unique dynamics and varied degrees of freedom, thereby guiding the diffusion process to generate dynamically rich and realistic motions. Subsequently, the Structure-Aligned Constraint Enforcer, is tasked with rigorously enforcing anatomical consistency throughout the generative process. It achieves this by directly integrating bone length constraints into the denoising procedure, leveraging stable structural features extracted from historical motion observations to ensure that generated poses adhere strictly to human biomechanics. Through the synergistic operation of these two modules, KinemaDiff enables direct and explicit control over physical realism and motion diversity, moving beyond the limitations of indirect structural priors and uniform noise application.

We extensively validate the effectiveness of our proposed diffusion model on Human3.6M and more challenging cross-dataset scenarios on AMASS. Our model outperformed previous models with multiple evaluation metrics on these datasets. Our contributions can be summarized as follows:

- We introduce KinemaDiff, a novel diffusion framework that integrates human skeletal structure and joint-specific motion dynamics directly within the diffusion process.
- We propose a learnable joint-adaptive noise generator to enhance motion diversity and a novel structural alignment mechanism to enforce anatomical consistency.
- We validate the effectiveness of our method through comprehensive experiments on Human 3.6M and the cross-dataset AMASS benchmark.

2 Related Work

Stochastic Human Motion Prediction. Early research in Human Motion Prediction (HMP) focused on deterministic forecasting using sequential models like RNNs and GCNs (Jain et al., 2016; Dang et al., 2021). To capture the inherent multimodality of human actions, the field of Human Motion Prediction (HMP) has shifted from deterministic forecasting to stochastic generation. Initial

generative approaches, primarily Variational Autoencoders (VAEs) (Walker et al., 2017) and Generative Adversarial Networks (GANs) (Barsoum et al., 2018), pioneered the generation of diverse futures but often struggled with long-term coherence and physical plausibility. More recently, Denoising Diffusion Models (DMs) have become the dominant paradigm, offering superior fidelity and diversity. Research in this area has seen rapid progress, from pioneering the paradigm with a two-stage framework Motiondiff (Wei et al., 2023) and simplifying the training pipeline via masked completion model HumanMAC (Chen et al., 2023), to introducing latent diffusion for more coherent, behavior-driven sampling model BeLFusion (Barquero et al., 2023) and integrating specialized network architectures like GCN-DCT to better capture spatio-temporal dynamics in CoMusion (Sun & Chowdhary, 2024). A concurrent line of work has also begun to adapt the diffusion process itself, such as SkeletonDiffusion (Curreli et al., 2025), which introduces anisotropic noise based on the skeleton's static structure. In contrast to these approaches, which largely treat the core diffusion mechanism as a fixed component or adapt it based on static priors, our work is the first to fundamentally reshape the denoising process to be dynamically aware of human kinematics. We achieve this by introducing a novel framework equipped with a structure-aligned constraint enforcer for anatomical consistency and a learnable, instance-adaptive Joint-Adaptive Noise Generator, offering a more flexible and physically grounded generative process.

Denoising Diffusion Probabilistic Models. Denoising Diffusion Probabilistic Models (DDPMs) (Ho et al., 2020; Nichol & Dhariwal, 2021; Song et al., 2020) have emerged as a powerful class of generative models, capable of synthesizing high-fidelity data by learning to reverse a progressive noising process. While these models have been successfully applied to a wide range of human-centric tasks (Gong et al., 2023; Shan et al., 2023), such applications have predominantly focused on innovating the denoiser's network architecture, while largely adopting a standard, generic diffusion process. However, a recent and promising research direction has begun to demonstrate significant gains by tailoring the diffusion process itself, particularly through task-specific noise designs (Sahoo et al., 2024; Huang et al., 2024). Following this trajectory, we introduce a novel diffusion framework specifically for human motion, which moves beyond architectural modifications. We propose to fundamentally reshape the denoising process with a joint-adaptive, structure-aligned mechanism that is inherently adapted to the kinematic properties of human skeleton data.

3 METHODOLOGY

3.1 PROBLEM FORMULATION

As illustrated in Fig.3, we denote the observed motion history of H frames as $\mathbf{x}^{(1:H)} = [\mathbf{x}^{(1)}; \mathbf{x}^{(2)}; \dots; \mathbf{x}^{(H)}] \in \mathbb{R}^{H \times 3J}$, where $\mathbf{x}^{(h)} \in \mathbb{R}^{3J}$ represents the joint coordinates at frame h, and J is the total number of joints. Given $\mathbf{x}^{(1:H)}$, the goal of Human Motion Prediction is to forecast the subsequent F frames $\mathbf{y}^{(1:F)} = \mathbf{x}^{(H+1:H+F)} = [\mathbf{x}^{(H+1)}; \mathbf{x}^{(H+2)}; \dots; \mathbf{x}^{(H+F)}] \in \mathbb{R}^{F \times 3J}$, where $\mathbf{y}^{(f)} \in \mathbb{R}^{J \times 3}$, and J is the number of body joints.

3.2 Preliminaries

Motion Diffusion. Let $\{y_t\}_{t=0}^T$ denote a Markov noising process, where y_0 represents the true data samples. For training, we progressively corrupt the target human motion sequence $\{y_t\}_{t=0}^T$ by adding noise. This forward diffusion process is represented as:

$$q(y_t \mid y_{t-1}) = \mathcal{N}(y_t; \sqrt{\alpha_t} y_{t-1}, (1 - \alpha_t)\mathbf{I}), \tag{1}$$

where $\{\alpha_t\}_{t=0}^T \in [0,1]$ controls the noise level. To approximate the underlying data distribution, the reverse diffusion process is formulated to iteratively remove noise from the corrupted samples y_t , starting from t=T down to t=1, as follows:

$$p_{\theta}(y_{t-1} \mid y_t) = \mathcal{N}(y_{t-1}; \mu_{\theta}(y_t, t), \sigma_{\theta}^2(y_t, t)\mathbf{I}).$$
(2)

In addition, following prior work (Chen et al., 2023; Barquero et al., 2023), we condition the model on the historical motion. The conditional reverse diffusion transition is then formulated as:

$$p_{\theta}(y_{t-1} \mid y_t, x) = \mathcal{N}(y_{t-1}; \mu_{\theta}(y_t, x, t), \sigma_{\theta}^2(y_t, x, t)\mathbf{I}).$$
(3)

where x denotes the motion history and y_t the noisy motion at step t.

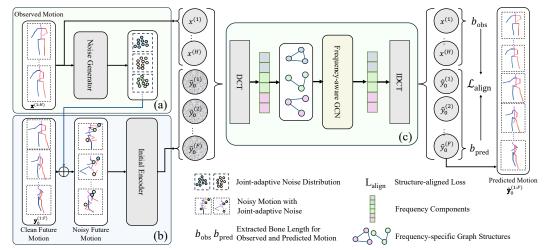


Figure 3: The overview of our proposed Kinemadiff. (a): Our Joint-adaptive noise generator. We learn joint-adaptive noise from the historical human joints and add it to the future human motions. (b): Initial motion reconstruction. The future human motion with injected noise is processed through a self-attention mechanism, which generates an initial prediction in the absence of external conditioning. (c): Structure-aligned constraint enforcer. The initial prediction is concatenated with the motion observations and then processed in the frequency domain through a frequency-aware GCN, and subsequently transformed back to the temporal domain for structural alignment.

Direct y₀ **prediction.** In existing diffusion models (Chen et al., 2023; Barquero et al., 2023), the objective is typically either to predict noise or to directly predict human motion. In this work, we adopt the latter, as it enables more effective optimization of the diffusion process through the incorporation of skeletal structural priors.

3.3 OVERALL ARCHITECTURE AND NETWORK DESIGN

Overall Architecture. As illustrated in Fig. 3, our network takes the observed motion history $\mathbf{x}^{(1:H)} \in \mathbb{R}^{H \times 3J}$ and the future frames $\mathbf{y}^{(1:F)} \in \mathbb{R}^{F \times 3J}$ as input. The future frames are perturbed with our joint-adaptive noise, which varies across joints and samples, yielding a sequence $\{y_t\}_{t=0}^T$ that follows the forward diffusion process. Our model is trained to reverse this process and ultimately reconstruct the predicted future motion $\hat{\mathbf{y}}_0 \in \mathbb{R}^{F \times 3J}$. First, we feed the noisy future frames \mathbf{y}_t into an encoder composed of several Transformer layers to obtain an initial reconstruction $\tilde{\mathbf{y}}_0 \in \mathbb{R}^{F \times 3J}$. To ensure structural alignment, we reshape the motion history as $\mathbf{x}^{(1:H)} \in \mathbb{R}^{H \times J \times 3}$ and the noisy future as $\tilde{\mathbf{y}}_0 \in \mathbb{R}^{F \times J \times 3}$. Next, we concatenate the initial reconstruction with the historical motion and process it with the Alignment Module to produce the denoised prediction.

Network Structure. As illustrated in Fig. 3, our network architecture consists of three main components: a joint-adaptive noise generator, an initial encoder, and a structure-aligned constraint enforcer. The first component generates joint-adaptive noise conditioned on motion history and joint characteristics using a few linear layers. The second component employs a temporal Transformer encoder to generate an initial reconstruction, providing a baseline prediction without conditioning. The third component is the proposed Structure-Aligned Constraint Enforcer, which models the full motion sequence in the frequency domain via DCT/IDCT. Unlike prior works, it employs GCNs with frequency-specific adjacency matrices, where each frequency band is associated with a tailored connectivity to capture its distinct motion patterns, thereby enabling more effective modeling of motion dynamics across different bands.

3.4 Joint-Adaptive Noise Generator.

Unlike previous diffusion-based methods that adopt a fixed schedule, our approach designs a learnable noise schedule. We introduce a multivariate noise schedule that assigns joint-adaptive noise rates, enabling the diffusion process to capture spatial variability in the human skeleton:

$$q(y_t \mid y_{t-1}) = \mathcal{N}(y_t; \alpha_t y_{t-1}, (1 - \alpha_t)\Sigma)$$
(4)

denoises the corrupted data samples y_t from t = T down to t = 1:

more realistic and coherent human motion generation.

predicted as clean human motion plus Gaussian noise:

compute the bone lengths of the motion sequence y from this mean:

STRUCTURE-ALIGNED CONSTRAINT ENFORCER

noise scaling can be expressed as:

where $\Sigma = \text{diag}(s_1^2, s_2^2, \dots, s_J^2)$, s_j denotes the noise scaling factor for the j-th joint. Specifically,

the scaling factor is determined by two aspects. First, it depends on the joint index, since different

joints exhibit distinct motion characteristics. Second, we further refine the scaling factor by condi-

tioning on the historical motion trajectories of each joint. Formally, the joint- and instance-specific

 $s_j = f_{\theta}(j, \mathbf{x}_i^{(1:H)})$

where j is the joint index, $\mathbf{x}_j^{(1:H)}$ denotes the observed motion history of joint j, and f_{θ} is a learnable function. Through this design, we are able to inject noise with varying intensities across different

joints, better reflecting their heterogeneous motion properties. The reverse diffusion process that

 $p_{\theta}(y_{t-1} \mid y_t) = \mathcal{N}(y_{t-1}; \mu_{\theta}(y_t, x, t), (1 - \alpha_t)\Sigma)$

where $\Sigma = \operatorname{diag}(s_1^2, s_2^2, \dots, s_J^2)$, s_j denotes the noise scaling factor for the j-th joint. Through the

above design, the diffusion process is able to inject joint- and instance-adaptive noise, leading to

Earlier diffusion-based approaches (Chen et al., 2023; Sun & Chowdhary, 2024) treated human motion prediction as a straightforward application of diffusion, overlooking the intrinsic structural

properties of human pose. For example, such methods fail to enforce structural constraints, such

as maintaining consistent bone lengths throughout motion generation. We proposed a Structure-

Aligned Constraint Enforcer, which effectively aligns the predicted human motion with historical

motion patterns based on the human skeletal structure during the diffusion process. Specifically, we

calculated the average bone length of connected joints from past human motions. Since the historical

motion is noise-free, it allows us to easily extract structural information. We regard the motion to be

 $y_t = \sqrt{\bar{\alpha}_t} y_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon$

where $\epsilon \sim \mathcal{N}(0, I)$ denotes Gaussian noise and y_0 is the clean motion. For a batch of y_t , we take

the mean across the batch. Since the noise ϵ follows a Gaussian distribution with zero mean, we can

 $\bar{y}_t = \frac{1}{B} \sum_{i=1}^{B} y_t^{(b)} \approx \sqrt{\bar{\alpha}_t} y_0,$

where B is the batch size. Specifically, we denote the set of skeletal connections as \mathcal{E} , where each

 $\ell_{i,j} = ||y^i - y^j||_2, \quad (i,j) \in \mathcal{E},$

As previously mentioned, our model predicts the future motion values directly at each step. When

t is relatively large, the direct prediction of y_0 tends to be inaccurate. To address this, we impose

 $(i,j) \in \mathcal{E}$ represents a connected joint pair. The corresponding bone length is defined as:

(5)

(6)

(7)

(8)

(9)

(10)

216 217

218

219

220 221

222 223

224 225

226

227 228

> 229 230

231 232

3.5

233 234

235 236

237 238 239

240 241

242 243

> 244 245 246

247 248

249 250

251 252

253 254

255 256

257 258 259

> 264 265

267

268 269

266

which encourages the predicted skeleton to preserve the structural scale of the observed motion,

thereby maintaining consistent bone proportions across time.

constraints on the bone lengths of the target human motion and perform an alignment with the structure of the historical human motion. Moreover, at each timestep, after the initial encoder,

we apply the same operation on \tilde{y}_0 to ensure that the human skeleton structure remains consistent

throughout the entire sequence. Specifically, for each bone connection $(i, j) \in \mathcal{E}$, we first compute

the average bone length over the observed history frames $\mathbf{x}^{(1:H)}$ as $\bar{b}_{\mathrm{obs}}^{(i,j)}$, and similarly obtain the average bone length over the predicted future frames $\hat{\mathbf{y}}_{0}^{(1:F)}$ as $\bar{b}_{\mathrm{pred}}^{(i,j)}$. The alignment loss is then defined as the mean disconverse between the context of \bar{b}_{pred} .

defined as the mean discrepancy between the two sets of averaged bone lengths:

where $y^i, y^j \in \mathbb{R}^3$ denote the 3D coordinates of joints i and j.

 $\mathcal{L}_{\mathrm{align}} = \frac{1}{|\mathcal{E}|} \sum_{(i,j) \in \mathcal{E}} \left| \bar{b}_{\mathrm{obs}}^{(i,j)} - \bar{b}_{\mathrm{pred}}^{(i,j)} \right|,$

3.6 Overall Learning Objectives

Our loss consists of two components: a reconstruction loss applied to the predicted results, and an alignment loss that enforces consistency between the predicted human motion and the observed human motion. For the reconstruction loss, we follow prior work (Sun & Chowdhary, 2024) and assign different weights to individual joints, which are weighted differently to reflect their relative importance in the motion context. The reconstruction loss is defined by:

$$\mathcal{L}_{\text{rec}} = \frac{1}{J} \sum_{j=1}^{J} \left(\gamma \cdot \left\| (x^j - \hat{x}^j) \cdot \lambda^j \right\|_1 + \left\| (y_0^j - \hat{y}_0^j) \cdot \lambda^j \right\|_1 \right), \tag{11}$$

where the superscript j denotes the joint index, λ^j is the weight assigned to each joint, and γ is a hyperparameter balancing the reconstruction of motion history and the prediction of future.

$$\mathcal{L}_{\text{total}} = \alpha \cdot \mathcal{L}_{\text{rec}} + \beta \cdot \mathcal{L}_{\text{align}}, \tag{12}$$

where α and β control the relative weight of the reconstruction and alignment losses.

4 EXPERIMENTS

We first introduce the experimental setup in §4.1. Then we assess the performance of our method across various settings, including intra-dataset forecasting on Human3.6M(§4.2), and more challenging cross-dataset generalization on AMASS(§4.3). Lastly, we provide ablative analyses in §4.4.

4.1 EXPERIMENTAL SETUP

Datasets. We conduct experiments on two widely used datasets, intra-dataset forecasting on Human 3.6M (Ionescu et al., 2013) and cross-dataset generalization on AMASS (Mahmood et al., 2019).

- **Human3.6M** is a seminal indoor dataset for 3D human motion analysis, widely utilized in stochastic Human Motion Prediction. It comprises 3.6 million frames, captured at 50Hz, documenting 11 subjects performing 15 common daily activities. Consistent with prior work (Barquero et al., 2023), we delineate subjects S1, S5, S6, S7, S8 for training, and subjects S9, S11 for evaluation.
- AMASS is a large-scale, highly diverse motion dataset for assessing cross-dataset generalization. It consolidates 24 distinct motion capture datasets, all standardized to the SMPL parameterization, accumulating over 9 million frames. Following prior research (Barquero et al., 2023), the dataset is partitioned into 11 training, 4 validation, and 7 testing constituent datasets.

Implementation Details. Our model, is trained end-to-end, following protocols as detailed below:

- **Diffusion Settings.** KinemaDiff employs a 10-step diffusion process with standard DDPM sampling, which is inherently augmented by our proposed kinematics-aware designs. Specifically, both the Structure-Aligned Constraint Enforcer and the Joint-Adaptive Noise Generator are integrated within the diffusion steps to ensure physical consistency and capture joint heterogeneity.
- Training Protocols. For both the Human3.6M and AMASS datasets, the model undergoes training for 500 epochs. We utilize the AdamW optimizer with a batch size of 128. The initial learning rate is set to 1e-4, which is subsequently decayed after the 200th epoch. These training parameters are consistent across both datasets.
- Dataset-Specific Protocols.
 - Human3.6M. Consistent with prior work (Barquero et al., 2023), we predict 100 future frames from 25 observed frames, using a 16-joint skeleton.
 - AMASS. Following the protocol established by (Barquero et al., 2023), the task involves forecasting 120 future frames (2s) based on 30 observed frames (0.5s).

Baselines. We compare our method against several representative diffusion-based methods.

- **BeLFusion**. Operating within a VAE-encoded latent space, BeLFusion employs a diffusion model to generate diverse and behavior-driven future motion predictions. Nevertheless, its physics-based consistency check is a post-processing step applied externally, rather than a constraint that guides the iterative denoising process internally.
- **CoMusion**. Integrating Graph Convolutional Networks (GCNs) within the Discrete Cosine Transform (DCT) space, CoMusion's denoiser effectively captures complex spatio-temporal dependencies. Though effective in modeling these dependencies, it focuses on advancing the architecture, rather than proposing adaptations to the core diffusion mechanism itself.

Table 1: Quantitative results on Human3.6M. The best results are highlighted in **bold**. The symbol '–' indicates not reported in the baseline work. For all metrics except for APD, lower is better.

Method	Reference_	Accuracy		Multimodality		Diversity	Realism	
Wethou		DE ↓	FDE ↓	$\overline{\text{MMADE}}$ \downarrow	MMFDE ↓	APD ↑	$\overline{\text{CMD}}\downarrow$	FID ↓
	GAN-Based							
HP-GAN (Barsoum et al., 2018)	[CVPRW2018] C).858	0.867	0.847	0.858	7.214	_	_
DeLiGAN (Gurumurthy et al., 2017	7) [CVPR2017] (0.483	0.534	0.520	0.545	6.509	_	_
VAE-Based								
TPK (Walker et al., 2017)	[ICCV2017] C).461	0.560	0.522	0.569	6.723	6.326	0.538
Motron (Salzmann et al., 2022)	[CVPR2022] C	0.375	0.488	0.509	0.539	7.168	40.796	13.743
DSF (Yuan & Kitani, 2019)	[ICLR2020] C).493	0.592	0.550	0.599	9.330	_	_
DLow (Yuan & Kitani, 2020)	[ECCV2020] C).425	0.518	0.495	0.531	11.741	4.927	1.255
GSPS (Mao et al., 2021)	[ICCV2021] C	0.389	0.496	0.476	0.525	14.757	10.758	2.103
DivSamp (Dang et al., 2022)	[ACMMM2022] C	0.370	0.485	0.475	0.516	15.310	11.692	2.083
DM-Based								
MotionDiff (Wei et al., 2023)	[AAAI2023] C).411	0.509	0.508	0.536	15.353	_	_
HumanMAC (Chen et al., 2023)	[ICCV2023] C	0.369	0.480	0.509	0.545	6.301	_	_
BeLFusion (Barquero et al., 2023)	[ICCV2023] C	0.372	0.474	0.473	0.507	7.602	5.988	0.209
CoMusion (Sun & Chowdhary, 202	24)[ECCV2024] (0.350	0.458	0.494	0.506	7.632	3.202	0.102
SkeletonDiff (Curreli et al., 2025)	[CVPR2025] C).344	0.450	0.487	0.512	7.249	4.178	0.123
Ours	- 0).331	0.449	0.500	0.520	6.912	4.60	0.083

• **SkeletonDiffusion**. Introducing a non-isotropic diffusion process, SkeletonDiffusion defines an anisotropic noise covariance matrix derived from the skeleton's static kinematic tree. Though effective in acknowledging joint heterogeneity, its noise characteristics remain fixed and are not adaptive to the unique dynamics of a given motion instance.

Evaluation Metrics. Following established practices, we evaluate our method across three key aspects: accuracy, diversity, and realism, using a comprehensive suite of standard metrics.

- Accuracy. We report Average Displacement Error (ADE) and Final Displacement Error (FDE), the mean ℓ_2 distance to the GT over the sequence and at the final frame, respectively.
- Diversity and Multimodality. We take Average Pairwise Distance (APD) to measure the variance
 among generated samples. we also report Multimodal ADE/FDE (MMADE/MMFDE), which
 assess multimodality by measuring the error to the best-matching ground-truth variant.
- Realism and Plausibility. We employ the Fréchet Inception Distance (FID) to assess the distributional similarity between generated and real motions. Additionally, following (Barquero et al., 2023), we use the Cumulative Motion Distribution (CMD) area for global plausibility to evaluate how realistically the model's generated diversity reflects that of the ground truth.

4.2 RESULTS OF INTRA-DATASET FORECASTING ON HUMAN 3.6M

The comparative performance on the Human3.6M dataset is systematically reported in Tab. 1. The results unequivocally demonstrate that our model establishes a new state-of-the-art result in forecasting accuracy and realism. Specifically, KinemaDiff achieves an ADE of **0.331** and an FDE of **0.449**, surpassing all prior methods. We attribute this superior accuracy to the Structure-Aligned Constraint Enforcer. By rigorously maintaining anatomical consistency at each step of the denoising process, this module prevents the accumulation of kinematic errors that can degrade long-term predictions, ensuring a physically grounded and accurate motion trajectory.

Furthermore, the effectiveness of our approach in generating high-fidelity motion is underscored by the FID score of **0.083**, a substantial 19% relative improvement over the previous leading model, CoMusion. A lower FID indicates that the distribution of generated motions is significantly closer to that of real human movements, not just in individual poses but in the naturalness of the entire sequence. This gain in realism is a direct result of the synergy between our two core modules: the Structure-Aligned Constraint Enforcer eliminates anatomically impossible poses, while the Joint-Adaptive Noise Generator sculpts more natural, heterogeneous joint movements, avoiding the robotic uniformity that can arise from conventional noise schedules. While other methods may achieve higher raw diversity scores (APD) or broader multimodal coverage (MMADE/FDE), Kine-

Table 2: Quantitative results for AMASS dataset. The best results are highlighted in **bold**. As AMASS does not contain class labels, the FID metric is not used for evaluation.

Method	Reference	Accuracy		Multimodality		Diversity	Realism
	Α	ADE ↓	FDE ↓	$MMADE\downarrow$	$\overline{MMFDE}\downarrow$	APD ↑	$\text{CMD} \downarrow$
VAE-Based							
TPK (Walker et al., 2017)	[ICCV2017] (0.656	0.675	0.658	0.674	9.283	17.127
DLow (Yuan & Kitani, 2020)	[ECCV2020]	0.590	0.612	0.618	0.617	13.170	15.185
GSPS (Mao et al., 2021)	[ICCV2021] (0.563	0.613	0.609	0.633	12.465	18.404
DivSampp (Dang et al., 2022)	[ACMMM2022] (0.564	0.647	0.623	0.667	24.724	50.239
DM-Based							
HumanMAC (Barquero et al., 2023) [ICCV2023] (0.511	0.554	0.593	0.591	9.321	_
BeLFusion (Barquero et al., 2023)	[ICCV2023] (0.513	0.560	0.569	0.585	9.376	16.995
CoMusion (Sun & Chowdhary, 202	4)[ECCV2024] (0.494	0.547	0.469	0.466	10.848	9.636
SkeletonDiff (Curreli et al., 2025)	[CVPR2025]	0.480	0.545	0.561	0.580	9.456	11.417
Ours	- (0.478	0.540	0.456	0.457	9.683	9.448

Table 3: Ablation of the main components in our Table 4: Experiment results on Humethod on Human3.6M.

Table 3: Ablation of the main components in our Table 4: Experiment results on Human3.6M with different Scheduler.

Encoder	J-Noise	Align	APD ↑	$ADE \downarrow$	FDE ↓	FID ↓	Scheduler	APD ↑	ADE ↓	FDE ↓	FID ↓
- ✓	-		19.601 9.600				Sqrt	6.837	0.342	0.457	0.108
-	-		6.214 7.243				Cosine	7.213	0.365	0.478	0.178
√		•	6.912				Variance	6.912	0.331	0.449	0.083

maDiff excels at ensuring that every generated sample possesses high physical fidelity, prioritizing the quality and plausibility of predictions as evidenced by its leading accuracy and realism metrics.

4.3 RESULTS OF CROSS-DATASET GENERALIZATION ON AMASS

To evaluate robustness and generalization, we present a comparative analysis on the diverse AMASS dataset in Tab. 2. In this challenging cross-dataset scenario, KinemaDiff demonstrates exceptional performance, achieving advanced results across the majority of metrics, including ADE (0.478), FDE (0.540), MMADE (0.456), MMFDE (0.457), and CMD (9.448). These results highlight the model's ability to learn fundamental principles of motion rather than dataset-specific artifacts. The superior generalization is primarily driven by the Structure-Aligned Constraint Enforcer, which learns intrinsic and invariant anatomical properties like bone lengths, making the model robust to the wide variety of novel motions present in AMASS. Concurrently, the results in multimodal metrics (MMADE/MMFDE) provide clear evidence for the efficacy of the Joint-Adaptive Noise Generator. On a diverse dataset like AMASS, where a single observation can lead to many valid future actions, the ability to generate instance-specific, heterogeneous noise allows the model to explore this rich possibility space more effectively than methods with static or uniform noise. It does not merely generate random variations but rather meaningful and plausible alternatives tailored to the input context. While some methods achieve a higher raw diversity score (APD) at the cost of accuracy, KinemaDiff achieves a competitive APD (9.683) while simultaneously delivering the best accuracy and multimodal coverage, demonstrating a superior balance between diversity and fidelity.

4.4 ABLATION STUDY

For in-depth analysis, we conduct ablative studies using Intra-Dataset Forecasting on Human3.6M.

Effect of our main components. To assess the contribution of our core components, we conducted an ablation study by removing the structure-aligned constraint enforcer, joint-adaptive noise, and the initial encoder. From Tab. 3, it can be seen that using only the baseline or only the initial encoder leads to low accuracy (ADE, FDE) and FID scores, despite high APD values. This suggests the model produces motions that are diverse yet largely implausible and unstructured. Incorporating the

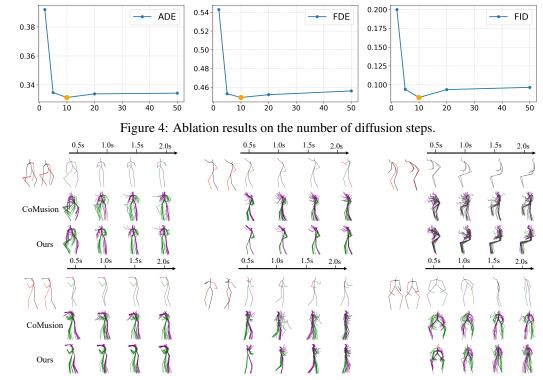


Figure 5: Visualization results. The red-black skeletons and green-purple skeletons denote the observed and predicted motions respectively.

structure-aligned constraint enforcer substantially improves accuracy metrics (ADE, FDE) and consistency metrics (FID) by leveraging structural cues aligned with historical motion. In addition, the joint-adaptive noise dynamically allocates noise to different joints, enabling more accurate modeling of human motion and further improving both accuracy (ADE, FDE) and consistency (FID) metrics.

Diffusion setting. We mainly investigate two aspects of the diffusion setting: the choice of scheduler and the number of diffusion steps. We evaluate multiple scheduler choices on the Human3.6M dataset in Tab. 4, and the results indicate that Variance Scheduler (Sun & Chowdhary, 2024) provides the best trade-off between stability and performance. In addition, we investigate the effect of the denoising step on model performance. We evaluate multiple metrics on Human3.6M under different timesteps. As shown in Fig. 4, the model achieves the best performance when the number of timesteps is 10. Moreover, choosing 10 timesteps ensures fast inference and accuracy.

Visualization results. In Fig. 5, we present a qualitative comparison by visualizing predicted motion sequences on Human3.6M. We use CoMusion as the baseline and select 15 predicted results for each frame. The visual analysis shows that our method generates more consistent and realistic human motions. Compared with the ground truth, our predictions are especially accurate for samples with smooth movements. In addition, unrealistic artifacts, such as sudden exaggerated leg lifts while walking, occur less frequently. Moreover, the diversity of generated motions better reflects realistic dynamics, with predictions concentrated near historical patterns and gradually diffusing over time.

5 Conclusion

In this work, we introduce a new diffusion model tailored for stochastic human motion prediction. Our algorithm integrates a Joint-Adaptive Noise Generator and a Structure-Aligned Constraint Enforcer directly within the diffusion process. The former enhances motion diversity with instance-aware noise, while the latter preserves anatomical consistency by embedding structural priors into each diffusion step. Results across multiple benchmarks demonstrate its effectiveness, owing to the unique integration of structural and dynamic priors within diffusion process.

ETHICS STATEMENT.

This work focuses on human motion prediction using publicly available benchmark datasets (Human3.6M and AMASS), which were collected and released under established research protocols. No personally identifiable or sensitive information is involved, and our methodology does not involve human subjects, sensitive attributes, or private data, posing no privacy or security concerns. All experiments follow ethical research practices, including proper citations, fair comparisons with prior works, and reproducibility efforts. All authors have read and will adhere to the ICLR Code of Ethics.

REPRODUCIBILITY STATEMENT.

We have made significant efforts to ensure the reproducibility of our work. A detailed description of our model architecture is provided in Section 3, while the evaluation protocols and training setup are presented in Section 4 of the main paper. Additional implementation details are included in the Appendix. All datasets used in our experiments are publicly available and widely adopted benchmarks, and our preprocessing steps strictly follow prior works to ensure consistency and comparability.

REFERENCES

- German Barquero, Sergio Escalera, and Cristina Palmero. Belfusion: Latent diffusion for behavior-driven human motion prediction. In <u>Proceedings of the IEEE/CVF international conference on computer vision</u>, pp. 2317–2327, 2023.
- Emad Barsoum, John Kender, and Zicheng Liu. Hp-gan: Probabilistic 3d human motion prediction via gan. In <u>Proceedings of the IEEE conference on computer vision and pattern recognition</u> workshops, pp. 1418–1427, 2018.
- Ling-Hao Chen, Jiawei Zhang, Yewen Li, Yiren Pang, Xiaobo Xia, and Tongliang Liu. Human-mac: Masked motion completion for human motion prediction. In <u>Proceedings of the IEEE/CVF</u> international conference on computer vision, pp. 9544–9555, 2023.
- Cecilia Curreli, Dominik Muhle, Abhishek Saroha, Zhenzhang Ye, Riccardo Marin, and Daniel Cremers. Nonisotropic gaussian diffusion for realistic 3d human motion prediction. In <u>Proceedings</u> of the Computer Vision and Pattern Recognition Conference, pp. 1871–1882, 2025.
- Lingwei Dang, Yongwei Nie, Chengjiang Long, Qing Zhang, and Guiqing Li. Msr-gcn: Multi-scale residual graph convolution networks for human motion prediction. In <u>Proceedings of the IEEE/CVF</u> international conference on computer vision, pp. 11467–11476, 2021.
- Lingwei Dang, Yongwei Nie, Chengjiang Long, Qing Zhang, and Guiqing Li. Diverse human motion prediction via gumbel-softmax sampling from an auxiliary space. In Proceedings of the 30th ACM international conference on multimedia, pp. 5162–5171, 2022.
- Jia Gong, Lin Geng Foo, Zhipeng Fan, Qiuhong Ke, Hossein Rahmani, and Jun Liu. Diffpose: Toward more reliable 3d pose estimation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 13041–13051, 2023.
- Liang-Yan Gui, Kevin Zhang, Yu-Xiong Wang, Xiaodan Liang, José MF Moura, and Manuela Veloso. Teaching robots to predict human motion. In <u>2018 IEEE/RSJ International Conference</u> on Intelligent Robots and Systems (IROS), pp. 562–567. IEEE, 2018.
- Swaminathan Gurumurthy, Ravi Kiran Sarvadevabhatla, and R Venkatesh Babu. Deligan: Generative adversarial networks for diverse and limited data. In <u>Proceedings of the IEEE conference on computer vision and pattern recognition</u>, pp. 166–174, 2017.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. <u>Advances in</u> neural information processing systems, 33:6840–6851, 2020.
- Xingchang Huang, Corentin Salaun, Cristina Vasconcelos, Christian Theobalt, Cengiz Oztireli, and Gurprit Singh. Blue noise for diffusion models. In <u>ACM SIGGRAPH 2024 conference papers</u>, pp. 1–11, 2024.

- Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. Human3. 6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. <u>IEEE transactions</u> on pattern analysis and machine intelligence, 36(7):1325–1339, 2013.
 - Ashesh Jain, Amir R Zamir, Silvio Savarese, and Ashutosh Saxena. Structural-rnn: Deep learning on spatio-temporal graphs. In <u>Proceedings of the ieee conference on computer vision and pattern recognition</u>, pp. 5308–5317, 2016.
 - Maosen Li, Siheng Chen, Zijing Zhang, Lingxi Xie, Qi Tian, and Ya Zhang. Skeleton-parted graph scattering networks for 3d human motion prediction. In <u>European conference on computer vision</u>, pp. 18–36. Springer, 2022.
 - Tiezheng Ma, Yongwei Nie, Chengjiang Long, Qing Zhang, and Guiqing Li. Progressively generating better initial guesses towards next stages for high-quality human motion prediction. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 6437–6446, 2022.
 - Naureen Mahmood, Nima Ghorbani, Nikolaus F Troje, Gerard Pons-Moll, and Michael J Black. Amass: Archive of motion capture as surface shapes. In <u>Proceedings of the IEEE/CVF</u> international conference on computer vision, pp. 5442–5451, 2019.
 - Wei Mao, Miaomiao Liu, and Mathieu Salzmann. Generating smooth pose sequences for diverse human motion prediction. In <u>Proceedings of the IEEE/CVF International Conference on Computer Vision</u>, pp. 13309–13318, 2021.
 - Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In International conference on machine learning, pp. 8162–8171. PMLR, 2021.
 - Brian Paden, Michal Čáp, Sze Zheng Yong, Dmitry Yershov, and Emilio Frazzoli. A survey of motion planning and control techniques for self-driving urban vehicles. <u>IEEE Transactions on intelligent vehicles</u>, 1(1):33–55, 2016.
 - Subham Sahoo, Aaron Gokaslan, Christopher M De Sa, and Volodymyr Kuleshov. Diffusion models with learned adaptive noise. <u>Advances in Neural Information Processing Systems</u>, 37:105730–105779, 2024.
 - Tim Salzmann, Marco Pavone, and Markus Ryll. Motron: Multimodal probabilistic human motion forecasting. In <u>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern</u> Recognition, pp. 6457–6466, 2022.
 - Wenkang Shan, Zhenhua Liu, Xinfeng Zhang, Zhao Wang, Kai Han, Shanshe Wang, Siwei Ma, and Wen Gao. Diffusion-based 3d human pose estimation with multi-hypothesis aggregation. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 14761–14771, 2023.
 - Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. <u>arXiv</u> preprint arXiv:2010.02502, 2020.
 - Jiarui Sun and Girish Chowdhary. Comusion: Towards consistent stochastic human motion prediction via motion diffusion. In <u>European Conference on Computer Vision</u>, pp. 18–36. Springer, 2024.
 - Jacob Walker, Kenneth Marino, Abhinav Gupta, and Martial Hebert. The pose knows: Video forecasting by generating pose futures. In Proceedings of the IEEE international conference on computer vision, pp. 3332–3341, 2017.
 - Dong Wei, Huaijiang Sun, Bin Li, Jianfeng Lu, Weiqing Li, Xiaoning Sun, and Shengxiang Hu. Human joint kinematics diffusion-refinement for stochastic motion prediction. In <u>Proceedings of</u> the AAAI Conference on Artificial Intelligence, volume 37, pp. 6110–6118, 2023.
 - Chenxin Xu, Robby T Tan, Yuhong Tan, Siheng Chen, Xinchao Wang, and Yanfeng Wang. Auxiliary tasks benefit 3d skeleton-based human motion prediction. In <u>Proceedings of the IEEE/CVF</u> international conference on computer vision, pp. 9509–9520, 2023.

Ye Yuan and Kris Kitani. Diverse trajectory forecasting with determinantal point processes. <u>arXiv</u> preprint arXiv:1907.04967, 2019.

Ye Yuan and Kris Kitani. Dlow: Diversifying latent flows for diverse human motion prediction. In European Conference on Computer Vision, pp. 346–364. Springer, 2020.

A APPENDIX

A.1 THE USE OF LARGE LANGUAGE MODELS (LLMS).

In preparing this manuscript, we used a large language model (LLM) only for grammar correction and improving the clarity of phrasing. The LLM was not involved in generating ideas, methods, experiments, analyses, or results. All scientific contributions, including the problem formulation, model design, and evaluation, are entirely the work of the authors.

A.2 ADDITIONAL VISUALIZATION.

Similarly, in the supplementary material, we provide additional qualitative visualizations using Co-Musion as the baseline, selecting 15 predicted results per sample, as shown in Fig. 6. The extended visual analysis demonstrates that our model's predictions are consistently closer to the ground truth, highlighting the improved accuracy. Furthermore, our method produces more coherent motions with fewer unrealistic poses. These visual improvements substantiate the effectiveness of our approach, which can be attributed to the tailored diffusion process incorporating the structure-aligned constraint enforcer and joint-adaptive noise generator.

A.3 ADDITIONAL HYPERPARAMETER SETTING.

Additional diffusion setting. We introduce three additional metrics—APD, MMADE, and CMD—to evaluate model performance across different timesteps. As shown in Fig. 7, when the number of timesteps is set to 10, the model achieves strong performance in terms of diversity, accuracy, and consistency. Therefore, we select this setting for our experiments.

Hyperparameter for loss function. In Section 3.6, we set $\gamma=0.1$ to balance the reconstruction of motion history and the prediction of future frames. The coefficients α and β are set to 1 and 2, respectively, to control the relative contributions of the reconstruction loss and the alignment loss. These values were selected based on experiments with several parameter configurations to identify the most effective setting.

Model setting. In the initial encoder, we stack two Transformer layers with a feature dimension of 512. In the Structure-Aligned Constraint Enforcer, we employ nine Frequency-aware GCN layers with a feature dimension of 125. The graph structure consists of N nodes, where N corresponds to the number of joints in the skeleton.

A.4 ADDITIONAL METRIC DESCRIPTIONS.

- APD (Average Pairwise Distance) measures the diversity of generated samples by computing the average distance between all pairs of generated motions.
- ADE (Average Displacement Error) computes the mean per-timestep distance between the predicted and ground-truth motions, reflecting overall accuracy.
- FDE (Final Displacement Error) measures the distance between the predicted and ground-truth motions at the final timestep, highlighting long-term prediction accuracy.
- MMADE and MMFDE extend ADE/FDE by comparing with clustered groundtruth variants, capturing a model's ability to generate multiple plausible outcomes.
- CMD (Conditional Motion Distance) quantifies global plausibility by comparing the areas under cumulative distributions of true and generated motion.
- FID (Fréchet Inception Distance for motion) computes the distance between the feature distributions of generated and ground-truth motions, reflecting realism at the distribution level.

As shown in Tab. 5, we present the formulas for the metrics used in the main text, with the definitions of relevant parameters as follows.

Notation. We denote the ground-truth motion sequence by $\mathbf{p}^{gt} = {\{\mathbf{p}_t^{gt}\}_{t=1}^T}$, and the k-th generated motion sequence by $\mathbf{p}^{(k)} = {\{\mathbf{p}_t^{(k)}\}_{t=1}^T}$, where T is the prediction horizon and K is the number of

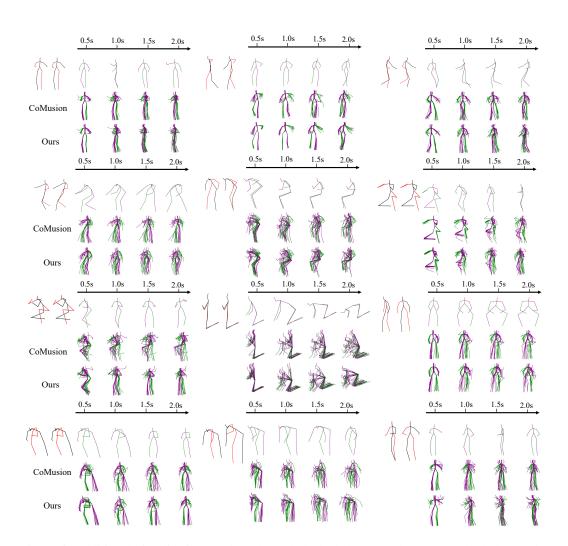


Figure 6: Additional visualization results. The red-black skeletons and green-purple skeletons denote the observed and predicted motions respectively.

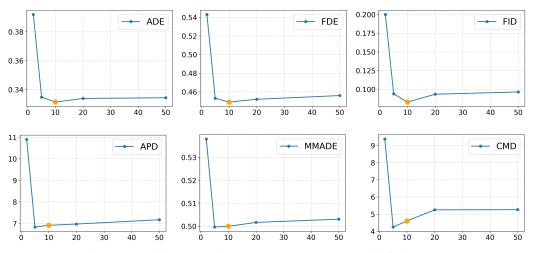


Figure 7: Additional ablation results on the number of diffusion steps.

FID*

Table 5: Evaluation metrics used for motion prediction.

	00
7	59
7	60
7	61
7	62
7	63

$$\begin{array}{lll} \textbf{Metric} & \textbf{Formula} \\ \\ \textbf{APD} & \frac{1}{K(K-1)} \sum_{i < j} \frac{1}{T} \sum_{t=1}^{T} \| \mathbf{p}_{t}^{(i)} - \mathbf{p}_{t}^{(j)} \|_{2} \\ \\ \textbf{ADE} & \frac{1}{T} \sum_{t=1}^{T} \| \hat{\mathbf{p}}_{t} - \mathbf{p}_{t}^{gt} \|_{2} \\ \\ \textbf{FDE} & \| \hat{\mathbf{p}}_{T} - \mathbf{p}_{T}^{gt} \|_{2} \\ \\ \textbf{CMD} & \sum_{t=1}^{T-1} (T-t) \| M_{t} - \bar{M} \|_{1} \\ \\ \textbf{FID*} & \| \mu_{g} - \mu_{r} \|_{2}^{2} + \text{Tr} \Big(\Sigma_{g} + \Sigma_{r} - 2(\Sigma_{g} \Sigma_{r})^{1/2} \Big) \\ \end{array}$$

generated samples. Here, $\mathbf{p}_t \in \mathbb{R}^{J \times 3}$ represents the 3D skeleton at timestep t, with J denoting the number of joints. For CMD, we compute the average displacement of all joints at frame t as M_t , and the average displacement over the whole test set as \bar{M} . (μ, Σ) are the mean and covariance of extracted motion features used for FID calculation.