# Deep Compositional Phase Diffusion for Long Motion Sequence Generation

**Ho Yin Au, Jie Chen, Junkun Jiang, Jingyu Xiang**

Department of Computer Science, Hong Kong Baptist University

# Intro: Long-term Compositional Generation

- Human Motion Generation
  - Short duration motion clip of a single semantics

- Long-term Compositional Generation
  - Multiple **sequentially connected clips**, each with single semantics.
    - Each clip is aligned to the **semantic** condition
    - **Transitions** between clips are smooth and natural
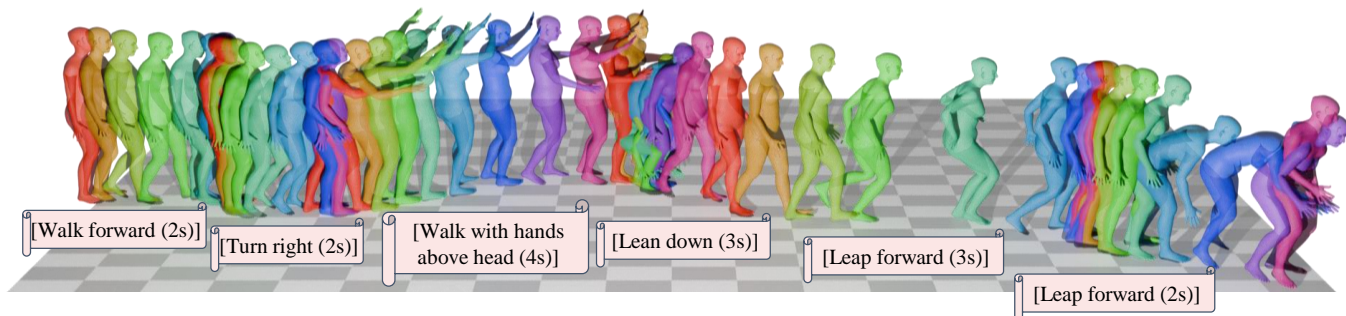  - Existing models <u>struggle to synthesize</u> natural and seamless transitions



[Walk forward (2s)]   [Turn right (2s)]   [Walk with hands above head (4s)]   [Lean down (3s)]   [Leap forward (3s)]   [Leap forward (2s)]

Illustration of Long-term Compositional Generation

# Related Works

- TEACH          (3DV 2022)
  - Autoregressively synthesizes motion clips, having minor discontinuities between clips
  - Blending transitions are generated using spherical linear interpolation (SLERP)

- Limitation
  - Blending transitions often appear <u>unrealistic</u>, reducing overall motion realism
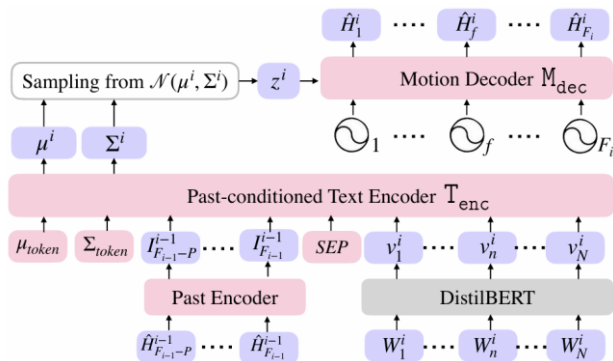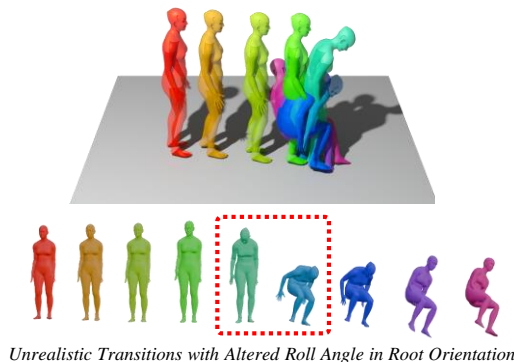


Illustration of TEACH's autoregressive
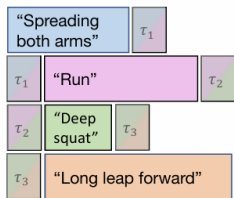variational encoder-decoder architecture



*Unrealistic Transitions with Altered Roll Angle in Root Orientation*

Result visualization of
[walk(2.4s), sit down(3.6s)]

▪ Athanasiou et al., Teach: Temporal action composition for 3d humans. In 3DV, pp. 414–423, 2022.
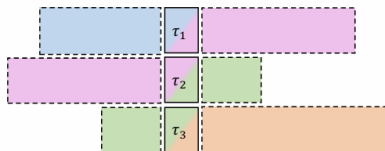
# Related Works

- PriorMDM        (ICLR 2024)
  - Synthesizes all motion clips in parallel, resulting in significant discontinuities between clips
  - Blending transitions are generated using the human motion diffusion model (MDM)
- Limitation
  - <u>Substantial discontinuities</u> between clips, and generated transitions fail to smooth them
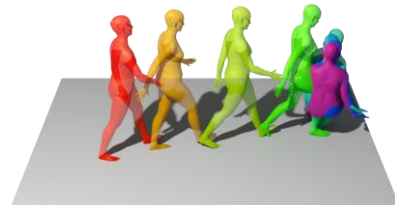


Illustration of priorMDM's
DoubleTake algorithm



*Rigid Transition with Foot Sliding and Abrupt Turn*

Result visualization of
[walk(2.4s), sit down(3.6s)]

▪ Shafir et al., Human motion diffusion as a generative prior, ICLR, 2024.

# Related Works

- DeepPhase      (SIGGRAPH 2022)
  - Encode motion into the periodic latent space using fixed-length convolutions and FFT
  - Excels at motion extrapolation and inbetweening, effectively capturing **motion dynamics**
- Limitation
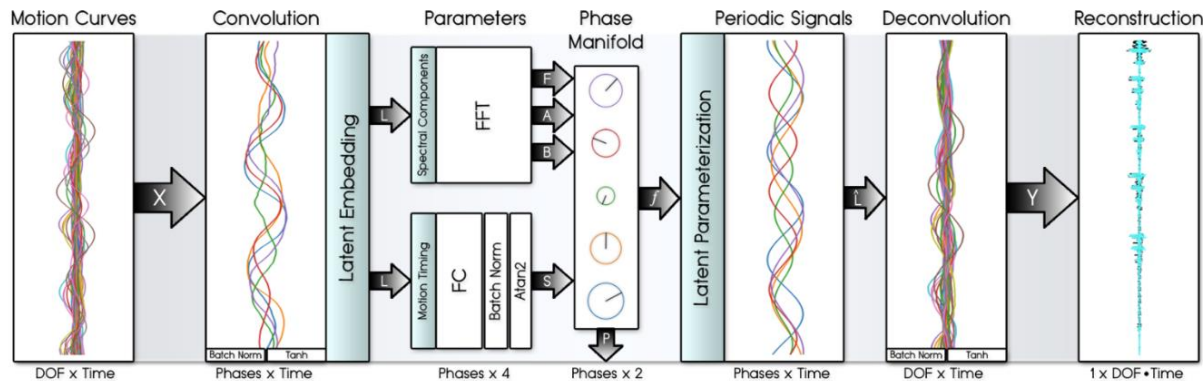  - Results in variable numbers of phase latents, more <u>difficult to learn text-motion alignment</u>.



Illustration of DeepPhase's periodic autoencoder architecture

Starke et al., Deepphase: Periodic autoencoders for learning motion phase manifolds. ACM Transactions on Graphics, 2022

# Framework Core Idea

- Integrates **motion dynamics from adjacent segments** into diffusion process
  - Minimizes discontinuity between consecutive motion segments
  - Utilizes these dynamics to generate realistic blending transitions

- Employs **phase mixing** to jointly integrate semantic and transitional conditions
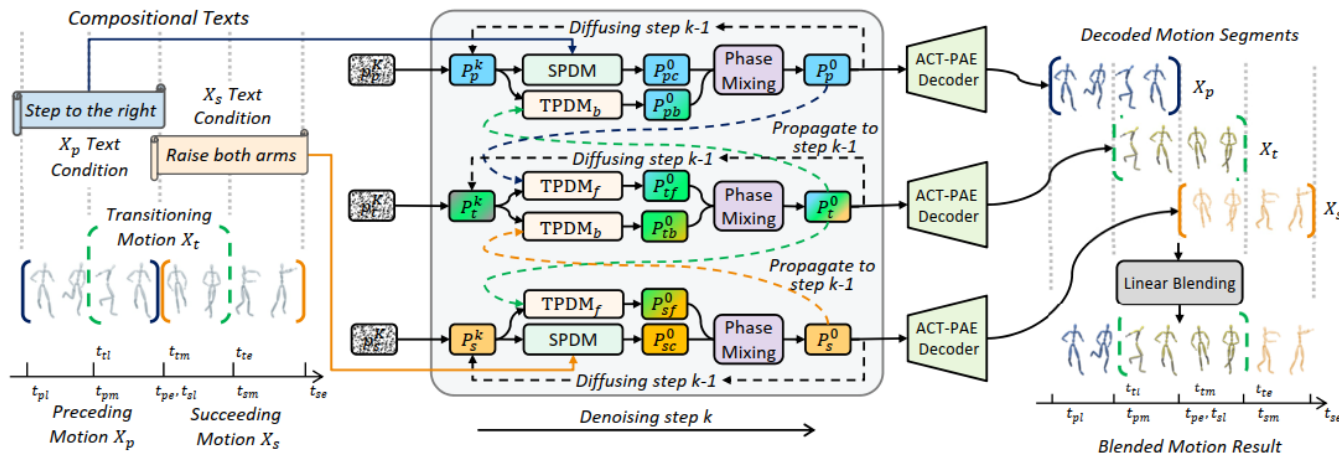


Illustration of our Compositional Motion Generation pipeline

# Framework Overview

- Three Components:
  - ACT-PAE: Encode human motion sequence into **phase latent space**
  - SPDM: Incorporates **semantic information** into the diffusion process
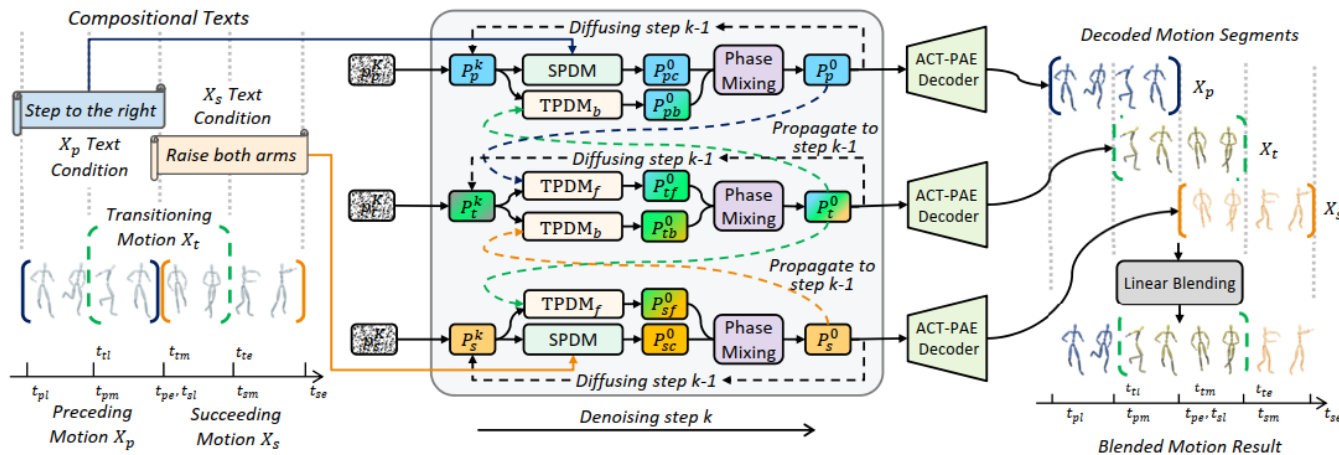  - TPDM: Integrates **adjacent motion dynamics** into the diffusion process



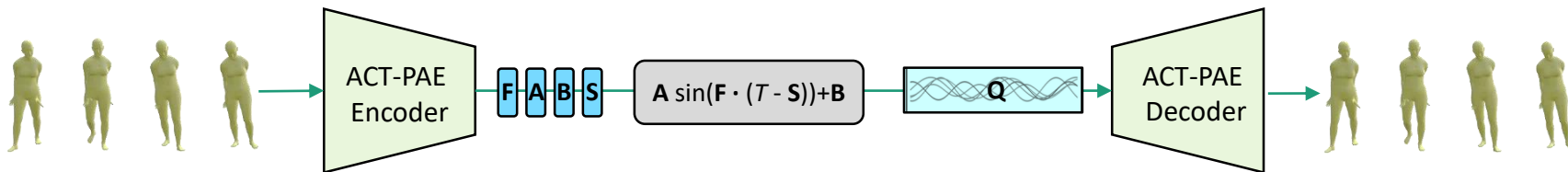Illustration of our Compositional Motion Generation pipeline

# ACT-PAE

- **A**ction **C**en**T**ric **P**eriodic **A**uto**E**ncoder (ACT-PAE):
  - Encode motion $X \in \mathbb{R}^{N \times E}$ into four phase latents $F, A, B, S \in \mathbb{R}^{Q}$
  - Reparameterizes latents into periodic signal $Q \in \mathbb{R}^{N \times Q}$ using:     $Q = A \sin(F \cdot (T - S)) + B$
  - Decode $Q$ back to reconstructed motion $\widehat{X}$

- Advantage over Traditional PAE:
  - Avoiding fixed-window motion slicing, enabling the capture of unified and semantic meaningful motion dynamics in the motion segment as a **cohesive unit**



Detail of the Action centric Periodic AutoEncoder (ACT-PAE)

# SPDM and TPDM

- **S**emantic **P**hase **D**iffusion **M**odule (SPDM):
  - Denoise phase latents $P_p^k$ utilizing semantic condition $\mathbf{C}$
  - Inputs include both phase latents $P_p^k = [F, A, B, S]$ and periodic signal $Q = A\sin(F \cdot (T - S)) + B$

- **T**ransitional **P**hase **D**iffusion **M**odule (TPDM):
  - Functions similarly to SPDM, but uses adjacent phase parameters $P_p^0$ as input

- Advantage over Latent Diffusion Model:
  - The periodic signal $Q$ captures **spatial-temporal context** within the phase latents $P$



Detail of the Semantic Phase Diffusion Module (SPDM)          Detail of the Transitional Phase Diffusion Module (TPDM)

# Compositional Phase Diffusion Pipelines

- Pipeline for **Compositional Motion Generation**:
  - Denoise phase latent using SPDM (for $P_c^0$) and two TPDMs (for $P_f^0$ and $P_b^0$)

  - Perform **phase mixing** on the denoised outputs:    $P^0 = \mathrm{r}\dfrac{P_f^0 + P_b^0}{2} + (1 - \mathrm{r})P_c^0$

  - Diffuse $P^0$ to the step k − 1 or decode via ACT-PAE to generate $X_\mathrm{p}, X_\mathrm{t}, X_\mathrm{s}$

  - **Linear blend** the transition segment $X_\mathrm{t}$ into the overlap region between $X_\mathrm{p}$ and $X_\mathrm{s}$
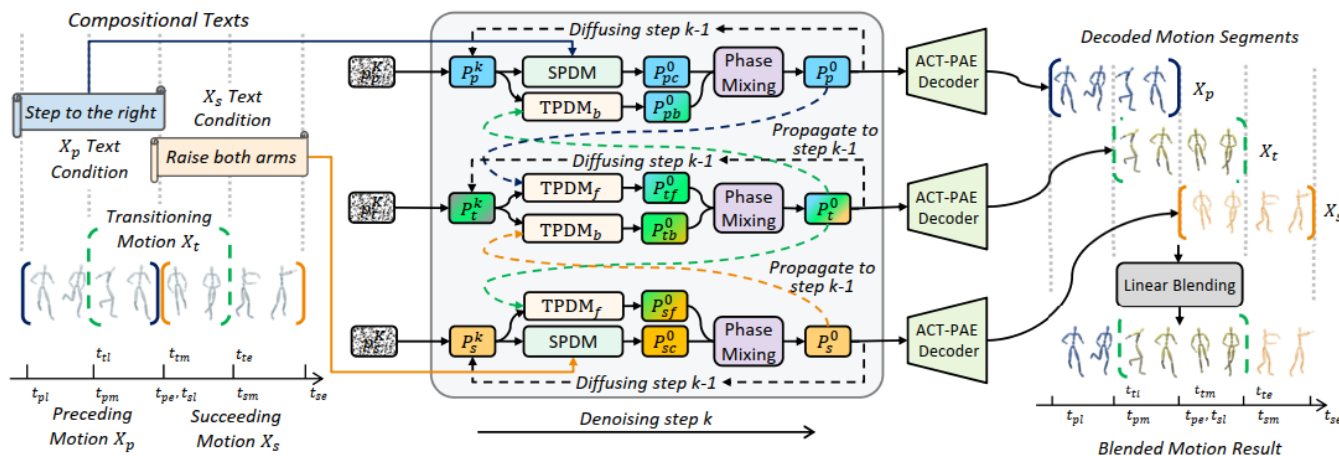


Illustration of our Compositional Motion Generation pipeline

# Compositional Phase Diffusion Pipelines

- Pipeline for **Motion Inbetweening**:
  - Encode user-provided motion segment into phase latents ($P_p^0$ and $P_s^0$) using ACT-PAE encoder
  - Denoise the inbetweening motion $X_i$ and the transitioning motions $X_{t_1}$ and $X_{t_2}$
  - **Linear blend** the transition segment $X_{t_1}$ and $X_{t_2}$ into the overlapping regions of $X_p$, $X_i$, $X_s$
  - Note that $X_i$ can be further conditioned with text input by incorporating an optional SPDM
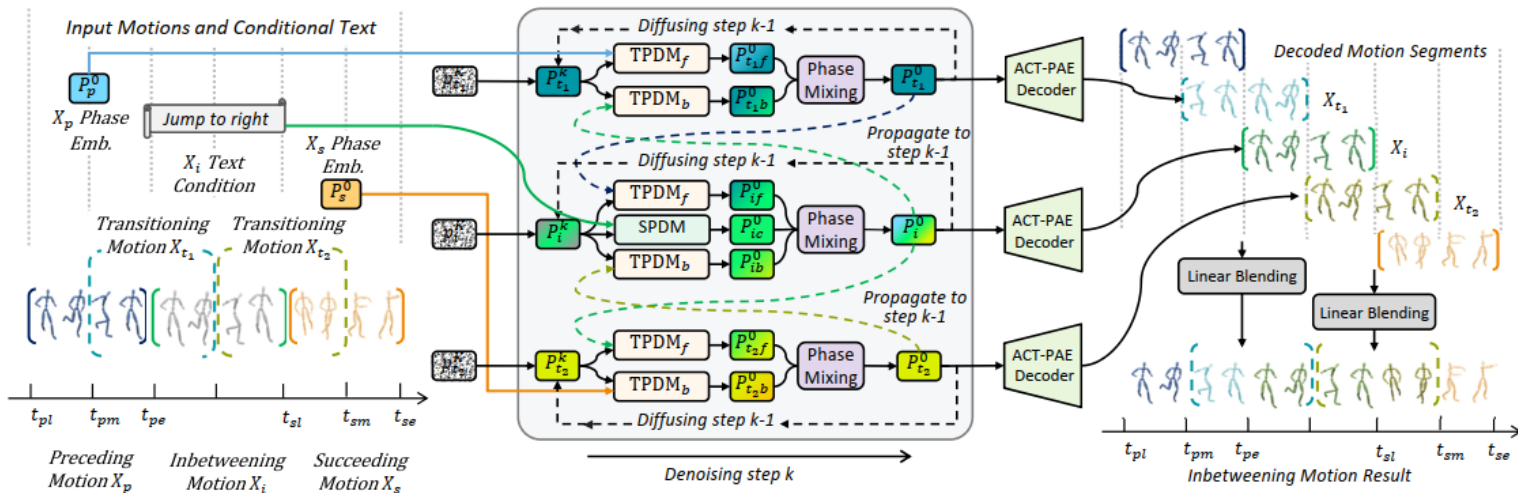


Illustration of our Motion Inbetweening pipeline

# Compositional Phase Diffusion Pipelines

- Advantages of the pipelines:
  - **Coherent Transitions**:
    Bidirectional TPDMs progressively <u>propagate phase information throughout the sequence</u>, ensuring smoother and more coherent motion generation
  - **Scalable and Flexible**:
    Pipeline supports an <u>arbitrary number of segments</u> by rearranging SPDM and TPDM modules, with <u>parallel processing</u> for efficient generation of long motion sequences.
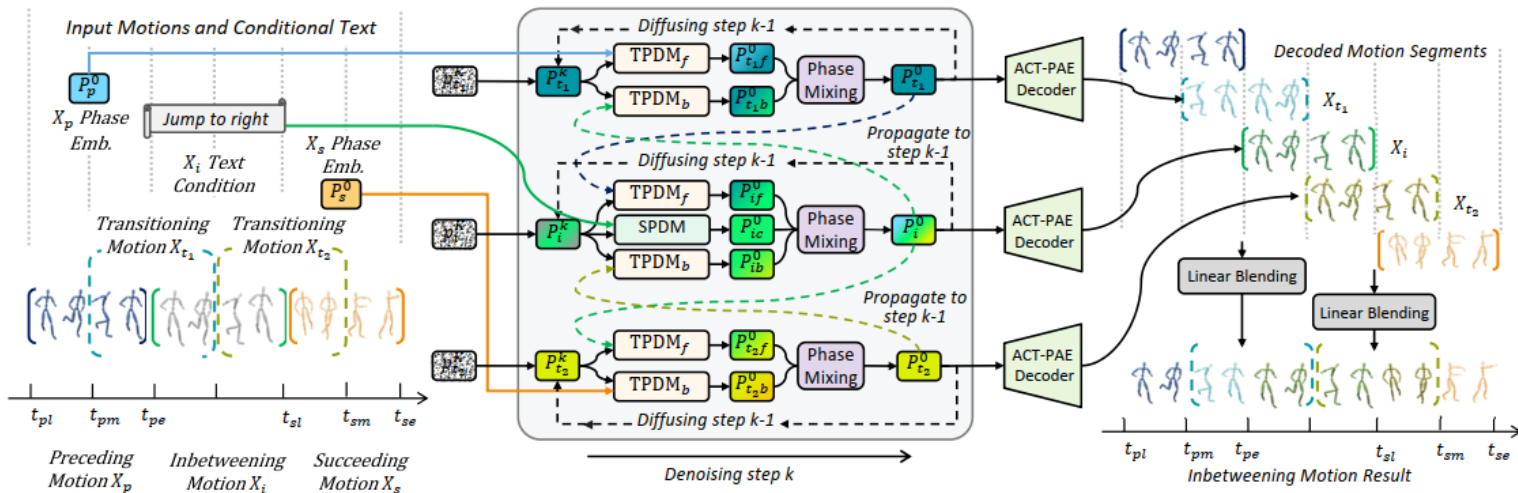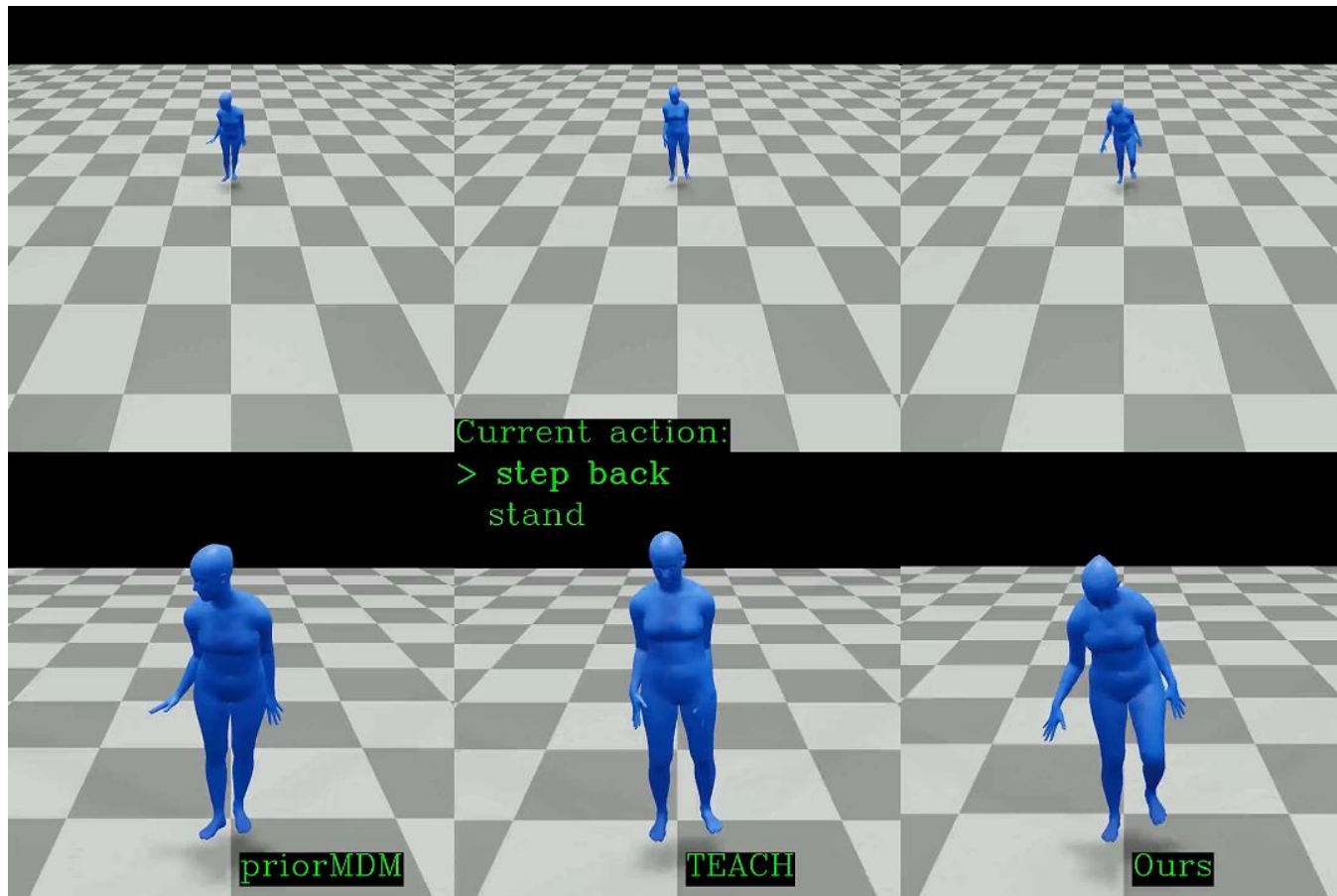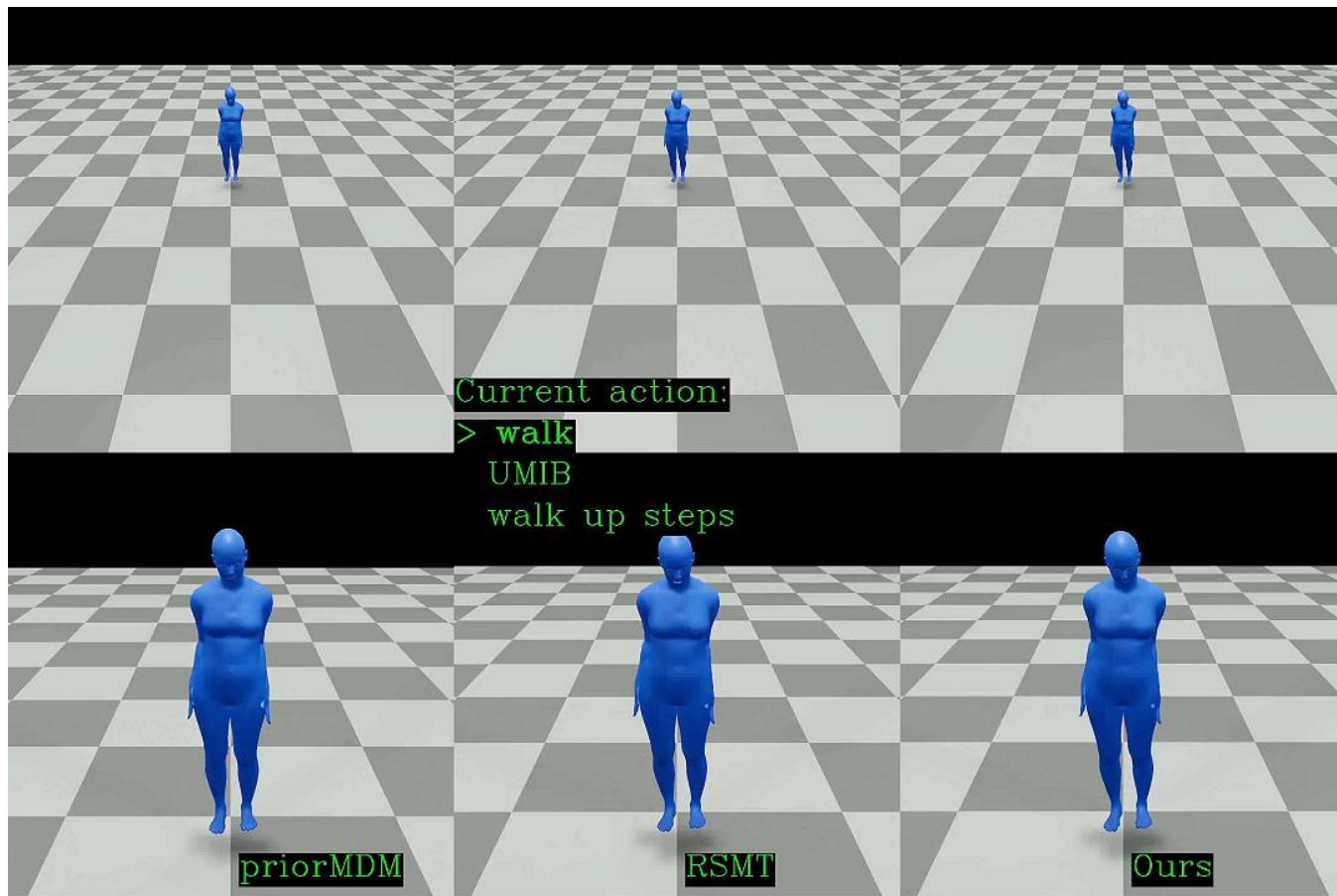


Illustration of our Motion Inbetweening pipeline

# Result: Compositional Motion Generation

# Result: Motion Inbetweening

# Summary

- Operating within a unified phase latent space facilitates alignment of transitional dynamics, enables <u>smoother transitions between motion clips</u>

- Scalable and efficient framework supports diverse motion generation tasks and allows <u>parallel processing of arbitrary number of motion segments</u>

**Github Code**

**arXiv Paper**