

Digital Humans – Winter 24/25

Lecture 13_2 – Motion Synthesis with Diffusion

Prof. Dr. Gerard Pons-Moll

University of Tübingen / MPI-Informatics

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN



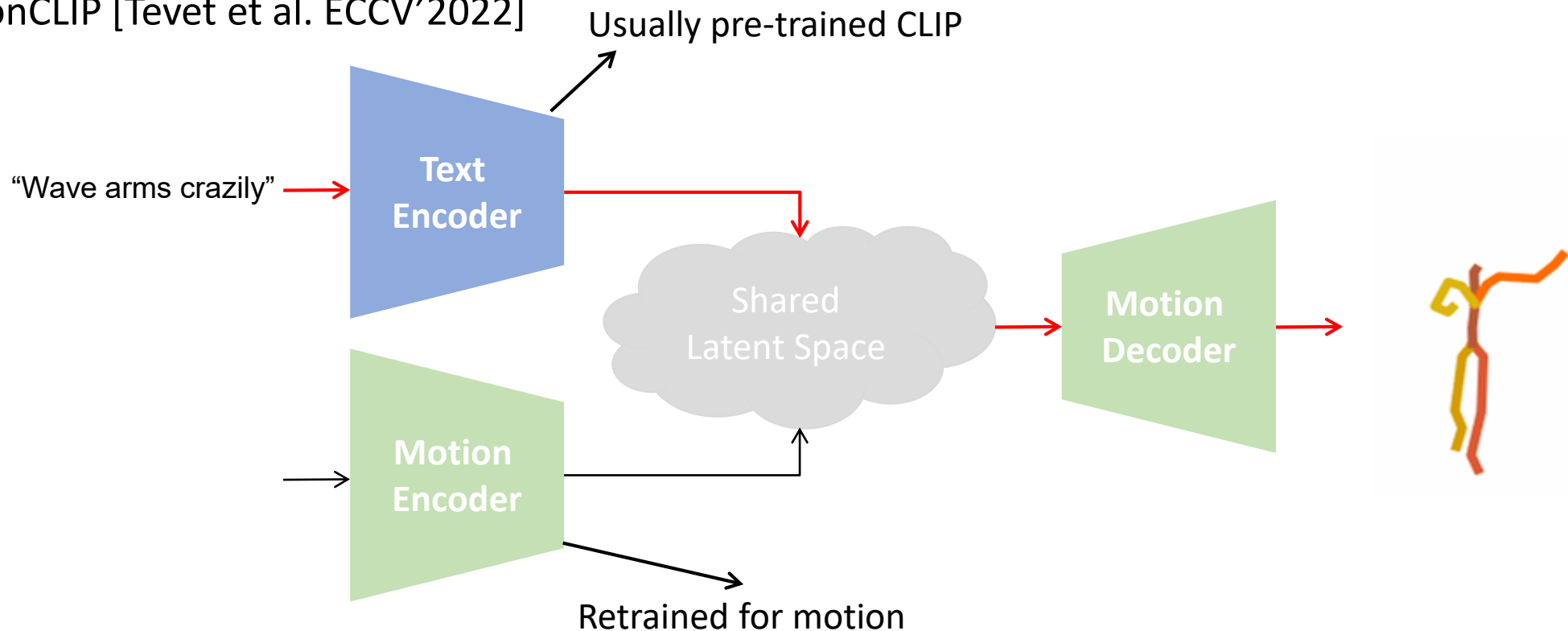
Main contents

- **Human motion diffusion model.**
 - **Text to motion generation.**
- Compositional motion generation with pretrained motion prior.
 - Long sequence.
 - Two-person interaction.
 - Trajectory control.
- Unified human motion synthesis and understanding with fine-grained semantics.
 - Bidirectional: text to motion and motion to text.
 - Hierarchical semantics: global and local text.

Classic way for text to motion generation

Text-to-motion using the VAE framework:

- TEMOS [Petrovich et al. ECCV'2022]
- T2M [Guo et al. CVPR'2022]
- MotionCLIP [Tevet et al. ECCV'2022]



Human Motion Diffusion Model

Guy Tevet

Sigal Raab

Brian Gordon

Yonatan Shafir

Daniel Cohen-Or

Amit H. Bermano

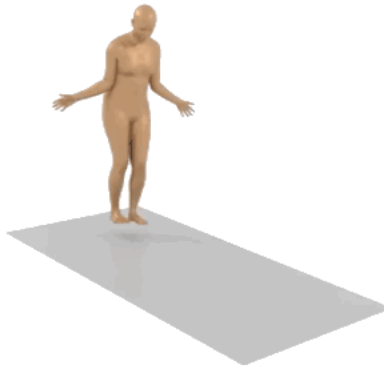
Tel Aviv University, Israel



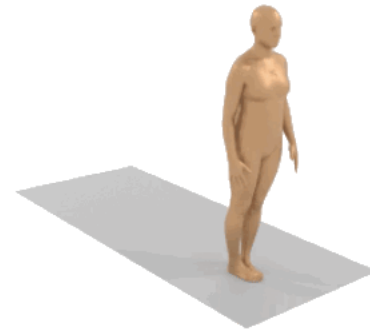
“a man kicks with something or someone with his left leg.”



“A person punches in a manner consistent with martial arts.”



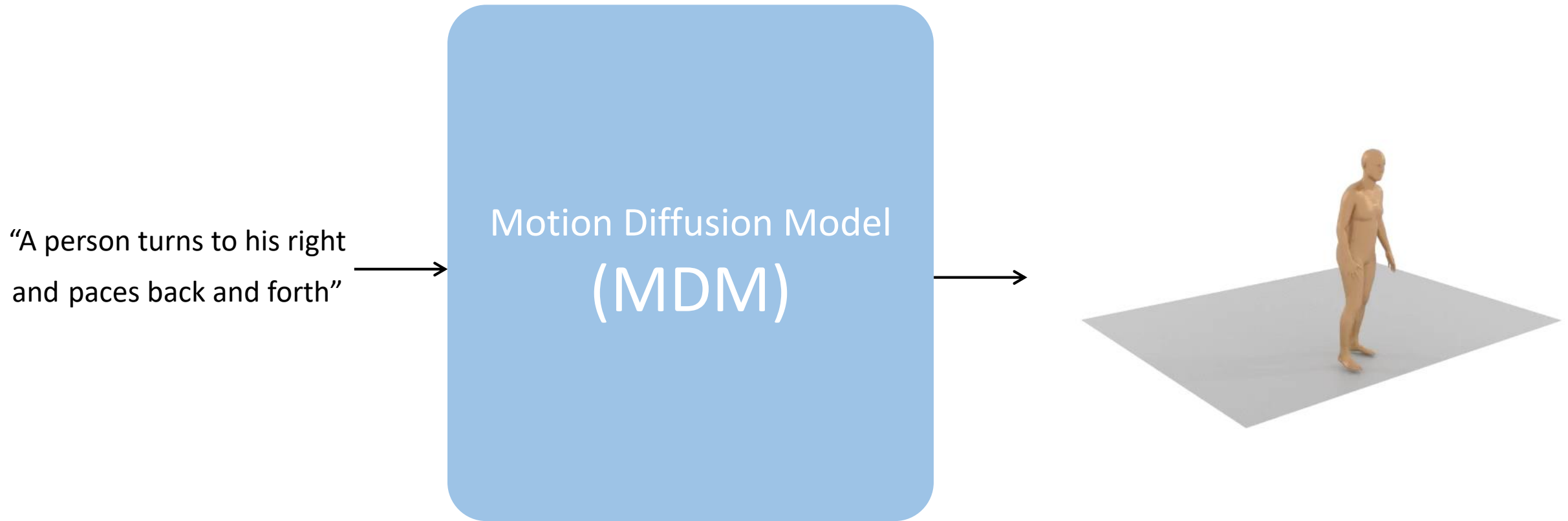
“A person is skipping rope.”



“a person walks backwards slowly.”

MDM: A Human Motion Framework

Goal: given text, generate plausible human motion.



MDM: A Human Motion Framework

High Quality

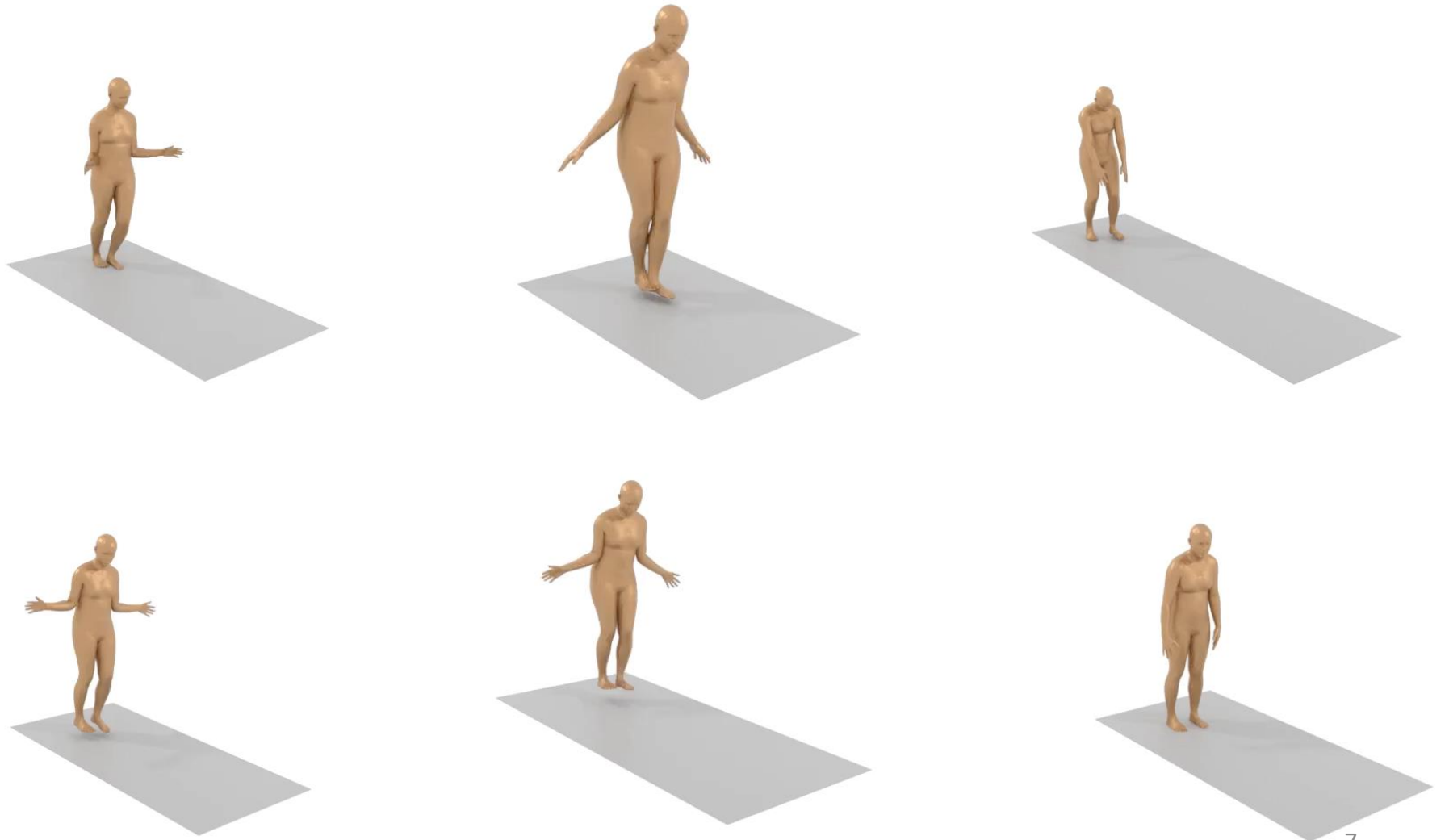


Global Position

MDM: A Human Motion Framework

Variability

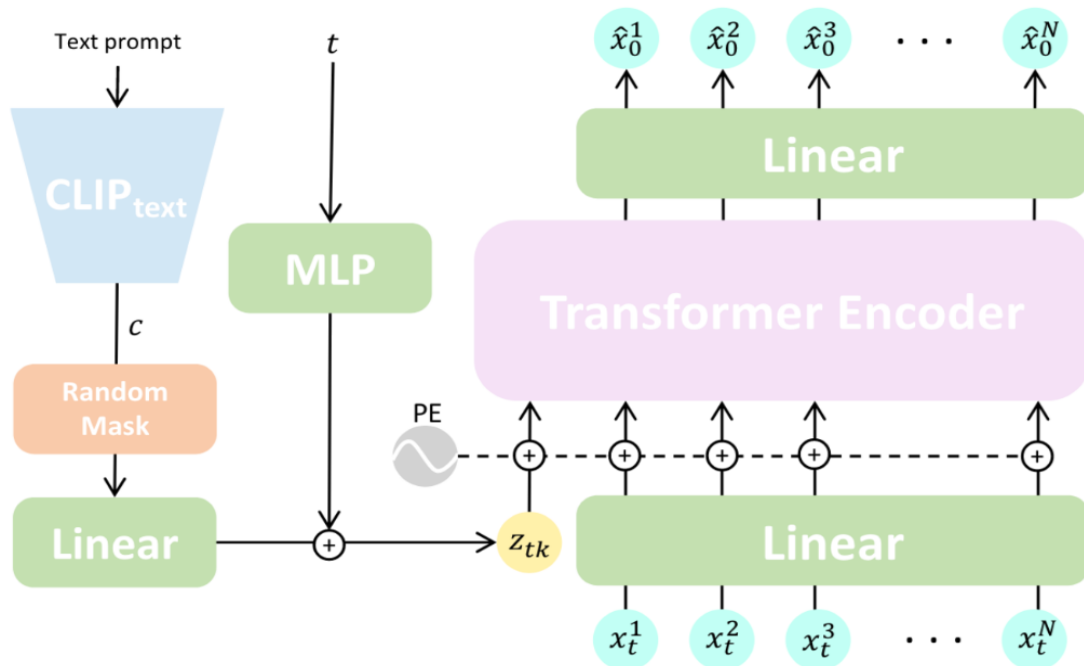
“A person is skipping rope.”



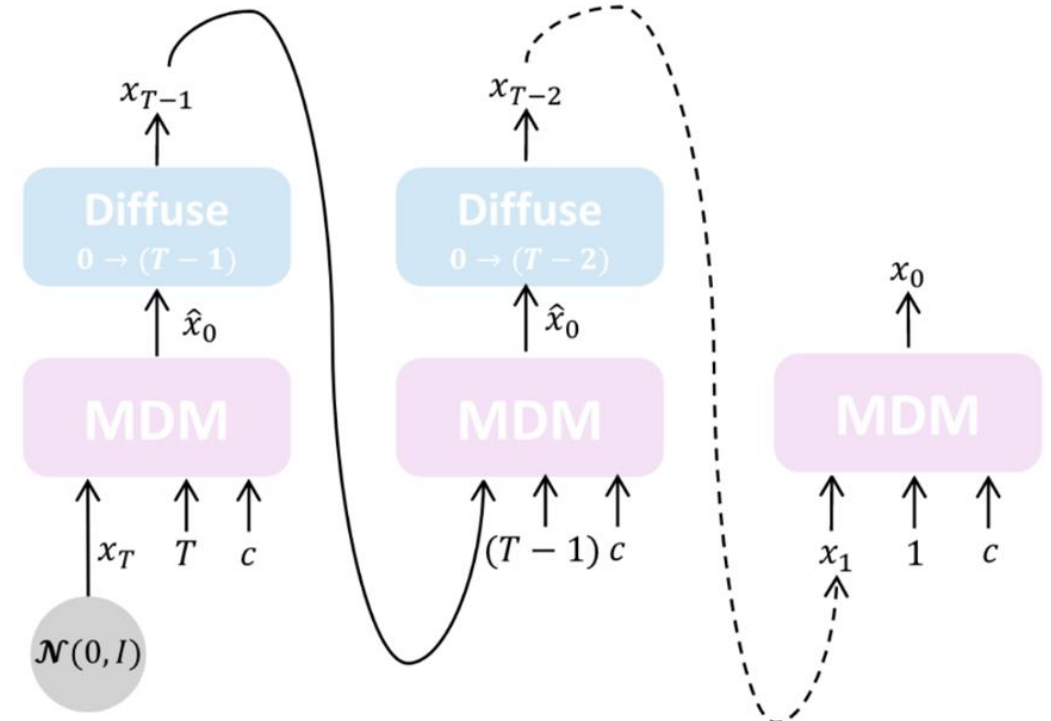
MDM framework

- CLIP text encoder + transformer based diffusion.
- Classifier-free guidance:
 - Training: randomly mask out text conditions.
 - Sampling: weighted combination of conditional and unconditional predictions.

MDM network architecture

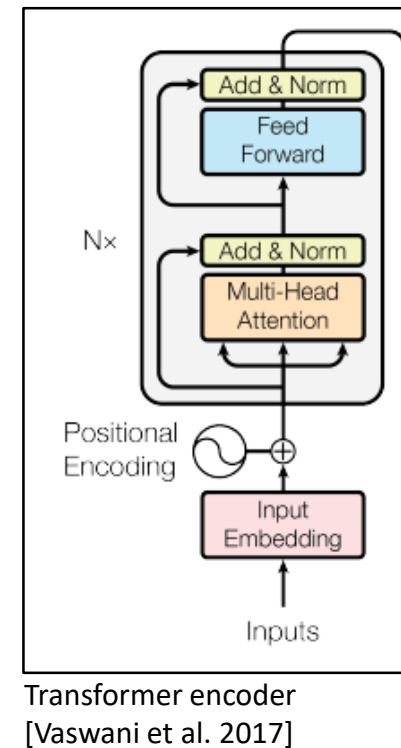
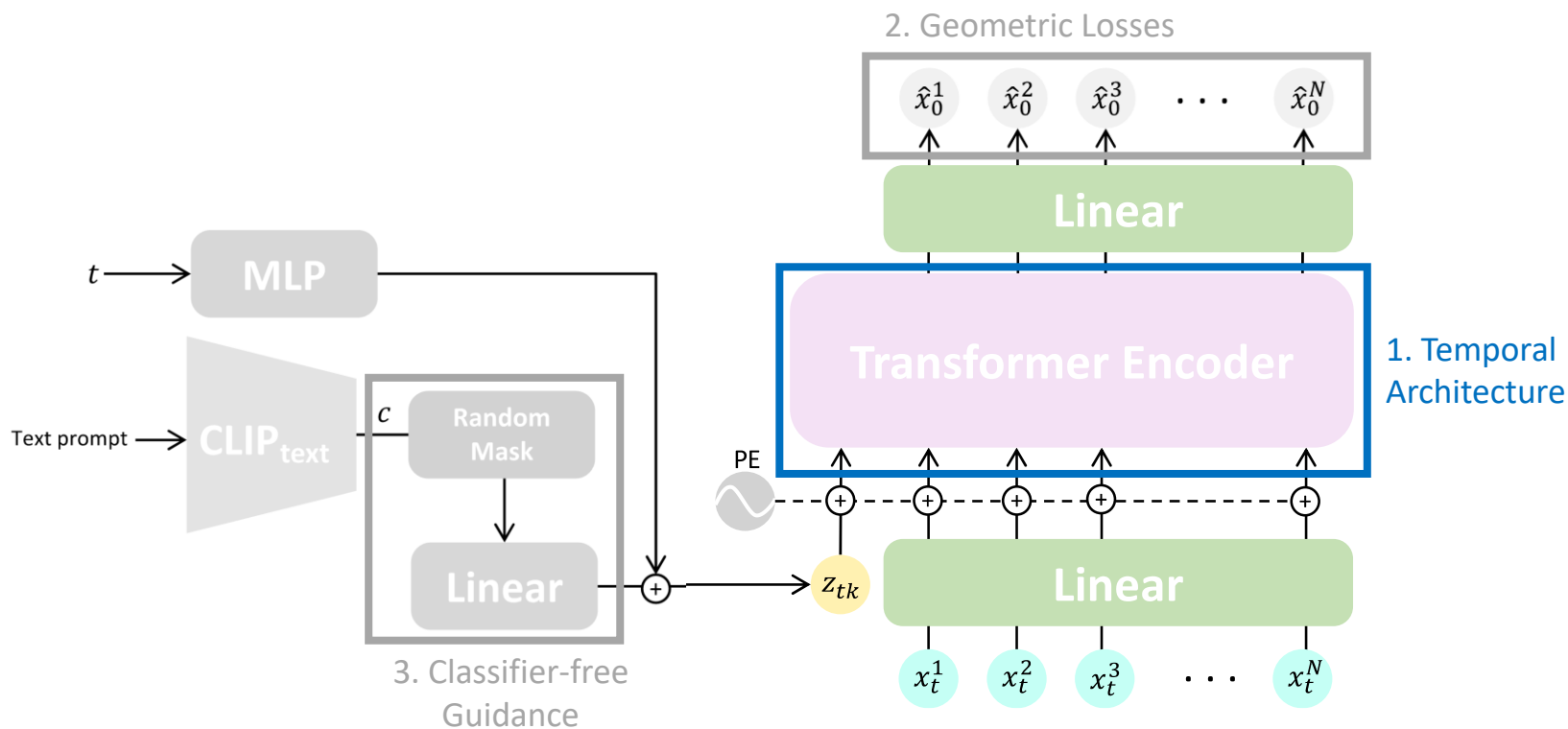


MDM reverse sampling



MDM architecture in more detail

- 1. Temporal architecture: encode time information via transformer.
- 2. Geometric representation $x_0^i \in \mathbb{R}^{J \times D}$: joint rotations or positions.



MDM training losses

- Diffusion loss on x_0 and geometric losses.

$$\mathcal{L} = \mathcal{L}_{\text{simple}} + \lambda_{\text{pos}} \mathcal{L}_{\text{pos}} + \lambda_{\text{vel}} \mathcal{L}_{\text{vel}} + \lambda_{\text{foot}} \mathcal{L}_{\text{foot}}$$

Geometric loss terms

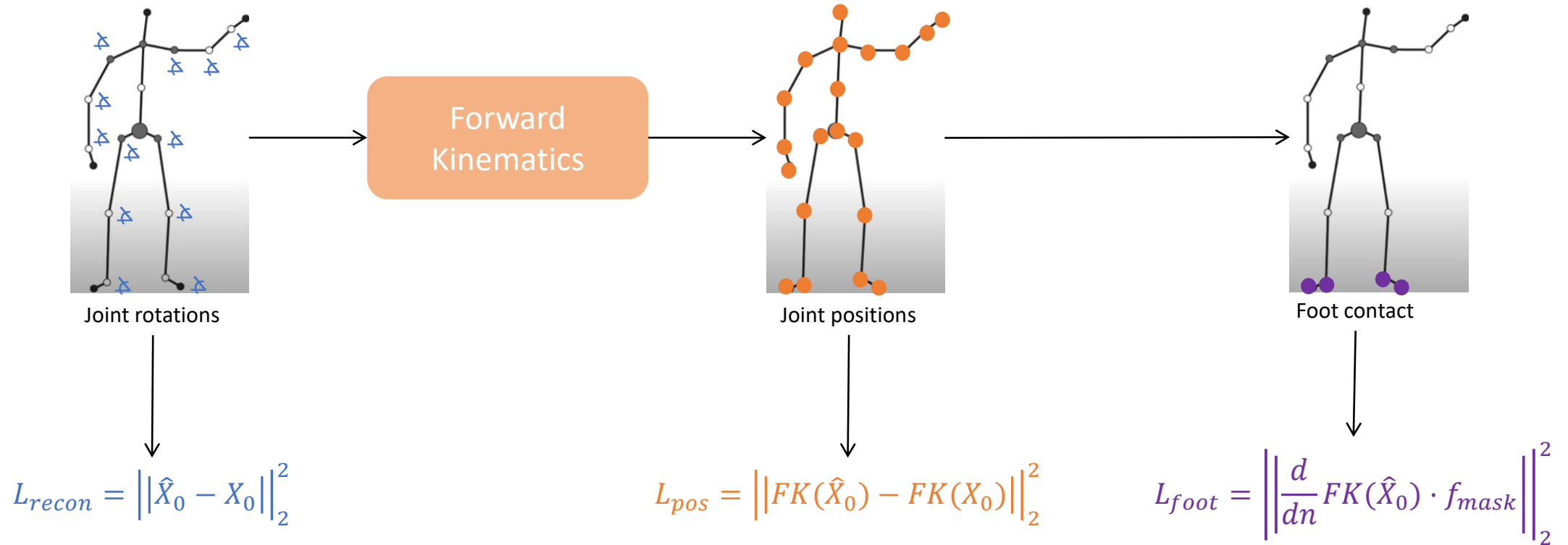
$$\mathcal{L}_{\text{simple}} = E_{x_0 \sim q(x_0|c), t \sim [1, T]} [\|x_0 - G(x_t, t, c)\|_2^2] \longrightarrow \text{Classic diffusion loss}$$

$$\mathcal{L}_{\text{pos}} = \frac{1}{N} \sum_{i=1}^N \|FK(x_0^i) - FK(\hat{x}_0^i)\|_2^2, \longrightarrow \text{Joint position after forward kinematics}$$

$$\mathcal{L}_{\text{foot}} = \frac{1}{N-1} \sum_{i=1}^{N-1} \|(FK(\hat{x}_0^{i+1}) - FK(\hat{x}_0^i)) \cdot f_i\|_2^2, \longrightarrow \text{Foot contact } f_i \in \{0, 1\}^J$$

$$\mathcal{L}_{\text{vel}} = \frac{1}{N-1} \sum_{i=1}^{N-1} \|(x_0^{i+1} - x_0^i) - (\hat{x}_0^{i+1} - \hat{x}_0^i)\|_2^2 \longrightarrow \text{Velocity}$$

Geometric losses visualization



Geometric Losses - Results

Warm-up



Warm-up



Sit



MDM
- without geometric
losses

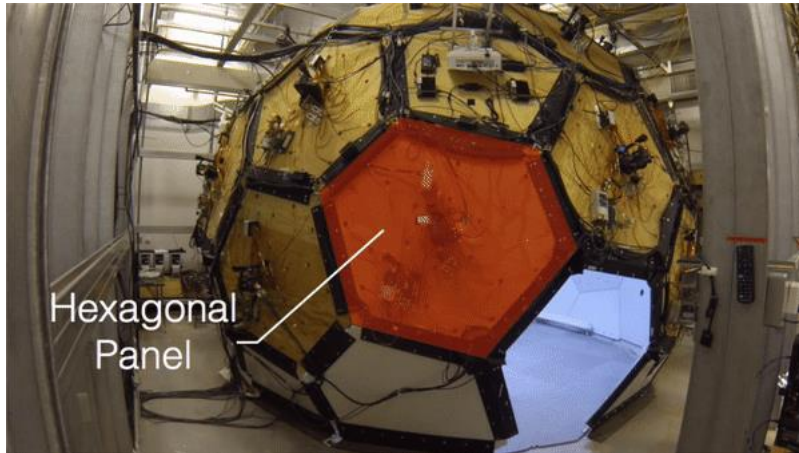
MDM sampling: classifier-free guidance

- Iterative reverse sampling with classifier-free guidance.
- Unconditional prediction: $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) = -\frac{1}{\sqrt{1-\bar{\alpha}_t}} \epsilon_{\theta}(\mathbf{x}_t, t)$
- With condition y : $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t, y) = \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) + \nabla_{\mathbf{x}_t} \log p(y|\mathbf{x}_t)$
- Guidance:
 - Training time: randomly mask out condition y . MDM: $y =$ the CLIP latent.

$$\begin{aligned}\nabla_{\mathbf{x}_t} \log p(y|\mathbf{x}_t) &= \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|y) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) \\ &= -\frac{1}{\sqrt{1-\bar{\alpha}_t}} \left(\underbrace{\epsilon_{\theta}(\mathbf{x}_t, t, y)}_{\text{Conditional}} - \underbrace{\epsilon_{\theta}(\mathbf{x}_t, t)}_{\text{Unconditional}} \right) \\ \bar{\epsilon}_{\theta}(\mathbf{x}_t, t, y) &= \epsilon_{\theta}(\mathbf{x}_t, t, y) - \sqrt{1-\bar{\alpha}_t} w \nabla_{\mathbf{x}_t} \log p(y|\mathbf{x}_t) \quad \text{Guidance scale } w \\ &= \epsilon_{\theta}(\mathbf{x}_t, t, y) + w(\epsilon_{\theta}(\mathbf{x}_t, t, y) - \epsilon_{\theta}(\mathbf{x}_t, t)) \\ &= (w+1)\epsilon_{\theta}(\mathbf{x}_t, t, y) - w\epsilon_{\theta}(\mathbf{x}_t, t)\end{aligned}$$

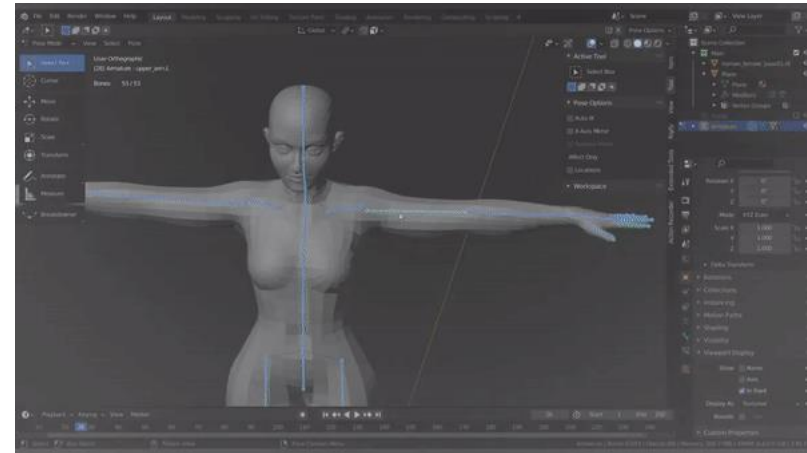
The Data: real capture or artist design.

CMU Panopticon – Motion Capture [Joo et al. 2015]



https://www.youtube.com/watch?v=zQt6g-Jel7M&ab_channel=HanbyulJoo

Real motion capture



https://www.youtube.com/watch?v=antc20EFh6k&t=24s&ab_channel=kfiraberman

Artist design

The Data

HumanML3D [Guo et al. 2022]

- 15K examples
- ~7 sec each



BABEL [Punnakkal et al. 2021]

MDM results: motion quality

Walk

Run

Jump

MDM



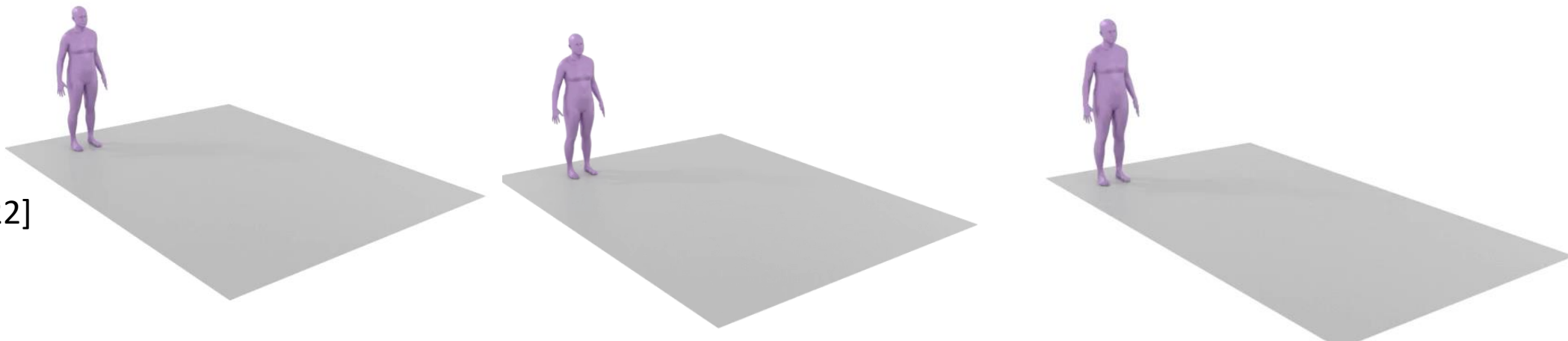
ACTOR
[Petrovich et al., 2021]



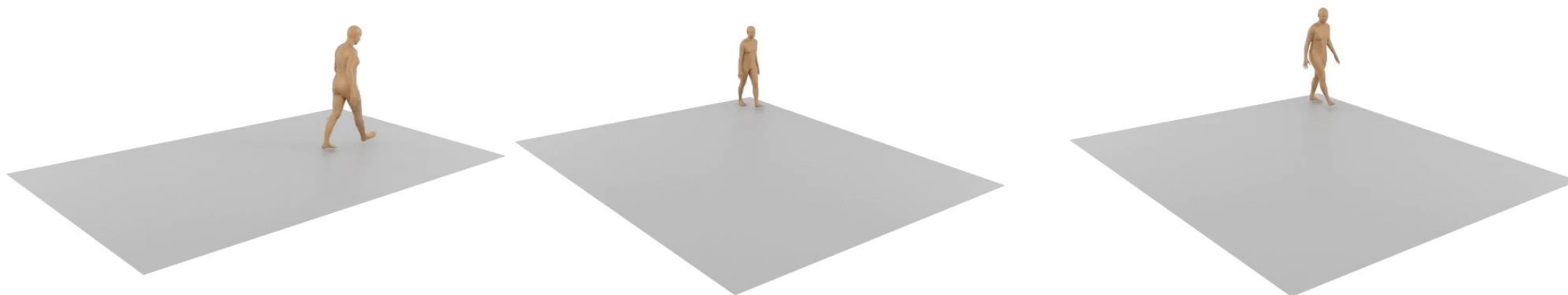
MDM results: motion diversity

“Person walking in an s shape”

T2M
[Guo et al. 2022]



MDM



Main contents

- Human motion diffusion model.
 - Text to motion generation.
- **Compositional motion generation with pretrained motion prior.**
 - Long sequence.
 - Two-person interaction.
 - Trajectory control.
- Unified human motion synthesis and understanding with fine-grained semantics.
 - Bidirectional: text to motion and motion to text.
 - Hierarchical semantics: global and local text.

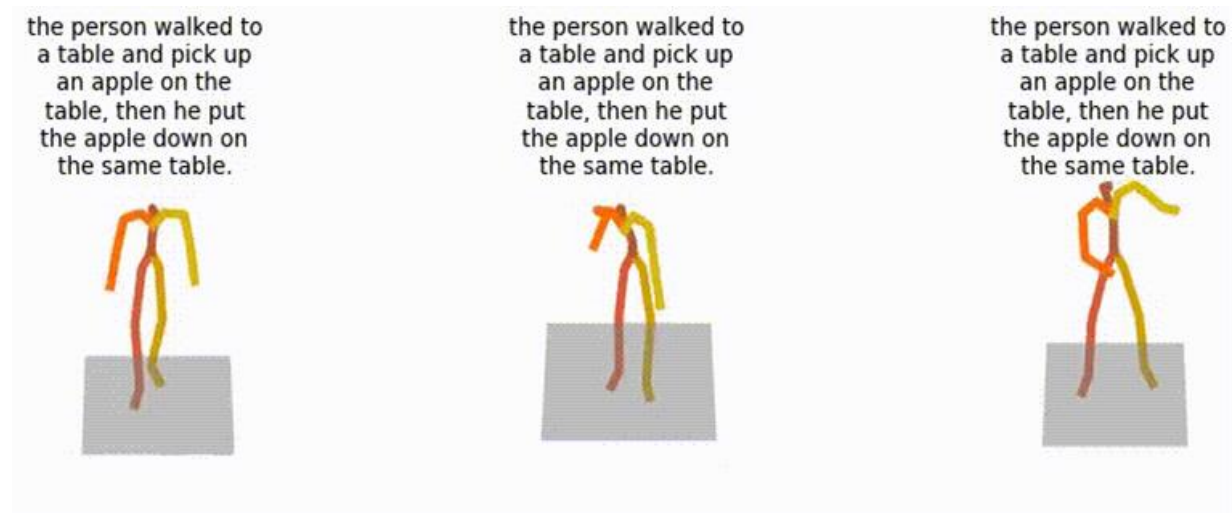
MDM limitations

- Maximum generation length 196 frames (9.8 seconds).
- No condition on the location.
- Results not satisfying for human-object-interaction prompts.

No spatial control



Cannot generate complex long motion

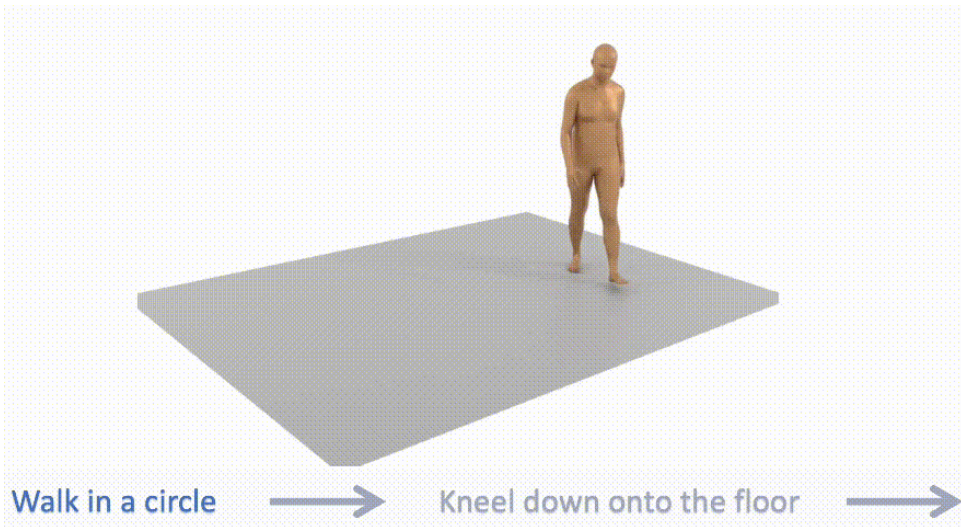
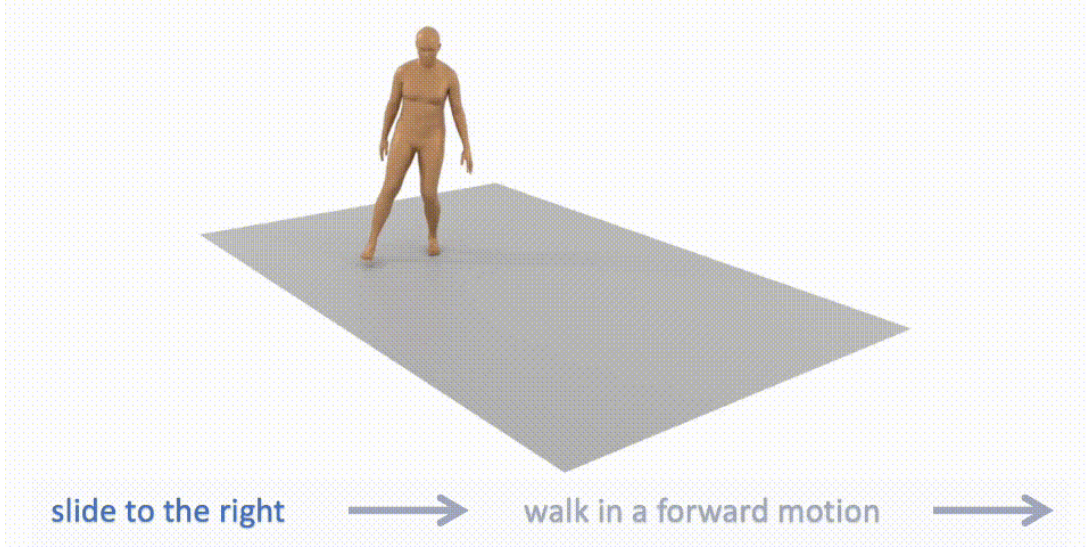
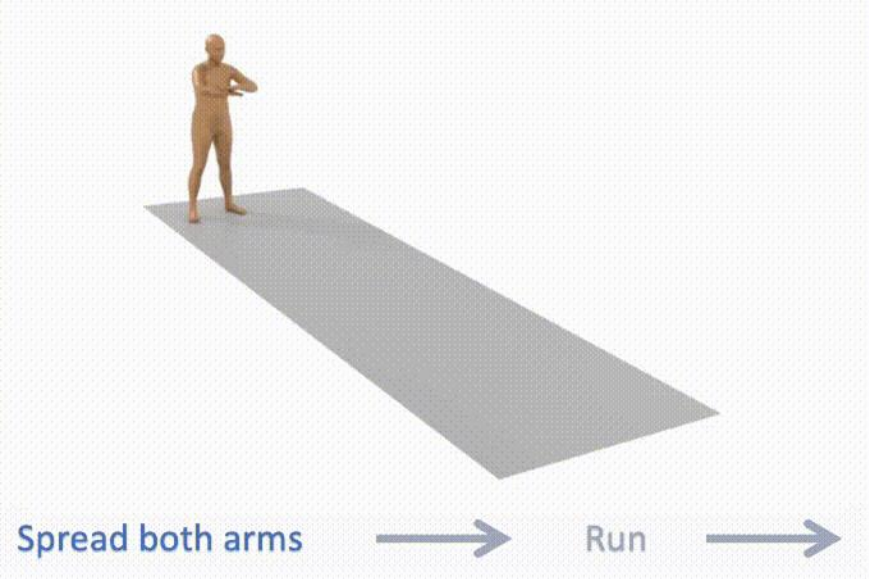


The data problem again

- Motion data is very expensive to obtain:
 - Real motion capture.
 - Designed by artists.
 - Training data for MDM: almost exclusively of short, single person sequences.
- Complex motions are compositional:
 - Very long motions.
 - Interaction between humans, or human and object.
 - Diverse control signals like text and spatial.
 - Can we use compose motions from pretrained motion priors?

PriorMDM: Human Motion Diffusion as a Generative Prior

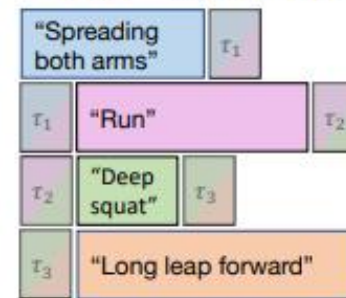
Long Motions



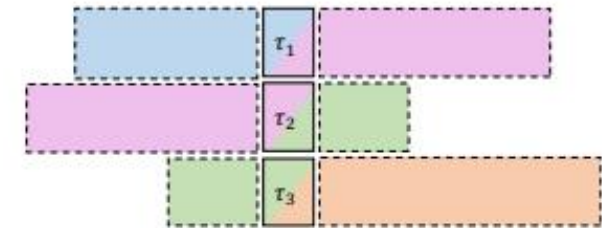
DoubleTake: composition in temporal domain

- Independent denoising + mixed diffusion in transition regions.
- Each temporal window is conditioned on different text.
- Communicate between segments via handshake ($\sim 1s/window$):
 - Inside the handshake τ : average of previous suffix and current prefix.
- Second take: add noise to handshake part and denoise back.

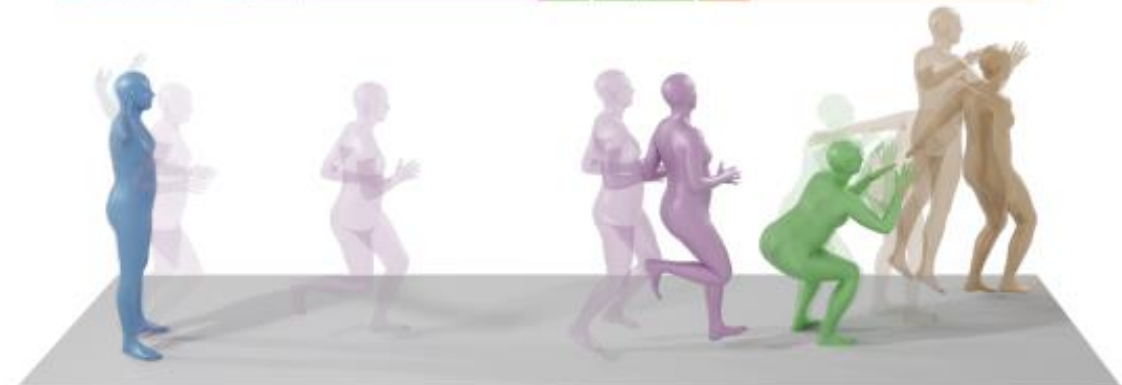
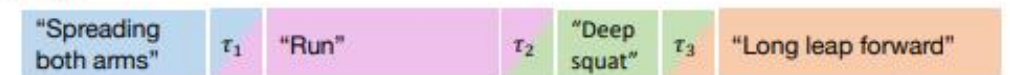
Take #1 – Handshake Generation



Take #2 – Transition Refinement



Unfolding

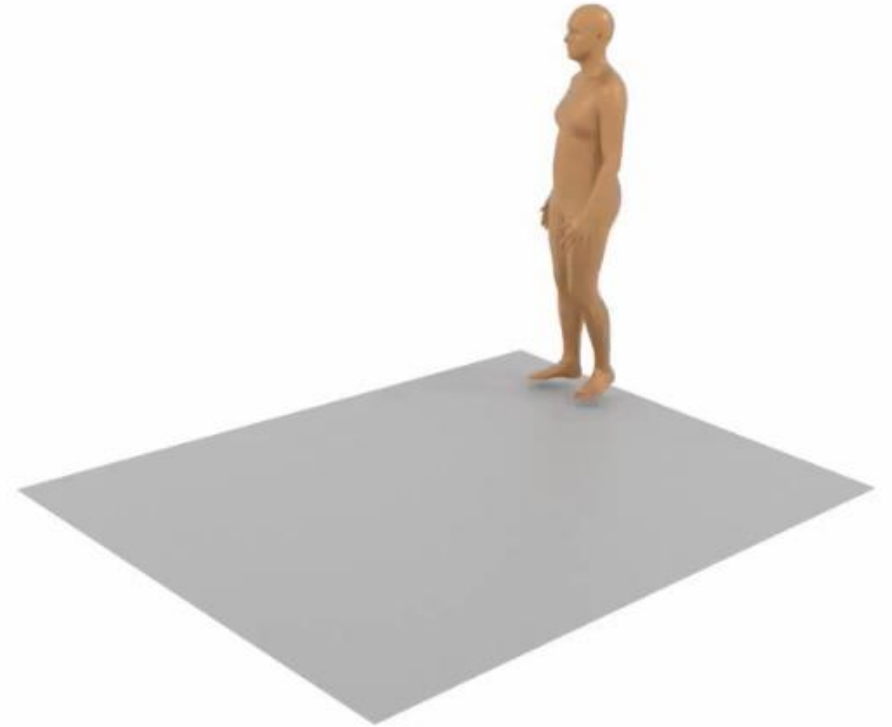


DoubleTake results

TEACH [Athanasiou et al. 2022]



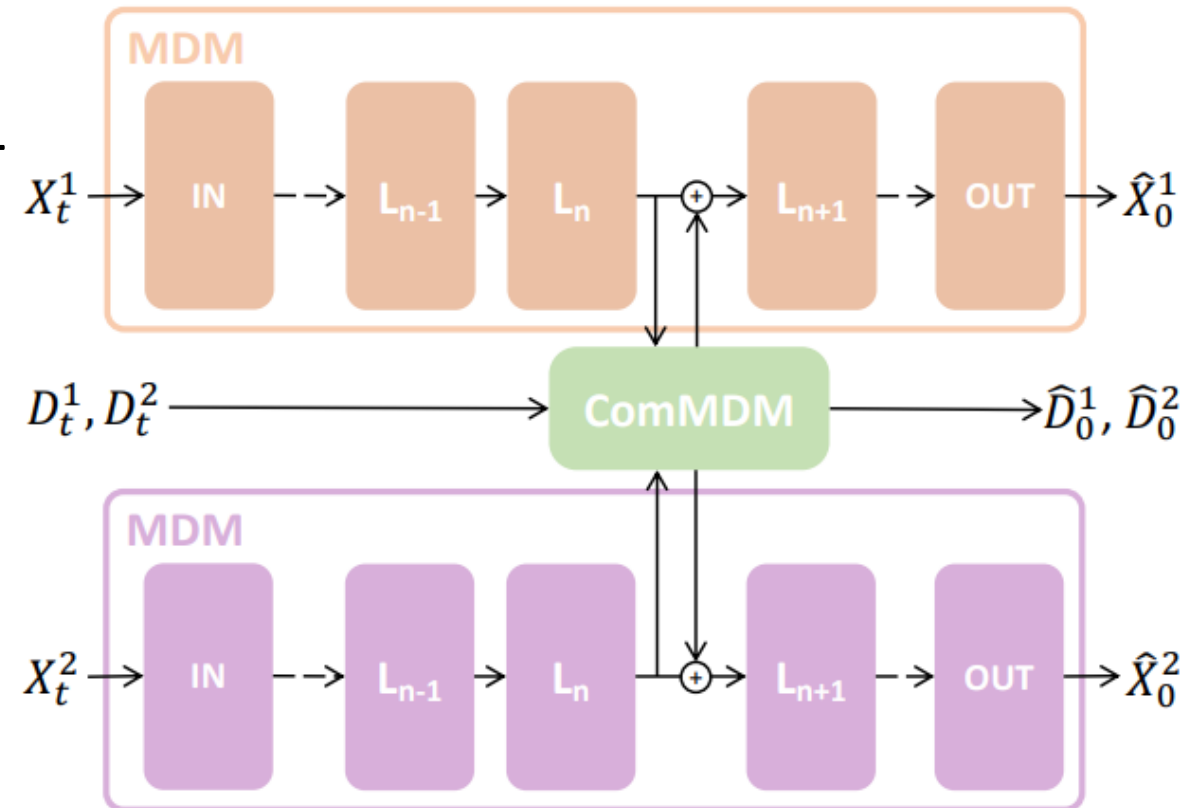
DoubleTake (Ours)



walk in a circle → stand → walk → reach out and shake right hand

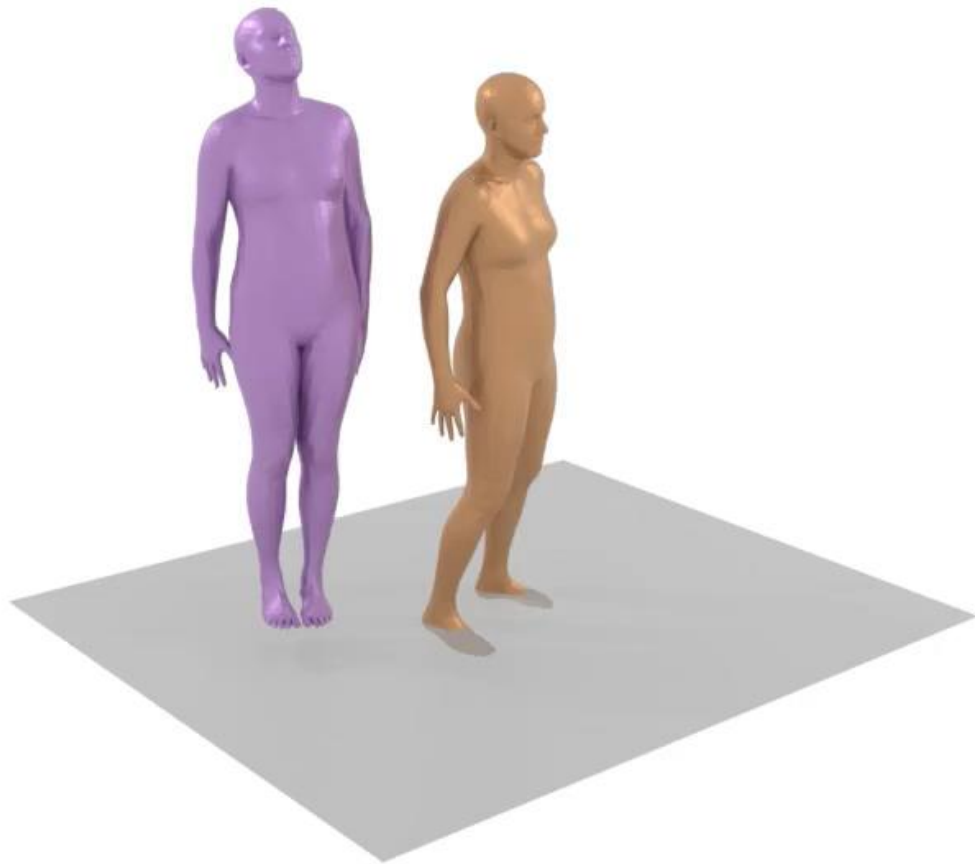
Two person motion generation

- Compositional: keep individual models and learn the interaction.
- ComMDM: learns to correct the output from individual models.
- Few-shot learning: trained on 55 two-person motion sequences.

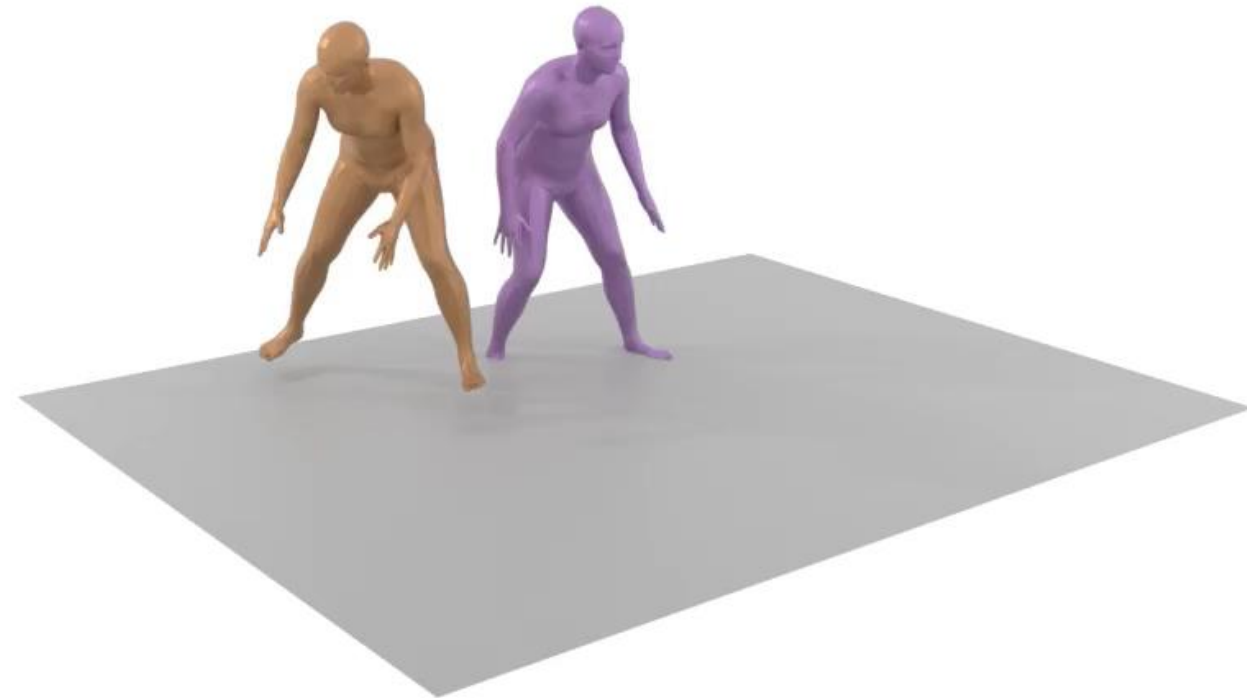


ComMDM for two person interaction

“A Capoeira practice. One is kicking and the other is avoiding the kick.”



“The two people are playing basketball, one with the ball the other is defending.”



MDM fine tuned for trajectory control

- Goal: control the motion with given trajectory.
- Training: fine tune the model to denoise only free body parts.
- Sampling: set hard constraint on the given trajectory.

Algorithm 1 Fine-tuning method

```
repeat  
   $x_0 \sim q(x_0)$   
   $t \sim \text{Uniform}(\{1, \dots, T\})$   
   $\epsilon \sim \mathcal{N}(0, I)$   
   $\epsilon[\textit{trajectory}] = 0$  ▷ Our addition  
  Take gradient descent step on:  
     $\nabla_{\theta} \|x_0 - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t)\|$   
until Converged
```

Algorithm 2 Sampling method

```
 $x_0^{(T)} = 0$   
for  $t = T, \dots, 0$  do  
   $x_0^{(t)}[\textit{trajectory}] = \textit{given trajectory}$  ▷  
  Original in-painting  
   $\epsilon \sim \mathcal{N}(0, I)$   
   $\epsilon[\textit{trajectory}] = 0$  ▷ Our addition  
   $x_0^{(t-1)} = \epsilon_{\theta}(\sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t)$   
end for
```

MDM fine tuned for trajectory control

MDM [Tevet et al . 2022]



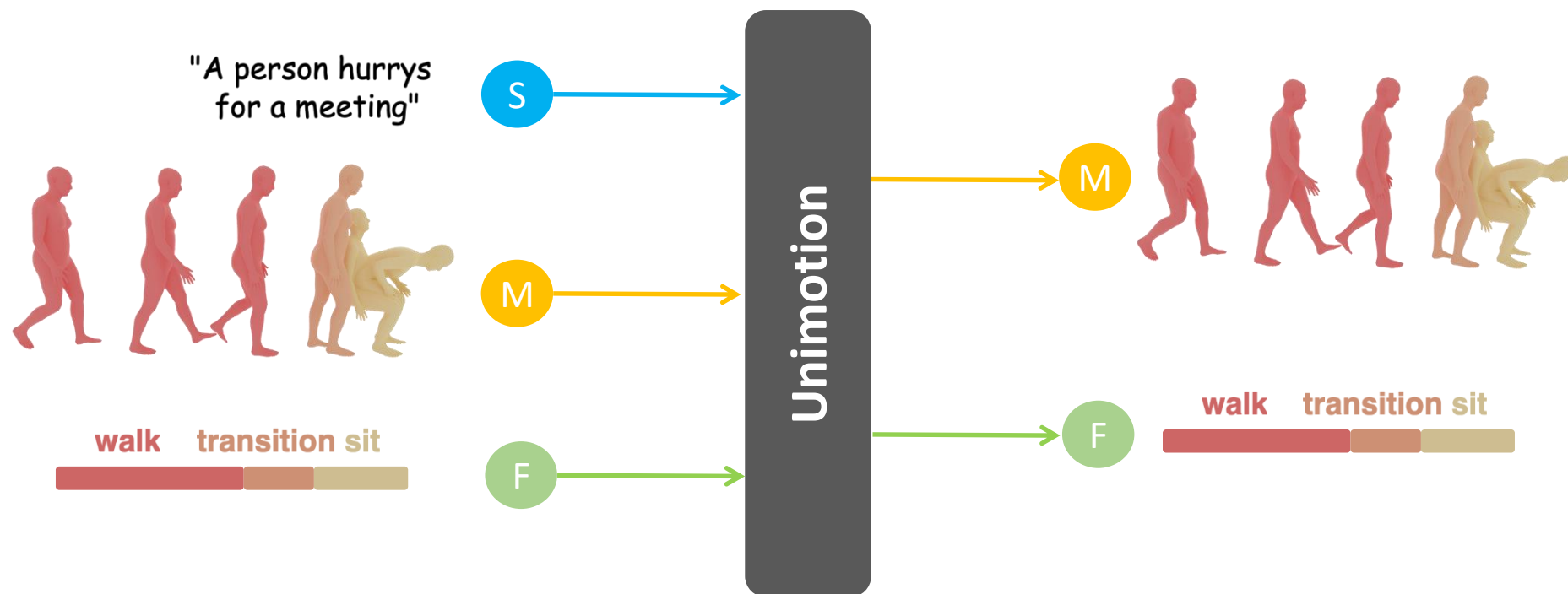
Fine-tuned MDM (Ours)



Main contents

- Human motion diffusion model.
 - Text to motion generation.
- Compositional motion generation with pretrained motion prior.
 - Long sequence.
 - Two-person interaction.
 - Trajectory control.
- **Unified human motion synthesis and understanding with fine-grained semantics.**
 - **Bidirectional: text to motion and motion to text.**
 - **Hierarchical semantics: global and local text.**

Goal: Unify **motion tasks** into a single model



Multi-tasks

- Sequence-level Text-to-Motion Generation

"A person walks forward, bends down to pick something up off the ground."

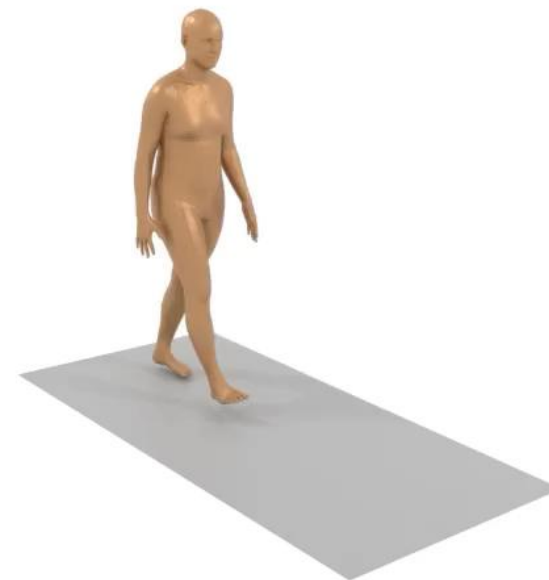


S

M



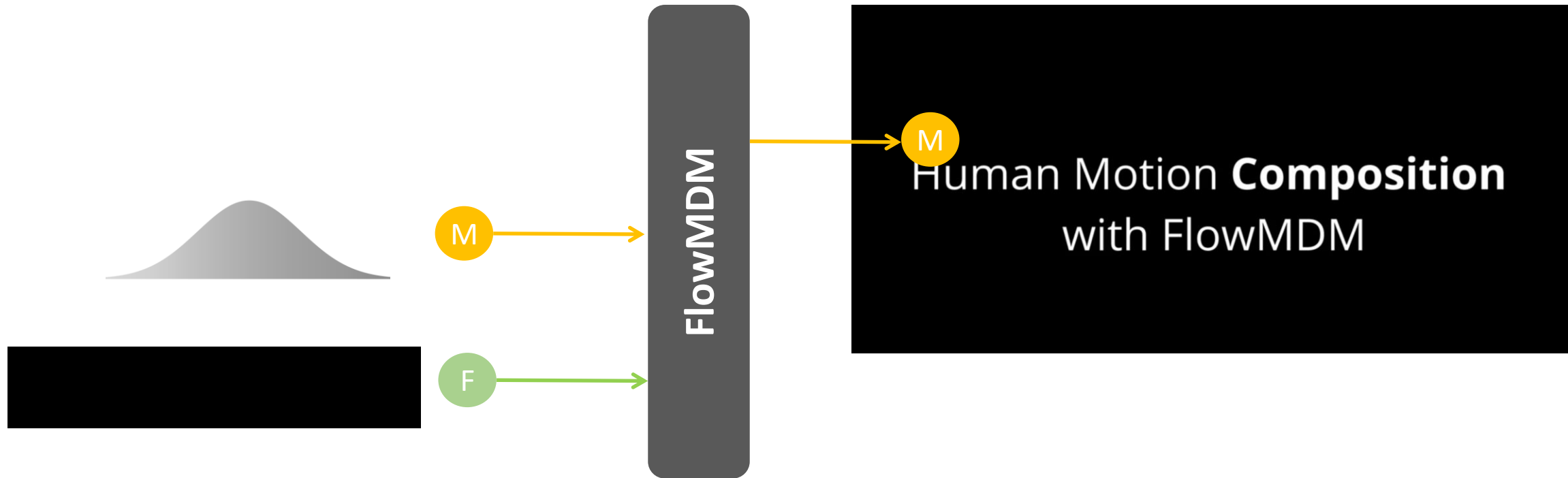
M



MDM: Tevet et al, ICLR' 23

Multi-tasks

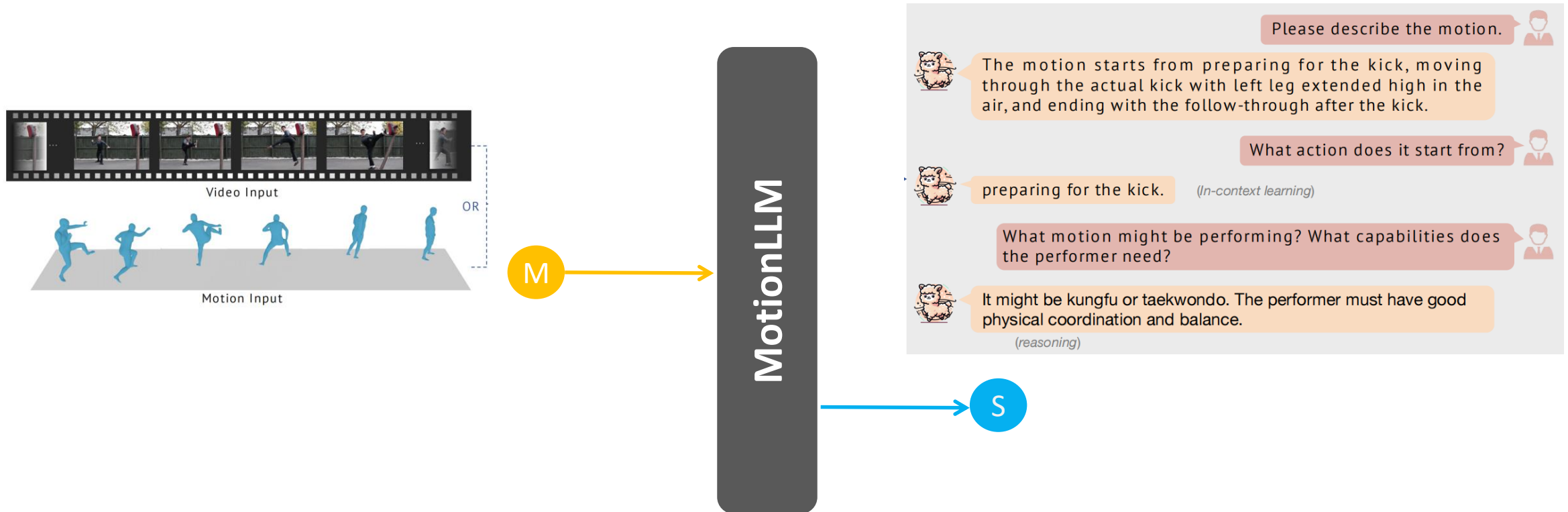
- Sequence-level Text-to-Motion Generation
- Frame-level Text-to-Motion Generation



FlowMDM: Barquero et al, CVPR' 24

Multi-tasks

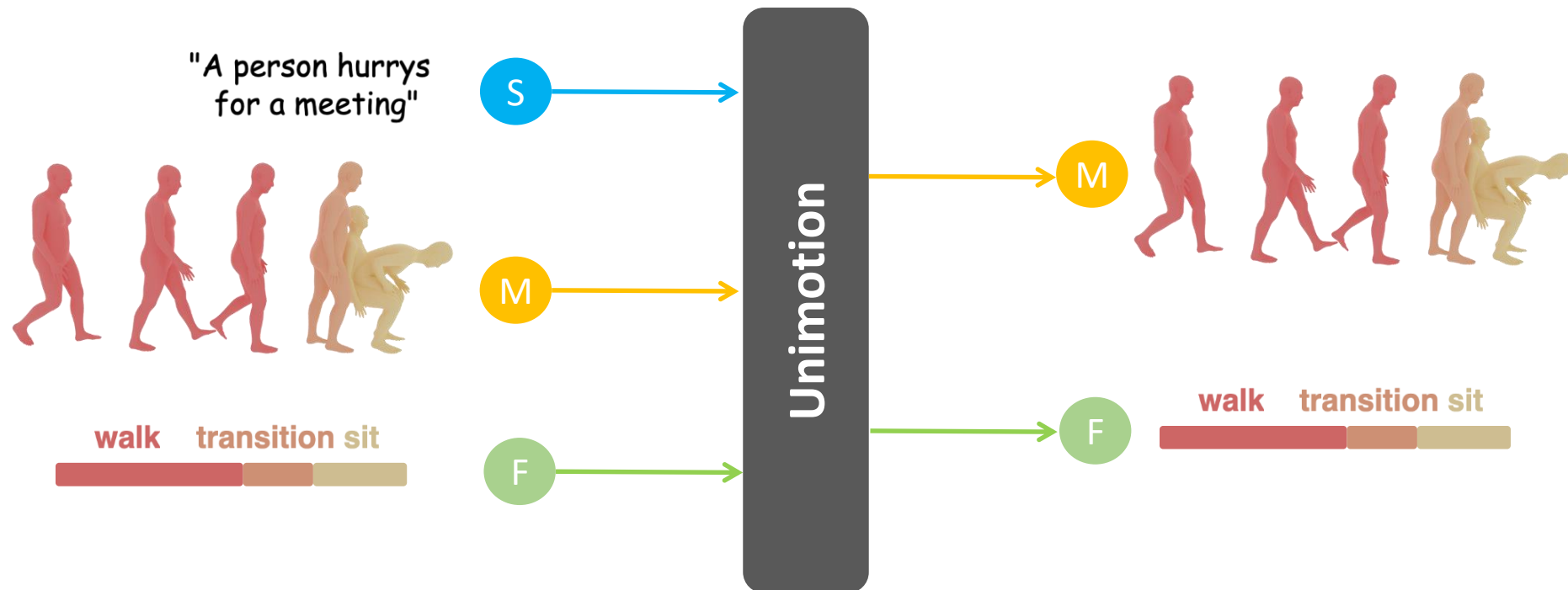
- Sequence-level Text-to-Motion Generation
- Frame-level Text-to-Motion Generation
- **Motion-to-Text Understanding**



MotionLLM: ArXiv'24

Goal: Unify **motion tasks** into a single model

Our model unifies all these tasks into a single model allowing for **multi-model input**.



Additionally, its flexibility allows for novel tasks, not yet performed by prior works.

Novel tasks

- Hierarchical Text-to-Motion Generation

"A person waves hands above head ."



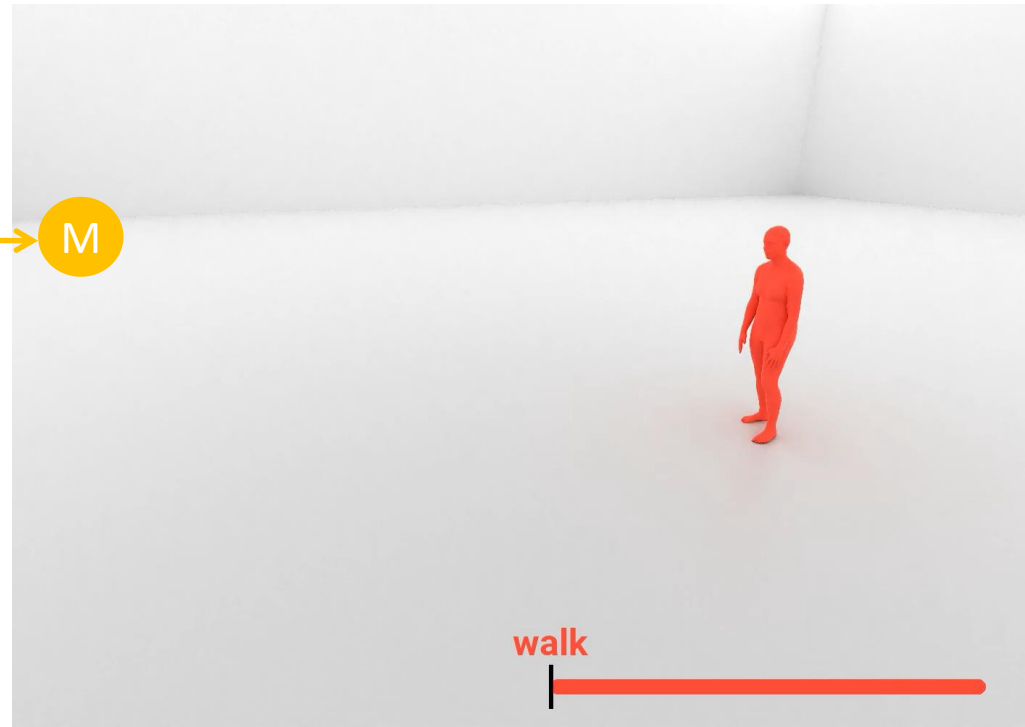
S

M

F

Unimotion

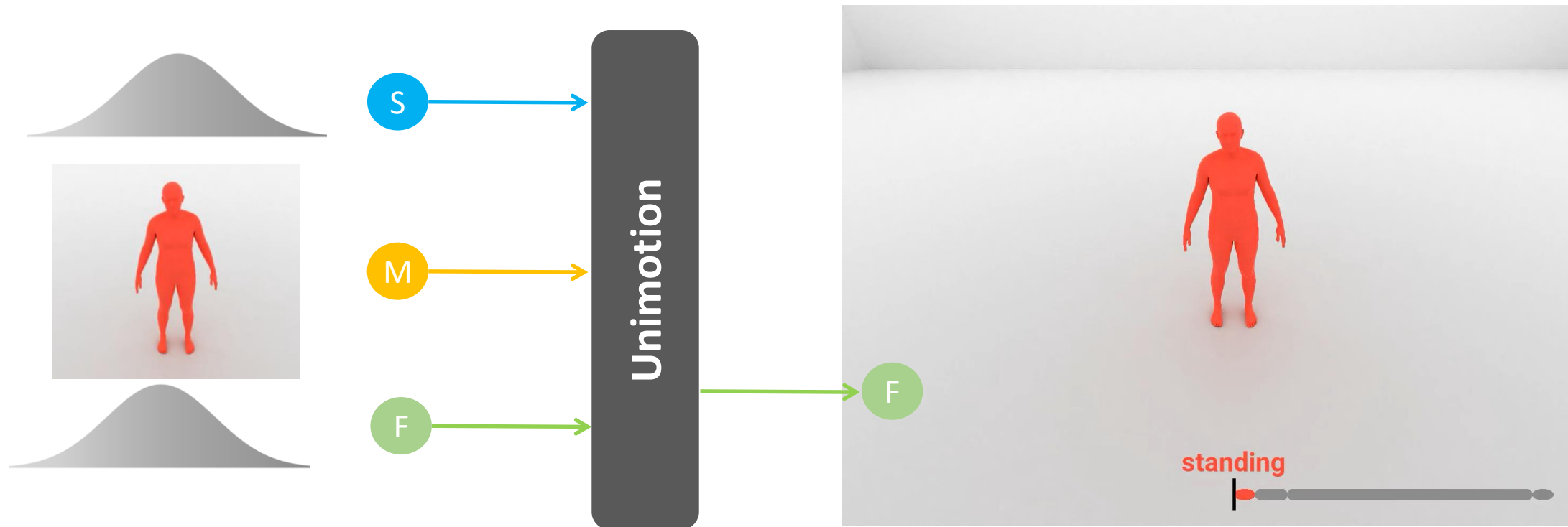
M



Novel tasks

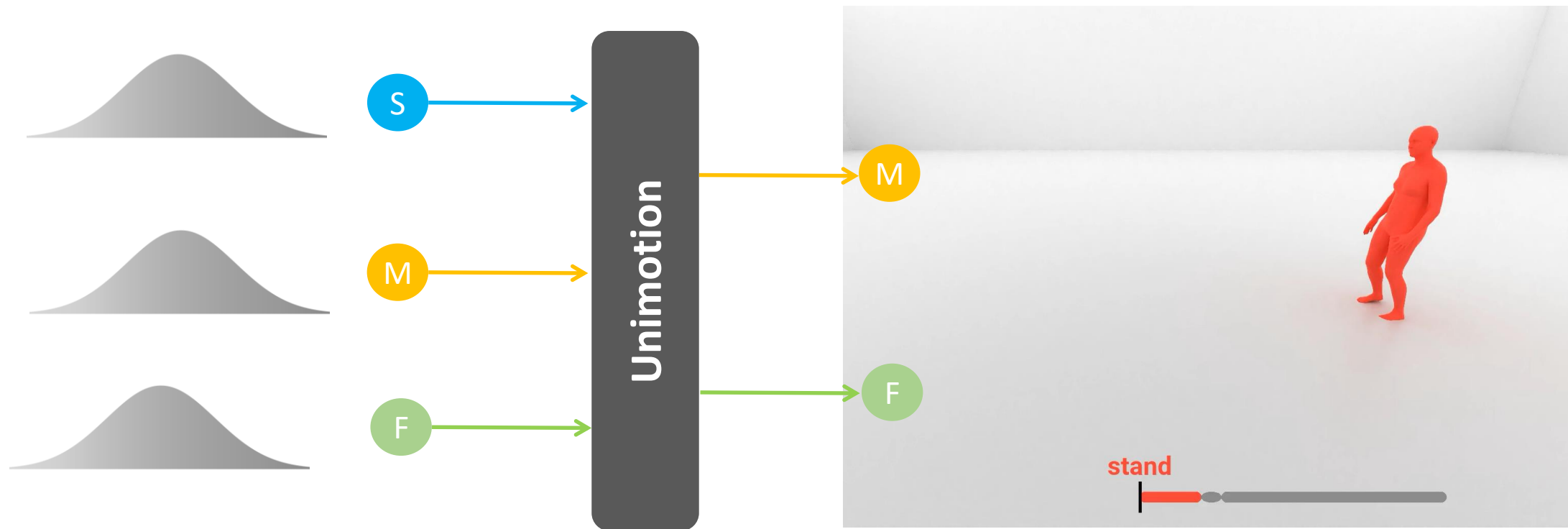
- Hierarchical Text-to-Motion Generation
- Motion-to-Text Understanding

Fine-grained with temporal alignment!



Novel tasks

- Hierarchical Text-to-Motion Generation
- Motion-to-Text Understanding
- Unconditional Joint Generation

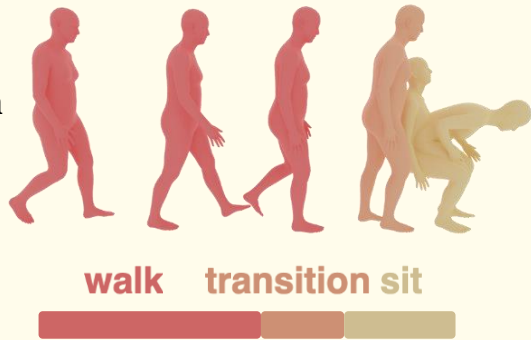


Method

Input

Noised (Motion, Frame-level text) +
Seq-level text

Clean motion



Frame-level text

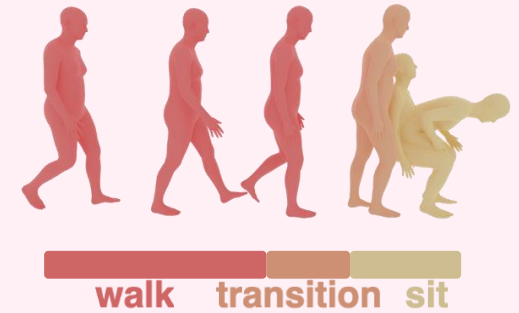
Seq-level text

"A person hurrys
for a meeting"

Model

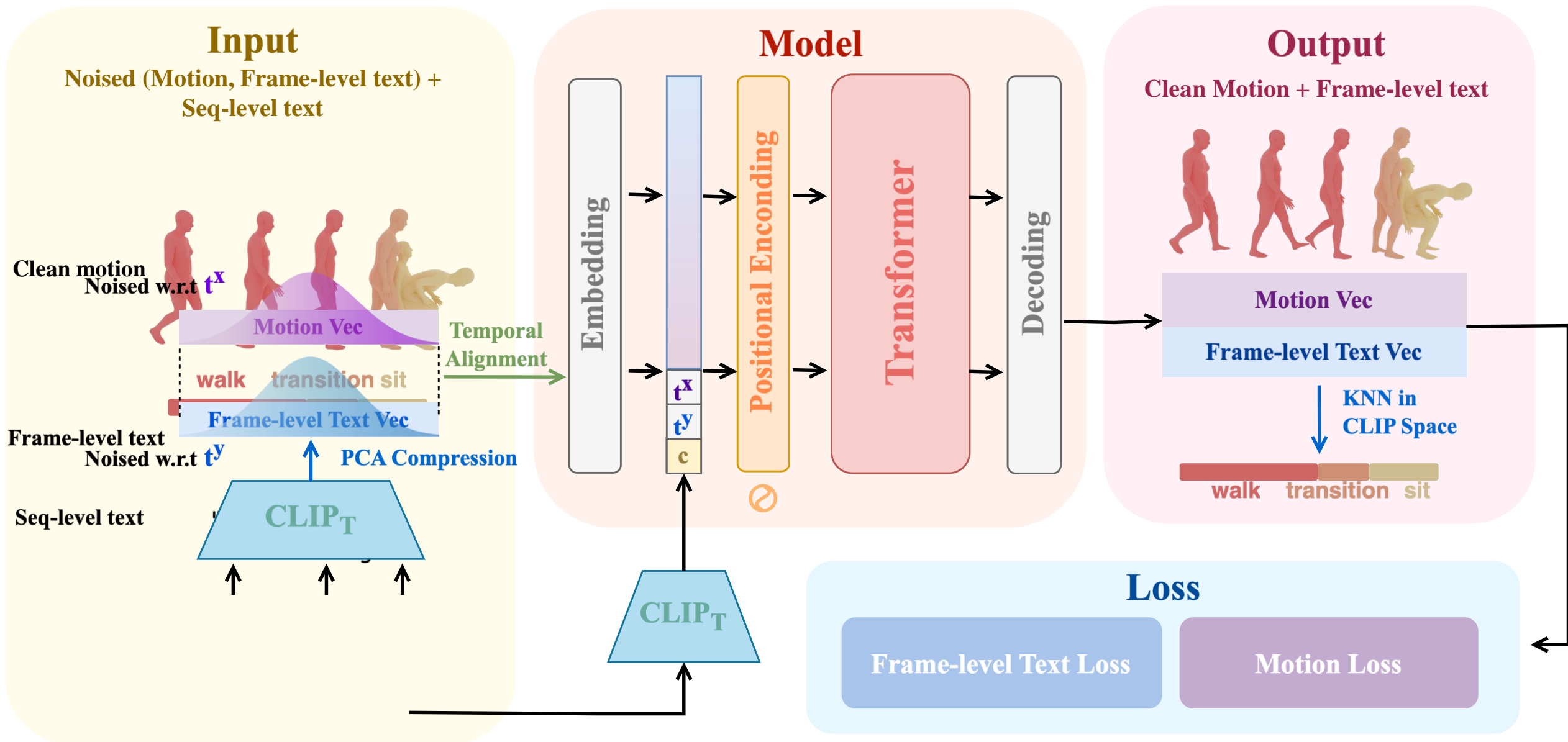
Output

Clean Motion + Frame-level text




Loss

Method




Multi-tasks

Hierarchical Text to Motion

Input: Seq-level + Frame-level Text
"A person waves hands above head"
3.sitting down + 1.walking
2.transition ↓

Output Motion

Motion to Text

Input Motion

Output: Frame-level Text ↓
4.walking 2.step up
3.step down 1.walking

Unconditional Joint Generation

Input: Noise Vec → 
+
2.jump forward
3.turn around 1.walking forward
Output: Motion + Frame-level Text

Motion Generation and Editing (Combined Task)

Input: Seq-level Text
"This person walks forward and picks something up then walks back" → 
Output Motion

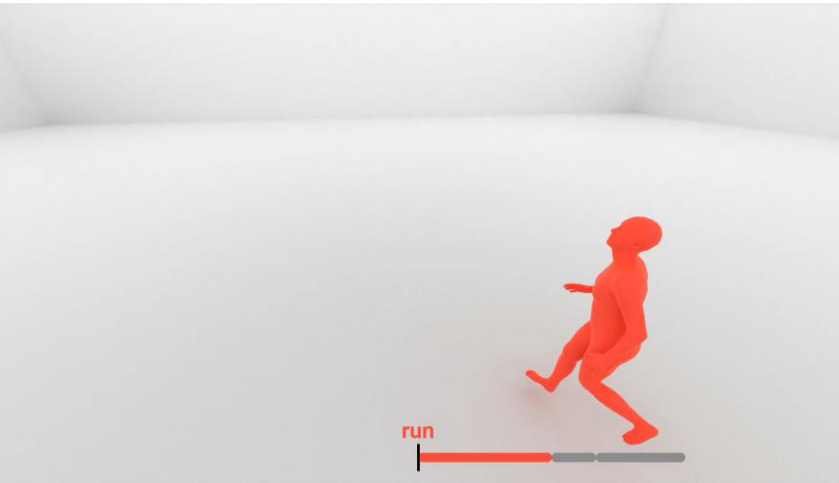
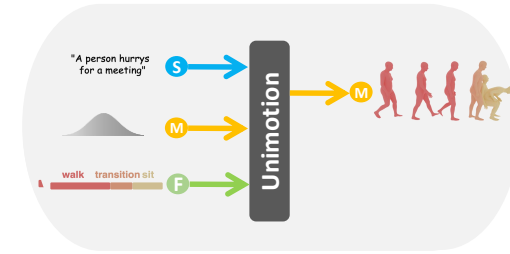
Output Text
+ 3.bring down right hand 5.walk forward 6.turn around 7.walk forward
1.turn around 2.walk forward 4.stand up

Edited Text
5.run forward
↓ **User Edit: Label Modification**

Edited Motion


Frame-level Text-to-Motion Generation

Our Method achieves flexible hierarchical control and stronger correspondence



FlowMDM

- ✗ Sequence-level control
- ✓ Frame-level control

STMC

- ✗ Sequence-level control
- ✓ Frame-level control

Ours

- ✓ Sequence-level control
- ✓ Frame-level control

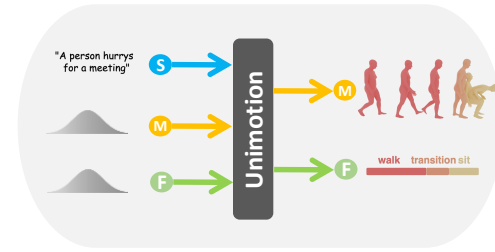
Li et al. UniMotion, 3DV'25

Barquero et al. FlowMDM, CVPR'24

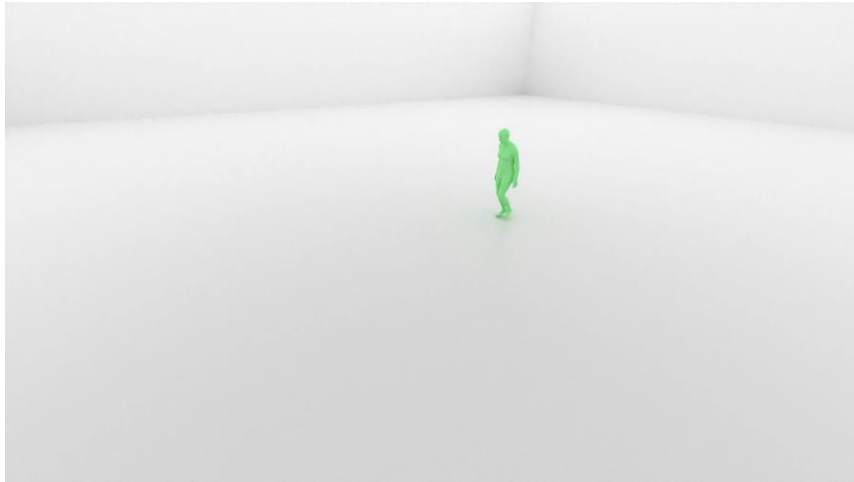
Petrovich et al. STMC, CVPRW24.

Sequence-level Text-to-Motion Generation

Our Method achieves motion synthesis and motion understanding at the same time

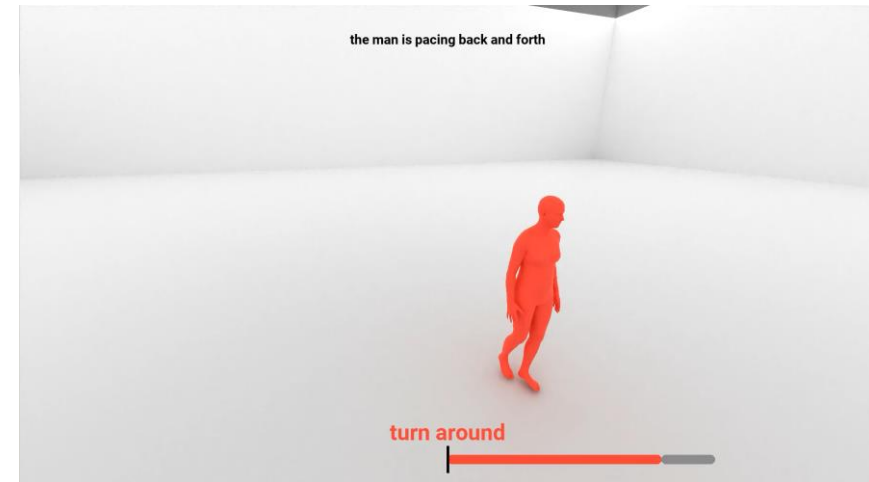


The man is pacing back and forth



MDM

- ✓ Motion generation
- ✗ Motion understanding

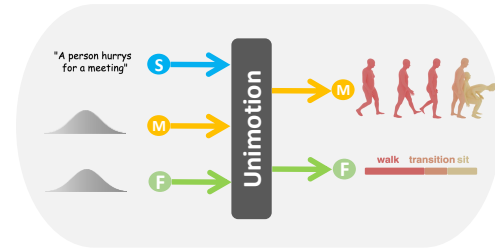


Ours

- ✓ Motion generation
- ✓ Motion understanding

Sequence-level Text-to-Motion Generation

Our Method achieves motion synthesis and motion understanding at the same time

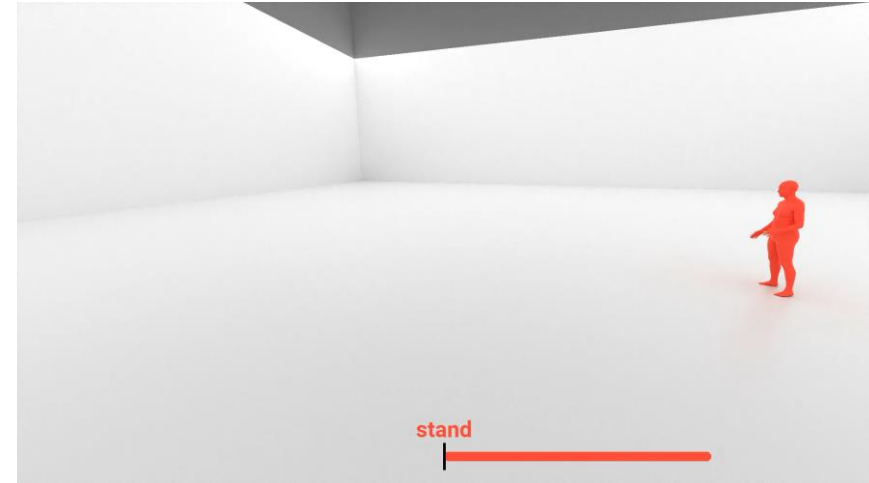


a person sprinting ahead, and then slowing down



MDM

- ✓ Motion generation
- ✗ Motion understanding



Ours

- ✓ Motion generation
- ✓ Motion understanding

Multi-tasks

Hierarchical Text to Motion

Input: Seq-level + Frame-level Text

"A person waves hands above head"

3.sitting down + 1.walking

2.transition ↓



Output Motion

Motion to Text

Input Motion



Output: Frame-level Text ↓

4.walking

2.step up

3.step down

1.walking

Unconditional Joint Generation

Input: Noise Vec →



+

2.jump forward

3.turn around

1.walking forward

Output: Motion + Frame-level Text

Motion Generation and Editing (Combined Task)

Input: Seq-level Text

"This person walks forward and picks something up then walks back"



Output Motion



Output Text +

1.turn around 2.walk forward 4.stand up

Edited Text

5.run forward

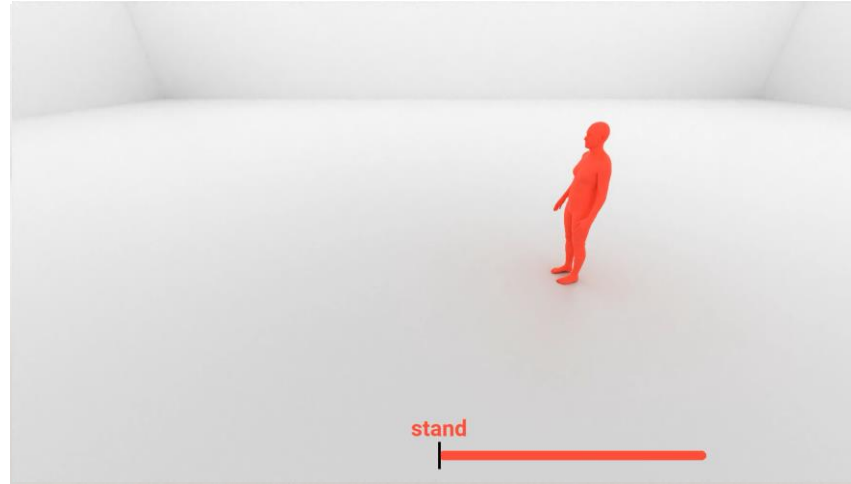
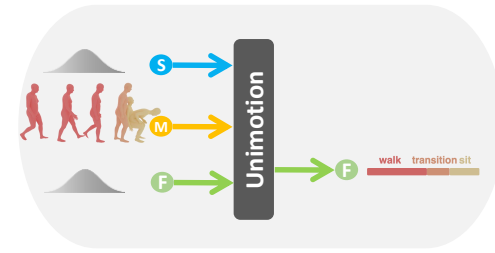
User Edit:
Label Modification

Edited Motion



Fine-grained Motion-to-text Understanding

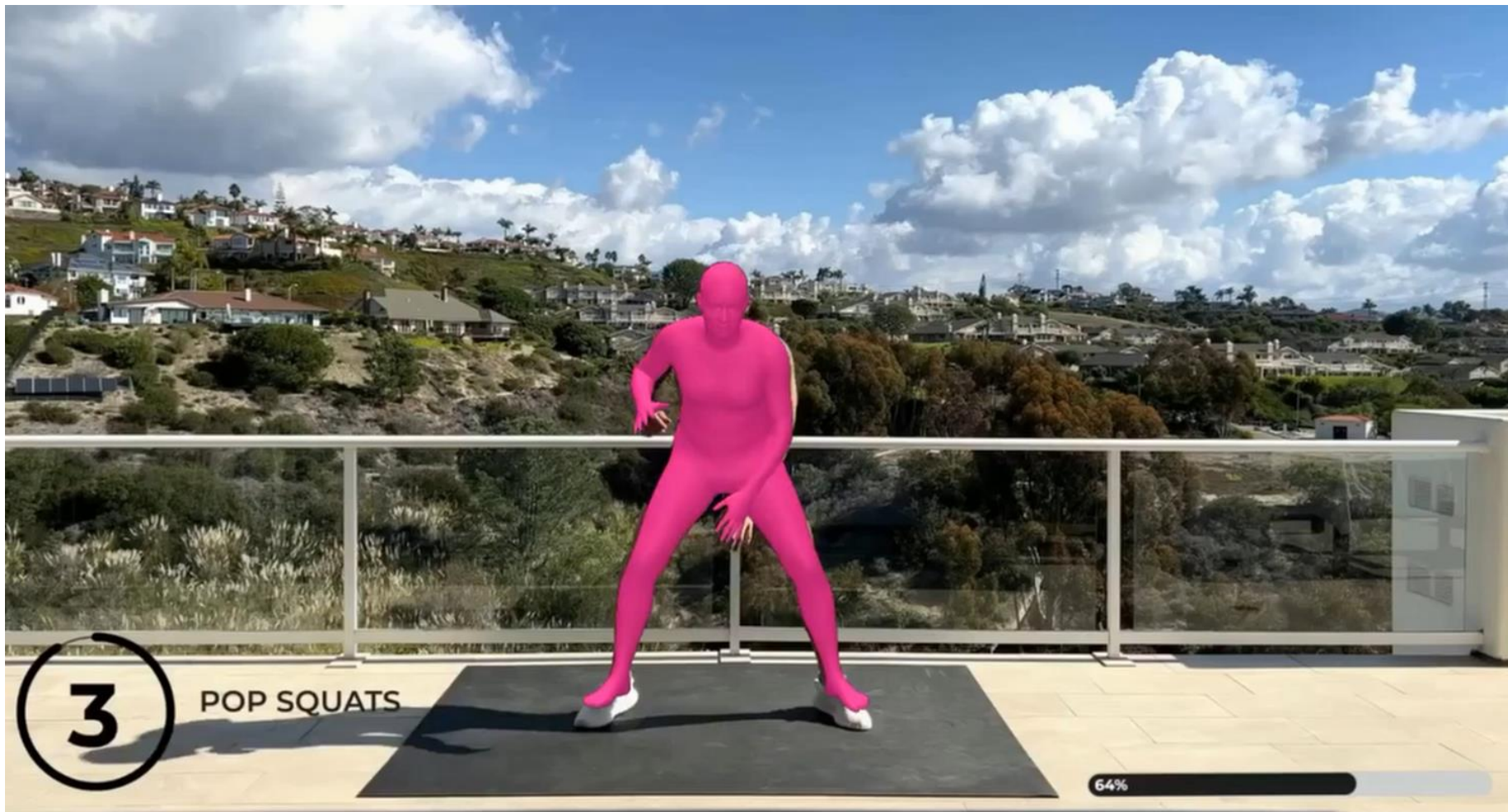
Our Method achieves motion understanding with fine-grained temporal awareness



Ours

- ✓ **Motion generation**
- ✓ **Motion understanding**

2D-video annotation



transition



- transition
- picking something up
- standing up
- standing and looking around
- standing right

Multi-tasks

Hierarchical Text to Motion

Input: Seq-level + Frame-level Text

"A person waves hands above head"

3.sitting down + 1.walking

2.transition ↓



Output Motion

Motion to Text

Input Motion



Output: Frame-level Text ↓

4.walking

2.step up

3.step down

1.walking

Unconditional Joint Generation

Input: Noise Vec →



+

2.jump forward

3.turn around

1.walking forward

Output: Motion + Frame-level Text

Motion Generation and Editing (Combined Task)

Input: Seq-level Text

"This person walks forward and picks something up then walks back"



Output Motion



Output Text +

1.turn around 2.walk forward 4.stand up

Edited Text

5.run forward

User Edit:
Label Modification

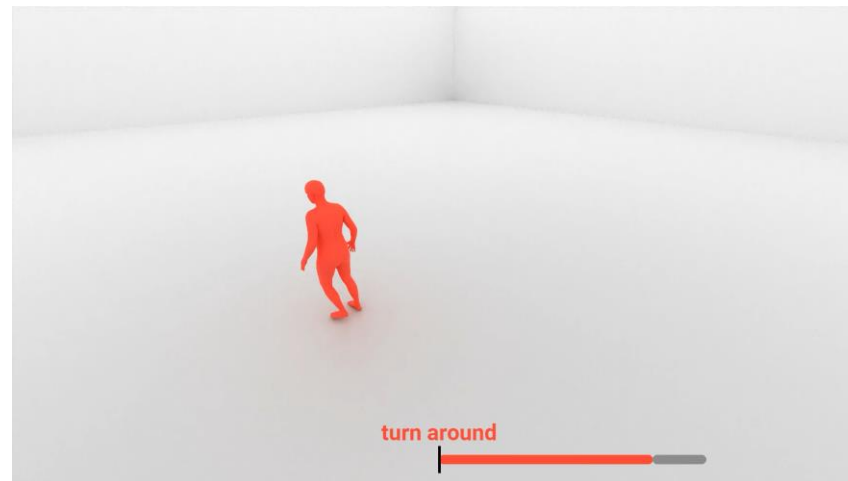
Edited Motion



Motion generation and Editing

User input: "The person walks forward and picks something up then walks back."

Model output:



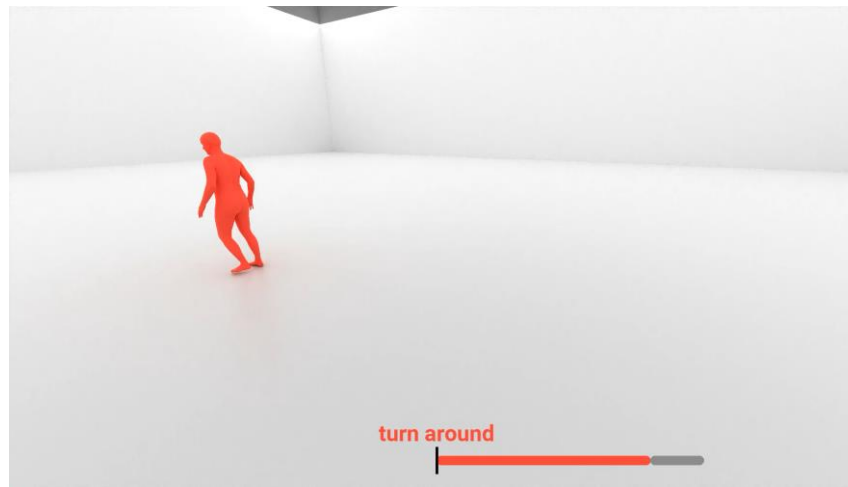
Motion generation and Editing

User input: "The person walks forward and picks something up then walks back."

Model output:



User edit:



Multi-tasks


Hierarchical Text to Motion

Input: Seq-level + Frame-level Text

"A person waves hands above head"

3.sitting down + 1.walking


2.transition ↓



Output Motion

Motion to Text

Input Motion




Output: Frame-level Text ↓

4.walking 2.step up

3.step down 1.walking

Unconditional Joint Generation

Input: Noise Vec →



+

2.jump forward

3.turn around 1.walking forward

Output: Motion + Frame-level Text

Motion Generation and Editing (Combined Task)

Input: Seq-level Text →

"This person walks forward and picks something up then walks back" →



Output Motion

Output Text +

3.bring down right hand 5.walk forward 6.turn around 7.walk forward

1.turn around 2.walk forward 4.stand up

Edited Text

5.run forward

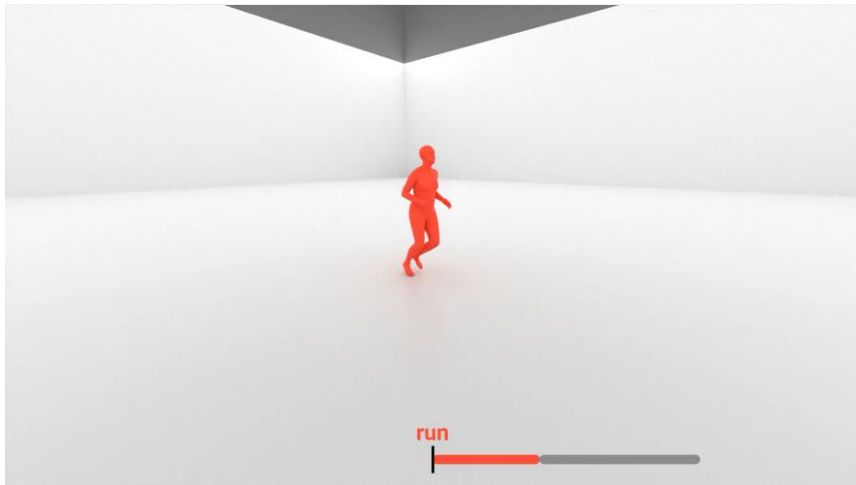
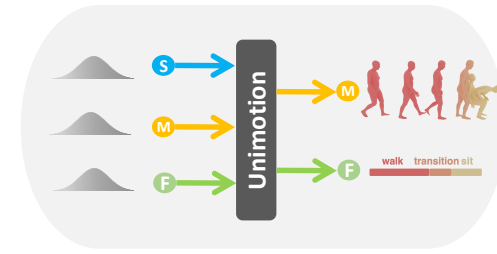
User Edit: Label Modification



Edited Motion

Unconditional Joint Generation

Our method is the first to do unconditional joint generation



Ours



Ours

- ✓ Motion generation
- ✓ Motion understanding

Take home messages

- Diffusion is also powerful for motion generation, but requires geometry constrain to train a good model.
- Pretrained motion prior is useful to generate complex compositional motions.
- A unified motion helps text and motion in both directions.

Slide credits

- Thanks to Guy Tevet for kindly providing the slides for diffusion based human motion models.

Thank you!

"A person is standing and waving goodbye!"

