# Uniform Text-Motion Generation and Editing via Diffusion Model

*Ruoyu Wang*, *Xiang Li*, *Tengjiao Sun*, *Yangfan He*, *TIANYU SHI*, *yitingxie*

*Genfun.ai* https://genlab3d.genfun.ai/

ty.shi@mail.utoronto.ca

## Motivation

**⬡ Limited to unimodal inputs and outputs ⬡**

▪ **Insufficient Guidance by Textual Instruction**

• Rely solely on textual instructions for guidance, lacking the capability to process motion or multimodal inputs.

• Textual instructions are often brief and ambiguous, making them insufficient to achieve the desired outcomes in many scenarios.

▪ **Restricted to motion generation**

• Prevent them from performing related tasks such as motion annotation or generating text-based descriptions of motions.
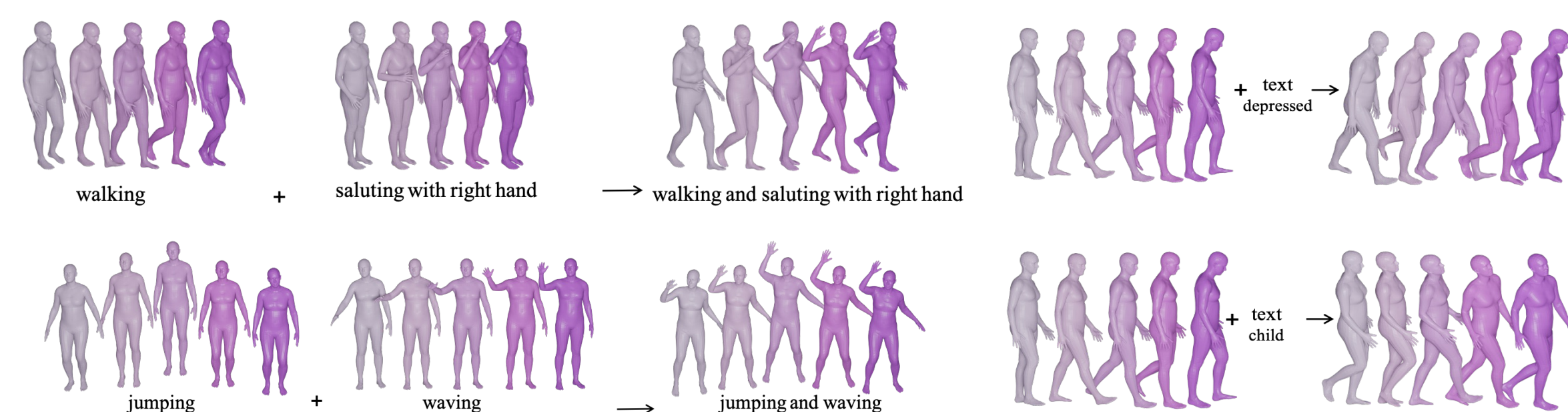
## Conclusion

▪ **Overcome the limitation of single-modal inputs and outputs**

• Demonstrate advanced effectiveness and generalization across multiple tasks, including text-driven motion generation, motion captioning, motion completion, and multimodal motion editing.
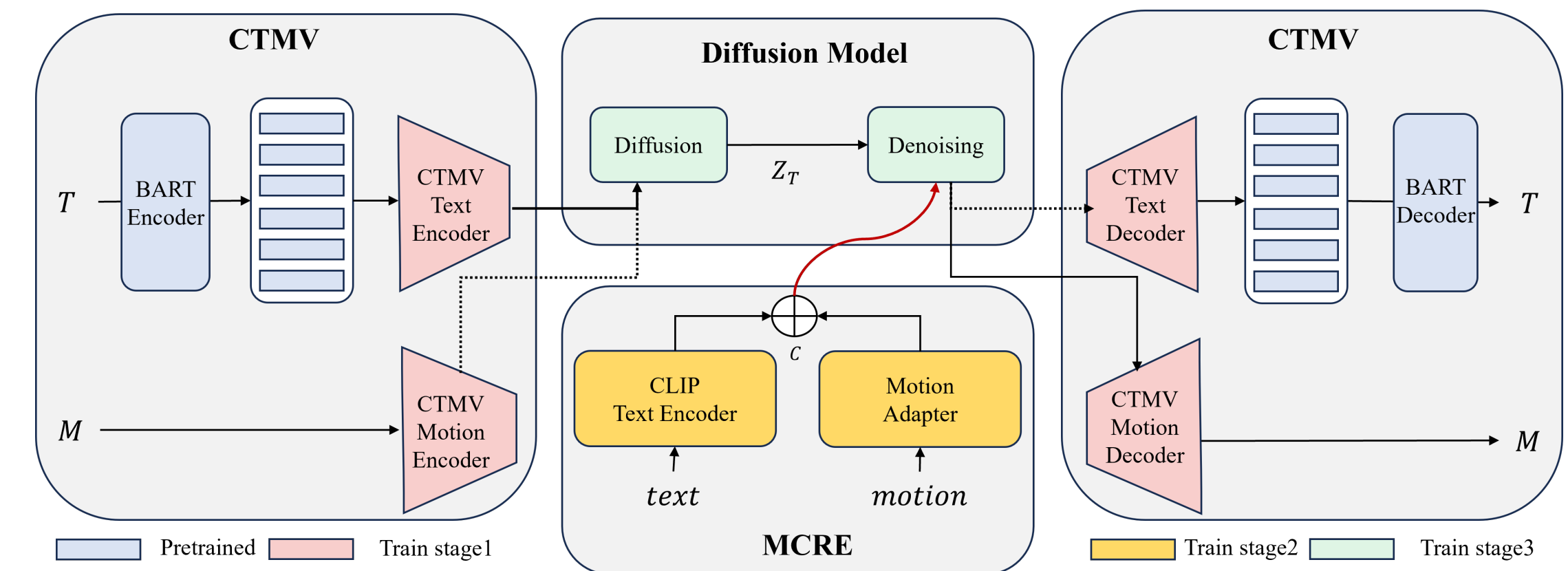
• Text-Driven Motion Generation



a man is throwing

a man is shooting basketball

• Multimodal Motion Editing



walking + saluting with right hand → walking and saluting with right hand

+ text depressed →

jumping + waving → jumping and waving

+ text child →

## Method



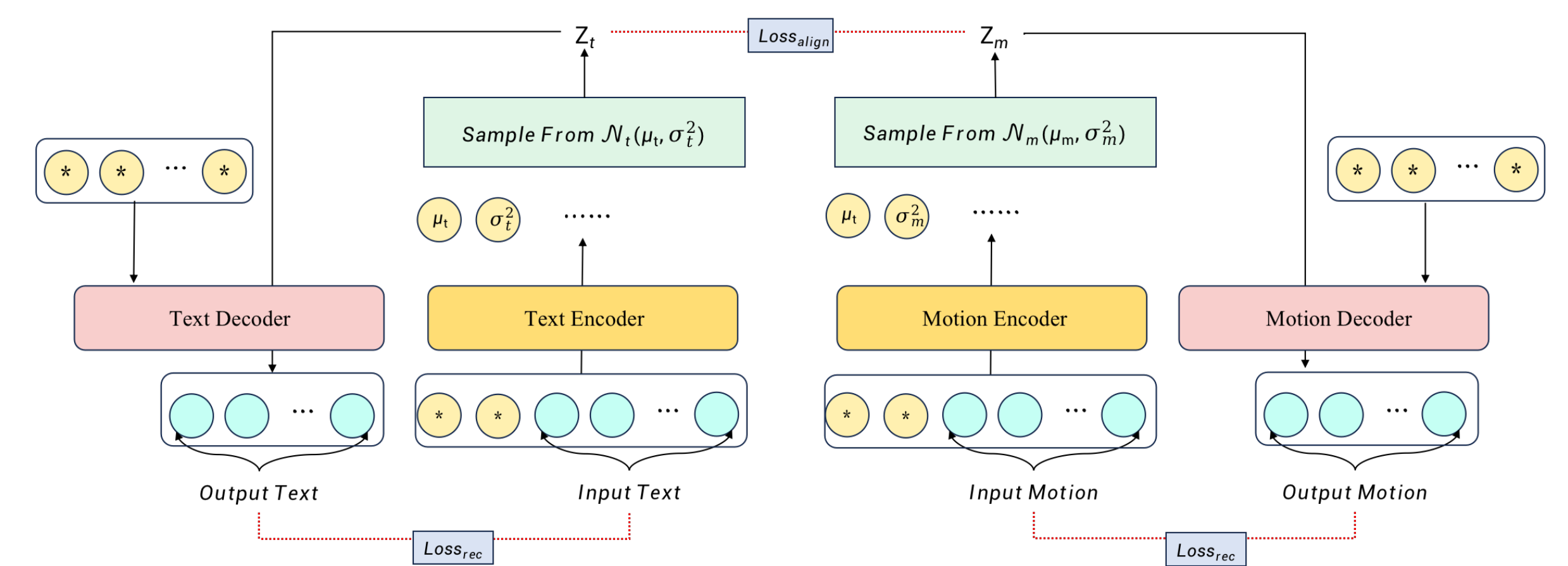▪ **Multimodal Conditional Representation and Editing (MCRE)**

• Designs a motion adapter to align motion with CLIP's text representations, leveraging its rich semantic understanding.



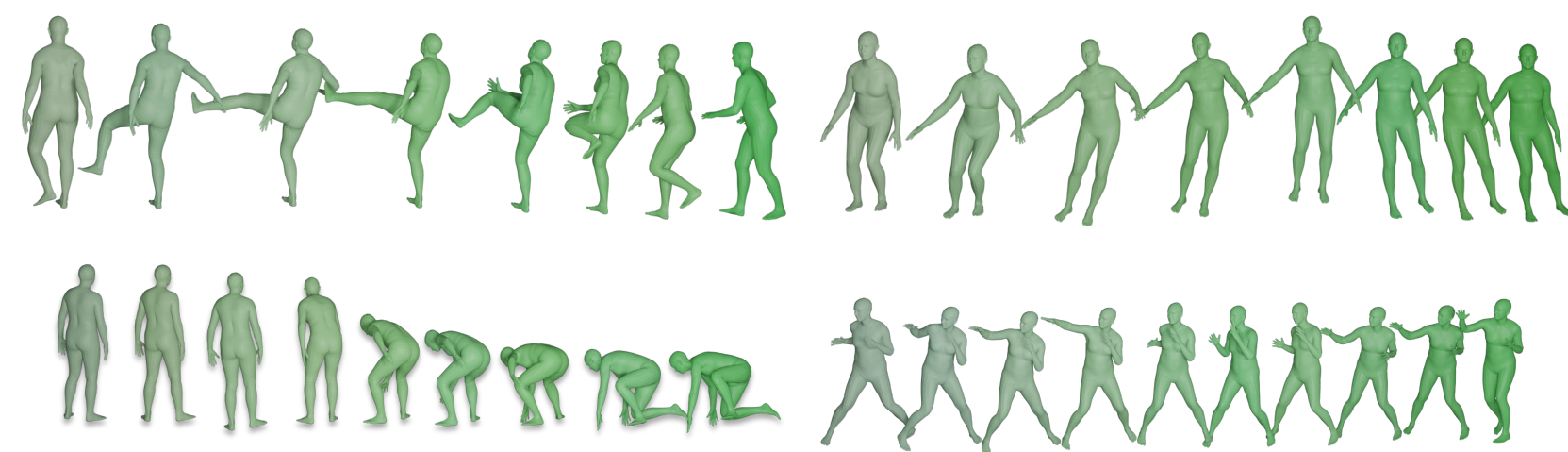▪ **Contrastive Text-Motion Variational Autoencoder (CTMV)**

• Aligns text-motion pairs into a shared latent space via transformer-based encoders leveraging contrastive learning.

• Offloads the task of generating high-frequency details of different modalities to the autoencoder

• Enables the diffusion to focus on the high-level semantics generation.

# Uniform Text-Motion Generation and Editing via Diffusion Model

▪ **Training Data--HumanML3D dataset**

• Comprising 14,616 motions and 44,970 descriptions.

• Data span various domains daily activities, exercise, and artistic performances.

• With an average duration of 7.1 seconds per action and an average description length of 12 words.



▪ **Training Strategy**

• Jointly training by training losses, which is composed of reconstruction and alignment losses.

▪ **Reconstruction Loss:**

• Cross-Entropy Loss for Text

• L2 loss for Motion

▪ **Alignment Losses**

• Cosine loss

• KL Loss.

▪ **Demonstrate advanced effectiveness and generalization across multiple tasks**

• Text-to-Motion

Table 1. Quantitative results of text-to-motion task on the HumanML3D and Motion-X test set

| Method | R-Precision↑ | | | FID↓ | MM-Dist↓ | Diversity→ | MModality↑ |
| | Top1 | Top2 | Top3 | | | | |
|---|---|---|---|---|---|---|---|
| *"HumanML3D Test set"* | | | | | | | |
| Real | 0.511±.003 | 0.703±.003 | 0.797±.002 | 0.002±.000 | 2.974±.008 | 9.503±.065 | - |
| MDM | 0.320±.005 | 0.498±.004 | 0.611±.007 | 0.544±.044 | 5.566±.027 | 9.559±.086 | 2.799±.072 |
| MotionDiffuse | 0.491±.001 | 0.681±.001 | 0.782±.001 | 0.630±.001 | 3.113±.001 | 9.410±.049 | 1.553±.042 |
| MLD | 0.481±.003 | 0.673±.003 | 0.772±.002 | 0.473±.013 | 3.196±.010 | 9.724±.082 | 2.413±.079 |
| Ours | 0.499±.003 | 0.683±.003 | 0.780±.002 | 0.339±.003 | 3.087±.003 | 9.527±.053 | 2.500±.081 |
| *"Motion-X Test set"* | | | | | | | |
| Real | 0.509±.004 | 0.702±.004 | 0.794±.003 | 0.003±.001 | 2.995±.009 | 9.508±.064 | - |
| MDM | 0.301±.006 | 0.477±.005 | 0.590±.008 | 0.645±.045 | 5.587±.028 | 9.558±.087 | 2.788±.073 |
| MotionDiffuse | 0.470±.002 | 0.660±.002 | 0.757±.002 | 0.638±.002 | 3.232±.002 | 9.411±.050 | 1.524±.043 |
| MLD | 0.475±.004 | 0.665±.004 | 0.764±.003 | 0.484±.014 | 3.255±.011 | 9.723±.083 | 2.312±.080 |
| Ours | 0.490±.004 | 0.677±.004 | 0.772±.003 | 0.355±.010 | 3.098±.003 | 9.560±.054 | 2.401±.081 |



a man is throwing

a man is shooting basketball

a man is dancing ballet

playing baseball

fancy skating

squat down

playing football

kneeling on one knee

playing the violin

• Motion-to-Text

Table 2. Quantitative results of motion-to-text task on the HumanML3D and Motion-X test set

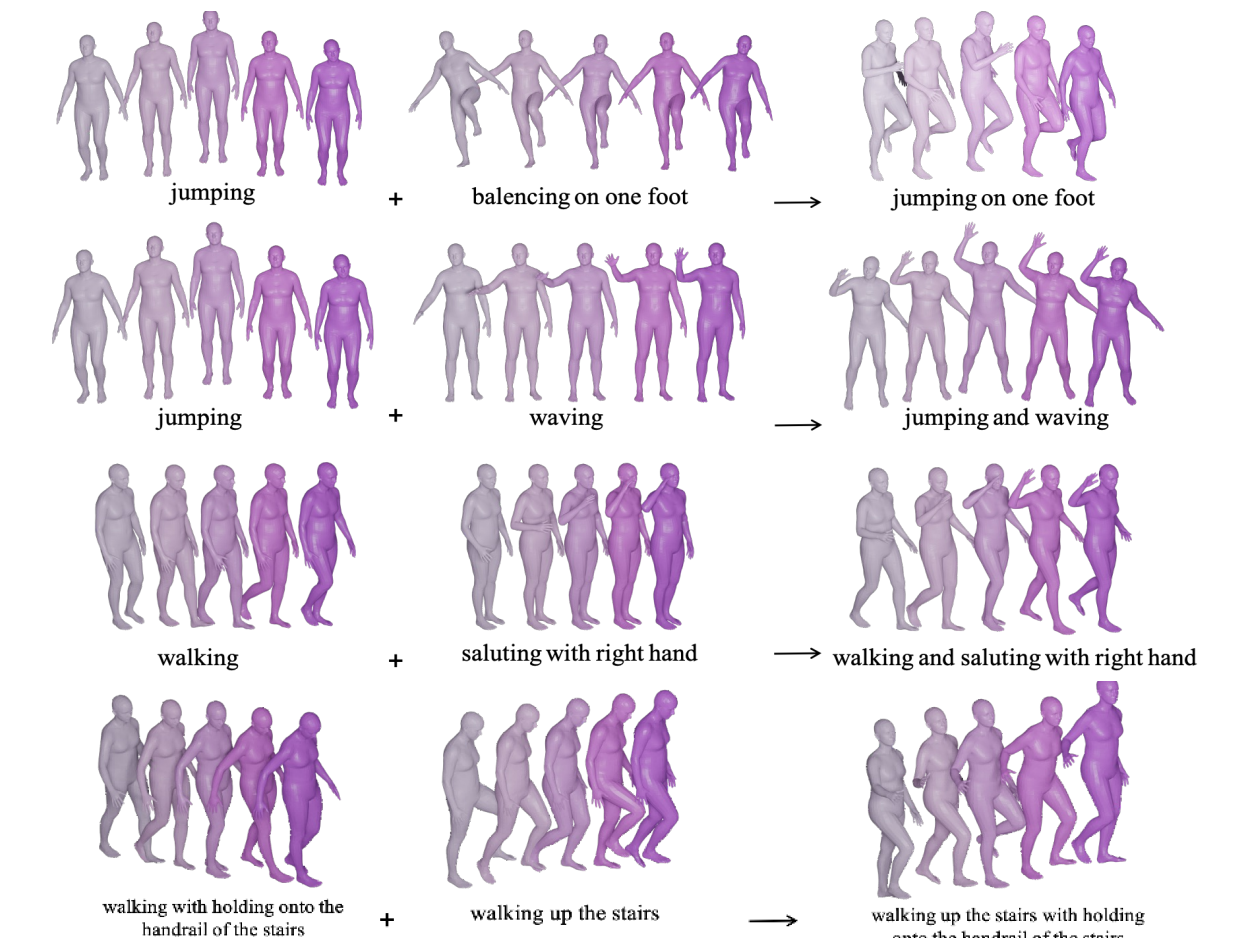| Method | R-Precision↑ | | MM-Dist↓ | $Length_{avg}$↑ | Bleu@1↑ | Bleu@4↑ | Rouge↑ | Cider↑ | BertScore↑ |
| | Top1 | Top3 | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| *"HumanML3D Test set"* | | | | | | | | | |
| Real | 0.523 | 0.828 | 2.901 | 12.750 | - | - | - | - | - |
| TM2T | 0.516 | 0.823 | 2.935 | 10.670 | 48.900 | 7.000 | 38.100 | 16.800 | 32.200 |
| MotionGPT | 0.543 | 0.827 | 2.821 | 13.040 | 48.200 | 12.470 | 37.400 | 29.200 | 32.400 |
| Ours | 0.520 | 0.825 | 2.878 | 13.000 | 49.300 | 11.700 | 35.000 | 30.300 | 32.800 |
| *"Motion-X Test set"* | | | | | | | | | |
| Real | 0.520 | 0.821 | 2.892 | 17.250 | - | - | - | - | - |
| TM2T | 0.484 | 0.803 | 2.975 | 11.901 | 46.500 | 6.781 | 35.903 | 15.201 | 29.909 |
| MotionGPT | 0.518 | 0.817 | 2.858 | 14.340 | 47.800 | 11.890 | 36.202 | 27.601 | 31.009 |
| Ours | 0.510 | 0.820 | 2.886 | 13.890 | 48.302 | 11.025 | 34.800 | 28.998 | 31.198 |

• Motion Completion

Table 4. Quantitative results of motion completion task on the HumanML3D and Motion-X test set

| Method | | Motion Prediction | | | | Motion-In-between | |
| | FID↓ | Diversity→ | ADE↓ | FDE↓ | FID↓ | Diversity↑ | ADE↓ |
|---|---|---|---|---|---|---|---|
| *"HumanML3D Test set"* | | | | | | | |
| Real | 0.002 | 9.503 | - | - | 0.002 | 9.503 | - |
| MDM | 6.031 | 7.813 | 5.446 | 8.561 | 2.698 | 8.420 | 3.787 |
| Ours | 1.702 | 9.001 | 4.740 | 6.670 | 1.203 | 9.600 | 3.669 |
| *"Motion-X Test set"* | | | | | | | |
| Real | 0.003 | 9.508 | - | - | 0.003 | 9.508 | - |
| MDM | 8.931 | 7.783 | 7.846 | 10.021 | 4.398 | 8.150 | 2.987 |
| Ours | 2.102 | 8.901 | 5.940 | 7.970 | 1.583 | 9.230 | 3.042 |

• Multimodal Motion Editing

1) motion based



jumping + balancing on one foot → jumping on one foot

jumping + waving → jumping and waving

walking + saluting with right hand → walking and saluting with right hand

walking with holding onto the handrail of the stairs + walking up the stairs → walking up the stairs with holding onto the handrail of the stairs

2) multimodal based



walking + text young woman → 

walking + text1 young man + text depressed → 

walking + text keep balance → keep balance while walking

swimming + text standing → swimming in standing position

walking + text limp → walking with a limp

swimming + text upside down → backswiming