HumanTOMATO: Text-aligned Whole-body Motion Generation

Page: https://lhchen.top/HumanTOMATO

Ling-Hao CHEN Tsinghua University, IDEA Research thu.lhchen@gmail.com

June 5, 2024

Road Map

Road Map

Human Motion

Road Map

Out-modality Driven source (text, audio, and so on)





Road Map

In-modality Driven source (spatial or temporal control)



Out-modality Driven source (text, audio, and so on)



Human Motion

Road Map

In-modality Driven source (spatial or temporal control)





□ Synthesizing text-aligned human motion (HumanTOMATO)

□ Zero-shot motion controlling (HumanMAC)

□ Tunable, real-time, and multi-modality motion controlling (MotionLCM)

□ Exploring on the scaling law for human motion (MotionLLM)

Road Map

In-modality Driven source (spatial or temporal control)



HumanMAC MotionLCM



Synthesizing text-aligned human motion The man is playing the ukulele, happily.

> Emphasizing one of my belief:

The source of intelligence is compression.

> Ilya Sutskever, previously from OpenAI:

Prediction is compression;

I find easier to think about compression when talking about unsupervised learning.

> Yi Ma *et al.* :

Closed-loop compressive architectures are ubiquitous for all intelligent beings and at all scales, from the brain to spinal circuits down to DNA.

 \succ Quick question

Q: Why is the compression so important for human motion?

A: Vinilla motion sapce is too noisy.



 $\frac{\text{vol}(\#\text{motion})}{\text{vol}(\#\text{whole})} \longrightarrow 0$

> Quick question

Q: Why is the compression so important for human motion?

A: Vinilla motion sapce is too noisy.

> Technical solution: Compression motion at first!



> Quick question

Q: Why is the compression so important for human motion?

A: Vinilla motion sapce is too noisy.

> Technical solution: Compression motion at first!

High-quality compression principal: Good reconstruction quality!



> High-quality compression principal: Good reconstruction quality!



> High-quality compression principal: Good reconstruction quality!



Text-aligned generation

Text-aligned generation



- Text-aligned generation
 - Q1: How generat codes?



- Text-aligned generation
 - Q1: How generat codes?

- Text-aligned generation
 - Q1: How generat codes?



- Text-aligned generation
 - Q1: How generat codes?
 - Q2: How to achieve text-motion alignment?







here we show some cases generated by human tomato

man swivels left to right with left hand up towards the face.



Body-level Comparison



Body-level





Vanilla VQ



A person walks clockwise-ly.



w/o TMA prior

w/ TMA prior

		Motion-X			GRAB	HumanML3D	
	All↓	Body \downarrow	Hand↓	All↓	Body \downarrow	Hand \downarrow	Body ↓
Vanilla VQ (512)	140.66	92.20	46.45	78.23	38.29	31.48	77.21
Vanilla VQ (1024)	139.33	91.77	46.40	76.01	37.34	29.89	71.34
RVQ (512×2)	<u>110.94</u>	73.97	<u>40.01</u>	<u>62.94</u>	31.12	27.28	63.05
$\mathrm{H}^{2}\mathrm{VQ}$ (512×2)	92.97	62.34	37.20	46.74	24.33	24.59	-

Table 7: Comparison of the motion reconstruction errors (MPJPE in mm) of different quantization methods on Motion-X, GRAB, and HumanML3D. Our H²VQ shows significant improvements.

embedding supervision		FID \downarrow	R-Precision ⁽³²⁾			TMA-R-Precision ⁽²⁵⁶⁾			Matching-	TMA-Matching-
			Top1 ↑	Top2 ↑	Top3 ↑	Top1 ↑	Top2 ↑	Top3 ↑	score \downarrow	score \downarrow
GT		0.002	0.500	0.708	0.814	0.407	0.578	0.673	2.888	0.768
CLIP	×	1.086	0.405	0.588	0.695	0.345	0.490	0.573	3.917	0.844
TMA	×	1.290	0.416	0.596	0.699	0.395	0.550	0.637	3.950	0.815
TMA	 ✓ 	1.174	0.416	0.603	0.703	0.399	0.555	0.638	3.894	0.809

Table 8: Ablation on a pre-trained text-motion-aligned model for motion generation on Motion-X. Both TMA embedding and text-motion alignment supervision help generate text-aligned motions.



□ Synthesizing text-aligned human motion (HumanTOMATO)

□ Zero-shot motion controlling (HumanMAC)

□ Tunable, real-time, and multi-modality motion controlling (MotionLCM)

□ Exploring on the scaling law for human motion (MotionLLM)

Road Map

In-modality Driven source (spatial or temporal control)



HumanMAC MotionLCM



MotionLCM: Real-time Controllable Motion Generation via Latent Consistency Model



MotionLCM: Real-time Controllable Motion Generation via Latent Consistency Model



MotionLCM: Real-time Controllable Motion Generation via Latent Consistency Model



Page: https://dai-wenxun.github.io/MotionLCM-page

MotionLLM: Understanding Human Behaviors from Human Motions and Videos



(a) Our MotionLLM takes motions or videos as inputs to understand human behaviors.





1. Begin in a supine position with arms at shoulder width. 2. Engage your core as you lift your legs to a 45 degree angle. 3. Bring your feet hip width apart. 4. Slide your hands down your thighs until your hands are on the floor. 5. Put down your feet to return to the supine position.

(b) An application: MotionLLM as your fitness coach based on its caption capability.

MotionLLM: Understanding Human Behaviors from Human Motions and Videos



(a) Our MotionLLM takes motions or videos as inputs to understand human behaviors.





1. Begin in a supine position with arms at shoulder width. 2. Engage your core as you lift your legs to a 45 degree angle. 3. Bring your feet hip width apart. 4. Slide your hands down your thighs until your hands are on the floor. 5. Put down your feet to return to the supine position.

(b) An application: MotionLLM as your fitness coach based on its caption capability.

Page: https://lhchen.top/MotionLLM

Thanks!