



stand and shush, angrily.

Yang-style 40 formTai Chi Competitionroutinestep34, happily.

The proposed HumanTOMATO can generate text-aligned whole-body motions with vivid and harmonious face, hand, and body motion.

## Key insight 1: Text-motion alignment



(b) Learning *image*-text aligned prior *implicitly* 

We clarify the importance of how to use the text and motion data to generate motions for the first time. We highlight the two method.

 $\checkmark$  Replace your **CLIP** text encoder with **TMA** text encoder. ✓ Introduce the text-motion alignment **supervision** to your motion generation model during training.

from OpenTMA import textencoder, motionencoder import torch

How to use? **motion = torch.randn(1, 64, 126)** # B = 1, T = , D = , need normalizationlengths = [64]t\_emb = textencoder(["a man is running"]).loc m\_emb = motionencoder(motion, lengths).loc

Shunlin Lu\*, Ling-Hao Chen\* Ailing Zeng, Jing Lin, Ruimao Zhang, Lei Zhang, Heung-Yeung Shum \*co-first author, random listing order Correspondence: R. Zhang, and H.-Y. Shum Project page: https://lhchen.top/HumanTOAMTO



Try TMA!



### (b) Facial Motion Tokenization (RVQ)





Auto-regressive motion generation. Both parts take textual description as input and predict tokens in an auto-regressive manner. The final whole-body motion is composed of both part motions decoded by the corresponding decoders.

# HumanTOMATO: Text-aligned Whole-body Motion Generation





## Results

		<b>D D</b> $(32)$			<b>TMA D D</b> (256)			16.4.1		MModality↑	Diversity↑
TA.	FID↓	R-Precision <sup>(02)</sup>			IMA-R-Precision <sup>(200)</sup>			Matching	IMA-Matching		
		Top1↑	Top2↑	Top3↑	Top1↑	Top2↑	Top3↑	Score *	Score *	initia during f	
Т	-	0.500	0.708	0.814	0.407	0.578	0.673	2.888	0.768	9. <u></u> 9	11.087
IOS	9.147	0.279	0.442	0.555	0.258	0.389	0.444	5.482	0.928	1.195	9.764
GPT	1.366	0.368	0.553	0.655	0.310	0.446	0.527	4.316	0.881	2.356	10.753
DM	3.800	0.352	0.547	0.634	0.310	0.430	0.530	4.050	0.840	2.530	11.400
LD	3.407	0.385	0.571	0.683	0.333	0.477	0.561	<u>3.901</u>	0.883	2.448	10.420
Diffuse	1.129	0.391	0.587	0.695	0.368	<u>0.493</u>	<u>0.584</u>	3.950	0.829	1.654	10.580
OMATO	1.174	0.416	0.603	0.703	0.399	0.555	0.638	3.894	0.809	1.732	10.812
Main results of motion generation on the Motion-X dataset.											

tizaction		Motion-X			GRAB	HumanML3D	
orko	All↓	Body $\downarrow$	Hand↓	All↓	Body $\downarrow$	Hand $\downarrow$	Body $\downarrow$
la VQ (512)	140.66	92.20	46.45	78.23	38.29	31.48	77.21
a VQ (1024)	139.33	91.77	46.40	76.01	37.34	29.89	71.34
$Q(512 \times 2)$	<u>110.94</u>	73.97	<u>40.01</u>	<u>62.94</u>	<u>31.12</u>	<u>27.28</u>	63.05
$^{\prime}\mathrm{Q}(512\times2)$	92.97	62.34	37.20	46.74	24.33	24.59	-

ng supervision		R-P	recision	(32)	TMA-R-Precision <sup>(256)</sup>			Matching-	TMA-Matching-
ing supervision		Top1 ↑	Top2 ↑	Top3 ↑	Top1 ↑	Top $2\uparrow$	Top3 ↑	score $\downarrow$	score $\downarrow$
GT	0.002	0.500	0.708	0.814	0.407	0.578	0.673	2.888	0.768
×	1.086	0.405	0.588	0.695	0.345	0.490	0.573	3.917	0.844
×	1.290	0.416	0.596	0.699	0.395	0.550	0.637	3.950	0.815
~	1.174	0.416	0.603	0.703	0.399	0.555	0.638	3.894	0.809
A blation on a protection d taxt motion aligned model for motion concretion									

Test-in-the-wild result (unseen text)

### Ling-Hao Chen