

Action-Conditioned 3D Human Motion Synthesis with Transformer VAE

Mathis Petrovich¹ Michael J. Black² Gül Varol¹

¹LIGM, École des Ponts, Univ Gustave Eiffel, CNRS, France

²Max Planck Institute for Intelligent Systems, Tübingen, Germany



MAX PLANCK INSTITUTE FOR INTELLIGENT SYSTEMS

2021 ICCV OCTOBER 11-17 VIRTUAL

<https://imagine.enpc.fr/~petrovim/actor/>

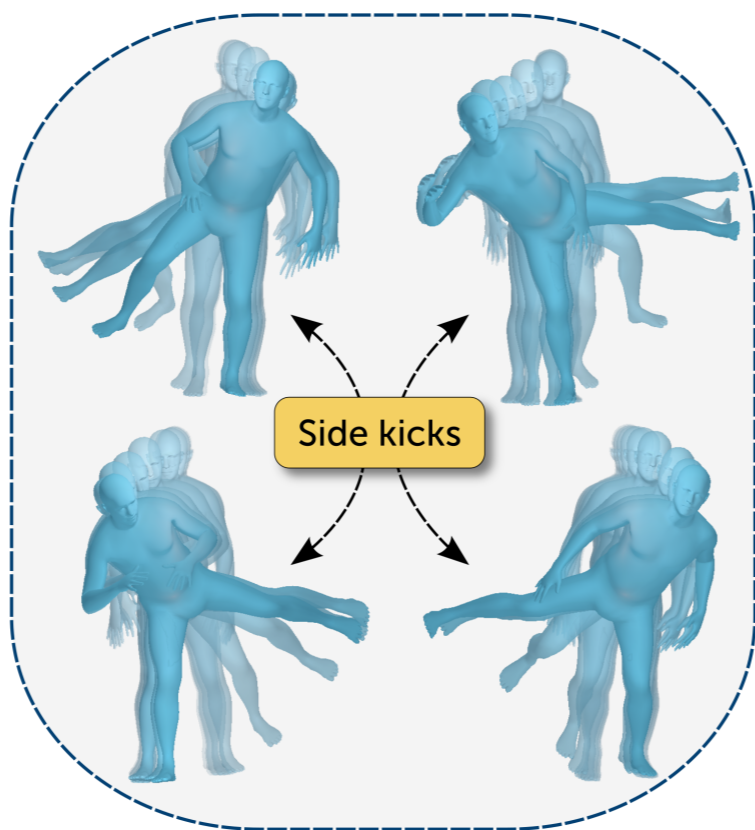
Introduction

Goal & contributions

- Generating synthetic but realistic and diverse human motion sequence given an **action label**
- Learning from noisy 3D body poses estimated from monocular action recognition datasets^{3,4,5}

Motivations

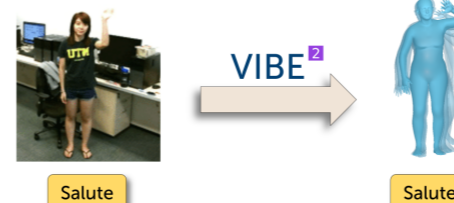
- Augmenting existing Mocap datasets, which are expensive and limited in size
- Serving as additional training data for motion recognition
- A compact action-aware latent space for human motions



Training data

NTU13⁴ (13 actions)

- RGB-D dataset, subset of NTU-120
- SMPL poses estimated with VIBE



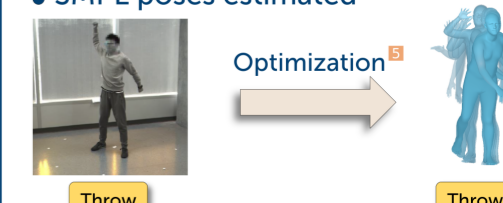
UESTC³ (40 actions)

- RGB-D dataset
- SMPL poses estimated with VIBE



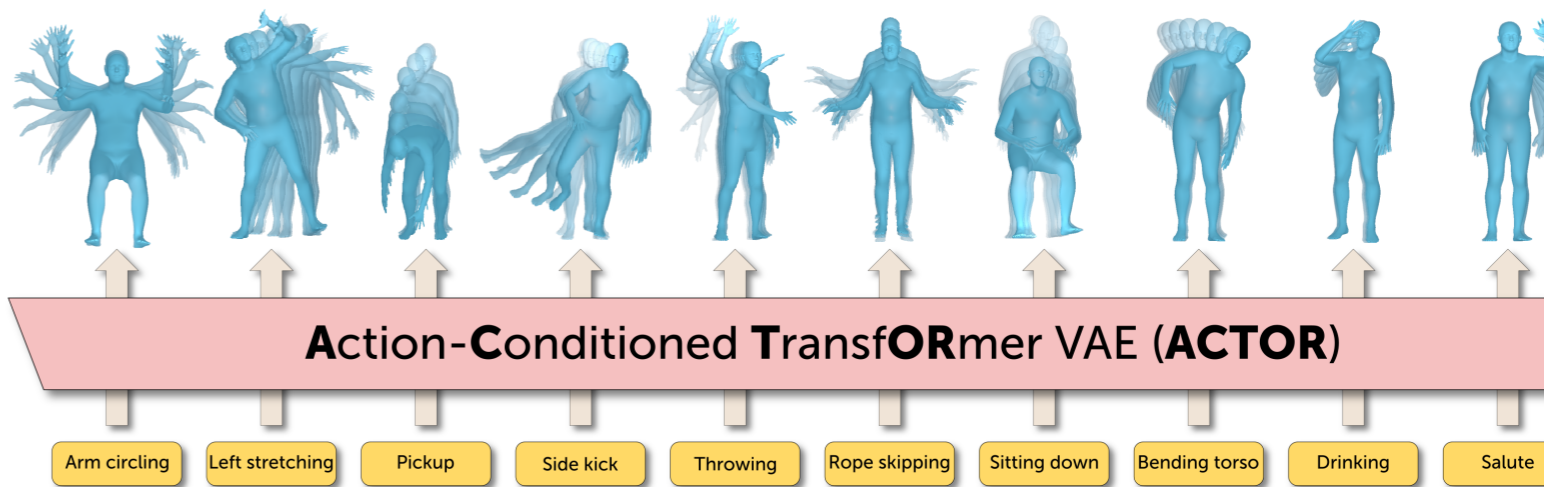
HumanAct12⁵ (12 actions)

- RGB-D + polarization images, subset of PHSPDataset
- SMPL poses estimated



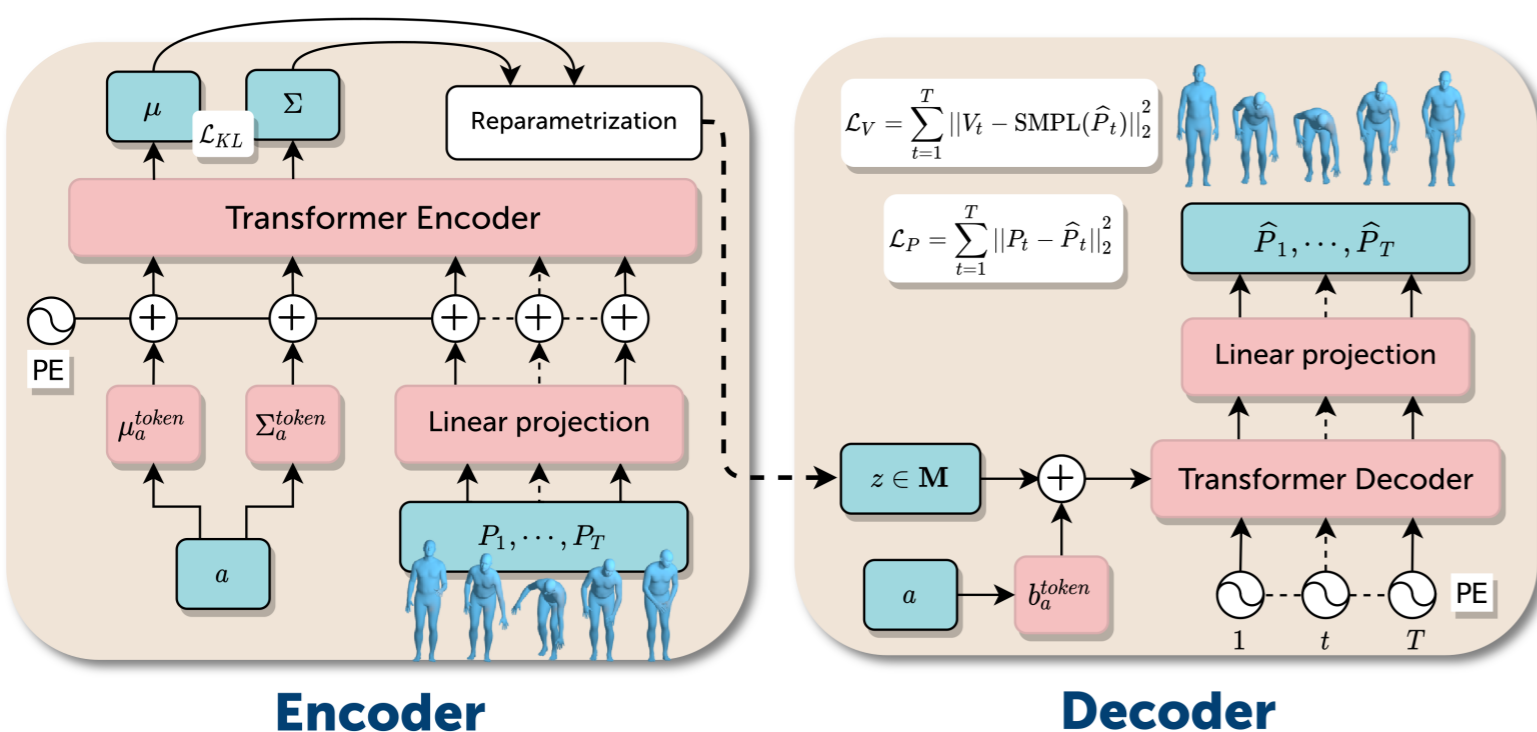
Qualitative results

The generated sequences are realistic, diverse and smooth.



ACTOR: Action-Conditioned TransfORmer VAE

- Non-autogressive Transformer architecture
- Sequence-level Variational autoencoder (VAE)
- Learnable tokens μ_a^{token} and Σ_a^{token}
- Loss terms on rotations and vertices (SMPL¹)
- Allows to generate variable length sequences with various body shapes



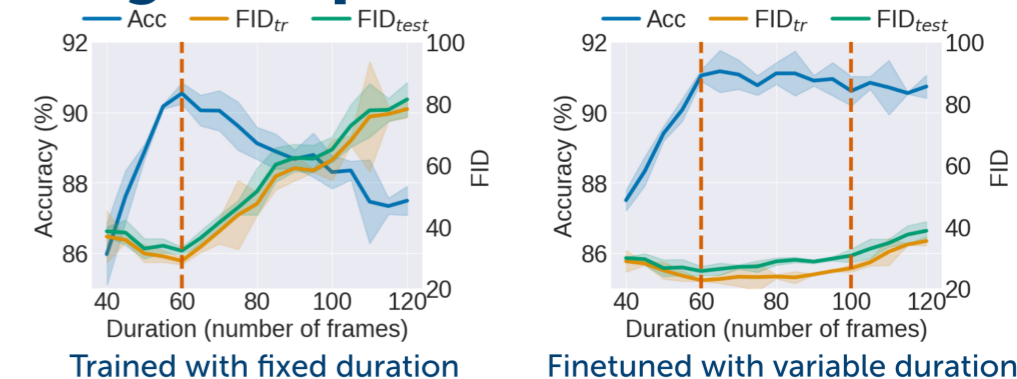
Generating variable-length sequences

We evaluate the capability of the models trained on UESTC³ with

- fixed size (60 frames)
- variable-size (between 60/100 frames)

on generating various durations.

The performance is overall improved when the model has previously seen duration variations in training.



Comparison with previous works

Metrics

Method	NTU-13 ⁴				HumanAct12 ⁵			
	FID _{tr} ↓	Acc. ↑	Div. →	Multimod. →	FID _{tr} ↓	Acc. ↑	Div. →	Multimod. →
FID: Fréchet Inception Distance (Similarity between GT distribution and the generation distribution)								
Real [Action2Motion]	0.03	99.9	7.11	2.19	0.09	99.7	6.85	2.45
Real*	0.02	99.8	7.07	2.25	0.02	99.4	6.86	2.60
Acc: Action recognition accuracy								
CondGRU	28.31	7.80	3.66	3.58	40.61	8.0	2.38	2.34
Two-stage GAN	13.86	20.2	5.33	3.49	10.48	42.1	5.96	2.81
Act-MoCoGAN	2.72	99.7	6.92	0.91	5.61	79.3	6.75	1.06
Action2Motion ¹	0.33	94.9	7.07	2.05	2.46	92.3	7.03	2.87
Multimod: Multimodality (Per-action diversity)								
ACTOR (ours)	0.11	97.1	7.08	2.08	0.12	95.5	6.84	2.53

References

- ¹Loper et al. SMPL: A skinned multi-person linear model 2015
- ²Liu et al. NTU RGB+D 120: A large-scale benchmark for 3D human activity understanding 2019
- ³Kocabas et al VIBE: Video inference for human body pose and shape estimation 2020
- ⁴Ji et al. A large-scale RGB-D database for arbitrary-view human action recognition 2018
- ⁵Liu et al. NTU RGB+D 120: A large-scale benchmark for 3D human activity understanding 2019
- ⁶Guo et al. Action2Motion: Conditioned generation of 3D human motions 2020