

Max Planck Institute for Intelligent Systems
Perceiving Systems

# DEEP REPRESENTATION LEARNING FOR HUMAN MOTION PREDICTION AND CLASSIFICATION

# MOTIVATION

Generative models of 3D human motion are often restricted to a small number of activities and can therefore not generalize well to novel movements or applications. In this work we propose a deep learning framework for human motion capture data that learns a generic representation from a large corpus of motion capture data and generalizes well to new, unseen, motions. Using an encoding-decoding network that learns to predict future 3D poses from the most recent past, we extract a feature representation of human motion. Most work on deep learning for sequence prediction focuses on video and speech. Since skeletal data has a different structure, we present and evaluate different network architectures that make different assumptions about time dependencies and limb correlations. To quantify the learned features, we use the output of different layers for action classification and visualize the receptive fields of the network units. Our method outperforms the recent state of the art in skeletal motion prediction even though these use action specific training data.

# MODELS



temporal encoder (S-TE)

temporal encoder (C-TE) c) hierarchical temporal encoder (H-TE)

We train the model on a large portion of the CMU mocap dataset, producing a generic representation.

### ACKNOWLEDGEMENT

This work was partly supported by the EU through the project socSMCs (H2020-FETPROACT-2014) and Swedish Foundation for Strategic Research.

#### REFERENCES

- [1] Fragkiadaki, Katerina and Levine, Sergey and Felsen, Panna and Malik, Jitendras. Recurrent network models for human dynamics. IEEE International Conference on Computer Vision (2015): 346-4354.
- Jain, Ashesh and Zamir, Amir R and Savarese, Silvio and Saxena, Ashutosh. Structural-RNN: Deep Learning on Spatio-Temporal Graphs. IEEE Conference on Computer Vision and Pattern Recognition (2016): 5308-5317.
- [3] Liu, Hailong and Taniguchi, Tadahiro. Feature extraction and pattern recognition for human motion by a deep sparse autoencoder. IEEE International Conference on Computer and Information Technology (2014): 173–181.

Judith Bütepage, Michael J. Black, Danica Kragic and Hedvig Kjellström butepage@kth.se, black@tuebingen.mpg.de, dani@kth.se, hedvig@kth.se



# CONCLUSION

In this work we develop an unsupervised representation learning scheme for long-term prediction of everyday human motion that is not confined to a small set of actions. We demonstrate that our learned low-dimensional representation can be used for action classification and that we outperform more complex deep learning models in terms of motion prediction. Our approach can be viewed as a generative model that has low computational complexity once trained, which makes it suitable for online tasks.

Projecting the activity of the units onto three dimensions reveals a low-dimensional representation of human motion dynamics.

#### CLASSIFICATION

Method	Classification rate								
Data $(1.6 \text{ s})$		0.76							
PCA		0.73							
	Lower Layer	Middle Layer	Upper Layer						
DSAE $[3]$	0.72	0.65	0.62						
S-TE	0.78	0.74	0.67						
C-TE	0.78	0.74	0.73						
H-TE	0.77	0.73	0.69						

We classify the actions walk, run, punching, boxing, jump, shake hands, laugh and drink based on the output of units of different layers. Deep sparse autoencoders (DSAE) reconstruct the input data and do not predict.

				— f	uture					
									A	
					A	A	T	T		
0ms ur ac	560	ms	Long T 1000m	$\frac{1}{1}$ erm	1600ms					
.53	0.	$\frac{1}{6}$	$\frac{0.67}{0.67}$		0.7					
.45	0.5	52	0.59		0.63			Com	nnarod	to re
.42	0.4	19	0.57		0.62			COII		
$\frac{.6M}{27}$			0.42		0.45	_		proa	ich is ai	ore to
.37 27		20 24	$\begin{array}{r} 0.43 \\ \hline 0.42 \end{array}$		0.43	_		Evei	n for lor	ıg-ter
.25	0.0	$\frac{3}{3}$	0.12	)	0.39	_		rema	ains sm	aller.
MU		I		I		-				
.33	0.3	35	0.37		0.37					
.25	0.3	30	0.33		0.35	_				
.24	0.2	28	0.31		0.33					

#### IEEE 2017 Conference on **Computer Vision and Pattern** Recognition





recurrent, action specific models, our apto generalize to unseen action types. term predictions, the motion prediction error