

# Contextually Plausible and Diverse 3D Human Motion Prediction (Supplementary Material)

Sadegh Aliakbarian\*  
Microsoft

Fatemeh Saleh  
ACRV, ANU

Lars Petersson  
Data61, CSIRO

Stephen Gould  
ACRV, ANU

Mathieu Salzmann  
CVLab, EPFL

## 1. LCP-VAE Architecture

Our motion prediction model follows the architecture depicted in Fig. 3 of the main paper. Below, we describe the architecture of each component in our model. Note that human poses, consisting of 32 joints in case of the Human3.6M dataset, are represented in 4D quaternion space. Thus, each pose at each time-step is represented with a vector of size  $1 \times 128$ . All the tensor sizes described below ignore the mini-batch dimension for simplicity.

The **observed motion encoder**, or CS-VAE motion encoder, is a single layer GRU [1] network with 1024 hidden units. If the observation sequence has length  $T_{obs}$ , the observed motion encoder maps  $T_{obs} \times 128$  into a single hidden representation of size  $1 \times 1024$ , i.e., the hidden state of the last time-step. This hidden state,  $h_t$ , acts as the condition to the LCP-VAE encoder and the direct input to the CS-VAE encoder.

**CS-VAE**, similarly to any variational autoencoder, has an encoder and a decoder. The CS-VAE encoder is a fully-connected network with ReLU non-linearities, mapping the hidden state of the motion encoder, i.e.,  $h_t$ , to an embedding of size  $1 \times 512$ . Then, to generate the mean and standard deviation vectors, we use two fully connected branches. They map the embedding of size  $1 \times 512$  to a mean vector of size  $1 \times 128$  and a standard deviation vector of size  $1 \times 128$ , where 128 is the length of the latent variable. Note that we apply a ReLU non-linearity to the vector of standard deviations to ensure that it is non-negative. We then use the reparameterization trick [7] to sample a latent variable of size  $1 \times 128$ . The CS-VAE decoder consists of multiple fully-connected layers, mapping the latent variable to a variable of size  $1 \times 1024$ , acting as the initial hidden state of the observed motion decoder. Note that we apply a Tanh non-linearity to the generated hidden state to mimic the properties of a GRU hidden state.

The **observed motion decoder**, or CS-VAE motion decoder, is similar to its motion encoder, except for the fact that it reconstructs the motion autoregressively. Additionally, it is initialized with the reconstructed hidden state, i.e.,

the output of the CS-VAE decoder. The output of each GRU cell at each time-step is then fed to a fully-connected layer, mapping the GRU output to a vector of size  $1 \times 128$ , which represents a human pose with 32 joints in 4D quaternion space. To decode the motions, we use a teacher forcing technique [10] during training. At each time-step, the network chooses with probability  $P_{tf}$  whether to use its own output at the previous time-step or the ground-truth pose as input. We initialize  $P_{tf} = 1$ , and decrease it linearly at each training epoch such that, after a certain number of epochs, the model becomes completely autoregressive, i.e., uses only its own output as input to the next time-step. Note that, at test time, the motions are generated completely autoregressively, i.e., with  $P_{tf} = 0$ .

Note that the future motion encoder and decoder have exactly the same architectures as the observed motion ones. The only difference is their input, where the future motion is represented by poses from  $T_{obs}$  to  $T_{end}$  in a sequence. In the following, we describe the architecture of LCP-VAE for motion prediction.

**LCP-VAE** is a conditional variational encoder. Its encoder’s input is a representation of future motion, i.e., the last hidden state of the future motion encoder,  $h_T$ , conditioned on  $h_t$ . The conditioning is done by concatenation, thus the input to the encoder is a representation of size  $1 \times 2048$ . The LCP-VAE encoder, similarly to the CS-VAE encoder, maps its input representation to an embedding of size  $1 \times 512$ . Then, to generate the mean and standard deviation vectors, we use two fully connected branches, mapping the embedding of size  $1 \times 512$  to a mean vector of size  $1 \times 128$  and a standard deviation vector of size  $1 \times 128$ , where 128 is the length of the latent variable. Note that we apply a ReLU non-linearity to the vector of standard deviations to ensure that it is non-negative. To sample the latent variable, we use our extended reparameterization trick, explained in the main paper. This unifies the conditioning and sampling of the latent variable. Then, similarly to CS-VAE, the latent variable is fed to the LCP-VAE decoder, which is a fully connected network that maps the latent representation of size  $1 \times 128$  to a reconstructed hidden state of size  $1 \times 1024$  for future motion prediction. Note that we apply

\*Work done while at the Australian National University.



Since  $\Sigma_c^{-1}\Sigma_c = I$ ,  $|\Sigma_c|$  will be cancelled out in the log term, which yields

$$\mathcal{L}_{prior}^{LCP-VAE} = -\frac{1}{2} \left[ \log \frac{1}{|\Sigma|} - d + \text{tr}\{\Sigma\} + (\mu_c - (\mu + \Sigma\mu_c))^T \Sigma_c^{-1} (\mu_c - (\mu + \Sigma\mu_c)) \right]. \quad (5)$$

## 4. Results on the Penn Action Dataset

As a complementary experiment, we evaluate our approach on the Penn Action dataset, which contains 2326 sequences of 15 different actions, where for each person, 13 joints are annotated in 2D space. Most sequences have less than 50 frames and the task is to generate the next 35 frames given the first 15. Results are provided in Table 1, where we compare our approach with the deterministic autoregressive (AR) counterpart. Note that the upper bound for the Context metric is 0.74, i.e., the classification performance given the Penn Action ground-truth motions.

Table 1. Quantitative evaluation on the Penn Action dataset. Note that a diversity of 1.21 is reasonably high for normalized 2D joint positions, i.e., values between 0 and 1, normalized with the width and the height of the image.

Method	Test MSE (KL) (Reconstructed)	Diversity (Sampled)	Quality (Sampled)	Context (Sampled)
LCP-VAE	0.034 (6.07)	1.21	0.46	0.70
AR Counterpart	0.048 (N/A)	0.00	0.46	0.51

## 5. Evaluating Sampling Quality

In Table 2, we compare our approach with the state-of-the-art deterministic motion prediction models [9, 6, 5, 2, 4] using the MAE metric in Euler space. To have a fair comparison, we generate one motion per observation by setting the latent variable to the distribution mode, i.e.,  $z = \mu_c$ . This allows us to generate a plausible motion without having access to the ground truth. To compare against the deterministic baselines, we follow the standard setting, and thus use 50 frames (i.e., 2sec) as observation to generate the next 25 frames (i.e., 1sec). Surprisingly, despite having a very simple motion decoder architecture (one-layer GRU network) with a very simple reconstruction loss function (MSE), this motion-from-mode strategy yields results that are competitive with those of the baselines that use sophisticated architectures and advanced loss functions. We argue that learning a good, context-preserving latent representation of human motion is the contributing factor to the success of our approach. This, however, could be used in conjunction with sophisticated motion decoders and reconstruction losses, which we leave for future research.

Table 2. Comparison with the state-of-the-art deterministic models for 4 actions of Human3.6M. Note that in our approach, we use  $z = \mu_c$  to generate a single motion.

Method	Walking						Eating					
	80	160	320	400	560	1000	80	160	320	400	560	1000
Zero Velocity	0.39	0.86	0.99	1.15	1.35	1.32	0.27	0.48	0.73	0.86	1.04	1.38
LSTM-3LR [2]	1.18	1.50	1.67	1.76	1.81	2.20	1.36	1.79	2.29	2.42	2.49	2.82
SRNN [6]	1.08	1.34	1.60	1.80	1.90	2.13	1.35	1.71	2.12	2.21	2.28	2.58
DAE-LSTM [3]	1.00	1.11	1.39	1.48	1.55	1.39	1.31	1.49	1.86	1.89	1.76	2.01
GRU [9]	0.28	0.49	0.72	0.81	0.93	1.03	0.23	0.39	0.62	0.76	0.95	1.08
AGED [4]	0.22	0.36	0.55	0.67	0.78	0.91	0.17	0.28	0.51	0.64	0.86	0.93
DCT-GCN [8]	<b>0.18</b>	<b>0.31</b>	0.49	0.56	0.65	<b>0.67</b>	<b>0.16</b>	0.29	0.50	0.62	0.76	1.12
LCP-VAE	0.20	0.34	<b>0.48</b>	<b>0.53</b>	<b>0.57</b>	0.71	0.20	<b>0.26</b>	<b>0.44</b>	<b>0.52</b>	<b>0.61</b>	<b>0.92</b>
Method	Smoking						Discussion					
	80	160	320	400	560	1000	80	160	320	400	560	1000
Zero Velocity	0.26	0.48	0.97	0.95	1.02	1.69	0.31	0.67	0.94	1.04	1.41	1.96
LSTM-3LR [2]	2.05	2.34	3.10	3.18	3.24	3.42	2.25	2.33	2.45	2.46	2.48	2.93
SRNN [6]	1.90	2.30	2.90	3.10	3.21	3.23	1.67	2.03	2.20	2.31	2.39	2.43
DAE-LSTM [3]	0.92	1.03	1.15	1.25	1.38	1.77	1.11	1.20	1.38	1.42	1.53	1.73
GRU [9]	0.33	0.61	1.05	1.15	1.25	1.50	0.31	0.68	1.01	1.09	1.43	1.69
AGED [4]	0.27	0.43	0.82	0.84	1.06	1.21	0.27	0.56	<b>0.76</b>	0.83	1.25	1.30
DCT-GCN [8]	0.22	<b>0.41</b>	0.86	0.80	0.87	1.57	<b>0.20</b>	<b>0.51</b>	0.77	0.85	1.33	1.70
LCP-VAE	<b>0.21</b>	0.43	<b>0.79</b>	<b>0.79</b>	<b>0.77</b>	<b>1.15</b>	0.22	0.55	0.79	<b>0.81</b>	<b>1.05</b>	<b>1.28</b>

## 6. Additional Qualitative Results

Here, we provide qualitative results on diverse human motion prediction on the Human3.6M dataset. As can be seen in Figures 1 to 6, the motions generated by our approach are diverse and natural, and mostly within the context of the observed motion. We also provide qualitative results as a video in a separate file.



Figure 1. Qualitative evaluation of the diversity in human motion. The first row illustrates the ground-truth motion. The first six poses of each row depict the observation (the condition) and the rest are sampled from our model. Each row is a randomly sampled motion (not cherry picked). As can be seen, all sampled motions are natural, with a smooth transition from the observed to the generated ones. The diversity increases as we increase the sequence length.

## References

- [1] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*, 2014. 1



Figure 2. Additional qualitative evaluation of the diversity in human motion.

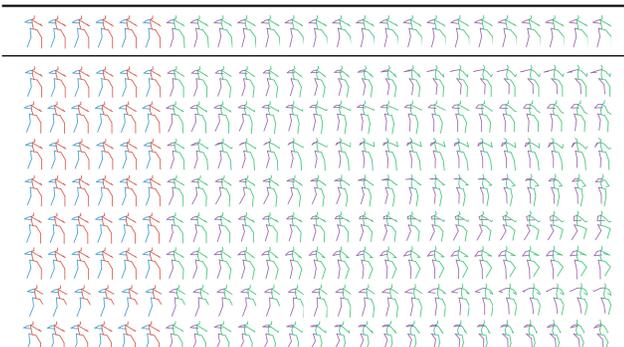


Figure 3. Additional qualitative evaluation of the diversity in human motion.

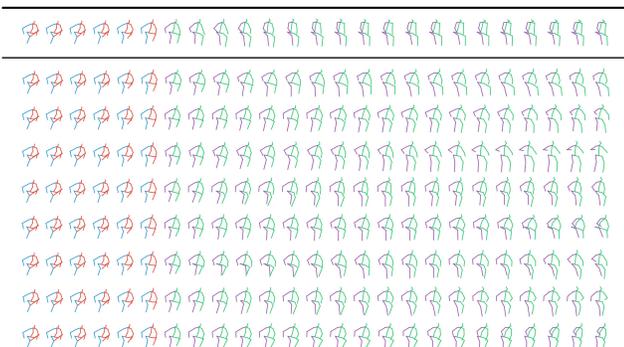


Figure 4. Additional qualitative evaluation of the diversity in human motion.

[2] Katerina Fragkiadaki, Sergey Levine, Panna Felsen, and Jitendra Malik. Recurrent network models for human dynamics. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4346–4354, 2015. 3

[3] Partha Ghosh, Jie Song, Emre Aksan, and Otmar Hilliges. Learning human motion models for long-term predictions. In *2017 International Conference on 3D Vision (3DV)*, pages 458–466. IEEE, 2017. 3

[4] Liang-Yan Gui, Yu-Xiong Wang, Xiaodan Liang, and José MF Moura. Adversarial geometry-aware human motion prediction. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 786–803, 2018. 3

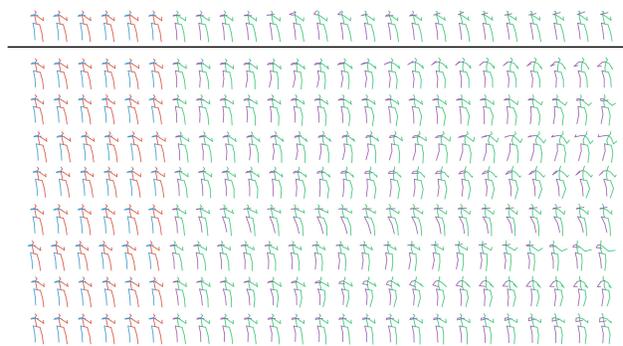


Figure 5. Additional qualitative evaluation of the diversity in human motion.

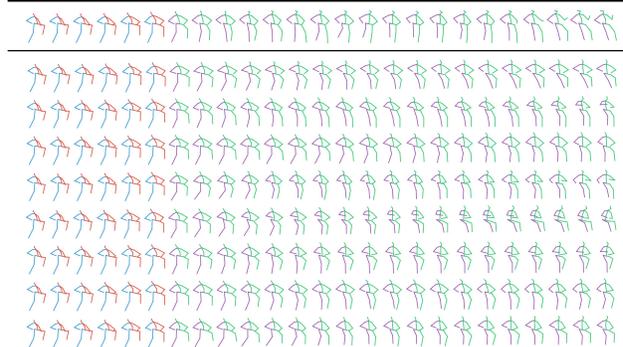


Figure 6. Additional qualitative evaluation of the diversity in human motion.

[5] Liang-Yan Gui, Yu-Xiong Wang, Deva Ramanan, and José MF Moura. Few-shot human motion prediction via meta-learning. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 432–450, 2018. 3

[6] Ashesh Jain, Amir R Zamir, Silvio Savarese, and Ashutosh Saxena. Structural-rnn: Deep learning on spatio-temporal graphs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5308–5317, 2016. 3

[7] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 1

[8] Wei Mao, Miaomiao Liu, Mathieu Salzmann, and Hongdong Li. Learning trajectory dependencies for human motion prediction. In *ICCV*, 2019. 3

[9] Julieta Martinez, Michael J Black, and Javier Romero. On human motion prediction using recurrent neural networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4674–4683. IEEE, 2017. 3

[10] Ronald J Williams and David Zipser. A learning algorithm for continually running fully recurrent neural networks. *Neural computation*, 1(2):270–280, 1989. 1