

# Learning Progressive Joint Propagation for Human Motion Prediction Supplementary

In this supplementary document, we provide materials not included in the main paper due to space constraints. Firstly, we provide more details about the DCT/IDCT process in Section 1. Secondly, we show the detailed architecture of our proposed network in Section 2. Lastly, detailed quantitative results are elaborated in Section 3.

## 1 Details of DCT/IDCT process

Following [3], we employ DCT to encode the temporal trajectories into the frequency domain. A key benefit is that, by discarding the high-frequencies, the DCT can provide a more compact representation to captures the smoothness of human motion. Specifically, given the trajectory  $\tilde{\mathbf{x}}_j = (\mathbf{x}_{j,1}, \mathbf{x}_{j,2}, \dots, \mathbf{x}_{j,t})$ , where  $j$  denotes the index of joint, the corresponding  $l^{th}$  DCT coefficient can be computed as:

$$C_{j,l} = \sqrt{\frac{2}{T}} \sum_{t=1}^T x_{j,t} \sqrt{\frac{1}{1 + \delta_{l,1}}} \cos\left(\frac{\pi}{2T}\right)(2t - 1)(l - 1) \quad (1)$$

where  $\delta_{m,n}$  denotes the *Kronecker* delta function :

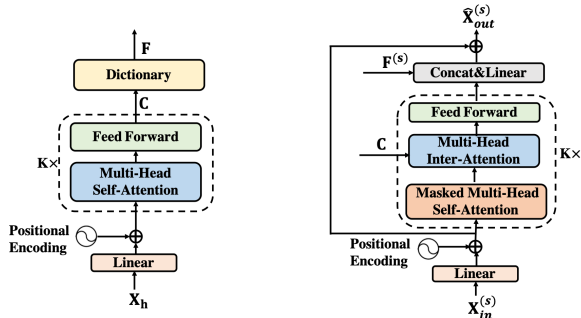
$$\delta_{m,n} = \begin{cases} 1 & \text{if } m = n \\ 0 & \text{if } m \neq n \end{cases} \quad (2)$$

Similarly, the final output can be transformed to temporal domain via the Inverse Discrete Cosine Transform (IDCT).

$$x_{j,t} = \sqrt{\frac{2}{T}} \sum_{l=1}^T C_{j,l} \sqrt{\frac{1}{1 + \delta_{l,1}}} \cos\left(\frac{\pi}{2T}\right)(2t - 1)(l - 1) \quad (3)$$

## 2 Network Architecture

Figure 1 illustrates the detailed architectures of our proposed network, including the transformer-encoder (left), dictionary (left), and the progressive decoder (right). Note that all operations (*i.e.* Linear, Feed Forward, Attention) are deployed for each joint with shared parameters. For more details of the progressive decoder and dictionary module, please kindly refer to Section 3.4 and 3.5 in our main paper.



**Fig. 1.** Detailed architecture of our proposed network. Left: **Encoder and Dictionary.**  $\mathbf{X}_h$  denotes the encoded historical trajectories of each joint.  $\mathbf{C}$  is the generated context feature and  $\mathbf{F}$  is the future dynamics summarized from the learned dictionary. Right: **Progressive decoder at stage  $s$ .**  $\mathbf{X}_{in}^{(s)}$ ,  $\hat{\mathbf{X}}_{out}^{(s)}$  are the input and corresponding output of the progressive decoder at stage  $s$ .  $\mathbf{F}^{(s)}$  is the future dynamics used in stage  $s$ . More details of the progressive decoding process can be found in Section 3.4 in our main paper.

### 3 More Quantitative Results

#### 3.1 Comparison with the State-of-the-art Methods

As mentioned in Section 4.3, we provide the full table of our proposed method compared with the state-of-the-art approaches [4, 3, 2, 5, 1] on Human3.6M dataset. As can be seen in Table 1 and Table 2, our results outperforms other methods on both MAE and MPJPE protocols.



## References

1. Gui, L.Y., Wang, Y.X., Liang, X., Moura, J.M.: Adversarial geometry-aware human motion prediction. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 786–803 (2018)
2. Li, C., Zhang, Z., Sun Lee, W., Hee Lee, G.: Convolutional sequence to sequence model for human dynamics. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 5226–5234 (2018)
3. Mao, W., Liu, M., Salzmann, M., Li, H.: Learning trajectory dependencies for human motion prediction. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 9489–9497 (2019)
4. Martinez, J., Black, M.J., Romero, J.: On human motion prediction using recurrent neural networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 2891–2900 (2017)
5. Wang, B., Adeli, E., Chiu, H.k., Huang, D.A., Niebles, J.C.: Imitation learning for human pose prediction. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 7124–7133 (2019)