

# LAL: Enhancing 3D Human Motion Prediction with Latency-aware Auxiliary Learning

## Supplementary Material

### 1. More Experimental Results

In main paper Table 3, we omit some prediction errors at certain timestamps and running time of CMU-Mocap [2] and 3DPW [6] due to limited space. Here, we provide the complete results of LTD [4], STS [5] and SPGSN [3] with LAL- $T_L$  in Table 1 and 2.

millisecond	CMU-Mocap average					running time (ms)
	80	160	320	400	1000	
LTD [4]	9.9	18.0	33.6	41.0	81.9	39.29
LTD-LAL-1	<b>9.3</b>	<b>17.4</b>	<b>32.1</b>	<b>39.8</b>	<b>81.8</b>	39.20
LTD-LAL-2	-	<b>16.7</b>	<b>30.9</b>	<b>39.2</b>	<b>80.7</b>	39.04
LTD-LAL-3	-	<b>15.6</b>	<b>29.4</b>	<b>38.5</b>	<b>79.7</b>	38.93
LTD-LAL-4	-	-	<b>28.7</b>	<b>38.3</b>	<b>78.4</b>	38.87
LTD-LAL-5	-	-	<b>28.1</b>	<b>37.8</b>	<b>76.5</b>	38.74
STS [5]	10.8	18.2	31.2	41.1	81.8	30.44
STS-LAL-1	11.4	<b>17.9</b>	<b>30.5</b>	<b>41.4</b>	<b>81.3</b>	30.39
STS-LAL-2	-	<b>16.8</b>	<b>28.4</b>	<b>39.5</b>	<b>80.7</b>	30.30
STS-LAL-3	-	<b>15.3</b>	<b>28.0</b>	<b>37.2</b>	<b>79.1</b>	30.18
STS-LAL-4	-	-	<b>27.2</b>	<b>37.6</b>	<b>76.2</b>	30.12
STS-LAL-5	-	-	<b>26.7</b>	<b>36.9</b>	<b>75.8</b>	30.06
SPGSN [3]	8.3	14.8	28.6	37.0	77.8	49.37
SPGSN-LAL-2	-	<b>13.4</b>	<b>27.1</b>	<b>35.7</b>	<b>76.8</b>	49.29
SPGSN-LAL-3	-	<b>12.5</b>	<b>25.8</b>	<b>34.9</b>	<b>75.1</b>	49.26
SPGSN-LAL-4	-	-	<b>25.1</b>	<b>34.8</b>	<b>74.9</b>	49.12
SPGSN-LAL-5	-	-	<b>24.5</b>	<b>33.9</b>	<b>74.2</b>	49.08

Table 1. Comparisons of MPJPE average errors in CMU-Mocap, where  $T_L = 1, 2, 3, 4, 5$ . Our LAL helps achieve progress without extra time cost during testing.

Similar to main paper Table 2 and 3, there is no LTD-LAL-1 for 3DPW, also no SPGSN-LAL-1 for CMU-Mocap and 3DPW, as their running time has already exceeded the *maximum* latency duration that it supports. Note that the sampling rate of 3DPW differs from the other two datasets, so the timestamps appear different (but we keep aligned with baselines to ensure fair comparison). Overall, LAL achieves improved results over baselines constantly, without the need for extra time cost during testing, which shows both the effectiveness and efficiency of our latency-aware prediction manner.

millisecond	3DPW Average					running time (ms)
	200	400	500	800	1000	
LTD [4]	35.6	67.5	80.2	106.8	117.8	41.66
LTD-LAL-2	<b>34.3</b>	<b>65.1</b>	<b>78.8</b>	<b>106.6</b>	<b>116.4</b>	41.50
LTD-LAL-3	<b>33.1</b>	<b>65.0</b>	<b>77.0</b>	<b>104.8</b>	<b>115.3</b>	41.44
LTD-LAL-4	<b>31.7</b>	<b>63.8</b>	<b>76.8</b>	<b>103.1</b>	<b>113.7</b>	41.29
LTD-LAL-5	<b>31.2</b>	<b>62.3</b>	<b>75.9</b>	<b>101.0</b>	<b>110.8</b>	41.20
STS [5]	37.8	67.5	77.3	106.7	112.2	31.84
STS-LAL-1	<b>37.2</b>	<b>66.8</b>	<b>76.5</b>	108.4	<b>111.4</b>	31.76
STS-LAL-2	<b>36.0</b>	<b>65.4</b>	<b>76.4</b>	<b>105.1</b>	<b>110.8</b>	31.75
STS-LAL-3	<b>33.8</b>	<b>64.8</b>	<b>75.3</b>	<b>103.6</b>	<b>109.0</b>	31.64
STS-LAL-4	<b>33.4</b>	<b>63.0</b>	<b>74.7</b>	<b>102.8</b>	<b>107.6</b>	31.55
STS-LAL-5	<b>34.6</b>	<b>61.9</b>	<b>73.1</b>	<b>101.6</b>	<b>106.6</b>	31.43
SPGSN [3]	32.9	64.5	76.2	104.0	111.1	50.61
SPGSN-LAL-2	<b>30.9</b>	<b>63.1</b>	<b>75.6</b>	<b>103.4</b>	<b>109.3</b>	50.46
SPGSN-LAL-3	<b>30.6</b>	<b>61.2</b>	<b>74.3</b>	<b>101.7</b>	<b>107.8</b>	50.39
SPGSN-LAL-4	<b>28.4</b>	<b>59.4</b>	<b>72.8</b>	<b>100.1</b>	<b>106.9</b>	50.35
SPGSN-LAL-5	<b>27.1</b>	<b>58.9</b>	<b>71.3</b>	<b>99.8</b>	<b>106.3</b>	50.28

Table 2. Comparisons of MPJPE average errors in 3DPW. The results and time consumption appear similar trend as in Table 1.

### 2. Ablation Experiments

#### 2.1. Training loss hyper-parameters

In main paper Figure 8, we show different  $(\lambda_P, \lambda_A)$  settings of Eq. (4), with STS [5] as the compared baseline and  $T_L = 3$ . Here, we provide the corresponding ablation with baseline LTD [4] and SPGSN [3] in Figure 1, and the improvement is most significant when  $(\lambda_P = 1.0, \lambda_A = 1.0)$  for both of them. Note that these experiments are all conducted on Human3.6M [1], and the hyper-parameter setting remain the same for each dataset and each  $T_L$  value.

#### 2.2. Ablations on refinement step

**Refinement loss hyper-parameters.** We provide the results of different  $(\lambda_Z, \lambda_U)$  weight settings in our main paper refinement loss Eq. (6). As shown in Figure 2 (left and middle), all of them can bring benefits more or less. We select  $\lambda_Z = 1.0, \lambda_U = 1.0$ , and maintain this setting for all baselines in our experiments.

**Why not merge the refinement loss into main training.**

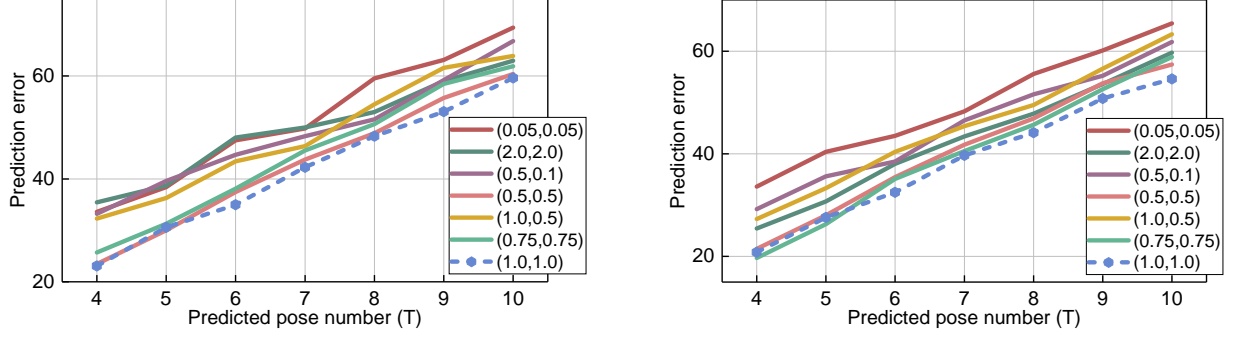


Figure 1. Results of different  $(\lambda_P, \lambda_A)$  settings of Eq. (4) in LTD-LAL-3 (left) and SPGSN-LAL-3 (right). The best is achieved when  $\lambda_P = 1.0, \lambda_A = 1.0$ , drawn in dashed blue line.

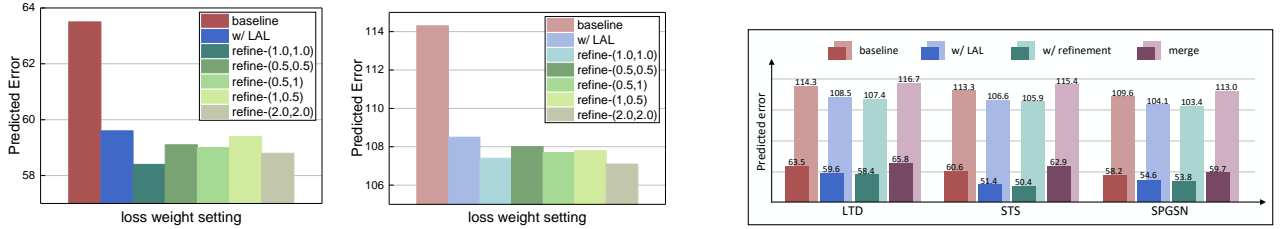


Figure 2. Left and middle: results of different  $(\lambda_Z, \lambda_U)$  settings in Eq. (6) for LTD-LAL-3. Left subfigure shows errors at 400 ms, while middle shows 1000 ms. We set  $\lambda_Z = 1.0, \lambda_U = 1.0$  that yields constant improvement. Right: merging refinement loss into LAL as an end-to-end training manner, showing no improvement compared to baselines. Dark colors indicate 400 ms errors while faded 1000 ms.

As we treat the refinement as an optional post-processing step, we should validate the effectiveness of using such separate training rather than end-to-end manner. We merge the feature statistics-based alignment loss into our main training, which let the primary prediction branch approximate the distribution of the auxiliary latent and final outputs from scratch, i.e., directly involving  $\mathcal{L}_{\text{fea}, Z}$  and  $\mathcal{L}_{\text{fea}, U}$  into our LAL. However, the auxiliary branch at the beginning is not well-trained that tends to generate wrong patterns of outputs (especially the final output  $U$ ), which may confuse the primary branch whether to fit towards the ground truth or the unstable auxiliary outputs. The results are shown in Figure 2 (right) purple bars, with no improvement and even higher errors than baselines, which verifies the above reasoning.

**Limitations.** In real-world scenarios, the intelligent system may confront incomplete observations due to the occlusion of human body or noise perturbation, which greatly impedes the prediction accuracy. Improving algorithm robustness against these factors is a worthy study that will further contribute to this task.

## References

[1] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments.

*IEEE transactions on pattern analysis and machine intelligence*, 36(7):1325–1339, 2013. 1

[2] Chen Li, Zhen Zhang, Wee Sun Lee, and Gim Hee Lee. Convolutional sequence to sequence model for human dynamics. In *CVPR*, pages 5226–5234, 2018. 1

[3] Maosen Li, Siheng Chen, Zijing Zhang, Lingxi Xie, Qi Tian, and Ya Zhang. Skeleton-parted graph scattering networks for 3d human motion prediction. In *ECCV*, pages 18–36. Springer, 2022. 1

[4] Wei Mao, Miaomiao Liu, Mathieu Salzmann, and Hongdong Li. Learning trajectory dependencies for human motion prediction. In *ICCV*, pages 9489–9497, 2019. 1

[5] Theodoros Sofianos, Alessio Sampieri, Luca Franco, and Fabio Galasso. Space-time-separable graph convolutional network for pose forecasting. In *ICCV*, pages 11209–11218, 2021. 1

[6] Timo Von Marcard, Roberto Henschel, Michael J Black, Bodo Rosenhahn, and Gerard Pons-Moll. Recovering accurate 3d human pose in the wild using imus and a moving camera. In *ECCV*, pages 601–617, 2018. 1