

# Appendix: Time-Aware Multi-Scale RNNs for Time Series Modeling

Zipeng Chen, Qianli Ma\* and Zhenxi Lin

School of Computer Science and Engineering,  
South China University of Technology, Guangzhou, China

zipengchencs@foxmail.com, qianlima@scut.edu.cn, zhenxi\_lin@foxmail.com

## A Detailed results on human motion prediction

For human motion prediction, the overall results of TP-RNN and TAMS-TP-RNN on the H3.6M data set are shown in Table 1. The comparison results show that the proposed TAMS-TP-RNN is significantly better than TP-RNN, yielding lower Euler angle errors and higher best rates, which strongly verifies the superiority of the proposed model.

## B More examples of music clips

For music genre recognition, more examples of music clips are shown in Fig. 1. Music series in the yellow region has small changes, so it's no need for the model to be updated frequently and a larger scale would be chosen. On the contrary, series in the purple region varies greatly, thus the choice of scale is more flexible. Therefore, at different time steps,

different temporal scales would be chosen to model the corresponding dynamic temporal patterns.

## C Ablation study on MTS data sets

We design five ablation models to verify the effectiveness of each component. The experiments are conducted on 15 MTS classification data sets and the results are shown in Table 2.

## D Computational efficiency analysis

The experiments are implemented using Python 2.7 and Tensorflow 1.11.0 on a personal computer with an Intel Core i7-8700K 3.70GHz CPU, a 64GB RAM, and a GeForce GTX 2080-Ti 11G graphics card. For "FaceDetection", the number of instances for training and testing are 5890 and 3524. The proposed model spent 12.77 seconds on training for each epoch with a batch size 64 and spent 2.52 seconds on testing for the whole testing set.

\*Qianli Ma is the corresponding author.

	TP-RNN						TAMS-TP-RNN					
Milliseconds	80	160	320	400	560	1000	80	160	320	400	560	1000
Walking	<b>0.25</b>	0.41	0.58	0.65	0.74	0.77	<b>0.25</b>	<b>0.39</b>	<b>0.57</b>	<b>0.63</b>	<b>0.73</b>	<b>0.76</b>
Eating	<b>0.20</b>	<b>0.33</b>	<b>0.53</b>	<b>0.67</b>	<b>0.84</b>	<b>1.14</b>	0.21	0.34	0.55	0.69	0.87	1.19
Smoking	0.26	<b>0.47</b>	<b>0.88</b>	0.90	0.98	1.66	<b>0.25</b>	<b>0.47</b>	<b>0.88</b>	<b>0.87</b>	<b>0.93</b>	<b>1.61</b>
Discussion	<b>0.30</b>	<b>0.66</b>	0.96	1.04	1.39	<b>1.74</b>	0.31	<b>0.66</b>	<b>0.95</b>	<b>1.03</b>	<b>1.34</b>	1.75
Directions	0.38	0.59	<b>0.75</b>	<b>0.83</b>	<b>0.95</b>	<b>1.38</b>	<b>0.36</b>	<b>0.58</b>	0.77	0.86	0.96	<b>1.38</b>
Greeting	0.51	0.86	1.27	1.44	1.72	1.81	<b>0.48</b>	<b>0.81</b>	<b>1.21</b>	<b>1.37</b>	<b>1.64</b>	<b>1.75</b>
Phoning	0.57	1.08	<b>1.44</b>	1.59	<b>1.47</b>	<b>1.68</b>	<b>0.56</b>	<b>1.06</b>	<b>1.44</b>	<b>1.58</b>	1.49	1.73
Posing	0.42	0.76	1.29	1.54	<b>1.75</b>	<b>2.47</b>	<b>0.32</b>	<b>0.54</b>	<b>1.10</b>	<b>1.38</b>	1.80	2.74
Purchases	<b>0.59</b>	<b>0.82</b>	<b>1.12</b>	<b>1.18</b>	<b>1.52</b>	<b>2.28</b>	0.63	0.87	1.13	1.20	1.53	2.30
Sitting	0.41	0.66	1.07	1.22	1.35	1.74	<b>0.38</b>	<b>0.61</b>	<b>1.00</b>	<b>1.16</b>	<b>1.27</b>	<b>1.67</b>
Sitting down	<b>0.41</b>	0.79	1.13	1.27	1.47	1.93	<b>0.41</b>	<b>0.77</b>	<b>1.09</b>	<b>1.22</b>	<b>1.42</b>	<b>1.90</b>
Taking photo	0.26	<b>0.51</b>	<b>0.80</b>	0.95	1.08	1.35	<b>0.25</b>	<b>0.51</b>	<b>0.80</b>	<b>0.92</b>	<b>1.00</b>	<b>1.28</b>
Waiting	<b>0.30</b>	<b>0.60</b>	<b>1.09</b>	<b>1.31</b>	1.71	<b>2.46</b>	0.32	0.62	1.10	<b>1.31</b>	<b>1.69</b>	2.47
Walking dog	0.53	0.93	<b>1.24</b>	<b>1.38</b>	1.73	1.98	<b>0.49</b>	<b>0.85</b>	1.39	1.45	<b>1.68</b>	<b>1.90</b>
Walking together	<b>0.23</b>	<b>0.47</b>	0.67	0.71	0.78	1.28	<b>0.23</b>	<b>0.47</b>	<b>0.66</b>	<b>0.70</b>	<b>0.77</b>	<b>1.26</b>
Average	0.37	0.66	0.99	1.11	1.30	<b>1.71</b>	<b>0.36</b>	<b>0.64</b>	<b>0.98</b>	<b>1.09</b>	<b>1.27</b>	<b>1.71</b>
No. best	7	7	8	5	5	7	<b>11</b>	<b>12</b>	<b>10</b>	<b>11</b>	<b>10</b>	<b>9</b>
Best rate	0.47	0.47	0.53	0.33	0.33	0.47	<b>0.73</b>	<b>0.80</b>	<b>0.67</b>	<b>0.73</b>	<b>0.67</b>	<b>0.60</b>

Table 1: Euler angle errors for 15 activities on the Human3.6M data set.

Data set	LSTM +MSFD(M)+TAFM	CW-LSTM +TAFM	LSTM +MSFD(M)	LSTM +MSFD(S)	CW-LSTM	LSTM
ArticulatoryWordRecognition	0.973	0.967	<b>0.976</b>	0.947	0.947	0.947
AtrialFibrillation	0.400	<b>0.467</b>	0.333	0.400	0.333	0.400
BasicMotions	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>	0.975
CharacterTrajectories	<b>0.994</b>	0.993	<b>0.994</b>	0.993	0.990	0.985
FaceDetection	<b>0.602</b>	0.596	0.595	0.562	0.594	0.573
HandMovementDirection	<b>0.473</b>	0.446	0.432	0.284	<b>0.473</b>	0.311
Heartbeat	<b>0.756</b>	0.712	0.751	0.732	0.722	0.698
MotorImagery	0.590	0.580	<b>0.600</b>	0.530	0.520	0.570
NATOPS	<b>0.956</b>	<b>0.956</b>	0.950	0.950	0.950	<b>0.956</b>
PEMS-SF	<b>0.890</b>	<b>0.890</b>	0.873	0.884	<b>0.890</b>	0.879
Pen Digits	<b>0.981</b>	0.979	0.977	0.976	0.976	0.972
Phoneme	<b>0.203</b>	0.186	0.194	0.178	0.178	0.189
SelfRegulationSCP2	<b>0.561</b>	0.544	<b>0.561</b>	0.544	0.544	0.528
SpokenArabicDigits	<b>0.990</b>	0.988	0.989	0.980	0.986	0.987
StandWalkJump	<b>0.400</b>	0.267	0.333	0.267	0.333	<b>0.400</b>

Table 2: Ablation study on 15 MTS classification data sets.

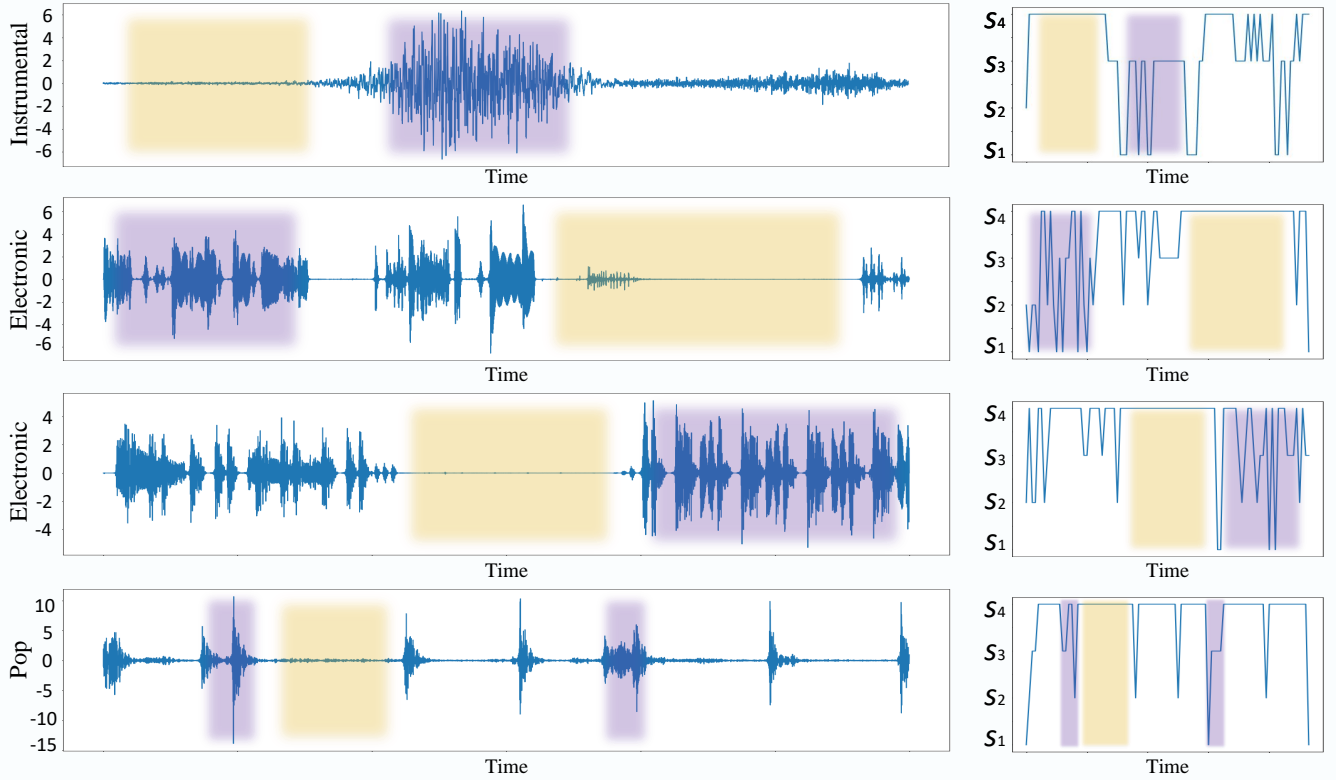


Figure 1: Five music clips and the corresponding weights assigned to the scales. Subfigures in the first column are five music clips of different genres, while subfigures in the second column are the scale with the largest weight at each time step  $t \text{ MOD } s_K = 0$  ( $s_K = 64$ ), corresponding to music clips in the first column. Series in the yellow region prefers the largest scale due to its small changes, while series in the purple region prefers multiple scales due to its sharp fluctuations.