

1 **Reply to Reviewer #1**

2 **Q1:** What other ways to generate fake sequences may be suitable for this problem?

3 A1: That is a good question. Although we here used a simple generation strategy (randomly shuffling some time steps), it has  
4 verified the improvement of the encoder’s ability will boost the performance of clustering. In the future, we will consider to use  
5 GAN to generate some more difficult fake sequences to further improve the ability of the encoder.

6 **Reply to Reviewer #2**

7 **Q1:** Comparison with other state-of-the-art deep clustering methods which are not designed for time-series.

8 A1: Following your suggestion, we compare our method with two state-of-the-art deep clustering methods (i.e., DEC (Xie et al.,  
9 2016 ICML) and IDEC (Guo et al., 2017 IJCAI) ). The results of DEC and IDEC are obtained by running the original authors  
10 implementation code and we use the Rand Index (RI) to evaluate performance, consistent with our paper. As shown in Table 1,  
11 our method still achieves the best performance with the highest average RI of 0.7714. We will try to add these comparison results  
to the main text or the appendix, modulo the page limit.

Table 1: Comparisons on 36 time series datasets (The No. of datasets is consistent with the one in Table 2 in main text)

Dataset	DEC(RI)	IDEC(RI)	DTCR(RI)	DTCR(NMI)	DTCR(ACC)	Dataset	DEC(RI)	IDEC(RI)	DTCR(RI)	DTCR(NMI)	DTCR(ACC)
1	0.5817	0.6210	<b>0.6868</b> (0.0026)	0.5513(0.0022)	0.6914(0.0028)	19	0.5423	0.5423	<b>0.5617</b> (0.0006)	0.1150(0.0005)	0.6733(0.0037)
2	0.5954	0.6276	<b>0.8046</b> (0.0018)	0.5613(0.0013)	0.5667(0.0013)	20	0.8590	0.8626	<b>0.8638</b> (0.0007)	0.5503(0.0006)	0.5940(0.0006)
3	0.4947	0.6053	<b>0.9000</b> (0.0001)	0.7610(0.0001)	0.8500(0.0001)	21	0.7435	0.7324	<b>0.7686</b> (0.0036)	0.4094(0.0041)	0.8395(0.0027)
4	0.4737	0.4789	<b>0.8105</b> (0.0033)	0.5310(0.0035)	0.8500(0.0031)	22	0.7484	0.7607	<b>0.7739</b> (0.0014)	0.2599(0.0010)	0.4430(0.0009)
5	0.6859	0.6870	<b>0.7501</b> (0.0022)	0.5021(0.0020)	0.6333(0.0017)	23	0.9447	0.9447	<b>0.9549</b> (0.0037)	0.9296(0.0033)	0.6286(0.0035)
6	0.5348	0.5350	<b>0.5357</b> (0.0011)	0.0468(0.0010)	0.8929(0.0015)	24	0.4263	<b>0.8091</b>	<b>0.8091</b> (0.0038)	0.5734(0.0029)	0.7805(0.0027)
7	0.4921	0.5767	<b>0.9286</b> (0.0016)	0.8122(0.0015)	0.9643(0.0016)	25	0.8189	<b>0.9030</b>	0.9023(0.0023)	0.6948(0.0017)	0.6300(0.0017)
8	0.9294	0.7347	<b>0.9682</b> (0.0032)	0.9418(0.0027)	0.9118(0.0028)	26	0.5732	0.6900	<b>0.8769</b> (0.0033)	0.9178(0.0027)	0.9285(0.0021)
9	0.7785	0.7786	<b>0.7825</b> (0.0008)	0.4553(0.0007)	0.8050(0.0008)	27	0.6514	0.6572	<b>0.8354</b> (0.0016)	0.6121(0.0017)	0.9056(0.0016)
10	0.5029	0.5330	<b>0.6075</b> (0.0024)	0.1180(0.0028)	0.7083(0.0025)	28	0.8837	0.8893	<b>0.9223</b> (0.0021)	0.6663(0.0019)	0.3816(0.0018)
11	0.6422	0.6233	<b>0.6648</b> (0.0034)	0.3691(0.0028)	0.8000(0.0026)	29	0.8841	0.8857	<b>0.9168</b> (0.0022)	0.8989(0.0018)	0.6492(0.0018)
12	0.5103	0.5114	<b>0.9638</b> (0.0032)	0.8056(0.0034)	0.8525(0.0039)	30	0.4984	0.5017	<b>0.5659</b> (0.0006)	0.3115(0.0008)	0.5658(0.0013)
13	0.4981	0.4974	<b>0.6398</b> (0.0011)	0.4200(0.0013)	0.7867(0.0011)	31	0.4991	0.4991	<b>0.8286</b> (0.0028)	0.3895(0.0033)	0.9000(0.0023)
14	<b>0.5963</b>	0.4956	0.5362(0.0035)	0.0989(0.0036)	0.6381(0.0031)	32	0.6293	0.6338	<b>0.6984</b> (0.0025)	0.4713(0.0019)	0.4253(0.0020)
15	0.5099	0.5099	<b>0.5759</b> (0.0017)	0.2248(0.0016)	0.7031(0.0015)	33	0.5007	0.5016	<b>0.7114</b> (0.0014)	0.4614(0.0009)	0.8622(0.0011)
16	0.5311	0.5519	<b>0.5913</b> (0.0016)	0.2289(0.0014)	0.7213(0.0017)	34	0.5679	0.5597	<b>0.7338</b> (0.0006)	0.0228(0.0001)	0.6775(0.0013)
17	0.6475	0.6220	<b>0.9763</b> (0.0016)	0.9653(0.0009)	0.9770(0.0017)	35	0.4913	0.5157	<b>0.6271</b> (0.0039)	0.2887(0.0038)	0.7963(0.0047)
18	0.7059	0.6800	<b>0.7982</b> (0.0028)	0.4661(0.0017)	0.7700(0.0023)	36	0.8893	0.8947	<b>0.8984</b> (0.0003)	0.5448(0.0003)	0.1630(0.0002)

12 **Q2:** Discussion on how to use the presented method if the number of clusters  $K$  is not known a priori.

14 A2: For most of the clustering algorithms, the choice of hyper-parameter  $K$  is indeed a tricky problem. We could use the popular  
15 "Elbow" method or Gap Statistic (Tibshirani et al., 2001 Royal Statistical Society) to choose  $K$ .

16 **Q3:** Whether the analysis and results in sections 4.3.2 and 4.3.3 are general properties of the algorithm?

17 A3: In fact, we conducted the experiments of sections 4.3.2 and 4.3.3 on all of the other datasets. The performance was consistent  
18 with the results in the article. We will add these results to the appendix.

19 **Reply to Reviewer #3**

20 **Q1:** The motivation behind the use of the k-means loss in this setting and how does this hold up compared to DTW?

21 A1: According to Zha et al. (2002 NIPS), K-means solves the minimum of a sum-of-squares cost function using coordinate  
22 descent, which is prone to local minima. Thus, they reformulate it as a trace maximization problem, which obtains optimal  
23 global solutions. Inspired by this, we reformulated K-means loss and integrated it into the autoencoder to guide representation  
24 learning. DTW may not be compatible with the trace maximization problem. This requires further investigation.

25 **Q2:** How robust is the method to hyper-parameters such as  $T$  and  $\lambda$ ?

26 A2: We performed experiments with  $T \in [10, 20, 30, 40]$  to explore the impact of  $T$  and verified our method is robust to  $T$ . We  
27 will add these results to the appendix.  $\lambda$  determines the weight of K-means loss, which affects the performance to some extent.  
28 As mentioned in the main text, we used a grid search approach with  $\lambda$  from  $[1, 1e-1, 1e-2, 1e-3]$ .

29 **Q3:** Do datasets contain a train test split? Is the entire dataset used for testing?

30 A3: In the UCR time series benchmark, each data set has a default training and test set. To make a fair comparison, we adopted  
31 the same protocol used in USSL (Zhang et al., 2018 TPAMI), which is trained on the training set and evaluated on the test set.

32 **Q4:** The experiment should be run a few times to show the impact of the random initialization.

33 A4: Following your suggestion, we run each experiment 5 times and record means and standard deviations in Table 1. The  
34 values in parentheses present standard deviations, which are all less than 0.01. DTCR still achieves the best performance with the  
35 lowest average rank of 2.9583 and the highest average RI of 0.7714. We will update these results to the main text.

36 **Q5:** The proposed method would also be more convincing if larger more complex datasets were evaluated.

37 A5: Thank you for your suggestion. Following the YADING paper (Ding et al., 2015 VLDB), which was designed for large-scale  
38 time series clustering, a larger and more complex dataset (StarLightCurves: 9, 236 samples, each sample’s length is 1, 024) is  
39 used to evaluate our method. We achieve NMI (normalized mutual information) metric of 0.6731, while NMI of YADING is  
40 0.6000, NMI of DEC is 0.6058, and NMI of IDEC is 0.6056. Since USSL did not provide source code, we can not obtain the  
41 result of it for comparison. We will try to add the comparison to the main text or the appendix, modulo the page limit.

42 **Q6:** Comparison with other state-of-the-art deep clustering methods which are not designed for time-series.

43 A6: Please see Table 1. As the results shown, our method outperforms DEC as well as IDEC significantly.

44 **Q7:** Accuracy and NMI also should be included to evaluate the method performance.

45 A7: Due to the limitations of the layout and comparison methods results, we only show the results of RI. In fact, we also recorded  
46 NMI and accuracy. The results are shown in Table 1. And our method also achieve the lowest average rank of 2.1944 while  
47 USSL achieve 2.2361 on the NMI metric. We will add them to the text, modulo the page limit.