

1 We thank the reviewers for their meaningful and valuable comments, which help to improve the quality of our work.

2 **R1 on paper motivation:** Figure 1 in our submission illustrates the limitation of using MSE as a loss for training deep
3 forecasting models. Since MSE is similar for the three predictions in (a), (b) and (c), a gradient-based optimization is
4 unable to produce training signals preferring predictions (b) and (c) over (a). By no means we wanted to claim that
5 certain existing approaches are unable to perform step prediction in this simple example. This work does not focus
6 on designing new forecasting models, but introduces the STDL loss function as an alternative to MSE. STDL is thus
7 model-agnostic and can be used for training various forecasting models - as shown in experiments and below.

8 **R1 on state-of-the-art methods:** We thank R1 for suggesting to compare our results to online regression with adaptive
9 parameters (forgetting recursive least squares, adaptive Kalman filters). Although these are historical approaches for
10 Bayesian inference in State Space Models (SSMs), their direct application to multi-step forecasting is not straightforward
11 because they require input data for adaptation at each time step. For this reason, several state-of-the-art multi-step
12 approaches combine SSMs and deep learning based on Seq2Seq architectures [1, 2, 3, 4].

13 To fulfill R1 requests, we perform additional experiments (shown in blue) on the Traffic dataset (Table 4 in submission).
14 The results of the Deep AR baseline¹ (obtained with GitHub code) is still outperformed by a simple Seq2Seq model
15 trained with STDL (results shown in submission, column 4 in black), and equivalent in temporal metrics. Training Deep
16 AR with STDL would be an interesting future exploration. Finally, we provide results when training the recent TT-RNN
17 (refs [48,49] in submission) with STDL, reaching the best shape and temporal performances. These new results further
18 highlight the importance of STDL ; we will be glad to add these comparisons in the final paper if accepted.

Eval loss		LSTNet-rec (MSE)	TT-RNN (MSE)	Deep AR (MSE)	Seq2Seq (STDL)	TT-RNN (STDL)
Euclidian	MSE (x100)	1.74 ± 0.11	0.840 ± 0.106	0.966 ± 0.068	1.00 ± 0.260	0.930 ± 0.09
Shape	DTW (x100)	42.0 ± 2.2	25.9 ± 1.99	27.8 ± 1.55	23.0 ± 1.62	21.4 ± 0.79
	Ramp (x10)	9.00 ± 0.577	6.71 ± 0.546	7.56 ± 0.42	5.93 ± 0.235	5.27 ± 0.27
Time	TDI (x10)	25.7 ± 4.75	17.8 ± 1.73	14.6 ± 0.94	14.4 ± 1.58	15.7 ± 1.02
	Hausdorff	2.34 ± 1.41	2.19 ± 0.12	2.04 ± 0.11	2.13 ± 0.514	1.88 ± 0.153

19 **R2 on more complex datasets:** As requested, we provide additional experiments on 2 more complex datasets:
20 household electricity consumption and solar energy. The former corresponds to a multivariate forecasting problem
21 involving 10 exogenous input variables (global intensity, voltage, sub-metering, date, *etc*), requiring the extraction of
22 complex interactions in data for spiky patterns prediction. The latter has very fine time granularity (10min vs 1h for
23 Traffic), needing to extract accurate time features. The results shown below again illustrate the superiority of training
24 Seq2Seq models with SDTL compared to MSE.

Method	Household electricity consumption			Solar energy		
	MSE (x10)	DTW	TDI	MSE (x1000)	DTW (x100)	TDI (x10)
Seq2Seq MSE	18.3 ± 2.5	4.54 ± 0.40	2.49 ± 0.26	13.7 ± 1.5	24.3 ± 3.4	12.9 ± 1.4
Seq2Seq STDL	19.9 ± 2.4	3.85 ± 0.26	2.30 ± 0.59	14.4 ± 0.57	20.9 ± 1.1	5.71 ± 0.83

25 **R2 on α tuning:** α is chosen on a validation set, by selecting the lowest value for which \mathcal{L}_{shape} gets comparable
26 performance than a reference DTW_{γ} trained model. This setup will be added in the final version if accepted.

27 **R3 on reporting STDL as evaluation metric:** these results will be added in our tables if accepted.

28 **R3 on training time:** 1 training epoch with our Seq2Seq GRU network takes about 0.5s for MSE vs 1.7s for SDTL on
29 Synthetic (1s vs 8s on ECG5000, 3s vs 33s on Traffic). The overhead is due to the sequential computation of the STDL
30 (dynamic programming in forward and backward passes). Note that no overhead is involved at test time.

31 **R3 on choice of δ :** We choose δ as the euclidean loss (paper 1. 102-103), which is common for computing DTW, but
32 any other distance (*e.g.* mean absolute error) could be employed.

33 **R3 on code sharing:** source code will be made available on GitHub after acceptance.

34 **R3 on ECG:** predicting the shape and time interval between heartbeats could be helpful for cardiologists, to detect
35 abnormal heartbeats such as 'premature ventricular contraction'.

36 **R3 on feature interpretation:** understanding the effects of our shape and time loss terms on the learned features is an
37 interesting but non trivial perspective. A possible direction to this end is to use feature visualization techniques⁵.

38 **R3 on confidence intervals:** we could use MC Dropout (Gal *et. al.*, ICML'16) to compute the predictive distribution
39 of trajectories ; or embed the STDL loss in a deep SSM architecture suited for probabilistic forecasting.

¹ D. Salinas, V. Flunkert, J. Gasthaus. "DeepAR: Probabilistic forecasting with autoregressive recurrent networks", ICML 2017

² S. Rangapuram, M. Seeger, J. Gasthaus, L. Stella, Y. Wang, "Deep state space models for time series forecasting", NeurIPS2018

³ X.Jin, S.Li, Y. Zhang, X.Yan, "Multi-step deep autoregressive forecasting with latent states", ICML 2019 Time Series Workshop

⁴ Y. Wang, A. Smola, D. C Maddix, J. Gasthaus, D. Foster, T. Januschowski. "Deep factors for forecasting", ICML 2019

⁵ Andrej Karpathy, Justin Johnson, and Li Fei-Fei. "Visualizing and understanding recurrent networks", ICLR 2016