

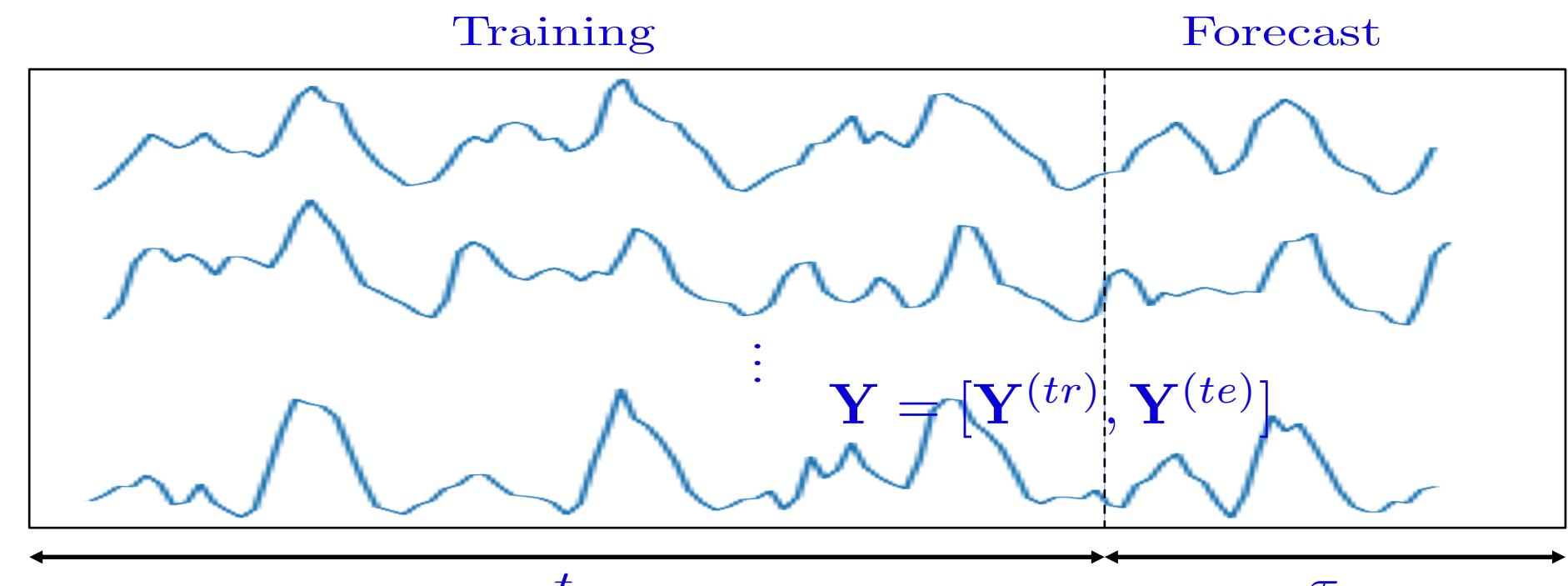
Think Globally, Act Locally: A Deep Neural Network Approach to High-Dimensional Time Series Forecasting

Rajat Sen[&], Hsiang-Fu Yu[&] and Inderjit Dhillon^{#,&}

^{&Amazon, #University of Texas at Austin}

High-Dimensional Time-Series Forecasting

- Applications like item demand forecasting, web-traffic forecasting can have millions of correlated time-series evolving together e.g forecasting future demand of items for a retailer like Amazon

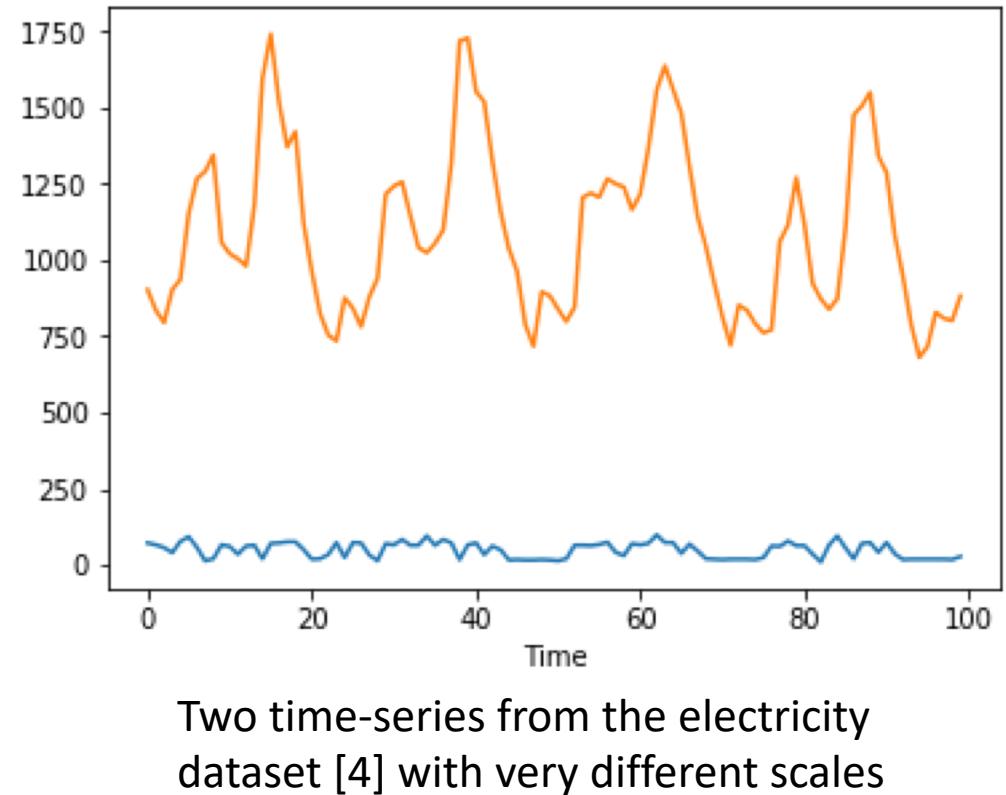


- Additionally there can be covariates at each time point, which are either common to all time-series (time of day, day of week etc.) or are unique to each time-series.

Key Challenges in Deep Temporal Models

Issues with Normalization

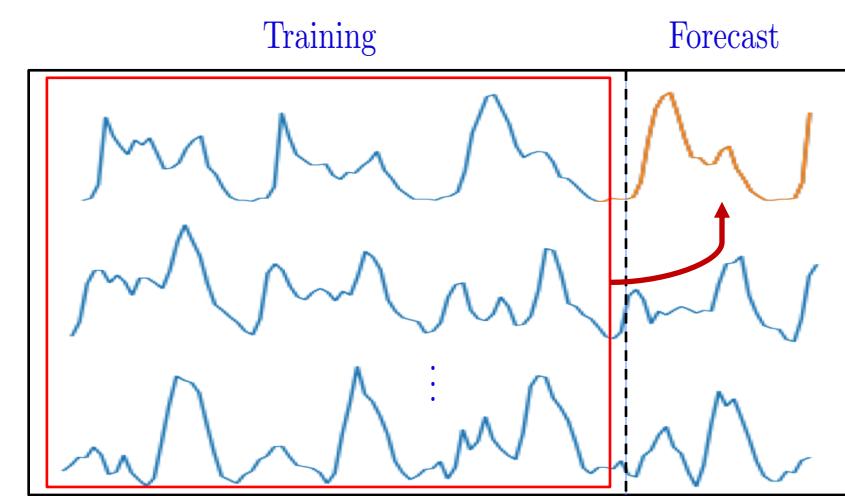
- The individual time-series may have vastly different scales. For instances, in demand forecasting some popular items may have orders of magnitude more demand, than niche items.
- This creates issues in training deep models based on LSTM or Temporal Convolution Networks (TCN). The datasets need to be normalized and the choice of scaling parameters can have impact on performance.
- For instance, [1] uses standard normalization, while [2] scales all the time-series with the value in their first time-index.



Two time-series from the electricity dataset [4] with very different scales

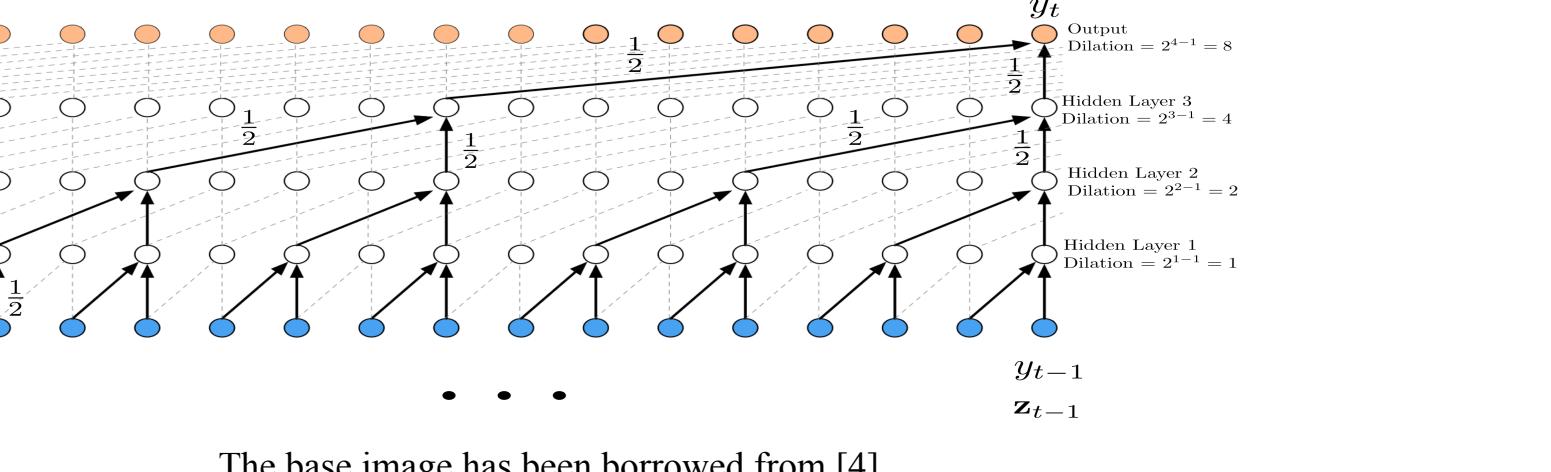
Using Local and Global during Prediction

- Most deep learning time-series models based on LSTM, Temporal Convolutions [1,2] use the whole dataset while training, but during forecasting of an individual time-series the predictions are a function of its local past.
- Global properties may be useful during prediction time. For instance, in retail demand forecasting, past values of similar items can be leveraged while predicting the future for a certain item.
- TRMF [3] can express all the time-series as linear combinations of basis time-series. These basis time-series can capture global patterns during prediction. However, TRMF can only model linear temporal dependencies.



LeveledInit: Handling Diverse Scales with TCN

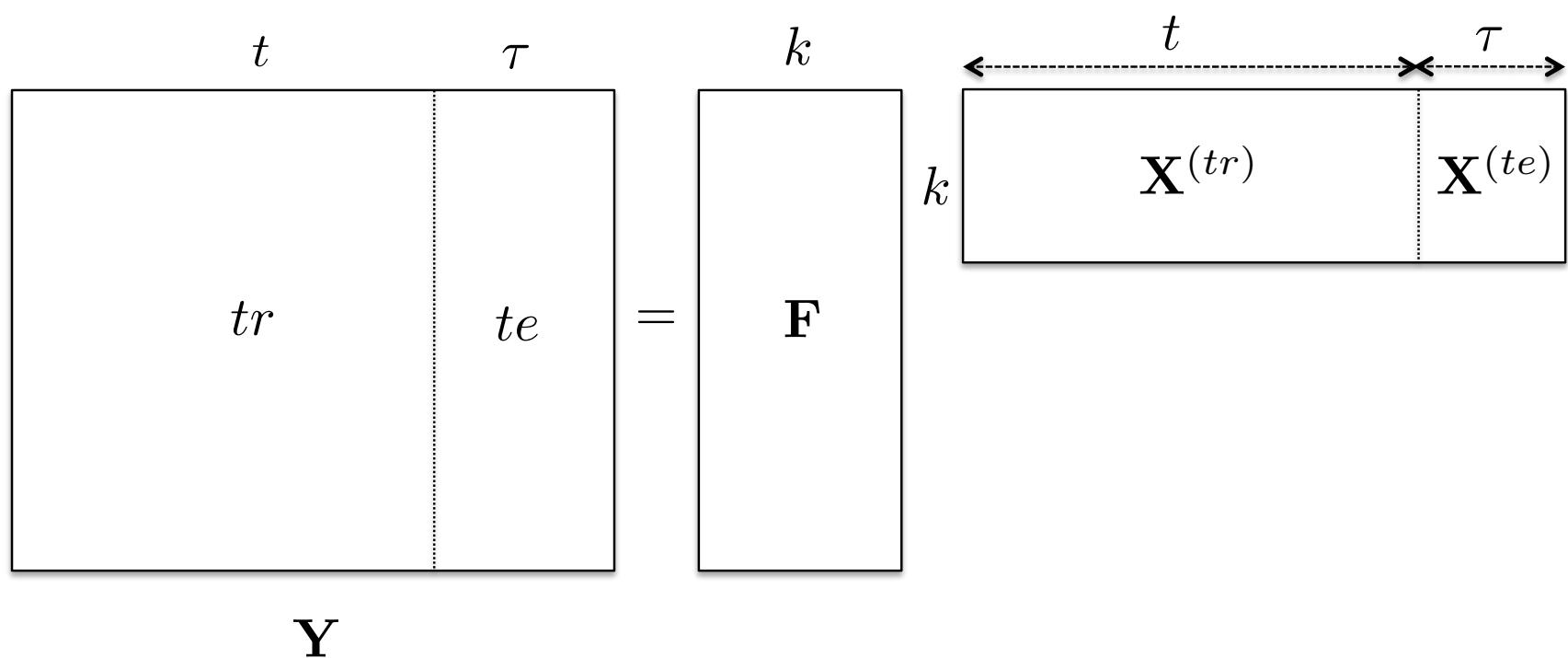
- In order to handle datasets, with a lot of variation in scale among different time-series, we propose a simple initialization scheme for Temporal Convolution Networks (TCN).
- LeveledInit:** We recommend setting the filter weights equal to $1/k$, where k is the filter size, in all layers except the input layer. In the input layer, all weights in the channel that takes in the original time-series are set to $1/k$, while all other weights are initialized to 0.
- We show that under some configurations of the TCN (filter size, dilation) LeveledInit leads to the network initially predicting the average of the past values in the dynamic range of the network, and then the network learns to predict variations around this mean value. The variations around the mean are relatively scale free and can alleviate the problems with scale.



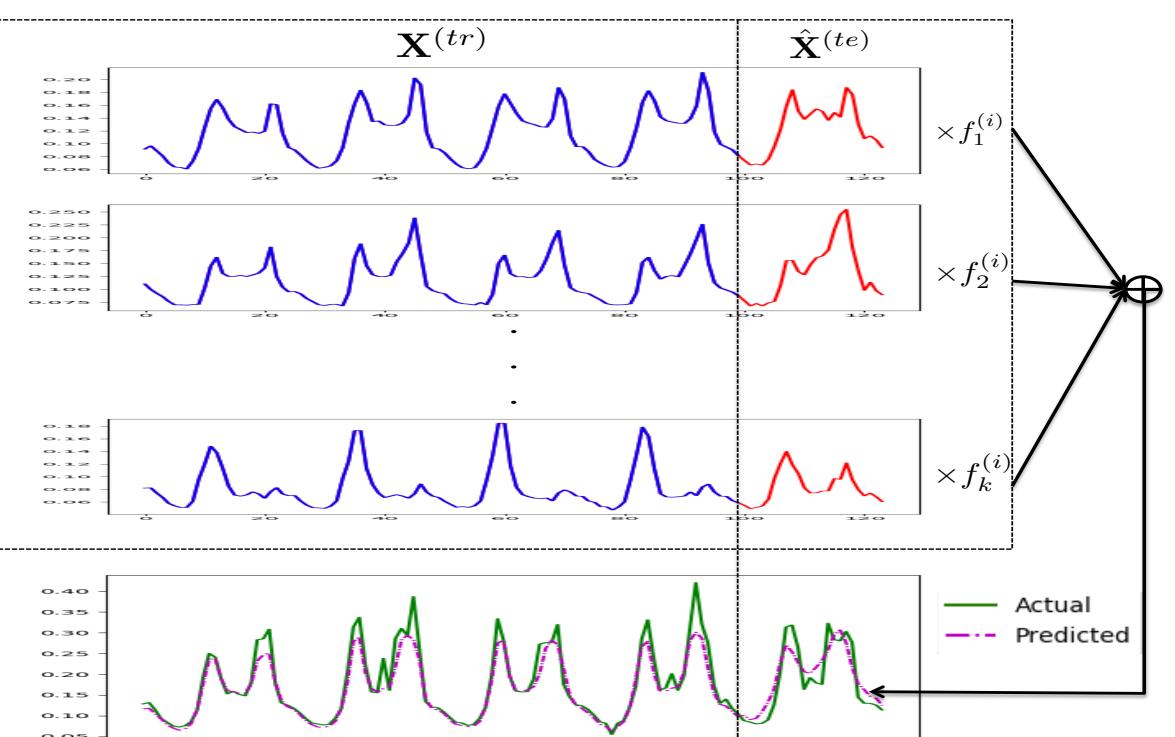
The base image has been borrowed from [4]

TCN-MF: Our Global Model

- Similar to [3], the idea is to decompose the training matrix $\mathbf{Y}^{(tr)}$ into factors \mathbf{F} and $\mathbf{X}^{(tr)}$ such that $\mathbf{X}^{(tr)}$ has temporal features. In other words $\mathbf{X}^{(tr)}$ consists of k basis time-series which can be predicted into the future, where the rank $k \ll n$.
- In order to compute non-linear temporal dependencies we use a TCN (trained concurrently) to regularize $\mathbf{X}^{(tr)}$ during the factorization. This encourages temporal structure in $\mathbf{X}^{(tr)}$.



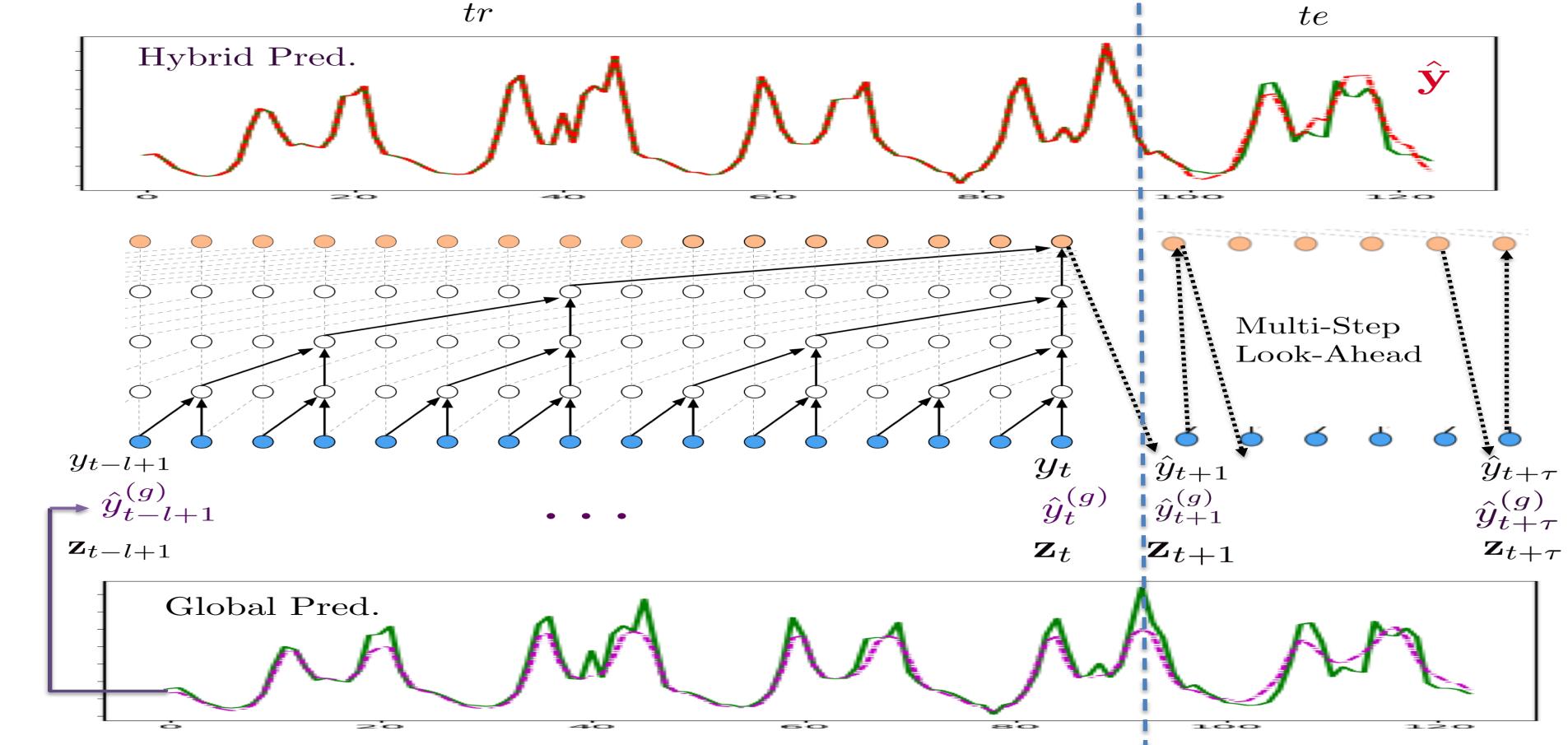
- Temporal regularization by TCN: $\mathcal{R}(\mathbf{X}^{(tr)} | \mathcal{T}_X(\cdot)) := \frac{1}{|\mathcal{T}|} \mathcal{L}(\mathbf{X}[:, \mathcal{T}], \mathbf{X}[:, \mathcal{T} - 1])$ where $\mathcal{T} = \{l, \dots, t\}$ and $\mathcal{L}(\cdot, \cdot)$ is the squared loss.
- Here, $\mathcal{T}_X(\cdot)$ is a TCN that is trained concurrently. Thus this regularization ensures $\mathbf{X}^{(tr)}$ is such that, there exists a TCN, that can predict $\mathbf{X}^{(tr)}$ based on past time-points.
- Thus the global loss being minimized is: $\mathcal{L}(\mathbf{Y}^{(tr)}, \mathbf{F}\mathbf{X}^{(tr)}) + \lambda \mathcal{R}(\mathbf{X}^{(tr)} | \mathcal{T}_X(\cdot))$
- We train the TCN and the factors alternately using mini-batch SGD (more details in paper)



In this figure, we show some of the basis time-series extracted from the traffic dataset, which can be combined linearly to yield individual original time-series. It can be seen that the basis series are highly temporal and can be predicted in the test range using the network $\mathcal{T}_X(\cdot)$

DeepGLO: Deep Global/Local Model

- Our final hybrid model, is a TCN, where the output of the global model is fed as a covariate, along with other covariates.
- Thus, DeepGLO can combine local per time-series properties along with the global predictions from the TCN-MF model.



The base image has been borrowed from [4]

Empirical Results

- We compare DeepGLO with other state of the art time-series models on rolling validation tasks on four public datasets depicted below.

Data	n	t	τ_w	n_w	$\text{std}(\{\mu(y_i)\})$	$\text{std}(\{\text{std}(y_i)\})$
electricity	370	25,968	24	7	1.19e4	7.99e3
traffic	963	10,392	24	7	1.08e-2	1.25e-2
wiki	115,084	747	14	4	4.85e4	1.26e4
PeMSD7(M)	228	11,232	9	160	3.97	4.42

The four datasets are electricity, traffic, wiki and pems. $\text{std}(\{\mu\})$ denotes the standard deviation among the means of all the time series in the data-set and the next column shows the same for St. deviation.

Algorithm	electricity $n = 370$		traffic $n = 963$		wiki $n = 115,084$	
	Normalized	Unnormalized	Normalized	Unnormalized	Normalized	Unnormalized
Proposed	0.1330/0.4530/0.162	0.0820/0.341/0.121	0.160/0.210/0.179	0.1480/0.1680/0.142	0.569/3.335/1.036	0.2370/0.4410/0.395
Local TCN (LeveledInit)	0.1430/0.3560/0.207	0.0920/0.237/0.126	0.1570/0.2010/0.156	0.1690/0.1770/0.169	0.2430/0.5450/0.431	0.2120/0.3160/0.296
Global TCN-MF	0.1440/0.4850/0.174	0.1060/0.2520/0.188	0.3390/0.4150/0.451	0.2260/0.2840/0.247	1.198/467.56	0.4330/1.590/0.509
Local-Only	0.1440/0.4850/0.174	0.1060/0.2520/0.188	0.3390/0.4150/0.451	0.2260/0.2840/0.247	1.198/467.56	0.4330/1.590/0.509
DeepAR	0.0860/0.259/0.141	0.0994/0.818/0.185	0.1400/0.2010/0.114	0.2110/0.3310/0.267	0.429/2.980/0.424	0.9938/120/1.475
TCN (no LeveledInit)	0.1470/0.4760/0.156	0.4230/0.769/0.523	0.2390/0.4250/0.281	0.336/0.322/0.497	0.5110/0.884/0.509	0.5110/0.884/0.509
Prophet	0.1970/0.3930/0.221	0.221/0.586/0.524	0.3130/0.600/0.420	0.3030/0.5590/0.403	0.429/0.6870/0.340	0.6973/51.086
Global-Only	0.1040/0.280/0.151	0.105/0.431/0.183	0.1590/0.226/0.181	0.210/0.322/0.275	0.309/0.847/0.451	0.3200/0.9380/0.503
TRMF (retrained)	0.2190/0.437/0.238	0.368/0.779/0.346	0.4680/0.841/0.581	0.3290/0.6870/0.340	0.6973/51.086	0.639/2.0000/0.893
Algorithm	PeMSD7(M) (MAE/MAPE/RMSE)					
DeepGLO (Unnormalized)	3.53/ 0.079 / 6.49					
DeepGLO (Normalized)	4.53/ 0.103 / 6.91					
STGCN(Cheb)	3.57/ 0.087/ 6.77					
STGCN(1 st)	3.79/ 0.091/ 7.03					

Comparison of algorithms on normalized and unnormalized versions of data-sets on rolling prediction tasks. The error metrics reported are WAPE/MAPE/SMAPE. TRMF is retrained before every prediction window, during the rolling predictions. All other models are trained once on the initial training set and used for further prediction for all the rolling windows. Prophet could not be scaled to the wiki dataset, even though it was parallelized on 32 core machine. Note that for DeepAR normalized setting corresponds to scalar=True and unnormalized setting is scalar=False, in the GluonTS implementation. In the bottom table DeepGLO is compared with the models in [5] on the pems dataset.

References

- Valentin Flunkert, David Salinas, and Jan Gasthaus. Deepar: Probabilistic forecasting with autoregressive recurrent networks. arXiv preprint arXiv:1704.04110, 2017.
- Anastasia Borovykh, Sander Bohte, and Cornelis W Oosterlee. Conditional time series forecasting with convolutional neural networks. arXiv preprint arXiv:1703.04691, 2017.
- Hsiang-Fu Yu, Nikhil Rao, and Inderjit Dhillon. Temporal regularized matrix factorization for highdimensional time series prediction. In Advances in neural information processing systems, pages 847–855, 2016.
- Aaron Van Den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. Wavenet: A generative model for raw audio. CoRR abs/1609.03499, 2016.
- Yu, Bing, Haoteng Yin, and Zhanxing Zhu. "Spatio-temporal graph convolutional networks: A deep learning framework for traffic forecasting." arXiv preprint arXiv:1709.04875 (2017).