## Qingyi Pan, Wenbo Hu, Ning Chen*

### Dept. of Comp. Sci. & Tech., Institute for AI, BNRist Lab, THBI Lab, Tsinghua University

Tsinghua University

IJCAI 2021 — MONTREAL

## Introduction

❑ It is important yet challenging to perform **accurate** and **interpretable** time series forecasting. Traditional parametric model are easy-to-interpret, but their predictive capabilities are limited. Deep architectures can boost the forecasting accuracy, they sacrifice interpretability. It is relatively unexplored to develop both accurate and interpretable methods for **multivariate time series forecasting.**

❑ **Existing work:**

1. Interpretation methods for general neural networks [Ribeiro et al., 2016; Shrikumar et al., 2017; Lundberg & Lee 2017]:

   Use gradient information to extract feature information for after the back-propagation

2. Transfer attention methods from the fields of language or vision [Bahadanau et al., 2014, Shih et al., 2019]:

   Attention values are calculated via the relative importance of the different time steps.

❑ **Key: Considering the time and feature dimensions in coherent manner**

- Represent the multivariate time series as a set of *window × feature* 2D series images.
- Each series image corresponds to a part of the multivariate time series within a given time window.
- Each row corresponds to one feature dimension.
- Follow the perturbation strategy in the smallest destroying region (SDR) principle [Dabkowski & Gal, 2017] for the reference series image.

$$\hat{x}_{t,i} = \begin{cases} x_{t,i} + \epsilon_{\sigma_1} & \text{noise} \\ g_{\sigma_2}(x_{t,i}) & \text{blur} \end{cases}$$

- Each series can be composed of three parts: trend, seasonality and residual [Hyandman & Athanasopoulos, 2018]
- Adding Gaussian blur extracts the trend information, and adding some noise enhances local information.

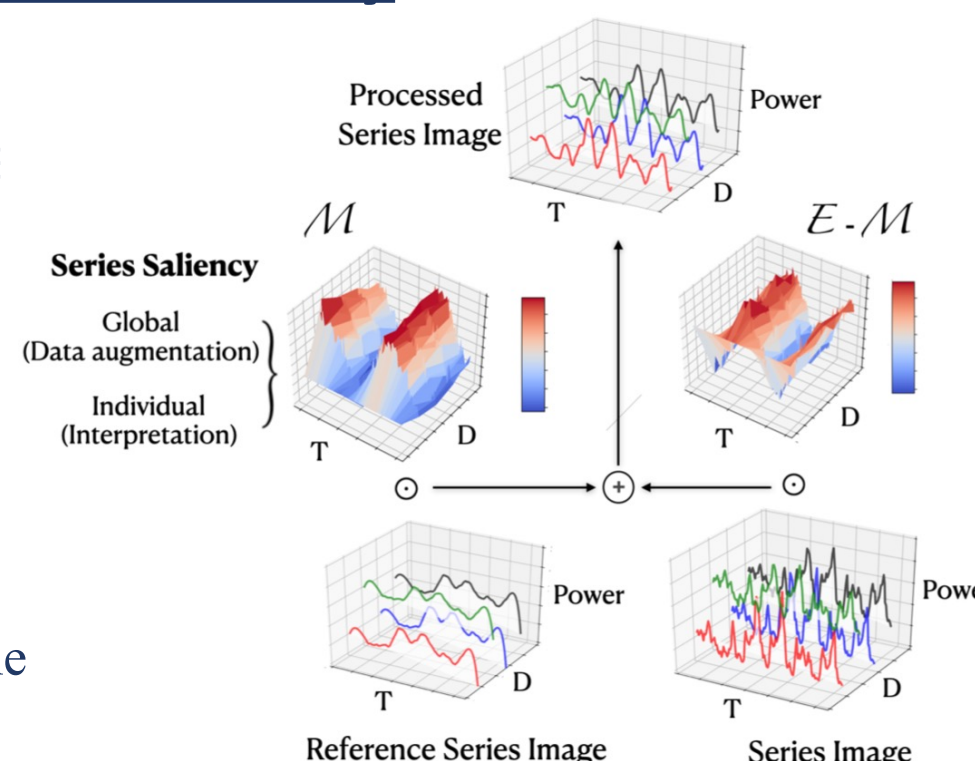❑ **Motivation: Sensitivity over perturbation for DNNs in time domain**

- When the amount of injected noise or blurring is small, the reference series image can be treated as data augmentation in time domain for deep models.
- If the perturbation is not set properly (e.g., too large), the blurring will introduce irregular roughness to cover the original series, making it difficult for DNNs to learn temporal patterns.

## Series Saliency

■ Series saliency module introduces a learnable mask $M \in [0,1]^{D \times T}$ and selectively combines the reference series images and original one:

$$\tilde{X} = M \odot \hat{X} + (E - M) \odot X,$$

■ Series saliency can generate data that cover the unexplored input space while maintaining the important characteristics of the original series image.

Mixup mechanism

## Methodology

❑ **Dual-path architecture**

- Scale of input data often changes in a non-periodic manner

$$\hat{y} = \underbrace{y^{(o)}}_{\text{Linear}} + \underbrace{y^{(r)}}_{\text{Non-Linear}}$$

❑ **Training with Series Saliency:**

$$\min_{\phi,M} \ell_1(\phi,M) + \lambda_1 \ell_m(M) + \lambda_2 \ell_r(M)$$

where $\lambda_1$ and $\lambda_2$ are the coefficient.

❑ **Training loss $\ell_1$:**

- $\ell(\phi,M) = \|f(\tilde{X}_i) - y_i\|^2 + \|f(X_i) - y_i\|^2$
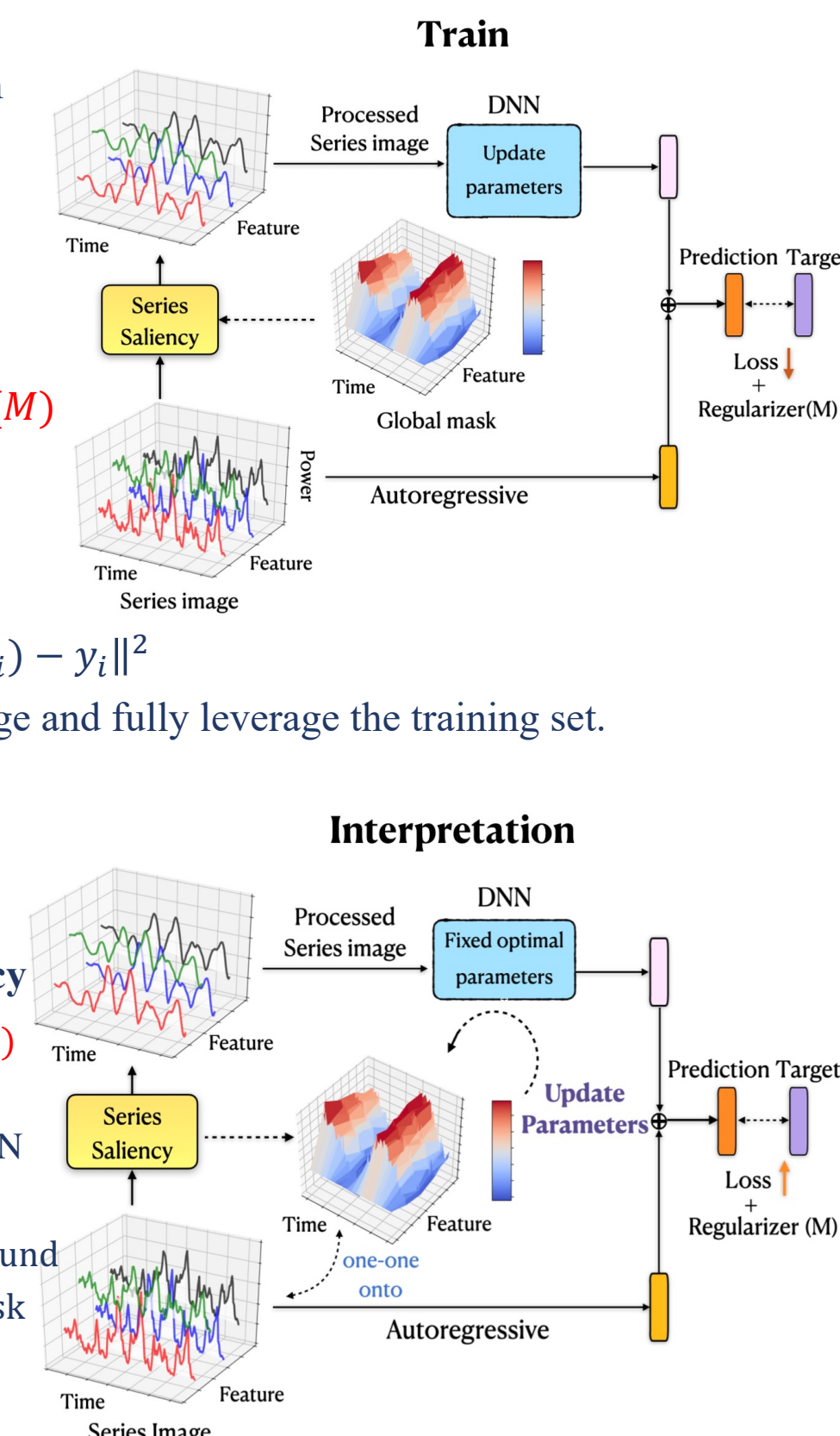- Generate the augmented series image and fully leverage the training set.

❑ **Regularization term:**

- $\ell_m(M) = \|M\|_2$
- $\ell_r(M) = \|MM^T - I\|_F$

❑ **Interpretation with Series Saliency**

$$\min_M -\|\hat{y}^* - y^*\| + \lambda_1 \ell_m(M) + \lambda_2 \ell_r(M)$$

- Fixed optimal parameters $\phi^*$ of the DNN and AR models after training.
- The most salient feature region are found by identifying the representative mask
- For AR part, the weights are easy-to-interpret because of linearity.
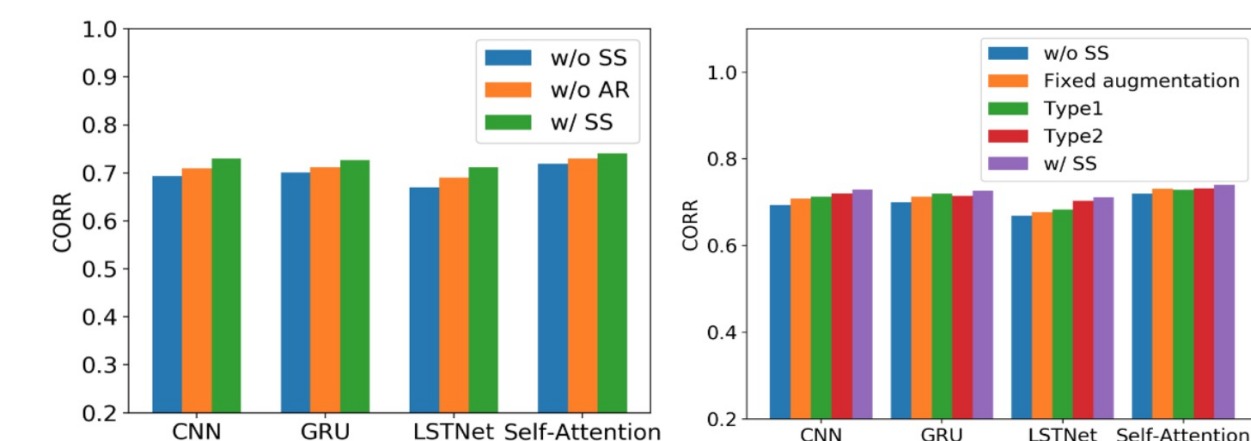
## Experiments

❑ **Results on Forecasting**

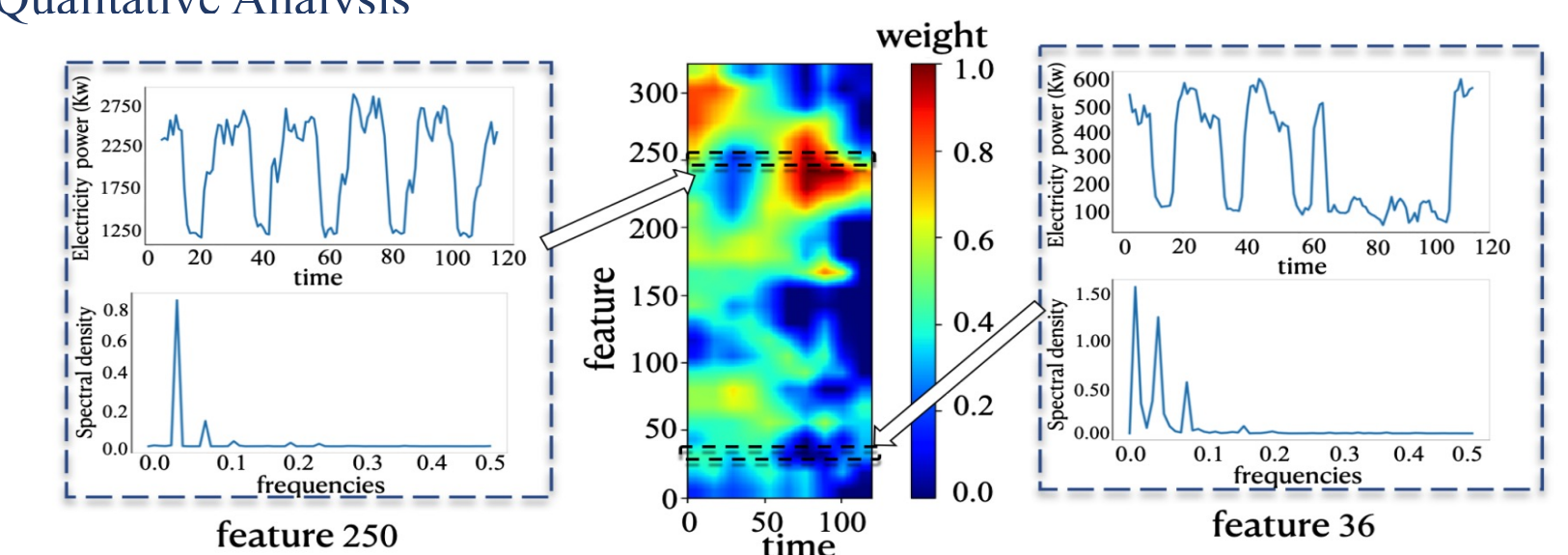➤ The three datasets are representative (from difficult to easy) and dimension from 13 to 321.

| Methods | Air Quality $\tau=3$ | $\tau=6$ | $\tau=12$ | Industry $\tau=3$ | $\tau=6$ | $\tau=12$ | Electricity $\tau=3$ | $\tau=6$ | $\tau=12$ |
|---|---|---|---|---|---|---|---|---|---|
| CNN | $0.775 \pm 0.003$ | $0.701 \pm 0.001$ | $0.636 \pm 0.001$ | $0.919 \pm 0.022$ | $0.909 \pm 0.019$ | $0.841 \pm 0.008$ | $0.883 \pm 0.004$ | $0.871 \pm 0.002$ | $0.866 \pm 0.004$ |
| GRU | $0.804 \pm 0.003$ | $0.712 \pm 0.002$ | $0.639 \pm 0.003$ | $0.953 \pm 0.003$ | $0.936 \pm 0.013$ | $0.904 \pm 0.011$ | $0.878 \pm 0.001$ | $0.877 \pm 0.003$ | $0.867 \pm 0.002$ |
| LSTNet | $0.777 \pm 0.001$ | $0.708 \pm 0.004$ | $0.624 \pm 0.004$ | $0.949 \pm 0.004$ | $0.934 \pm 0.003$ | $0.876 \pm 0.011$ | $0.922 \pm 0.004$ | $0.913 \pm 0.002$ | $0.906 \pm 0.002$ |
| SA | $0.813 \pm 0.002$ | $0.722 \pm 0.003$ | $0.643 \pm 0.003$ | $0.961 \pm 0.002$ | $0.942 \pm 0.005$ | $0.905 \pm 0.009$ | $0.919 \pm 0.007$ | $0.907 \pm 0.001$ | $0.902 \pm 0.003$ |
| CNN w/ SS | $0.779 \pm 0.005$ | $0.723 \pm 0.009$ | $0.641 \pm 0.007$ | $0.941 \pm 0.006$ | $0.927 \pm 0.004$ | $0.881 \pm 0.001$ | $0.898 \pm 0.004$ | $0.893 \pm 0.002$ | $0.892 \pm 0.007$ |
| GRU w/ SS | $0.809 \pm 0.003$ | $0.716 \pm 0.012$ | $0.649 \pm 0.003$ | $0.955 \pm 0.001$ | $0.935 \pm 0.002$ | $0.912 \pm 0.003$ | $0.905 \pm 0.004$ | $0.889 \pm 0.008$ | $0.878 \pm 0.003$ |
| LSTNet w/ SS | $0.794 \pm 0.008$ | $0.724 \pm 0.002$ | $0.641 \pm 0.003$ | $0.959 \pm 0.004$ | $0.938 \pm 0.001$ | $0.901 \pm 0.002$ | $\mathbf{0.928 \pm 0.003}$ | $\mathbf{0.918 \pm 0.003}$ | $0.907 \pm 0.001$ |
| SA w/ SS | $\mathbf{0.819 \pm 0.003}$ | $\mathbf{0.732 \pm 0.009}$ | $\mathbf{0.658 \pm 0.001}$ | $\mathbf{0.965 \pm 0.003}$ | $\mathbf{0.955 \pm 0.016}$ | $\mathbf{0.916 \pm 0.004}$ | $0.923 \pm 0.003$ | $0.915 \pm 0.001$ | $\mathbf{0.911 \pm 0.002}$ |

➤ Ablation study

❑ **Results on Interpretation**

➤ Qualitative Analysis

➤ Quantitative Comparison

| Methods | Industry | Air Quality | Electricity |
|---|---|---|---|
| Grad | $0.214 \pm 0.007$ | $0.297 \pm 0.007$ | $0.199 \pm 0.008$ |
| DeepLift | $0.211 \pm 0.008$ | $0.241 \pm 0.006$ | $0.174 \pm 0.003$ |
| Ablation | $0.204 \pm 0.008$ | $0.225 \pm 0.006$ | $0.213 \pm 0.009$ |
| Occlusion | $0.124 \pm 0.004$ | $0.221 \pm 0.011$ | $0.142 \pm 0.005$ |
| Shapley | $0.145 \pm 0.006$ | $0.211 \pm 0.010$ | $0.171 \pm 0.005$ |
| Attention | $0.141 \pm 0.007$ | $0.203 \pm 0.007$ | $0.139 \pm 0.003$ |
| Type1 | $0.131 \pm 0.005$ | $0.205 \pm 0.008$ | $0.143 \pm 0.004$ |
| Type2 | $0.122 \pm 0.002$ | $0.201 \pm 0.005$ | $0.135 \pm 0.001$ |
| w/ SS | $\mathbf{0.117 \pm 0.002}$ | $\mathbf{0.192 \pm 0.003}$ | $\mathbf{0.131 \pm 0.002}$ |

## Conclusion

❑ We propose Series Saliency to boost both **accuracy** and **interpretability** for multivariate time series **forecasting.**

❑ Series saliency module acts as an **adaptive data augmentation** method for training deep models while can be optimized for **interpretable** forecasting in both feature and time dimensions.