

# An Image Inpainting Approach to Short-Term Load Forecasting

Yanzhu Liu , Shreya Dutta, Adams Wai Kin Kong , and Chai Kiat Yeo 

**Abstract**—In current power systems, electrical energy is generated whenever there is a demand for it. Therefore, load forecasting, which estimates the active load in advance, is imperative for power system planning and operations. Based on the time horizon, load forecasting is classified as very short-term (below one day), short-term (a day to two weeks), medium-term (two weeks to three years) and long-term (over three years). This paper focuses on the short-term forecasting. The complex multi-level seasonality of load series (e.g., the load in a given hour is not only dependent on load in the previous hour, but also on the previous days load in the same hour, and on the previous week's load in the same day-of-the-week and hour) makes this task challenging, especially when the load data is represented in 1d numerical series. However, in multi-channel images, the patterns in spatial neighbourhood of one channel and the patterns in the neighbourhood along the channel dimension are able to be captured by 3d image processing operations. Hence, this study proposes to transform electrical load data from 1d series to 3d images and transform the problem from future series forecasting to missing patch inpainting. Furthermore, it proposes a recurrent neural network to model the temporal trends in the series by convolutional operations on the spatial neighbourhood in the images. Experimental results demonstrate the effectiveness of the proposed method on two benchmarks and show the capability of inferring the future load from related series if there is a lack of history.

**Index Terms**—Deep learning, short-term load forecasting, univariate time-series.

## I. INTRODUCTION

**E**LECTRIC load forecasting plays a crucial role in power system planning and operations [1]. In terms of forecasting horizon, it can be classified into four categories [2]: very short-term (seconds to several hours), short-term (a day to two weeks), medium-term (weeks up to three years) and long-term (over three years). Very short-term load forecasting is used for power system transients and dynamics; medium-term load

forecasting helps scheduling of fuel supplies and maintenance program and long-term forecasting is necessary for strategic planning of power networks [3]. Short-term load forecasting is very important for power system planning and operations, it provides the basis for power system planning to avoid over and under utilization of generating capacity and guides the power system operations such as automatic generation control and security assessment [4]–[7]. Moreover, short-term electricity forecasting helps in optimizing energy usage in urban buildings, realizing energy-saving operations and improving deployment strategies to avoid electricity outages during peak times [8].

In terms of forecasting object, electric load forecasting can be classified either based on the load levels of power system, such as system load, distribution feeder load, and building & residential load, or based on the area's population and living standard of people, such as domestic load, commercial load, industrial load, and agriculture load [9]. However, recent rapid development of AI research provides load type-agnostic solutions. AI models are designed to automatically learn the underlying patterns of load series from historical data, which highly depends on the data availability. Before enough historical load records are accumulated, most data-driven forecasting models are not applicable. To make use of other load data with similar characteristics as training samples, it requires models with powerful generalization capability. Generalization refers to extracting underlying information that is shared among multiple related load series and is suitable to apply on the target series.

Besides historical data, electric load is subject to various exogenous factors, such as weather conditions, calendar effect and types of consumers. *Multivariate models* in load forecasting integrate these variables for prediction. *Univariate models* are those that do not rely on any exogenous variables. It has been reported that for short-term forecasting, the accuracies of multivariate models are generally better than univariate models [2]. However, it cannot be guaranteed that exogenous data is always available and reliable. The proposed model in this paper is a univariate model that does not make use of any information beyond historical series.

The challenge of short-term load forecasting arises from the complex multi-level seasonality of load series. For example, the load in a given hour is not only dependent on the load in the previous hour, but also on the previous day's load in the same hour, and on the previous week's load in the same day-of-the-week and hour. From the data perspective, the only input of univariate load forecasting model is a 1d series of numbers. The models have to capture the multi-level seasonality

Manuscript received 6 June 2021; revised 20 October 2021 and 12 February 2022; accepted 26 February 2022. Date of publication 15 March 2022; date of current version 22 December 2022. Paper no. TPWRS-00899-2021. (Corresponding author: Yanzhu Liu.)

Yanzhu Liu is with the Nanyang Technological University, Singapore 639798. She is now with Institute for Infocomm Research, A-Star, Singapore 138632 (e-mail: liu\_yanzhu@i2r.a-star.edu.sg).

Shreya Dutta is with the Energy Research Institute @ Nanyang Technological University (ERI@N), Singapore 639798 (e-mail: shreya@ntu.edu.sg).

Adams Wai Kin Kong and Chai Kiat Yeo are with the School of Computer Science and Engineering, Nanyang Technological University, Singapore 639798 (e-mail: adamskong@ntu.edu.sg; ASCKYEO@ntu.edu.sg).

Color versions of one or more figures in this article are available at <https://doi.org/10.1109/TPWRS.2022.3159493>.

Digital Object Identifier 10.1109/TPWRS.2022.3159493

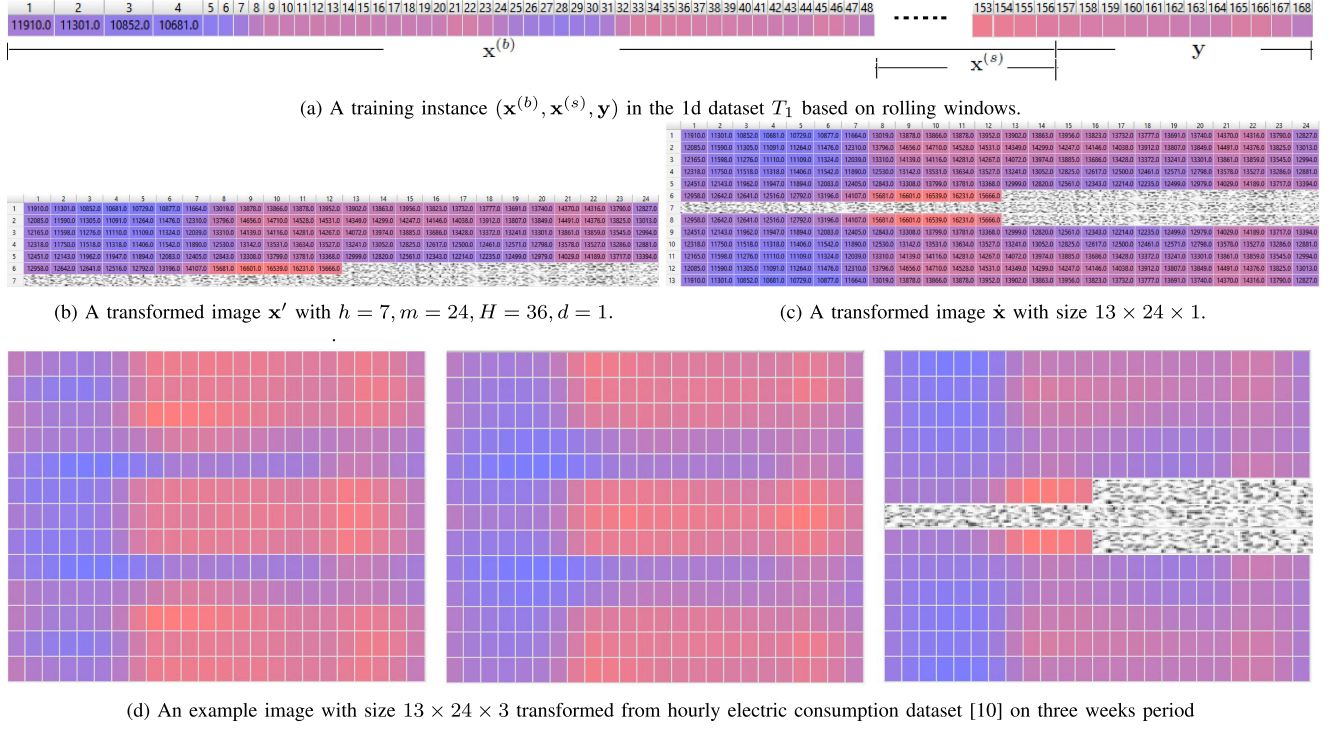


Fig. 1. Examples of the transformed dataset set.

from long look-back horizons along the 1d series. However, if the series numbers are properly mapped to pixels of an image, a week-long half-hourly series (i.e.,  $7 \times 48$  time stamps) corresponds to a tiny patch with  $7 \times 48$  pixels. Therefore, longer pattern of series is able to be explored in local regions of images. And under proper re-organization (e.g., Fig. 1), the multi-level seasonality can be encoded in spatial patterns, which motivates the authors to use a computer vision approach to solve the short-term load forecasting problem.

In computer vision, image inpainting is a task to fill missing pixels in a damaged image. If the series are transformed into images and a section of the future values to be predicted is transformed as an unobserved patch, estimating future values in the load series, in fact, becomes similar to generating pixels in the missing area of the images. To adapt image inpainting models for load forecasting, this paper firstly encodes a series including historical energy meter values and a length of unknown future values to a multi-channel image, each channel of which is built from series values of one week period and the values to be predicted are kept as a missing patch in the image. Based on that, a recurrent neural network (RNN) integrated with 3D convolutions in the memory cells is proposed to generate the missing patches. It takes a  $\tau$ -channel image with an unobserved patch in the last channel as input and produces a 1-channel image with the filled patch as output. The recurrent network models the sequential property between channels of images. The 3d convolution operations encapsulated in LSTM cell capture shift-invariant local patterns in spatial neighbourhoods of both intra and inter channels. Because the proposed network is learnt from instances converted from segments of load series, the

training set can be constructed by segments of multiple historical series. For a newly installed sensor or a new building with little history, multiple load series of sensors located in nearby regions can be used as historical data to train the proposed network. If the few historical values of the target sensor are able to construct a testing instance, the forecasting can be queried from the network by inputting the testing instance. Due to similar visual characteristics can be shared between different images, a model capturing the visual patterns of an image can be applied to related images. This property provides the basis that the information learnt by the proposed model from related load series can be applied to the others.

The rest of this paper is organized as follows. Section II reviews the literature of load forecasting. Section III describes the proposed method and its network architecture. Section IV reports the experimental results, and section V concludes the paper.

## II. RELATED WORK

Numerous approaches and their variants are developed for short-term load forecasting [2]. Many of them are based on the following classical techniques. 1) *Regression analysis* estimates the relationship between load as the dependent variable and independent variables such as weather conditions and events. The prominent approaches include multiple linear regression [11], semi-parametric additive regression [12], and fuzzy regression [13]. 2) *Statistical time series* techniques, represented by exponential smoothing [14] and ARIMA (Auto-regressive Integrated Moving Average) [15], are univariate models which do

not rely on any exogenous variables. ARIMA models can be generalized to include explanatory factors. 3) *Machine learning* techniques, represented by SVR (support vector regression) [16] and ANNs (artificial neural networks) [5], learn the nonlinear mapping from load and exogenous variables to future load values in the supervised fashion. Generally speaking, regression analysis and machine learning based models are multivariate models, while statistical time series models are univariate. Because electric load is highly dependent on weather conditions, multivariate models incorporating the explanatory variables have potential advantages in terms of accuracy over univariate models [2], [17]. In contrast, due to no additional information requirements, univariate models can be widely and fairly applied on more datasets with only historical load series. It is worth mentioning that Hong *et al.* [2] provided a comprehensive review on them. Readers can refer to it for more details.

In recent years, deep learning technique has achieved great success in wide areas. Several deep neural network approaches have also been developed for short-term load forecasting. Singh *et al.* [18] provided a very recent review. From the network architecture perspective, three predominant networks for load forecasting are Stacked Auto-encoder [19], Recurrent Neural Network (RNN) [20] and Convolution Neural Network (CNN) [21]. From the learning strategy perspective, two different supervised manners are applied: end-to-end training, which directly learns from raw load series and outputs the predicted series in the future horizon by the neural networks [20], [22]; and hybrid training, where neural layers work as automatic feature extractor capturing the nonlinear abstract features from load series linked with another forecaster (e.g., SVR, fuzzy regressor or another type of neural network) to produce the estimated series based on the extracted features [23], [24]. From the data format perspective, the input of the neural networks can be 1d univariate load series, 2d transformed image-like grids, and 2d original multivariate series which can be presented as a matrix in which rows are for the time dimension and columns are for the load and exogenous variables. The output of the neural networks is either 1d future load series, or one real number for the load value in the immediate next time stamp mainly used in multi-step forecasting.

Generally speaking, CNNs consist of convolutional layers performing the convolutional operation on input data and fully-connected layers computing the final output. Because CNNs are powerful to capture local correlation and extract abstract features, but not specially designed for sequential data, most of the CNNs based load forecasting approaches only keep convolutional layers for feature extraction and replace fully-connected layers by other forecasters. As CNNs are developed for grid data, if univariate load series data are viewed as 1d grids, CNNs can be used for short-term load forecasting with 1d convolution [25]. Plenty of CNNs based approaches have been proposed to transform load series into images and perform 2d convolutional operation on them. The outputs of convolutional layers are the extracted features, which are passed into a predictor such as multilayer neural network, K-means cluster and so on [24], [26], [27] to get the final estimation. As recurrent networks are designed for processing sequential data, LSTM has

been used as the forecaster integrated with the 2d convolutional layers [28]–[31].

The proposed neural network in this paper is a LSTM equipped with 3d convolution, but it is very different from the convolution and RNN hybrid approaches in the literature. It takes  $\tau$ -channel images ( $\tau > 3$ ) as inputs and produces 1-channel images as outputs, which is trained end-to-end. The recurrent component models the sequential property among the channels of images, not among the series values at time stamps as traditional RNN based forecasting models. The Eidetic-3d [32] cell is used as the LSTM cell. The output of an Eidetic-3d cell at a certain layer for a time is computed from cell input, hidden state, cell memory and global memory. All the four inputs are filtered with 3d convolutional operations before being passed into the gates. Therefore, 3d convolution filters are encapsulated inside the LSTM cell, but not the separate convolutional layers installed on top of the whole RNN as existing load forecasting approaches.

The proposed approach is motivated by image inpainting idea. In computer vision, image inpainting or image completion aims to fill missing region with visually realistic and semantically reasonable content. With power of rich representation capability and being able to automatically extract abstract knowledge, deep learning techniques reduce barriers of interdisciplinary applications. Under reconstruction of electric load series to images with unknown future series as a missing patch, inpainting models are possible to provide accurate performance for load forecasting. In literature, numerous deep inpainting models have been proposed [33]–[36]. An extensive review about them is offered by Zhen *et al.* in [37]. Deep neural networks (DNNs) based approaches are the state-of-the-art [38], [39] for image inpainting. Yeh *et al.* [40] proposed a constrained semantic image inpainting model to generate the missing content by conditioning on training images. Yu *et al.* [41] introduced a contextual attention mechanism into a generative adversarial network to explicitly utilize known surrounding background information as references to generate missing patches. In [42], they extended inpainting model to free shape of missing regions and made use of guidance learnt from millions of unlabeled images. However, to the best knowledge of the authors, there is no existing work using such techniques on time series forecasting, because they are very different tasks on the different data types. However, deep learning, which is able to automatically extract abstract features provides an opportunity to apply the techniques across domains.

### III. IMAGE BASED ENERGY CONSUMPTION SERIES FORECASTING

For the sake of clear presentation, the problem setting and notations are given first. Given a length- $S$  observed series history  $[y_1, y_2, \dots, y_s]$  and a forecast horizon  $H$  ( $H \ll S$ ), the task is to predict the vector of future values  $[y_{s+1}, y_{s+2}, \dots, y_{s+H}]$ , where  $y_i \in \mathbb{R}$  is the energy consumption at the  $i$ -th time step of the series. This study focuses on higher frequency series than daily, e.g., in hour and half hour. The term “image” in this paper refers to as a multi-channel matrix in  $\mathbb{R}^{h \times w \times d}$ , where



$h$ ,  $w$  are height and width, and  $d$  is the number of channels and not necessarily equals to 3 like a RGB color image. The term “frame” refers to a slice of an image along the channel dimension, i.e., it is a matrix in  $\mathbb{R}^{h \times w}$ .

### A. Time Series Encoding

Following the strategy of deep learning methods in traditional time series forecasting, the training set is constructed based on rolling windows. To encode both short-term and long-term temporal information, an instance is built as a pair of sub-series determined by a big and a small window respectively. Let the big window size  $L_b = 7 \times m \times d - H$ , where  $m$  is the frequency per day of the series, (e.g.,  $m = 24$  in hourly series and  $m = 48$  in half-hourly series) and  $L_b < S$ . Thus, a sub-series with size  $L_b$  consists of contiguous time stamp values within  $d$  weeks. Let the small window size  $L_s = h \times m \times d' - H$ , where  $h = \lceil \frac{H}{m} \rceil$  and  $\lceil \cdot \rceil$  is the ceiling function. Thus, a sub-series with size  $L_s$  consists of contiguous time stamp values within  $h \times d'$  days. Denote  $T_1 = \{(X_b, X_s, Y)\}$  as the training set with  $n$  labeled instances:  $X_b = \{\mathbf{x}_i^{(b)} | \mathbf{x}_i^{(b)} = [y_{i-(L_b-1)}, y_{i-(L_b-2)}, \dots, y_{i-1}, y_i], L_b < i \leq S\}$ ,  $X_s = \{\mathbf{x}_i^{(s)} | \mathbf{x}_i^{(s)} = [y_{i-(L_s-1)}, y_{i-(L_s-2)}, \dots, y_{i-1}, y_i], L_s < i \leq S\}$ , and  $Y = \{\mathbf{y}_i | \mathbf{y}_i = [y_{i+1}, y_{i+2}, \dots, y_{i+H}], i + H < S\}$ . Therefore, for the  $i$ -th instance in the training set,  $(\mathbf{x}_i^{(b)}, \mathbf{x}_i^{(s)})$  is the pair of input series and  $\mathbf{y}_i$  is its label series. To simplify the notations, when referring to any input series, we omit the subscript  $i$  for  $\mathbf{x}_i^{(b)}$ ,  $\mathbf{x}_i^{(s)}$  and  $\mathbf{y}_i$ , and rewrite  $\mathbf{x}^{(b)} = [x_1, x_2, \dots, x_{L_b}]$ ,  $\mathbf{x}^{(s)} = [x_1, x_2, \dots, x_{L_s}]$  and  $\mathbf{y} = [y_1, y_2, \dots, y_H]$ . The target of the problem is to predict  $\mathbf{y}^*$  for a testing input  $(\mathbf{x}^{*(b)}, \mathbf{x}^{*(s)})$ . Fig. 1(a) shows an example of an instance  $(\mathbf{x}^{(b)}, \mathbf{x}^{(s)}, \mathbf{y})$  in this traditional 1d time series dataset  $T_1$ , where  $m = 24$ ,  $d = 1$ ,  $d' = 6$  and  $H = 12$ . The numbers in the rectangles are energy values, and the brighter color indicates higher consumption in that time stamp.

Each instance  $(\mathbf{x}^{(b)}, \mathbf{x}^{(s)}, \mathbf{y})$  in  $T_1$  will be encoded into two images based on the same strategy. Let  $\mathbf{x}$  denote the input series with the window size  $L = h \times m \times d$ , (i.e.,  $\mathbf{x} = \mathbf{x}^{(b)}$  if  $h = 7$  and  $\mathbf{x} = \mathbf{x}^{(s)}$  if  $h = \lceil \frac{H}{m} \rceil$ ). To keep the unknown values in the forecasting horizon as a patch in the transformed image, a length- $H$  random noise series  $\mathbf{z} = [z_1, z_2, \dots, z_H]$  is appended to each input series  $\mathbf{x}$ , and (1) is used to produce the input image  $\mathbf{x}' \in \mathbb{R}^{h \times m \times d}$ :

$$\mathbf{x}'_{p,q,r} = \begin{cases} x_{(r-1) \times m \times h + (p-1) \times m + q} & \text{if } (r-1) \times m \times h + (p-1) \times m + q \leq L \\ z_{(r-1) \times m \times h + (p-1) \times m + q - L} & \text{otherwise} \end{cases} \quad (1)$$

where  $p \in \{1, 2, \dots, h\}$  and  $q \in \{1, 2, \dots, m\}$  are the indexes of rows and columns respectively and  $r \in \{1, 2, \dots, d\}$  is the index of channels. Denote  $I_{\mathbf{x}} = \{(p, q, r) | (r-1) \times m \times h + (p-1) \times m + q \leq L\}$  as the index set of the region with  $\mathbf{x}$  and  $I_{\mathbf{z}} = \{(p, q, r) | (r-1) \times m \times h + (p-1) \times m + q > L\}$  as the index set of the region with  $\mathbf{z}$ . According to (1), the input series  $\mathbf{x}$  with appended  $\mathbf{z}$  is arranged in left-to-right then top-to-bottom order as  $d$  matrices of  $h \times m$  size.  $m$  is the number of time slots per day and therefore, e.g. for  $\mathbf{x}^{(b)}$ , such  $h \times m$  ( $h =$

7) matrix consists of contiguous time series values within one week. These matrices are concatenated together as a  $d$ -channel image  $\mathbf{x}'$  and  $\mathbf{z}$  is the last  $H$  pixels in the last channel if  $H < h \times m$ , which is the unobserved patch to be predicted. If  $d = 1$ , there is only one frame in the image, and Fig. 1(b) shows its 2d representation in a transformed image with  $h = 7$ ,  $m = 24$ ,  $H = 36$ .

Based on the above reorganization strategy, a window of  $7 \times m \times d - H$  series is constructed into an image with the size of  $7 \times m$  and  $d$ -channel. The images in the new training set are very tiny and narrow compared to most image application datasets where CNN is applicable. Moreover, the targets we aim to predict are the values in  $I_{\mathbf{z}}$  region, which are located at the boundary of the images and hence there are no spatial neighbors from future time stamps. Therefore, a simple trick is used to augment the images further by (2):

$$\dot{\mathbf{x}}_{p,q,r} = \begin{cases} \mathbf{x}'_{p,q,r} & p \leq h, r \leq d \\ \mathbf{x}'_{2h-p,q,r} & p > h, r \leq d \end{cases} \quad (2)$$

where  $\dot{\mathbf{x}} \in \mathbb{R}^{2h-1 \times m \times d}$ . According to (2), each frame of the image is flipped down along the bottom boundary as a mirror reflection, as shown in Fig. 1(c).  $\dot{I}_{\mathbf{x}}$  and  $\dot{I}_{\mathbf{z}}$  are the index sets of pixels with original values from  $\mathbf{x}$  and  $\mathbf{z}$  respectively. Note that  $\mathbf{y}$  consists of the corresponding ground truth of pixels in  $\dot{I}_{\mathbf{z}}$ . This constructed image training set is denoted as  $T_2 = \{(\dot{\mathbf{x}}^{(b)}, \dot{\mathbf{x}}^{(s)}, \mathbf{y})\}$ .

As per the traditional rolling window based time series forecasting problem, the testing input is the segment in the last time window, i.e., the testing input  $\mathbf{x}^{*(b)} = [y_{S-L_b+1}, y_{S-L_b+2}, \dots, y_S]$ ,  $\mathbf{x}^{*(s)} = [y_{S-L_s+1}, y_{S-L_s+2}, \dots, y_S]$ , and the testing target  $\mathbf{y}^* = [y_{S+1}, y_{S+2}, \dots, y_{S+H}]$  in the original series. Following the same transform strategy, the testing instance  $(\mathbf{x}^{*(b)}, \mathbf{x}^{*(s)}, \mathbf{y}^*)$  is converted to  $(\dot{\mathbf{x}}^{*(b)}, \dot{\mathbf{x}}^{*(s)}, \mathbf{y}^*)$  by using (1) and (2). Table I lists the notations defined in this paper.

The reason why univariate time series are suitable to be explicitly organized as images in this way is because of the specific characteristics of energy consumption data: it has more obvious seasonality than time series in other domains. In the transformed image  $\dot{\mathbf{x}}^{(b)}$ , pixels in adjacent columns (e.g.,  $\dot{\mathbf{x}}_{p,q-1,r}$  and  $\dot{\mathbf{x}}_{p,q,r}$ ) correspond to contiguous time slots of the same day (e.g., 10am and 11am in day  $q$ ); pixels in adjacent rows (e.g.,  $\dot{\mathbf{x}}_{p-1,q,r}$  and  $\dot{\mathbf{x}}_{p,q,r}$ ) correspond to the same time slots of consecutive days (e.g., 10am in day  $p-1$  and day  $p$ ), and pixels in adjacent channels (e.g.,  $\dot{\mathbf{x}}_{p,q,r-1}$  and  $\dot{\mathbf{x}}_{p,q,r}$ ) correspond to the same time slots of consecutive weeks (e.g., Monday 10am in week  $r-1$  and week  $r$ ). Therefore, the techniques used to capture relationships in the spatial neighborhood can be adapted to extract features for the temporal pattern in the original time series. Fig. 1(d) provides an image example which is transformed from a three-week hourly energy consumption series [10]. The periodic characteristic of energy consumption is visualized in the figure with the brighter color indicating more consumption in the time stamp.

TABLE I  
TABLE OF NOTATIONS

Symbol	Domain	Description
$S$	$\mathbb{R}$	Length of observed historical series
$H$	$\mathbb{R}$	Forecast horizon
$\mathbf{x}^{(b)}$	$\mathbb{R}^{L_b}$	$\mathbf{x}^{(b)} = [x_1, x_2, \dots, x_{L_b}]$ , long input series with length $L_b$
$\mathbf{x}^{(s)}$	$\mathbb{R}^{L_s}$	$\mathbf{x}^{(s)} = [x_1, x_2, \dots, x_{L_s}]$ , short input series with length $L_s$
$\mathbf{y}$	$\mathbb{R}^H$	Ground truth of forecast series (i.e., $\mathbf{y} = [y_1, y_2, \dots, y_H]$ )
$\mathbf{z}$	$\mathbb{R}^H$	Noise series (i.e., $\mathbf{z} = [z_1, z_2, \dots, z_H]$ )
$m$	$\mathbb{R}$	Frequency per day of the target series (e.g., $m = 24$ in hourly series and $m = 48$ in half-hourly series)
$h$	$\mathbb{R}$	Height of the images transformed by Eq.1
$d$	$\mathbb{R}$	The number of channels of the images transformed from $\mathbf{x}^{(b)}$
$d'$	$\mathbb{R}$	The number of channels of the images transformed from $\mathbf{x}^{(s)}$
$\tau$	$\mathbb{R}$	The number of image channels inputted into LSTM cell
$\mathbf{x}'$	$\mathbb{R}^{h \times m \times d}$	Image transformed by Eq.1
$\dot{\mathbf{x}}^{(b)}$	$\mathbb{R}^{2h-1 \times m \times d}$	Image transformed from $\mathbf{x}^{(b)}$ by Eq.2
$\dot{\mathbf{x}}^{(s)}$	$\mathbb{R}^{2h-1 \times m \times d'}$	Image transformed from $\mathbf{x}^{(s)}$ by Eq.2

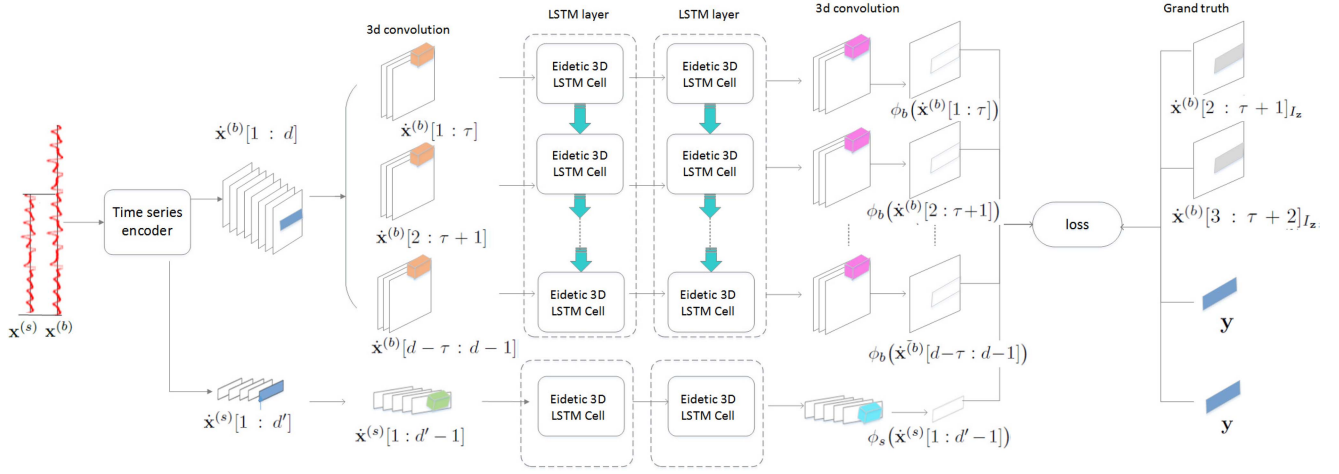


Fig. 2. The architecture of tsRNN. The notations in the figure are defined in Table I.

### B. A Recurrent Neural Network for Patch Generation

Based on the time series encoding in section III-A, the problem setting of the univariate energy consumption can be stated as an image inpainting problem: given an image training set  $T_2 = \{(\dot{\mathbf{x}}^{(b)}, \dot{\mathbf{x}}^{(s)}, \mathbf{y})\}$ , the task is to learn a model to predict  $\mathbf{y}^*$  in the unobserved patch located by  $I_z$  for a testing input  $(\dot{\mathbf{x}}^{*(b)}, \dot{\mathbf{x}}^{*(s)})$ . A recurrent neural network is proposed to generate the unobserved patches. The proposed architecture - recurrent neural network for time series forecasting (tsRNN) is shown in Fig. 2. tsRNN consists of two LSTMs integrated together by the objective function, one of which extracts long-term temporal features on  $\mathbf{x}^{(b)}$ , and the other explores short-term features on  $\mathbf{x}^{(s)}$ . For the sake of convenience, denote the frames of any image  $\mathbf{x} \in \mathbb{R}^{h \times w \times d}$  from channels  $a$  to  $c$  as  $\mathbf{x}[a : c]$ , where  $a, c \in \{1, \dots, d\}$  and  $a < c$ . Thus,  $\mathbf{x}$  is equivalent to  $\mathbf{x}[1 : d]$ . As shown in Fig. 2, time series  $\mathbf{x}^{(b)}, \mathbf{x}^{(s)}$  are transformed into two images  $\dot{\mathbf{x}}^{(b)}[1 : d]$  and  $\dot{\mathbf{x}}^{(s)}[1 : d']$  with noise patches highlighted in blue by the time series encoder described in section III-A first, and then  $\dot{\mathbf{x}}^{(b)}[1 : d]$  is split into  $d - \tau$   $\tau$ -channel images  $\dot{\mathbf{x}}^{(b)}[1 : \tau], \dot{\mathbf{x}}^{(b)}[2 : \tau + 1], \dots, \dot{\mathbf{x}}^{(b)}[d - \tau : d - 1]$ , where  $\tau$  is the predefined temporal depth as per the standard recurrent neural networks. For each  $\tau$ -channel image  $\dot{\mathbf{x}}^{(b)}[1 : \tau], \dots, \dot{\mathbf{x}}^{(b)}[d -$

$\tau : d - 2]$ , the patch region located by  $I_z$  in the frames next to the last channel (i.e.,  $\dot{\mathbf{x}}^{(b)}[2 : \tau + 1]_{I_z}, \dot{\mathbf{x}}^{(b)}[3 : \tau + 2]_{I_z}, \dots$  respectively) is the ground truth. For the last image  $\dot{\mathbf{x}}^{(b)}[d - \tau : d - 1]$ ,  $\mathbf{y}$  is the ground truth. As shown in the last column in Fig. 2, grey regions are ground truth patches and blue regions are the final target regions we want to generate. In tsRNN, only one input image  $\dot{\mathbf{x}}^{(s)}[1 : d' - 1]$  is used for  $\dot{\mathbf{x}}^{(s)}[1 : d']$ .

As discussed in section III-A, due to the periodic property of energy consumption data, dependencies exist between pixels in different frames of these transformed images. Long short-term memory (LSTM) is a powerful tool to capture the time dependencies. Wang *et al.* [32] proposed a new LSTM cell (Eidetic-3 d) for video future frame generation, which is a memory cell interacting with its historical records via a gate-controlled self-attention module. We equip Eidetic-3 d as the basic cell in tsRNN. As shown in Fig. 2, tsRNN consists of two LSTMs. The  $\tau$ -channel images  $\dot{\mathbf{x}}^{(b)}[1 : \tau]$  to  $\dot{\mathbf{x}}^{(b)}[d - \tau : d - 1]$  are input into the top LSTM in order. Each of them is propagated through Eidetic-3 d cell layer by layer, and the cells in the same layer (indicated by dotted rectangles in the figure) share weights. Similarly,  $\dot{\mathbf{x}}^{(s)}[1 : d']$  is fed into the bottom LSTM. A following 3 d convolutional layer is used to output feature

maps as the generated frames  $\phi_b(\dot{\mathbf{x}}^{(b)}[1 : \tau])$ ,  $\phi_b(\dot{\mathbf{x}}^{(b)}[2 : \tau + 1])$ ,  $\dots$ ,  $\phi_b(\dot{\mathbf{x}}^{(b)}[d - \tau : d - 1])$ , and  $\phi_s(\dot{\mathbf{x}}^{(s)}[1 : d' - 1])$ . The generated patches located by the index set  $I_z$  (highlighted by grey rectangles in the figure) are supervised by the loss function in (3):

$$\begin{aligned} \mathcal{L} = & \frac{1}{2} \sum_{(\dot{\mathbf{x}}^{(b)}, \dot{\mathbf{x}}^{(s)}, \mathbf{y}) \in T_2} \left( \|\phi_s(\dot{\mathbf{x}}^{(s)}[1 : d' - 1])_{I_z} - \mathbf{y}\|_2^2 \right. \\ & + \sum_{i=1}^{d-\tau-1} \|\phi_b(\dot{\mathbf{x}}^{(b)}[i : i + \tau - 1])_{I_z} \\ & - \dot{\mathbf{x}}^{(b)}[i + 1 : i + \tau]_{I_z}\|_2^2 \\ & \left. + \|\phi_b(\dot{\mathbf{x}}^{(b)}[d - \tau : d - 1])_{I_z} - \mathbf{y}\|_2^2 \right) \end{aligned} \quad (3)$$

where  $\phi_b(\cdot)$  and  $\phi_s(\cdot)$  are the mapping functions representing the layer operations of the top LSTM and the bottom LSTM in tsRNN respectively, and operator  $[\cdot]_{I_z}$  collects values at each position  $(p, q, r)$  in the index set  $I_z = \{(p, q, r)\}$  in order of increasing  $p, q, r$ .  $\|\cdot\|_2^2$  is the  $L_2$  norm.

The learning of tsRNN is summarized in Algorithm 1. The input is the history series  $[y_1, y_2, \dots, y_S]$  with frequency  $m$  and the horizon  $H$  to be predicted. The training set  $T_2$  is built from this series based on the horizon  $H$ . Given different forecasting horizon  $H$ , different images are constructed correspondingly due to the different  $I_z$  even if the same history series is provided. Fig. 2 shows an example that  $H$  is shorter than one week. However, tsRNN supports longer horizons without any changes, in that case the generated frames and ground truth images are more than one channel.

#### IV. EXPERIMENTAL RESULTS

The proposed tsRNN is comprehensively evaluated on benchmarks of two different types of power load - system load and residential load with different frequencies and different value scales. The first benchmark has two sub datasets of system load, which is used to evaluate the predictive performance of tsRNN compared with the baseline methods for different forecasting horizons. The second benchmark is residential load, which is used to evaluate the generalization capability of tsRNN. Furthermore, since the data-driven univariate load forecasting methods predict future series only based on historical records, which do not rely on any other information related to load types, these two benchmarks with different types of load are also used to evaluate the performances of load type-agnostic models. The benchmarks used in the experiments are summarized as follows:

- PJM system load data [10] is collected from the website of PJM Interconnection LLC, which is a regional transmission organization in the United States operating an electricity transmission system serving 14 states of US. The experiments are performed on two sub datasets - AEP and Dayton, which record the electricity consumption in Michigan and a part of Ohio served by the company AEP and Dayton under PJM respectively. Both consist of one hourly series from 2 Oct 2004 00:00 to 2 Aug 2018 23:00

---

#### Algorithm 1: Pseudo code of learning in tsRNN.

---

- Input:** An energy consumption series  $[y_1, y_2, \dots, y_S]$ , where  $y_i \in \mathbb{R}$ , the series frequency  $m$ , and the forecast horizon  $H$ .
- Output:** Predictions of future values  $\tilde{\mathbf{y}}^* = [\tilde{y}_{S+1}, \tilde{y}_{S+2}, \dots, \tilde{y}_{S+H}]$ .
- 1: Initialize integers  $d, d'$  and  $\tau$  ( $\tau < d$ ). Initialize a noise vector  $\mathbf{z} = [z_1, z_2, \dots, z_H]$ .
  - 2: Build the training set  $T_1 = \{(\mathbf{x}^{(b)}, \mathbf{x}^{(s)}, \mathbf{y})\}$  from  $[y_1, \dots, y_S]$  based on the definition.
  - 3: Transform  $T_1 = \{(\mathbf{x}^{(b)}, \mathbf{x}^{(s)}, \mathbf{y})\}$  to  $T_2 = \{(\dot{\mathbf{x}}^{(b)}, \dot{\mathbf{x}}^{(s)}, \mathbf{y})\}$  based on (1) and (2).
  - 4: **for**  $epoch = 1$  to  $MAX_{epoch}$  **do**
  - 5:   Shuffle  $T_2$  and divide it into mini-batches  $T_s$ .
  - 6:   **for each**  $T_s$  **do**
  - 7:     Forward propagate instances through the three layers of tsRNN as shown in Fig. 2.
  - 8:     Calculate loss by (3).
  - 9:     Backward propagate and update the weights.
  - 10:   **end for**
  - 11: **end for**
  - 12: Construct the testing instance  $(\dot{\mathbf{x}}^{(b)*}, \dot{\mathbf{x}}^{(s)*}, \mathbf{y}^*)$  from  $([y_{S-L_b+1}, y_{S-L_b+2}, \dots, y_S], [y_{S-L_s+1}, y_{S-L_s+2}, \dots, y_S], [y_{S+1}, y_{S+2}, \dots, y_{S+H}])$ .
  - 13: Forward propagate  $(\dot{\mathbf{x}}^{(b)*}, \dot{\mathbf{x}}^{(s)*})$ , and output  $\tilde{\mathbf{y}}^* = \phi_s(\dot{\mathbf{x}}^{(s)*}[1 : d' - 1])_{I_z}$ .
- 

with 121,272 time steps in megawatts (MW). The min and max consumption values in AEP series and Dayton series are (9581.0, 25695.0) and (982.0, 3746.0).

- London residential load data [43] is collected from London smart meter data store, which contains the energy consumption readings of 5,567 households in 112 blocks between November 2011 and February 2014 in kilowatt-hours (kWh). This benchmark consists of 5,567 half-hourly series, and the starting and ending dates are not the same for all series. The min and max consumption values of different series are different, and all the consumption values are in the range of [0 to 10] kWh.

##### A. Comparing to Baseline Methods

To comprehensively evaluate the predictive performance of tsRNN, statistical univariate time series forecasting methods, traditional machine learning methods and deep learning methods are compared on the PJM benchmark. The baseline methods include seven traditional models: Seasonal Auto Regressive Integrated Moving Average (SARIMA, [44]), Exponential Smoothing (ExpSm, [45]), Theta [46], Trigonometric Box-Cox ARMA Trend Seasonal (TBATS, [47]), and Dynamic Harmonic Regression (DHR, [48]), CatBoost [49], Pooled Regression [50]; and six deep neural network methods: multilayer perception (MLP) and long short-term memory (LSTM), DeepAR [51], N-BEATS [52], WaveNet [53] and Transformer [54]. In the first category, CatBoost and Pooled Regression are global methods,

TABLE II  
EVALUATION METRICS

Name	Equations
RMSE (root mean squared error)	$\sqrt{\frac{1}{H} \sum_{i=1}^H (y_i - \tilde{y}_i)^2}$
MAE (mean absolute error)	$\frac{1}{H} \sum_{i=1}^H  y_i - \tilde{y}_i $
MAPE (mean absolute percent error)	$\frac{1}{H} \sum_{i=1}^H \left  \frac{y_i - \tilde{y}_i}{y_i} \right $
sMAPE (symmetric mean absolute percent error)	$\frac{1}{H} \sum_{i=1}^H \frac{2 \times  y_i - \tilde{y}_i }{ y_i  +  \tilde{y}_i }$
ND (normalized deviation)	$\frac{\sum_{i=1}^H  y_i - \tilde{y}_i }{\sum_{i=1}^H  y_i }$

which predict the future series based on all training series samples. However, other four statistical baselines are local methods, which predict the future values of the target series only based on its own historical series. In the deep learning category, WaveNet and Transformer are CNN-based methods; MLP and N-BEATS are fully-connected networks; and LSTM and DeepAR are recurrent networks. Six metrics [55] are used as performance indexes. Table II lists their definitions, where  $\tilde{y}_i$  is the prediction for  $i$ -th time stamp in the forecasting horizon,  $y_i$  is the ground truth and  $H$  is the horizon length.

AEP and Dayton data, each of which contains 121,272 continuous time steps, are split into a history segment for training and a future segment for testing based on the 8:2 ratio (i.e., the first 98,720 time stamps are for training and the remaining 24,192 slots of 144-week steps are for testing). For all deep learning methods - tsRNN, MLP, LSTM, DeepAR, N-BEATS, WaveNet and Transformer - the segment  $[y_1 \cdots y_S]$  ( $S=98,720$ ) is used to construct the training sets based on rolling windows, and the 144-week testing segment  $[y_{S+1} \cdots y_{|Y|}]$  ( $|Y|$  is the total number of time stamps in  $Y$ ) is split into 144 segments  $\{[y_{S+i*168-L_b}, \dots, y_{S+i*168}] | i = 0, \dots, 143\}$  and converted into the testing samples  $\{(\mathbf{x}_i^{(b)*}, \mathbf{x}_i^{(s)*}, \mathbf{y}_i^*) | i = 0, \dots, 143\}$ . All the methods are compared for two forecasting horizons:  $h = 48$  (two days) and  $h = 168$  (one week). For each  $h$ , only one model of tsRNN is trained on  $[y_1 \cdots y_S]$ , which is applied to predict for all 144 testing samples, likewise for all deep learning baselines. However, local statistical methods, including SARIMA, ExpSm, Theta, TBATS, and DHR, are only able to predict the future slots next to the end of training slots. For example, if the forecasting target is  $[y_{S+1} \cdots y_{S+h}]$ , the training segment must end at  $y_S$ . As a result, to obtain the predictions for testing samples, 144 models of local methods are learnt from  $[y_1 \cdots y_S]$ ,  $[y_1 \cdots y_{S+168}]$ ,  $[y_1 \cdots y_{S+2 \times 168}]$ ,  $\dots$ ,  $[y_1 \cdots y_{|Y|-168}]$  respectively. This is a limitation of them compared to deep learning methods.

By using grid search on  $\{0, 1, 2\}$  to set the optimal combination for the six standard parameters of SARIMA

TABLE III  
HYPER-PARAMETERS OF THE DEEP LEARNING MODELS FOR THE PJM BENCHMARK

$h = 48$		$h = 168$
tsRNN		
#layers	2 Eidetic LSTM layers	
#nodes	4 cells with [2,3,3] conv (1 <sup>st</sup> layer) 4 cells with [2,3,3] conv (2 <sup>nd</sup> layer)	
window size	$d = 5, d' = 5, \tau = 2$	
minibatch size	192	
#iterations	10000	
learning rate	0.0001	
MLP		
#layers	2 fully-connected layers + Relu	
#nodes	224 neurons (1 <sup>st</sup> layer) 192 neurons (2 <sup>nd</sup> layer)	
window size	128	256
minibatch size	64	
#iterations	75,000	
learning rate	0.001	
LSTM		
#layers	1 LSTM layer + Relu + Dropout 1 fully-connected layer	
#nodes	128 units(1 <sup>st</sup> layer) 48 neurons(2 <sup>nd</sup> layer)	192 units(1 <sup>st</sup> layer) 168 neurons(2 <sup>nd</sup> layer)
window size	128	256
minibatch size	64	
#iterations	75,000	
learning rate	0.001	

(auto-regressive, integrated, moving average, seasonal auto-regressive, seasonal integrated and seasonal moving average), it is observed that for the first ten models learnt from  $[y_1 \cdots y_S]$ ,  $[y_1 \cdots y_{S+1}]$ ,  $\dots$ ,  $[y_1 \cdots y_{S+10}]$ , the optimal combinations of all models are  $[1, 1, 1, 1, 1]$ . Therefore, in the experiments, the six parameters of SARIMA are set as 1. For standard MLP and LSTM models, five sets of features are extracted. For instance, for the time stamp  $i$ , the features include history consumption values  $[y_{i-L} \cdots y_i]$  within a window  $L$ , consumption differences between every two adjacent time stamps  $[y_{i-L+1} - y_{i-L}, \dots, y_i - y_{i-1}]$ , codes indicating whether the consumption values are increasing (i.e., 1 if  $y_{j+1} > y_j$ ; -1 if  $y_{j+1} < y_j$ ; 0 otherwise.  $j = i - L, \dots, i$ ), one binary code indicating whether the time stamp  $i$  is in a workday, 24 one-hot codes indicating that the time stamp  $i$  is in which hour of the day (e.g., if the time stamp  $i$  belongs to 00:00-00:59, the codes are  $[1 \ 0 \ 0 \ 0 \ 0 \dots]$ ). Settings of other statistical baselines follow the default parameters in Monash Time Series Forecasting Archive [56].

The hyper-parameters used in the experiments including window size, architecture information and the training parameters of MLP, LSTM and tsRNN are listed in Table III. Hyper-parameters of other baselines follow the default settings in Monash Time Series Forecasting Archive [56]. Table IV and V summarize the comparison results on the PJM benchmark. Due to the different value scales of AEP and Dayton, the training and testing segments are normalized by  $\frac{y_i - \min}{\max - \min}$  where  $\min$  and  $\max$  are minimum and maximum values of the training segments, and the prediction results are computed on the normalized data. The numbers in Table IV and Table V are the mean errors on the 144 testing samples and  $\pm$  represents standard deviation. The

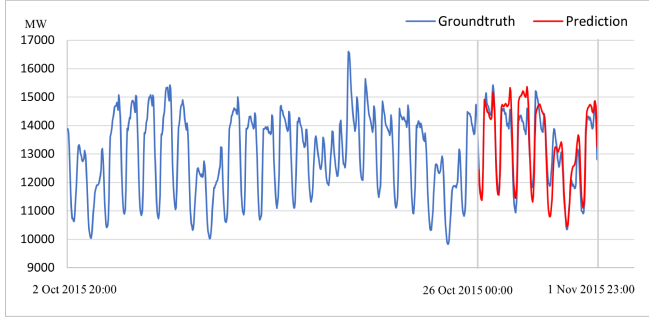


TABLE IV  
RESULTS ON PJM - AEP

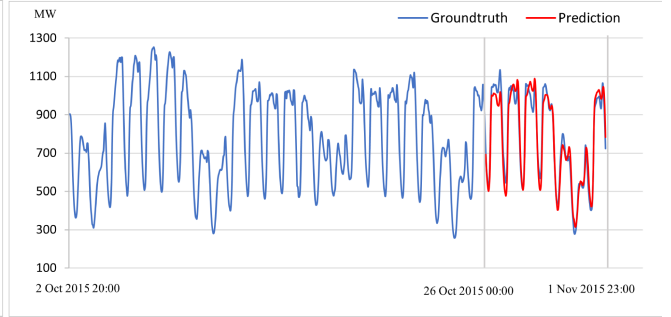
	$h = 48$					$h = 168$				
	RMSE	MAE	MAPE	sMAPE	ND	RMSE	MAE	MAPE	sMAPE	ND
SARIMA	.057±.04	.047±.03	.151±.11	.146±.09	.138±.09	.113±.07	.094±.06	.433±.33	.311±.17	.302±.19
ExpSm	.095±.06	.077±.05	.235±.14	.226±.10	.228±.13	.126±.06	.101±.05	.431±.30	.322±.12	.332±.17
CatBoost	.056±.03	<u>.046±.03</u>	.149±.13	<u>.139±.08</u>	.138±.09	.089±.04	.073±.04	.323±.28	.259±.11	.237±.12
PooledRegression	.057±.04	.048±.03	.152±.13	.146±.09	.141±.09	.090±.04	.074±.04	.333±.30	.259±.11	.239±.12
Theta	.111±.04	.096±.04	.381±.20	.312±.11	.290±.10	.134±.05	.113±.04	.655±.47	.386±.12	.371±.12
TBATS	.096±.05	.080±.05	.272±.17	.263±.19	.240±.14	.194±.21	.161±.17	.771±.95	.491±.35	.528±.56
DHR	.210±.08	.179±.07	.648±.34	.618±.20	.536±.17	.223±.06	.186±.05	1.051±.82	.618±.15	.630±.24
MLP	.122±.08	.108±.08	.375±.31	.322±.20	.327±.25	.127±.06	.107±.06	.437±.28	.367±.16	.340±.16
DeepAR	.086±.05	.069±.05	.196±.12	.220±.14	.197±.10	.092±.04	.076±.04	.380±.35	.262±.12	.256±.15
N-BEATS	<u>.055±.03</u>	<u>.046±.03</u>	<b>.143±.12</b>	.140±.08	<b>.133±.07</b>	<u>.088±.04</u>	<u>.072±.04</u>	<u>.304±.25</u>	<u>.257±.10</u>	<u>.231±.10</u>
Transformer	.073±.04	.061±.04	.204±.17	.180±.10	.188±.12	.138±.06	.120±.06	.417±.16	.503±.19	.371±.12
WaveNet	<b>.053±.03</b>	<b>.044±.03</b>	<u>.145±.12</u>	<b>.134±.07</b>	<u>.134±.08</u>	<u>.088±.04</u>	.073±.03	.363±.32	.257±.12	.251±.14
LSTM	.114±.04	.100±.07	.346±.25	.303±.17	.294±.19	.121±.05	.102±.05	.493±.34	.352±.12	.335±.15
tsRNN-2d-half	.075±.03	.065±.03	.227±.14	.199±.07	.191±.08	.135±.04	.120±.04	.655±.05	.400±.18	.444±.27
tsRNN-2d	.065±.03	.055±.03	.190±.13	.174±.08	.168±.08	.135±.04	.120±.04	.657±.05	.401±.18	.445±.27
tsRNN	.060±.03	.051±.03	.176±.14	.157±.08	.159±.09	<b>.083±.04</b>	<b>.069±.03</b>	<b>.294±.25</b>	<b>.246±.09</b>	<b>.224±.09</b>

TABLE V  
RESULTS ON PJM - DAYTON.

	$h = 48$					$h = 168$				
	RMSE	MAE	MAPE	sMAPE	ND	RMSE	MAE	MAPE	sMAPE	ND
SARIMA	.053±.04	.044±.03	.110±.08	.110±.08	.108±.07	.113±.07	.093±.06	.310±.23	.247±.14	.254±.17
ExpSm	.069±.04	.057±.03	.137±.07	.139±.07	.140±.08	.103±.05	.084±.05	.266±.16	.225±.10	.230±.12
CatBoost	<u>.051±.03</u>	<u>.042±.03</u>	<u>.100±.06</u>	<u>.100±.06</u>	<u>.101±.07</u>	<u>.076±.03</u>	<u>.062±.03</u>	<u>.178±.09</u>	<u>.175±.07</u>	<u>.165±.07</u>
PooledRegression	.053±.03	.044±.03	.105±.07	.106±.07	.106±.07	.080±.04	.066±.04	.198±.12	.183±.09	.177±.09
Theta	.103±.04	.090±.04	.248±.08	.233±.07	.222±.07	.124±.05	.105±.04	.365±.19	.290±.09	.287±.10
TBATS	.094±.05	.079±.05	.208±.12	.204±.13	.194±.11	.168±.15	.140±.13	.464±.47	.394±.36	.392±.38
DHR	.194±.07	.165±.06	.415±.13	.464±.14	.403±.11	.202±.06	.168±.05	.552±.22	.461±.10	.463±.13
MLP	.094±.05	.082±.05	.217±.15	.206±.12	.204±.13	.103±.04	.087±.04	.251±.11	.255±.11	.233±.09
DeepAR	.067±.03	.056±.03	.185±.13	.171±.08	.168±.09	.088±.05	.072±.04	.236±.16	.197±.10	.200±.12
N-BEATS	.055±.03	.046±.03	.152±.13	.141±.07	.137±.08	.076±.04	<u>.062±.03</u>	.184±.12	.176±.08	.166±.08
Transformer	.064±.03	.054±.03	.176±.14	.173±.09	.160±.09	.127±.05	.110±.05	.296±.08	.360±.13	.293±.09
WaveNet	<b>.050±.03</b>	<b>.041±.03</b>	<b>.098±.06</b>	<b>.098±.06</b>	<b>.099±.06</b>	.114±.07	.097±.07	.303±.21	.241±.12	.270±.18
LSTM	.088±.05	.077±.05	.196±.11	.190±.09	.185±.10	.096±.04	.080±.04	.238±.10	.220±.08	.212±.08
tsRNN-2d-half	.062±.03	.052±.02	.136±.06	.130±.05	.129±.05	.100±.03	.086±.03	.264±.12	.237±.08	.242±.10
tsRNN-2d	.058±.03	.049±.02	.127±.06	.122±.05	.122±.06	.099±.03	.086±.03	.263±.12	.236±.08	.242±.10
tsRNN	.056±.03	.047±.03	.123±.07	.117±.06	.119±.07	<b>.071±.03</b>	<b>.058±.03</b>	<b>.176±.09</b>	<b>.163±.07</b>	<b>.159±.07</b>



(a) AEP series segment from 02 Oct 2015 to 01 Nov 2015.



(b) DAYTON series segment from 02 Oct 2015 to 01 Nov 2015.

Fig. 3. Examples of series and predictions of tsRNN.

best performances are highlighted in bold, and the second best results are underlined.

As shown in the tables, for the forecasting horizon  $h = 48$ , WaveNet and N-BEATS perform the best on AEP dataset. For Dayton dataset, WaveNet is the leading model and CatBoost follows behind. However, when horizon  $h = 168$ , the performance of WaveNet drops significantly compared to that on the

shorter horizon. Similar characteristic is observed on local statistical methods (e.g., SARIMA) which are good at short future forecasting, but their performances change a lot with increasing length of the forecasting horizon. The proposed tsRNN achieves the best results for the forecasting horizon  $h = 168$  on both AEP and Dayton datasets, and its performance is affected less by the length of forecasting horizon. Fig. 3 shows examples



TABLE VI  
RESULTS ON THE LONDON BENCHMARK.

	$h = 48$			$h = 168$		
	RMSE	MAE	sMAPE	RMSE	MAE	sMAPE
SARIMA	.126±.07	.089±.05	.806±.27	.126±.08	.085±.07	.835±.29
ExpSm	.124±.08	.087±.06	.844±.36	.131±.23	.093±.21	.900±.37
MLP	.139±.19	.092±.09	.805±.26	.133±.14	.089±.05	.824±.25
LSTM	.139±.19	.093±.09	.809±.27	.132±.14	.088±.05	.819±.25
tsRNN	<b>.094±.05</b>	<b>.063±.04</b>	<b>.711±.27</b>	<b>.102±.05</b>	<b>.067±.03</b>	<b>.718±.24</b>

from AEP and Dayton dataset. The blue curves are the true series data, and red curves are the one week-ahead predictions performed by tsRNN. It can be observed that the main errors occur at the time stamps when the electricity load values change significantly.

### B. Ablation Studies

To evaluate the benefits of the data reconstruction and augmentation, the proposed method is tested on every intermediate transformed data shown in Fig. 1. For the original one-dimensional series data (i.e., Fig. 1(a)) before reconstruction, the proposed idea degenerates to a standard LSTM, whose performance has been reported in Tables IV and V. Because tsRNN is proposed for 3D data in terms of two key aspects: the supervised learning is to infer frame  $d$  based on frames 1 to  $d-1$  of the images; and 3D convolution kernel works on adjacent multiple frames, it is not able to be performed directly on 2D reconstructions. Therefore, we replicate the target frame (as shown in Fig. 1(b) and Fig. 1(c)) to enable the above 3D operations, i.e., in Fig. 2,  $\dot{\mathbf{x}}^{(b)}[1:d]$  is with  $d = 2$  and  $\dot{\mathbf{x}}^{(b)}[1] = \dot{\mathbf{x}}^{(b)}[2]$ . Under this replication, all frames of  $\dot{\mathbf{x}}^{(s)}[1:d]$  with  $\dot{\mathbf{x}}^{(s)}[1] = \dot{\mathbf{x}}^{(s)}[2]$  in Fig. 2 are the noises replacing the target regions. Hence, the part of loss on  $\dot{\mathbf{x}}^{(s)}$  of (3) is omitted in the experiments for 2D data.

TsRNN-2d-half in Tables IV and V is the tsRNN model performed on the replicated 2D data of Fig. 1(b), and tsRNN-2 d is that on Fig. 1(c). As indicated by the results of all metrics for all four dataset settings in Tables IV and V, tsRNN-2d-half and tsRNN-2 d perform much better than LSTM but worse than tsRNN, which verified the benefits of the data reconstruction. Moreover, comparing to tsRNN-2d-half, the inputs of tsRNN-2 d enjoy spatial augmentation by flipping. For 48-horizon forecasting on both datasets, tsRNN-2 d outperforms tsRNN-2d-half benefiting from the data augmentation. However, for 168-horizon forecasting, tsRNN-2 d and tsRNN-2d-half produce very similar results, but perform much worse than tsRNN, which demonstrates that when models are out of capability to capture the underlying patterns, the spatial augmentation of data is not able to help much.

### C. Generalization Capability of tsRNN

To evaluate the generalization capability of tsRNN, the London benchmark is divided on blocks as training/testing partition (blocks 1-90 for training and blocks 91-112 for testing). The energy consumption of each household in a block is stored as

TABLE VII  
HYPER-PARAMETERS OF THE DEEP LEARNING MODELS FOR THE LONDON BENCHMARK

	$h = 48$	$h = 168$
tsRNN		
#layers	2 Eidetic LSTM layers	
#nodes	4 cells with [2,3,3] conv (1 <sup>st</sup> layer) 4 cells with [2,3,3] conv (2 <sup>nd</sup> layer)	
window size	$d = 5, d' = 5, \tau = 2$	
minibatch size	64	
#iterations	1000	
learning rate	0.0001	
MLP		
#layers	2 fully-connected layers + Relu	
#nodes	128 neurons(1 <sup>st</sup> layer) 96 neurons(2 <sup>nd</sup> layer)	224 neurons(1 <sup>st</sup> layer) 192 neurons(2 <sup>nd</sup> layer)
window size	128	256
minibatch size	48	
#iterations	100,000	
learning rate	0.001	
LSTM		
#layers	1 LSTM layer + Relu + Dropout 1 fully-connected layer	
#nodes	128 units(1 <sup>st</sup> layer) 48 neurons(2 <sup>nd</sup> layer)	192 units(1 <sup>st</sup> layer) 168 neurons(2 <sup>nd</sup> layer)
window size	128	256
minibatch size	64	
#iterations	85,000	
learning rate	0.001	

one time series, and therefore, in this experiment, the models are learnt on multiple series from the training blocks and applied on the series from the testing blocks. As the starting and ending time of different series are different, we use all slots from all training series to construct the transformed images, and use the last horizon ( $h = 48, 168$ ) of each testing series as the forecasting target. The aims of this experiment are to evaluate the capability of a method to model the common patterns of energy consumption of households in these blocks based on an assumption that these blocks have high similarity in terms of energy consumption.

The number of households in the training blocks (i.e., blocks 1-90) is 4,397, and the average number of training samples of tsRNN to be constructed from one household series is 500 by the encoder of Section III-A. Therefore, the number of training samples is more than 2,100,000. To deal with the large scalability of this dataset, in the training phase of tsRNN, for each epoch, 100 households are randomly selected and those series are converted to images and split as minibatches to train the model. When one epoch is finished, another 100 random households are

chosen. There are abnormal points in the series as the values in a whole day are zeros. In the preprocessing, the days with more than 15 zero values are removed from the training series. When there are zeros in the testing series, and the corresponding MAPE and ND will be infinity for any predictions (see the definitions in Table II). Therefore, the three matrices RMSE, MAE and sMAPE are evaluated in this experiment. Consumption values are normalized by minimum and maximum values of each series separately using  $\frac{y_i - \min}{\max - \min}$ . Table VI shows the results. Benefiting from making use of series from the training blocks, tsRNN outperforms all baselines significantly. It shows that given time series under the assumption of high similarity, tsRNN is able to capture their common characteristics and generalized to other related series. The architecture and hyper-parameters of deep learning models in this experiments are listed in Table VII.

In the same way as the experiments in Section IV-A, separate SARIMA and ExpSm models are learnt for each of the testing series (the number of testing series is 1,038), whereas only one deep learning model is trained from all of the training series. For tsRNN and the other deep learning baselines, following the traditional strategy, the segment before the target horizon in the testing series is used to construct the testing input. As SARIMA and ExpSm are only able to predict the future stamps based on the previous history in the same time series, in each testing series, the exact segment converted to the testing input of deep learning models is used to learn the SARIMA and ExpSm model. Therefore, for different forecasting horizon  $h = 48$  and  $h = 168$ , the segments in a same testing series used to learn SARIMA and ExpSm models are different. For example, given a testing series  $y$ , if  $h = 48$ , the segment used to learn SARIMA and ExpSm is  $[y_{|y|-48-L_b} \cdots y_{|y|-48}]$ ; if  $h = 168$ , the segment will be  $[y_{|y|-168-L_b} \cdots y_{|y|-168}]$ . From the experimental results listed in Table VI, we can observe that the performances of SARIMA and ExpSm for  $h = 48$  and  $h = 168$  are close, which is obviously different from the phenomenon in the experiments of Section IV-A. In those experiments, the predictions for the next 48 or 168 time stamps come from the same SARIMA and ExpSm model learnt from the same segment. Therefore, it can be concluded that the performances of SARIMA and ExpSm are highly affected by the length of training segment, particularly when the length of training segment is not significantly longer than the forecasting horizon.

## V. CONCLUSION

Benefiting from the rich representation capability of images, this paper tackles the long-term prediction and generalization challenges of short-term load forecasting problem by transforming a consumption series into multi-channel images. Based on the transformation, a forecasting problem for univariant load series is converted to an image inpainting problem. This paper proposes an adapted recurrent deep neural network with 3d convolutions to generate the unobserved patch in the transformed images, which is the prediction of future time stamps of the consumption series. Experimental results have shown its effectiveness and generalization capabilities compared to both statistical models and standard deep learning methods.

## ACKNOWLEDGMENT

The authors would like to thank Siemens Smart Infrastructure, Singapore for providing their domain expertise in the field of building energy and Energy Research Institute @ Nanyang Technological University (ERI@N) to conduct this research as a part of the EcoCampus Flagship Programme which focused on campus-wide sustainability.

## REFERENCES

- [1] T. Hong *et al.*, "Energy forecasting: Past, present, and future," *Foresight: Int. J. Appl. Forecasting*, vol. 32, pp. 43–48, 2014.
- [2] T. Hong and S. Fan, "Probabilistic electric load forecasting: A tutorial review," *Int. J. Forecasting*, vol. 32, no. 3, pp. 914–938, 2016.
- [3] N. Phuangpornpitak and W. Prommee, "A study of load demand forecasting models in electric power system operation and planning," *GMSARN Int. J.*, vol. 10, pp. 19–24, 2016.
- [4] W. Xifan and J. R. McDonald, *Modern Power System Planning*. New York, NY, USA: McGraw-Hill, 1994.
- [5] H. S. Hippert, C. E. Pedreira, and R. C. Souza, "Neural networks for short-term load forecasting: A review and evaluation," *IEEE Trans. Power Syst.*, vol. 16, no. 1, pp. 44–55, Feb. 2001.
- [6] K. Metaxiotis, A. Kagiannas, D. Askounis, and J. Psarras, "Artificial intelligence in short term electric load forecasting: A state-of-the-art survey for the researcher," *Energy Convers. Manage.*, vol. 44, no. 9, pp. 1525–1534, 2003.
- [7] T. Hong, *Short Term Electric Load Forecasting*. Chapel Hill, NC, USA: North Carolina State Univ., 2010.
- [8] S. Humeau, T. K. Wijaya, M. Vasirani, and K. Aberer, "Electricity load forecasting for residential customers: Exploiting aggregation and correlation between households," in *Proc. Sustain. Internet ICT Sustainability*, 2013, pp. 1–6.
- [9] Circuit Globe, "Electrical Load," 2019. [Online]. Available: <https://circuitglobe.com/electrical-load.html>
- [10] Rob Mulla, "Hourly Energy Consumption," 2018. [Online]. Available: <https://www.kaggle.com/robikscube/hourly-energy-consumption>
- [11] N. Charlton and C. Singleton, "A refined parametric model for short term load forecasting," *Int. J. Forecasting*, vol. 30, no. 2, pp. 364–368, 2014.
- [12] S. Fan and R. J. Hyndman, "Short-term load forecasting based on a semi-parametric additive model," *IEEE Trans. Power Syst.*, vol. 27, no. 1, pp. 134–141, Feb. 2012.
- [13] V. Hinojosa and A. Hoese, "Short-term load forecasting using fuzzy inductive reasoning and evolutionary algorithms," *IEEE Trans. Power Syst.*, vol. 25, no. 1, pp. 565–574, Feb. 2010.
- [14] J. W. Taylor, "Short-term electricity demand forecasting using double seasonal exponential smoothing," *J. Oper. Res. Soc.*, vol. 54, no. 8, pp. 799–805, 2003.
- [15] S.-J. Huang and K.-R. Shih, "Short-term load forecasting via ARMA model identification including non-Gaussian process considerations," *IEEE Trans. Power Syst.*, vol. 18, no. 2, pp. 673–679, May 2003.
- [16] B.-J. Chen *et al.*, "Load forecasting using support vector machines: A study on EUNITE competition 2001," *IEEE Trans. Power Syst.*, vol. 19, no. 4, pp. 1821–1830, Nov. 2004.
- [17] M. Singh and R. Maini, "Various electricity load forecasting techniques with pros and cons," *Int. J. Recent Technol. Eng.*, vol. 8, no. 6, pp. 220–229, 2020.
- [18] K. N. Singh *et al.*, "A review on deep learning models for short-term load forecasting," *Appl. Artif. Intell. Mach. Learn.*, vol. 778, pp. 705–721, 2021.
- [19] L. Wang, Z. Zhang, and J. Chen, "Short-term electricity price forecasting with stacked denoising autoencoders," *IEEE Trans. Power Syst.*, vol. 32, no. 4, pp. 2673–2681, Jul. 2017.
- [20] H. Shi, M. Xu, and R. Li, "Deep learning for household load forecasting—A novel pooling deep RNN," *IEEE Trans. Smart Grid*, vol. 9, no. 5, pp. 5271–5280, Sep. 2017.
- [21] Z. Deng, B. Wang, Y. Xu, T. Xu, C. Liu, and Z. Zhu, "Multi-scale convolutional neural network with time-cognition for multi-step short-term load forecasting," *IEEE Access*, vol. 7, pp. 88058–88071, 2019.
- [22] Ö. F. Ertugrul, "Forecasting electricity load by a novel recurrent extreme learning machines approach," *Int. J. Elect. Power Energy Syst.*, vol. 78, pp. 429–435, 2016.

- [23] C. Tong, J. Li, C. Lang, F. Kong, J. Niu, and J. J. Rodrigues, "An efficient deep model for day-ahead electricity load forecasting with stacked denoising auto-encoders," *J. Parallel Distrib. Comput.*, vol. 117, pp. 267–273, 2018.
- [24] X. Dong, L. Qian, and L. Huang, "Short-term load forecasting in smart grid: A combined CNN and K-means clustering approach," in *Proc. IEEE Int. Conf. Big Data Smart Comput.*, 2017, pp. 119–125.
- [25] K. Amarasinghe, D. L. Marino, and M. Manic, "Deep neural networks for energy load forecasting," in *Proc. IEEE 26th Int. Symp. Ind. Electron.*, 2017, pp. 1483–1488.
- [26] I. Koprinska, D. Wu, and Z. Wang, "Convolutional neural networks for energy time series forecasting," in *Proc. Int. Joint Conf. Neural Netw.*, 2018, pp. 1–8.
- [27] L. Li, K. Ota, and M. Dong, "Everything is image: CNN-based short-term electrical load forecasting for smart grid," in *Proc. 14th Int. Symp. Pervasive Syst. Algorithms Netw., 11th Int. Conf. Frontier Comput. Sci. Technol., 3rd Int. Symp. Creative Comput.*, 2017, pp. 344–351.
- [28] C. Tian, J. Ma, C. Zhang, and P. Zhan, "A deep neural network model for short-term load forecast based on long short-term memory network and convolutional neural network," *Energies*, vol. 11, no. 12, 2018, Art. no. 3493.
- [29] W. He, "Load forecasting via deep neural networks," *Procedia Comput. Sci.*, vol. 122, pp. 308–314, 2017.
- [30] M. Cai, M. Pipattanasomporn, and S. Rahman, "Day-ahead building-level load forecasts using deep learning vs. traditional time-series techniques," *Appl. Energy*, vol. 236, pp. 1078–1088, 2019.
- [31] L. Han, Y. Peng, Y. Li, B. Yong, Q. Zhou, and L. Shu, "Enhanced deep networks for short-term and medium-term load forecasting," *IEEE Access*, vol. 7, pp. 4045–4055, 2018.
- [32] Y. Wang, L. Jiang, M.-H. Yang, L.-J. Li, M. Long, and L. Fei-Fei, "Eidetic 3D LSTM: A model for video prediction and beyond," in *Proc. Int. Conf. Learn. Representations*, 2019.
- [33] K. Nazeri, E. Ng, T. Joseph, F. Qureshi, and M. Ebrahimi, "Edge-Connect: Structure guided image inpainting using edge prediction," in *Proc. IEEE/CVF Int. Conf. Comput. Vis. Workshops*, 2019, pp. 3265–3274.
- [34] G. Liu, F. A. Reda, K. J. Shih, T.-C. Wang, A. Tao, and B. Catanzaro, "Image inpainting for irregular holes using partial convolutions," in *Proc. Eur. Conf. Comput. Vis.*, 2018, pp. 85–100.
- [35] Y. Wang, X. Tao, X. Qi, X. Shen, and J. Jia, "Image inpainting via generative multi-column convolutional neural networks," in *Proc. Int. Conf. Neural Inf. Process. Syst.*, vol. 31, 2018, pp. 329–338.
- [36] Z. Yan, X. Li, M. Li, W. Zuo, and S. Shan, "Shift-net: Image inpainting via deep feature rearrangement," in *Proc. Eur. Conf. Comput. Vis.*, 2018, pp. 1–17.
- [37] Z. Qin, Q. Zeng, Y. Zong, and F. Xu, "Image inpainting based on deep learning: A review," *Displays*, vol. 69, 2021, Art. no. 102028.
- [38] Y. Zeng, J. Fu, H. Chao, and B. Guo, "Learning pyramid-context encoder network for high-quality image inpainting," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 1486–1494.
- [39] C. Zheng, T.-J. Cham, and J. Cai, "Pluralistic image completion," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 1438–1447.
- [40] R. A. Yeh, C. Chen, T. Yian Lim, A. G. Schwing, M. Hasegawa-Johnson, and M. N. Do, "Semantic image inpainting with deep generative models," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 5485–5493.
- [41] J. Yu, Z. Lin, J. Yang, X. Shen, X. Lu, and T. S. Huang, "Generative image inpainting with contextual attention," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 5505–5514.
- [42] J. Yu, Z. Lin, J. Yang, X. Shen, X. Lu, and T. Huang, "Free-form image inpainting with gated convolution," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2019, pp. 4471–4480.
- [43] D. J. Michael, "Smart meters in London," 2019. [Online]. Available: <https://www.kaggle.com/jeanmidev/smart-meters-in-london>
- [44] J. Durbin and S. J. Koopman, *Time Series Analysis by State Space Methods*. London, U.K.: Oxford Univ. Press, 2012.
- [45] R. Hyndman, A. B. Koehler, J. K. Ord, and R. D. Snyder, *Forecasting With Exponential Smoothing: The State Space Approach*. Berlin, Germany: Springer, 2008.
- [46] V. Assimakopoulos and K. Nikolopoulos, "The theta model: A decomposition approach to forecasting," *Int. J. Forecasting*, vol. 16, no. 4, pp. 521–530, 2000.
- [47] A. M. D. Livera, R. J. Hyndman, and R. D. Snyder, "Forecasting time series with complex seasonal patterns using exponential smoothing," *J. Amer. Stat. Assoc.*, vol. 106, no. 496, pp. 1513–1527, 2011.
- [48] R. J. Hyndman and G. Athanasopoulos, *Forecasting: Principles and Practice*. Australia: OTexts, 2018.
- [49] L. Prokhorenkova, G. Gusev, A. Vorobev, A. V. Dorogush, and A. Gulin, "Catboost: Unbiased boosting with categorical features," in *Proc. Int. Conf. Neural Informat. Process. Syst.*, 2018, pp. 6639–6649.
- [50] J. Trapero, N. Kourentzes, and R. Fildes, "On the identification of sales forecasting models in the presence of promotions," *J. Oper. Res. Soc.*, vol. 66, pp. 299–307, 2014.
- [51] D. Salinas, V. Flunkert, J. Gasthaus, and T. Januschowski, "DeepAR: Probabilistic forecasting with autoregressive recurrent networks," *Int. J. Forecasting*, vol. 36, no. 3, pp. 1181–1191, 2020.
- [52] B. N. Oreshkin, D. Carpo, N. Chapados, and Y. Bengio, "N-beats: Neural basis expansion analysis for interpretable time series forecasting," in *Proc. Int. Conf. Learn. Representations*, 2020. [Online]. Available: <https://openreview.net/forum?id=r1ecqn4YwB>
- [53] A. I. Borovykh, S. M. Bohte, and C. W. Oosterlee, "Conditional time series forecasting with convolutional neural networks," in *Lecture Notes Comput. Science/Lecture Notes Artif. Intell.*, 2017, pp. 729–730.
- [54] A. Vaswani *et al.*, "Attention is all you need," in *Proc. Adv. Neural Inf. Process. Syst.*, 2017, pp. 5998–6008.
- [55] A. Botchkarev, "Performance metrics (error measures) in machine learning regression, forecasting and prognostics: Properties and typology," in *IJIKM* 14, 2019, pp. 45–79.
- [56] P. Montero-Manso, "Monash time series forecasting archive," in *35th Conf. Neural Informat. Process. Syst. Datasets Benchmarks Track (Round 2)*, 2021.