

Parallel Spatio-Temporal Attention-Based TCN for Multivariate Time Series Prediction

Jin Fan, Ke Zhang, Yipan Huang, Yifei Zhu, Baiping Chen*

School of Computer Science and Technology

Hangzhou Dianzi University

Hangzhou, China

{fanjin, ke.zhang, yipan.huang, zhuyifei, chenbp}@hdu.edu.cn

Abstract—As industrial systems become more complex and monitoring sensors for everything from surveillance to our health become more ubiquitous, multivariate time series prediction is taking an important place in the smooth-running of our society. A recurrent neural network with attention to help extend the prediction windows is the current-state-of-the-art for this task. However, we argue that their vanishing gradients, short memories, and serial architecture make RNNs fundamentally unsuited to long-horizon forecasting with complex data. Temporal convolutional networks (TCNs) do not suffer from gradient problems and they support parallel calculations, making them a more appropriate choice. Additionally, they have longer memories than RNNs, albeit with some instability and efficiency problems. Hence, we propose a framework, called PSTA-TCN, that combines a parallel spatio-temporal attention mechanism to extract dynamic internal correlations with stacked TCN backbones to extract features from different window sizes. The framework makes full use parallel calculations to dramatically reduce training times, while substantially increasing accuracy with stable prediction windows up to 13 times longer than the status quo.

Index Terms—multivariate time series prediction, spatio-temporal attention, parallel stacked TCN

I. INTRODUCTION

Complex systems are commonplace in today’s manufacturing plants [1], health monitoring [2] and ensuring these systems run smoothly inevitably involves constant monitoring of numerous diverse data streams, from temperature and pressure sensors to image and video feeds to CPU usage levels, biometric data, etc. [6], [7], [21]–[24]. However, rather than merely watching for sensor readings to approach certain thresholds, today’s smart analytics systems must look to predict eventualities based on historical patterns. And, generally speaking, the more historical data that can be considered in a prediction, the better the chances of capturing patterns in different variables, and the more accurate the prediction. Presently, recurrent neural networks (RNNs) are the go-to approach for multivariate time series prediction [15], [16]. However, we argue that RNNs are fundamentally ill-suited to this task. They are plagued by issues with vanishing gradients, and techniques like LSTM and GRUs only lessen the problem,

they do not solve it. Even with attention to try and focus on the most important information, RNNs still struggle to capture a sufficient amount of temporal context for highly accurate predictions. Further, because calculations for the current time step need to be completed before starting the next, RNNs tend to spend an excessive amount of time inefficiently waiting for results. Temporal convolutional networks (TCN) [13] suffer from none of these problems. Unlike RNNs, TCNs do not have gradient issues; they support layer-wise computation, which means every weight in every layer can be updated in every time step simultaneously; and, although not excessively, their memories are longer than RNNs. Hence, TCNs have three very significant advantages over RNNs. However, conventional TCNs give every feature equal weight, which results in the limitation of accuracy because every feature has different importance.

Our solution is, therefore, a feedforward network architecture to combine the advantages of a TCN while avoiding the disadvantages of an RNN. Generally, given the target time series y_1, y_2, \dots, y_{T-1} with $y_t \in \mathbb{R}$, the objective is to predict y_T . Such forecasting method ignore the affect of exogenous series x_1, x_2, \dots, x_T with $x_t \in \mathbb{R}^n$. So we choose to combine the exogenous series with the target series as input, i.e., $y_T = F(x_1, x_2, \dots, x_T, y_1, y_2, \dots, y_{T-1})$, $F(\cdot)$ is a nonlinear mapping need to learn.

Moreover, inspired by attention mechanism [14], which has less parameters compared with CNNs and RNNs, and the demand for arithmetic is even smaller. More importantly, attention mechanism does not depend on the calculation results of the previous steps, so it can be processed in parallel with TCN, which makes the model improve performance without consuming too much time. Therefore, we propose a novel attention mechanism comprising both spatial and temporal attention running in parallel to even further improve accuracy and stability. The spatial attention stream gives different weights to the various exogenous features, while the temporal attention stream extracts the correlations between all time steps within the attention window. We also provide an exhaustive interpretation for the fluctuation of single-step prediction in different historical window sizes. Hence, the key advancement made by this work is a framework for multivariate time series prediction that consists of a parallel spatio-temporal attention

This work was supported by a grant from The National Natural Science Foundation of China(No.U1609211), National Key Research and Development Project(2019YFB1705100). The corresponding author is Baiping Chen.

mechanism (PSTA) that can extract internal correlations from exogenous series in parallel branches and two stacked TCN backbones. Our experiments show PSTA-TCN framework has three distinct advantages over the current alternatives:

- **Speed:** PSTA-TCN trains 14 times faster than the current state-of-the-art DSTP [35] and 12 times faster than DSTP’s predecessor DARNN [34].
- **Stability:** Our proposed parallel mechanism has improved the stability of TCN in long-term and long history time series prediction.
- **Accuracy:** our method is verified to perform better than the most advanced time series forecasting methods in both single-step and multi-step predictions.

II. RELATED WORKS

Time series prediction is fundamental to the human condition. No area of activity escapes our desire to prepare, profit or prevent through forecasting, be it finance forecasting [3], weather forecasting [4], [5], human activity detection [8], energy consumption prediction [10], industrial fault diagnosis [41] and etc..

Our explorations into the domain of sequence modeling to generate these forecasts, i.e., time-series predictions, have taken us from statistical engines to multi-layer perceptrons (MLPs) to recursive models [26], [27].

Traditional statistical methods of time-series analysis, such as ARIMA [28] and SVR [29], date back as far as 1970. These are lightweight methods, but they cannot balance spatial correlations with temporal dependencies. MLPs were, arguably, the first post-NN solution to sequence modeling. They are fairly simplistic networks that operate linearly and do not share parameters. Although still relevant to many applications where time series prediction is needed, MLPs quickly become unwieldy with large numbers of input parameters, as is common with today’s complex monitoring systems. With advances in deep learning, RNNs came to be the default scheme for time series modeling [11], [12]. RNNs share parameters in each time step, and each time step is a function of its previous time step, which means, in theory, RNNs have unlimited memory [13], [30]. However, RNNs suffer from issues with vanishing gradient when a data sequence becomes too long [31]. Long short-term memory (LSTM) [32] and gated recurrent units (GRU) [33] can lessen this problem, but not to the extent that long short-term becomes long term. The field of vision, both forwards and backwards, is still limited.

The current state-of-the-art RNN solutions both involve attention. Qin et al. [34] developed DARNN, a dual-staged attention-based recurrent network, in 2017. And after that, Lieu et al. [35] published an improved version of DARNN, called DSTP (dual-stage two-phase attention), which employs multiple attention layers to jointly select the most relevant input characteristics and capture long-term temporal dependencies.

Although attention-based RNNs have many strengths, they have some inherent flaws that cannot be overcome. As mentioned, one is serial calculation, i.e., the calculation for the

current time step must be completed before the calculation for the next time step can begin. Hence, processes like training and testing cannot be parallelized [36].

TCNs support parallel computing and, further, with a feed-forward model, they can be used for sequence modeling [13]. Further, unlike RNNs, the hierarchical structure of TCNs makes it possible to capture long-range patterns. Although, we find that predictions with very long sequences (e.g., the length is 32) are not particularly efficient or stable. Hence, we designed a novel spatio-temporal attention mechanism to address these issues. The result is a framework for multivariate time series prediction that leverages the best thinking from both TCN and RNN-based strategies, as outlined in the next section.

III. SPATIO-TEMPORAL ATTENTION BASED TCN

PSTA-TCN comprises a parallel spatio-temporal attention mechanism and two stacked TCN backbones. In this section, we provide an overview of the network architecture and details of these two main systems, beginning with the problem statement and notations.

A. Notation and Problem Statement

Consider a multivariate exogenous series $X = \langle X^{(1)}, X^{(2)}, \dots, X^{(n)} \rangle \in \mathbb{R}^{n \times T}$, where n denotes the dimensions of the exogenous series, and T is the length of the window size. The i -th exogenous series $X^{(i)}$ is denoted as $X^{(i)} = \langle X_1^{(i)}, X_2^{(i)}, \dots, X_T^{(i)} \rangle \in \mathbb{R}^T$, where the length of $X^{(i)}$ is also T . The target series is defined as $Y = \langle y_1, y_2, \dots, y_T \rangle \in \mathbb{R}^T$, also with a length of T . Typically, given the previous exogenous series $X = \langle X^{(1)}, X^{(2)}, \dots, X^{(n)} \rangle$ and target series $Y = \langle y_1, y_2, \dots, y_T \rangle$, we aim to predict the future value of $\hat{Y} = \langle \hat{y}_{T+1}, \hat{y}_{T+2}, \dots, \hat{y}_{T+\tau} \rangle \in \mathbb{R}^\tau$, where τ is the time step to predict. The objective is formulated as:

$$\hat{y}_{T+1}, \hat{y}_{T+2}, \dots, \hat{y}_{T+\tau} = F(X_1, X_2, \dots, X_T, Y) \quad (1)$$

where $F(\cdot)$ is the nonlinear mapping we aim to learn.

B. Model

Fig.1 shows the architecture of our proposed PSTA-TCN model. The input, which is a multivariate time series comprising both exogenous and target series, is fed into two parallel backbones simultaneously. One backbone begins with a spatial attention block for extracting the spatial correlations between the exogenous and target series. The other begins a temporal attention block to capture the temporal dependencies between all time steps in the window. The output of these blocks is then transmitted through two identical stacked TCN backbones. After processing dilated convolutions and residual connections, the results are delivered to a dense layer and then summed to produce the final prediction.

Parallel Spatial-Temporal Attention

Inspired by multi-stage attention models, we employ a spatial attention block to extract spatial correlations between exogenous series and historical target series. Meanwhile, we

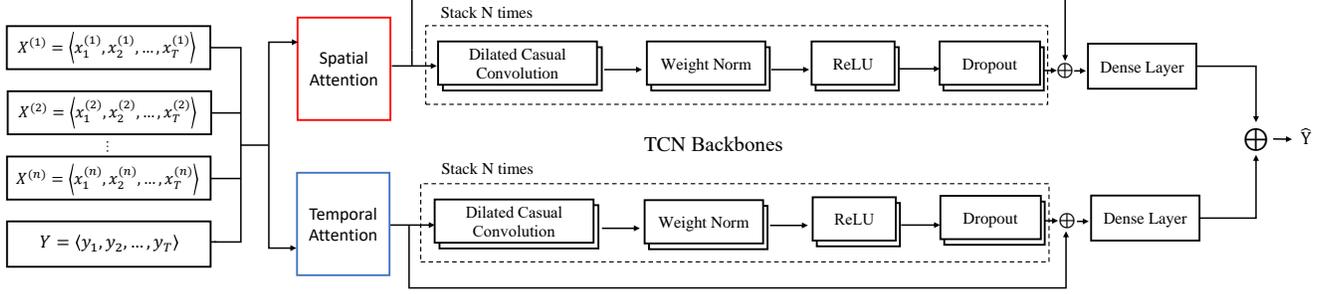


Fig. 1. Overview of the PSTA-TCN architecture. PSTA-TCN comprises input layers, attention blocks, TCN backbones and dense layers. The model input is $\langle X^{(1)}, X^{(2)}, \dots, X^{(n)}, Y \rangle$, and the output is $\hat{Y} = \langle \hat{y}_{T+1}, \hat{y}_{T+2}, \dots, \hat{y}_{T+\tau} \rangle$, where τ is the future time step to predict, T is the window size, n is the dimension of exogenous series X .

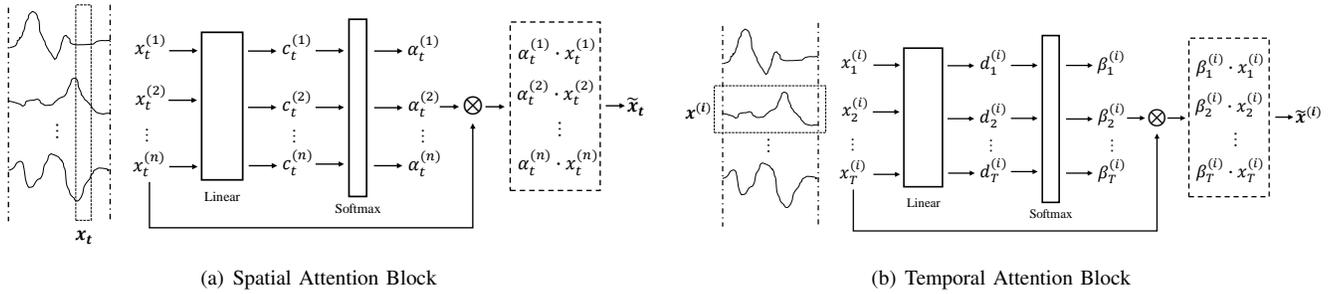


Fig. 2. Two inter-layer transformation diagram. (a) Transformation details in spatial attention block. Input $x_t = \langle X_t^{(1)}, X_t^{(2)}, \dots, X_t^{(n)} \rangle$ has been framed by vertical dashed lines in the left, c_t is the intermediate weight obtained after the linear transformation, α_t is normalized with a softmax operation, and \tilde{x}_t represents for the output already weighted. (b) Transformation details in temporal attention block. Input $x^i = \langle X_1^{(i)}, X_2^{(i)}, \dots, X_T^{(i)} \rangle$ has been framed by horizontal dashed lines in the left, d^i is the intermediate weight obtained after linear transformation, β^i is normalized with a softmax operation, and \tilde{x}^i is the final weighted output.

use a temporal attention block to obtain long history temporal dependencies across window size T . Fig. 2 shows the inter-layer transformations in temporal attention block and spatial attention block, respectively. We omit the description about the processing of input Y to be succinct. Fig. 2(a) shows the workflow for the spatial attention block. The input is formulated as $x_t = \langle X_t^{(1)}, X_t^{(2)}, \dots, X_t^{(n)} \rangle$, where n indicates the dimensions of the full exogenous series, and t indicates a time step in the current window. First, a spatial attention weight vector c_t is generated to represent the importance of each feature in time step t a by applying a linear transformation to the original input:

$$c_t = W_c^\top x_t + b_c \quad (2)$$

where $W_c \in \mathbb{R}^{n \times 1}$, $b_c \in \mathbb{R}$ are the parameters to learn. Next, the weighted vector c_t is normalized with a softmax function to ensure all the attentions sum to 1, resulting in vector α_t :

$$\alpha_t^{(k)} = \frac{\exp(c_t^{(k)})}{\sum_{i=1}^{n+1} \exp(c_t^{(i)})} \quad (3)$$

Fig. 2(b) shows the process for calculating temporal attention. The input takes the form $x^{(i)} = \langle X_1^{(i)}, X_2^{(i)}, \dots, X_T^{(i)} \rangle$,

where i indicates the i -th exogenous series and T is the window size. Again, a linear transformation of the original input produces a temporal attention weight vector $d^{(i)}$ reflecting the importance of i -th exogenous series among all time steps.

$$d^{(i)} = W_d^\top x^{(i)} + b_d \quad (4)$$

where $W_d \in \mathbb{R}^{T \times 1}$, $b_d \in \mathbb{R}$ are the parameters to learn. And the vector d^i is normalized with a softmax function.

$$\beta_t^{(i)} = \frac{\exp(d_t^{(i)})}{\sum_{t=1}^T \exp(d_t^{(i)})} \quad (5)$$

where the current time step $t \in [1, T]$.

Stacked TCN backbones

As a new exploration in sequence modeling, TCN benefits from CNNs (i.e. convolutional network [17] based models) with stronger parallelism and more flexible receptive fields than RNNs, and requires less memory when facing long sequences. As shown in Fig. 1, we use generic TCN as basic backbone, and stack one TCN backbone for N times to provide N levels.

Convolution layers in TCN is causal which means there is no "information leakage", i.e., when calculating the output

at time step t , only the states at or before time step t are convolved.

Dilated convolution stops the network from growing too deep when dealing with long sequences by forcing the receptive field of each layer to grow exponentially, as a larger receptive field with fewer parameters and fewer layers is more beneficial. The effective history in each layer of TCN is $(k-1)d$, where k is the kernel size, and d is the dilated factor. For the purpose of controlling the amount of parameters, we choose a fixed size of k , and each layer increases the value of d exponentially, i.e., $d = 2^i$ where i means the level of the network.

However, when faced with ultra-long sequences, dilated convolution will not be enough. A deeper network will need to be trained to make the model sufficiently powerful, which we do using residual connections to avoid the issue of vanishing gradients. The residual connections can be defined by adding up X and $F(X)$:

$$\text{Output} = \text{ReLU}(X + F(X)) \quad (6)$$

where X represent for the original input, $F(\cdot)$ means the processing of one TCN backbone.

IV. EXPERIMENTS AND RESULTS

A. Datasets

To test PSTA-TCN, we compared its performance in a bespoke prediction task against 5 other methods: 2 RNNs, 2 RNNs with attention (the current state-of-the-arts), and 1 vanilla TCN as a baseline. The experimental scenario was human activity, and the task was to make long-term motion prediction.

To collect the data, we attached four wearable micro-sensors [38] to 10 participants and asked them to perform five sessions of 10 squats. The sensors (configured with the master on the left arm and slaves on the right arm and each knee) measure acceleration and angular velocity data along three axes and visualize it in a mobile app connected by Bluetooth. Fig. 3 pictures the wearable microsensors, one of the participants fitted with the devices and the mobile app interface. Sampling 50 times per second for the duration of the exercise (approx 0.02 seconds), we gathered 81,536 data points in each of 24 data series, i.e., 4 sensors * 3 axes * 2 dimensions (acceleration and angular velocity) for each participant to constitute a multivariate time series of 1.96 million data. For clarity, we list a sample of acceleration and angular velocity data from our dataset in Table I. Our prediction task on self-designed dataset can be formulated as:

$$\begin{aligned} \hat{A}_{T+1}, \hat{A}_{T+2}, \dots, \hat{A}_{T+\tau} &= F(A_{X1}, \dots, A_{XT}, A_{Y1}, \dots, A_{YT}, \\ &A_{Z1}, \dots, A_{ZT}, \hat{A}_1, \dots, \hat{A}_T, V_{X1}, \dots, V_{XT}, V_{Y1}, \dots, V_{YT}, \\ &V_{Z1}, \dots, V_{ZT}, \hat{V}_1, \dots, \hat{V}_T) \end{aligned} \quad (7)$$

where $A_X = (A_{X1}, \dots, A_{XT})$, $A_Y = (A_{Y1}, \dots, A_{YT})$ and $A_Z = (A_{Z1}, \dots, A_{ZT})$ are a window size of the acceleration data in X-axis, Y-axis and Z-axis, respectively. Likewise,

TABLE I
A SAMPLE OF ACCELERATION AND ANGULAR VELOCITY DATA.

	Acceleration			Angular Velocity		
	X	Y	Z	X	Y	Z
Master	-8.00155	0.08966	-0.71372	84.5	150.8	40.4
Slave-1	11.62156	-1.68806	-0.61927	-83.7	-179.8	162.7
Slave-2	-9.81514	-0.19487	-1.71795	82.0	-70.5	151.7
Slave-3	11.16128	0.78904	-0.61688	-78.4	178.6	14.9

V_X , V_Y and V_Z are a window size of the angular velocity data in X-axis, Y-axis and Z-axis, respectively. \hat{A}_t and \hat{V}_t represent for the resultant acceleration and resultant Angular velocity at a historical time step t separately. Meanwhile, $\hat{A}_{T+1}, \hat{A}_{T+2}, \dots, \hat{A}_{T+\tau}$ is the target series we need to predict and τ represents for the number of prediction steps. $F(\cdot)$ is the nonlinear mapping we aim to learn.

In our experiment, the dataset treats 1.96 million data as a whole which is chronologically split into training set and test set by a ratio of 4:1. Additionally, we segmented each dataset into windows using the sliding window method [25] and, to avoid overfitting, we randomly shuffled all the windows. The specific parameter settings will be introduced in the next Section IV-C.

B. Baseline methods

LSTM [32]: LSTM was designed to solve the gradient vanishing problem in standard RNNs. It uses gated units to selectively retain or remove information in time series data, capturing long-term dependencies in the process.

GRU [33]: As a variant of LSTM, GRU merges different gated units in LSTM, and also combines the cell state and the hidden state, making the model lighter and suitable for scenarios with smaller amounts of data.

DARNN [34]: The first of the state-of-the-art methods, DARNN is a single-step predictor. It uses dual-stage attention to capture dependencies in both input exogenous data and encoder hidden states.

DSTP [35]: DSTP is the second of the state-of-the-arts. Its basic structure is similar to DARNN, but it involves an additional phase of attention so as to process the exogenous series and the target series separately.



Fig. 3. The wearable microsensors; a participant wearing the devices; the app interface and data visualization

TABLE II
SINGLE-STEP PREDICTION AMONG DIFFERENT WINDOW SIZE

Window size	Metrics	Methods					
		LSTM	GRU	DARNN	DSTP	TCN	PSTA-TCN
32	RMSE	0.0821	0.0842	0.0767	0.0777	0.0629	0.0579
	MAE	0.0507	0.0524	0.0241	0.0223	0.0293	0.0238
64	RMSE	0.0863	0.0872	0.0781	0.0786	0.0659	0.0612
	MAE	0.0532	0.0549	0.0331	0.0250	0.0316	0.0306
128	RMSE	0.0942	0.0922	0.0762	0.0804	0.0735	0.0706
	MAE	0.0631	0.0576	0.0509	0.0239	0.0429	0.0425
256	RMSE	0.1006	0.1084	0.8681	0.0796	0.0701	0.0682
	MAE	0.0735	0.0640	0.0540	0.0505	0.0394	0.0388

TABLE III
MULTI-STEP PREDICTION AMONG DIFFERENT PREDICTING STEPS

Prediction step	Metrics	Methods					
		LSTM	GRU	DARNN	DSTP	TCN	PSTA-TCN
2	RMSE	0.0947	0.1278	0.0863	0.1013	0.0850	0.0842
	MAE	0.0461	0.0634	0.0468	0.0372	0.0473	0.0505
4	RMSE	0.1423	0.1785	0.1158	0.1403	0.1036	0.0893
	MAE	0.0638	0.0887	0.0683	0.0697	0.0662	0.0598
8	RMSE	0.2568	0.2340	0.2089	0.1897	0.1268	0.1060
	MAE	0.1221	0.1035	0.1393	0.1162	0.0840	0.0673
16	RMSE	0.3567	0.3398	0.3166	0.3091	0.1216	0.1094
	MAE	0.2534	0.1676	0.2347	0.2099	0.0758	0.0773
32	RMSE	0.5957	0.5012	0.4705	0.4484	0.2090	0.1122
	MAE	0.3624	0.2785	0.3512	0.3172	0.1496	0.0697

TCN [13]: This is a vanilla TCN consisting of causal convolution, residual connection and dilation convolution. The receptive fields are flexible, and parallel calculations are supported.

C. Hyperparameter setting and evaluation metrics

We conducted two main sets of experiments – first single-step predictions, then multi-step predictions. During the training process, we set the batch size to 64 and the initial learning rate to 0.001. With the single-step predictions, we tested the performance of each model with different window sizes $T \in \{32, 64, 128, 256\}$, i.e., with different amounts of historical information. With the multi-step predictions, we fixed the window size to $T = 32$, and varied the prediction steps $\tau \in \{2, 4, 8, 16, 32\}$ to verify the impact of different prediction steps. To be fair, we conducted a grid search for all models to find the best hyperparameter settings. Specifically, we set $m = p = 128$ for DARNN, $m = p = q = 128$ for DSTP. As for TCN and our model, we set the kernel size to 7 and level to 8. To ensure the reproducibility of

experimental results, we set the random seeds to an integer for all experiments, which is 1111.

We chose the two most commonly used assessment metrics in the field of time series forecasting for the evaluation: root mean squared error (RMSE) and mean absolute error (MAE). The specific formulations used were:

$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^N (\hat{y}_t^i - y_t^i)^2} \quad (8)$$

$$MAE = \frac{1}{N} \sum_{i=1}^N |\hat{y}_t^i - y_t^i| \quad (9)$$

where y_t is the ground truth at time step t and \hat{y}_t is the predicted value at time step t . Lower rates of both reflect better accuracy.

D. Results

The results for the single-step predictions are shown in Table II, and the multi-step predictions are provided in Table III. Fig. 4 represents the results as a line chart. Across all tests, PSTA-TCN consistently achieved the lowest RMSE and MAE scores by a substantial margin.

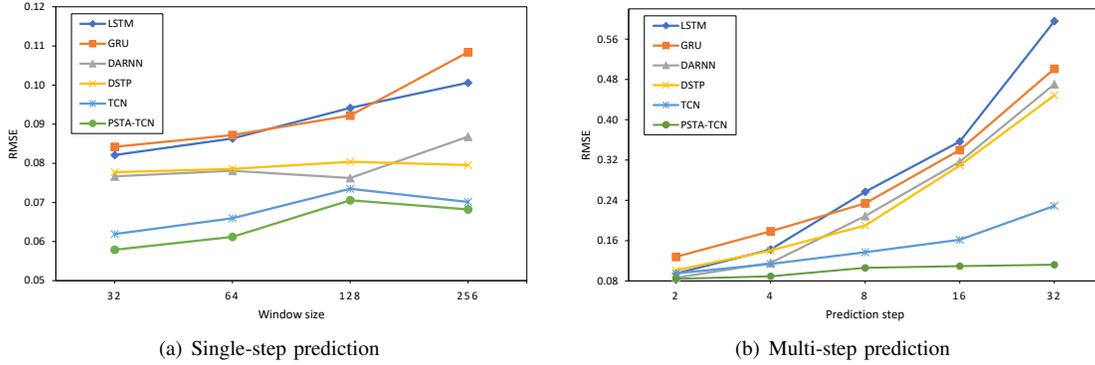


Fig. 4. Performance of single-step prediction and multi-step prediction. All baseline methods are compared with our proposed methods

The results in Table II show accuracy with different amounts of historical information. LSTM and GRU are relatively old models. They do not have attention, which means they have no effective way of screening past information so, as expected, their performance was sub-par. There was little difference between DARNN and DSTP in terms of prediction quality, with DSTP doing marginally better due to multiple attention. However, Fig. 5 does show some significant differences in training time depending on the window size T , which is discussed further in the next section.

TCN and PSTA-TCN were significantly more accurate, plus accuracy began to increase again after a nadir as the window size passed 128. We would expect the RMSE to decline as the window size expands and more historical information is considered in the prediction. However, what we find is a fluctuation, particularly for the two TCN methods. Upon further analysis, we find two reasons to explain this phenomenon: 1) When the historical information increases, spatio-temporal attention does not capture enough of the long-term dependencies; hence, the network needs to deepen to consider more parameters; and 2) as the input data extends, the load on the model increases significantly, making it harder to train. Therefore, although a larger window size brings more reference information, it also increases the difficulty of training the model, and the resulting cycle of deepening the network and training the parameters manifests as fluctuations in the final accuracy.

In terms of the multi-step predictions (Table III and Fig. 4(b)), the clearest observation is that the accuracy of the RNN-based methods declines significantly more as the number of prediction steps increases, relative to the TCN-based methods. Notably, PSTA-TCN remained remarkably accurate, even when predicting very long sequences. In contrast to the RNNs, PSTA-TCN was much more stable and was better able to extract the spatio-temporal dependencies from historical information.

In comparison to the baseline form of TCN, the addition of parallel attention meant PSTA-TCN was able to maintain a high level of accuracy well beyond the 32 steps where TCN began to obviously decline. We speculate that the reason is, in

the long-term prediction, our proposed spatio-temporal attention mechanism extracts more hierarchical feature information from the original data, which makes our model have more reference to do long-term prediction under the same historical window size compared with vanilla TCN. Overall, these results demonstrate PSTA-TCN to be a very promising strategy for improving stability and extending the longevity of network memory for multivariate time series prediction.

V. FURTHER EXPERIMENTS

A. Time complexity

Fig. 5 compares the training time of each model with different window sizes T at a training batch size of 64. What is clear is that the calculation time for both DARNN and DSTP increases significantly as the window size increases. This is due to the serial nature of the underlying RNNs and the complexity of the attention mechanisms. At $T = 256$, DSTP takes 46 times longer to train than vanilla TCN, and 14 times longer than PSTA-TCN. DARNN's complexity is not much better at 42 times TCN and 13 times PSTA-TCN. Hence, we find that in the face of more historical information, both DARNN and DSTP begin to lose their luster.

In practice, RNNs tend to spend an excessive amount of time waiting for the calculation results of the previous time step, whereas TCNs leverage parallel computing to radically reduce the amount of training time required. Our strategy is

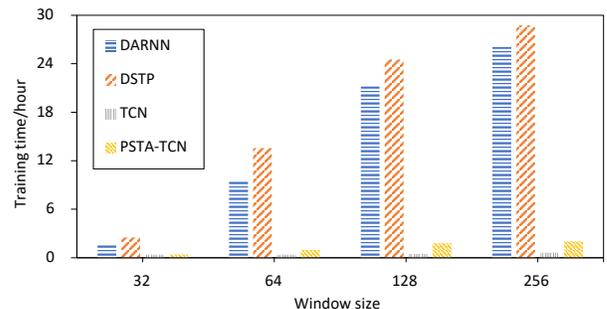


Fig. 5. Training time comparison on single-step prediction among different window sizes

to sacrifice part of this reduction in favor of a spatio-temporal attention mechanism to leverage long sequences and maximize accuracy. As a result, PSTA-TCN is more stable with long sequences than a standard TCN, and faster and more accurate than an RNN.

B. Ablation studies

To explore the contribution of each module in PSTA-TCN, we compared PSTA-TCN with its variants as follows:

- **P-TCN**: Remove all attention, leaving only parallel TCN backbones.
- **PSA-TCN**: Remove temporal attention, leaving only spatial attention module.
- **PTA-TCN**: Remove spatial attention, leaving only temporal attention module.

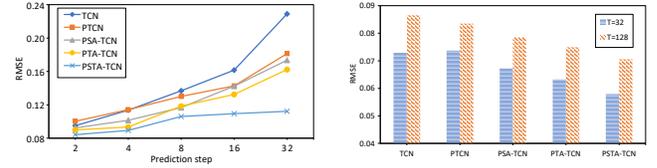
Fig. 6(a) shows the stepwise results for the multi-step prediction experiments, and Fig. 6(b) shows the single-step prediction performance for each model with window sizes of $T = \{32, 128\}$ to reflect a relatively short history and a relatively long history.

From Fig. ?? we can observe that: 1) our model outperformed PTA-TCN and PSA-TCN by a considerable margin, neither spatial attention, temporal attention, nor the parallel backbones are primarily responsible for PSTA-TCN’s performance improvement over TCN. In fact, it is only when all three are combined that we see accuracy improve by a considerable margin.

2) Parallel TCN, as a model combination method we proposed, provides additional information to improve overall performance. In multi-step forecasting, the performance of P-TCN was significantly more accurate than vanilla TCN, especially as the prediction horizon grew longer. The reason we presume is that the parallel TCN backbones extend the vanilla TCN with much more parameters, so that our model has stronger expression ability and a better performance in the high difficulty of long-term prediction task. The innovative application of P-TCN is one of the reasons why PSTA-TCN is able to maintain stability as the number of prediction steps increases.

C. Influence of Hyperparameters

Finally we investigate the influence of hyperparameters in stacked TCN backbones, which are hidden dimension, the number of levels and kernel size. The results of single-step prediction are in Fig. 7(a)7(b)7(c). As we can see, the RMSE curve of the model falls first and then rises, which means we do have the optimal choice of hyperparameters to strengthen our model. For instance, $H=12$ for the hidden dimensions, $L=8$ for the number of layers, $K=7$ for the kernel size. And we also would like to understand the influence of window size when we make multi-step prediction. We have window size $T \in \{8, 16, 32, 64, 128, 256\}$, and the prediction step equals to 32. We control other conditions to do the experiment and get the results as shown in the figure. As we can observed from Fig.7(d), the RMSE curve shows the trend of turbulence, so the window size still have non-negligible effect on the final



(a) The result of multi-step prediction (b) The result of single-step prediction

Fig. 6. Performance comparison among different vairants

prediction, the optimal value of which is 32. Firstly, prediction is poor when the window size is small ($T = 8, 16$). Such a phenomenon mainly comes from, the smaller the window size is, the less historical information the model can use, and the timing characteristics cannot be completely captured. Especially when the historical data are less than the number of steps to be predicted, the prediction effect is very poor. While on the other hand, when the window size is larger than the number of prediction steps, we can find that the accuracy of prediction decreases instead of increasing. We speculate the reason as, the performance improvement of our spatio-temporal attention module is limited. Limited by the window size, when the window size is larger than a threshold, the importance evaluation in the window will be distorted.

VI. CONCLUSION

In this paper, we proposed a novel parallel spatio-temporal attention based TCN (PSTA-TCN), which consists of parallel spatio-temporal attention and stacked TCN backbones. On the basis of the TCN backbone, we makes full use of the parallelism of TCN model to speed up training times while the gradient problems associated with RNNs.

We apply spatial and temporal attention in two different branches to efficiently capture spatial correlations and temporal dependencies, respectively. With the help of this attention mechanism, our proposed PSTA-TCN improves stability over long-term predictions, outperforming the current state-of-the-art by a large margin. Although designed for time series

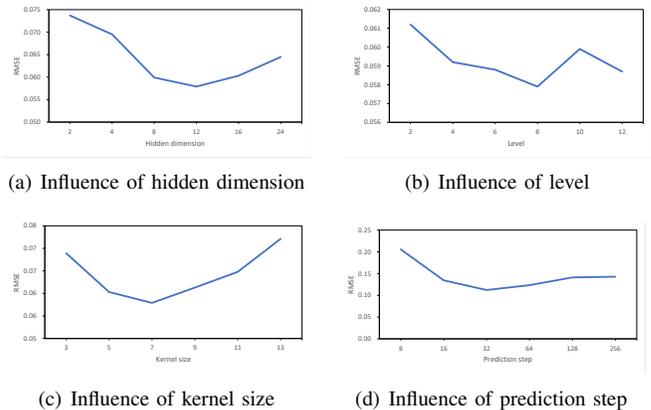


Fig. 7. Influence of different hyperparameters in PSTA-TCN.

forecasting, PSTA-TCN also has potential as a general feature extraction tool in the fields of industrial data mining [40] and fault diagnosis [41]. In the future, we plan to compress PSTA-TCN to adapt to other resource-restrained edge devices while maintaining the original accuracy as much as possible. And we also want to explore an effective combination of CNN and RNN using attention mechanism as the connection module.

ACKNOWLEDGMENT

This work was supported by a grant from The National Natural Science Foundation of China(No.U1609211), National Key Research and Development Project(2019YFB1705102).

CONFLICT OF INTEREST STATEMENT

We declare that we have no financial and personal relationships with other people or organizations that can inappropriately influence our work, there is no professional or other personal interest of any nature or kind in any product, service and/or company that could be construed as influencing the position presented in, or the review of, the manuscript entitled.

REFERENCES

- [1] F. Liu, S. Xue, J. Wu, C. Zhou, W. Hu, C. Paris, S. Nepal, J. Yang, and P. S. Yu, "Deep learning for community detection: Progress, challenges and opportunities," *arXiv preprint arXiv:2005.08225*, 2020.
- [2] R. Zhao, R. Yan, Z. Chen, K. Mao, P. Wang, and R. X. Gao, "Deep learning and its applications to machine health monitoring," *Mechanical Systems and Signal Processing*, vol. 115, pp. 213–237, 2019.
- [3] H. Li, Y. Shen, and Y. Zhu, "Stock price prediction using attention-based multi-input lstm," in *Asian Conference on Machine Learning*, 2018, pp. 454–469.
- [4] E. Soares, P. Costa Jr, B. Costa, and D. Leite, "Ensemble of evolving data clouds and fuzzy models for weather time series prediction," *Applied Soft Computing*, vol. 64, pp. 445–453, 2018.
- [5] F. Zamora-Martínez, P. Romeu, P. Botella-Rocamora, and J. Pardo, "Online learning of indoor temperature forecasting models towards energy efficiency," *Energy and Buildings*, vol. 83, pp. 162–172, 2014.
- [6] J. Wu, X. Zhu, C. Zhang, and S. Y. Philip, "Bag constrained structure pattern mining for multi-graph classification," *Ieee transactions on knowledge and data engineering*, vol. 26, no. 10, pp. 2382–2396, 2014.
- [7] J. Wu, S. Pan, X. Zhu, and Z. Cai, "Boosting for multi-graph classification," *IEEE transactions on cybernetics*, vol. 45, no. 3, pp. 416–429, 2014.
- [8] M. Cornacchia, K. Ozcan, Y. Zheng, and S. Velipasalar, "A survey on activity detection and classification using wearable sensors," *IEEE Sensors Journal*, vol. 17, no. 2, pp. 386–403, 2016.
- [9] J. Wu, S. Pan, X. Zhu, C. Zhang, and X. Wu, "Multi-instance learning with discriminative bag mapping," *IEEE Transactions on Knowledge and Data Engineering*, vol. 30, no. 6, pp. 1065–1080, 2018.
- [10] L. M. Candanedo, V. Feldheim, and D. Deramaix, "Data driven prediction models of energy use of appliances in a low-energy house," *Energy and buildings*, vol. 140, pp. 81–97, 2017.
- [11] M. Han and M. Xu, "Laplacian echo state network for multivariate time series prediction," *IEEE transactions on neural networks and learning systems*, vol. 29, no. 1, pp. 238–244, 2017.
- [12] S. Sivakumar and S. Sivakumar, "Marginally stable triangular recurrent neural network architecture for time series prediction," *IEEE Transactions on Cybernetics*, no. 99, pp. 1–15, 2017.
- [13] S. Bai, J. Z. Kolter, and V. Koltun, "An empirical evaluation of generic convolutional and recurrent networks for sequence modeling," *arXiv preprint arXiv:1803.01271*, 2018.
- [14] R. Hübner, M. Steinhauser, and C. Lehle, "A dual-stage two-phase model of selective attention," *Psychological review*, vol. 117, pp. 759–84, 07 2010.
- [15] Y. Li, Z. Zhu, D. Kong, H. Han, and Y. Zhao, "Ea-lstm: Evolutionary attention-based lstm for time series prediction," *Knowledge-Based Systems*, vol. 181, p. 104785, 2019.
- [16] Y. Hua, Z. Zhao, R. Li, X. Chen, Z. Liu, and H. Zhang, "Deep learning with long short-term memory for time series prediction," *IEEE Communications Magazine*, vol. 57, no. 6, pp. 114–119, 2019.
- [17] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel, "Backpropagation applied to handwritten zip code recognition," *Neural computation*, vol. 1, no. 4, pp. 541–551, 1989.
- [18] N. Kalchbrenner, L. Espeholt, K. Simonyan, A. v. d. Oord, A. Graves, and K. Kavukcuoglu, "Neural machine translation in linear time," *arXiv preprint arXiv:1610.10099*, 2016.
- [19] J. Gehring, M. Auli, D. Grangier, and Y. N. Dauphin, "A convolutional encoder model for neural machine translation," *arXiv preprint arXiv:1611.02344*, 2016.
- [20] Y. N. Dauphin, A. Fan, M. Auli, and D. Grangier, "Language modeling with gated convolutional networks," in *International conference on machine learning*, 2017, pp. 933–941.
- [21] M. Christ, A. W. Kempa-Liehr, and M. Feindt, "Distributed and parallel time series feature extraction for industrial big data applications," *arXiv preprint arXiv:1610.07717*, 2016.
- [22] H. Yan, J. Wan, C. Zhang, S. Tang, Q. Hua, and Z. Wang, "Industrial big data analytics for prediction of remaining useful life based on deep learning," *IEEE Access*, vol. 6, pp. 17 190–17 197, 2018.
- [23] L. Hou and N. W. Bergmann, "Novel industrial wireless sensor networks for machine condition monitoring and fault diagnosis," *IEEE transactions on instrumentation and measurement*, vol. 61, no. 10, pp. 2787–2798, 2012.
- [24] Y. Xu, Y. Sun, J. Wan, X. Liu, and Z. Song, "Industrial big data for fault diagnosis: Taxonomy, review, and applications," *IEEE Access*, vol. 5, pp. 17 368–17 380, 2017.
- [25] S. Huang, D. Wang, X. Wu, and A. Tang, "Dsanet: Dual self-attention network for multivariate time series forecasting," in *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*, 2019, pp. 2129–2132.
- [26] Y. Liang, S. Ke, J. Zhang, X. Yi, and Y. Zheng, "Geoman: Multi-level attention networks for geo-sensory time series prediction," in *IJCAI*, 2018, pp. 3428–3434.
- [27] H. Hao, Y. Wang, Y. Xia, J. Zhao, and F. Shen, "Temporal convolutional attention-based network for sequence modeling," *arXiv preprint arXiv:2002.12530*, 2020.
- [28] G. E. Box and D. A. Pierce, "Distribution of residual autocorrelations in autoregressive-integrated moving average time series models," *Journal of the American statistical Association*, vol. 65, no. 332, pp. 1509–1526, 1970.
- [29] T. Van Gestel, J. A. Suykens, D.-E. Baestaens, A. Lambrechts, G. Lanckriet, B. Vandaele, B. De Moor, and J. Vandewalle, "Financial time series prediction using least squares support vector machines within the evidence framework," *IEEE Transactions on neural networks*, vol. 12, no. 4, pp. 809–821, 2001.
- [30] I. Goodfellow, Y. Bengio, and A. Courville, *Deep learning*. MIT press, 2016.
- [31] Y. Bengio, P. Simard, and P. Frasconi, "Learning long-term dependencies with gradient descent is difficult," *IEEE Transactions on Neural Networks*, vol. 5, no. 2, pp. 157–166, 1994.
- [32] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [33] K. Cho, B. Van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, "Learning phrase representations using rnn encoder-decoder for statistical machine translation," *arXiv preprint arXiv:1406.1078*, 2014.
- [34] Y. Qin, D. Song, H. Chen, W. Cheng, G. Jiang, and G. Cottrell, "A dual-stage attention-based recurrent neural network for time series prediction," *arXiv preprint arXiv:1704.02971*, 2017.
- [35] Y. Liu, C. Gong, L. Yang, and Y. Chen, "Dstp-rnn: A dual-stage two-phase attention-based recurrent neural network for long-term and multivariate time series prediction," *Expert Systems with Applications*, vol. 143, p. 113082, 2020.
- [36] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in neural information processing systems*, 2017, pp. 5998–6008.
- [37] S. Merity, N. S. Keskar, and R. Socher, "Regularizing and optimizing lstm language models," *arXiv preprint arXiv:1708.02182*, 2017.
- [38] J. Fan, H. Wang, Y. Huang, K. Zhang, and B. Zhao, "Aedmts: an attention-based encoder-decoder framework for multi-sensory time series analytic," *IEEE Access*, vol. PP, pp. 1–1, 02 2020.

- [39] R. Dai, L. Minciullo, L. Garattoni, G. Francesca, and F. Bremond, "Self-attention temporal convolutional network for long-term daily living activity detection," in *2019 16th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS)*. IEEE, 2019, pp. 1–7.
- [40] J. Zhu, Z. Ge, Z. Song, and F. Gao, "Review and big data perspectives on robust data mining approaches for industrial process modeling with outliers and missing data," *Annual Reviews in Control*, vol. 46, pp. 107–133, 2018.
- [41] Y. Wang and H. Li, "Industrial process time-series modeling based on adapted receptive field temporal convolution networks concerning multi-region operations," *Computers & Chemical Engineering*, p. 106877, 2020.