# AA-Forecast: Anomaly-Aware Forecast for Extreme Events

**Ashkan Farhangi**
University of Central Florida
Orlando, USA
ashkan.farhangi@ucf.edu

**Jiang Bian**
Baidu Inc.
Beijing, China
bianjiang03@baidu.com

**Arthur Huang**
University of Central Florida
Orlando, USA
arthur.huang@ucf.edu

**Haoyi Xiong**
Baidu Inc.
Beijing, China
xionghaoyi@baidu.com

**Jun Wang**
University of Central Florida
Orlando, USA
jun.wang@ucf.edu

**Zhishan Guo**
University of Central Florida
Orlando, USA
zhishan.guo@ucf.edu

## Abstract

Time series models often deal with extreme events and anomalies, both prevalent in real-world datasets. Such models often need to provide careful probabilistic forecasting, which is vital in risk management for extreme events such as hurricanes and pandemics. However, it is challenging to automatically detect and learn to use extreme events and anomalies for large-scale datasets, which often require manual effort. Hence, we propose an anomaly-aware forecast framework that leverages the previously seen effects of anomalies to improve its prediction accuracy during and after the presence of extreme events. Specifically, the framework automatically extracts anomalies and incorporates them through an attention mechanism to increase its accuracy for future extreme events. Moreover, the framework employs a dynamic uncertainty optimization algorithm that reduces the uncertainty of forecasts in an online manner. The proposed framework demonstrated consistent superior accuracy with less uncertainty on three datasets with different varieties of anomalies over the current prediction models.

## 1 Introduction

Climate change is increasing the severity of natural disasters. Compared to the 1990s, natural disasters quadrupled in terms of economic damage in the U.S. alone [1]. Time series forecasting during the presence of such extreme events (e.g., hurricanes) is critical for resource allocation and resilience planning [2, 3, 4].

Intuitively, the high accuracy and low uncertainty of the forecasts are critical insights for uncovering the influence of external shocks and events on large-scale time series data [5]. To provide sustainable economic development and resilience planning, it is crucial to understand how different industries are influenced by and recover from such extreme events over time [6]. However, it remains a challenge to develop reliable and accurate prediction models as the real-world dataset often contains anomalies that tend to be rare and random. Hence, it is crucial to develop a forecast model that can leverage previously seen extreme events and anomalies for future forecasts.

Although there have been considerable achievements in machine learning-based models, existing methods tend to overlook anomalies' special effects on real-world time series data. For instance, LSTMs [7] address the vanishing gradient problem via gate mechanism and have the ability to capture complex temporal dependencies [8]. Yet, Khandelwal et al. [9] show that even LSTMs have a limited ability to capture long-term dependencies, and their awareness of context degrades as the length of the input sequence increases. Consequently, making them inefficient to capture and learn from rare occurrences or extreme events.

As an alternative, Li et al. [10] considered the use of transformers for time series prediction. Transformers use a self-attention mechanism that allows each observation in the feature sequence to attend independently to every other feature in the sequence. However, they have considerable computational and memory requirements that grow quadratically with respect to sequence length, making it computationally rigorous to train large-scale data [10]. Such deficiency makes them computationally unsuitable for extreme events that often appear in longer sequences than the transformer's inputs. Moreover, it was not even clear from the design itself that transformers can be as effective as RNNs, whereas Zaheer et al. [11] reported that the attention mechanism in transformers does not even obey the sequence order of time steps which is essential for the time series domain. What is more, non-transformer architectures (i.e., MLP) when designed and trained properly, can perform competitively with transformers [12].

This lack of systematic strategy in handling anomalies and not providing predictions with non-transparent uncertainty levels makes current forecasting methods unreliable during the presence of extreme events. As a result, a key aspect of our knowledge in developing time series models for critical moments of extreme events will remain a puzzle unless the long-term effects of anomalies are well captured and utilized.

**Contribution.** This work proposes a novel and generalized anomaly-aware prediction framework, AA-Forecast, which automatically extracts and uses anomalies to optimize its probabilistic forecasting. Specifically,

- AA-Forecast extracts anomalies through a novel decomposition method and leverages them through an attention mechanism designed to optimize its probabilistic forecasting during extreme events. Also, AA-Forecast is able to perform zero-shot prediction for unseen time series and does not suffer from quadratic computational time and memory complexity of transformers.
- An online optimization procedure is proposed to minimize the prediction uncertainties of AA-Forecast framework, which features applying the optimal dropout probability at every time step during testing.
- Extensive experimental studies are conducted on three real-world datasets prone to extreme events and anomalies. The comparisons with state-of-the-art models illustrate the higher accuracy and less uncertainty in the AA-Forecast's prediction.

## 2  Problem Formulation

In this study, we are interested in the task of time series forecast under the influence of extreme events and anomalies. Mathematically, given a dataset $\mathbf{D} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \ldots, \mathbf{x}^{(K)}\}$ with $K$ univariate time series, $\mathbf{x}^{(k)} = \{x_1^{(k)}, x_2^{(k)}, \ldots, x_T^{(k)}\}$ denotes a time series instance with length $T$, where $\mathbf{x}^{(k)} \in \mathbb{R}^T$. For every time step, the corresponding extreme events are aligned and labeled as covariates $\mathbf{e}^{(k)} = \{e_1^{(k)}, e_2^{(k)}, \ldots, e_T^{(k)}\}$. Extreme events are considered as the influence of external events that promote a dynamic occurrence within a limited time steps [13]. Specifically, $e_t^{(k)} \in \mathbb{R}$ indicates the level of extreme event (e.g., hurricane category) at time $t$, otherwise, $e_t^{(k)} = 0$ indicates a non-extreme event condition for periods outside of the event. To this end, we denote the data with extreme events as a series of tuples $\widehat{\mathbf{x}}^{(k)} \triangleq \{(x_1^{(k)}, e_1^{(k)}), (x_2^{(k)}, e_2^{(k)}), \ldots, (x_T^{(k)}, e_T^{(k)})\}$. Particularly, given the previous $\tau$ observations $\widehat{\mathbf{x}}_{t-\tau+1:t}^{(k)} = \{(x_{t-\tau+1}^{(k)}, e_{t-\tau+1}^{(k)}), (x_{t-\tau+2}^{(k)}, e_{t-\tau+2}^{(k)}), \ldots, (x_t^{(k)}, e_t^{(k)})\}$, we aim to model the conditional distribution of the next observation:

$$p(x_{t+1}^{(k)} \,|\, \widehat{\mathbf{x}}_{t-\tau+1:t}^{(k)}; \mathbf{\Phi}), \tag{1}$$
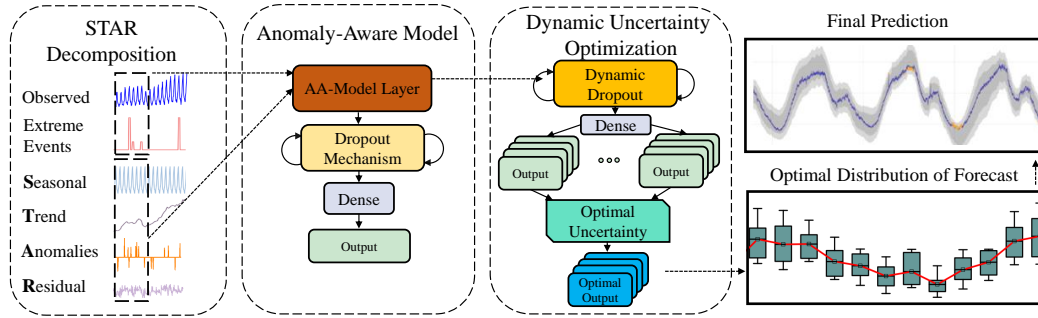
2

Figure 1: Main components of AA-Forecast: (i) **STAR Decomposition** to automatically extract essential features such as anomalies, (ii) an **Anomaly-Aware Model** to leverage such extracted features, and (iii) a **Dynamic Uncertainty Optimization** to reduce the uncertainty of the network. The final predicted series contains confidence intervals with the least uncertainty.

where $\mathbf{\Phi}$ denotes the parameters of a nonlinear prediction model. We are also interested in reducing the uncertainty of predictions in an online setting, whereas uncertainty of prediction can be viewed as the variability of the distribution. Therefore, the optimization problem during the online settings is defined as:

$$\mathbf{\Phi}_{\text{on}}^* = \text{argmin}_{\mathbf{\Phi}} \ \mathcal{V} \left( p(x_{t+1}^{(k)} \mid \widehat{\mathbf{x}}_{t-\tau+1:t}^{(k)}; \mathbf{\Phi}) \right), \tag{2}$$

where $\mathcal{V}(\cdot)$ represents the variability of the probability distribution and $\mathbf{\Phi}_{\text{on}}^*$ is the optimal online parameters of the nonlinear prediction model that produces the least amount of uncertainty in each time step.

## 3 AA-Forecast Framework

The proposed AA-Forecast framework consists of three main components. Section 3.1 proposes a novel anomaly decomposition method that automatically extracts the anomalies and essential features of the time series data. Then, the extracted anomalies are fed into an anomaly-aware model detailed in Section 3.2. Specifically, it leverages an attention mechanism on anomalies and extreme events to produce the distribution of the forecasts. To further reduce the forecast uncertainty in an online manner, Section 3.3 proposes a dynamic uncertainty optimization algorithm.

### 3.1 STAR Decomposition

STAR decomposition is used as a strategy to not only extract the anomalies and sudden changes of data but also decompose the complex time series to its essential components. Unfortunately, widely popular decomposition method such as STL [14] does not extract anomalies. Although recent works such as STR [15] and RobustSTL [16] are designed to be robust to the extreme effect of anomalies in their decomposition, they are not used to explicitly extract anomalies from the residual component.

To alleviate these issues, we propose STAR decomposition that are decomposes the original time series $\mathbf{x}^{(k)}$ in a multiplicative manner to its **s**easonal ($\mathbf{s}^{(k)}$), **t**rend ($\mathbf{t}^{(k)}$), **a**nomalies ($\mathbf{a}^{(k)}$), and **r**esidual ($\mathbf{r}^{(k)}$) components:

$$\mathbf{x}^{(k)} = \mathbf{s}^{(k)} \times \mathbf{t}^{(k)} \times \mathbf{a}^{(k)} \times \mathbf{r}^{(k)} \tag{3}$$

Such decomposition not only increases the dimensions of the original data but also supplies us with the extracted anomalies. As shown in Figure 1, we begin the decomposition by approximating the trend line $\mathbf{t}^{(k)}$ with the locally weighted scatterplot smoothing (i.e., LOESS [17]). Then, we divide the original data $\mathbf{x}^{(k)}$ by the approximated trend line to derive the detrended time series[1].

---

[1] We use the log transform of $\mathbf{x}^{(k)}$ to handle the situation that specific values of original data are zero.

We then partition the detrended time series into periods of cyclic sub-series where the cycle size is determined by the time interval of the dataset. As an example, the cycle size for a monthly dataset would be 12 (one year as a cycle). Then, we obtain the seasonal component ($\mathbf{s}^{(k)}$) by grouping the detrended series in each period and deriving the average value of each period across the time series. Subsequently, the residual component ($\mathbf{r}^{(k)}$) is derived by dividing the seasonal and trend segments from the original series.

Note that the anomaly component ($\mathbf{a}^{(k)}$) can be considered as the oddities of the dataset, which do not follow the extracted trend or seasonal components. Intuitively, anomalies are spread out through residual components, which also contain noise and other real-world effects. To distinguish the anomalies from residual components, statistical metrics such as mean and variance are not the appropriate measure as they are highly sensitive to the severity level of anomaly values. As one expects, the severity of the anomalies can change the mean and variance values which are unwanted. To resolve this issue, we leverage the median of the residuals, which is immune to the severity of the outliers in the residual components. Next, we define robustness score $\rho_t^{(k)}$ for each observation at time $t$ as:

$$\rho_t^{(k)} = \frac{|r_t^{(k)} - \dot{r}^{(k)}|}{\sqrt{\frac{\sum_{t=1}^{T} |r_t^{(k)} - \dot{r}^{(k)}|}{T-1}}} \tag{4}$$

where $\rho_t^{(k)}$ stands for the strength of the anomalies, $r_t^{(k)}$ is the residual at time step $t$ and $\dot{r}^{(k)}$ is the median of the residuals. Note that the larger $\rho_t$ indicates that a drastic change has occurred in the trend and seasonal components. We then extract the anomalies from residuals as below:

$$\mathbf{a}_t^{(k)} = \begin{cases} 1, & \rho_t^{(k)} < \rho_c^{(k)} \\ r_t^{(k)}, & \rho_t^{(k)} > \rho_c^{(k)} \end{cases} \tag{5}$$

where $\rho_c^{(k)}$ is the constant threshold given by the value of a robustness score ranked in the $p$-value $0.05^2$ while the values of elements in $\rho^{(k)}$ are ranked in descending order from large to small. Notably, when the value of the anomaly component ($\mathbf{a}^{(k)}$) deviates further from the value 1, it indicates an abrupt change in the trend and the seasonal component (no sign of anomalies). On the contrary, when both anomaly and residual values are equal 1 ($r_t^{(k)} = 1$ and $\mathbf{a}_t^{(k)} = 1$), it indicates that the observed signal at time $t$ explicitly follows the trend and seasonal component. Note that such important information might not be automatically inferred when additive decomposition methods are being used. This is due to the fact that the values of residual components can differ from one dataset to another which requires manual effort in their detection.

A sample result of anomaly decomposition is shown in the left-most part of Figure 1, where the observed time series data is decomposed into its seasonal, trend, anomalies, and residual components respectfully. Each of these components holds essential information about the characteristics of the time series and will be leveraged to train the forecast model. To this end, we concatenate the derived decomposed vector of time series with the input, which includes the observed time series and its labeled extreme event. Specifically, STAR decomposition concatenates the original time series to $\widetilde{\mathbf{x}}^{(k)} = (\mathbf{x}^{(k)}, \mathbf{e}^{(k)}, \mathbf{s}^{(k)}, \mathbf{t}^{(k)}, \mathbf{a}^{(k)}, \mathbf{r}^{(k)})$ which can be leveraged by the anomaly-aware model described in the next section.

## 3.2 Anomaly-Aware Model

The Anomaly-Aware model is designed to explicitly incorporate extracted anomalies $\mathbf{a}^{(k)}$ and extreme event covariates $\mathbf{e}^{(k)}$ into the prediction. As these features are rare in the whole time series, feeding them directly into a regular RNN like LSTM [7] can be potentially ignored during the training of the model. Note that the extracted anomalies and previously experienced external events hold valuable information regarding the effect of extreme events that should be handled carefully.

Recent robust prediction models rely on the LSTMs or transformers architecture to provide robustness in their prediction. Even though LSTMs are designed to obtain long-term dependencies, their

---

[2]Adopted based on the choice of the $p$-value (0.05) which is used as a standard level of statistical significance.

capacity to pay different degrees of attention to sub-window features within large time steps is inadequate [11]. As an example, Khandelwa et al. [9] showed that even though the LSTM model can have an effective sequence size of 200 observations, they are only able to sharply distinguish the 50 closest observations. This indicates that even LSTMs struggle to capture long-term dependencies. On the other hand, conventional transformers suffer from quadratic computation and memory requirements, which limits their ability to process long input sequences. Even though such memory bottlenecks have been improved by using sparse-attention algorithms [10], their performance improvement is not significant compared to a full-attention mechanism for real-world datasets [18]. Given that extreme events and anomalies are rare and can appear at very long distances from each other, it is computationally infeasible to increase the input sequence to provide attention to all previously seen anomalies and extreme events.

To address such problems, one must pay attention to all of the anomalies and extreme events throughout the dataset, no matter how far they have occurred. Intuitively, due to their rare nature, they hold greater importance in learning, given that the trend and seasonal patterns are often easier to predict by statistical or deep learning models. Ergo, we developed a novel attention mechanism explicitly for extreme events and anomalies, which are considered the crucial time steps of time series data and often cause the biggest error in prediction.

**Architecture Design of AA-Model.** LSTMs and GRUs are suitable for predicting the recurring patterns with a fairly low computational time and memory complexity which suffer from the quadratic complexity of full-attention transformers. However, we enhance the long-term dependencies of these models through an attention mechanism that retains the effect of anomalies and extreme events for future predictions. Such a decision in architecture allows the model not only to be computationally feasible for handling large-scale datasets but also to take the critical moments of extreme events and anomalies into consideration.

Given the past $\tau$ time steps of observations as $\tilde{\mathbf{x}}_{t:t-\tau+1}$[3], we derive the hidden states of an RNN that deals with vanishing gradient problem (e.g., LSTM or GRU) as:

$$\mathbf{h}_{t:t-\tau+1} = \text{RNN}\left(\tilde{\mathbf{x}}_{t:t-\tau+1}\right), \tag{6}$$

where $h_t$ is the hidden layer of RNN at time step $t$. Note that we are only giving attention to anomalies and extreme events which are naturally rare and belong to a small population of observations. Moreover, both could have different impacts on the prediction and based on the type of dataset, can be challenging to model. Hence, we design the attention mechanism to automatically incorporate extreme events and anomalies during their occurrence:

$$J = \{t \in \mathbb{Z}^+ | e_t \neq 0 \vee a_t \neq 1\}, \tag{7}$$

where $J$ is the set of time steps including two possible circumstances: presence of extreme events covariates ($e_t \neq 0$) or anomalies ($a_t \neq 1$). We then gather all the previous hidden states of the RNNs for all critical time steps in $J$ and regularize them by the weights generated from the attention layer as $v_t$ which follows:

$$v_t = \tanh(\mathbf{w}_\alpha^\top \mathbf{h}_t + b_\alpha), \quad \forall t \in J \tag{8}$$

where $\mathbf{w}_\alpha$ and $b_\alpha$ are the attention layer's weight and bias. Then, we derive the attention weights of all previous values as:

$$\alpha_t = \text{Softmax}\left(v_1, v_2, \ldots, v_t\right), \quad \forall t \in J \tag{9}$$

where $\alpha_t$ is the attention weight at the critical time steps. The generated attention weights are then used in the AA-Forecast layer as:

$$\mathcal{A}_t = \begin{cases} \mathbf{h}_t, & \forall t \notin J \\ \sum_{t \in J} \alpha_t \cdot \mathbf{h}_t, & \forall t \in J \end{cases} \tag{10}$$

where the attention values are only calculated in the presence of anomalies and extreme events as shown in Figure 2. The value of the next time step is calculated through a dense layer:

$$y_{t+1} = \mathbf{w}_d(\mathcal{A}_{t:t-\tau+1}) + b_d, \tag{11}$$

where $\mathbf{w}_d$ and $b_d$ are the weights and biases of the dense layer. To train the network, we minimize the prediction loss $\mathcal{L}$ which is defined as follows:

$$\mathbf{\Phi}_{\text{off}}^* = \text{argmin}_{\mathbf{\Phi}} \, \mathcal{L}\left(\mathcal{F}_{\mathbf{\Phi}}(\widetilde{\mathbf{x}}), y\right), \tag{12}$$

---

[3]To reduce the ambiguity of the AA-Forecast layer, we are omitting the superscript (k) from this section
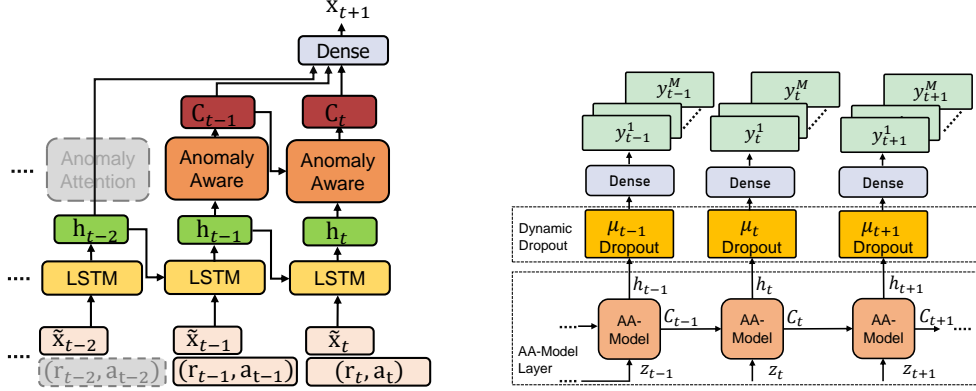
Figure 2: **Left:** AA-Model Architectures. **Right:** Dynamic dropout $\mu_t$ determines the optimal probability of dropout at each time step during the online settings (i.e., inference). The output $\hat{y}$ consists of a distribution of predicted test values. The dropout optimization improves the certainty and accuracy at each time step $t$ by determining how relevant the previous hidden state is for the next time step prediction.

where $\mathcal{F}_{\Phi}$ is the anomaly-aware model and $y$ is the training label, which is the ground truth of the next time step prediction. Note that $\Phi_{\text{off}}^*$ represents the optimal model parameters after the offline training phase.

## 3.3 Dynamic Uncertainty Optimization

Although Monte Carlo (MC) dropout [19] probability is treated as a static hyperparameter in previous studies [20, 21], it plays an important role in the prediction outcome and can be leveraged to reduce the uncertainty of the prediction during the testing phase [22]. Therefore, we rely on an automatic selection mechanism for the optimal dropout in online settings. Such selection is based on the uncertainty of the prediction produced during the testing phase (Figure 2).

Note that the model's uncertainty is desired to be the lowest and as stable as possible in real-world settings. Therefore, it is essential to optimize further the uncertainty of the model prediction both during the offline training and online testing phase. Specifically, we apply a dropout operation after every AA-Forecast layer with a specific probability ($p$).

AA-Forecast not only reports the prediction distribution but also provides the point prediction (average of the distribution) and the prediction uncertainty (variability of the distribution). Specifically, by producing $M$ forecast for every time step in an online manner (test data $\widetilde{\mathbf{x}}^*$) from the previously trained model $\mathcal{F}_{\Phi}(\widetilde{\mathbf{x}})$, we obtain $M$ outputs $y^*$ as a from the prediction distribution $\left\{ y_{(1)}^*, \ldots, y_{(M)}^* \right\}$. Then, the average of the distribution is calculated as $\bar{y}^* = \frac{1}{M} \sum_{m=1}^{M} y_{(m)}^*$.

We represent uncertainty as to the variability of the prediction distribution —- the standard deviation (SD) of the probability distribution of future observations conditional on the information available at the time of forecasting. We further optimize the uncertainty of the framework by deriving the optimal dropout probability $p$ at each time step. We derive the prediction error for the probability $p$ between 0 and 1 with 0.1 increments. Notably, without such probability (i.e., $p = 0$) the model prediction deviates from probabilistic forecasting and does not provide a level of uncertainty in its prediction for each time step. The optimal uncertainty $\mu_t$ is then reported when it results in a minimal variance (i.e., SD) of the predicted values, thereby reducing the prediction uncertainty to its minimum during the testing phase. To this end, the prediction uncertainty is formulated as:

$$\sigma^2 \left( \mathcal{F}_{\Phi}(\widetilde{\mathbf{x}}^*) \right) = \sqrt{ \frac{1}{M} \sum_{m=1}^{M} \left( y_{(m)}^* - \bar{y}^* \right)^2 }. \tag{13}$$

6

**Algorithm 1** Psuedecode for AA-Forecast

---

**Input:** data $\widetilde{\mathbf{x}}^{(k)} = (\mathbf{x}^{(k)}, \mathbf{e}^{(k)}, \mathbf{s}^{(k)}, \mathbf{t}^{(k)}, \mathbf{a}^{(k)}, \mathbf{r}^{(k)})$

1: Initialize parameters $\boldsymbol{\Phi}$
2: **for** $k = 1$ to $K_{\text{train}}$ **do**
3:     Sample $(\tilde{\mathbf{x}}^k, y^k)$ from training data:
4:     **for** $b = 1$ to $B$ **do**
5:             $\boldsymbol{\Phi}_{e+1} \leftarrow \boldsymbol{\Phi}_e$ - $\xi \cdot \nabla \, \mathcal{L}(\mathcal{F}_{\boldsymbol{\Phi}}(\tilde{\mathbf{x}}^k), y^k)$
6:         Update the optimal parameters:
             $\boldsymbol{\Phi} = \text{argmin}_{\boldsymbol{\Phi}} \mathcal{L}(\mathcal{F}_{\boldsymbol{\Phi}}(\tilde{\mathbf{x}}^k), y^k)$
7:     **end for**
8: **end for**
9: Dynamic Uncertainty optimization: $\boldsymbol{\Phi}^* \leftarrow \boldsymbol{\Phi}$
10: **for** $\delta = 0.1$ to $0.9$ increment by $0.1$ **do**
11:     Update the optimal uncertainty:
         $\boldsymbol{\Phi}^* = \text{argmin}_{\boldsymbol{\Phi}} \mathcal{V}(\mathcal{F}_{\boldsymbol{\Phi}}(x^k))$
12: **end for**

---

Algorithm 1 presents the pseudocode for AA-Forecast. Specifically, we sample $(\tilde{\mathbf{x}}^k, y^k)$ as a driving example which includes extracted anomalies $\mathbf{a}^{(k)}$ and extreme events $\mathbf{r}^{(k)}$. Next, we train the model by maximizing the overall prediction accuracy. Upon testing, the network leverage dynamic uncertainty optimization further optimizes the prediction uncertainty automatically in online testing so that it would not require any further training.

Note that the network's predictions during the testing phase cannot benefit from the supervised training. However, the control of variability is possible and ensures that the prediction uncertainty is minimal in each step of future predictions, regardless of whether the labels are provided or not. Additionally, the algorithm testing time complexity is similar to other RNN-based models due to the use of dynamic uncertainty optimization during the test phase solely. This allows the model to provide the least amount of uncertainty during the presence of anomalies or extreme events where critical online decisions are being made.
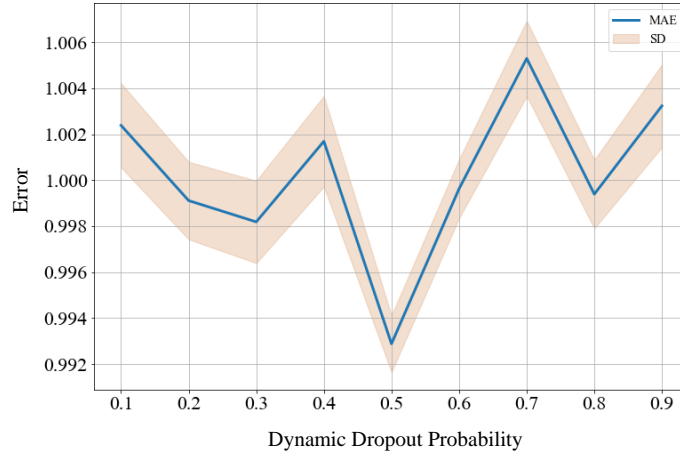


Figure 3: Effects of dynamic uncertainty optimization on prediction error and uncertainty during the occurrence of an anomaly. The method automatically selects the optimal probability that yields the lowest uncertainty.

As an example, Figure 3 shows that the optimal uncertainty results can occur when the standard deviation is the lowest. Intuitively, the network at $p = 0.5$ shows the highest confidence in its prediction (i.e., the lowest uncertainty) where unnecessary neurons are dropped out from the network. Therefore, the network automatically selects and reports the $p = 0.5$ probability as the best choice for this time step in the testing phase.

# 4 Experiments

This section reports multiple experiments comparing the proposed AA-Forecast framework with baseline models using different types of large-scale time series datasets.

## 4.1 Dataset and Experimental Settings

Three real-world time series with diverse structures and domains are gathered[4]. The detailed description and data collection procedure are as follows:

- We provide a new spatio-temporal benchmark dataset (`Hurricane`), which is suited for forecasting during extreme events and anomalies. The dataset is gathered through the Florida Department of Revenue which provides the monthly sales revenue (2003-2020) for the tourism industry in the 67 counties of Florida which are prone to annual hurricanes. We further enriched and aligned the raw time series with the history of hurricane categories for each region upon impact. More precisely, the hurricane category indicates the maximum sustained wind speed which can result in catastrophic damages [23].

- The second dataset (`COVID-19`) showcases the changes in the number of employees based on one million employees active in the US during the COVID-19 pandemic and is gathered from Homebase [24]. We further enriched the data with the state-level policies as an indication of extreme events (e.g., the state's business closure order).

- The third dataset (`Electricity`) is a publicly available benchmark dataset that contains the electricity consumption of 370 consumers on an hourly basis from 2011 to 2014. Note that this benchmark dataset is anonymized and does not contain extreme event labels, yet AA-Forecast is able to automatically extract the anomalies, indicating abrupt changes in trend and seasonality.

Table 1: Summary statistics of the datasets.

| Dataset | Hurricane | COVID-19 | Electricity |
|---|---|---|---|
| Time step | Monthly | Daily | Hourly |
| # Unique time series | 9,876 | 15,312 | 370 |
| # Observation | 9,876 | 15,312 | 11,952,480 |
| # Train | 7,900 | 12,250 | 9,561,984 |
| # Test | 1,975 | 3,062 | 2,390,496 |
| # Regions | 48 | 50 | 370 |
| # Extreme events | 88 | 100 | - |
| # Anomalous points | 102 | 124 | 672 |

We propose two sets of experiments for all baseline models. The first experiment follows a standard 80-20 dividing of the dataset to training and testing sets and $\tau = 12$ for window length. The second experiment evaluates the zero-shot prediction capability of the model based on various window search ranges in $\{3, 6, 12, 24\}$, and thus is more applicable for real-world settings when the newly added time series cannot train on a newly added time series. Hence, the second experiment evaluates the prediction accuracy of all models on a set of completely unseen time series.

The models are implemented using Python 3.7 and tested on a cloud workstation with two Intel Xeon 2.3 GHz CPUs, 64 GB RAM, and one Nvidia Tesla A100 GPU. We conduct a grid search over all tunable hyperparameters on the held-out validation set for baseline methods and our framework. To provide a fair evaluation, all baseline models benefit from the essential features extracted by AA-Forecast except the ARIMA model which does not benefit from multidimensional features. Moreover, future known information is not included in all the models.

We kept training to 40 iterations for all experiments. The reported values are the average of the observed error five times during the test stage. The hyperparameters of all baseline methods are tuned based on a grid search.

---

[4]All datasets are publicly available at https://github.com/ashfarhangi/AA-Forecast

Table 2: Hyperparameters of AA-Forecast used for each dataset.

| Parameter | Hurricane | COVID-19 | Electricity |
|---|---|---|---|
| Batch size | 128 | 64 | 64 |
| Learning rate | $1 \times 10^{-5}$ | $3 \times 10^{-5}$ | $5 \times 10^{-5}$ |
| Weight decay | $1 \times 10^{-6}$ | $1 \times 10^{-5}$ | $1 \times 10^{-4}$ |
| Number of epochs | 40 | 40 | 40 |
| Static dropout | 0.5 | 0.4 | 0.6 |

## 4.2 Methods for Comparison

The baseline methods for comparison include:

- ARIMA [25]: A traditional autoregressive integrated moving average method for time series prediction and often used as a baseline.

- AE-LSTM [26]: An LSTM network that uses an autoencoder for deep feature extraction and provides a deterministic prediction.

- SARIMAX [27]: An autoregressive model that can handle seasonality and exogenous features of time series.

- UberNN [8]: An LSTM-based model that uses Monte Carlo dropout to provide uncertainty and is able to extract deep features of time series through autoencoders.

- TSE-SC [28]: was recently proposed as a Transformer-based Deep Learning model that can forecast abrupt changes accurately. (i) STAR Decomposition to automatically ex- tract essential features such as anomalies, (ii) an Anomaly-Aware Model to leverage such extracted features, and (iii) a Dynamic Uncertainty Optimization to reduce the uncertainty of the network. The final predicted

- AA-Forecast (LSTM) is our proposed model with LSTM cells.

- AA-Forecast (GRU) is our proposed model with GRU cells.

## 4.3 Metrics

For providing a comprehensive evaluation, we adopted three different evaluation metrics. The first evaluation metric is the Continuous Ranked Probability Score (CRPS), which evaluates probabilistic forecasting. Formally defined as $\mathbf{CRPS} = \int_{-\infty}^{\infty} (F(y) - \mathbb{1}(y - \hat{y}))^2 \, \mathrm{d}y$ where $F$ is the cumulative distribution function of its forecast distribution and $\mathbb{1}$ is the Heaviside step function. We also report the root mean square error (RMSE). Formally defined as $\mathbf{RMSE} = \sqrt{\frac{1}{N} \sum_{i=1}^{N} \left( y_{t,(i)} - \hat{y}_{t,(i)} \right)^2}$ where $y_t$ is the mean of the predicted distribution at time $t$ and $\hat{y}_t$ is the observed value at time $t$. The third evaluation metric is the standard deviation (SD) that is correlated to the uncertainty of the prediction and is denoted as $\mathbf{SD} = \sqrt{\frac{1}{N} \sum_{i=1}^{N} \left( y_{t,(i)} - \tilde{y}_t \right)^2}$ where $\bar{y}$ is the mean of the predicted distribution.

## 4.4 Experimental Results

We provide two comprehensive comparisons and evaluations of the proposed AA-Forecast framework: the aforementioned 80-20 testing where 20% of the data are unseen, as well as the testing on zero-shot prediction where the whole time series is unseen. In both cases, we calculate the CRPS, RMSE, and SD. Lastly, we provided an ablation study to discuss the effectiveness of different AA-Forecast components.

**The $80 - 20$ testing.** We first used the 'older' 80% of each time series in training and tested the accuracy of prediction on the rest of 20%. Table 3 reports the loss of the networks under such $80 - 20$ testing, where the SD of AA-Forecast (GRU) method is lower than all baseline methods, showing the model's high confidence in the forecasts.

Among the baseline methods, UberNN and TSE-SC have shown good accuracy but suffer from higher SD (uncertainty) compared to the AA-Forecast (LSTM-GRU) models. Considering that the

Table 3: Performance comparison of our proposed framework and baseline models under $80-20$ testing.

| Methods | Metrics | Dataset | | |
|---|---|---|---|---|
| | | Electricity | COVID-19 | Hurricane |
| ARIMA [25] | CRPS | 1.150 | 0.103 | 0.761 |
| | RMSE | 1.520 | 0.114 | 0.802 |
| | SD | 0.225 | 0.011 | 0.106 |
| AE-LSTM [26] | CRPS | 0.895 | 0.086 | 0.531 |
| | RMSE | 1.296 | 0.087 | 0.576 |
| | SD | 0.215 | 0.009 | 0.102 |
| SARIMAX [27] | CRPS | 0.911 | 0.098 | 0.532 |
| | RMSE | 1.285 | 0.108 | 0.578 |
| | SD | 0.195 | 0.009 | 0.093 |
| UberNN [8] | CRPS | 0.633 | 0.071 | 0.442 |
| | RMSE | 1.015 | 0.081 | 0.453 |
| | SD | 0.134 | 0.007 | 0.073 |
| TSE-SC [28] | CRPS | 0.583 | 0.062 | 0.384 |
| | RMSE | 0.983 | 0.072 | 0.423 |
| | SD | 0.146 | 0.007 | 0.092 |
| **AA-Forecast** | CRPS | 0.546 | **0.059** | 0.237 |
| **(LSTM)** | RMSE | 0.949 | **0.068** | 0.274 |
| | SD | 0.095 | **0.003** | 0.060 |
| **AA-Forecast** | CRPS | **0.493** | 0.063 | **0.216** |
| **(GRU)** | RMSE | **0.894** | 0.073 | **0.253** |
| | SD | **0.081** | 0.003 | **0.051** |

extracted features are available for all the baseline methods, we believe the higher uncertainty of SD is due to their static dropout probability that is constant for all time steps. Therefore, the two proposed models, AA-Forecast (LSTM-GRU), consistently outperform state-of-the-art methods. Considering all three evaluation metrics, AA-Forecast (GRU) is the best-suited framework for our dataset as it provides higher accuracy and confidence.

**Zero-shot prediction.** Table 4 demonstrates the zero-shot prediction abilities for the selected models. Both AA-Forecast (LSTM-GRU) predictions follow the observed time series in general. The prediction errors are comparably low during the presence of extreme events (i.e., hurricanes). This is mainly due to the anomaly attention mechanism developed to further reduce the prediction error during extreme events. Moreover, extracted anomalies from STAR decomposition led to the recall of the hurricane effects on previously seen regions, thus providing predictions for unseen time series data with a lower error given the presence of anomalies. Figure 4 showcases a sample of these predictions for each model where for every time step, the prediction uncertainty is the least.
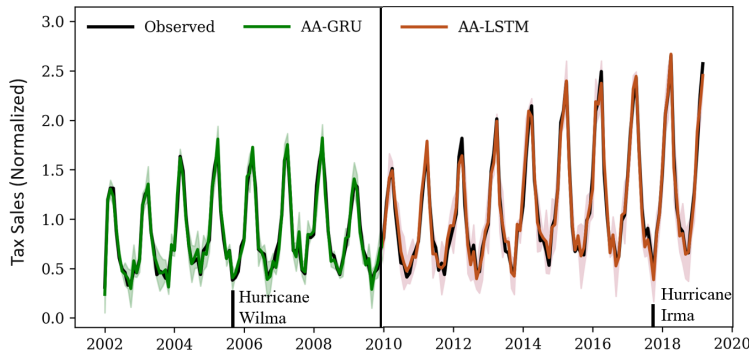


Figure 4: Zero-shot prediction for hotel tax sales of Collier County, Florida, U.S. Both variations of AA-Forecast are concatenated for demonstration.

Table 4: Performance comparisons of zero-shot prediction abilities of models using ten randomly selected counties' sales tax data where they have not been used in training entirely.

| Methods | Metrics | Input time window | | | |
|---|---|---|---|---|---|
| | | 3 | 6 | 12 | 24 |
| ARIMA [25] | CRPS | 0.893 | 0.891 | 0.861 | 0.831 |
| | RMSE | 0.934 | 0.932 | 0.922 | 0.872 |
| | SD | 0.119 | 0.1154 | 0.115 | 0.113 |
| AE-LSTM [26] | CRPS | 0.663 | 0.661 | 0.651 | 0.601 |
| | RMSE | 0.708 | 0.706 | 0.696 | 0.646 |
| | SD | 0.115 | 0.112 | 0.111 | 0.109 |
| SARIMAX [27] | CRPS | 0.664 | 0.662 | 0.662 | 0.602 |
| | RMSE | 0.714 | 0.712 | 0.712 | 0.652 |
| | SD | 0.106 | 0.102 | 0.102 | 0.100 |
| UberNN [8] | CRPS | 0.547 | 0.545 | 0.535 | 0.485 |
| | RMSE | 0.585 | 0.583 | 0.573 | 0.523 |
| | SD | 0.086 | 0.082 | 0.082 | 0.08 |
| TSE-SC [28] | CRPS | 0.766 | 0.764 | 0.754 | 0.704 |
| | RMSE | 0.795 | 0.793 | 0.783 | 0.733 |
| | SD | 0.105 | 0.102 | 0.101 | 0.099 |
| **AA-Forecast** | CRPS | 0.362 | 0.361 | 0.351 | 0.301 |
| **(LSTM)** | RMSE | 0.406 | 0.404 | 0.394 | 0.344 |
| | SD | 0.073 | 0.071 | 0.069 | 0.067 |
| **AA-Forecast** | CRPS | **0.348** | **0.346** | **0.336** | **0.286** |
| **(GRU)** | RMSE | **0.385** | **0.383** | **0.373** | **0.323** |
| | SD | **0.064** | **0.060** | **0.062** | **0.058** |

Given that the network did not train on the selected time series directly, it's able to transfer its knowledge from previously seen extreme events (i.e., the effect of cat 4 hurricanes) and provide more accurate prediction when not provided with such ability.

## 4.5 Ablation Study

In this section, we provide an extensive analysis of the performance of AA-Forecast, as well as the impact of each component on the performance of AA-Forecast. The results are shown in Table 5 where we removed each component and reported the changes in accuracy and uncertainty.

**Influence of anomaly-aware decomposition.** To demonstrate that the anomaly-aware decomposition can aid in improving the time series prediction, we fed the input series to the prediction model directly. This modification results in the worst performance in our ablation study. Note that AA-Forecast (GRU) still benefits from dynamic dropout optimization and extreme event labels, and the predicted uncertainty is optimized. However, the accuracy of AA-Forecast prediction (GRU) drops because of the limited number of features, indicating that the neural network does not have a strong ability to capture complex and nonlinear information. This can highlight the role of auxiliary features such as decomposed anomalies and extreme events for forecasting.

**Influence of uncertainty optimization.** We also used a static dropout throughout the experiments at every time step, which caused a substantial increase in SD. Uncertainty optimization of dropout plays a critical role in reducing the uncertainty of the forecast intervals. Such modification also caused a higher error in the forecast, which is the model's inability to forecast with higher confidence.

**Influence of anomaly attention.** We conducted experiments to demonstrate the effectiveness of anomaly-awareness through the network's attention mechanism. Specifically, we directly fed the extreme events and anomalies without the anomaly-attention mechanism described in Section 3.2. Such change makes limits AA-Forecast's knowledge about hurricanes and the severity of their effects. As an example, in Figure 5 (right), the results show that the network's error during the presence of harder-to-predict time points (anomalies and extreme events) weakens. Thus, removing the attention mechanism for anomalous/extreme event points of the dataset will reduce the performance of the model during the critical months of extreme events such as hurricanes. Simply relying on

Table 5: Ablation study on AA-Forecast (GRU) model using the sales tax dataset to show the effectiveness of its components.

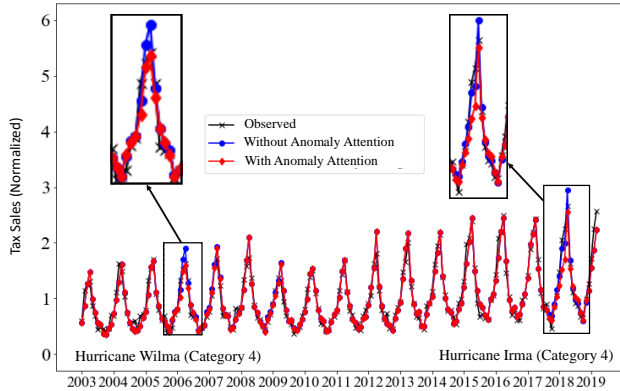| AA-Forecast (GRU) | Metrics | Time window | | | |
|---|---|---|---|---|---|
| | | 3 | 6 | 12 | 24 |
| w/o STAR Decomposition | CRPS | 0.493 | 0.446 | 0.445 | 0.457 |
| | RMSE | 0.512 | 0.464 | 0.463 | 0.494 |
| | SD | 0.074 | 0.071 | 0.070 | 0.070 |
| w/o Uncertainty Optimization | CRPS | 0.429 | 0.431 | 0.43 | 0.367 |
| | RMSE | 0.466 | 0.471 | 0.467 | 0.404 |
| | SD | 0.088 | 0.088 | 0.087 | 0.083 |
| w/o Anomaly Attention | CRPS | 0.379 | 0.380 | 0.367 | 0.317 |
| | RMSE | 0.416 | 0.417 | 0.404 | 0.354 |
| | SD | 0.067 | 0.067 | 0.063 | 0.061 |
| AA-Forecast (GRU) | CRPS | **0.348** | **0.346** | **0.336** | **0.286** |
| | RMSE | **0.385** | **0.383** | **0.373** | **0.323** |
| | SD | **0.064** | **0.060** | **0.060** | **0.058** |



Figure 5: Influence of anomaly attention on hurricanes. Two Category 4 hurricanes (Wilma and Irma) have caused similar annual sales losses. Anomaly-attention activation occurs during the presence of extreme events which makes it computationally efficient compared to the full-attention mechanism in transformers.

the previously seen dataset will not allow the network to handle external events and sudden changes effectively.

## 4.6 Discussion

**Interpretation.** The benefits of providing optimal uncertainty in prediction are twofold: first, it provides a systematic way to aid in resource allocation. Second, it further prepares the domain for interventions. For example, if one region receives more catastrophic extreme events, the resources can be transferred to that region. Moreover, government and industries can provide better-informed interventions and decisions (e.g., financial aid relief during COVID-19). As shown in the ablation study, including additional features such as extreme events and anomalous points can improve accuracy and better prime the model to handle predictions than deviate from trend or seasonality. Moreover, as shown in Figure 5 without proper attention to these points, they result in a large amount of error in forecasting. Given that such critical moments are of high importance during extreme events such as natural disasters, the performance of the model during critical time steps can be improved. Hence, it is essential to provide a thorough learning objective in our time series models to not only improve the overall performance but take critical moments into more consideration. Furthermore, allowing the model to provide its level of uncertainty establishes transparency and builds a level of trust for the users.

**Limitations & future directions.** Although the dynamic dropout mechanism guarantees the least uncertainty in predictions, it cannot provide guarantees to do the same for prediction accuracy.

This is due to the randomness nature of the dropout which we left as a future work where the dropout can appear for a predetermined distribution of the neurons. Therefore, maximizing the useful information contained in the multidimensional model serves to predict time series in extreme events. When it is not available, it's more reasonable to suggest methods that extract potential critical time steps such as anomalous points (e.g., STAR decomposition).

## 5 Related Works

Anomalies in time series data often produce a high variance of uncertainty prediction that is difficult to predict, thus becoming a challenge for reliable model design [8, 29]. To provide a more reliable forecast during the presence of anomalies, probabilistic forecasting methods are often studied, which can report a level of uncertainty [30].

The majority of Bayesian Neural Networks in probabilistic forecasting requires specific training and optimization methods and require additional model parameters that result in a larger amount of computation. Hence, MC dropout is preferred due to its practicability and its out-of-the-box solution [8]. Applying standard dropout to Bayesian Neural Networks often results in poor performance on account of dropout noise preventing the network from maintaining long-term memory [31]. Gal and Ghahramani [19] proposed the MC dropout, in which the dropout can be interpreted as a sampling method that is equivalent to a variational approximation of a deep Gaussian process. MC dropout that is used for recurrent layers has proved to be successful and is commonly used in practice by applying dropout to recurrent connections in a way that can preserve long-term memory [31]. In previous studies, static MC dropout was used throughout their experiments, which suffers the model's robustness toward the effect of anomalies. Given that probabilistic models still require an overall great accuracy of their forecasts, optimizing the uncertainty in prediction intervals remains a challenging question

## 6 Conclusion

We propose an anomaly-aware time series prediction framework, namely AA-Forecast, to capture and leverage the effect of extreme events and anomalies for the time series prediction task. It features a novel anomaly decomposition method that also extracts the essential features of the data. We also proposed an anomaly-aware model to leverage the extracted anomalies through an attention mechanism. Moreover, we reduced the uncertainty of the network without any further training so that the prediction uncertainty is minimal through the testing state. We compare our framework with several statistical and deep learning models using three real-world time series datasets. The results show that the AA-Forecast framework outperforms these models in prediction error and uncertainty. For future work, the prediction performance could be further improved if we target specific groups of neurons (e.g., the neurons containing unnecessary details of the time series dynamics) for dynamic dropout optimization.

## References

[1] Adam B Smith. 2021 us billion dollar weather and climate disasters in historical context including new county-level exposure, vulnerability and projected damage mapping. In *102nd American Meteorological Society Annual Meeting*. AMS, 2022.

[2] Min Jing, Kok Yew Ng, Brian Mac Namee, Pardis Biglarbeigi, Rob Brisk, Raymond Bond, Dewar Finlay, and James McLaughlin. Covid-19 modelling by time-varying transmission rate associated with mobility trend of driving via apple maps. *Journal of Biomedical Informatics*, 122:103905, 2021.

[3] Carlos Santos-Burgoa, John Sandberg, Erick Suárez, Ann Goldman-Hawes, Scott Zeger, Alejandra Garcia-Meza, Cynthia M Pérez, Noel Estrada-Merly, Uriyoan Colón-Ramos, Cruz María Nazario, et al. Differential and persistent risk of excess mortality from hurricane maria in puerto rico: a time-series analysis. *The Lancet Planetary Health*, 2(11):e478–e488, 2018.

[4] Asif Khan, Sughra Bibi, Jiaying Lyu, Abdul Latif, and Ardito Lorenzo. Covid-19 and sectoral employment trends: assessing resilience in the us leisure and hospitality industry. *Current Issues in Tourism*, 24(7):952–969, 2021.

[5] Linara Adilova, Siming Chen, and Michael Kamp. Novelty detection in sequential data by informed clustering and modeling. *arXiv preprint arXiv:2103.03943*, 2021.

[6] David Rolnick, Priya L Donti, Lynn H Kaack, Kelly Kochanski, Alexandre Lacoste, Kris Sankaran, Andrew Slavin Ross, Nikola Milojevic-Dupont, Natasha Jaques, Anna Waldman-Brown, et al. Tackling climate change with machine learning. *arXiv preprint arXiv:1906.05433*, 2019.

[7] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.

[8] Lingxue Zhu and Nikolay Laptev. Deep and confident prediction for time series at uber. In *2017 IEEE International Conference on Data Mining Workshops (ICDMW)*, pages 103–110. IEEE, 2017.

[9] Urvashi Khandelwal, He He, Peng Qi, and Dan Jurafsky. Sharp nearby, fuzzy far away: How neural language models use context. *arXiv preprint arXiv:1805.04623*, 2018.

[10] Shiyang Li, Xiaoyong Jin, Yao Xuan, Xiyou Zhou, Wenhu Chen, Yu-Xiang Wang, and Xifeng Yan. Enhancing the locality and breaking the memory bottleneck of transformer on time series forecasting. In *Advances in Neural Information Processing Systems*, pages 5243–5253, 2019.

[11] Manzil Zaheer, Guru Guruganesh, Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. Big bird: Transformers for longer sequences. *arXiv preprint arXiv:2007.14062*, 2020.

[12] Ilya O Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas Unterthiner, Jessica Yung, Andreas Steiner, Daniel Keysers, Jakob Uszkoreit, et al. Mlp-mixer: An all-mlp architecture for vision. *Advances in Neural Information Processing Systems*, 34:24261–24272, 2021.

[13] Lisa Hanna Broska, Witold-Roger Poganietz, and Stefan Vögele. Extreme events defined—a conceptual discussion applying a complex systems approach. *Futures*, 115:102490, 2020.

[14] Robert B Cleveland, William S Cleveland, Jean E McRae, and Irma Terpenning. Stl: A seasonal-trend decomposition. *Journal of official statistics*, 6(1):3–73, 1990.

[15] Alexander Dokumentov and Rob J Hyndman. Str: A seasonal-trend decomposition procedure based on regression. *arXiv preprint arXiv:2009.05894*, 2020.

[16] Qingsong Wen, Jingkun Gao, Xiaomin Song, Liang Sun, Huan Xu, and Shenghuo Zhu. Robuststl: A robust seasonal-trend decomposition algorithm for long time series. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 5409–5416, 2019.

[17] William S Cleveland. Robust locally weighted regression and smoothing scatterplots. *Journal of the American statistical association*, 74(368):829–836, 1979.

[18] Bryan Lim, Sercan O Arik, Nicolas Loeff, and Tomas Pfister. Temporal fusion transformers for interpretable multi-horizon time series forecasting. *arXiv preprint arXiv:1912.09363*, 2019.

[19] Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059, 2016.

[20] David Salinas, Valentin Flunkert, Jan Gasthaus, and Tim Januschowski. Deepar: Probabilistic forecasting with autoregressive recurrent networks. *International Journal of Forecasting*, 36(3):1181–1191, 2020.

[21] Nikolay Laptev, Jason Yosinski, Li Erran Li, and Slawek Smyl. Time-series extreme event forecasting with neural networks at uber. In *International Conference on Machine Learning*, volume 34, pages 1–5, 2017.

[22] Hud Wahab, Vivek Jain, Alexander Scott Tyrrell, Michael Alan Seas, Lars Kotthoff, and Patrick Alfred Johnson. Machine-learning-assisted fabrication: Bayesian optimization of laser-induced graphene patterning using in-situ raman analysis. *Carbon*, 167:609–619, 2020.

[23] Nhc data archive. `https://www.nhc.noaa.gov/data/`. (Accessed on 08/20/2022).

[24] Alexander W Bartik, Marianne Bertrand, Feng Lin, Jesse Rothstein, and Matt Unrath. Labor market impacts of covid-19 on hourly workers in small-and medium-sized businesses: Four facts from homebase data. *Institute for Research on Labor and Employment*, 2020.

[25] George EP Box and David A Pierce. Distribution of residual autocorrelations in autoregressive-integrated moving average time series models. *Journal of the American statistical Association*, 65(332):1509–1526, 1970.

[26] Alaa Sagheer and Mostafa Kotb. Unsupervised pre-training of a deep lstm-based stacked autoencoder for multivariate time series forecasting problems. *Scientific Reports*, 9(1):1–16, 2019.

[27] Agostino Tarsitano and Ilaria L Amerise. Short-term load forecasting using a two-stage sarimax model. *Energy*, 133:108–114, 2017.

[28] Ling Cai, Krzysztof Janowicz, Gengchen Mai, Bo Yan, and Rui Zhu. Traffic transformer: Capturing the continuity and periodicity of time series for traffic forecasting. *Transactions in GIS*, 24(3):736–755, 2020.

[29] Jingyue Pang, Datong Liu, Yu Peng, and Xiyuan Peng. Anomaly detection based on uncertainty fusion for univariate monitoring series. *Measurement*, 95:280–292, 2017.

[30] Longyuan Li, Junchi Yan, Xiaokang Yang, and Yaohui Jin. Learning interpretable deep state space model for probabilistic time series forecasting. In *IJCAI*, pages 2901–2908, 2019.

[31] Alex Labach, Hojjat Salehinejad, and Shahrokh Valaee. Survey of dropout methods for deep neural networks. *arXiv preprint arXiv:1904.13310*, 2019.