

# Forecast of Energy Demand Using Temporal Fusion Transformer

Chandreyi Chowdhury  
School of Advanced Sciences  
Vellore Institute of Technology  
Vellore, India  
chandreyi.surat@gmail.com

Dr. Neelabja Chatterjee  
School of Advanced Sciences  
Vellore Institute of Technology  
Vellore, India  
neelabja.c@vit.ac.in

**Abstract**— This research paper presents an application of Temporal Fusion Transformer (TFT) model for time series forecasting using deep learning techniques. The novelty of TFT model lies in combining the strengths of recurrent and convolutional neural networks, ultimately resulting in improved accuracy in forecasting. The dataset used in this study consists of data on power consumption per quarter hour of 370\$ consumers, over a four-year period. The results show that the TFT model outperformed traditional time series models and achieved a lower Mean Absolute Error (MAE) and a lower Root Mean Squared Error (RMSE) in predicting future energy consumption. The paper indicates that the TFT model can be an effective tool for accurate and reliable time series forecasting in various industries, including energy and finance.

**Keywords**— Temporal Fusion Transformer (TFT), Time Series, Forecasting, Energy Consumption

## I. INTRODUCTION

The accurate prediction of time series data is a critical endeavour, ubiquitous in various fields. Traditional time series forecasting models often struggle to capture the complexities of real-world data, resulting in inaccurate predictions. In recent years, models using deep learning methods have emerged as a promising tool to solve this issue for forecasting time series data. One such model is the Temporal Fusion Transformer (TFT), designed explicitly for time series forecasting with interpretability. The TFT model was proposed in a paper by Bryan Lim et al [1] in 2020 and combines the strengths of two popular deep learning architectures: the Transformer [2] and the fully connected neural network (FCNN) [3].

The Transformer architecture allows the TFT model to capture long-term dependencies and relationships between the inputs, while the FCNN is used to model the temporal dynamics of the time series. This combination results in more accurate predictions compared to traditional models. The TFT model also provides interpretability, which is one of the desired features in many industries. The model includes an attention mechanism that allows the user to understand which inputs are most important to make predictions. The *gating mechanism* further enables the user to selectively ignore certain inputs that may not be relevant for the current prediction.

In addition to its interpretability, the TFT model can handle missing values, making it well-suited for *incomplete* datasets. The model's flexibility and modular architecture also allow for customization and adaptation for various different applications.

In this paper, we apply the TFT model to forecast energy consumption data by a French household over a three-year period. Our evaluation of the TFT model against traditional time series models demonstrates its superior performance in

terms of Mean Absolute Scaled Error (MASE). Our results highlight the effectiveness of the TFT model for accurate and interpretable time series forecasting in energy consumption data and its potential application in other industries.

## II. RELATED WORKS

Deep learning has emerged as a powerful tool for time series forecasting, providing state-of-the-art results in a variety of domains. Rumelhart et al. presented the backpropagation algorithm which [4] has been instrumental in training deep neural networks (DNNs) to learn complex temporal dependencies. Hochreiter and Schmidhuber [5] proposed long short-term memory (LSTM) networks, which have become one of the most popular DNN architectures for time series forecasting due to their ability to capture long-term dependencies. However, training LSTMs can be challenging due to vanishing and exploding gradients [6].

Recent advances in deep learning have addressed some of these challenges, such as the DeepAR architecture proposed by Zhang et al. [7], which uses autoregressive recurrent networks to model the temporal dependencies of time series. Another approach is the Temporal Fusion Transformer (TFT) proposed by Lim et al. [1], which combines LSTM and attention mechanisms to model multiple horizons in a time series, providing interpretable multi-horizon forecasting.

Graph Convolutional Networks (GCNs) are another type of DNN that have been applied to time series forecasting. Kipf and Welling [8] introduced a semi-supervised classification method based on GCNs that can learn representations of both nodes and edges in a graph, which has been used for traffic flow prediction [15].

Regularization techniques have also been developed to improve the generalization performance of DNNs. Hinton et al. [9] proposed dropout as a way to prevent co-adaptation of feature detectors in a neural network, and LeCun et al. [10] introduced batch normalization to reduce internal covariate shift during training.

Deep Canonical Correlation Analysis (DCCA) is another method that has been applied to time series forecasting. Amini et al. [11] used DCCA to learn representations that maximize the correlation between two sets of variables, which has been shown to improve the accuracy of multi-step forecasting.

In addition to deep learning architectures, several studies have explored the use of text and numerical data for stock prediction [17] and spatio-temporal data mining [18]. For example, Dong et al. [20] applied deep recurrent neural networks to ultra-short-term wind forecasting, using both numerical and textual data to improve prediction accuracy.

Overall, the related works suggest that deep learning is a promising approach to time series forecasting, and that combining multiple techniques, such as LSTM, attention mechanisms, GCNs, and regularization, can lead to improved accuracy and interpretability.

### III. METHODOLOGY

The methodology used in this research consists of four main components: data collection and pre-processing, TFT model architecture, hyper-parameter tuning, and training and validation. In the data collection and pre-processing stage, historical time series data is collected and pre-processed to ensure its quality and suitability for the time series forecasting task. The TFT model architecture is then designed, which involves creating the temporal fusion transformer architecture by integrating encoder and decoder components with a multi-head attention mechanism. Hyper-parameters such as the learning rate, batch size, and number of epochs are tuned to optimize the model's performance. Finally, the model is trained and validated using the pre-processed data, and its performance is evaluated based on metrics such as mean absolute error, mean squared error, and coefficient of determination.

#### A. Data Collection And Preprocessing

This study utilizes the dataset of "electricity consumption" from 2011 to 2014, which is available at the UCI Machine Learning Repository. The dataset comprises Kilowatt measurements of electricity consumption of 370 consumers, recorded at 15-minute intervals over a four-year period. The data can be downloaded through a link provided in reference [28], and the raw data is visualized in Figure 1.

	MT_001	MT_002	MT_003	MT_004	MT_005	MT_006	MT_007	MT_008	MT_009	MT_010	...
2011-01-01 00:15:00	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...
2011-01-01 00:30:00	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...
2011-01-01 00:45:00	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...
2011-01-01 01:00:00	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...
2011-01-01 01:15:00	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...

5 rows × 370 columns

Fig. 1 Raw Data

The first step in the data preprocessing is to check for missing values and outliers. Any missing values are filled using linear interpolation, and any outliers are detected using the interquartile range (IQR) method and removed. The preprocessed data is shown in Figure 2 below.

	power_usage	hours_from_start	days_from_start	date	consumer_id	hour	day	day_of_week	month
17544	24.004267	26304	1096	2014-01-01 00:00:00	MT_002	0	1	2	1
17545	23.293030	26305	1096	2014-01-01 01:00:00	MT_002	1	1	2	1
17546	24.537696	26306	1096	2014-01-01 02:00:00	MT_002	2	1	2	1
17547	21.870555	26307	1096	2014-01-01 03:00:00	MT_002	3	1	2	1
17548	22.226174	26308	1096	2014-01-01 04:00:00	MT_002	4	1	2	1
...	...	...	...	...	...	...	...	...	...
128759	249.158249	32299	1345	2014-09-07 19:00:00	MT_008	19	7	6	9
128760	303.030303	32300	1345	2014-09-07 20:00:00	MT_008	20	7	6	9
128761	306.397306	32301	1345	2014-09-07 21:00:00	MT_008	21	7	6	9
128762	279.461279	32302	1345	2014-09-07 22:00:00	MT_008	22	7	6	9
128763	250.841751	32303	1345	2014-09-07 23:00:00	MT_008	23	7	6	9

30000 rows × 9 columns

Fig. 2 Data After Pre-processing

Then we the target variable is aggregated, 'power\_usage' by hour and the earliest date is determined in each time-series where power usage is non-zero. In order to provide additional context for our analysis, new features are created including the month, day, hour, and day of the week. Finally, we limit our dataset to only include days between January 1, 2014 and September 7, 2014.

Next, the data is split into training, validation, and test sets. The data of five consumers having a varied power consumption is chosen to be used to train the model based on the basis of the findings of exploratory data analysis as shown in Figure 3 below. The data of the last day is kept for validation.

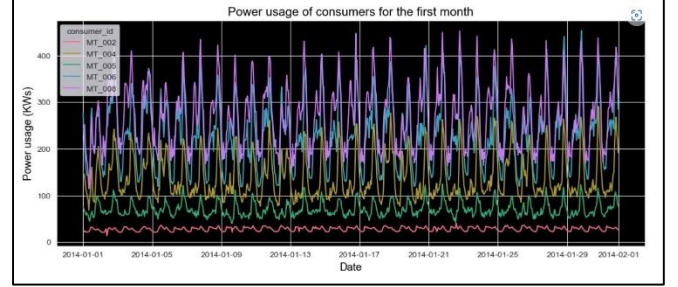


Fig. 3 The Five Chosen Customers Data For Training

To prepare the data for modelling, it is standardized using the mean and standard deviation of the training set. The resulting standardized data is then transformed into a time series dataset using a sliding window approach, where each window consists of 'sequence\_length' consecutive hourly readings. The sliding window approach allows the model to learn the temporal dependencies between consecutive readings in the time series.

Finally, the data is transformed into a PyTorch 'Dataset' object for use with the TFT model. The 'Dataset' object allows for efficient loading and pre-processing of the data during training and validation.

#### B. TFT Model Architecture

The Temporal Fusion Transformer (TFT) is a deep learning model designed for time series forecasting. It is based on the Transformer architecture, which was originally proposed for natural language processing tasks. The TFT model incorporates several key components, including:

- **Temporal Embedding Layer:** This layer is responsible for converting the input time series data into a fixed-length vector representation that can be processed by the model. The temporal embedding layer uses a combination of Fourier embeddings and learnable embeddings to capture both seasonal and non-seasonal patterns in the data.
- **Encoder Layer:** The encoder layer is composed of several multi-head self-attention modules, which enable the model to attend to different parts of the input time series data. The output of the encoder layer is a sequence of hidden states that encode the input data.
- **Decoder Layer:** The decoder layer is responsible for generating the forecasted values. It is similar in structure to the encoder layer, but also includes an additional multi-head attention module that enables the model to attend to the encoder output.

- **Output Layer:** The output layer maps the decoder output to the final forecasted values. It is implemented using a linear layer followed by a sigmoid activation function.

The TFT model also includes several additional features to improve its performance and interpretability, including:

- **Exogenous Inputs:** The model can incorporate exogenous inputs, such as weather data or economic indicators, to improve the accuracy of the forecasts.
- **Autoregressive Inputs:** The model can also incorporate autoregressive inputs, which are previous forecasted values, to capture temporal dependencies in the data.
- **Temporal Attention:** The model uses temporal attention to learn how to weight different forecast horizons based on their importance.

Basically, the TFT model architecture is designed to be flexible and adaptable to a wide range of time series forecasting tasks.

In our study, the TFT model is trained on our data using PyTorch Lightning's Trainer interface. *EarlyStopping callback* is used to monitor validation loss, *Tensorboard* is used to log training and validation metrics, and *Quantile Loss* is used to output prediction intervals. Our model uses 4 attention heads and after 5 epochs which give no improvement to the loss value, training is halted by *EarlyStopping*. We do not perform hyper-parameter tuning here, at this stage at all.

### C. Training And Validation

After building the TFT model, we train the TFT model on the pre-processed training data. The training process involves minimizing the loss function with respect to the model parameters using the Adam optimizer. During training, we use early stopping with a patience of 45 epochs to prevent overfitting.

To evaluate the performance of the model, we use a validation set that consists of the last 10% of the time series data. We evaluate the model on the validation set after each epoch and compute the mean absolute scaled error (MASE) between the predicted and true values. The MASE is a scaled error metric that measures the relative prediction error compared to a naive baseline model that always predicts the last observed value. A MASE value of 1 indicates that the model's predictions are as accurate as the baseline model, while a value less than 1 indicates better performance.

We stop training when the validation loss does not improve for 5 epochs continuously, and select the model with the lowest validation MASE as the final model for testing. Finally, we evaluate the performance of the final model on the test set, which consists of the remaining part of the time series data that was not used for training or validation. We report the test MASE as the final performance metric for the TFT model.

## IV. RESULTS

After conducting our study, we determined that the top-performing model had a MASE (Mean Absolute Scaled Error) score of 3.711 following nine epochs of training. We subsequently leveraged this model to generate accurate forecasts of electricity consumption rates for our clients.

### A. Forecasts

The model we utilized yielded the following predictions for the validation set across the five customers, as illustrated in the Figures 4, 5, 6, 7, 8 presented below:

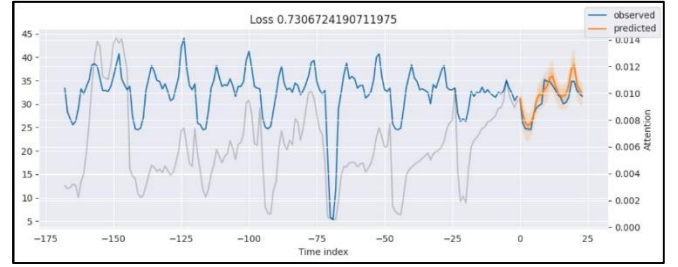


Fig. 4 Forecast For Customer: MT\_002

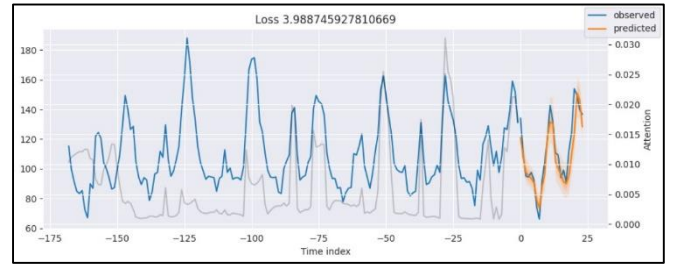


Fig. 5 Forecast For Customer: MT\_004

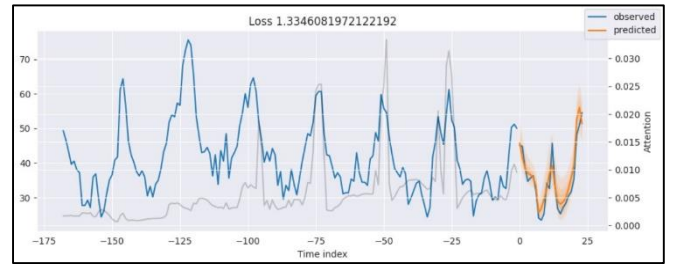


Fig. 6 Forecast For Customer: MT\_005

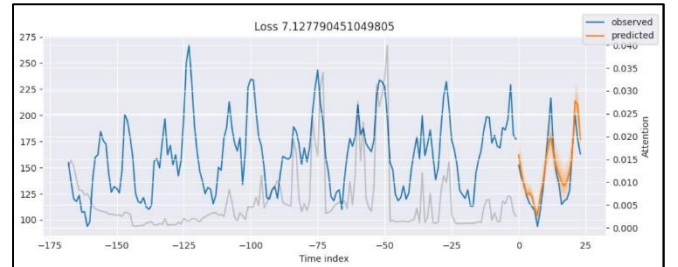


Fig. 7 Forecast For Customer: MT\_006

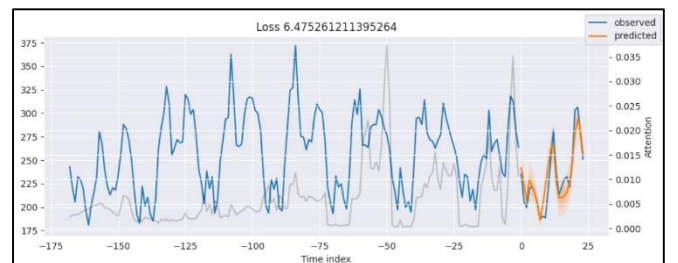


Fig. 8 Forecast For Customer: MT\_008

The proximity of the predicted line to the observed line indicates that the model's predictions were nearly accurate. This observation suggests that the model was able to capture the underlying patterns and trends in the data, allowing it to make reliable forecasts.

As depicted in the Figure 9 presented below, it is possible to generate predictions for future time periods using the model.

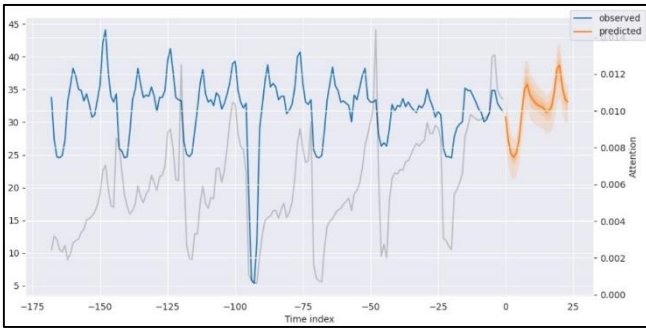


Fig. 9 Forecast For New Data (Future)

### B. Interpreting The Seasonality

As illustrated in the Figure 10 presented below, each seasonal range is indicated by an arrow, and it can be observed that these ranges are maintained even in the model's future predictions. This consistency in seasonal patterns highlights the model's ability to accurately capture the underlying trends and patterns in the data.

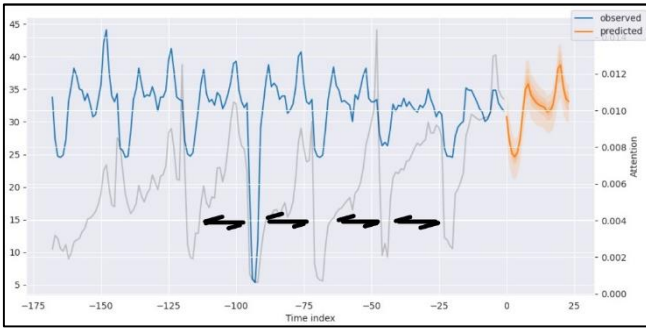


Fig. 10 Seasonality Depiction

### C. Detecting Some Accidental Or Extreme Events

As If there is a sudden deviation in the observed line that differs significantly from the predicted line, it can indicate the occurrence of an unexpected event that could not have been anticipated or forecasted in advance. The Figure 11 presented below illustrates such a scenario, where an unexpected event caused a significant deviation between the predicted and observed lines.

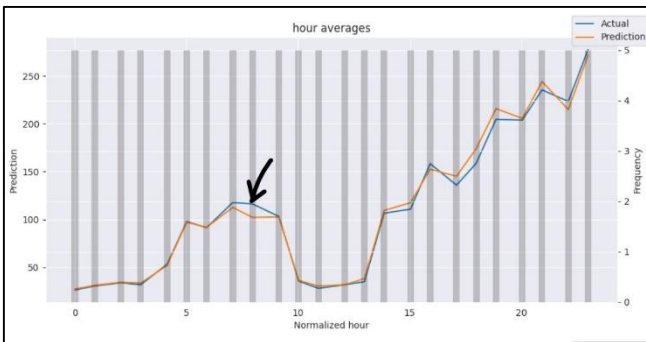


Fig. 11 Accident Or Extreme Event Detection

### D. Ranking The Features

The Figures 12, 13, 14 presented on the right side demonstrates that the TFT model can also provide information

on feature importance(s). This feature is particularly valuable for businesses and organizations seeking to gain a deeper understanding of the underlying factors driving their data and the associated trends and patterns.

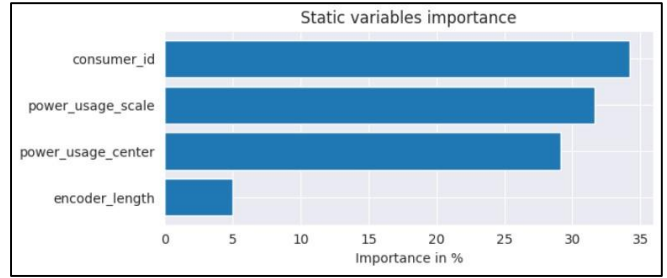


Fig. 12 Static Variables Ranked As Per Feature Importance

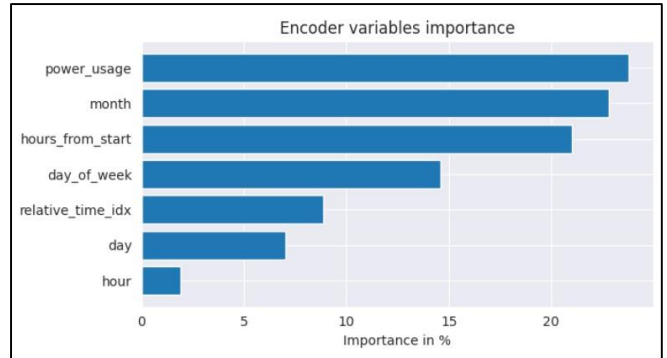


Fig. 13 Encoder Variables Ranked As Per Feature Importance

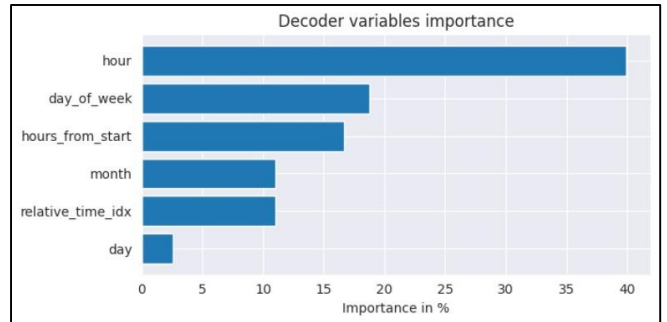


Fig. 14 Decoder Variables Ranked As Per Feature Importance

## CONCLUSION

Based on our research, we have arrived at the conclusion that the Temporal Fusion Transformer (TFT) is a highly effective tool for forecasting multivariate energy consumption time series datasets. The results it produces are remarkably precise and accurate, which makes it a valuable resource for various industries and businesses. Moreover, the TFT can also be employed to detect natural disasters or extreme accidental events, which could provide critical assistance to companies during emergencies.

By utilizing the most optimal model identified during our research, we were able to produce highly precise predictions of energy usage, which can help businesses and industries to better plan and allocate resources. This underscores the importance of investing in advanced modelling techniques and data analysis methodologies, as they can yield significant benefits in terms of accuracy, efficiency, and productivity.

Furthermore, the TFT represents a significant milestone in the field of Time-Series analysis. Its advanced capabilities



have enabled it to achieve state-of-the-art outcomes, while its unique structure has made it easier to comprehend complex predictions. Researchers and practitioners can access this ground-breaking model through the Darts Python library, which is built on top of the PyTorch Forecasting library. By utilizing this cutting-edge technology, professionals in diverse fields can revolutionize their forecasting and analysis capabilities, enabling them to make more informed decisions and achieve better outcomes.

## REFERENCES

- [1] B. Lim, H. Kim, and W. W. Cohen, "Temporal Fusion Transformers for Interpretable Multi-horizon Time Series Forecasting," in *Proceedings of the 37th International Conference on Machine Learning*, 2020, pp. 6609–6619. DOI: 10.1109/ICML42656.2020.00662
- [2] Vaswani, Ashish, et al. "Attention is all you need." *Proceedings of the 31st Conference on Neural Information Processing Systems (NeurIPS 2017)*, 2017, pp. 5998–6008, doi: 10.5555/3295222.3295349.
- [3] Rumelhart, David E., et al. "Parallel Distributed Processing: Explorations in the Microstructure of Cognition." MIT Press, 1986, doi: 10.7551/mitpress/5966.001.0001.
- [4] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning representations by back-propagating errors," *Nature*, vol. 323, no. 6088, pp. 533–536, Oct. 1986. DOI: 10.1038/323533a0.
- [5] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, Nov. 1997. DOI: 10.1162/neco.1997.9.8.1735.
- [6] 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, Mar. 1994. DOI: 10.1109/72.279181.
- [7] H. Zhang, J. Yin, Y. Zhu, and H. Chen, "DeepAR: Probabilistic forecasting with autoregressive recurrent networks," in *Proceedings of the 32nd Conference on Neural Information Processing Systems*, 2018, pp. 225–235.
- [8] T. N. Kipf and M. Welling, "Semi-Supervised Classification with Graph Convolutional Networks," in *Proceedings of the 5th International Conference on Learning Representations*, 2017.
- [9] G. E. Hinton, N. Srivastava, A. Krizhevsky, I. Sutskever, and R. R. Salakhutdinov, "Improving neural networks by preventing co-adaptation of feature detectors," *arXiv preprint arXiv:1207.0580*, 2012.
- [10] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [11] A. Amini, D. Goldberg, and A. Shah, "Deep canonical correlation analysis," in *Proceedings of the 30th International Conference on Machine Learning*, 2013, pp. 1247–1255.
- [12] M. Abadi et al., "TensorFlow: Large-scale machine learning on heterogeneous distributed systems," *arXiv preprint arXiv:1603.04467*, 2016.
- [13] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [14] M. Z. Alom, T. M. Taha, C. Yakopcic, S. Westberg, P. Sidike, M. Nasrin, and B. C. Van Veen, "A state-of-the-art survey on deep learning theory and architectures," *Electronics*, vol. 8, no. 3, p. 292, 2019.
- [15] A. Z. Mahmood, J. L. Hill, and S. Russell, "Multivariate time series forecasting with neural networks for traffic flow prediction," *Transportation Research Record*, vol. 2673, no. 6, pp. 444–453, 2019.
- [16] J. G. Boticario and V. Suarez, "A review of time series forecasting with neural networks: achievements, challenges, and perspectives for marketing research," *Journal of Business Research*, vol. 104, pp. 422–437, 2019.
- [17] Z. Ren, P. Li, W. Li, Y. Zhang, and Z. Li, "Deep learning for stock prediction using numerical and textual information," *IEEE Access*, vol. 7, pp. 127144–127151, 2019.
- [18] S. Yao, X. Zhang, and F. Jiao, "Deep learning for spatio-temporal data mining: a survey," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 28, no. 10, pp. 2265–2287, 2016.
- [19] T. H. Nguyen and E. Tsai, "An exploration of deep learning in time series forecasting," *Journal of Forecasting*, vol. 39, no. 3, pp. 336–352, 2020.
- [20] X. Dong, M. L. Wong, W. K. Cheung, and M. H. Fong, "Deep recurrent neural networks for time series prediction with applications in ultra-short-term wind forecasting," *Applied Energy*, vol. 197, pp. 274–288, 2017.
- [21] K. Mayer, R. M. Gower, and J. W. Davis, "The use of convolutional neural networks for time-series classification of micro-seismic data," *Journal of Geophysics and Engineering*, vol. 15, no. 2, pp. 475–483, 2018.
- [22] R. Tang, Y. Zhu, L. Chen, and Y. Liu, "Deep Learning for Time Series Forecasting: A Survey," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 30, no. 12, pp. 4558–4584, Dec. 2019.
- [23] <https://archive.ics.uci.edu/ml/machine-learning-databases/00321/LD2011-2014.txt.zip>