

**To what extent can the incorporation of Temporal Convolution into Transformer Improve
the Efficiency of Time Series Prediction?**

Word Count: 3467

Computer Science

To what extent can the incorporation of Temporal Convolution into Transformer Improve the Efficiency of Time Series Prediction?

Contents

Introduction	4
Related work	4
Transformer architecture	4
Convolution-based architectures	5
Convolution attention combined architectures	6
Inspiration	6
TCformer	10
Structure Overview	10
Convolution extraction	12
Sequence length reduction	13
Experiments	13
Experiments for single-step prediction	13
Dataset	14
Experiment setup	14
Results	14
Experiments for sequence to sequence	17
Dataset	17
Experiment setup	18
Results	19
Efficiency	21
Conclusionand Evaluation	23

Future works 23

Jupyter notebook code for experiment on Transformer's single-step-prediction task . . .	34
Jupyter notebook code for experiment on LSTM's single-step-prediction task	43
The TCformer using 1D CNN for basic Transformer on seq2seq task	52
The TCformer using 2D CNN for basic Transformer on seq2seq task	57
The TCformer using 1D CNN for iTransformer on seq2seq task	61

Introduction

Time series prediction is a critical task in various domains such as finance, healthcare, weather forecasting, and transportation. It uses historical data to predict the future. In recent years, the attention mechanism was well known after being implemented in the transformer architecture [1]. Since 2017, there have been many works around this model in many tasks such as natural language processing [2], and computer vision [3] given its performance. In the field of time series prediction, many researchers found that the standard Transformer architecture may not fully exploit the temporal relationships and dependencies present in time series data. So there many variants of transformers appeared, including informer [4], patchTST [5], that all aim to improve its performance. In this work, we propose a generalized architecture, where the input data can be partially extracted. Then, we used an example of using CNN and Transformer in this architecture, the TCformer, and experimented it's efficiency and performance.

Related work

Transformer architecture

Despite the impressive works above, the transformer architecture comes with its great size. Many other works in training large-scale transformer models appeared with up to 600 billion weights that cost for 2048 TPU v3 cores [6]. Recent work [7] has reduced the size with compression with the cost of precision. Currently, researchers have no way to get around with the bulky model without compromises. Seeing this flaw, we proposed a new hybrid architecture TCformer that both uses CNN and Transformer to achieve an improvement in precision with the reduction in size for multivariate time series prediction tasks.

Considering the architecture of the Transformer, different from convolution layers, its parameter size is inherently dependent on the size of the input. In Transformer, the scaled dot-product attention [1] computes the output as:

$$Attention(Q, K, V) = softmax(\frac{QK^T}{\sqrt{d_k}}V) \quad (1)$$

Where Q represents the query matrix; K is the key matrix, and V value matrix. In the time series prediction task, those matrices are 3D tensors that have the shape of:

$$Q, K, V \in \mathbb{R}^{[B \times L \times (D/H)]} \quad (2)$$

where B is the batch size, H the number of heads, L the sequence length, D the numbers of features. Thus, they are all dependent on the input size. Even though it can be addressed through a linear layer that maps the dimension lower, it's not a common approach due to the huge information loss through this process which decreases the model's precision. This attention mechanism itself already made this model large, not to mention the parameters in the FFN (Feed Forward Network).

Due to this full attention mechanism, this leads to the quadratic dependency of Transformer-based models [8], which makes the computational and memory dependent on the sequence length. With that being said, reducing the dimension while maintaining the precision of the attention mechanism became the key to reducing the model's size.

Flash attention has made a significant breakthrough in reducing the space complexity to linear and using SRAM to reduce IO complexity. Different from the algorithmic improvement, we made the model more efficient by changing the architecture without changing the specific layers. These two approaches don't conflict and can be used simultaneously. However, in our work, we used the most basic encoder layer for simplicity.

Convolution-based architectures

Though Transformer has gained the most focus after its success in Large Language Models (LLMs) and was mostly researched in doing various tasks including times series, convolution has also been applied to time series [9]. From the two base mechanisms, the time series models have diverged into two groups as more models are proposed based on either CNN or Transformer. Moreover, recent work [10] has shown that CNN is more accurate than the Transformer model, reaching 92% accuracy compared to the 80% accuracy of the Transformer.

However, pure convolutional architectures are hard to deal with time series data. Each kernel, relatively small to the whole sequence, can only see a very limited range of data and can have an overview of a long series only when the network is deep enough and the pieces of traits converge at the top of the network. Thus it's hard for those architectures to expand their capacity in processing long sequences. Moreover, CNN has many more hyperparameters on each layer, such as kernel size, stride and padding. This makes it hard to build a model that's at its maximum performance. In our work, the CNN is used for simply extracting general features in the past data which only requires one to two layers of it, and leaves the rest of the prediction task to the Transformer encoder.

Convolution attention combined architectures

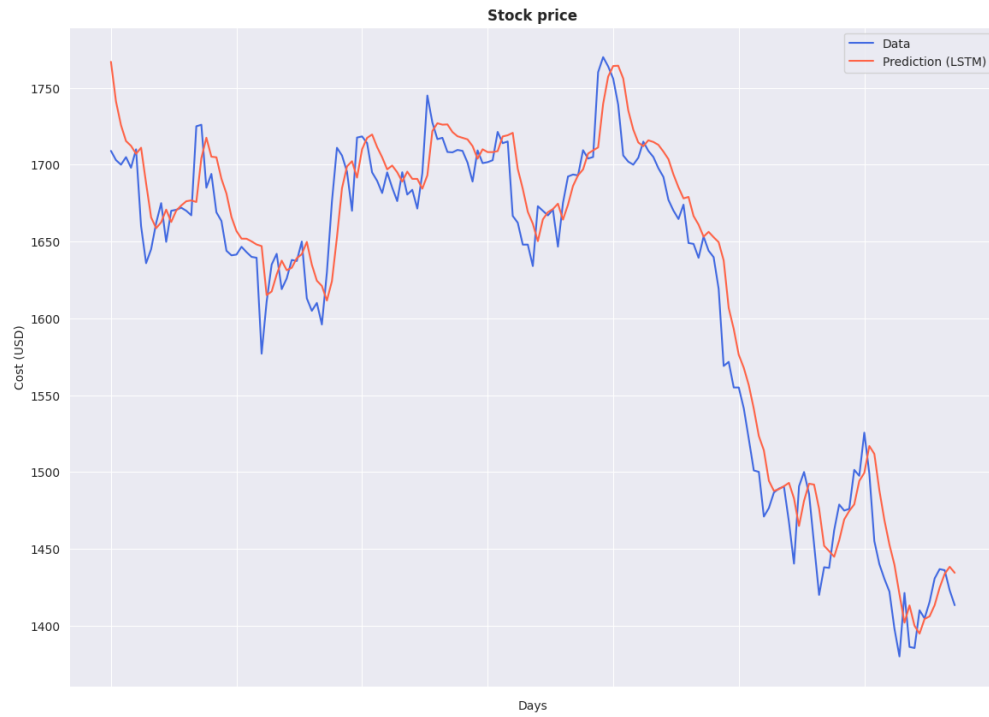
The combination of convolution layers with attention mechanism isn't a rare approach. In ConTNet, researchers have embedded the encoding block in between the convolution layers [11] for computer vision tasks. In time series prediction, Informer [4] has changed the encoder by adding convolution layers between the attention layers for handling longer time series sequences.

Unlike previous works, the TCformer we proposed neither embed a convolutional layer in between the model nor applies it to the whole input.

Inspiration

We originally used Transformer solely on stock prediction task and were surprised how poorly it performed in the prediction compared to other models.

(a) *LSTM*



(b) *Transformer*

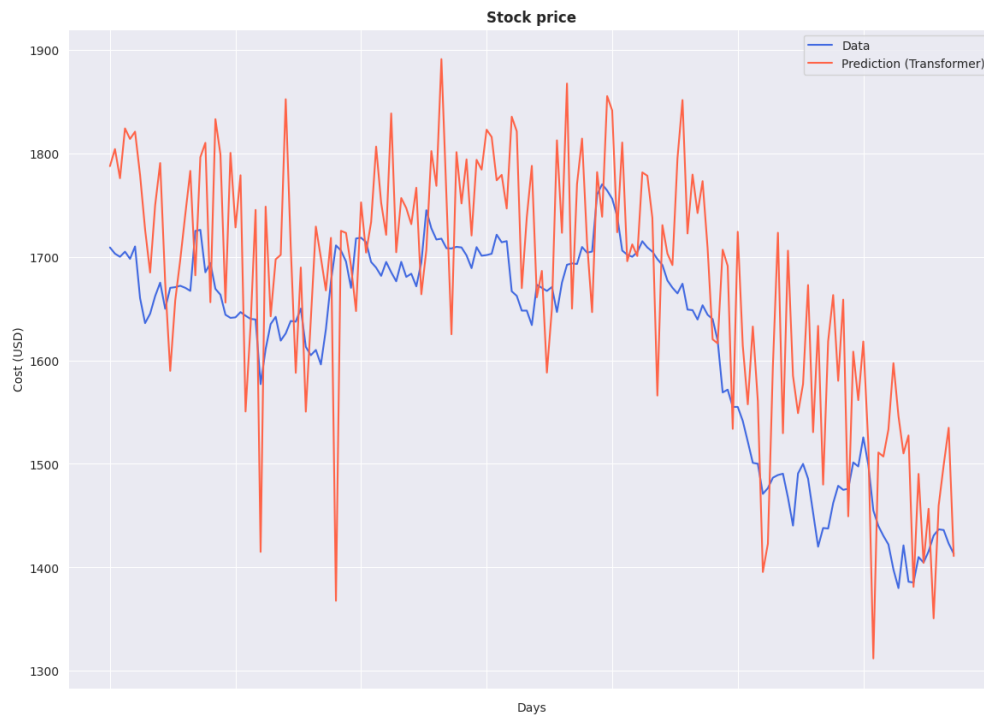


Figure 1

Stock price prediction comparison

We ran both LSTM and Transformer on a Chinese stock with code 600519. On the left of Figure 1, LSTM is predicted relatively closely and shows a smooth trend. However, on the right, the Transformer's prediction had shown a significant variation with sharp turning points that occurred frequently. Numerically, LSTM's MSE test loss is about 91.84% lower than Transformer's. We reflected on why the prevalent Transformer model failed in such a task. After trying a range of hyperparameters, the Transformer's performance had some minor fluctuation, but the sharp fluctuation remains. By examining the prediction graph of the Transformer, we found out that the general trend for the Transformer is correct. Thus, the high loss is attributed to the sharp peaks and troughs. Conceptually, we guessed it's the attention mechanism that led to this issue.

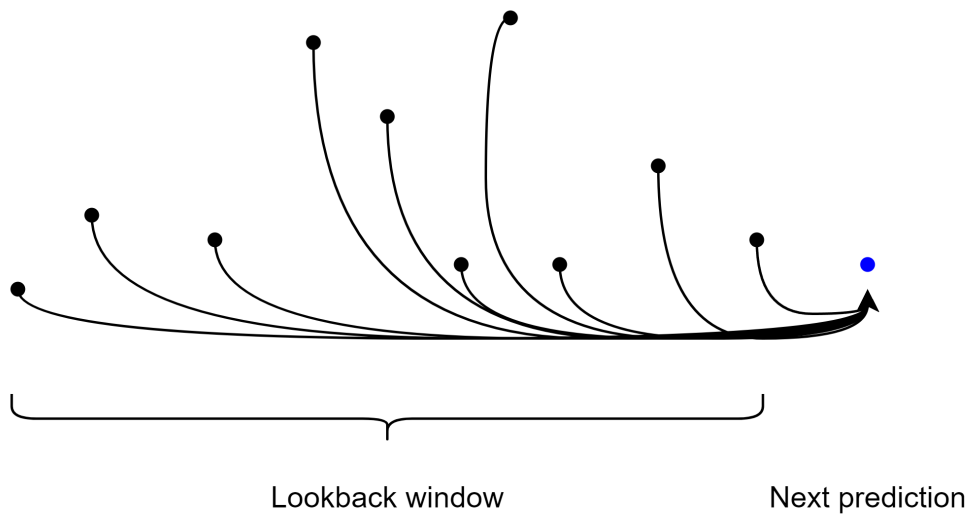


Figure 2

Attention in time series

In Figure 2, it illustrates the attention mechanism in time series. Every point in the lookback window (black) will have a weighted contribution to the probability of the next prediction (blue). While this might be successful when the data points are word tokens in Natural Language Processing (NLP), this use of information is too detailed and specific in time series tasks. The model will need to consider every single data point that might fluctuate, leading to increased fluctuation in prediction. As the series increases, there will be more irrelevant

fluctuations captured by the Transformer while failing to capture the general trend.

To address this problem, we used CNN to "summarize" and extract the input features. This can compress a range k data points into a general one, where k is the CNN kernel size while leaving the recent data (last part of data) uncompressed.

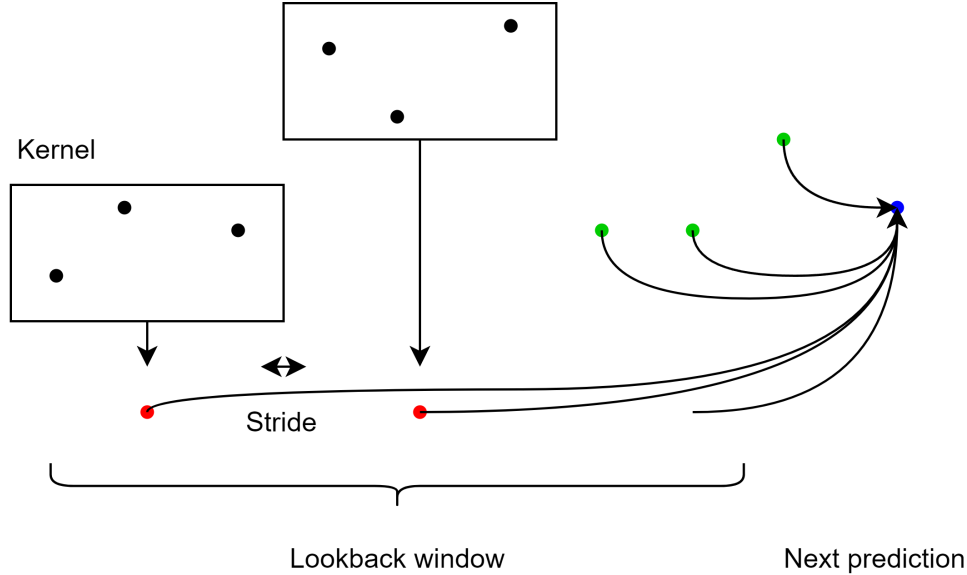


Figure 3

TCformer attention for time series

In Figure 3, the kernel of shape 3 is extracting the feature of 3 data points and mapping it into one (red). The distance the kernel moves forward is dependent on the stride.

The extracted data points (red) along with the raw (green) then applied attention in Transformer. This not only reduces the sequence length but also captures the trend in previous data.

However, this design is based on an important assumption: more recent data in the dataset needs to be more important. This allows us to split the data. This applies to many different datasets, such as traffic, solar, electricity, stock and so on. Conversely, in language models where the input might refer to words in a very early context would not be expected to have an improved precision.

TCformer

The architecture we propose is not a specific one but a general idea of how current models can preprocess the input tensor for better efficiency. In this work, we used a simple Encoder for the base model and CNN to preprocess the input. For the most common single-step forecasting task, where the model predicts the next timestamp X_{n+1} by observing the historical data X_1, \dots, X_n .

Structure Overview

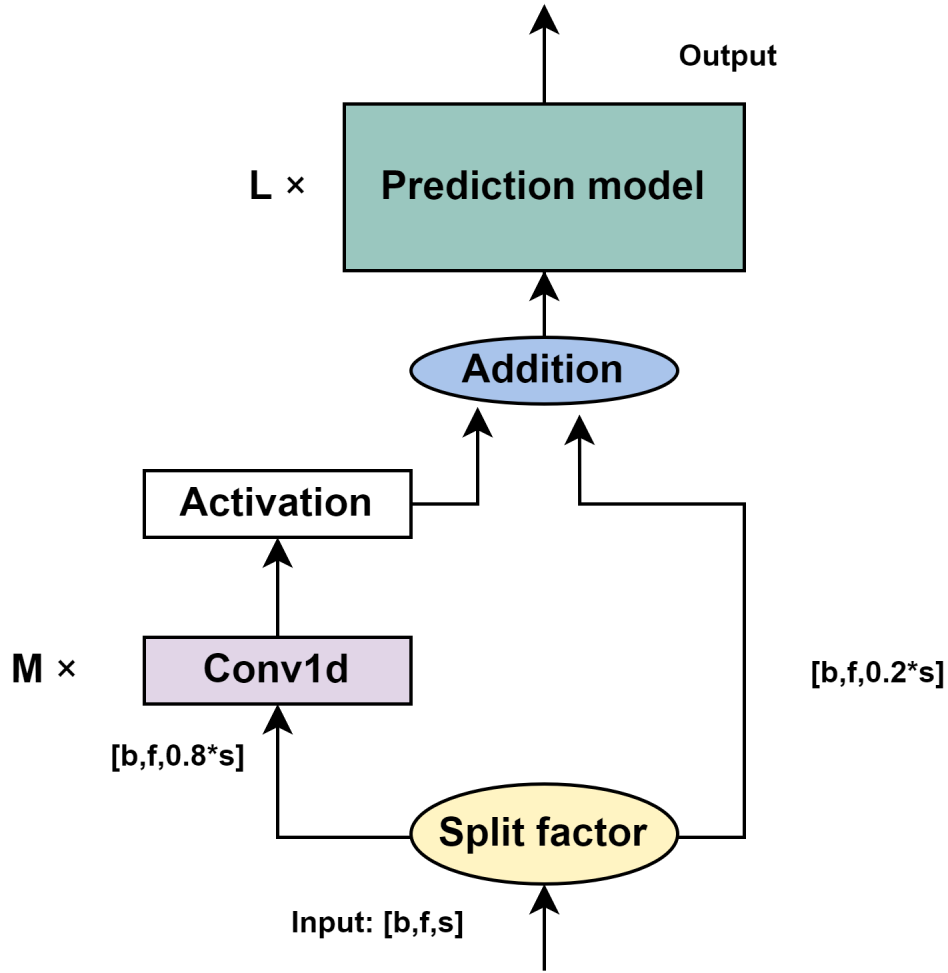


Figure 4

TCformer sample structure

The TCformer has a simple and flexible structure. In general, we use CNN for extracting the first $p\%$ of the input data, leaving the rest raw, and feeding it to the prediction model. In

Figure 4, the split factor is a scalar p that splits the dataset into two on the dimension of sequence length. Note that the split factor is a hyperparameter that needs to be manually assigned. For example, a scalar $p = 0.8$ would split the input of shape $X \in \mathbb{R}^{b,f,s}$ into $X \in \mathbb{R}^{b,f,0.8 \cdot s}$ and $X \in \mathbb{R}^{b,f,0.2 \cdot s}$ where:

1. b is the batch size
2. f is the dimension of the feature
3. s is the sequence length

In Figure 4, there can be M Conv1d layer. They are 1-dimensional convolution layers that extract traits from the $p\%$ of the input and reduce the sequence length, where p is the split factor and it's set to 0.8 in this illustration. There could be multiple layers added flexibly. It can be replaced by a Conv2d (needs to ensure that the dimension of the model doesn't reduce). The right branch leaves the raw data that will be joined back with the left branch in the concat component. In the concatenated component, the two branches' tensors will be merged on the third dimension s . Lastly, the concatenated tensor will be the input for the prediction model.

Convolution extraction

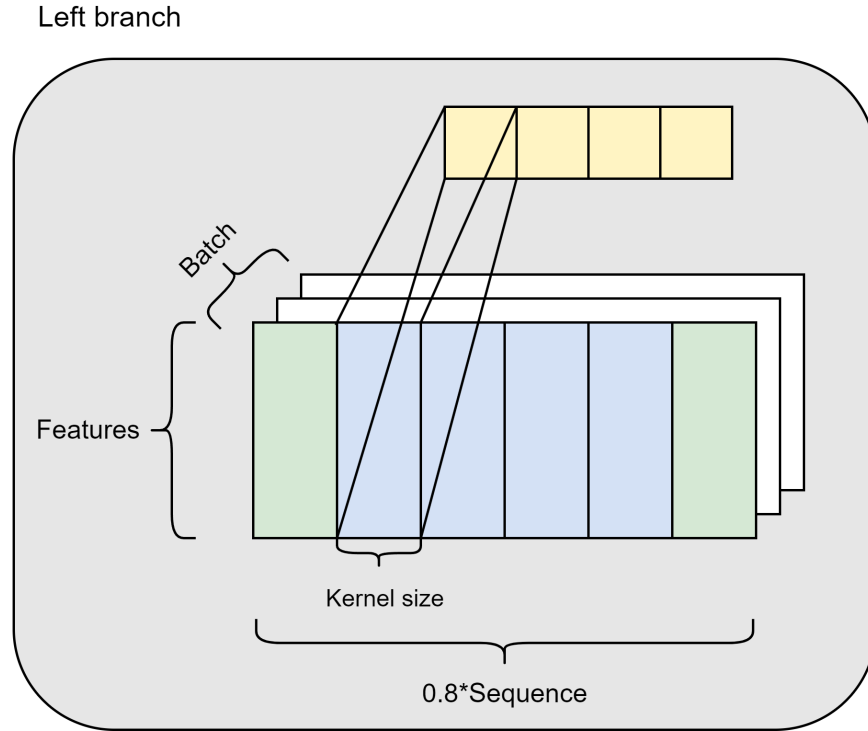


Figure 5

Convolutional sequence summary on left branch

In the left branch illustrated in 5, the Conv1d layer is outputting summaries for the first 80% of earlier data, leaving the rest 20% raw for the Transformer encoder to process. This will cause the input X to change its third dimension's shape s from:

$$s_1 = 0.8 \cdot s \quad (3)$$

$$s_2 = \lfloor \frac{0.8 \cdot s - k}{r} \rfloor + 1 \quad (4)$$

1. s_1 is the sequence length before the convolution layer
2. k is kernel size
3. s_2 is the sequence length after the convolution layer

4. r is the stride

On the right branch, the raw data is not processed. Given that the Transformer has already a good performance on raw data, applying convolution is not needed and might reduce the information perceived in the encoder block.

Sequence length reduction

After the addition, the input for the encoder block will have a sequence length of:

$$s_3 = s_2 + 0.2 \cdot s \quad (5)$$

$$= \lfloor \frac{0.8 \cdot s - k}{r} \rfloor + 1 + 0.2 \cdot s \quad (6)$$

The reduction of sequence length s_3 compared to s can be calculated:

$$l = s - s_3 \quad (7)$$

$$= s - \lfloor \frac{0.8 \cdot s - k}{r} \rfloor - 1 - 0.2 \cdot s \quad (8)$$

Notice l is positively related with k and r , changing the stride and kernel size can give a flexible control on the model size. In the attention block, a reduction of l will reduce:

$$l \times b \times f \quad (9)$$

in each of the Q, K, V matrices, largely reducing its parameters.

Experiments

We experimented with the model on both single-step prediction and sequence-to-sequence (seq2seq) tasks.

Experiments for single-step prediction

We conducted experiments on 3 real-world datasets in the experiment, including weather used in Autoformer [12] and AMD stock data.

Dataset

The weather dataset included 52696 rows of data with 21 covariates.

The AMD stock dataset is accessed through Baostock API [13]. It includes 10995 rows of data with 5 covariates. The data spans from the 2nd of January 1981 to the 14th of August 2024.

Experiment setup

The transformer model has a hidden dimension of 8, 2 layers, and 2 heads for attention. For all models, the optimizer uses Adam, with $lr=1e-3$, $wd=1e-4$, $dropout=0.1$, training in 5 epochs with a batch size of 512 with train test 9 : 1 splitting. The TCformer uses a split factor of 0.8 with the CNN using kernel size 6 with a stride of 6. The criteria uses MSE (Mean Squared Error).

Results

TCformer1d	transformer	seq_len	pct_improve(%)
0.018490	0.073744	16	74.93
0.015125	0.061805	32	75.53
0.012830	0.053022	64	75.80
0.017049	0.069303	128	75.40
0.016919	0.070958	256	76.16
0.017332	0.072824	512	76.20

Table 1

AMD stock data test MSE loss

In Table 1, for AMD stock data, the TCformer with 1d on average has improved performance by about 75% while having reduced size of the model. During the experiment, the TCformer was successfully trained on the dataset when the sequence length is 1024 however the Transformer model failed due to GPU memory overflow, so the sequence length for the

experiment is limited up to 512.

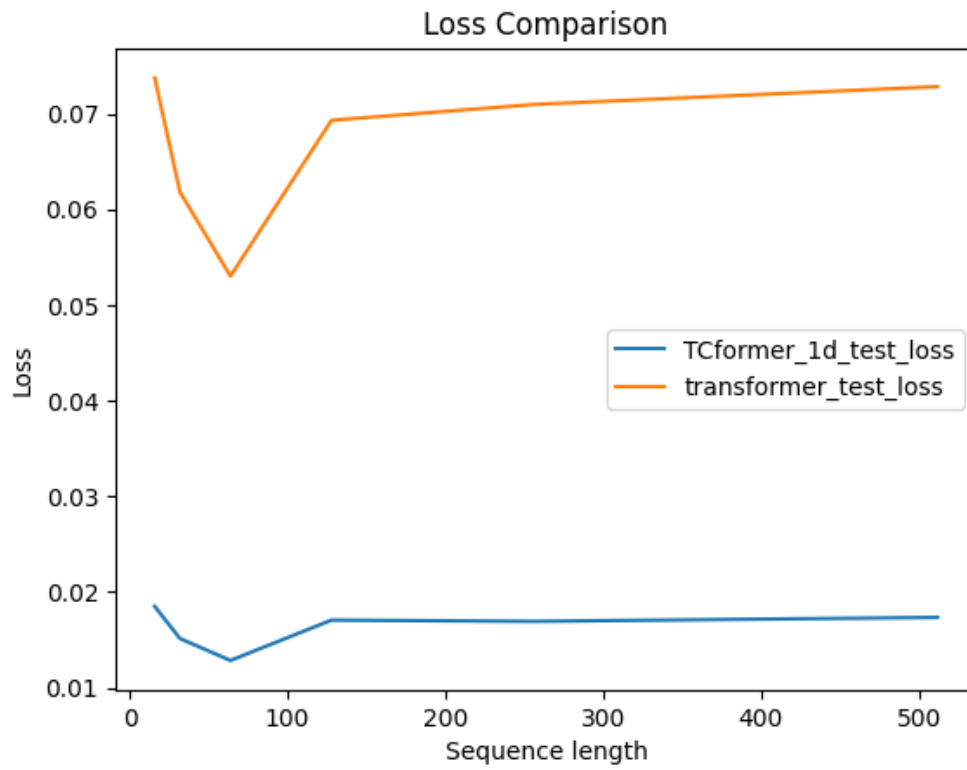


Figure 6

AMD loss curve

In Figure 6, the test loss for both models has its minimum at a sequence length of 64. For every sequence length, the TCformer has a much lower loss compared to Transformer. The Transformer loss shows an increase as sequence length increases from 128 to 512 while this is relatively not as significant in TCformer's loss. This shows that the sequence reduction has enabled TCformer to handle longer sequences while maintaining similar precision.

TCformer1d	transformer	seq_len	pct_improve (%)
0.026144	0.104435	16	74.97
0.026144	0.104253	32	74.92
0.026188	0.103928	64	74.80
0.025542	0.101652	128	74.87
0.028019	0.107837	256	74.02
0.026371	0.104923	512	74.87

Table 2

Weather data test MSE loss

Like the results in Table 1, the percentage improved is around 74% in Table 2. Given that Transformer has a poor performance, this great improvement is within expectation.

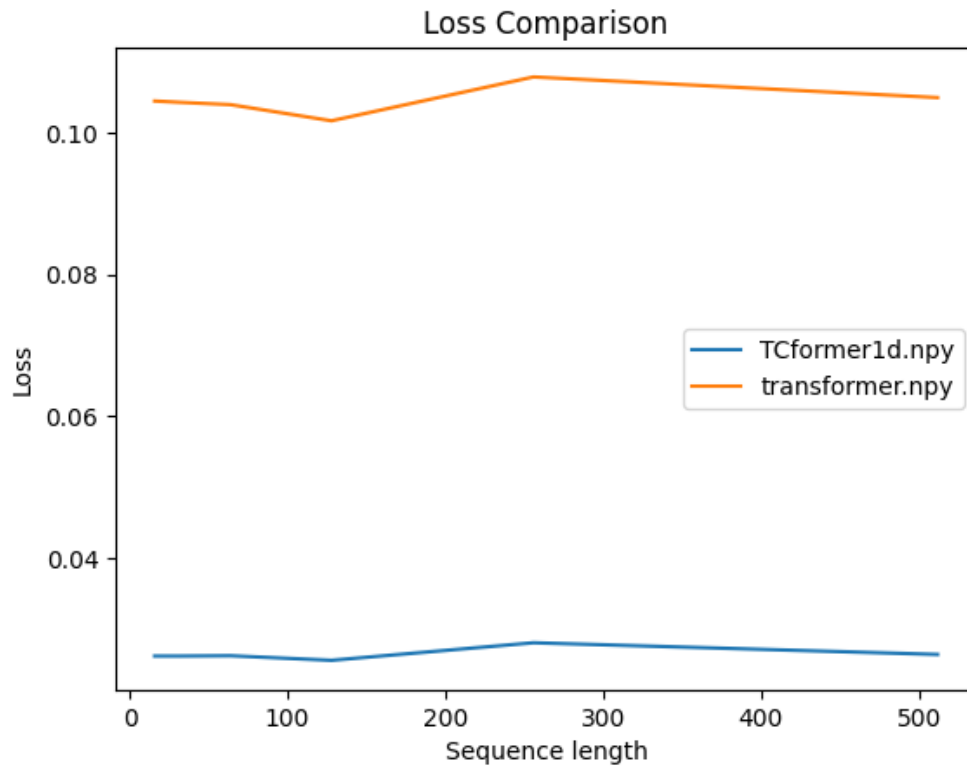


Figure 7

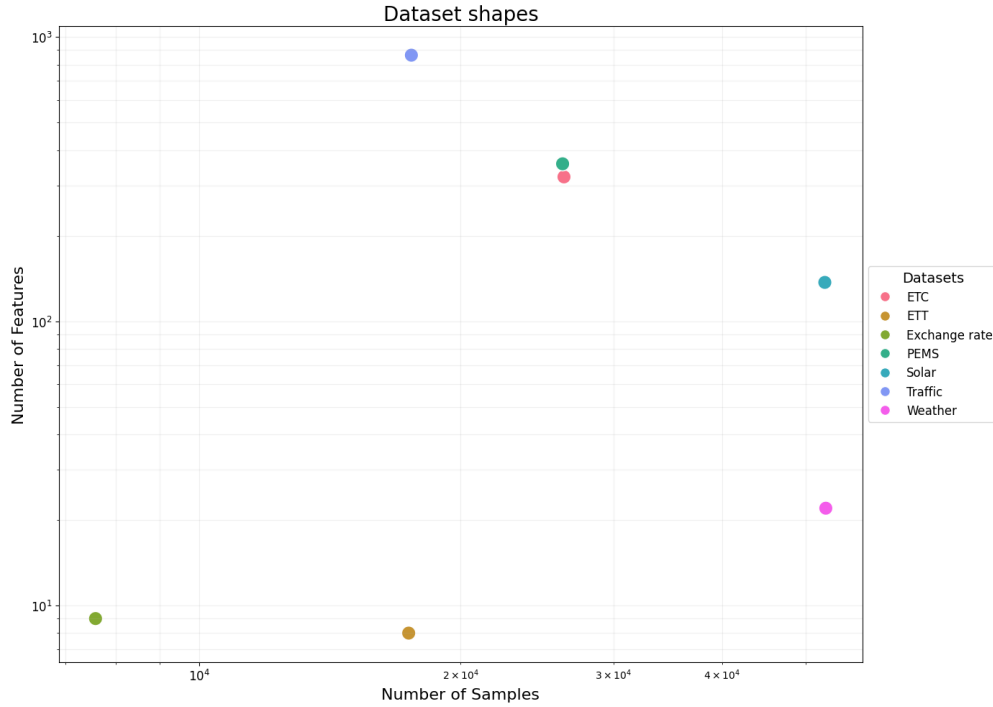
Weather loss curve

Experiments for sequence to sequence

Dataset

For sequence to sequence, we used 7 datasets (ECL, ETT subset 1, Exchange, Traffic, Weather, and Solar energy), the same as the experiment for iTransformer [14]. For each dataset, we ran Transformer, iTransformer, and TCiTransformer (iTransformer with CNN). Since iTransformer currently has the best accuracy on different datasets, we used it as our prediction model, simply replacing the TCformer’s Transformer Encoder with the iTransformer Encoder.

The dataset covers a range of samples and feature sizes:

**Figure 8***Dataset shapes****Experiment setup***

For each model, we set the prediction length to $P = \{96, 192, 336\}$ and the lookback length to 96. Since PEMS is too large, we set the lookback 48 and prediction length to $P = \{12, 24, 48\}$. We set two layers of CNN, and split factor to 0.8 for the iTransformer with CNN. The kernel of the first layer is 5, stride 2 and the second layer is 3, stride 1. We experimented with the model with different hyperparameters such as the CNN kernel size, stride, split factors and so on. The impact on the performance isn't obvious, so we set two layers to have the kernel capturing 9 timestamps at each time. This ensures that most of the dataset doesn't cover a whole period but extracts the trends that exist in a period and then gives it to the prediction model. The kernel isn't set too small to ensure that it can reduce the sequence length.

Results

Dataset	Model	MSE	MAE	Input length	Output length
ECL	Transformer	0.249665	0.349949	96	96
ECL	iTransformer	0.167462	0.256914	96	96
ECL	TCiTransformer	0.180147	0.280937	96	96
ETTh1	Transformer	0.796016	0.702380	96	96
ETTh1	iTransformer	0.387790	0.406386	96	96
ETTh1	TCiTransformer	0.393983	0.415327	96	96
Exchange	Transformer	0.677474	0.637843	96	96
Exchange	iTransformer	0.086314	0.206543	96	96
Exchange	TCiTransformer	0.086183	0.207386	96	96
PEMS03	Transformer	0.105133	0.204023	48	12
PEMS03	iTransformer	0.073994	0.179749	48	12
PEMS03	TCiTransformer	0.066601	0.170845	48	12
solar	Transformer	0.198320	0.232168	96	192
solar	iTransformer	0.212609	0.245238	96	96
solar	TCiTransformer	0.224425	0.245268	96	96
traffic	Transformer	0.665396	0.365834	96	96
traffic	iTransformer	0.465135	0.323030	96	96
traffic	TCiTransformer	0.546155	0.343953	96	192
weather	Transformer	0.221558	0.310034	96	96
weather	iTransformer	0.183019	0.224638	96	96
weather	TCiTransformer	0.167566	0.213238	96	96

Table 3

Summary table for loss (selecting lowest MSE)

For each model on each dataset, we selected it's the best-performing results from the varying output sequence length. In Table 3, iTransformer1dSplit is the iTransformer with CNN. For the seq2seq task, the TCiTransformer didn't have much increase compared to iTransformer in most datasets. In ECL, ETTh1, solar and traffic, the model's precision is lower than iTransformer. However, in PEMS03, Exchange and weather, TCiTransformer had a higher precision. The full result is in Appendix .

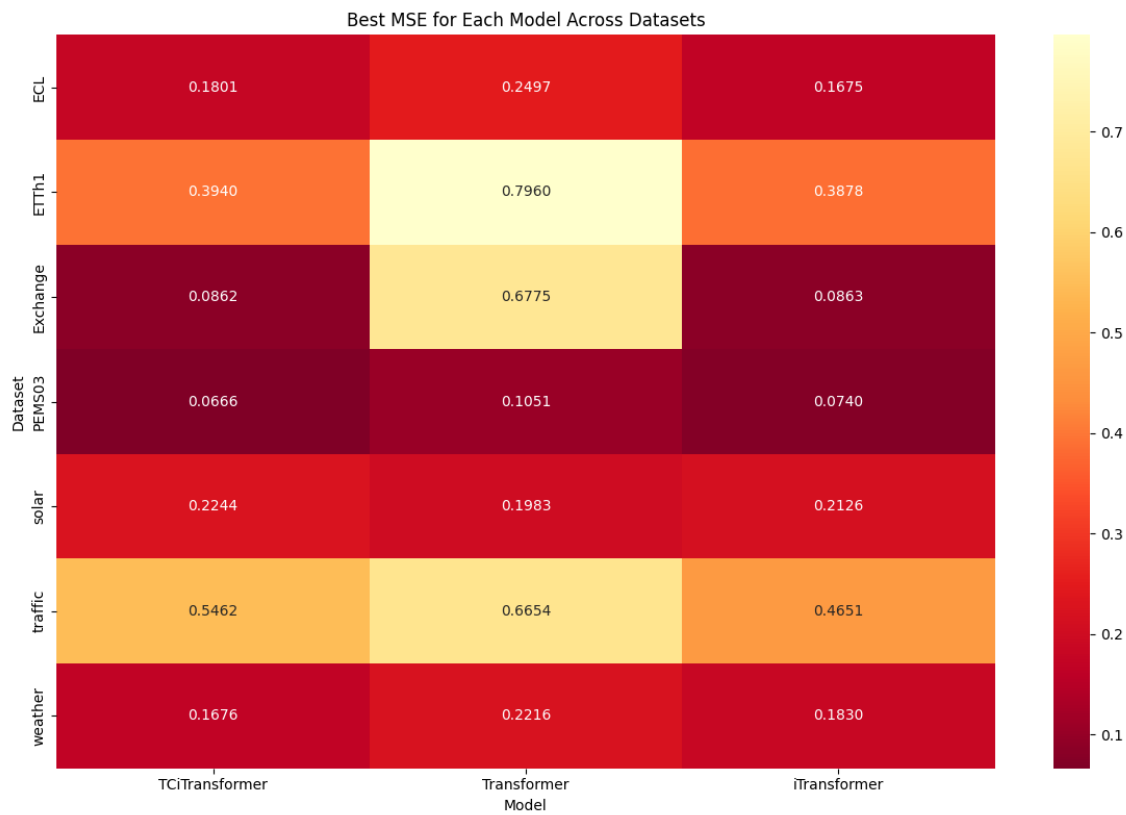


Figure 9

MSE heatmap

This is illustrated clearer in Figure 9. The darker means the better the model is.

We selected the prediction of our model on traffic, which it's performing worse, and weather which it performs better:

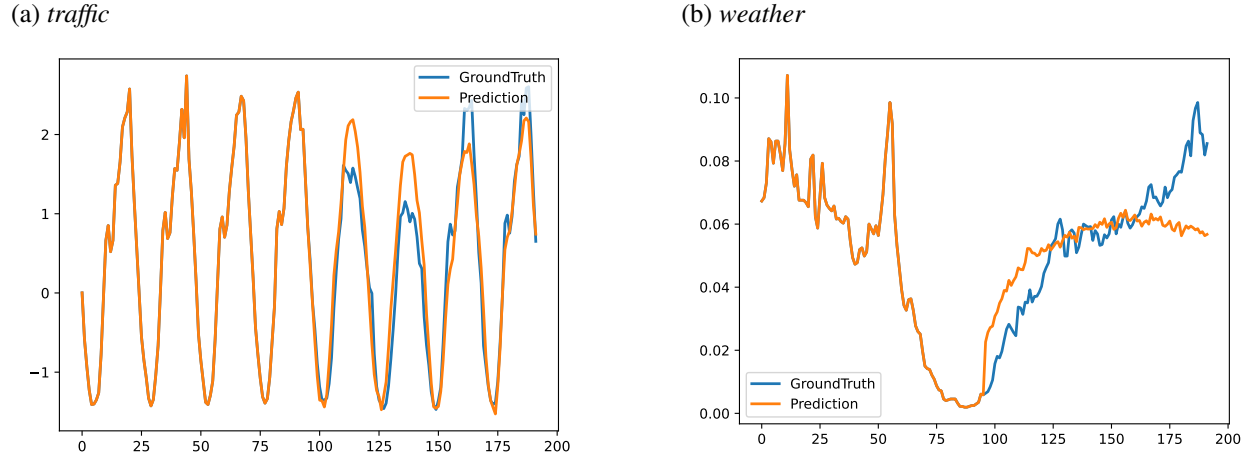


Figure 10

Prediction on traffic and Weather for both input and output length is 96

The traffic data is much more periodic with a generally regular pattern, while the weather data is more irregular. By observing other datasets, we found out that TCiTransformer is better in non-periodic time series forecasting. This might be because the CNN cannot be trained to capture the increase and decrease over a small range of time, while it gives the general trend in datasets such as weather.

Efficiency

Although the sequence length can be reduced, as mentioned in section , the CNN requires time for operation. Depending on the kernel size and stride, the time for TCformer architecture can vary. In the experiments above, we recorded the training time for traffic data, using the same kernel size and stride:

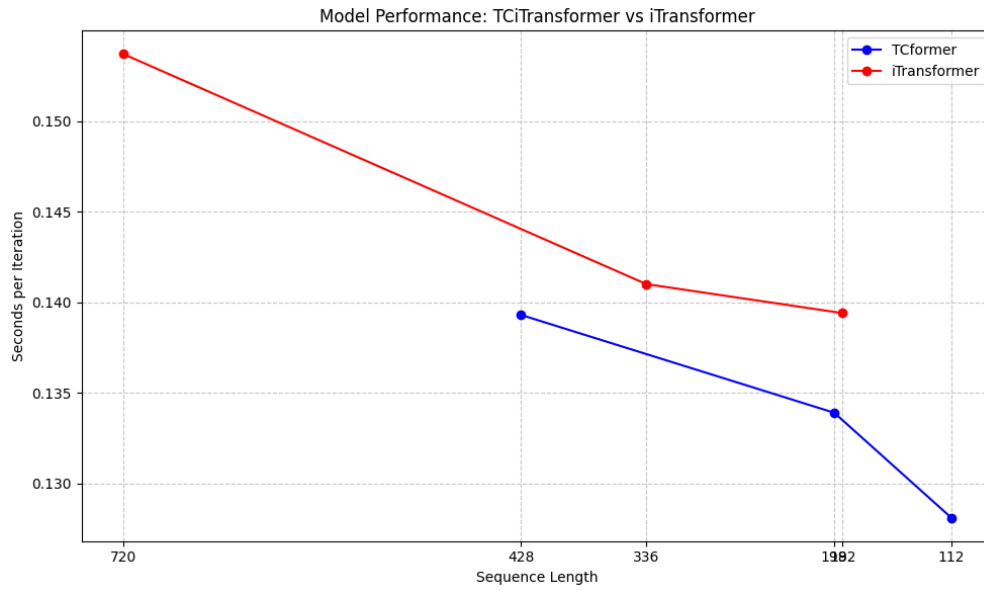


Figure 11

Speed comparison on different sequence length scale

Fig.11 shows the speed of TCiTransformer against iTransformer. Because TCformer architecture reduces the sequence length, so there is a right shift of the blue line on the sequence length. However, even at the same sequence length, TCiTransformer outperformed iTransformer.

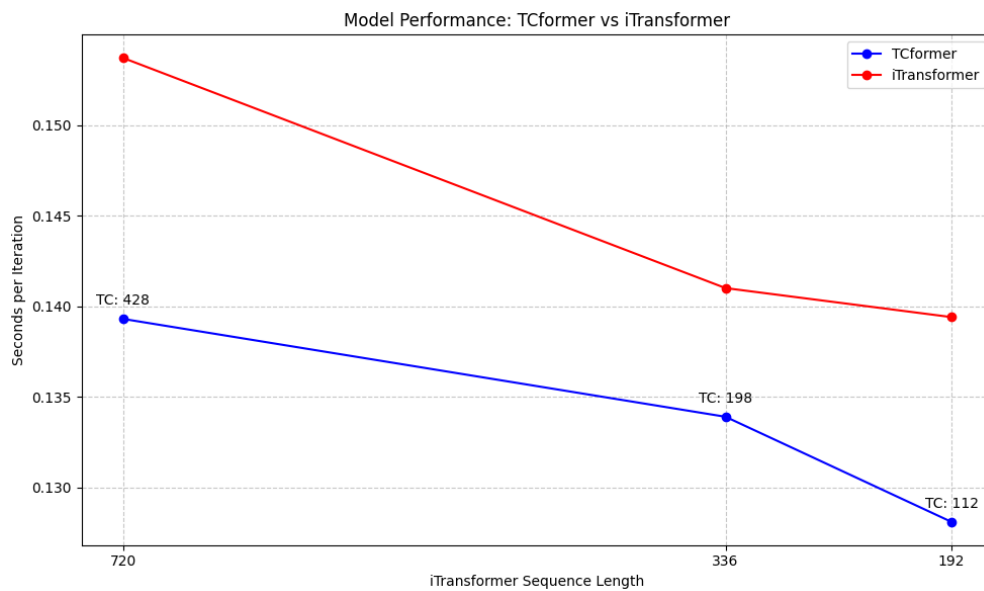


Figure 12

Speed comparison on same sequence length scale

Fig.12 plots the speed of the model on the same sequence length scale. It shows at the same input sequence length, the reduced encoder block input will lead to an increase in speed, despite the cost of CNN calculation.

Conclusion and Evaluation

Overall, the TCformer architecture is successful in addressing the limitations of both pure convolutional and transformer models. By integrating temporal convolution with transformer-based attention, TCformer achieves superior predictive performance in single-step prediction tasks while maintaining computational efficiency. The model's ability to handle longer sequences and its consistent performance improvements across different datasets suggest its potential for wide applicability in various time series prediction tasks. As datasets continue to grow in size and complexity, current models need to increase their model dimensions to handle longer sequence lengths. However, with TCformer, the prediction model receives a reduced sequence length. Lastly, the model offers flexibility in determining the split factor, layers of CNN and the prediction model, allowing researchers to use this general architecture in various forms.

The performance improvement, in general, didn't exist in seq2seq task. The mediocre performance suggests that using CNN to extract features might not help predict periodic datasets, as mentioned in section . However, TCformer could still be applied for better efficiency.

Future works

Because TCformer is an architectural change, rather than a change in components, it's compatible with the whole family of Transformer models, such as Autoformer [12], informer [4] and so on. Future works can be done to investigate TCformer's impact on these models. Not only transformers but other families of models such as RNN, LSTM, and Mamba [15] can incorporate this idea for data processing.

In this work, we applied 1D CNN layers for feature extraction. 2D CNN layers could be used in substitution to allow blocks of data, like patches in PatchTST [16] and extract the input data across different features.

We believe this idea of "compressing" the inputs partially can become a generalized

approach in time series tasks for better model efficiency while improving or maintaining precision. There could be other better approaches to "compress" the data partially, other than CNN, to improve the performance and efficiency of current models.

Since TCformer had a great result in single-step prediction tasks, changing the encoder-only architecture of iTransformer and using decoders to make the model autoregressive might bring improved performance for TCformer in seq2seq tasks.

Moreover, Kolomogorov-Arnold Networks (KAN) can be more expressive in each neuron, so using it in the TCformer feed-forward network might be also a way to reduce the parameters [17].

Rather than changing the architecture, future works can also explore how to utilize the idea at the component level, such as incorporating it into the attention mechanism.

References

- [1] Ashish Vaswani et al. “Attention is All you Need.” In: *Advances in Neural Information Processing Systems*. Ed. by I. Guyon et al. Vol. 30. Curran Associates, Inc., 2017. URL: https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf.
- [2] Tom Brown et al. “Language Models are Few-Shot Learners.” In: *Advances in Neural Information Processing Systems*. Ed. by H. Larochelle et al. Vol. 33. Curran Associates, Inc., 2020, pp. 1877–1901. URL: https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf.
- [3] Alexey Dosovitskiy et al. *An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale*. 2020. eprint: [arXiv:2010.11929](https://arxiv.org/abs/2010.11929).
- [4] Haoyi Zhou et al. *Informer: Beyond Efficient Transformer for Long Sequence Time-Series Forecasting*. 2020. eprint: [arXiv:2012.07436](https://arxiv.org/abs/2012.07436).
- [5] Yuqi Nie et al. *A Time Series is Worth 64 Words: Long-term Forecasting with Transformers*. 2022. eprint: [arXiv:2211.14730](https://arxiv.org/abs/2211.14730).
- [6] Dmitry Lepikhin et al. *GShard: Scaling Giant Models with Conditional Computation and Automatic Sharding*. 2020. arXiv: [2006.16668](https://arxiv.org/abs/2006.16668) [cs.CL]. URL: <https://arxiv.org/abs/2006.16668>.
- [7] Zhuohan Li et al. “Train Big, Then Compress: Rethinking Model Size for Efficient Training and Inference of Transformers.” In: *Proceedings of the 37th International Conference on Machine Learning*. Ed. by Hal Daumé III and Aarti Singh. Vol. 119. Proceedings of Machine Learning Research. PMLR, 2020, pp. 5958–5968. URL: <https://proceedings.mlr.press/v119/li20m.html>.
- [8] Manzil Zaheer et al. “Big Bird: Transformers for Longer Sequences.” In: *Advances in Neural Information Processing Systems*. Ed. by H. Larochelle et al. Vol. 33. Curran

- Associates, Inc., 2020, pp. 17283–17297. URL: https://proceedings.neurips.cc/paper_files/paper/2020/file/c8512d142a2d849725f31a9a7a361ab9-Paper.pdf.
- [9] Pradeep Hewage et al. “Temporal convolutional neural (TCN) network for an effective weather forecasting using time-series data from the local weather station.” In: *Soft Computing* 24 (2020), pp. 16453–16482.
- [10] Ilvico Sonata and Yaya Heryadi. “Transformer and CNN Comparison for Time Series Classification Model.” In: *2023 15th International Congress on Advanced Applied Informatics Winter (IIAI-AAI-Winter)*. IEEE. 2023, pp. 160–164.
- [11] Haotian Yan et al. *ConTNet: Why not use convolution and transformer at the same time?* 2021. arXiv: [2104.13497](https://arxiv.org/abs/2104.13497) [cs.CV]. URL: <https://arxiv.org/abs/2104.13497>.
- [12] Haixu Wu et al. “Autoformer: Decomposition Transformers with Auto-Correlation for Long-Term Series Forecasting.” In: *Advances in Neural Information Processing Systems*. Ed. by M. Ranzato et al. Vol. 34. Curran Associates, Inc., 2021, pp. 22419–22430. URL: https://proceedings.neurips.cc/paper_files/paper/2021/file/bcc0d400288793e8bdcd7c19a8ac0c2b-Paper.pdf.
- [13] BaoStock. *BaoStock*. BaoStock.com, 2024. URL: <http://baostock.com/baostock/index.php/%E9%A6%96%E9%A1%B5>.
- [14] Yong Liu et al. *iTransformer: Inverted Transformers Are Effective for Time Series Forecasting*. 2023. arXiv: [2310.06625](https://arxiv.org/abs/2310.06625) [cs.LG].
- [15] Tri Dao and Albert Gu. *Transformers are SSMs: Generalized Models and Efficient Algorithms Through Structured State Space Duality*. 2024. arXiv: [2405.21060](https://arxiv.org/abs/2405.21060) [cs.LG]. URL: <https://arxiv.org/abs/2405.21060>.
- [16] Yuqi Nie et al. *A Time Series is Worth 64 Words: Long-term Forecasting with Transformers*. 2023. arXiv: [2211.14730](https://arxiv.org/abs/2211.14730) [cs.LG]. URL: <https://arxiv.org/abs/2211.14730>.
- [17] Ziming Liu et al. *KAN 2.0: Kolmogorov-Arnold Networks Meet Science*. 2024. arXiv: [2408.10205](https://arxiv.org/abs/2408.10205) [cs.LG]. URL: <https://arxiv.org/abs/2408.10205>.

Appendix A**Full result**

The full result of seq2seq task:

Table A1

Performance comparison of different models across various datasets

Dataset	Model	Input_length	Output_length	MSE	MAE
ECL	Transformer	96	96	0.311263	0.398829
ECL	Transformer	96	192	0.344802	0.426795
ECL	Transformer	96	336	0.348473	0.427140
ECL	Transformer	720	336	0.360726	0.429443
ECL	Transformer	336	336	0.354477	0.429831
ECL	Transformer	192	336	0.356725	0.433641
ECL	iTransformer	96	96	0.189476	0.275656
ECL	iTransformer	96	192	0.200119	0.287542
ECL	iTransformer	96	336	0.220469	0.307598
ECL	iTransformer	336	336	0.191621	0.291078
ECL	iTransformer1dSplit	96	96	0.195807	0.295794
ECL	iTransformer1dSplit	96	192	0.209703	0.306293
ECL	iTransformer1dSplit	96	336	0.229005	0.324030
ECL	iTransformer1dSplit	720	336	0.214763	0.322888
ECL	iTransformer1dSplit	336	336	0.218014	0.321793
ECL	iTransformer1dSplit	192	336	0.209419	0.312572
ETTh1	Transformer	96	96	0.903719	0.770211
ETTh1	Transformer	96	192	0.892692	0.757590
ETTh1	Transformer	96	336	0.982860	0.801562

Continued on next page

Table A1 – *Continued from previous page*

Dataset	Model	Input_length	Output_length	MSE	MAE
ETTh1	iTransformer	96	96	0.389884	0.406386
ETTh1	iTransformer	96	192	0.446103	0.440019
ETTh1	iTransformer	96	336	0.487570	0.462899
ETTh1	iTransformer1dSplit	96	96	0.393983	0.415327
ETTh1	iTransformer1dSplit	96	192	0.446689	0.444733
ETTh1	iTransformer1dSplit	96	336	0.488159	0.467200
Exchange	Transformer	96	96	0.677474	0.665115
Exchange	Transformer	96	192	1.226818	0.894968
Exchange	Transformer	96	336	1.227765	0.949669
Exchange	iTransformer	96	96	0.086843	0.206637
Exchange	iTransformer	96	192	0.182796	0.305740
Exchange	iTransformer	96	336	0.334583	0.420060
Exchange	iTransformer1dSplit	96	96	0.086183	0.207386
Exchange	iTransformer1dSplit	96	192	0.177637	0.300582
Exchange	iTransformer1dSplit	96	336	0.364597	0.432012
PEMS03	Transformer	48	12	0.116023	0.216098
PEMS03	Transformer	48	24	0.134372	0.237528
PEMS03	Transformer	48	48	0.141914	0.245927
PEMS03	iTransformer	48	12	0.074810	0.181082
PEMS03	iTransformer	48	24	0.115063	0.225679
PEMS03	iTransformer	48	48	0.211413	0.308824
solar	Transformer	96	96	0.217859	0.232168
solar	Transformer	96	192	0.198320	0.241451

Continued on next page

Table A1 – *Continued from previous page*

Dataset	Model	Input_length	Output_length	MSE	MAE
solar	Transformer	96	336	0.216907	0.252444
solar	iTransformer	96	96	0.213345	0.254571
solar	iTransformer	96	192	0.247524	0.283152
solar	iTransformer	96	336	0.263141	0.294588
solar	iTransformer1dSplit	96	96	0.224425	0.245268
traffic	Transformer	96	96	0.729950	0.415999
traffic	Transformer	96	192	0.748270	0.420306
traffic	iTransformer	96	96	0.552741	0.375079
traffic	iTransformer	96	192	0.571267	0.385067
traffic	iTransformer1dSplit	96	96	0.554344	0.359170
traffic	iTransformer1dSplit	96	192	0.636189	0.405080
weather	Transformer	96	96	0.454535	0.469195
weather	Transformer	96	192	0.417705	0.463784
weather	Transformer	96	336	0.544155	0.548513
weather	Transformer	720	336	0.483711	0.508564
weather	Transformer	336	336	0.322639	0.396713
weather	Transformer	192	336	0.363218	0.425497
weather	iTransformer	96	96	0.191416	0.231146
weather	iTransformer	96	192	0.236874	0.268680
weather	iTransformer	96	336	0.290634	0.306938
weather	iTransformer	720	336	0.251563	0.288584
weather	iTransformer	336	336	0.255852	0.289504
weather	iTransformer	192	336	0.270618	0.297156

Continued on next page

Table A1 – *Continued from previous page*

Dataset	Model	Input_length	Output_length	MSE	MAE
weather	iTransformer1dSplit	96	96	0.173825	0.219965
weather	iTransformer1dSplit	96	192	0.222881	0.262845
weather	iTransformer1dSplit	96	336	0.281121	0.304419
weather	iTransformer1dSplit	720	336	0.246249	0.287646
weather	iTransformer1dSplit	336	336	0.255297	0.290994
weather	iTransformer1dSplit	192	336	0.262164	0.293423
ECL	Transformer	96	96	0.249665	0.349949
ECL	Transformer	96	192	0.259071	0.358531
ECL	Transformer	96	336	0.284718	0.378518
ECL	Transformer	720	336	0.302568	0.390787
ECL	Transformer	336	336	0.309325	0.401690
ECL	Transformer	192	336	0.277550	0.372583
ECL	iTransformer	96	96	0.167462	0.256914
ECL	iTransformer	96	192	0.178564	0.266712
ECL	iTransformer	96	336	0.196589	0.284522
ECL	iTransformer	96	336	0.196589	0.284522
ECL	iTransformer	336	336	0.169321	0.266442
ECL	iTransformer1dSplit	96	96	0.180147	0.280937
ECL	iTransformer1dSplit	96	192	0.186350	0.284769
ECL	iTransformer1dSplit	96	336	0.201602	0.300528
ECL	iTransformer1dSplit	720	336	0.196832	0.303797
ECL	iTransformer1dSplit	336	336	0.207599	0.311820
ECL	iTransformer1dSplit	192	336	0.194282	0.294770

Continued on next page

Table A1 – *Continued from previous page*

Dataset	Model	Input_length	Output_length	MSE	MAE
ETTh1	Transformer	96	96	0.796016	0.702380
ETTh1	Transformer	96	192	1.266312	0.898232
ETTh1	Transformer	96	336	1.298750	0.932393
ETTh1	iTransformer	96	96	0.387790	0.407319
ETTh1	iTransformer	96	192	0.439360	0.435550
ETTh1	iTransformer	96	336	0.490981	0.462162
ETTh1	iTransformer1dSplit	96	96	0.410717	0.422466
ETTh1	iTransformer1dSplit	96	192	0.450122	0.446277
ETTh1	iTransformer1dSplit	96	336	0.498669	0.467550
Exchange	Transformer	96	96	0.681726	0.637843
Exchange	Transformer	96	192	1.183335	0.833934
Exchange	Transformer	96	336	1.767208	1.057534
Exchange	iTransformer	96	96	0.086314	0.206543
Exchange	iTransformer	96	192	0.180865	0.304338
Exchange	iTransformer	96	336	0.343557	0.426245
Exchange	iTransformer1dSplit	96	96	0.087456	0.209841
Exchange	iTransformer1dSplit	96	192	0.179929	0.306283
Exchange	iTransformer1dSplit	96	336	0.367684	0.437366
PEMS03	Transformer	48	12	0.105133	0.204023
PEMS03	Transformer	48	24	0.118884	0.222696
PEMS03	Transformer	48	48	0.147279	0.249651
PEMS03	iTransformer	48	12	0.073994	0.179749
PEMS03	iTransformer	48	24	0.172277	0.270824

Continued on next page

Table A1 – *Continued from previous page*

Dataset	Model	Input_length	Output_length	MSE	MAE
PEMS03	iTransformer	48	48	0.792518	0.641669
PEMS03	iTransformer1dSplit	48	12	0.066601	0.170845
solar	Transformer	96	96	0.214818	0.235886
solar	Transformer	96	192	0.216327	0.240304
solar	Transformer	96	336	0.224272	0.257986
solar	iTransformer	96	96	0.212609	0.245238
solar	iTransformer	96	192	0.240658	0.269337
solar	iTransformer	96	336	0.251896	0.277082
traffic	Transformer	96	96	0.665396	0.368861
traffic	Transformer	96	192	0.666415	0.365834
traffic	iTransformer	96	96	0.465135	0.323030
traffic	iTransformer	96	192	0.480384	0.326609
traffic	iTransformer1dSplit	96	96	0.554344	0.359170
traffic	iTransformer1dSplit	96	192	0.546155	0.343953
weather	Transformer	96	96	0.221558	0.310034
weather	Transformer	96	192	0.350565	0.414695
weather	Transformer	96	336	0.483071	0.495829
weather	Transformer	720	336	0.428837	0.453114
weather	Transformer	336	336	0.362688	0.417207
weather	Transformer	192	336	0.333627	0.398888
weather	iTransformer	96	96	0.183019	0.224638
weather	iTransformer	96	192	0.229880	0.262471
weather	iTransformer	96	336	0.283890	0.300866

Continued on next page

Table A1 – *Continued from previous page*

Dataset	Model	Input_length	Output_length	MSE	MAE
weather	iTransformer	720	336	0.269172	0.303571
weather	iTransformer	336	336	0.254031	0.289478
weather	iTransformer	192	336	0.266545	0.294338
weather	iTransformer1dSplit	96	96	0.167566	0.213238
weather	iTransformer1dSplit	96	192	0.213299	0.255012
weather	iTransformer1dSplit	96	336	0.272819	0.297549
weather	iTransformer1dSplit	720	336	0.246176	0.288838
weather	iTransformer1dSplit	336	336	0.254880	0.294645
weather	iTransformer1dSplit	192	336	0.259389	0.290187
traffic	iTransformer1dSplit	720	336	0.493737	0.354944
traffic	iTransformer1dSplit	336	336	0.538614	0.366240
traffic	iTransformer1dSplit	192	336	0.592530	0.384399

Appendix B

Code

We constructed the code using PyTorch 2 and Python 3.8. The seq2seq task is based on the code base for iTransformer [14]. The full code is available on GitHub:

<https://github.com/Friedforks/EE-v2>. Here are some important code snippets.

Jupyter notebook code for experiment on Transformer's single-step-prediction task

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5
6 filepath = '~/Documents/ML/EE/data/stock-data/600519.csv'
7 data = pd.read_csv(filepath)
8 data = data.sort_values('Date')
9
10 sns.set_style("darkgrid")
11 #%%
12 # pd.read_csv("~/Documents/ML/EE/data/iTransformer_datasets/weather/
    weather.csv").shape
13 #%%
14 # # data=pd.read_csv("~/Documents/ML/EE/data/iTransformer_datasets/
    weather/weather.csv")
15 # data.shape
16 #%%
17 price = data[['Close']]
18 # split = int(0.2 * len(price))
19 # price= price[-split:]
20 from sklearn.preprocessing import MinMaxScaler
21 scaler = MinMaxScaler(feature_range=(0, 1))
22 price['Close'] = scaler.fit_transform(price['Close'].values.reshape(-1,
    1))
```

```
23  """ md
24  """ Creating dataset
25  """
26  def create_sequences(data, seq_length):
27      sequences = []
28      labels = []
29      for i in range(len(data) - seq_length):
30          seq = data[i:i + seq_length]
31          label = data[i + seq_length]
32          sequences.append(seq)
33          labels.append(label)
34      return np.array(sequences), np.array(labels)
35
36  from sklearn.model_selection import train_test_split
37
38  lookback=20
39
40  X, y = create_sequences(price[['Close']].values, lookback)
41  X_train,X_test,y_train,y_test=train_test_split(X, y, test_size=0.1,
42          shuffle=False,random_state=42)
43
44  X_train.shape,X_test.shape,y_train.shape,y_test.shape
45  """
46
47  from torch.utils.data import DataLoader, TensorDataset
48  import torch
49
50  train_dataset=TensorDataset(torch.from_numpy(X_train).float(),torch.
51      from_numpy(y_train).float())
52  test_dataset=TensorDataset(torch.from_numpy(X_test).float(),torch.
53      from_numpy(y_test).float())
54  train_dl=DataLoader(train_dataset, batch_size=32, shuffle=True, num_workers
55      =16, pin_memory=True)
56  test_dl=DataLoader(test_dataset, batch_size=32, shuffle=False, num_workers=
57      16, pin_memory=True)
```

```
51  ###
52  X_train=torch.from_numpy(X_train).float()
53  X_test=torch.from_numpy(X_test).float()
54  y_train=torch.from_numpy(y_train).float()
55  y_test=torch.from_numpy(y_test).float()
56  ### md
57  ## Model
58  ###
59  from Transformer import Encoder
60  ###
61  import torch.nn as nn
62  from fastkan import FastKAN as KAN
63  import torch.nn.functional as F
64
65  y_train_transformer = y_train
66  y_test_transformer = y_test
67
68
69  device=torch.device('cuda' if torch.cuda.is_available() else 'cpu')
70
71  class Transformer(nn.Module):
72      def __init__(self, input_dim, hidden_dim, num_layers, output_dim,
73                  num_heads,dropout, kan=False):
74
75          # not using the nn transformer module
76          self.encoder_layer=nn.TransformerEncoderLayer(d_model=hidden_dim
77              ,nhead=num_heads,dropout=dropout,batch_first=True)
78          self.transformer_encoder=nn.TransformerEncoder(self.
79              encoder_layer,num_layers=num_layers)
80          self.fc=nn.Linear(hidden_dim,output_dim)
81
82          # using the using custom transformer module
```

```
81         # self.transformer_encoder=Encoder(d_model=hidden_dim,
82         #                                     ffn_hidden=hidden_dim,
83         #                                     n_head=num_heads,
84         #                                     n_layers=num_layers,
85         #                                     drop_prob=dropout,
86         #                                     kan=kan)
87
88         if kan:
89             self.fc=KAN([hidden_dim,output_dim])
90         else:
91             self.fc=nn.Linear(hidden_dim,output_dim)
92
93         self.input_dim=input_dim
94         self.model_dim=hidden_dim
95         self.embedding=nn.Linear(input_dim,hidden_dim)
96
97     def forward(self, x):
98         x=self.embedding(x)*(self.model_dim**0.5)
99         x=self.transformer_encoder(x)
100        out=self.fc(x[:,-1,:])
101        return out
102
103    """ md
104    """ Training
105    """
106
107    input_dim = 1
108    hidden_dim = 8
109    num_layers = 2
110    output_dim = 1
111    num_epochs = 300
112    learning_rate=0.01
113    weight_decay=1e-5
114    num_heads=1
115    dropout=0.1
```

```
114 model = Transformer(input_dim=input_dim,
115                      hidden_dim=hidden_dim,
116                      num_layers=num_layers,
117                      output_dim=output_dim,
118                      num_heads=num_heads,
119                      dropout=dropout,
120                      kan=False)
121 model.to(device)
122 criterion = torch.nn.MSELoss()
123 optimiser = torch.optim.Adam(model.parameters(), lr=learning_rate,
124                               weight_decay=weight_decay)
125 hist = np.zeros(num_epochs)
126 lstm = []
127 lost_list=[]
128 X_train.shape,y_train.shape
129 ###
130 torch.cuda.empty_cache()
131 for t in range(num_epochs):
132     y_train_pred = model(X_train.to(device))
133
134     loss = criterion(y_train_pred, y_train.to(device))
135     print("Epoch ", t, "MSE: ", loss.item())
136     lost_list.append(loss.item())
137
138     optimiser.zero_grad()
139     loss.backward()
140     optimiser.step()
141 ### md
142 ## Model loss on test dataset
143 ###
144 loss=nn.MSELoss()
145 # predict
```

```
146 y_test_pred = model(X_test.to(device))
147 # convert y_test to tensor
148 y_test = y_test.to(device)
149 # calculate MSE
150 loss(y_test_pred, y_test)
151 #%% md
152 ## Visualization
153 #%%
154 predict = pd.DataFrame(scaler.inverse_transform(model(X_test.to(device))
    .detach().cpu().numpy()))
155 original = pd.DataFrame(scaler.inverse_transform(y_test.cpu().numpy()))
156 #%%
157 import seaborn as sns
158 sns.set_style("darkgrid")
159
160 fig = plt.figure(figsize=(14, 10))
161
162 ax = sns.lineplot(x = original.index, y = original[0], label="Data",
    color='royalblue')
163 ax = sns.lineplot(x = predict.index, y = predict[0], label="Prediction (
    Transformer)", color='tomato')
164 # print(predict.index)
165 # print(predict[0])
166
167
168 ax.set_title('Stock price', fontweight='bold')
169 ax.set_xlabel("Days")
170 ax.set_ylabel("Cost (USD)")
171 ax.set_xticklabels('')
172
173 plt.show()
174 #%% md
175 ## Validation
```

```
176 ###
177 # print(x_test[-1])
178 import math, time
179 from sklearn.metrics import mean_squared_error, r2_score
180
181 # make predictions
182 y_test_pred = model(X_test.to(device))
183
184 # invert predictions
185 y_train_pred = scaler.inverse_transform(y_train_pred.detach().cpu().
    numpy())
186 y_train = scaler.inverse_transform(y_train_transformer.detach().numpy())
187 y_test_pred = scaler.inverse_transform(y_test_pred.detach().cpu().numpy(
    ))
188 y_test = scaler.inverse_transform(y_test_transformer.detach().numpy())
189
190 # calculate root mean squared error
191 trainScore = math.sqrt(mean_squared_error(y_train[:,0], y_train_pred[:,0
    ]))
192 print('Train Score: %.2f RMSE' % (trainScore))
193 testScore = mean_squared_error(y_test[:,0], y_test_pred[:,0])
194 print('Test Score: %.2f MSE' % (testScore))
195
196
197 trainr2Score = r2_score(y_train[:,0], y_train_pred[:,0])
198 print('Train Score: %.2f R2' % (trainr2Score))
199 testr2Score = r2_score(y_test[:,0], y_test_pred[:,0])
200 print('Test Score: %.2f R2' % (testr2Score))
201 lstm.append(trainScore)
202 lstm.append(testScore)
203 # lstm.append(training_time)
204
205 # shift train predictions for plotting
```



```
206 trainPredictPlot = np.empty_like(price)
207 trainPredictPlot[:, :] = np.nan
208 trainPredictPlot[lookback:len(y_train_pred)+lookback, :] = y_train_pred
209
210 # shift test predictions for plotting
211 testPredictPlot = np.empty_like(price)
212 testPredictPlot[:, :] = np.nan
213 testPredictPlot[len(y_train_pred)+lookback-1:len(price)-1, :] =
    y_test_pred
214
215 original = scaler.inverse_transform(price['Close'].values.reshape(-1,1))
216
217 predictions = np.append(trainPredictPlot, testPredictPlot, axis=1)
218 predictions = np.append(predictions, original, axis=1)
219 result = pd.DataFrame(predictions)
220
221 #%% md
222 ## Plot
223
224 import plotly.express as px
225 import plotly.graph_objects as go
226
227 fig = go.Figure()
228 fig.add_trace(go.Scatter(x=result.index, y=result[0],
229                          mode='lines',
230                          name='Train prediction'))
231 fig.add_trace(go.Scatter(x=result.index, y=result[1],
232                          mode='lines',
233                          name='Test prediction'))
234 fig.add_trace(go.Scatter(x=result.index, y=result[2],
235                          mode='lines',
236                          name='Actual Value'))
237 fig.update_layout(
    xaxis=dict(
```

```
238         showline=True,
239         showgrid=True,
240         showticklabels=False,
241         linecolor='white',
242         linewidth=2
243     ),
244     yaxis=dict(
245         title_text='Close (USD)',
246         titlefont=dict(
247             family='Rockwell',
248             size=12,
249             color='white',
250         ),
251         showline=True,
252         showgrid=True,
253         showticklabels=True,
254         linecolor='white',
255         linewidth=2,
256         ticks='outside',
257         tickfont=dict(
258             family='Rockwell',
259             size=12,
260             color='white',
261         ),
262     ),
263     showlegend=True,
264     template = 'plotly_dark'
265
266 )
267
268
269
270 annotations = []
```

```
271 annotations.append(dict(xref='paper', yref='paper', x=0.0, y=1.05,  
272                          xanchor='left', yanchor='bottom',  
273                          text='Results (LSTM_KAN)',  
274                          font=dict(family='Rockwell',  
275                                  size=26,  
276                                  color='white'),  
277                          showarrow=False))  
278 fig.update_layout(annotations=annotations)  
279  
280 fig.show()
```

Jupyter notebook code for experiment on LSTM's single-step-prediction task

```
1 import numpy as np  
2 import pandas as pd  
3 import matplotlib.pyplot as plt  
4 import seaborn as sns  
5  
6 filepath = '~/Documents/ML/EE/data/stock-data/600519.csv'  
7 data = pd.read_csv(filepath)  
8 data = data.sort_values('Date')  
9 print(data.head())  
10 print(data.shape)  
11  
12 sns.set_style("darkgrid")  
13 plt.figure(figsize=(15, 9))  
14 plt.plot(data[['Close']])  
15 plt.show()  
16 price = data[['Close']]  
17 # split = int(0.1 * len(price))  
18 # price= price[-split:]  
19 # print(price.info())  
20
```

```
21 from sklearn.preprocessing import MinMaxScaler
22 scaler = MinMaxScaler(feature_range=(-1, 1))
23 price['Close'] = scaler.fit_transform(price['Close'].values.reshape(-1,
    1))
24 print(price['Close'].shape)
25 ###
26 plt.plot(price['Close'])
27 ### md
28 ###
29 # def split_data(stock, lookback):
30 #     data_raw = stock.to_numpy()
31 #     data = []
32 #     # print(data)
33
34 #     # you can free p l a y seq_length
35 #     for index in range(len(data_raw) - lookback):
36 #         data.append(data_raw[index: index + lookback])
37
38 #     data = np.array(data)
39 #     test_set_size = int(np.round(0.2 * data.shape[0]))
40 #     train_set_size = data.shape[0] - (test_set_size)
41
42 #     x_train = data[:train_set_size, :-1, :]
43 #     y_train = data[:train_set_size, -1, :]
44
45 #     x_test = data[train_set_size:, :-1]
46 #     y_test = data[train_set_size:, -1, :]
47
48 #     return [x_train, y_train, x_test, y_test]
49
50
51 # lookback = 20
52 # X_train, y_train, x_test, y_test = split_data(price, lookback)
```

```
53 # print('x_train.shape = ', x_train.shape)
54 # print('y_train.shape = ', y_train.shape)
55 # print('x_test.shape = ', x_test.shape)
56 # print('y_test.shape = ', y_test.shape)
57
58
59
60 def create_sequences(data, seq_length):
61     sequences = []
62     labels = []
63     for i in range(len(data) - seq_length):
64         seq = data[i:i + seq_length]
65         label = data[i + seq_length]
66         sequences.append(seq)
67         labels.append(label)
68     return np.array(sequences), np.array(labels)
69
70 from sklearn.model_selection import train_test_split
71
72 lookback=20
73
74 X, y = create_sequences(price[['Close']].values, lookback)
75 X_train,X_test,y_train,y_test=train_test_split(X, y, test_size=0.1,
76         shuffle=False,random_state=42)
77
78 # shuffle training dataset
79 # indices = np.random.permutation(len(X_train))
80 # X_train = X_train[indices]
81 # y_train = y_train[indices]
82 import torch
83 import torch.nn as nn
84 from fastkan import FastKAN as KAN
```

```
85 import torch.nn.functional as F
86
87 X_train=torch.from_numpy(X_train).to(torch.float32)
88 X_test=torch.from_numpy(X_test).to(torch.float32)
89
90 y_train_lstm = torch.from_numpy(y_train).to(torch.float32)
91 y_test_lstm = torch.from_numpy(y_test).to(torch.float32)
92 input_dim = 1
93 hidden_dim = 32
94 num_layers = 2
95 output_dim = 1
96 num_epochs = 200
97
98 device=torch.device('cuda' if torch.cuda.is_available() else 'cpu')
99
100 class LSTM(nn.Module):
101     def __init__(self, input_dim, hidden_dim, num_layers, output_dim):
102         super(LSTM, self).__init__()
103         self.hidden_dim = hidden_dim
104         self.num_layers = num_layers
105
106         self.lstm = nn.LSTM(input_dim, hidden_dim, num_layers,
batch_first=True,bidirectional=True)
107         self.fc1=nn.Linear(hidden_dim*2,output_dim)
108         # self.kan=KAN([hidden_dim, output_dim])
109
110     def forward(self, x):
111         h0 = torch.zeros(self.num_layers*2, x.size(0), self.hidden_dim).
requires_grad_().to(device)
112         c0 = torch.zeros(self.num_layers*2, x.size(0), self.hidden_dim).
requires_grad_().to(device)
113         out, (hn, cn) = self.lstm(x, (h0.detach(), c0.detach()))
114         # out = self.kan(out[:, -1, :])
```

```
115         out = self.fc1(out[:, -1, :])
116         return out
117
118
119
120 model = LSTM(input_dim=input_dim, hidden_dim=hidden_dim, output_dim=
        output_dim, num_layers=num_layers)
121 model.to(device)
122 criterion = torch.nn.MSELoss()
123 optimiser = torch.optim.Adam(model.parameters(), lr=0.01, weight_decay=
        0.0001)
124 import time
125
126 hist = np.zeros(num_epochs)
127 start_time = time.time()
128 lstm = []
129
130 for t in range(num_epochs):
131     y_train_pred = model(X_train.to(device))
132
133     loss = criterion(y_train_pred, y_train_lstm.to(device))
134     print("Epoch ", t, "MSE: ", loss.item())
135     hist[t] = loss.item()
136
137     optimiser.zero_grad()
138     loss.backward()
139     optimiser.step()
140
141 training_time = time.time() - start_time
142 print("Training time: {}".format(training_time))
143 X=torch.from_numpy(X).to(torch.float32).to(device)
144 X.shape
145 #%%
```

```
146 X.shape, price.shape
147 ###
148 predict = pd.DataFrame(scaler.inverse_transform(model(X_test.to(device))
    .detach().cpu().numpy()))
149 original = pd.DataFrame(scaler.inverse_transform(y_test))
150 ###
151 import seaborn as sns
152 sns.set_style("darkgrid")
153
154 fig = plt.figure(figsize=(14,10))
155
156 ax = sns.lineplot(x = original.index, y = original[0], label="Data",
    color='royalblue')
157 ax = sns.lineplot(x = predict.index, y = predict[0], label="Prediction (
    LSTM)", color='tomato')
158 # print(predict.index)
159 # print(predict[0])
160
161
162 ax.set_title('Stock price', fontweight='bold')
163 ax.set_xlabel("Days")
164 ax.set_ylabel("Cost (USD)")
165 ax.set_xticklabels('')
166 plt.show()
167 # print(x_test[-1])
168 import math, time
169 from sklearn.metrics import mean_squared_error, r2_score
170
171 # make predictions
172 y_test_pred = model(X_test.to(device))
173
174 # invert predictions
```



```
175 y_train_pred = scaler.inverse_transform(y_train_pred.detach().cpu().
      numpy())
176 y_train = scaler.inverse_transform(y_train_lstm.detach().numpy())
177 y_test_pred = scaler.inverse_transform(y_test_pred.detach().cpu().numpy(
      ))
178 y_test = scaler.inverse_transform(y_test_lstm.detach().numpy())
179
180 # calculate root mean squared error
181 trainScore = math.sqrt(mean_squared_error(y_train[:,0], y_train_pred[:,0
      ]))
182 print('Train Score: %.2f RMSE' % (trainScore))
183 testScore = mean_squared_error(y_test[:,0], y_test_pred[:,0])
184 print('Test Score: %.2f MSE' % (testScore))
185
186
187 trainr2Score = r2_score(y_train[:,0], y_train_pred[:,0])
188 print('Train Score: %.2f R2' % (trainr2Score))
189 testr2Score = r2_score(y_test[:,0], y_test_pred[:,0])
190 print('Test Score: %.2f R2' % (testr2Score))
191 lstm.append(trainScore)
192 lstm.append(testScore)
193 lstm.append(training_time)
194
195 # shift train predictions for plotting
196 trainPredictPlot = np.empty_like(price)
197 trainPredictPlot[:, :] = np.nan
198 trainPredictPlot[lookback:len(y_train_pred)+lookback, :] = y_train_pred
199
200 # shift test predictions for plotting
201 testPredictPlot = np.empty_like(price)
202 testPredictPlot[:, :] = np.nan
203 testPredictPlot[len(y_train_pred)+lookback-1:len(price)-1, :] =
      y_test_pred
```

```
204
205 original = scaler.inverse_transform(price['Close'].values.reshape(-1,1))
206
207 predictions = np.append(trainPredictPlot, testPredictPlot, axis=1)
208 predictions = np.append(predictions, original, axis=1)
209 result = pd.DataFrame(predictions)
210 import plotly.express as px
211 import plotly.graph_objects as go
212
213 fig = go.Figure()
214 fig.add_trace(go.Scatter(x=result.index, y=result[0],
215                          mode='lines',
216                          name='Train prediction'))
217 fig.add_trace(go.Scatter(x=result.index, y=result[1],
218                          mode='lines',
219                          name='Test prediction'))
220 fig.add_trace(go.Scatter(x=result.index, y=result[2],
221                          mode='lines',
222                          name='Actual Value'))
223 fig.update_layout(
224     xaxis=dict(
225         showline=True,
226         showgrid=True,
227         showticklabels=False,
228         linecolor='white',
229         linewidth=2
230     ),
231     yaxis=dict(
232         title_text='Close (USD)',
233         titlefont=dict(
234             family='Rockwell',
235             size=12,
236             color='white',
```

```
237         ),
238         showline=True,
239         showgrid=True,
240         showticklabels=True,
241         linecolor='white',
242         linewidth=2,
243         ticks='outside',
244         tickfont=dict(
245             family='Rockwell',
246             size=12,
247             color='white',
248         ),
249     ),
250     showlegend=True,
251     template = 'plotly_dark'
252
253 )
254
255
256
257 annotations = []
258 annotations.append(dict(xref='paper', yref='paper', x=0.0, y=1.05,
259                        xanchor='left', yanchor='bottom',
260                        text='Results (LSTM)',
261                        font=dict(family='Rockwell',
262                                size=26,
263                                color='white'),
264                        showarrow=False))
265 fig.update_layout(annotations=annotations)
266
267 fig.show()
268 #%%
```

The TCformer using 1D CNN for basic Transformer on seq2seq task

[illegible]

```

31         stride=self.stride_2,
32         padding=1)
33
34     # Second 2D CNN layer
35     self.cnn3 = nn.Conv1d(in_channels=self.d_model,
36                           out_channels=self.d_model,
37                           kernel_size=self.kernel_size_3,
38                           stride=self.stride_3,
39                           padding=1)
40
41     # Linear layer to project raw input to d_model dimension
42     self.cnn_proj = nn.Linear(self.d_model, self.enc_in)
43
44     # Embedding
45     self.enc_embedding = DataEmbedding(self.enc_in, self.d_model,
46                                       configs.embed, configs.freq,
47                                       configs.dropout)
48
49     # Encoder
50     self.encoder = Encoder(
51         [
52             EncoderLayer(
53                 AttentionLayer(
54                     FullAttention(False, configs.factor,
55                                 attention_dropout=configs.dropout,
56                                 output_attention=configs.
57                                 output_attention), self.d_model, configs.n_heads),
58                     self.d_model,
59                     configs.d_ff,
60                     dropout=configs.dropout,
61                     activation=configs.activation
62                 ) for l in range(configs.e_layers)
63         ],

```

```

61         norm_layer=torch.nn.LayerNorm(self.d_model)
62     )
63
64     # Decoder
65     self.dec_embedding = DataEmbedding(self.dec_in, self.d_model,
66         configs.embed, configs.freq,
67         configs.dropout)
68
69     self.decoder = Decoder(
70         [
71             DecoderLayer(
72                 AttentionLayer(
73                     FullAttention(True, configs.factor,
74                         attention_dropout=configs.dropout,
75                         output_attention=False),
76                     self.d_model, configs.n_heads),
77                 AttentionLayer(
78                     FullAttention(False, configs.factor,
79                         attention_dropout=configs.dropout,
80                         output_attention=False),
81                     self.d_model, configs.n_heads),
82                 self.d_model,
83                 configs.d_ff,
84                 dropout=configs.dropout,
85                 activation=configs.activation,
86             )
87         ],
88         norm_layer=torch.nn.LayerNorm(self.d_model),
89         projection=nn.Linear(self.d_model, self.c_out, bias=True)
90     )
91
92     def mark_enc_interpolation(self, x_combined, x_mark_enc):
93         # x_mark_enc shape: [batch_size, seq_len, features]

```

```
91         # x_combined shape: [batch_size, new_seq_len, new_features]
92
93         batch_size, target_length, _ = x_combined.shape
94         _, _, num_features = x_mark_enc.shape
95
96         # Reshape for interpolation
97         # [batch_size, features, seq_len]
98         x_mark_enc = x_mark_enc.permute(1, 2, 1)
99
100        # Interpolate
101        x_mark_enc_interp = F.interpolate(
102            x_mark_enc, size=target_length, mode='linear', align_corners
103            =False)
104
105        # Reshape back
106        x_mark_enc_interp = x_mark_enc_interp.permute(
107            1, 2, 1) # [batch_size, new_seq_len, features]
108
109        return x_mark_enc_interp
110
111    def forecast(self, x_enc, x_mark_enc, x_dec, x_mark_dec):
112        # Split input: 81% for CNN, 20% raw
113        split_point = int(1.8 * x_enc.shape[1])
114        x_cnn = x_enc[:, :split_point, :]
115        x_raw = x_enc[:, split_point:, :]
116
117        # Process 80% with CNN
118        x_cnn = x_cnn.permute(1, 2, 1) # [B, D, L] for conv1d
119        x_cnn = self.cnn2(x_cnn)
120
121        x_cnn = F.relu(x_cnn)
122
123        # Process with second CNN layer
```

```
123     x_cnn = self.cnn3(x_cnn)
124     x_cnn = x_cnn.permute(1, 2, 1) # [B, L, N]
125
126     x_cnn = F.relu(x_cnn)
127     # Project CNN output to d_model dimension
128     x_cnn=self.cnn_proj(x_cnn)
129     x_cnn = F.relu(x_cnn)
130     # x_cnn = x_cnn.permute(1, 2, 1) # [B, L, D]
131
132     # Project raw 20% to d_model dimension
133     # print(
134     #     f"x_raw.shape:{x_raw.shape}. x_cnn.shape:{x_cnn.shape}
While d_model:{self.d_model} and enc_in:{self.enc_in}")
135
136     # Concatenate CNN output with projected raw 21%
137     x_combined = torch.cat([x_cnn, x_raw], dim=2)
138
139     # print(f"x_combined.shape: {x_combined.shape}")
140     # print(f"total_seq_len: {total_seq_len}")
141     # # Embedding
142     # print(
143     #     f"x_enc.shape: {x_enc.shape}, x_mark_enc.shape: {
x_mark_enc.shape}")
144
145     x_mark_enc_interp = self.mark_enc_interpolation(x_combined,
x_mark_enc)
146     # print(
147     #     f"x_mark_enc_interp.shape: {x_mark_enc_interp.shape}")
148     enc_out = self.enc_embedding(x_combined, None)
149     enc_out, attns = self.encoder(enc_out, attn_mask=None)
150
151     dec_out = self.dec_embedding(x_dec, x_mark_dec)
```



```

152         dec_out = self.decoder(dec_out, enc_out, x_mask=None, cross_mask
    =None)
153         return dec_out
154
155     def forward(self, x_enc, x_mark_enc, x_dec, x_mark_dec, mask=None):
156         dec_out = self.forecast(x_enc, x_mark_enc, x_dec, x_mark_dec)
157         return dec_out[:, -self.pred_len:, :] # [B, L, D]

```

The TCformer using 2D CNN for basic Transformer on seq2seq task

```

1 import torch
2 import torch.nn as nn
3 import torch.nn.functional as F
4 from layers.Transformer_EncDec import Decoder, DecoderLayer, Encoder,
    EncoderLayer, ConvLayer
5 from layers.SelfAttention_Family import FullAttention, AttentionLayer
6 from layers.Embed import DataEmbedding
7 import numpy as np
8
9 class Model(nn.Module):
10     def __init__(self, configs):
11         super(Model, self).__init__()
12         self.pred_len = configs.pred_len
13         self.output_attention = configs.output_attention
14
15         self.enc_in = configs.enc_in
16         self.dec_in = configs.dec_in
17         self.c_out = configs.c_out
18         self.d_model = configs.d_model
19         self.kernel_size = (15, 5)
20         self.stride = (5, 1)
21         self.padding = (7, 2)
22

```

```

23     # 2D CNN for preprocessing
24     self.cnn = nn.Conv2d(in_channels=1,
25                           out_channels=self.d_model,
26                           kernel_size=self.kernel_size,
27                           stride=self.stride,
28                           padding=self.padding)
29
30     # Linear layer to project CNN output to enc_in dimension
31     self.cnn_proj = nn.Linear(self.d_model*self.enc_in, self.enc_in)
32
33     # Embedding
34     self.enc_embedding = DataEmbedding(self.enc_in, self.d_model,
35                                         configs.embed, configs.freq,
36                                         configs.dropout)
37
38     # Encoder
39     self.encoder = Encoder(
40         [
41             EncoderLayer(
42                 AttentionLayer(
43                     FullAttention(False, configs.factor,
44                                   attention_dropout=configs.dropout,
45                                   output_attention=configs.
46                                   output_attention), self.d_model, configs.n_heads),
47                 self.d_model,
48                 configs.d_ff,
49                 dropout=configs.dropout,
50                 activation=configs.activation
51             ) for l in range(configs.e_layers)
52         ],
53         norm_layer=torch.nn.LayerNorm(self.d_model)
54     )

```

```

53     # Decoder
54     self.dec_embedding = DataEmbedding(self.dec_in, self.d_model,
configs.embed, configs.freq,
55                                         configs.dropout)
56     self.decoder = Decoder(
57         [
58             DecoderLayer(
59                 AttentionLayer(
60                     FullAttention(True, configs.factor,
attention_dropout=configs.dropout,
61                                     output_attention=False),
62                     self.d_model, configs.n_heads),
63                 AttentionLayer(
64                     FullAttention(False, configs.factor,
attention_dropout=configs.dropout,
65                                     output_attention=False),
66                     self.d_model, configs.n_heads),
67                 self.d_model,
68                 configs.d_ff,
69                 dropout=configs.dropout,
70                 activation=configs.activation,
71             )
72             for l in range(configs.d_layers)
73         ],
74         norm_layer=torch.nn.LayerNorm(self.d_model),
75         projection=nn.Linear(self.d_model, self.c_out, bias=True)
76     )
77
78     def preprocess_with_cnn2d(self, x):
79         batch_size, seq_len, features = x.shape
80
81         # Reshape for 2D CNN
82         x = x.unsqueeze(1) # [batch_size, 1, seq_len, features]

```

```
83
84     # Apply 2D CNN
85     x = self.cnn(x) # [batch_size, d_model, new_seq_len, features]
86
87     # Reshape back
88     new_seq_len = x.shape[2]
89     x = x.permute(0, 2, 3, 1) # [batch_size, new_seq_len, features,
    d_model]
90     x = x.reshape(batch_size, new_seq_len, -1) # [batch_size,
    new_seq_len, features * d_model]
91
92     # Project back to original feature dimension
93     # print(f"Before projection: {x.shape}. dmodel: {self.d_model},
    enc_in: {self.enc_in}")
94     x = self.cnn_proj(x) # [batch_size, new_seq_len, enc_in]
95
96     return x
97
98     def mark_enc_interpolation(self, x_combined, x_mark_enc):
99         batch_size, target_length, _ = x_combined.shape
100         _, _, num_features = x_mark_enc.shape
101
102         x_mark_enc = x_mark_enc.permute(0, 2, 1)
103         x_mark_enc_interp = F.interpolate(
104             x_mark_enc, size=target_length, mode='linear', align_corners
    =False)
105         x_mark_enc_interp = x_mark_enc_interp.permute(0, 2, 1)
106
107         return x_mark_enc_interp
108
109     def forecast(self, x_enc, x_mark_enc, x_dec, x_mark_dec):
110         # Preprocess with 2D CNN
111         x_processed = self.preprocess_with_cnn2d(x_enc)
```

```

112
113     # Interpolate mark_enc to match new sequence length
114     x_mark_enc_interp = self.mark_enc_interpolation(x_processed,
115     x_mark_enc)
116
117     # Embedding
118     enc_out = self.enc_embedding(x_processed, x_mark_enc_interp)
119     enc_out, attns = self.encoder(enc_out, attn_mask=None)
120
121     dec_out = self.dec_embedding(x_dec, x_mark_dec)
122     dec_out = self.decoder(dec_out, enc_out, x_mask=None, cross_mask
123     =None)
124     return dec_out
125
126     def forward(self, x_enc, x_mark_enc, x_dec, x_mark_dec, mask=None):
127         dec_out = self.forecast(x_enc, x_mark_enc, x_dec, x_mark_dec)
128         return dec_out[:, -self.pred_len:, :] # [B, L, D]

```

The TCformer using 1D CNN for iTransformer on seq2seq task

```

1 import torch
2 import torch.nn as nn
3 import torch.nn.functional as F
4 from layers.Transformer_EncDec import Encoder, EncoderLayer
5 from layers.SelfAttention_Family import FullAttention, AttentionLayer
6 from layers.Embed import DataEmbedding_inverted
7
8
9 class Model(nn.Module):
10     def __init__(self, configs):
11         super(Model, self).__init__()
12         self.seq_len = configs.seq_len
13         self.pred_len = configs.pred_len

```

```

14         self.output_attention = configs.output_attention
15         self.use_norm = configs.use_norm
16         self.d_model = configs.d_model
17         self.enc_in = configs.enc_in
18
19         # CNN parameters
20         self.kernel_size_1 = 5
21         self.stride_1 = 2
22         self.kernel_size_2 = 3
23         self.stride_2 = 1
24
25         self.split_factor=0.8
26         combined_seq_len = self.cnn_seq_calc(self.cnn_seq_calc(int(
27             self.split_factor*configs.seq_len), self.kernel_size_1, self
28             .stride_1), self.kernel_size_2, self.stride_2)+configs.seq_len-int(
29             self.split_factor*configs.seq_len)
30
31         # First 1D CNN layer
32         self.cnn1 = nn.Conv1d(in_channels=self.enc_in,
33                               out_channels=self.d_model,
34                               kernel_size=self.kernel_size_1,
35                               stride=self.stride_1,
36                               padding=0)
37
38         # Second 1D CNN layer
39         self.cnn2 = nn.Conv1d(in_channels=self.d_model,
40                               out_channels=self.d_model,
41                               kernel_size=self.kernel_size_2,
42                               stride=self.stride_2,
43                               padding=0)
44
45         self.dropout1 = nn.Dropout(p=0.2)
46
47         # Linear layer to project raw input to d_model dimension

```

```

45         self.cnn_proj = nn.Linear(self.d_model, self.enc_in)
46
47         # combined_seq_len = (int(
48         #     0.8*configs.seq_len) - self.kernel_size) // self.stride +
1 + configs.seq_len - int(0.8*configs.seq_len)
49         print(f"combined_seq_len: {combined_seq_len}")
50
51         # Embedding
52         self.enc_embedding = DataEmbedding_inverted(combined_seq_len,
configs.d_model, configs.embed, configs.freq,
53                                                     configs.dropout)
54
55         self.class_strategy = configs.class_strategy
56
57         # Encoder-only architecture
58         self.encoder = Encoder(
59             [
60                 EncoderLayer(
61                     AttentionLayer(
62                         FullAttention(False, configs.factor,
attention_dropout=configs.dropout,
63                                     output_attention=configs.
output_attention), configs.d_model, configs.n_heads),
64                     configs.d_model,
65                     configs.d_ff,
66                     dropout=configs.dropout,
67                     activation=configs.activation
68                 ) for l in range(configs.e_layers)
69             ],
70             norm_layer=torch.nn.LayerNorm(configs.d_model)
71         )
72         self.projector = nn.Linear(
73             configs.d_model, configs.pred_len, bias=True)

```

```

74     def cnn_seq_calc(self, seq_len, kernel_size, stride):
75         return (seq_len - kernel_size) // stride + 1
76
77     def forecast(self, x_enc, x_mark_enc, x_dec, x_mark_dec):
78         if self.use_norm:
79             # Normalization
80             means = x_enc.mean(1, keepdim=True).detach()
81             x_enc = x_enc - means
82             stdev = torch.sqrt(
83
84                 torch.var(x_enc, dim=1, keepdim=True, unbiased=False) +
1e-5)
85             x_enc /= stdev
86
87         B, L, N = x_enc.shape
88
89         split_point = int(self.split_factor * L)
90         x_cnn = x_enc[:, :split_point, :]
91         x_raw = x_enc[:, split_point:, :]
92
93         # Process 80% with first CNN layer
94         x_cnn = x_cnn.permute(0, 2, 1) # [B, N, L] for conv1d
95         x_cnn = self.cnn1(x_cnn)
96
97         x_cnn = F.relu(x_cnn)
98
99         # Process with second CNN layer
100        x_cnn = self.cnn2(x_cnn)
101        x_cnn = x_cnn.permute(0, 2, 1) # [B, L, N]
102
103        x_cnn = F.relu(x_cnn)
104        # Project CNN output to d_model dimension
105        x_cnn=self.cnn_proj(x_cnn)

```



```

106         x_cnn = F.relu(x_cnn)
107         # Concatenate CNN output with raw 20%
108         x_combined = torch.cat([x_cnn, x_raw], dim=1)
109
110         # Interpolate x_mark_enc to match x_combined length
111         # x_mark_enc_interp = self.mark_enc_interpolation(x_combined,
x_mark_enc)
112
113         # Embedding
114         # print(f"x_combined: {x_combined.shape} and x_mark_enc_interp:{
x_mark_enc_interp.shape}")
115         enc_out = self.enc_embedding(x_combined, None)
116
117         # Encoder
118         enc_out, attns = self.encoder(enc_out, attn_mask=None)
119
120         # Projection
121         dec_out = self.projector(enc_out).permute(0, 2, 1)[: , : , :N]
122
123         if self.use_norm:
124             # De-Normalization
125             dec_out = dec_out * \
126                 (stdev[:, 0, :].unsqueeze(1).repeat(1, self.pred_len, 1)
)
127             dec_out = dec_out + \
128                 (means[:, 0, :].unsqueeze(1).repeat(1, self.pred_len, 1)
)
129
130         return dec_out
131
132     def mark_enc_interpolation(self, x_combined, x_mark_enc):
133         batch_size, target_length, _ = x_combined.shape
134         _, _, num_features = x_mark_enc.shape

```

```
135
136     x_mark_enc = x_mark_enc.permute(0, 2, 1)
137     x_mark_enc_interp = F.interpolate(
138         x_mark_enc, size=target_length, mode='linear', align_corners
139         =False)
140     x_mark_enc_interp = x_mark_enc_interp.permute(0, 2, 1)
141
142     return x_mark_enc_interp
143
144 def forward(self, x_enc, x_mark_enc, x_dec, x_mark_dec, mask=None):
145     dec_out = self.forecast(x_enc, x_mark_enc, x_dec, x_mark_dec)
146     return dec_out[:, -self.pred_len:, :]
```