

METATST: ESSENTIAL TRANSFORMER COMPONENTS FOR TIME SERIES ANALYSIS

Anonymous authors

Paper under double-blind review

ABSTRACT

1 This paper presents MetaTST, a versatile time series Transformer architecture that
2 combines standard Transformer components with time series-specific features,
3 omitting the traditional token mixer in favor of non-parametric pooling opera-
4 tors. The study’s two primary contributions include defining the MetaTST ar-
5 chitecture and showcasing its empirical success across forecasting, classification,
6 imputation, and anomaly detection tasks. These results establish MetaTST as a
7 robust and adaptable foundation for future time series Transformer designs, rais-
8 ing important questions about the necessity of attention mechanisms in time series
9 analysis.

10 1 INTRODUCTION

11 Time series analysis techniques is widely used in real world applications. In recent years, deep learn-
12 ing for time series analysis has received great interests. Many classical models, such as MLP, CNN
13 and RNN, have found their variations for time series analysis. Transformer (Vaswani et al., 2017),
14 which is designed for NLP tasks, is now becoming popular in many areas such as CV (Dosovitskiy
15 et al., 2021) and time series analysis. Benefits from its self-attention mechanism, Transformers can
16 capture dependencies of long sequence. This lead to the success of Transformers in many areas.

17 In those time series transformers, Autoformer (Wu et al., 2021), FEDformer (Zhou et al., 2022) are
18 among the best variants successfully applied to time series data. One of the main challenges they
19 all trying to solve is the computation/memory bottleneck brought by the quadratic complexity of at-
20 tention mechanism. With the insight that attention on time series often turns out to be sparse (Zhou
21 et al., 2021), they adopt various substitute attention block specially designed for time series which
22 can capture new time series features and have lower complexity. For example, the auto-correlation
23 (Wu et al., 2021) replaces self-attention with series-wise connections that can be calculated effi-
24 ciently via FFT (Fast Fourier Transform) with $O(L \log L)$ complexity. FEDformer use FFT and
25 Wavelet Transform to capture the features in frequency demain. Along this line of research, the
26 success of these models are mainly attributed to their newly devised attention substitution.

27 Although the performance of time series Transformers grows, its effectiveness is questioned by a re-
28 cent work (Zeng et al., 2023). The authors demonstrate that a simple linear projection with seasonal-
29 trend decomposition can outperform most Transformer variants, putting question on the effective-
30 ness of Transformer architecture and attention mechanism for time series analysis, especially in the
31 LTSF (Long-term Time Series Forecasting) task. As a fight-back, PatchTST (Nie et al., 2023) im-
32 proves the capacity of Transformer architecture by introducing patching and channel-independence.
33 Moreover, in CV, Metaformer (Yu et al., 2022a) provides a strong baseline for vision Transformers.
34 It uses a simple pooling operator as the token mixer (which is traditionally implemented by attention
35 mechanism) to aggregates information among tokens and achieves reasonable performance, thereby
36 attributes the model capacity to the Transformer architecture itself.

37 With all these observations, this paper aims to explore what is really useful for time series trans-
38 formers. We abstract the essential parts of time series Transformers as MetaTST (**Meta Time Series**
39 **Transformer**). MetaTST contain time series tailored components such as decomposition, instance
40 norm as well as patching technique. Meanwhile, it does not specify concrete token mixer. By im-
41 plementing the token mixer with simple non-parametric operator pooling, we demonstrate that the
42 MetaTST architecture can bring promising performance through extensive experiments on 4 time
43 series analysis tasks.

44 The contributions of this paper are two-fold. Firstly, this paper summarize the time series transform-
 45 ers into a general architecture MetaTST, and empirically demonstrate that general transformer archi-
 46 tecture plus with time series tailored components can achieve promising performance. Secondly, this
 47 paper evaluates the proposed MetaTST on different time series tasks including forecasting, classifi-
 48 cation, imputation and anomaly detection. MetaTST performs on par with other well-acknowledge
 49 time series Transformers. Thus, MetaTST can serve as a good start base for future time series
 50 Transformer design.

51 2 RELATED WORK

52 Transformer (Vaswani et al., 2017) is first proposed for NLP tasks and then rapidly become popular
 53 in many various tasks such as computer vision (Dosovitskiy et al., 2021) and time series (Li et al.,
 54 2019; Zhou et al., 2021). Along the line of transformers for time series analysis, the main challenge
 55 of time series Transformer is the quadratic complexity of dot-product attention in self-attention
 56 mechanism. In order to tackle this problem, (Zhou et al., 2021) points out that the attention score
 57 is sparsely distributed, thereby it is possible to reduce the complexity of attention mechanism while
 58 maintaining most information. For example, Autoformer (Wu et al., 2021) propose auto-correlation
 59 that can seamlessly replace multi-head attention and be able to capture series-wise dependence of
 60 time series. Fedformer (Zhou et al., 2022) capture frequency domain information with Fourier
 61 Transform.

62 The other line of research provides methods on how to incorporate insights of time series into deep
 63 learning models especially for Transformers. Multi-level seasonal-trend decomposition is proposed
 64 by (Wu et al., 2021) and proved to be a useful design by (Zeng et al., 2023). (Nie et al., 2023)
 65 proposes patching to enable the model to directly capture series-wise dependence and keep channel
 66 independent. (Kim et al., 2022) and (Liu et al., 2022) notice the problem of distribution shift between
 67 training and testing dataset. Similar instance normalization is proposed to solve this problem.

68 However, as questioned by (Zeng et al., 2023), are Transformers effective for time series forecast-
 69 ing? They show that a simple linear model with decomposition can beat many complex Transformer-
 70 based models on long-term time series forecasting task. Metaformer (Yu et al., 2022a;b) points out
 71 that complex token-mixer (attention) in Transformer can be replaced by a light-weight and simple
 72 pooling module while maintaining most of performance. What really matters is the Metaformer
 73 architecture that consists of input-embedding, residual connection, arbitrary token-mixer, channel-
 74 mixer. This paper, however, aims at verifying is similar hypothesis holds in time series forecasting
 75 task: Metaformer plus with add-on time series adopted tricks are all you need for time series fore-
 76 casting.

77 3 METHOD

78 3.1 THE METATST FRAMEWORK

79 Figure.1 shows the overall framework of MetaTST. MetaTST is an abstracted general architecture
 80 based on transformer with time series related modifications. Note that the token mixer, which is often
 81 implemented by various attention mechanisms, is not specified, meaning that any token/time-wise
 82 aggregation modules can be applied. Given the input I steps multivariate time series $\mathbf{X} \in \mathbb{R}^{I \times C}$
 83 of C variables, the input is first processed by instance norm module to mitigate the influence of
 84 distribution shift between training and testing sets. Then the positional encoding is added and the
 85 whole sequence is transformed by patching to make it suitable for Transformers.

86 After that, the input time series is decomposed into seasonal part and trend part, then fed into the
 87 MetaTST encoder stacks. Each stack contains a token mixer to gather time-wise information and a
 88 feed forward layer module to gather channel-wise information. Two series decomposition modules
 89 are also included to gradually decompose the time series so that it can be processed better by next
 90 module. Note that only the seasonal part goes through these modules, the decomposed trend are
 91 aggregated together and added with the seasonal part at the end of the encoder stack. Finally, the
 92 extracted features are fed into the projection head, which could be different between generative tasks
 93 such as forecasting and identify tasks such as classification. If it is for generative tasks, the output
 94 has to be denormalized.

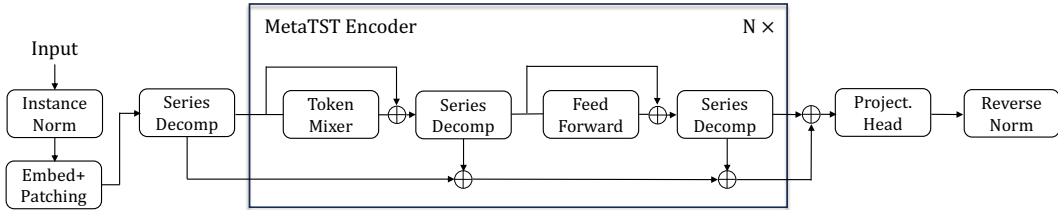


Figure 1: The overall framework of MetaTST.

95 3.2 ESSENTIAL COMPONENTS FOR TIME SERIES TRANSFORMER

96 **Pooling as Token Mixer.** Token mixer is often implemented by various attention mechanism, such
 97 as vanilla attention (Vaswani et al., 2017), autocorrelation (Wu et al., 2021), frequency enhanced
 98 block (Zhou et al., 2022) and so on. This line of work often attributes their model capacity to the
 99 elaborately designed attention mechanism. In this paper, we use a simple parameter-free operator,
 100 i.e. average pooling, to replace the attention. Compared with other attention mechanisms, pooling
 101 is extremely simple and the computation cost is rather low. As a token mixer, the receptive field of a
 102 single pooling operator cannot cover the whole sequence. Thus, the pooling size is set to be rather
 103 large to increase the receptive filed of each pooling layer.

104 **Decomposition.** Time series often consists of components with different dynamics. For example the
 105 house price may grow with years and fluctuate within a year. Thus it is useful to decompose those
 106 patterns and process for them respectively. Seasonal-Trend Decomposition Zeng et al. (2023); Wu
 107 et al. (2021) has been used in several time series forecasting models. And it is of great importance for
 108 their accurate forecasting. Formally, given the input series \mathbf{X} , the decomposition module divide it
 109 into seasonal part \mathbf{X}_s and trend part \mathbf{X}_t . This procedure can be implemented simply via AvgPool1d
 110 in PyTorch. Formally,

$$\begin{aligned} \mathbf{X}_t &= \text{AvgPool1d}(\mathbf{X}) & (1) \\ \mathbf{X}_s &= \mathbf{X} - \mathbf{X}_t & (2) \end{aligned}$$

111 A time series may contain complicated patterns that cannot be decomposed with only one operation.
 112 Thus, it is necessary to do multiple decomposition operation. In MetaTST, global decomposition is
 113 conducted firstly to filter out global trend part, so that the encoders only handle the seasonal part.

114 **Patching.** Patching is first introduced in vision Transformers (Dosovitskiy et al., 2021). It split a
 115 input 2D image into local patches so that they can be treated as a sequence by Transformer. Back into
 116 time series, this technique is also useful since it can significantly reduces nominal sequence length
 117 and eliminate the memory constaints hindering Time Series Transformers to handle long sequences
 118 (Nie et al., 2023). Given the original input series $\mathbf{X} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)}\}$, for each univariate time
 119 series $\mathbf{x}^{(i)}$, it is splitted into 2D patches with patch length P and stride S . Then the patches sequence
 120 is $x_p^i \in \mathbb{R}^{P \times N}$ and $N = \frac{L-P}{S} + 2$ is the number of patches. However, with batches and multivariate
 121 setting, this proecess generates a 4D tensor $\mathbf{X}_p \in \mathbb{R}^{B \times C \times P \times N}$. We merge the first two dimension of
 122 \mathbf{X}_p and then get $\mathbf{X}'_p \in \mathbb{R}^{(B * C) \times P \times N}$ so that it can be processed by Transformer models.

123 **Instance Normalization.** The data distribution between training and test set can be different, lead-
 124 ing to degradation of a well trained model performance on test set. The instance norm Kim et al.
 125 (2022); Liu et al. (2022) can tackle this problem to some extend. By normalize each input time
 126 series instance, and denormalize back the model outputs, it stablizes the value to comply with the
 127 distribution of the test set. Thereby increase the performance on generative tasks such as forecasting,
 128 imputation, and anomaly detection (the observation outliers compared with prediction are regarded
 129 as anomaly). MetaTST adopts a RevIN layer which makes extra learnable affine transform of the
 130 normalized data. Formally, for k -th instance, each point $x_{kt}^{(i)}$ in input series at step t is normalized
 131 as:

$$\hat{x}_{kt}^{(i)} = \gamma_k \left(\frac{x_{kt}^{(i)} - \mathbb{E}_t [x_{kt}^{(i)}]}{\sqrt{\text{Var} [x_{kt}^{(i)}] + \epsilon}} \right) + \beta_k \quad (3)$$

132 and final prediction is denormalized as:

$$\hat{y}_{kt}^{(i)} = \sqrt{\text{Var} [x_{kt}^{(i)}] + \epsilon} \cdot \left(\frac{\tilde{y}_{kt}^{(i)} - \beta_k}{\gamma_k} \right) + \mathbb{E}_t [x_{kt}^{(i)}] \quad (4)$$

133 where γ_k and β_k can be fixed or learnable parameters.

134 4 EXPERIMENTS

135 **Baselines.** Since this paper aims to summarize the effective components in time series analysis,
 136 we compare the performance of MetaTST with several well-acknowledged Transformer-based time
 137 series models, including Autoformer (Wu et al., 2021), FEDformer (Zhou et al., 2022), Pyraformer
 138 (Liu et al., 2021). Beside, to verify the effectiveness of MetaTST architecture, vanilla Transformer
 139 (Vaswani et al., 2017) is taken as baseline as well.

140 **General Setup.** The model is trained with the ADAM (Kingma & Ba, 2014) optimizer with an initial
 141 learning rate of 10^{-3} . Batch size is set to 32, shrunk if the model runs out of GPU memory under
 142 large batch size. The training process is early stopped within 10 epochs for generative tasks including
 143 forecasting, imputation and anomaly detection, implemented in PyTorch Paszke et al. (2019) with
 144 codebase from (Wu et al., 2022) and conducted on NVIDIA RTX 3090 24GB GPUs. Generally, the
 145 time series Transformers have 2 encoder layers and 1 decoder layer. Since the MetaTST does not
 146 contain a decoder, for fair comparison, the number of encoder layer in MetaTST is set to 3.

147 4.1 FORECASTING

148 **Setup.** In order to verify the hypothesis, we conduct empirical experiments of long term forecasting
 149 task on ETTm1, Traffic, Weather and ECL datasets Zhou et al. (2021); Wu et al. (2021), as well as
 150 short term forecasting task on M4 dataset (Makridakis et al., 2018). Loss function is Mean Squared
 151 Error (MSE).

152 **Results.** Table. 1 and Table. 2. shows the long-term forecasting results and short-term forecasting re-
 153 sults respectively. Surprisingly, MetaTST achieve most of the best performance on these benchmarks.
 154 For the M4 dataset, MetaTST outperforms all other models, showing that the proposed framework
 155 suits the forecasting tasks very well.

156 Pooling operator aggregates nearly tokens evenly. Thus it is an extremely simple token mixer. How-
 157 ever, the experiment results show that with that kind of simple token mixing operator, MetaTST still
 158 obtain competitive performance compared with other Transformer-based model. Fig. 2 gives show
 159 cases of forecasting results on ECL and ETTm1 dataset. Although they are difference quantitatively
 160 on MSE metric, the actual prediction shows no significant difference. This findings conveys that
 161 the MetaTST is the base-stone for Transformer models to achieve reasonable performance on time
 162 series forecasting task.

163 4.2 IMPUTATION

164 **Setup.** Missing values often appear in real world time series data due to the malfunction of
 165 data collector. To facilitate downstream tasks, it is necessary to recover the original data with the
 166 partially missing data. To verify the performance of MetaTST on imputation task, three typical
 167 datasets ETTm1, ECL and Weather are selected. In order to compare the model capacity under
 168 different proportions of missing data, the ratio we randomly masked in the experiment varies in
 169 12.5%, 25%, 37.5%, 50%.

170 **Results.** As shown in Table. 3, the MetaTST performs on par with other Transformer-based models.
 171 Revealing that the MetaTST architecture is suitable for imputation task.

172 4.3 ANOMALY DETECTION

173 **Setup.** Detecting anomalies from monitoring data is an important application for various areas.
 174 Since anomalies are often hidden in large amounts of data, it is hard to find those anomalies by peo-
 175 ple. Here we focus on unsupervised time series anomaly detection. The experiments are conducted

Table 1: Results of the long-term forecasting task

Dataset	Length	Autoformer		FEDformer		Pyraformer		MetaTST	
		MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
ettm1	96	0.438	0.446	0.419	0.452	0.604	0.513	0.329	0.367
	192	0.484	0.470	0.447	0.456	0.651	0.559	0.374	0.390
	336	0.464	0.475	0.443	0.456	0.779	0.653	0.402	0.409
	720	0.464	0.479	0.539	0.508	0.896	0.701	0.463	0.443
traffic	96	0.602	0.384	0.590	0.365	0.867	0.468	0.512	0.336
	192	0.605	0.371	0.600	0.369	0.869	0.467	0.509	0.332
	336	0.684	0.432	0.643	0.406	0.881	0.469	0.523	0.336
	720	0.650	0.395	0.653	0.400	0.896	0.473	0.559	0.353
weather	96	0.270	0.346	0.218	0.304	0.194	0.276	0.186	0.223
	192	0.305	0.369	0.275	0.347	0.227	0.312	0.230	0.260
	336	0.352	0.395	0.406	0.439	0.304	0.366	0.283	0.298
	720	0.456	0.458	0.453	0.462	0.395	0.418	0.344	0.344
electricity	96	0.234	0.342	0.193	0.310	0.386	0.449	0.170	0.259
	192	0.215	0.324	0.212	0.326	0.378	0.443	0.178	0.266
	336	0.291	0.389	0.233	0.350	0.376	0.443	0.193	0.282
	720	0.296	0.391	0.268	0.377	0.376	0.445	0.233	0.315

Table 2: Results of the short-term forecasting task in the M4 dataset.

Period	Metric	Autoformer	FEDformer	Pyraformer	Transformer	PatchTST	MetaTST
Year	SMAPE	69.522	17.974	13.604	14.694	13.564	13.396
	MASE	18.142	4.062	3.075	3.304	3.050	3.005
	OWA	4.409	1.061	0.803	0.865	0.799	0.788
Quarterly	SMAPE	73.760	14.485	10.610	11.506	10.791	10.805
	MASE	13.282	1.872	1.246	1.375	1.299	1.305
	OWA	8.192	1.340	0.936	1.024	0.964	0.966
Monthly	SMAPE	69.837	18.235	13.887	15.589	14.540	13.262
	MASE	11.164	1.592	1.053	1.209	1.139	1.005
	OWA	7.670	1.381	0.976	1.109	1.039	0.932
Others	SMAPE	106.379	6.721	4.804	5.829	6.350	4.778
	MASE	82.033	4.793	3.238	4.034	4.020	3.268
	OWA	24.129	1.463	1.016	1.249	1.302	1.018
Average	SMAPE	72.533	16.699	12.581	13.915	13.006	12.279
	MASE	16.821	2.388	1.674	1.872	1.761	1.650
	OWA	7.072	1.240	0.901	1.002	0.940	0.884

176 on five anomaly detection benchmarks including: SMD, MSL, SMap, SWaT and PSM, covering
 177 different applications. Following previous work on this task Xu et al. (2021); Wu et al. (2022),
 178 the dataset is splited into consecutive non-overlapping segments by sliding window. And only the
 179 classical reconstruction error is regarded as the shared anomaly criterion for all experiments.

180 **Results.** As shown in Table 4, MetaTST achieves a reasonable performance in anomaly detection
 181 task with the mose simple token-mixer. The performance can be attributed to the MetaTST archi-
 182 tecture.

183 5 CONCLUSION AND FUTURE WORK

184 This paper summarizes recent research on time series Transformers by proposing a abstract model
 185 architecture called MetaTST. It contains essential components for time series Transformers includ-
 186 ing the overall architecture, instance normalization, decomposition and patching. Compared with
 187 other time series Transformers, MetaTST uses a simple pooling operation but can still achieve com-
 188 petitive results, showing that the capacity of time series Transformers attributes a lot to the whole
 189 time-series-adopted architecture. Thus, the hypothesis proposed by Metaformer perhaps holds in
 190 time series analysis area. Our work reveals where the capacity of time series Transformers come
 191 from. Thus, MetaTST has the potential to be the base model for future model design and serve as

Table 3: Imputation results on Weather, ETTm1 and ECL datasets.

Dataset	Mask Ratio	Transformer		Autoformer		FEDformer		Pyraformer		MetaTST	
		MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
Weather	0.125	0.033	0.087	0.357	0.438	0.044	0.107	0.030	0.074	0.031	0.057
	0.250	0.035	0.086	0.144	0.252	0.055	0.128	0.036	0.089	0.033	0.057
	0.375	0.039	0.097	0.135	0.239	0.076	0.159	0.039	0.091	0.034	0.058
	0.500	0.042	0.094	0.180	0.281	0.116	0.211	0.041	0.092	0.038	0.063
ETTm1	0.125	0.023	0.107	0.718	0.699	0.034	0.130	0.032	0.128	0.046	0.143
	0.250	0.028	0.117	0.526	0.573	0.053	0.163	0.035	0.132	0.055	0.150
	0.375	0.035	0.130	0.350	0.443	0.083	0.202	0.041	0.140	0.060	0.159
	0.500	0.044	0.145	0.313	0.402	0.133	0.260	0.048	0.152	0.067	0.167
ECL	0.125	0.150	0.278	0.191	0.328	0.185	0.323	0.190	0.303	0.059	0.163
	0.250	0.157	0.282	0.198	0.309	0.207	0.340	0.216	0.346	0.072	0.183
	0.375	0.168	0.290	0.216	0.346	0.225	0.355	0.195	0.305	0.088	0.203
	0.500	0.180	0.297	0.234	0.360	0.251	0.372	0.207	0.312	0.108	0.227

Table 4: Results of anomaly detection task.

	Transformer			Autoformer			FEDformer			Pyraformer			PatchTST			MetaTST		
	P	R	F1	P	R	F1	P	R	F1	P	R	F1	P	R	F1	P	R	F1
MSL	89.98	73.79	81.09	90.53	74.96	82.01	90.71	75.41	82.35	89.01	70.84	78.90	88.31	70.77	78.57	88.51	71.64	79.19
PSM	99.36	83.20	90.56	99.99	78.96	88.24	99.98	81.94	90.07	98.53	88.36	93.17	98.84	93.54	96.12	98.73	90.91	94.66
SMAP	90.96	62.28	73.94	91.47	67.66	77.79	89.96	55.47	68.62	89.56	54.54	67.80	90.63	55.51	68.85	90.17	53.75	67.35
SMD	78.48	65.27	71.26	78.41	65.06	71.12	78.44	64.98	71.08	79.16	93.54	73.23	87.26	82.12	84.61	87.15	77.53	82.06
SWAT	99.70	66.08	79.48	99.96	65.55	79.18	99.96	65.55	79.18	99.94	65.56	79.18	91.34	83.31	87.14	91.45	84.23	87.69
Avg F1		79.27			79.67			78.26			78.46			83.06				82.19

192 a baseline for new Transformer-based models. Each part of MetaTST is proven to be effective by
 193 extensive experiments.

194 REFERENCES

- 195 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
 196 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, and Sylvain Gelly. An image
 197 is worth 16x16 words: Transformers for image recognition at scale. In *International Conference*
 198 *on Learning Representations*, 2021.
- 199 Taesung Kim, Jinhee Kim, Yunwon Tae, Cheonbok Park, Jang-Ho Choi, and Jaegul Choo. Re-
 200 versible instance normalization for accurate time-series forecasting against distribution shift. In
 201 *International Conference on Learning Representations*, 2022.
- 202 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
 203 *arXiv:1412.6980*, 2014.
- 204 Shiyang Li, Xiaoyong Jin, Yao Xuan, Xiyong Zhou, Wenhui Chen, Yu-Xiang Wang, and Xifeng
 205 Yan. Enhancing the locality and breaking the memory bottleneck of transformer on time series
 206 forecasting. *Advances in neural information processing systems*, 32, 2019.
- 207 Shizhan Liu, Hang Yu, Cong Liao, Jianguo Li, Weiyao Lin, Alex X Liu, and Schahram Dustdar.
 208 Pyraformer: Low-complexity pyramidal attention for long-range time series modeling and fore-
 209 casting. In *International conference on learning representations*, 2021.
- 210 Yong Liu, Haixu Wu, Jianmin Wang, and Mingsheng Long. Non-stationary transformers: Exploring
 211 the stationarity in time series forecasting. *Advances in Neural Information Processing Systems*,
 212 35:9881–9893, 2022.
- 213 Spyros Makridakis, Evangelos Spiliotis, and Vassilios Assimakopoulos. The m4 competition: Re-
 214 sults, findings, conclusion and way forward. *International Journal of Forecasting*, 34(4):802–808,
 215 2018.
- 216 Yuqi Nie, Nam H Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth 64
 217 words: Long-term forecasting with transformers. In *The Eleventh International Conference on*
 218 *Learning Representations*, 2023. <https://github.com/yuqinie98/PatchTST/tree/main>.

- 219 A Paszke, S Gross, F Massa, A Lerer, Jea PyTorch Bradbury, G Chanan, T Killeen, Z Lin,
220 N Gimselshin, L Antiga, et al. An imperative style, high-performance deep learning library.
221 *Adv. Neural Inf. Process. Syst.*, 32:8026, 2019.
- 222 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
223 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural informa-*
224 *tion processing systems*, 30, 2017.
- 225 Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Autoformer: Decomposition trans-
226 formers with auto-correlation for long-term series forecasting. *Advances in Neural Information*
227 *Processing Systems*, 34:22419–22430, 2021.
- 228 Haixu Wu, Tengge Hu, Yong Liu, Hang Zhou, Jianmin Wang, and Mingsheng Long. Timesnet:
229 Temporal 2d-variation modeling for general time series analysis. In *The Eleventh International*
230 *Conference on Learning Representations*, 2022.
- 231 Jiehui Xu, Haixu Wu, Jianmin Wang, and Mingsheng Long. Anomaly transformer: Time series
232 anomaly detection with association discrepancy. In *International Conference on Learning Repre-*
233 *sentations*, 2021.
- 234 Weihao Yu, Mi Luo, Pan Zhou, Chenyang Si, Yichen Zhou, Xinchao Wang, Jiashi Feng, and
235 Shuicheng Yan. Metaformer is actually what you need for vision. In *Proceedings of the IEEE/CVF*
236 *conference on computer vision and pattern recognition*, pp. 10819–10829, 2022a.
- 237 Weihao Yu, Chenyang Si, Pan Zhou, Mi Luo, Yichen Zhou, Jiashi Feng, Shuicheng Yan, and Xin-
238 chao Wang. Metaformer baselines for vision. *arXiv preprint arXiv:2210.13452*, 2022b.
- 239 Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are transformers effective for time series
240 forecasting? In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, pp.
241 11121–11128, 2023. ISBN 2374-3468.
- 242 Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang.
243 Informer: Beyond efficient transformer for long sequence time-series forecasting. In *Proceedings*
244 *of the AAAI conference on artificial intelligence*, volume 35, pp. 11106–11115, 2021. ISBN
245 2374-3468.
- 246 Tian Zhou, Ziqing Ma, Qingsong Wen, Xue Wang, Liang Sun, and Rong Jin. Fedformer: Frequency
247 enhanced decomposed transformer for long-term series forecasting. In *International Conference*
248 *on Machine Learning*, pp. 27268–27286. PMLR, 2022. ISBN 2640-3498.