

PatchMixer: A Patch-Mixing Architecture for Long-Term Time Series Forecasting

Zeying Gong, Yujin Tang and Junwei Liang^{*}

Hong Kong University of Science and Technology (Guangzhou)

zgong313@connect.hkust-gz.edu.cn, {yujintang,junweiliang}@hkust-gz.edu.cn

Abstract

Recently, transformers incorporating patch-based representations have set new benchmarks in long-term time series forecasting. This naturally raises an important question: Is the impressive performance of patch-based transformers primarily due to the use of patches rather than the transformer architecture itself? To explore this, we introduce **PatchMixer**, a patch-based CNN that enhances accuracy and efficiency through depthwise separable convolution. Our experimental results on seven time-series forecasting benchmarks indicate that PatchMixer achieves relative improvements of 3.9%, 11.6%, and 21.2% in comparison to state-of-the-art Transformer, MLP, and CNN models, respectively. Additionally, it demonstrates **2-3** training and inference times faster than the most advanced method. We also found that optimizing the patch embedding parameters and enhancing the objective function enables PatchMixer to better adapt to different datasets, thereby improving the generalization of the patch-based approach. Code is available at: <https://github.com/Zeying-Gong/PatchMixer>.

1 Introduction

Long-term time series forecasting (LTSF) remains a crucial challenge within machine learning, entailing the prediction of future data points based on historically observed sequences. LTSF applications across various domains include traffic flow estimation, energy management, and weather forecasting.

In the LTSF domain, prevalent methods are typically Transformer and MLP-based models, which have a global receptive field due to features like self-attention in Transformers and global operations in MLPs. Informer [Zhou *et al.*, 2021] and various Transformer variants [Wu *et al.*, 2021; Zhou *et al.*, 2022; Liu *et al.*, 2022b] have been developed for time series analysis. However, the simple MLP networks in [Zeng *et al.*, 2023] surprisingly outperformed previous models, prompting a shift to MLPs and a reevaluation of Transformers' ability to capture long-term temporal relationships.

^{*}Corresponding author

Recently, PatchTST [Nie *et al.*, 2023], a patch-based Transformer, demonstrated the continued dominance of Transformers in LTSF with superior predictive performance. This raises a key question: Is the impressive performance of patch-based Transformers due to the Transformer architecture itself or the patch-based input representation?

To answer this question, we might as well turn our focus to Convolutional Neural Networks (CNNs), which only have the local receptive field. Therefore, it is common to consider this architecture unsuitable for LTSF tasks. However, this contrast raises an intriguing possibility: achieving comparable or superior forecasting performance with patch presentation and convolutional core module could highlight the importance of data preprocessing techniques in the time series.

In this paper, we introduce a novel patch-based model known as PatchMixer. It is based on a depthwise separable convolutional module and distinguished by its “patch-mixing” design, which (i) processes time series data in patches, retaining the sequential structure, and (ii) captures dependencies within and across variables through shared weights, emphasizing the role of patch-based preprocessing in striking a balance between efficiency and performance. The main contributions are as follows:

- PatchMixer surpasses the state-of-the-art (SOTA) Transformer, MLP, and CNN models in Mean Squared Error (MSE), achieving improvements of 3.9%, 11.6%, and 21.2% respectively, across seven prominent long-term forecasting benchmarks.
- PatchMixer balances efficiency and performance, demonstrating a **3x** faster inference and **2x** faster training speed compared to the current SOTA models under the same configurations.
- In our exploration of patch-based parameter optimization, we enhance the patch-mixing method's predictive performance and generalization capabilities by improving the patch embedding and the loss function.

2 Related Work

CNNs for long-term context. Most temporal CNNs adopt complex architectures to extend their limited receptive fields. For example, the TCN model [Bai *et al.*, 2018] utilized dilated causal convolutions to expand the field. MICN [Wang

et al., 2023] introduced a novel multi-scale hybrid decomposition and isometric convolution. TimesNet [Wu *et al.*, 2023] applied a parallel convolutional structure in the Inception [Szegedy *et al.*, 2015] style. Outside the realm of time series analysis, there are also a few CNN-based studies for long-term context [Gu *et al.*, 2021; Romero *et al.*, 2021; Poli *et al.*, 2023].

Depthwise Separable Convolution. This technique can be abbreviated as **DWConv**. It was first introduced in 2014 [Sifre and Mallat, 2014], and widely adopted in computer vision. It gained prominence in the Inception V1 and V2 models [Szegedy *et al.*, 2015; Ioffe and Szegedy, 2015], supported Google’s MobileNet [Howard *et al.*, 2017] for mobile efficiency, and was scaled up in the Xception network [Chollet, 2017]. ConvMixer [Trockman and Kolter, 2022] highlighted patch representation in computer vision task via it.

3 The Patch-Mixing Design

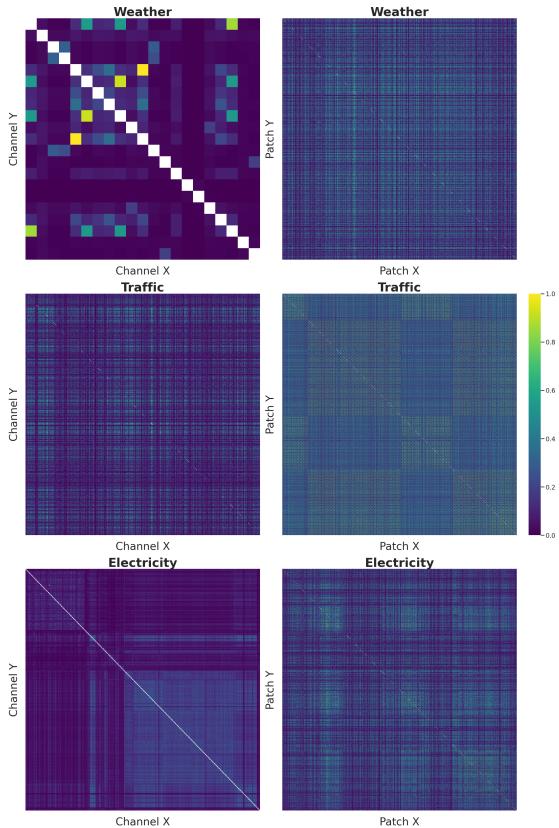


Figure 1: Channel Dependency vs. Patch Dependency analysis was conducted on the 3 largest datasets: Traffic, Electricity, and Weather. The left panel shows the normalized mutual information among different variables, revealing sparse correlations. The right panel illustrates a pronounced intra-variable dependency within single-variable temporal patches, indicated by the more intense coloration.

In the LTSF field, we address the following task: Given a set of multivariate time series instances with a historical look-back window L : $X^{M \times L} = (x_1^M, \dots, x_L^M)$, where each

x_t^M at time step t represents a time series vector x of M variables. Our objective is to make predictions for the subsequent T time steps, resulting in the prediction sequence $\hat{X}^{M \times T} = (\hat{x}_{L+1}^M, \dots, \hat{x}_{L+T}^M)$.

From a signal processing perspective, a multivariate time series is a signal with multiple channels. “Channel mixing” integrates information across variables, while “channel independence” treats each variable separately. This strategy was first used in [Zeng *et al.*, 2023] and further validated in [Nie *et al.*, 2023].

For the patch-based methods, a sliding window of size P and stride S unfolds across the extended series, generating N patches of length S (see Figure 2, left). This process starts by padding each univariate series, replicating the final value S times to preserve edge information, thus $N = \lfloor \frac{(L-P)}{S} \rfloor + 2$.

Therefore, our PatchMixer abstracts the problem as follows. The multivariate time series $X^{M \times L}$ is split into M univariate series $x^{(i)} \in \mathbb{R}^{1 \times L}$, where $i = 1, \dots, M$. In patch embedding, they are transferred to temporal patches $x_p^{(i)} = (x_{1:P}^{(i)}, x_{1+S:P+S}^{(i)}, \dots, x_{L-P+1:L}^{(i)})$, P stands for patch length and S for patch step. The model outputs a prediction for each series, denoted by $\hat{x}^{(i)} \in \mathbb{R}^{1 \times T}$ over the forecast horizon T . By aggregating these individual forecasts, we compose the final multivariate prediction $\hat{X}^{M \times T} = (\hat{x}_{L+1}^M, \dots, \hat{x}_{L+T}^M)$.

We propose “patch mixing”, which uses patching and channel independence to extract cross-patch temporal features. To validate our hypothesis, we measure the interdependence of channels and patches using **Normalized Mutual Information (NMI)**, which quantifies the shared information between two variables [Shannon, 1948], defined as:

$$NMI(X; Y) = \frac{2I(X; Y)}{H(X) + H(Y)} \quad (1)$$

where X and Y are two random variables, $I(X; Y)$ is the mutual information, while $H(X)$ and $H(Y)$ are the respective entropies of X and Y .

As Figure 1 shows, intra-variable temporal patterns hold significant mutual information in the main datasets. The improvement of patch mixing on prediction performance is evident from the experimental results in Table 3.

4 The PatchMixer Model

4.1 Model Structure

The overall architecture of PatchMixer is depicted in Figure 2. We employ a patch-based DWConv module called PatchMixer Block. This design holds efficient parameter usage while maintaining a wide receptive field via patch embedding. Besides, we also devise dual forecasting heads, ensuring a holistic feature representation.

4.2 PatchMixer Block

PatchMixer Block utilizes DWConv as its core. It is composed of the following:

Depthwise Convolution: We utilize depthwise convolution with the kernel size $K = 8$. Assuming that $x^{N \times D}$ represents one of the univariate series, σ denotes an activation function GELU [Hendrycks and Gimpel, 2016], BN represents

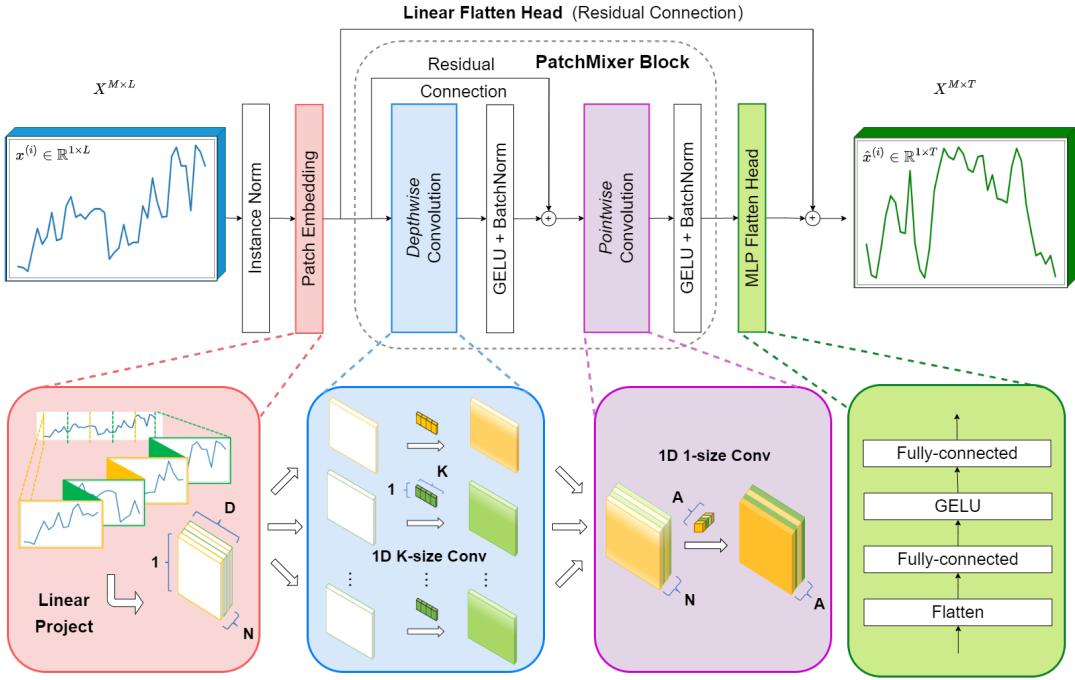


Figure 2: PatchMixer overview.

the BatchNorm operation and l represents the current layer, this process is as follows:

$$x_l^{N \times D} = \text{BN} \left(\sigma(\text{Conv}_{N \rightarrow N}(x_{l-1}^{N \times D}), \text{kernel} = \text{step} = K) \right) \quad (2)$$

Pointwise Convolution: It is applied to capture inter-patch feature correlations, enhanced by a residual connection.

$$x_{l+1}^{A \times D} = \text{BN} \left(\sigma(\text{Conv}_{N \rightarrow N}(x_l^{N \times D}), \text{kernel} = \text{step} = 1) \right) \quad (3)$$

The above equations 3 demonstrate the process of the univariate series $x^{N \times D}$ in layer l passing through the pointwise convolution kernel in layer $(l+1)$, while $\text{Conv}_{N \rightarrow N}$ represents input and output channels are both N .

4.3 Dual Forecasting Heads

PatchMixer introduces a dual-head mechanism. The residual connection spotlights linear trends, and the MLP head addresses nonlinear dynamics. The combination effectively captures a broad spectrum of temporal patterns in fine-to-coarse temporal modeling [Zhou *et al.*, 2023].

4.4 Instance Normalization

This technique [Ulyanov *et al.*, 2016; Kim *et al.*, 2021] is applied before patching, with the original mean and standard deviation reintegrated post-forecasting to preserve the initial scale and distribution of the input time series.

5 Experiments

5.1 Multivariate Long-term Forecasting

Datasets. We use 7 popular multivariate datasets provided in [Wu *et al.*, 2021] for forecasting. Detailed statistical data on

the size of the datasets are as follows.

	Variables	Timesteps	Frequencies
Weather	21	52696	10 Minutes
Traffic	862	17544	1 Hour
Electricity	321	26304	1 Hour
ETTh1	7	17420	1 Hour
ETTh2	7	17420	1 Hour
ETTm1	7	69680	15 Minutes
ETTm2	7	69680	15 Minutes

Table 1: Statistics of popular datasets used for benchmarking.

- Weather: This dataset collects 21 meteorological indicators in Germany, such as temperature and humidity.
- Traffic: This dataset records the road occupancy rates from different sensors on San Francisco freeways.
- Electricity: This dataset describes 321 customers' hourly electricity consumption.
- ETT (Electricity Transformer Temperature): It contains data from two electric transformers (1 and 2), each with two resolutions (15 minutes and 1 hour), resulting in 4 subsets: *ETTm1*, *ETTm2*, *ETTh1*, and *ETTh2*.

Implementation details. For a fair comparison, we adopt most of PatchTST's configurations. However, several differences are noteworthy: (i) The dimension of the Feed Forward Layer in PatchTST is $F = 256$ in most datasets. Our embedding dimension is set directly to $D = 256$ across all datasets. (ii) PatchTST has three layers of vanilla Transformer encoders, while our model achieves better performance with

just one PatchMixer Block. (iii) In our model’s dual heads, the MLP head has two linear layers with GELU activation: one projects the hidden representation to $D = 2 \times T$ for the forecasting length T , and the other projects it to the final prediction target $D = T$. The linear head projects the embed vector directly from $N \times D$ to T . (iv) Our model uses Mean Squared Error (MSE) and Mean Absolute Error (MAE) in a 1:1 ratio as the loss function. Details are in Section 5.2.

Baselines and metrics. We choose the most representative LTSF models as our baselines, including Transformer-based models like PatchTST (2023), FEDformer (2022), Autoformer (2021), Informer (2021), in addition to two CNN-based models containing MICN (2023) and TimesNet (2023), with the significant MLP-based model DLinear (2023). We employ commonly used evaluation metrics: MSE and MAE.

Results. Table 2 presents the multivariate long-term forecasting results. Our model outperforms all baselines on the three largest benchmarks: Traffic, Electricity, and Weather. On other datasets, it achieves the best performance for most prediction lengths. Quantitatively, PatchMixer shows improvements over the SOTA transformer PatchTST with reductions of 3.9% in MSE and 3.0% in MAE. Compared to the leading MLP model DLinear, it reduces MSE by 11.6% and MAE by 9.4% and achieves reductions of 21.2% in MSE and 12.5% in MAE over the most advanced CNN model TimesNet.

5.2 Ablation Study

Ablation of PatchMixer Block. We conducted an ablation study by deleting or substituting the corresponding part, including patching, DWConv, and dual heads, with identical configurations for each variant.

As indicated in Table 3, our findings highlight several key observations: (i) The patch technique enhances performance across various dataset sizes. (ii) The convolutional module outperforms the attention mechanism across all datasets. (iii) The linear head is more effective on small datasets, while the MLP head excels on larger ones; the dual-head combination provides balanced performance. (iv) Position encoding tends to reduce accuracy in patch-based models for LTSF tasks, with better performance observed without it in two out of three datasets.

Varying Look-back Windows. Figure 3 shows the impact of historical length on forecasting accuracy. Key trends are: (i) Traditional CNN-based methods show fluctuating performance with increasing input length, while PatchMixer’s loss decreases steadily, utilizing longer look-back windows effectively. (ii) Across varying input lengths, PatchMixer consistently outperforms PatchTST, showing the patch-based approach, rather than the Transformer architecture, plays a key role in enhancing performance.

5.3 Efficiency Analysis

Comparison of Model Complexity We analyze the computational complexity of PatchMixer, contrasting it with the previous SOTA method PatchTST using the same configuration, with the number of patches N and the embedding dimension D . The convolutional approach has a kernel size of K , with $D \gg N > K$. Table 4 shows the complexity results and

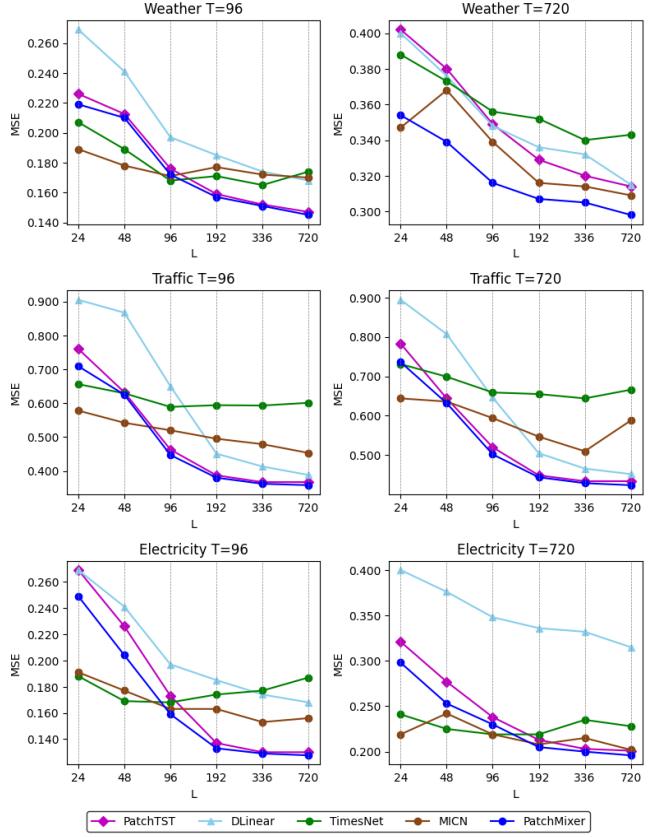


Figure 3: MSE scores with varying look-back windows on top 3 largest datasets. We report the top 5 methods for better observation. The look-back windows $L = [24, 48, 96, 192, 336, 720]$, and the prediction horizons $T = [96, 720]$.

multiply-accumulate operations (MACs) for this configuration. For comparison, we also include the computational cost of a standard convolution to emphasize the significant computational savings mainly due to the DWConv module.

Training and Inference Efficiency. As shown in Figure 4, we conducted experiments on the ETTm1 dataset using a batch size of 8 across 7 variables with identical configurations. We evaluated look-back lengths from 96 to 2880.

Our results show two key improvements: PatchMixer is 3 times faster in inference and twice as fast in training as PatchTST. Additionally, PatchTST’s performance is sensitive to look-back length, while PatchMixer shows fewer fluctuations with increasing input length, demonstrating its stability.

5.4 Patch Embedding and Loss Optimization

We further explore the parametric properties of the patch-mixing architecture. The predictive performance of patch-based methods can be significantly improved by optimizing the patch embedding and enhancing the objective function.

Varying Patch Length. Figure 5 illustrates how patch lengths affect performance. For patch-based methods, losses generally decrease or show minor fluctuations with larger patches. Notably, PatchMixer consistently achieves greater

Models		PatchMixer (Ours)		PatchTST 2023		DLinear 2023		MICN 2023		TimesNet 2023		FEDformer 2022		Autoformer 2021		Informer 2021	
Metrics		MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
Weather	96	0.151	0.193	<u>0.152</u>	<u>0.199</u>	0.176	0.237	0.172	0.240	0.165	0.222	0.238	0.314	0.249	0.329	0.354	0.405
	192	0.194	0.236	<u>0.197</u>	<u>0.243</u>	0.220	0.282	0.218	0.281	0.215	0.264	0.275	0.329	0.325	0.370	0.419	0.434
	336	0.225	0.267	<u>0.249</u>	<u>0.283</u>	0.265	0.319	0.275	0.329	0.274	0.304	0.339	0.377	0.351	0.391	0.583	0.543
	720	0.305	0.323	<u>0.320</u>	<u>0.335</u>	0.323	0.362	0.314	0.354	0.339	0.349	0.389	0.409	0.415	0.426	0.916	0.705
Traffic	96	0.363	0.245	<u>0.367</u>	<u>0.251</u>	0.410	0.282	0.479	0.295	0.593	0.321	0.576	0.359	0.597	0.371	0.733	0.410
	192	0.384	0.254	<u>0.385</u>	<u>0.259</u>	0.423	0.287	0.482	0.297	0.617	0.336	0.610	0.380	0.607	0.382	0.777	0.435
	336	0.393	0.258	<u>0.398</u>	<u>0.265</u>	0.436	0.296	0.492	0.297	0.629	0.336	0.608	0.375	0.623	0.387	0.776	0.434
	720	0.429	0.283	<u>0.434</u>	<u>0.287</u>	0.466	0.315	0.510	0.309	0.640	0.350	0.621	0.375	0.639	0.395	0.827	0.466
Electricity	96	0.129	0.221	<u>0.130</u>	<u>0.222</u>	0.140	0.237	0.153	0.264	0.168	0.272	0.186	0.302	0.196	0.313	0.304	0.393
	192	0.144	0.237	<u>0.148</u>	<u>0.240</u>	0.153	0.249	0.175	0.286	0.184	0.289	0.197	0.311	0.211	0.324	0.327	0.417
	336	0.164	0.257	<u>0.167</u>	<u>0.261</u>	0.169	0.267	0.192	0.303	0.198	0.300	0.213	0.328	0.214	0.327	0.333	0.422
	720	0.200	0.289	<u>0.202</u>	<u>0.291</u>	0.203	0.301	0.215	0.323	0.220	0.320	0.233	0.344	0.236	0.342	0.351	0.427
ETTh1	96	0.353	0.381	<u>0.375</u>	<u>0.399</u>	0.375	0.399	0.405	0.430	0.384	0.402	0.376	0.415	0.435	0.446	0.941	0.769
	192	0.373	0.394	0.414	0.421	<u>0.405</u>	<u>0.416</u>	0.447	0.468	0.436	0.429	0.423	0.446	0.456	0.457	1.007	0.786
	336	0.392	0.414	<u>0.431</u>	<u>0.436</u>	0.439	0.443	0.579	0.549	0.491	0.469	0.444	0.462	0.486	0.487	1.038	0.784
	720	0.445	0.463	<u>0.449</u>	<u>0.466</u>	0.472	0.490	0.699	0.635	0.521	0.500	0.469	0.492	0.515	0.517	1.144	0.857
ETTh2	96	0.225	0.300	<u>0.274</u>	<u>0.336</u>	0.289	0.353	0.349	0.401	0.340	0.374	0.332	0.374	0.332	0.368	1.549	0.952
	192	0.274	0.334	<u>0.339</u>	<u>0.379</u>	0.383	0.418	0.442	0.448	0.402	0.414	0.407	0.446	0.426	0.434	3.792	1.542
	336	0.317	0.368	<u>0.331</u>	<u>0.380</u>	0.480	0.465	0.652	0.569	0.452	0.452	0.400	0.471	0.477	0.479	4.215	1.642
	720	0.393	0.426	<u>0.379</u>	<u>0.422</u>	0.605	0.551	0.800	0.652	0.462	0.468	0.412	0.469	0.453	0.490	3.656	1.619
ETTm1	96	0.291	0.340	<u>0.290</u>	<u>0.342</u>	0.299	0.343	0.302	0.352	0.340	0.377	0.326	0.390	0.505	0.475	0.626	0.560
	192	0.325	0.362	<u>0.332</u>	<u>0.369</u>	0.335	<u>0.365</u>	0.342	0.380	0.374	0.387	0.365	0.415	0.553	0.496	0.725	0.619
	336	0.353	0.382	<u>0.366</u>	<u>0.453</u>	0.369	<u>0.386</u>	0.381	0.403	0.392	0.413	0.392	0.425	0.621	0.537	1.005	0.741
	720	0.413	0.413	<u>0.420</u>	<u>0.533</u>	0.425	<u>0.421</u>	0.434	0.447	0.433	0.436	0.446	0.458	0.671	0.561	1.133	0.845
ETTm2	96	0.174	<u>0.256</u>	0.165	0.255	<u>0.167</u>	0.260	0.188	0.286	0.183	0.271	0.180	0.271	0.255	0.339	0.355	0.462
	192	0.227	<u>0.295</u>	0.220	0.292	<u>0.224</u>	0.303	0.236	0.320	0.242	0.309	0.252	0.318	0.281	0.340	0.595	0.586
	336	0.266	0.323	<u>0.278</u>	<u>0.329</u>	0.281	0.342	0.295	0.355	0.304	0.348	0.324	0.364	0.339	0.372	1.270	0.871
	720	0.344	0.372	0.367	<u>0.385</u>	0.397	0.421	0.422	0.445	0.385	0.400	0.410	0.420	0.433	0.432	3.001	1.267

Table 2: Multivariate long-term forecasting results with our model PatchMixer. The prediction lengths $T \in \{96, 192, 336, 720\}$ for all datasets. We report $L = 336$ for PatchMixer and PatchTST, and the best result in $L = 24, 48, 96, 192, 336$ for the other baseline models by default. Thus it could be a strong baseline. The best results are in **bold** and the second best results are in underlined.

Dataset	Electricity (Large)	Weather (Medium)	ETTh2 (Small)
Patch + DWConv + Dual Heads	<u>0.200</u>	0.305	<u>0.393</u>
DWConv + Dual Heads	0.211	0.322	0.415
Patch + Linear Head	0.210	0.329	0.383
Patch + MLP Head	0.199	0.321	0.403
Patch + Dual Heads	<u>0.200</u>	0.320	0.400
Patch + Attention + Dual Heads	0.203	0.319	0.396
Patch + Attention + Dual Heads (wo. pos.)	0.202	0.320	0.395

Table 3: Ablation Study across diverse dataset sizes, under $L = 336$ and $T = 720$ on each dataset. MSE evaluates the results, while the best are in **bold** and the second best results are in underlined.

performance gains at larger patch sizes.

Varying Patch Strides. In Figure 6, as the patch stride increases, the MSE oscillates irregularly within a range of ± 0.002 . Besides, the lowest loss function points differ across datasets, suggesting that the patch stride parameter doesn't significantly aid in optimization.

Loss Function. We study the effects of different loss functions in Table 5. LTSF tasks have mainly used MSE, while a few models like [Liu *et al.*, 2022a] utilized MAE. Recent research has introduced SmoothL1Loss [Lin *et al.*, 2023], which combines aspects of both MSE and MAE. This observation leads us to explore a better loss approach.

Network Type	Computational Complexity	MACs
Attention Mechanism	$O(N \cdot D^2 + N^2 \cdot D)$	293.63M
Standard Convolution	$O(N^2 \cdot D \cdot K)$	175.57M
PatchMixer Block	$O(N^2 \cdot D + N \cdot D \cdot K)$	66.32M

Table 4: Comparison of efficiency between patch-based Transformer and CNN models under $L = 336$ and $T = 720$ on the ETTm1 dataset. Note that $D >> N > K$.

6 Conclusion

In LTSF, since intra-variable mutual information in time series benchmark datasets is denser than that among variables, we hypothesize that patch mixing outperforms channel mixing in predictive performance. We validate this through PatchMixer, a CNN-based model. Our experiments demonstrate that even with models of limited receptive field, the patch-mixing architecture can improve predictive performance while accelerating training and inference. Lastly, we emphasize the importance of data optimization in LTSF tasks, which can enhance predictive performance through patch embedding and loss function optimization.

7 Acknowledgements

This work was supported by the National Natural Science Foundation of China (No. 62306257). The views and conclu-

Models		PatchMixer								PatchTST									
		MSE+MAE		MSE		MAE		SmoothL1loss		MSE+MAE		MSE		MAE		SmoothL1loss			
Metric		MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE		
Weather	96	0.149	0.193	0.154	0.196	0.152	0.190	0.151	0.193	0.151	0.192	0.152	0.199	0.152	0.186	0.150	0.191		
	192	0.191	0.233	0.197	0.237	0.193	0.231	0.194	0.237	0.195	0.234	0.197	0.243	0.197	0.230	0.196	0.236		
	336	0.225	0.269	0.225	0.267	0.227	0.265	0.229	0.272	0.247	0.275	0.249	0.283	0.250	0.273	0.250	0.278		
	720	0.307	0.324	0.302	0.322	0.308	0.321	0.309	0.326	0.321	0.328	0.320	0.335	0.322	0.326	0.321	0.330		
Traffic	96	0.362	0.242	0.370	0.252	0.369	0.237	0.367	0.252	0.364	0.240	0.367	0.251	0.379	0.231	0.382	0.234		
	192	0.382	0.252	0.388	0.258	0.388	0.244	0.386	0.256	0.382	0.245	0.385	0.259	0.398	0.239	0.403	0.245		
	336	0.392	0.257	0.400	0.266	0.398	0.246	0.400	0.267	0.396	0.253	0.398	0.265	0.411	0.246	0.419	0.257		
	720	0.428	0.282	0.436	0.288	0.429	0.266	0.435	0.290	0.434	0.277	0.434	0.287	0.443	0.265	0.460	0.296		
Electricity	96	0.128	0.221	0.128	0.221	0.217	0.130	0.224	0.131	0.223	0.130	0.222	0.131	0.224	0.131	0.223	0.131	0.223	
	192	0.144	0.237	0.142	0.236	0.143	0.233	0.145	0.240	0.148	0.239	0.148	0.240	0.149	0.241	0.148	0.240	0.149	0.240
	336	0.164	0.257	0.163	0.255	0.162	0.252	0.166	0.260	0.165	0.256	0.167	0.261	0.165	0.257	0.167	0.257	0.167	0.257
	720	0.201	0.290	0.199	0.289	0.199	0.284	0.204	0.293	0.208	0.293	0.202	0.291	0.207	0.290	0.207	0.290	0.207	0.290
ETTh1	96	0.355	0.383	0.354	0.384	0.353	0.379	0.356	0.384	0.376	0.401	0.375	0.399	0.367	0.392	0.376	0.400	0.376	0.400
	192	0.373	0.394	0.376	0.397	0.376	0.392	0.375	0.394	0.411	0.418	0.414	0.421	0.411	0.416	0.412	0.418	0.412	0.418
	336	0.391	0.410	0.397	0.421	0.396	0.410	0.394	0.411	0.429	0.432	0.431	0.436	0.431	0.427	0.430	0.431	0.430	0.431
	720	0.446	0.463	0.446	0.462	0.437	0.450	0.444	0.462	0.445	0.462	0.449	0.466	0.443	0.455	0.442	0.460	0.442	0.460
ETTh2	96	0.220	0.298	0.226	0.300	0.224	0.296	0.222	0.298	0.275	0.334	0.274	0.336	0.277	0.331	0.276	0.334	0.277	0.334
	192	0.267	0.332	0.276	0.335	0.272	0.331	0.270	0.333	0.340	0.375	0.339	0.379	0.343	0.374	0.341	0.375	0.343	0.375
	336	0.304	0.363	0.319	0.368	0.311	0.364	0.307	0.364	0.329	0.378	0.331	0.380	0.333	0.378	0.331	0.378	0.331	0.378
	720	0.375	0.417	0.395	0.427	0.380	0.416	0.377	0.417	0.378	0.419	0.379	0.422	0.382	0.417	0.380	0.419	0.380	0.419
ETTm1	96	0.290	0.340	0.292	0.341	0.290	0.334	0.289	0.339	0.290	0.338	0.290	0.342	0.294	0.330	0.294	0.330	0.294	0.330
	192	0.325	0.361	0.326	0.362	0.328	0.357	0.327	0.362	0.334	0.365	0.332	0.369	0.339	0.359	0.337	0.358	0.337	0.358
	336	0.353	0.382	0.354	0.382	0.355	0.377	0.355	0.382	0.359	0.382	0.366	0.392	0.361	0.378	0.362	0.378	0.361	0.378
	720	0.413	0.413	0.417	0.413	0.415	0.409	0.416	0.413	0.421	0.420	0.420	0.424	0.415	0.414	0.415	0.414	0.415	0.414
ETTm2	96	0.164	0.251	0.168	0.253	0.165	0.249	0.164	0.251	0.165	0.250	0.165	0.255	0.164	0.246	0.164	0.246	0.164	0.246
	192	0.220	0.291	0.224	0.291	0.219	0.285	0.219	0.289	0.219	0.289	0.220	0.292	0.215	0.283	0.218	0.285	0.218	0.285
	336	0.264	0.322	0.265	0.320	0.265	0.318	0.261	0.317	0.275	0.326	0.278	0.329	0.270	0.320	0.270	0.323	0.270	0.323
	720	0.342	0.375	0.343	0.370	0.347	0.370	0.345	0.371	0.365	0.382	0.367	0.385	0.355	0.374	0.363	0.380	0.363	0.380
Avg.		0.292	0.316	0.296	0.318	0.294	0.312	0.294	0.318	0.305	0.322	0.306	0.327	0.307	0.318	0.309	0.322	0.309	0.322

Table 5: Ablation study of loss functions for training in PatchMixer. 4 cases are included: (a) both MSE and MAE are included in loss function; (b) MSE; (c) MAE; (d) SmoothL1loss. The best results are in **bold**.

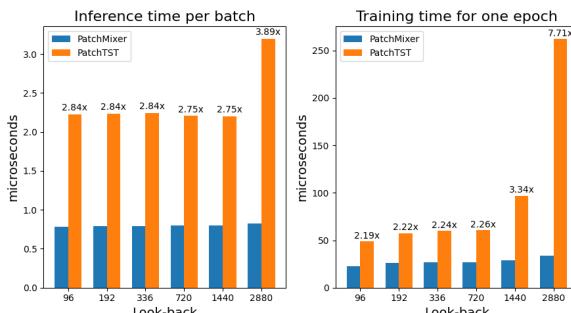


Figure 4: PatchMixer vs. PatchTST: Comparison of Efficiencies.

sions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of the National Natural Science Foundation.

References

- [Bai *et al.*, 2018] Shaojie Bai, J Zico Kolter, and Vladlen Koltun. An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. *arXiv preprint arXiv:1803.01271*, 2018.
- [Chollet, 2017] François Chollet. Xception: Deep learning with depthwise separable convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1251–1258, 2017.
- [Gu *et al.*, 2021] Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured state spaces. *arXiv preprint arXiv:2111.00396*, 2021.
- [Hendrycks and Gimpel, 2016] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016.
- [Howard *et al.*, 2017] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mo-

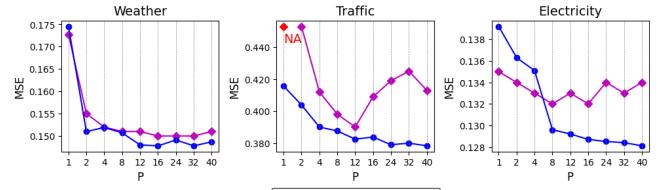


Figure 5: MSE scores with varying patch length on the top 3 largest datasets. The patch length $P = [1, 2, 4, 8, 12, 16, 24, 32, 40]$, where $L = 336$ and $T = 96$. “NA” means the setup runs out of GPU memory (NVIDIA GTX4090 24GB) even with batch size 1.

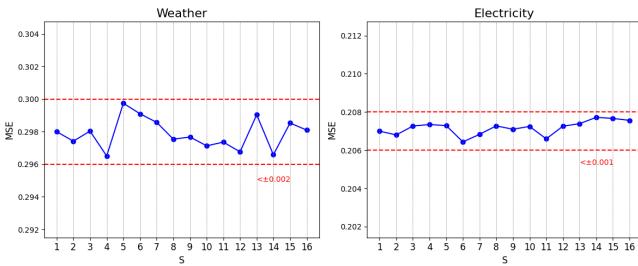


Figure 6: MSE scores with varying patch strides S from 1 to 16 where the look-back window is 336 and the prediction length is 720.

bile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.

[Ioffe and Szegedy, 2015] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. pmlr, 2015.

[Kim *et al.*, 2021] Taesung Kim, Jinhee Kim, Yunwon Tae, Cheonbok Park, Jang-Ho Choi, and Jaegul Choo. Reversible instance normalization for accurate time-series forecasting against distribution shift. In *International Conference on Learning Representations*, 2021.

[Lin *et al.*, 2023] Shengsheng Lin, Weiwei Lin, Wentai Wu, Songbo Wang, and Yongxiang Wang. Petformer: Long-term time series forecasting via placeholder-enhanced transformer. *arXiv preprint arXiv:2308.04791*, 2023.

[Liu *et al.*, 2022a] Minhao Liu, Ailing Zeng, Muxi Chen, Zhijian Xu, Qixia Lai, Lingna Ma, and Qiang Xu. Scinet: Time series modeling and forecasting with sample convolution and interaction. *Advances in Neural Information Processing Systems*, 35:5816–5828, 2022.

[Liu *et al.*, 2022b] Shizhan Liu, Hang Yu, Cong Liao, Jian-guo Li, Weiyao Lin, Alex X Liu, and Schahram Dustdar. Pyraformer: Low-complexity pyramidal attention for long-range time series modeling and forecasting. In *International Conference on Learning Representations*, 2022.

[Nie *et al.*, 2023] Yuqi Nie, Nam H. Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth 64 words: Long-term forecasting with transformers. In *International Conference on Learning Representations*, 2023.

[Poli *et al.*, 2023] Michael Poli, Stefano Massaroli, Eric Nguyen, Daniel Y Fu, Tri Dao, Stephen Baccus, Yoshua Bengio, Stefano Ermon, and Christopher Ré. Hyena hierarchy: Towards larger convolutional language models. *arXiv preprint arXiv:2302.10866*, 2023.

[Romero *et al.*, 2021] David W Romero, Anna Kuzina, Erik J Bekkers, Jakub M Tomczak, and Mark Hoogendoorn. Ckconv: Continuous kernel convolution for sequential data. *arXiv preprint arXiv:2102.02611*, 2021.

[Shannon, 1948] Claude Elwood Shannon. A mathematical theory of communication. *The Bell system technical journal*, 27(3):379–423, 1948.

[Sifre and Mallat, 2014] Laurent Sifre and Stéphane Mallat. Rigid-motion scattering for texture classification. *arXiv preprint arXiv:1403.1687*, 2014.

[Szegedy *et al.*, 2015] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1–9, 2015.

[Trockman and Kolter, 2022] Asher Trockman and J Zico Kolter. Patches are all you need? *arXiv preprint arXiv:2201.09792*, 2022.

[Ulyanov *et al.*, 2016] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Instance normalization: The missing ingredient for fast stylization. *arXiv preprint arXiv:1607.08022*, 2016.

[Wang *et al.*, 2023] Huiqiang Wang, Jian Peng, Feihu Huang, Jince Wang, Junhui Chen, and Yifei Xiao. Micn: Multi-scale local and global context modeling for long-term series forecasting. In *International Conference on Learning Representations*, 2023.

[Wu *et al.*, 2021] Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Autoformer: Decomposition transformers with Auto-Correlation for long-term series forecasting. In *Advances in Neural Information Processing Systems*, 2021.

[Wu *et al.*, 2023] Haixu Wu, Tengge Hu, Yong Liu, Hang Zhou, Jianmin Wang, and Mingsheng Long. Timesnet: Temporal 2d-variation modeling for general time series analysis. In *International Conference on Learning Representations*, 2023.

[Zeng *et al.*, 2023] Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are transformers effective for time series forecasting? In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, pages 11121–11128, 2023.

[Zhou *et al.*, 2021] Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang. Informer: Beyond efficient transformer for long sequence time-series forecasting. In *The Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Virtual Conference*, volume 35, pages 11106–11115. AAAI Press, 2021.

[Zhou *et al.*, 2022] Tian Zhou, Ziqing Ma, Qingsong Wen, Xue Wang, Liang Sun, and Rong Jin. FEDformer: Frequency enhanced decomposed transformer for long-term series forecasting. In *Proc. 39th International Conference on Machine Learning*, 2022.

[Zhou *et al.*, 2023] Jiaming Zhou, Kun-Yu Lin, Yu-Kun Qiu, and Wei-Shi Zheng. Twinformer: Fine-to-coarse temporal modeling for long-term action recognition. *IEEE Transactions on Multimedia*, pages 1–14, 2023.