# PHASEFORMER: FROM PATCHES TO PHASES FOR EF-FICIENT AND EFFECTIVE TIME SERIES FORECASTING

**Anonymous authors**Paper under double-blind review

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

023

025

026027028

029

031

033

034

035

037

038

040

041

042 043

044

046

047

048

049

051

052

#### **ABSTRACT**

Periodicity is a fundamental characteristic of time series data and has long played a central role in forecasting. Recent deep learning methods strengthen the exploitation of periodicity by treating patches as basic tokens, thereby improving predictive effectiveness. However, their efficiency remains a bottleneck due to large parameter counts and heavy computational costs. This paper provides, for the first time, a clear explanation of why patch-level processing is inherently inefficient, supported by strong evidence from real-world data. To address these limitations, we introduce a phase perspective for modeling periodicity and present an efficient yet effective solution, PhaseFormer. PhaseFormer features phasewise prediction through compact phase embeddings and efficient cross-phase interaction enabled by a lightweight routing mechanism. Extensive experiments demonstrate that PhaseFormer achieves state-of-the-art performance with around 1k parameters, consistently across benchmark datasets. Notably, it excels on large-scale and complex datasets, where models with comparable efficiency often struggle. This work marks a significant step toward truly efficient and effective time series forecasting. Code is available at this repository: https: //anonymous.4open.science/r/ICLR26-PhaseFormer-17678.

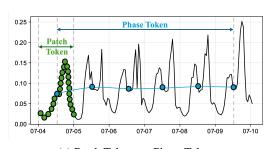
## 1 Introduction

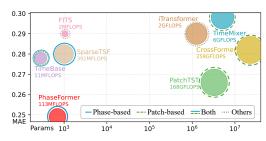
Time series forecasting underpins decision-making across diverse domains such as finance, energy, climate science, and healthcare, playing a pivotal role in tasks including weather forecasting (Qureshi et al., 2025; Wu et al., 2021a), energy consumption planning (Lai et al., 2018; Alvarez et al., 2010; Cheng et al., 2021), traffic scheduling (Cirstea et al., 2022; 2021; Wu et al., 2021b). In recent years, deep learning has demonstrated promising potential in this field by leveraging end-to-end modeling and powerful representational capacity to extrapolate from history to future trends.

A central inductive bias in forecasting models is periodicity—the recurring temporal structure inherent in many real-world time series. Recent advances exploited this property by segmenting sequences into patch tokens, potentially aligned with cycles, prior to processing by the crafted models (Nie et al., 2023; Zhang & Yan, 2023; Huang et al., 2025; Tang & Zhang, 2025). For instance, Nie et al. (2023) applied Transformer to tokenized time series to capture temporal correlations within and between cycles, while Zhang & Yan (2023) extended this paradigm by modeling cross-dimension dependencies and cross-scale interactions.

Despite their effectiveness, patch-based approaches struggle to scale efficiently to large and complex datasets (Nie et al., 2023; Zhang & Yan, 2023; Tang & Zhang, 2025). We attribute this poor scalability to the substantial variability of cycle patterns in real-world scenarios. This variability stems from dynamic external factors, which continuously shift the cycle patterns. For instance, traffic flow patterns may evolve as new infrastructure is introduced, while electricity demand can change with adjustments in work schedules. This variability forces models to construct a high-dimensional representation space to faithfully accommodate the broadened distribution, which inevitably inflates both parameter counts and computational costs (Nie et al., 2023; Zhang & Yan, 2023). Additionally, these methods also struggle to generalize under such varying behavior, resulting in unreliable forecasts for samples beyond training data.

To address this challenge, we introduce a novel phase-based perspective that focuses on values aligned at the same offset across successive cycles. From this perspective, the dynamics of a time





- (a) Patch Token vs. Phase Token.
- (b) Model Accuracy and Efficiency Comparison.

Figure 1: Comparison between patch-based and phase-based representations for time-series fore-casting. (a) illustrates the difference in tokenization. (b) jointly evaluates model accuracy, parameter scale, and computational overhead on the Traffic dataset, where marker size indicates FLOPS.

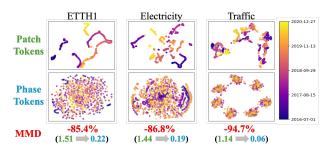
series are characterized by the cross-period trends of each phase–captured as phase tokens—while disregarding the full cyclic behavior. As illustrated in Fig. 1a, phase tokens exhibit significantly lower variability than patch tokens, enabling more efficient and generalizable representation. Importantly, excluding cycle patterns has minimal impact on forecasting effectiveness, since the cyclic behaviors remain locally stable and thus require little effort to predict. We study and verify these properties in depth in Sec. 3 using real-world data, showing the stationarity and compactness of the feature space offered by phase tokenization.

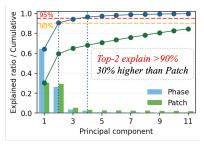
Building on these insights, we propose Phase-based Routing Transformer, abbreviated as **Phase-Former**, which reframes time series as a collection of phase tokens and casts step-wise prediction as phase-wise prediction. Specifically, PhaseFormer (i) aligns and extracts phase tokens from the input sequence and maps them into a shared low-dimensional latent space, (ii) employs a lightweight routing mechanism to enable efficient communication across phases, and (iii) applies a shared predictor to project the latent representations into forecasts for each phase. Extensive experiments demonstrate that, compared with PatchTST (Nie et al., 2023) and Crossformer (Zhang & Yan, 2023), PhaseFormer achieves over **99.9%** reduction in both parameter count and computational cost, while delivering consistent improvements in prediction accuracy across all seven benchmark datasets, as illustrated by the Traffic dataset in Fig. 1b. Moreover, in contrast to methods with comparable efficiency such as SparseTSF (Lin et al., 2024) and TimeBase (Huang et al., 2025), PhaseFormer significantly enhances predictive effectiveness, particularly on large and complex datasets. Finally, we conduct a comprehensive analysis of different configurations to reveal the necessity of the constructed components and the effects of various hyperparameters. Our contributions are as follows:

- 1. We introduce a phase-based perspective that aligns values across cycles for the characterization of long-term time series, empirically and theoretically demonstrating improved feature stationarity and compactness over the patch-based perspective.
- 2. We propose PhaseFormer, a lightweight forecasting model that reframes time series as phase tokens, maps them into a shared latent space, and employs a routing mechanism with a shared predictor to enable efficient phase-wise forecasting.
- Extensive experiments are conducted to showcase that PhaseFormer achieves substantial efficiency gains while consistently improving forecasting accuracy, establishing a superior efficiency-effectiveness trade-off across diverse benchmarks.

#### 2 RELATED WORKS

**Transformer-Based Forecasting Architectures.** Early Transformer-based models for long sequence forecasting often overlooked the periodicity in time series (Zhou et al., 2021; Li et al., 2019). Subsequent research introduced domain-specific priors that better understand recurring temporal structures. Autoformer (Wu et al., 2021a) and FEDformer (Zhou et al., 2022) incorporated decomposition strategies and frequency-domain modeling enabling explicit representation of seasonal–trend patterns. Pyraformer (Liu et al., 2021) and Crossformer (Zhang & Yan, 2023) further enriched temporal modeling by embedding multi-scale hierarchies and cross-variable dependencies,





- (a) Temporal shift comparison across multiple datasets.
- (b) PCA for phase tokens on Traffic.

Figure 2: Visualization of phase tokenization and its advantages. (a) Phase tokenization yields more stable representations than patch-based embeddings. (b) Phase tokens exhibit clear low-dimensionality compared with patch tokens.

while Liu et al. (2022) explicitly accounted for distributional shifts. More recently, PatchTST (Nie et al., 2023) reframed time series as patch sequences to enable more accurate characterization of sequence-level semantics, followed by an extension to jointly consider spatial and temporal correlations (Huo et al., 2025). Generally speaking, these models embed progressively stronger temporal biases, though often at the cost of massive parameter counts and heavy computation.

Efficiency-Oriented Forecasting Models. A growing body of research emphasizes efficiency, aiming to design lightweight forecasting architectures. Patch-based MLP variants such as xPatch (Stitsyuk & Choi, 2025), TimeMixer (Wang et al., 2024), and PITS (Lee et al., 2024) exploited compact tokenization or hierarchical dependencies to reduce parameter counts while maintaining accuracy. Beyond patches-based methods, frequency-based counterparts leverage spectral representations for compression and denoising. FreTS (Yi et al., 2023) applied MLPs in the frequency domain, Yi et al. (2024) learned frequency filters to improve noise robustness, and FITS (Xu et al., 2024) achieved strong accuracy with only 10k parameters. Deng et al. (2024) demonstrated that selective decomposition can deliver both parsimony and capability. More recently, SparseTSF (Lin et al., 2024) and TimeBase (Huang et al., 2025) highlighted the importance of cross-period correlation, sharing a similar motivation with ours. Despite their impressive computational efficiency, these methods still fall short on forecasting accuracy for large and complex datasets such as Traffic and Electricity (Lin et al., 2024; Huang et al., 2025). Moreover, they lack systematic analysis to answer the fundamental question: Why can phase tokens serve as an efficient alternative to patch tokens?

# 3 MOTIVATIONS

To motivate our approach, we conduct a comparative analysis of the geometric structures of patch and phase tokens across three widely used datasets. As illustrated in Fig. 1a, a patch token is composed of adjacent observations within a local period, whereas a phase token is constructed by extracting values at identical offsets across consecutive periods. We gain the following two important insights from the thorough analysis.

**Insight 1: Phase tokens are globally stationary, while patch tokens are locally stationary.** To provide an intuitive overview of their geometric structures, we project both types of tokens into two-dimensional spaces using t-SNE (van der Maaten & Hinton, 2008). As shown in Fig. 2a, the distributions of patch tokens drift continuously over time but exhibit local coherence, indicating *local stationarity* and supporting the minimal impact of excluding cycle patterns from intensive processing. In contrast, phase tokens form compact and coherent clusters that remain stable over the long term, reflecting strong *global stationarity*. To rigorously quantify the long-term drift, we compute the average discrepancy distance between each subsequent week and the initial week. Specifically, we adopt the Maximum Mean Discrepancy (MMD) metric (Ouyang & Key, 2021), a statistical measure of distributional divergence:

$$MMD^{2}(P,Q) = \mathbb{E}_{x,x' \sim P}[k(x,x')] + \mathbb{E}_{y,y' \sim Q}[k(y,y')] - 2 \,\mathbb{E}_{x \sim P,y \sim Q}[k(x,y)], \tag{1}$$

where P and Q denote tokens collected from two different weeks, respectively, and  $k(\cdot,\cdot)$  is the RBF kernel function. As two distributions become closer, their MMD value approaches zero. The results at the bottom of Fig. 2a show that the average MMD distance of the phase token space is significantly smaller than that of the patch token space. Taken together, both qualitative and quantitative analyses demonstrate that phase tokenization exhibits substantially lower temporal distribution divergence, thereby facilitating better generalization across the time axis.

**Insight 2: Phase tokens reside in a lower-dimensional subspace than patch tokens.** To measure the effective dimensionality of the token space, we perform principal component analysis (PCA) on it. Surprisingly, as illustrated in Fig. 2b, two dimensions are already sufficient to explain over 90% of the variance of phase tokens, whereas patch tokens require more than eleven dimensions to achieve the same degree of explanation, owing to their drifting behavior observed in Fig. 2a. Consequently, phase information resides in a low-dimensional subspace, *providing a principled basis for parameter- and computation-efficient modeling*.

We further establish, based on perturbation theory, that phase tokenization remains stable under perturbations of cycle patterns, whereas patch tokenization undergoes structural shifts. Due to space limitations, only the core theorem is presented here, while the detailed proof is provided in Sec. A.6.

**Theorem 1 (Phase Tokenization Stability)** Let  $X = AG^{\top} + N \in \mathbb{R}^{D \times H}$  with  $\operatorname{rank}(A) = \operatorname{rank}(G) = r \ll \min(D, H)$ , and consider the transformed data

$$X' = XS^{\top} + R,\tag{2}$$

where  $||N'||_2 \le ||S||_2 ||N||_2$ ,  $||R||_2 \le \varepsilon (||M||_F + ||N||_F)$ , and let  $\delta_{\min} > 0$  denote the minimal spectral separation. Then there exists a universal constant C > 0 such that:

1. For phase tokenization and corresponding subspace  $U_r$ , there exists:

$$d(\mathcal{U}_r(X), \mathcal{U}_r(X')) \le C \frac{\|N\|_2 + \|N'\|_2 + \|R\|_2}{\delta_{\min}},$$
 (3)

with exact invariance in the noiseless case (N = R = 0).

2. For patch tokenization and corresponding subspace  $V_r$ , there exists:

$$d(\mathcal{V}_r(X), \mathcal{V}_r(X')) \ge d(\text{Col}(G), \text{Col}(SG)) - C \frac{\|N\|_2 + \|N'\|_2 + \|R\|_2}{\delta_{\min}}.$$
 (4)

**Takeaways.** Phase tokenization is structurally invariant under the cycle pattern change S and only subject to perturbations from noise and small day-to-day mismatches. In contrast, patch tokenization generally suffers from a non-vanishing structural offset. Hence, phase tokenization is more robust and consistent under cycle pattern drifts.

## 4 METHODOLOGY

Given the focus on periodicity, we adopt the channel-independent paradigm (Nie et al., 2023; Zeng et al., 2023) and omit the channel dimension throughout the remainder of this paper. The objective of forecasting is to predict the future trajectory  $\mathbf{Y} \in \mathbb{R}^{L_{\text{out}}}$  from an input sequence  $\mathbf{X} \in \mathbb{R}^{L_{\text{in}}}$ , where  $L_{\text{in}}$  and  $L_{\text{out}}$  denote the input and output lengths, respectively. In the following sections, we describe the data preprocessing procedure, present the proposed network architecture, and finally analyze the computational complexity of the method.

#### 4.1 DATA PRE-PROCESSING

**Normalization and De-Normalization.** Following Kim et al. (2021), we normalize inputs with their estimated mean and standard deviation, and de-normalize predictions to the original scale.

**Phase Tokenization and De-Tokenization.** Phase tokenization transforms the one-dimensional input sequence into a two-dimensional phase–period matrix for the following processing. Conversely, phase de-tokenization reconstructs the predicted phase–period matrix back into a one-dimensional

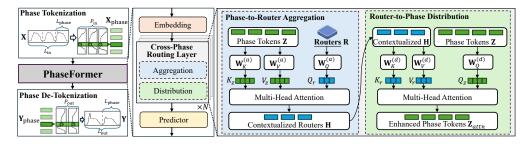


Figure 3: The overview of PhaseFormer.

output sequence. Let  $L_{\rm phase}$  denote the period length, which can be estimated using autocorrelation analysis.

To ensure that the input sequence length is a multiple of  $L_{\rm phase}$ , we circularly pad the sequence to length  $P_{\rm in}*L_{\rm phase}$ , where  $P_{\rm in}=\left\lceil\frac{L_{\rm in}}{L_{\rm phase}}\right\rceil$ . As illustrated in Fig. 3, the padded sequence  ${\bf X}$  is then reshaped into a phase–period matrix  ${\bf X}_{\rm phase}\in\mathbb{R}^{L_{\rm phase}\times P_{\rm in}}$ , where each entry  ${\bf X}_{\rm phase}[\ell,p]$  corresponds to the observation at the  $\ell^{\rm th}$  phase of the  $p^{\rm th}$  period. In the de-tokenization process, the predicted phase–period matrix is mapped back to the temporal domain by reversing the transformation, thereby reconstructing the final one-dimensional forecast sequence.

#### 4.2 Phase-based Routing Transformer

The phase–period matrix is fed into our proposed phase-based routing Transformer, termed Phase-Former, to capture and extrapolate temporal dynamics at the phase level in an efficient and effective way. As illustrated in Fig. 3, PhaseFormer first applies an embedding layer to the phase tokens, then refines them through multiple cross-phase routing layers, and finally maps them to the target via a shared predictor. Next, we elaborate on the design of these modules in detail.

## 4.2.1 Embedding Layer

The embedding layer projects the phase tokens  $\mathbf{X}_{\text{phase}}$  into a low-dimensional representation space, allowing the informative components to be extracted from raw observations that are often contaminated by perturbations. Formally, for each phase index  $\ell \in \{1,\ldots,L_{\text{phase}}\}$ , the corresponding phase token  $\mathbf{X}_{\text{phase}}[\ell,:]$  is mapped into a d-dimensional representation through a linear function  $f_{\theta}$ , parameterized by  $\theta \in \mathbb{R}^{P_{\text{in}} \times d}$ :

$$\mathbf{Z} = f_{\theta}(\mathbf{X}_{\text{phase}}) \in \mathbb{R}^{L_{\text{phase}} \times d} \tag{5}$$

To better capture the temporal ordering among phases, we introduce a set of learnable positional embeddings  $\mathbf{E}_{pos} \in \mathbb{R}^{L_{phase} \times d}$  to distinguish the relative position of each phase, following Liu et al. (2023a). These embeddings are added to  $\mathbf{Z}$  in a phase-wise manner, so that each phase representation is enriched with its positional information:

$$\tilde{\mathbf{Z}} = \mathbf{Z} + \mathbf{E}_{pos}. \tag{6}$$

The resulting  $\dot{\mathbf{Z}}$  is then forwarded to the cross-phase routing layers for higher-level feature interaction and forecasting.

## 4.2.2 Cross-Phase Routing Layer

Directly modeling full pairwise interactions among phase representations via self-attention is computationally expensive. To handle this, we introduce a set of learnable routers  $\mathbf{R} \in \mathbb{R}^{M \times d}$  to mediate information exchange across phases, drawing inspiration from previous methods (Jaegle et al., 2021; Zhang & Yan, 2023). This design substantially reduces the quadratic cost of self-attention while preserving rich cross-phase dependencies.

Cross-phase routing consists of two steps: (i) *phase-to-router aggregation*, which selectively compresses information from phase representations into the compact set of routers; and (ii) *router-to-*

*phase distribution*, which selectively propagates the aggregated cross-phase information from the routers back to the phase representations. Both steps are implemented via cross-attention, allowing the model to scale efficiently while preserving strong representational capacity.

**Phase-to-Router Aggregation.** The routers attend to the phase representations to extract contextual information, yielding contextualized router embeddings  $\mathbf{H} \in \mathbb{R}^{M \times d}$ . Specifically, the routers act as queries while the phases provide keys and values. The projection matrices  $\mathbf{W}_Q^{\mathrm{agg}}, \mathbf{W}_K^{\mathrm{agg}}, \mathbf{W}_V^{\mathrm{agg}} \in \mathbb{R}^{d \times d}$  map the representations into query, key, and value spaces, respectively:

$$\mathbf{Q}_r = \mathbf{R} \mathbf{W}_O^{\text{agg}}, \quad \mathbf{K}_z = \tilde{\mathbf{Z}} \mathbf{W}_K^{\text{agg}}, \quad \mathbf{V}_z = \tilde{\mathbf{Z}} \mathbf{W}_V^{\text{agg}}.$$
 (7)

The aggregated router embeddings are then obtained via multi-head attention (MHA) with  $d_h$  heads:

$$\mathbf{H} = \mathrm{MHA}(\mathbf{Q}_r, \mathbf{K}_z, \mathbf{V}_z). \tag{8}$$

**Router-to-Phase Distribution.** The aggregated information in the routers is subsequently redistributed to the phase representations, thereby enabling cross-phase information flow. In this step, the phase representations serve as queries while the routers provide keys and values, yielding refined phase representations  $\mathbf{Z}_{\text{attn}}$ . The projection matrices  $\mathbf{W}_Q^{\text{dist}}, \mathbf{W}_K^{\text{dist}}, \mathbf{W}_V^{\text{dist}} \in \mathbb{R}^{d \times d}$  are used for this distribution:

$$\mathbf{Q}_{z} = \tilde{\mathbf{Z}} \mathbf{W}_{Q}^{\text{dist}}, \quad \mathbf{K}_{r} = \mathbf{H} \mathbf{W}_{K}^{\text{dist}}, \quad \mathbf{V}_{r} = \mathbf{H} \mathbf{W}_{V}^{\text{dist}},$$
(9)

$$\mathbf{Z}_{\text{attn}} = \text{MHA}(\mathbf{Q}_z, \mathbf{K}_r, \mathbf{V}_r). \tag{10}$$

This mechanism restores phase-level resolution while simultaneously enforcing coherence across phases through the contextualized routers. Ultimately, each phase representation attends to all others through a two-stage routing pathway.

#### 4.2.3 PREDICTOR

The predictor produces multi-step forecasts of length  $P_{\text{out}}$  for all phases simultaneously, based on their refined representations. Taking as input the refined phase representations  $\mathbf{Z}_{\text{attn}} \in \mathbb{R}^{L_{\text{phase}} \times d}$  from the final cross-phase routing layer, the predictor is realized as a linear mapping  $g_{\phi}$ , parameterized by  $\phi \in \mathbb{R}^{d \times P_{\text{out}}}$ :

$$\mathbf{Y}_{\text{phase}} = q_{\phi}(\mathbf{Z}_{\text{attn}}) \in \mathbb{R}^{L_{\text{phase}} \times P_{\text{out}}}.$$
 (11)

All phases share the same predictor parameters, which enforces consistency across phases and reduces the number of trainable parameters. This not only improves efficiency but also regularizes learning, thereby enhancing generalization. Finally, the predicted phase–period matrix  $\mathbf{Y}_{phase}$  is passed through de-tokenization and de-normalization to produce the final forecast  $\mathbf{Y}$ .

#### 4.3 Complexity of PhaseFormer

For each variable, the overall complexity of PhaseFormer can be summarized as follows: the phase embedding layer requires  $O(L_{\rm phase}P_{\rm in}d)$  time and  $O(L_{\rm phase}d)$  memory. The cross-phase routing layer, which dominates computation, incurs  $O((L_{\rm phase}+M)d^2+ML_{\rm phase}d)$  time and  $O(HML_{\rm phase}+(L_{\rm phase}+M)d)$  memory. Finally, the predictor costs  $O(L_{\rm phase}dP_{\rm out})$  time and  $O(L_{\rm phase}P_{\rm out})$  memory. Aggregating across N blocks, the end-to-end time complexity is:

$$O(N((L_{\text{phase}} + M)d^2 + ML_{\text{phase}}d) + L_{\text{phase}}d(P_{\text{in}} + P_{\text{out}})).$$

Substituting  $P_{\rm in} = \lceil L_{\rm in}/L_{\rm phase} \rceil$  and  $P_{\rm out} = \lceil L_{\rm out}/L_{\rm phase} \rceil$  into the above expression gives:

$$O(N((L_{\text{phase}} + M)d^2 + ML_{\text{phase}}d) + d(L_{\text{in}} + L_{\text{out}})).$$

As investigated in Sec. 3, the phase token space exhibits a inherently low-dimensional structure, which allows M and d to be chosen as fixed and small numbers. Thus, the computational cost grows in a linear manner with both the input length  $L_{\rm in}$  and the output horizon  $L_{\rm out}$ .

Table 1: Main results for long-term forecasting. The input sequence length is  $L_{\text{input}} = 720$ , and results are averaged over forecast horizons  $L_{\text{out}} \in \{96, 192, 336, 720\}$ . The best results are shown in **bold**, and the second-best in <u>underline</u>.

Dataset	PhaseFormer		PatchTST		iTransformer		Crossformer		FEDformer		TimeBase		SparseTSF		FITS		TimeMixer	
Dataset	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
ETTh1	0.403	0.415	0.420	0.439	0.453	0.467	0.517	0.512	0.523	0.523	0.404	0.416	0.406	0.418	0.419	0.435	0.452	0.474
ETTh2	0.346	0.388	0.344	0.390	0.392	0.422	1.468	0.867	0.428	0.469	0.347	0.397	0.345	0.383	0.334	0.382	0.386	0.425
ETTm1	0.346	0.374	0.354	0.383	0.370	0.401	0.390	0.417	0.438	0.465	0.356	0.380	0.362	0.383	0.359	0.382	0.383	0.413
ETTm2	0.250	0.313	0.251	0.319	0.278	0.337	0.392	0.426	0.401	0.452	0.250	0.314	0.252	0.316	0.285	0.336	0.314	0.367
Electricity	0.160	0.250	0.169	0.265	0.165	0.263	0.180	0.273	0.235	0.348	0.167	0.258	0.168	0.263	0.172	0.270	0.171	0.273
Traffic	0.386	0.249	0.394	0.266	0.406	0.290	0.545	0.282	0.638	0.400	0.418	0.278	0.413	0.280	0.410	0.290	0.421	0.298
Weather	0.223	0.260	0.223	0.264	0.233	0.273	0.255	0.304	0.354	0.393	0.227	0.262	0.243	0.285	0.241	0.283	0.237	0.281

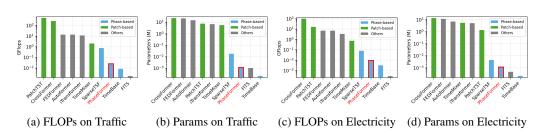


Figure 4: Comparison of FLOPs and parameter counts across models on the Traffic and Electricity. Patch-based models are shown in green, phase-based models in blue, and other models in gray.

## 5 EXPERIMENTS

#### 5.1 Long-term Time Series Forecasting

We conduct a joint evaluation of model efficiency and predictive accuracy. The comparative analysis highlights that the proposed *PhaseFormer* establishes an improved effectiveness-efficiency tradeoff in terms of parameter scale and error metrics. We also provide the code in https://anonymous.4open.science/r/ICLR26-PhaseFormer-17678.

**Datasets and Setup.** Experiments are performed on seven widely used long-term time series fore-casting datasets: *ETTh1*, *ETTh2*, *ETTm1*, *ETTm2*, *Weather*<sup>2</sup>, *Electricity*<sup>3</sup>, and *Traffic*<sup>4</sup>, covering a diverse range of real-world scenarios. The details of the datasets are provided in Sec. A.2. Following prior works (Nie et al., 2023; Zhang & Yan, 2023; Huang et al., 2025), we adopt a 6:2:2 split for the ETT datasets and a 7:1:2 split for the other datasets. For PhaseFormer, we report the average results over three random seeds, while for the other baselines we follow their official implementations and released code. We evaluate the forecasting **accuracy** of all tested models using mean squared error (MSE) and mean absolute error (MAE), and assess **efficiency** in terms of floating-point operations (FLOPs) and the number of parameters (Params).

Baselines and Implementation Details. We evaluate our approach against eight competitive baselines, encompassing both state-of-the-art Transformer architectures and efficiency-oriented forecasting models. We compare our method with PatchTST(2023), iTransformer(2023b), Crossformer(2023), FEDformer(2022), TimeBase(2025), SparseTSF(2024), FITS(2024), and TimeMixer(2024). Among these, PatchTST, Crossformer, and TimeMixer are patch-based; SparseTSF is phase-based; TimeBase integrates patch and phase paradigms; FITS and FEDformer are frequency-domain; and iTransformer models the full sequence directly. For all baselines, we adopt the recommended configurations provided in their official implementations. The model is optimized using the Adam optimizer with a fixed learning rate of  $1 \times 10^{-3}$ . Following the settings from efficiency-oriented works (Huang et al., 2025; Lin et al., 2024; Xu et al., 2024), the look-back length is set to 720 time steps. More implementation details are provided in Sec. A.2.

https://github.com/zhouhaoyi/ETDataset

<sup>&</sup>lt;sup>2</sup>https://www.bgc-jena.mpg.de/wetter/

<sup>3</sup>https://archive.ics.uci.edu/ml/datasets

<sup>4</sup>https://pems.dot.ca.gov/

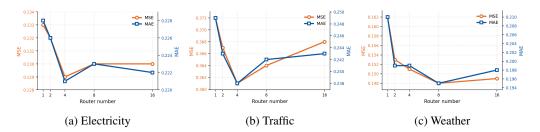


Figure 5: Effect of varying the number of routers M on forecasting performance on three datasets.

Table 2: Cross-Phase Routing layer ablation. Each cell reports MSE, MAE, and FLOPs. Lower is better for all metrics. FLOPs are reported in millions (MFLOPs). The best results are highlighted with **Bold**, and the second-best results with <u>Underlined</u>.

Dataset	P	haseForn	ier	w/	FullAtter	tion	w/ l	LinearMi	xing	w/o Routing			
	MSE	MAE	FLOPs	MSE	MAE	FLOPs	MSE	MAE	FLOPs	MSE	MAE	FLOPs	
Weather	0.1503	0.1971	3.119	0.1527	0.2005	3.202	0.1700	0.2226	0.920	0.1907	0.2406	0.783	
Electricity	0.1290	0.2209	42.213	0.1295	0.2217	48.951	0.1403	0.2334	14.068	0.1423	0.2365	11.972	
Traffic	0.3721	0.2475	113.356	0.3791	0.2513	131.452	0.3842	0.2532	37.776	0.3892	0.2584	32.149	

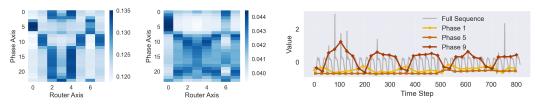
Main Results. We evaluate the predictive accuracy of PhaseFormer and the baseline methods on seven datasets. Tab. 1 reports the average prediction errors across four forecasting horizons, with detailed results provided in Sec. A.3.1. Overall, PhaseFormer consistently achieves superior performance on nearly all datasets, with particularly notable gains on complex and dynamic benchmarks such as Weather, Electricity, and Traffic. For example, on the largest dataset, Traffic, PhaseFormer surpasses the second-best method, PatchTST, by 6.3% and outperforms TimeBase by 10.4%, underscoring its robustness on large-scale and heterogeneous data. The only exception is ETTh2, where PhaseFormer ranks second to FITS while still maintaining highly competitive accuracy. A closer examination reveals that patch-based baselines, including PatchTST, Crossformer, and TimeMixer, exhibit performance degradation on the Electricity, likely due to stronger distributional shifts. In contrast, PhaseFormer delivers stable and consistently superior results under these challenging conditions, highlighting the robustness of its phase-based design.

Efficiency Comparison. We evaluate the computational overhead of all models, with detailed results in Sec. A.3.2. Fig. 4 shows the FLOPs and the number of parameters of all tested models on the Electricity and Traffic. Overall, phase-based models incur lower overhead than patch-based ones. On the Traffic dataset, PhaseFormer achieves an extraordinary FLOPs reduction of about 99.99% over PatchTST and Crossformer. Beyond patch-based baselines, it also outperforms other phase-based models like SparseTSF, consistently delivering high efficiency. This stems from the lower variety of phase tokens over time (Sec. 3), making them inherently more efficient to process. Taken together with the previous accuracy evaluations, these results clearly demonstrate that PhaseFormer provides an efficient yet effective solution, delivering superior performance on complex datasets.

## 5.2 ABLATION STUDIES AND ANALYSIS

Varying the Number of Routers. We systematically evaluate the impact of different number of routers M on model performance, with results summarized in Fig. 5. The experiments indicate that across three datasets, the model's prediction error generally decreases as the number of routers M increases, before eventually stabilizing or slightly rising. It is worth noting that the best performance is usually achieved when  $M \in \{4, 8\}$ , which is much smaller compared to the actual number of phase tokens,  $L_{\rm phase} = 24$ . This observation indicates that the phase token spans an inherently low-dimensional space, so only a small number of routers is sufficient to effectively capture and represent its underlying structure. More detailed results are provided in Sec. A.3.4.

**Effectiveness of Cross-Phase Routing.** To assess the contribution of the cross-phase routing layer, we compare four variants of the model: **PhaseFormer**, which adopts the original cross-phase routing layer; **w/ FullAttention**, which substitutes the cross-phase routing layer with a full attention mechanism; and **w/ LinearMixing**, which replaces the cross-phase routing layer with a linear layer;



(a) Aggregation Weights

- (b) Distribution Weights
- (c) Visualization of three selected phase tokens

Figure 6: Case study on a sample from the Traffic dataset. (a) Attention weight matrix during Phase-to-Router aggregation. (b) Attention weight matrix during Router-to-Phase distribution. Both matrices capture the association between 8 routers and 24 input phases. (c) Visualization of three representative phases (1, 5, and 9), each represents a distinct attentive pattern with routers.

and **w/o Routing**, which directly projects each phase into its own future. All other experimental settings are kept identical across these variants.

As summarized in Tab. 2, PhaseFormer consistently outperforms **w/ LinearMixing** and **w/o Routing**, indicating that explicit cross-phase routing is crucial for modeling periodic dynamics. Moreover, PhaseFormer not only incurs less computational and memory overhead, but also achieves lower prediction error than **w/ FullAttention**, showing that the routing layer is both efficient and effective. We attribute these gains to operating in a low-dimensional phase token space, which concentrates informative interactions and reduces cost.

#### 5.3 CASE STUDY

We select one sample from Traffic dataset, comprising an input sequence of length 720 and an output sequence of length 96 (816 time steps in total). The input sequence is fed into PhaseFormer, and we record the attention-weight matrices at the first cross-phase routing layer during both aggregation (Phase—Router) and distribution (Router—Phase). As shown in Fig. 6a and Fig. 6b, both attention patterns exhibit clear local similarity: adjacent phases tend to be assigned to the same routers and to receive attention from similar routers. This indicates that the routing mechanism captures temporally consistent phase relationships. Meanwhile, the attention weights reveal that certain phases share similar attentive patterns. To analyze this further, we focus on three phases with distinct attentive patterns and visualized them in Fig. 6c. These phases display different temporal behaviors: Phase 5 remains relatively stable over long horizons, whereas Phase 9 and Phase 1 both exhibit a pronounced 7-day periodicity but with opposite trends. The differing patterns of these phase tokens suggest that the router structure not only distinguishes among phase tokens but also effectively models their periodicity and trend characteristics.

## 6 CONCLUSION

This work identifies the inefficiencies of patch-based forecasting and presents PhaseFormer, a phase-centric model that captures periodicity via compact phase representations and lightweight cross-phase routing. Both theoretical analysis and empirical validation converge on the same conclusion that phase representations remain both more robust and more efficient than patch-based approaches under cycle pattern shifts. Consequently, PhaseFormer maintains high predictive accuracy while remaining lightweight compared to patch-based methods. More broadly, these results provide a practical pathway for building lightweight yet powerful forecasting models that retain accuracy without heavy and complex architectures.

However, the approach assumes locally stable periodicity across the input and output horizons; under highly irregular or non-repetitive cycles, phase representations may fail to capture meaningful dynamics. Future work will relax this assumption by modeling non-stationarity and complex drifts, aiming to develop more resilient phase representations and further establish PhaseFormer as a benchmark for long-term time-series forecasting.

# 7 ETHICS STATEMENT

This study focuses on methodological advances in time-series forecasting and does not involve human subjects, personally identifiable data, or sensitive private information. All experiments use publicly available benchmark datasets that are widely adopted in the research community, and their use complies with the terms of release. We do not employ proprietary or confidential data, and no conflicts of interest exist. The contributions are purely technical in nature and do not promote harmful applications. All authors affirm adherence to fairness, research integrity, and relevant legal and ethical standards, in line with the ICLR Code of Ethics.

## 8 REPRODUCIBILITY STATEMENT

We make substantial efforts to ensure reproducibility. All datasets used in our experiments are publicly accessible, with links provided in Sec. 5. Detailed dataset statistics, preprocessing steps, and partitioning procedures appear in Sec. A.2. Model architectures, hyperparameters, and training procedures (including optimizer choice, learning rate, look-back window length, and router configuration) are described in Sec. 5 and the Sec. A.2.

For fair comparison, we follow the official implementations of all baseline models and provide references to their sources. Comprehensive experimental results, including ablation studies, efficiency analyses, and visualizations, appear in Sec. 5 and Sec. A.3. Theoretical analyses supporting our design choices also appear in the Sec. A.6.

Finally, to facilitate independent verification, we release anonymized source code and experiment scripts in a public repository at https://anonymous.4open.science/r/ICLR26-PhaseFormer-17678. Collectively, these measures ensure that our reported results are reliably reproducible and extensible by the research community.

## REFERENCES

- Francisco Martinez Alvarez, Alicia Troncoso, Jose C. Riquelme, and Jesus S. Aguilar Ruiz. Energy time series forecasting based on pattern sequence similarity. *IEEE Transactions on Knowledge and Data Engineering (TKDE)*, 2010.
- Xu Cheng, Fan Shi, Xiufeng Liu, Meng Zhao, and Shengyong Chen. A novel deep class-imbalanced semi-supervised model for wind turbine blade icing detection. *IEEE Transactions on Neural Networks and Learning Systems (TNNLS)*, 2021.
- Razvan-Gabriel Cirstea, Tung Kieu, Chenjuan Guo, Bin Yang, and Sinno Jialin Pan. Enhancenet: Plugin neural networks for enhancing correlated time series forecasting. In *IEEE International Conference on Data Engineering (ICDE)*, 2021.
- Razvan-Gabriel Cirstea, Bin Yang, Chenjuan Guo, Tung Kieu, and Shirui Pan. Towards spatio-temporal aware traffic time series forecasting. In *IEEE International Conference on Data Engineering (ICDE)*, 2022.
- Jinliang Deng, Feiyang Ye, Du Yin, Xuan Song, Ivor Tsang, and Hui Xiong. Parsimony or capability? decomposition delivers both in long-term time series forecasting. *Advances in Neural Information Processing Systems (NeurIPS)*, 2024.
- Qihe Huang, Zhengyang Zhou, Kuo Yang, Zhongchao Yi, Xu Wang, and Yang Wang. Timebase: The power of minimalism in efficient long-term time series forecasting. In *International Conference on Machine Learning (ICML)*, 2025.
- Menghao Huo, Kuan Lu, Yuxiao Li, and Qiang Zhu. Ct-patchtst: Channel-time patch time-series transformer for long-term renewable energy forecasting. *arXiv preprint:2501.08620*, 2025.
- Andrew Jaegle, Felix Gimeno, Andy Brock, Oriol Vinyals, Andrew Zisserman, and Joao Carreira. Perceiver: General perception with iterative attention. In *International Conference on Machine Learning (ICML)*, 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 (ICLR)*, 2021.
  - Guokun Lai, Wei-Cheng Chang, Yiming Yang, and Hanxiao Liu. Modeling long- and short-term temporal patterns with deep neural networks. In *ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*, 2018.
  - Seunghan Lee, Taeyoung Park, and Kibok Lee. Learning to embed time series patches independently. In *International Conference on Learning Representations (ICLR)*, 2024.
  - Shiyang Li, Xiaoyong Jin, Yao Xuan, Xiyou Zhou, Wenhu Chen, Yu-Xiang Wang, and Xifeng Yan. Enhancing the locality and breaking the memory bottleneck of transformer on time series forecasting. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2019.
  - Shengsheng Lin, Weiwei Lin, Wentai Wu, Haojun Chen, and Junjie Yang. Sparsetsf: Modeling long-term time series forecasting with 1k parameters. In *International Conference on Machine Learning (ICML)*, 2024.
  - Hangchen Liu, Zheng Dong, Renhe Jiang, Jiewen Deng, Jinliang Deng, Quanjun Chen, and Xuan Song. Spatio-temporal adaptive embedding makes vanilla transformer sota for traffic forecasting. In ACM International Conference on Information and Knowledge Management (CIKM), 2023a.
  - Shizhan Liu, Hang Yu, Cong Liao, Jianguo 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 (ICLR)*, 2021.
  - Yong Liu, Haixu Wu, Jianmin Wang, and Mingsheng Long. Non-stationary transformers: Exploring the stationarity in time series forecasting. In *Advances in Neural Information Processing Systems* (*NeurIPS*), 2022.
  - Yong Liu, Tengge Hu, Haoran Zhang, Haixu Wu, Shiyu Wang, Lintao Ma, and Mingsheng Long. itransformer: Inverted transformers are effective for time series forecasting. In *International Conference on Learning Representations (ICLR)*, 2023b.
  - 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 (ICLR)*, 2023.
  - Liwen Ouyang and Aaron Key. Maximum mean discrepancy for generalization in the presence of distribution and missingness shift. In *NeurIPS Workshop on Distribution Shifts: Connecting Methods and Applications*, 2021.
  - Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. *arXiv preprint:1912.01703*, 2019.
  - Md. Mahin Uddin Qureshi, Amrin Binte Ahmed, Adisha Dulmini, Mohammad Mahboob Hussain Khan, and Rumana Rois. Developing a seasonal-adjusted machine-learning-based hybrid timeseries model to forecast heatwave warning. *Scientific Reports*, 2025.
  - Artyom Stitsyuk and Jaesik Choi. xpatch: Dual-stream time series forecasting with exponential seasonal-trend decomposition. In *AAAI Conference on Artificial Intelligence (AAAI)*, 2025.
  - Peiwang Tang and Weitai Zhang. Unlocking the power of patch: Patch-based mlp for long-term time series forecasting. In *AAAI Conference on Artificial Intelligence (AAAI)*, 2025.
  - Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of Machine Learning Research (JMLR)*, 2008.

- Shiyu Wang, Haixu Wu, Xiaoming Shi, Tengge Hu, Huakun Luo, Lintao Ma, James Y. Zhang, and Jun Zhou. Timemixer: Decomposable multiscale mixing for time series forecasting. In *International Conference on Learning Representations (ICLR)*, 2024.
  - Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Autoformer: Decomposition transformers with auto-correlation for long-term series forecasting. Advances in Neural Information Processing Systems (NeurIPS), 34:22419–22430, 2021a.
  - Xinle Wu, Dalin Zhang, Chenjuan Guo, Chaoyang He, Bin Yang, and Christian S. Jensen. Autocts: Automated correlated time series forecasting. *Proceedings of the VLDB Endowment (VLDB)*, 2021b.
  - Zhijian Xu, Ailing Zeng, and Qiang Xu. Fits: Modeling time series with 10k parameters. In *International Conference on Learning Representations (ICLR)*, 2024.
  - Kun Yi, Qi Zhang, Wei Fan, Shoujin Wang, Pengyang Wang, Hui He, Ning An, Defu Lian, Longbing Cao, and Zhendong Niu. Frequency-domain mlps are more effective learners in time series forecasting. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.
  - Kun Yi, Jingru Fei, Qi Zhang, Hui He, Shufeng Hao, Defu Lian, and Wei Fan. Filternet: Harnessing frequency filters for time series forecasting. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2024.
  - Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are transformers effective for time series forecasting? In AAAI Conference on Artificial Intelligence (AAAI), 2023.
  - Yunhao Zhang and Junchi Yan. Crossformer: Transformer utilizing cross-dimension dependency for multivariate time series forecasting. In *International Conference on Learning Representations* (ICLR), 2023.
  - 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 *AAAI Conference on Artificial Intelligence (AAAI)*, 2021.
  - Tian Zhou, Ziqing Ma, Qingsong Wen, Xue Wang, Liang Sun, and Rong Jin. Fedformer: Frequency enhanced decomposed transformer for long-term series forecasting. In *International Conference on Machine Learning (ICML)*, 2022.

# A APPENDIX

#### A.1 DETAILS ABOUT BASELINES

In our experiments, we incorporated a diverse set of time series forecasting models, with particular emphasis on approaches based on **Patch Tokenization** and efficient forecasting models. The details of these models are as follows:

- PatchTST A channel-independent Transformer that treats each variable as an individual channel and segments the time series into patches as tokens. This design reduces the complexity of the attention mechanism and enables the utilization of longer historical sequences, thereby improving long-term forecasting accuracy.
- 2. **iTransformer** A channel-dependent Transformer that models variables themselves as tokens to capture inter-variable relationships, while simultaneously accounting for nonlinear temporal variations within each variable.
- 3. Crossformer A multi-scale Transformer that performs patching or segmentation along the temporal dimension and employs a two-stage attention mechanism (within-time and cross-variable). This design effectively captures both temporal dependencies and intervariable correlations, making it particularly suitable for datasets characterized by strong inter-variable coupling and mixed long- and short-term patterns.
- 4. FEDformer A model that integrates trend-seasonal decomposition with frequency-domain analysis. It extracts a small number of significant frequency components to enhance periodic forecasting performance while maintaining controlled model complexity in long-term forecasting tasks.
- SparseTSF A lightweight model that reduces temporal complexity through periodic down-sampling or subsequence selection, aiming to achieve competitive periodic forecasting performance with minimal resource consumption.
- 6. FITS A lightweight model that leverages frequency-domain features and interpolation operations to reconstruct the predicted sequences. With fewer parameters and low computational overhead compared with other models, it demonstrates strong performance on time series with distinct spectral structures.
- 7. TimeBase A model that constructs temporal bases (via patching or segmentation strategies) to represent historical and future variations. Its objective is to maintain satisfactory forecasting accuracy while reducing computational and parameter costs.
- 8. **TimeMixer** An patch-based forecasting model fully based on MLPs. It employs Past-Decomposable-Mixing to decouple seasonal and trend components across different scales (fine and coarse) and utilizes Future-Multipredictor-Mixing to aggregate multi-scale predictions. This design achieves a balance of efficiency and accuracy in both short-term and long-term forecasting tasks.

# A.2 IMPLEMENTATION DETAILS

Dataset	Var	Length	T	L	Freq	Scale
ETTh1	7	14,400	720	96~720	1hour	0.1M
ETTh2	7	14,400	720	$96 \sim 720$	1hour	0.1M
ETTm1	7	57,600	720	$96 \sim 720$	15mins	0.4M
ETTm2	7	57,600	720	$96 \sim 720$	15mins	0.4M
Weather	21	52,696	720	$96 \sim 720$	10mins	1.1M
Electricity	321	26,304	720	$96 \sim 720$	1hour	8.1M
Traffic	862	17,544	720	$96 \sim 720$	1hour	15.0M

Table 3: Dataset statistics used in experiments.

We present detailed statistics of the datasets in Tab. 3. The data loading and preprocessing procedures follow prior works (Nie et al., 2023; Huang et al., 2025).

All baseline methods are implemented based on their original papers or official code. For cases where fixed random seeds are not specified, each experiment is repeated three times to ensure stability. All experiments are conducted using PyTorch (Paszke et al., 2019) on a single NVIDIA A100 24GB GPU.

For model configuration, the primary period is determined via frequency-domain analysis by selecting the dominant component, while the number of routers is chosen through grid search. We mainly use a single-layer model with 8 routers for ETT datasets, a two-layer model with 4 routers for Traffic and Electricity datasets, and a 3-layer model with 8 routers for the Weather dataset. Please refer to the released code for complete training details at https://anonymous.4open.science/r/ICLR26-PhaseFormer-17678.

#### A.3 FULL RESULTS

## A.3.1 THE DETAILED FORECASTING ACCURACY RESULTS

We present detailed forecasting results across all prediction horizons on the test sets in Tab. 5, with the input length fixed to 720. PhaseFormer consistently delivers strong and stable performance across most datasets and forecasting lengths. The only notable exception is ETTh2, where FITS slightly outperforms our model. This highlights the robustness of PhaseFormer across diverse scenarios, even though some simple datasets may still favor specialized baselines. It is also worth noting that TimeBase, which adopts a phase-based strategy, achieves competitive results on the relatively simple ETT datasets. In contrast, PhaseFormer demonstrates its advantage primarily on Traffic and Electricity, which are more complex and challenging datasets. This distinction illustrates that while phase-inspired models may be effective in straightforward settings, PhaseFormer generalizes better and excels in more demanding real-world contexts.

#### A.3.2 THE DETAILED FORECASTING EFFICIENCY RESULTS

We further provide the efficiency comparison of PhaseFormer against all baselines in terms of FLOPs and number of parameters, with the input length set to 720 and the output length fixed at 96. The results in Fig. 4 reveal that PhaseFormer achieves a favorable trade-off between accuracy and efficiency. Despite its stronger predictive performance, PhaseFormer maintains moderate model size and computational cost, often comparable to or even lower than other transformer-based models: On complex datasets such as Traffic, PhaseFormer outperforms large baselines like PatchTST with substantially fewer FLOPs; On simpler datasets, even when specialized models such as TimeBase or FITS show competitive accuracy, their efficiency advantage diminishes when considering scalability to larger, real-world datasets. These findings underscore that PhaseFormer is not only accurate but also efficient, making it more suitable for deployment in resource-constrained or latency-sensitive environments.

## A.3.3 THE DETAILED RESULTS OF PCA VISUALIZATION

We present PCA visualization results on the ETTh1, ETTh2, ETTm1, ETTm2, Electricity, and Weather datasets in Figure 7, in addition to Figure 2b. The findings are consistent with those observed on the Traffic dataset: phase tokenization yields a significantly more compact space compared to patch tokenization.

#### A.3.4 THE DETAILED RESULTS OF VARYING ROUTER NUMBERS

We further provide detailed results on the effect of varying the number of routers (1,2,4,8,16) across three datasets: Traffic, Electricity, and Weather. The input window was fixed at 720, and the output length was set to 96.

Our observations show that the number of routers does influence model performance, but the optimal configuration typically involves a relatively small number of routers. Specifically, the best performance was achieved with 8 routers on the Weather dataset, and with 4 routers on both the Electricity and Traffic datasets. Since routers serve as the foundation for aggregation and distribution in the phase token space, these results provide supporting evidence that the phase token space captures low-dimensional features, allowing strong performance even with fewer routers.

Table 4: Parameters and FLOPS across models for different datasets.

Model	Traffic		Weather		Electricity		ETTh1		ETTh2		ETTm1		ETTm2	
WIOGCI	Params	FLOPS	Params	FLOPS	Params	FLOPS	Params	FLOPS	Params	FLOPS	Params	FLOPS	Params	FLOPS
PhaseFormer	1.156K	13.9	308	0.15	1.156K	5.18	1.156K	0.11	1.156K	0.11	1.156K	0.11	1.156K	0.11
PatchTST	7.589M	498,577.49	1.373M	1,054.77	1.373M	16,122.93	587.68K	51.29	587.68K	51.29	587.68K	51.29	587.68K	51.29
iTransformer	6.731M	11,652.34	5.153M	257.54	5.153M	3,347.97	369.9K	8.12	304.1K	6.68	304.1K	7.29	304.1K	7.29
Crossformer	22.954M	259,209.90	158.34K	84.09	13.537M	96,564.63	2.069M	544.20	2.069M	544.20	2.069M	544.20	2.069M	544.20
FEDformer	21.206M	13,679.70	5.828M	2,757.24	11.861M	6,904.61	5.792M	2,734.51	5.792M	2,734.51	5.793M	2,734.95	5.793M	2,734.95
TimeBase	214	8.44	214	0.21	214	3.14	214	0.07	214	0.07	704	0.23	704	0.23
SparseTSF	17.949K	751.31	4.509K	5.14	4.509K	78.61	4.509K	1.71	4.509K	1.71	4.509K	1.71	4.509K	1.71
FITS	1.054K	1.76	272	0.01	462	0.28	272	0.004	272	0.004	2.646K	0.04	2.646K	0.04
TimeMixer	5.697M	2,026.53	5.562M	205.40	5.584M	739.64	4.024M	125.91	4.024M	125.91	4.024M	125.95	4.024M	125.95

Table 5: Full results across datasets and prediction lengths. Each entry reports MAE and MSE. The input length is set to 720. The best results are marked with **Bold**, and the second-best results are marked with **Underlined**.

Dataset	Horizon	Phasel	Former	Patch	TST	iTrans	former	Cross	former	FEDf	ormer	Time	Base	Spars	eTSF	FI	TS	Time	Mixer
Dataset	HOHZOH	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
	96	0.359	0.382	0.377	0.408	0.389	0.421	0.408	0.442	0.485	0.500	0.365	0.387	0.362	0.389	0.380	0.402	0.410	0.441
ETTh1	192	0.397	0.404	0.413	0.431	0.424	0.446	0.472	0.496	0.481	0.498	0.403	0.409	0.404	0.412	0.415	0.424	0.448	0.465
EIIII	336	0.425	0.424	0.436	0.444	0.456	0.469	0.480	0.486	0.522	0.521	0.409	0.419	0.435	0.426	0.449	0.460	0.475	0.490
	720	0.431	0.450	0.455	0.475	0.545	0.532	0.710	0.616	0.604	0.575	0.440	0.448	0.426	0.448	0.433	0.457	0.475	0.500
	96	0.275	0.338	0.276	0.339	0.305	0.361	1.164	0.744	0.401	0.451	0.292	0.350	0.294	0.346	0.271	0.336	0.315	0.380
ETTh2	192	0.341	0.376	0.342	0.385	0.405	0.421	1.414	0.830	0.425	0.464	0.339	0.387	0.340	0.377	0.332	0.374	0.383	0.415
E11112	336	0.369	0.405	0.364	0.405	0.411	0.436	1.220	0.794	0.427	0.471	0.394	0.420	0.360	0.398	0.355	0.396	0.415	0.436
	720	0.402	0.436	0.395	0.434	0.448	0.470	2.074	1.103	0.462	0.493	0.400	0.448	0.383	0.425	0.378	0.423	0.432	0.471
	96	0.293	0.344	0.298	0.352	0.315	0.369	0.306	0.353	0.406	0.441	0.311	0.351	0.314	0.359	0.313	0.357	0.332	0.384
ETTm1	192	0.323	0.361	0.335	0.373	0.349	0.388	0.341	0.385	0.450	0.477	0.338	0.371	0.348	0.376	0.339	0.369	0.362	0.398
EIIIII	336	0.358	0.381	0.366	0.389	0.381	0.409	0.383	0.420	0.436	0.466	0.364	0.386	0.368	0.386	0.367	0.385	0.386	0.413
	720	0.412	0.410	0.420	0.421	0.437	0.439	0.532	0.512	0.462	0.479	0.413	0.414	0.419	0.413	0.417	0.417	0.452	0.457
	96	0.163	0.256	0.165	0.260	0.179	0.274	0.244	0.338	0.339	0.406	0.162	0.256	0.167	0.259	0.166	0.256	0.192	0.285
ETTm2	192	0.219	0.293	0.219	0.298	0.239	0.314	0.350	0.412	0.397	0.452	0.218	0.293	0.219	0.297	0.271	0.328	0.307	0.362
ETTIMZ	336	0.269	0.326	0.268	0.333	0.309	0.356	0.400	0.431	0.418	0.452	0.270	0.328	0.271	0.330	0.352	0.380	0.380	0.412
	720	0.351	0.379	0.352	0.386	0.387	0.407	0.574	0.525	0.451	0.499	0.352	0.380	0.353	0.380	0.352	0.380	0.380	0.412
	96	0.148	0.195	0.149	0.199	0.159	0.212	0.151	0.210	0.289	0.342	0.146	0.198	0.174	0.231	0.176	0.232	0.163	0.223
Weather	192	0.193	0.237	0.193	0.243	0.203	0.252	0.220	0.273	0.340	0.403	0.185	0.241	0.216	0.267	0.203	0.256	0.201	0.255
weather	336	0.242	0.278	0.240	0.281	0.253	0.291	0.287	0.342	0.370	0.408	0.263	0.281	0.260	0.299	0.261	0.299	0.258	0.300
	720	0.309	0.332	0.312	0.334	0.317	0.337	0.362	0.393	0.420	0.421	0.314	0.331	0.325	0.345	0.325	0.346	0.329	0.348
	96	0.129	0.221	0.141	0.240	0.135	0.233	0.140	0.237	0.226	0.341	0.139	0.231	0.139	0.239	0.147	0.253	0.142	0.247
Electricity	192	0.148	0.238	0.156	0.256	0.155	0.253	0.165	0.259	0.220	0.336	0.153	0.245	0.155	0.250	0.159	0.256	0.164	0.273
Electricity	336	0.165	0.257	0.172	0.267	0.169	0.267	0.190	0.286	0.224	0.337	0.169	0.262	0.171	0.265	0.169	0.270	0.171	0.260
	720	0.201	0.285	0.208	0.299	0.204	0.301	0.227	0.312	0.271	0.378	0.207	0.294	0.208	0.300	0.214	0.302	0.209	0.313
	96	0.361	0.238	0.363	0.250	0.374	0.273	0.512	0.265	0.664	0.431	0.394	0.267	0.389	0.272	0.374	0.273	0.404	0.293
Traffic	192	0.373	0.243	0.382	0.258	0.393	0.283	0.528	0.271	0.613	0.382	0.407	0.270	0.399	0.272	0.393	0.282	0.404	0.292
Hanne	336	0.385	0.248	0.399	0.268	0.409	0.292	0.543	0.281	0.612	0.379	0.417	0.278	0.417	0.279	0.423	0.292	0.425	0.293
	720	0.428	0.270	0.432	0.289	0.450	0.314	0.598	0.314	0.664	0.410	0.456	0.298	0.449	0.299	0.450	0.314	0.453	0.314

Table 6: Impact of router number R on prediction accuracy. Each entry reports MSE and MAE. The input length is set to 720. The best results are marked with **Bold**, and the second-best results are marked with **Underlined**.

Dataset	R:	=1	R:	=2	R:	=4	R:	=8	R=16		
	MSE	MAE									
Weather	0.162	0.210	0.153	0.199	0.151	0.199	0.148	0.195	0.149	0.198	
Traffic	0.372	0.249	0.367	0.243	0.361	0.238	0.364	0.242	0.368	0.243	
Electricity	0.133	0.228	0.132	0.226	0.129	0.221	0.130	0.223	0.130	0.222	

## A.4 ADDITIONAL HYPER PARAMETERS ANALYSIS

## A.4.1 IMPACT OF MODEL PARAMETER SCALE ON PERFORMANCE

We conducted comparative experiments on three variants of the PhaseFormer model with different parameter scales across the Electricity, Traffic, and Weather datasets. The three configurations are: a single-layer model with latent dimension 8 ( $\approx$ 1.72K parameters), a single-layer model with latent dimension 16 ( $\approx$ 5.48K parameters), and a two-layer model with latent dimension 32 ( $\approx$ 37.1K parameters). Fig. 7 summarizes the results in terms of MSE, MAE and FLOPs.

Overall, the effect of model scale on performance is not consistent. On the Traffic dataset, larger models yield slight improvements, whereas on the Electricity and Weather datasets, the smaller and medium-sized models perform better. These findings indicate that PhaseFormer achieves a favorable balance between computational efficiency and predictive accuracy at relatively small parameter scales, and increasing model size does not lead to uniform gains across all tasks.

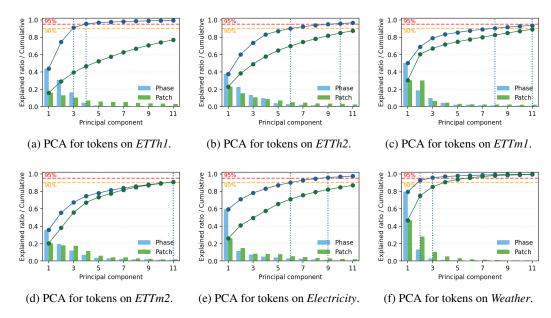


Figure 7: Visualization of phase tokenization across six datasets: ETTh1, ETTh2, ETTm1, ETTm2, Electricity, and Weather.

Table 7: Comparison of PhaseFormer variants. Each cell reports MSE, MAE, and FLOPs. The input length is fixed as 720 steps and the output length is fixed as 96 steps. The best results are marked with **Bold**, and the second-best results are marked with <u>Underlined</u>.

Dataset	Pha	seForme	r-1.7K	Pha	aseForm	er-5K	PhaseFormer-37K				
	MSE	MAE	FLOPs	MSE	MAE	FLOPs	MSE	MAE	FLOPs		
Electricity	0.129	0.220	9.41M	0.129	0.221	31.97M	0.131	0.223	221.05M		
Traffic	0.361	0.241	25.27M	0.366	0.243	85.84M	0.360	0.236	593.61M		
Weather	0.150	0.199	0.62M	0.151	0.194	2.09M	0.174	0.217	14.46M		

#### A.4.2 IMPACT OF INPUT LENGTH ON PERFORMANCE

We examine how the input window size affects the prediction accuracy and computational cost of PhaseFormer. Throughout this section, L denotes using the most recent L time steps as model input. The output length is fixed as 96 steps. As summarized in Fig. 8, increasing L reduces MSE and MAE across datasets, indicating that PhaseFormer benefits from longer historical context for modeling long-range temporal dependencies.

In terms of efficiency, the parameter count and FLOPs per forward pass remain nearly constant as L increases, with only modest growth (see Fig. 9) attributable primarily to the embedding stage. This behavior arises because the sequence length processed by the core encoder/decoder is governed by the number of *phases*, which depends on the data's learned periodic structure rather than by the raw input length. Consequently, scaling L mainly affects the embedding computations, whose cost is relatively small compared to the phase-based modules.

## A.5 SHOWCASES

To provide a clearer comparison of predictive performance across different models, we present the results of PhaseFormer, PatchTST, FITS, and TimeBase on the Traffic dataset. PhaseFormer demonstrates strong predictive performance, as reflected by both the shape of its forecasts and the actual prediction errors.

Compared with PhaseFormer, PatchTST produces lower peak values within each cycle, failing to fully match the true curve. This discrepancy is likely due to phase shifts in the periodic pattern that reduce peak amplitudes. FITS, which performs prediction in the frequency domain with frequency

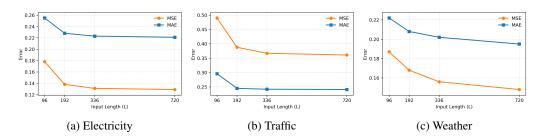


Figure 8: Prediction error test results across datasets under different input lengths.

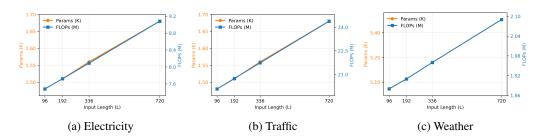


Figure 9: Efficiency evaluation of PhaseFormer across datasets with varying input lengths.

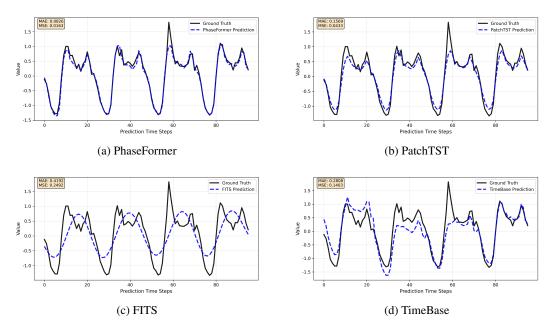


Figure 10: Visualization of forecasting results on Traffic dataset. The black lines stand for the ground truth and the blue lines stand for forecasting results.

band partitioning, tends to overlook high-frequency information. As a result, its predictions deviate from the ground truth, though the forecasts still preserve a periodic structure. TimeBase aligns well with the general cyclical pattern but fails to capture true variations across cycles, a limitation stemming from its patch-based basis construction mechanism.

# A.6 THEORETICAL ANALYSIS FOR PHASE TOKENIZATION

We model the data matrix as

$$X = AG^{\top} + N \in \mathbb{R}^{D \times H},$$

where  $A \in \mathbb{R}^{D \times r}$ ,  $G \in \mathbb{R}^{H \times r}$  are column full rank with  $\operatorname{rank}(A) = \operatorname{rank}(G) = r$ , and N is noise. The true signal is  $M = AG^{\top}$ . We assume  $r \ll \min(D, H)$ .

Patch tokenization corresponds to the row space Row(X) (the right singular r-subspace), while phase tokenization corresponds to the column space Col(X) (the left singular r-subspace).

A shared transformation applies  $S \in \mathbb{R}^{H \times H}$  on the hourly dimension:

$$X' = XS^{\top} = A(SG)^{\top} + N'.$$

For a matrix Y, define the spectral separation as

$$\operatorname{sep}_r(Y) := \min_{i \le r, j > r} |\sigma_i(Y) - \sigma_j(Y)|.$$

When rank(Y) = r, this equals  $\sigma_r(Y)$ . In particular,

$$\delta = \sigma_r(M), \quad \delta' = \sigma_r(MS^\top), \quad \delta_{\min} = \min(\delta, \delta').$$

Assume S is invertible on Col(G), i.e. rank(SG) = r. Define

$$\kappa := \sigma_{\min}(S|_{\operatorname{Col}(G)}) > 0.$$

Moreover, since  $\sigma_r(M) \geq \sigma_r(A)\sigma_r(G)$  and  $\sigma_r(MS^\top) \geq \kappa \sigma_r(A)\sigma_r(G)$ , we have the useful bound  $\delta_{\min} \geq \min(1, \kappa) \sigma_r(A)\sigma_r(G)$ .

For two r-dimensional subspaces  $\mathcal{U}$ ,  $\mathcal{V}$ , their distance is

$$d(\mathcal{U}, \mathcal{V}) := \|P_{\mathcal{U}} - P_{\mathcal{V}}\|_2 = \sin \Theta_{\max}(\mathcal{U}, \mathcal{V}),$$

where  $P_{\mathcal{U}}$  is the orthogonal projector onto  $\mathcal{U}$ . This metric satisfies the triangle inequality.

**Lemma 1 (Column space preservation)** Let  $M = AG^{\top}$ . If rank(SG) = r, then

$$\operatorname{Col}(MS^{\top}) = \operatorname{Col}(M) = \operatorname{Col}(A).$$

If rank(SG) < r, then the column space shrinks.

Lemma 2 (Row space change) For  $M = AG^{\top}$ ,

$$\operatorname{Row}(M) = \operatorname{Col}(G), \quad \operatorname{Row}(MS^{\top}) = \operatorname{Col}(SG).$$

As a result:

$$d(\operatorname{Col}(G), \operatorname{Col}(SG)) > 0 \iff S(\operatorname{Col}(G)) \neq \operatorname{Col}(G).$$

**Lemma 3 (Wedin's**  $\sin \Theta$  **theorem)** Let  $\widehat{M} = M + E$  and  $\delta = \sup_r(M) > 0$ . Then

$$d(\mathcal{U}_r(M), \mathcal{U}_r(\widehat{M})) \leq C \frac{\|E\|_2}{\delta},$$

where  $U_r(M)$  denotes the leading left singular r-subspace of M (the right case is analogous). Here C is an absolute constant (often  $C \in [2, 2\sqrt{2}]$ ). The condition  $\delta > 0$  is necessary.

**Theorem 2** Assume rank(SG) = r and  $\delta, \delta' > 0$ . Then:

1. For phase tokenization,

$$d\left(\mathcal{U}_r(X), \mathcal{U}_r(X')\right) \leq C\left(\frac{\|N\|_2}{\delta} + \frac{\|N'\|_2}{\delta'}\right) \leq C\frac{\|N\|_2 + \|N'\|_2}{\delta_{\min}}.$$

In the noiseless case, Lemma 1 ensures exact invariance, so the distance is 0.

2. For patch tokenization,

$$d\big(\mathcal{V}_r(X),\mathcal{V}_r(X')\big) \geq d_0 - C\Big(\frac{\|N\|_2}{\delta} + \frac{\|N'\|_2}{\delta'}\Big),$$
where  $d_0 := d(\operatorname{Col}(G),\operatorname{Col}(SG))$ . If  $S(\operatorname{Col}(G)) \neq \operatorname{Col}(G)$ , then  $d_0 > 0$ .

**Proof 1** For phase tokenization, apply the triangle inequality:

$$d(\mathcal{U}_r(X),\mathcal{U}_r(X')) \le d(\mathcal{U}_r(X),\mathcal{U}_r(M)) + d(\mathcal{U}_r(M),\mathcal{U}_r(MS^\top)) + d(\mathcal{U}_r(MS^\top),\mathcal{U}_r(X')).$$

 By Lemma 1 the middle term vanishes, and the two boundary terms are bounded by Wedin's theorem, yielding the stated inequality.

For patch tokenization, we have

the bound.

$$d(\mathcal{V}_r(X), \mathcal{V}_r(X')) \ge d(\mathcal{V}_r(M), \mathcal{V}_r(MS^\top)) - d(\mathcal{V}_r(M), \mathcal{V}_r(X)) - d(\mathcal{V}_r(MS^\top), \mathcal{V}_r(X')).$$
  
By Lemma 2 the first term equals  $d_0$ , and the other two are controlled by Wedin's theorem, proving

In real-world scenarios, slight variations in timing or conditions occur from day to day, so the daily transformations are not exactly identical. We model this systematic inconsistency by introducing a small perturbation  $\Delta_d$ , which captures the mismatch between the ideal linear transformation S and the actual data-generating process. Suppose each day's transform is  $S_d = S + \Delta_d$  with  $\|\Delta_d\|_2 \le \varepsilon$ . Then

$$X' = XS^{\top} + R, \qquad R_{d.:} = X_{d.:}\Delta_d^{\top}.$$

Bounding row by row gives  $||R_{d,:}||_2 \le \varepsilon ||X_{d,:}||_2$ , hence

$$||R||_F \le \varepsilon ||X||_F \quad \Rightarrow \quad ||R||_2 \le \varepsilon (||M||_F + ||N||_F).$$

**Theorem 3 (Stability under day-wise perturbations)** Under the relaxed model, each day uses  $S_d = S + \Delta_d$  with  $\|\Delta_d\|_2 \le \varepsilon$ , so that

$$X' = XS^{\top} + R, \qquad R_{d.:} = X_{d.:}\Delta_d^{\top}.$$

Let X = M + N with  $M = AG^{\top}$ ,  $\operatorname{rank}(A) = \operatorname{rank}(G) = r$ , and assume  $\operatorname{rank}(SG) = r$  so that  $\delta = \sigma_r(M) > 0$  and  $\delta' = \sigma_r(MS^{\top}) > 0$ . Define  $\delta_{\min} = \min(\delta, \delta')$  and  $d_0 = d(\operatorname{Col}(G), \operatorname{Col}(SG))$ . Then there exists an absolute constant  $C \in [2, 2\sqrt{2}]$  such that:

1. Phase tokenization (left r-subspace):

$$d(\mathcal{U}_{r}(X), \mathcal{U}_{r}(X')) \leq C \left( \frac{\|N\|_{2}}{\delta} + \frac{\|N'\|_{2}}{\delta'} + \frac{\|R\|_{2}}{\delta'} \right)$$

$$\leq C \frac{\varepsilon(\|M\|_{F} + \|N\|_{F}) + \|N\|_{2} + \|N'\|_{2}}{\delta_{\min}}$$

2. Patch tokenization (right r-subspace):

$$d(\mathcal{V}_{r}(X), \mathcal{V}_{r}(X')) \geq d_{0} - C\left(\frac{\|N\|_{2}}{\delta} + \frac{\|N'\|_{2}}{\delta'} + \frac{\|R\|_{2}}{\delta'}\right)$$

$$\geq d_{0} - C\frac{\varepsilon(\|M\|_{F} + \|N\|_{F}) + \|N\|_{2} + \|N'\|_{2}}{\delta_{\min}}$$

In particular, if  $S(\operatorname{Col}(G)) = \operatorname{Col}(G)$  then  $d_0 = 0$  and patch tokenization is also preserved up to the same perturbation scale.

Moreover, using  $\delta_{\min} \ge \min(1, \kappa) \, \sigma_r(A) \sigma_r(G)$  with  $\kappa = \sigma_{\min}(S|_{\operatorname{Col}(G)}) > 0$  makes the role of signal strength explicit.

**Proof 2** By row-wise control,  $||R_{d,:}||_2 \le \varepsilon ||X_{d,:}||_2$ , hence

$$||R||_F \le \varepsilon ||X||_F \le \varepsilon (||M||_F + ||N||_F), \qquad ||R||_2 \le ||R||_F \le \varepsilon (||M||_F + ||N||_F).$$

Insert the chain

$$X \to M \to MS^{\top} \to MS^{\top} + N' \to X' = MS^{\top} + N' + R.$$

For phase subspace  $U_r$ , according to the triangle inequality,

$$d(\mathcal{U}_r(X), \mathcal{U}_r(X')) \le d(\mathcal{U}_r(X), \mathcal{U}_r(M)) + d(\mathcal{U}_r(M), \mathcal{U}_r(MS^\top)) + d(\mathcal{U}_r(MS^\top), \mathcal{U}_r(MS^\top + N')) + d(\mathcal{U}_r(MS^\top + N'), \mathcal{U}_r(X')).$$

The middle term vanishes by Column space preservation (Lemma 1). Applying Wedin's  $\sin \Theta$  theorem (Lemma 3) to the remaining three perturbations  $E \in \{N, N', R\}$  yields  $C\|N\|_2/\delta + C\|N'\|_2/\delta' + C\|R\|_2/\delta'$ . Use  $\delta_{\min} \leq \delta, \delta'$  and Step 1 to obtain Item 1.

For patch subspace  $V_r$ , we use the reverse triangle inequality:

$$d(\mathcal{V}_r(X), \mathcal{V}_r(X')) \ge d(\mathcal{V}_r(M), \mathcal{V}_r(MS^\top)) - d(\mathcal{V}_r(M), \mathcal{V}_r(X)) - d(\mathcal{V}_r(MS^\top), \mathcal{V}_r(MS^\top + N')) - d(\mathcal{V}_r(MS^\top + N'), \mathcal{V}_r(X')).$$

The first term equals  $d_0$  by Row space change (Lemma 2). Apply Wedin's theorem to the other three terms to obtain the final results.

## A.7 THE USE OF LARGE LANGUAGE MODELS

In this work, large language models (specifically ChatGPT-5) are used solely for polishing the writing, identifying grammatical issues, and performing proofreading.