# Multi-Modal View Enhanced Large Vision Models for Long-Term Time Series Forecasting

ChengAo Shen<sup>1</sup>, Wenchao Yu<sup>2</sup>, Ziming Zhao<sup>1</sup>, Dongjin Song<sup>3</sup>, Wei Cheng<sup>2</sup>, Haifeng Chen<sup>2</sup>, Jingchao Ni<sup>1</sup>

## **Abstract**

Time series, typically represented as numerical sequences, can also be transformed into images and texts, offering multi-modal views (MMVs) of the same underlying signal. These MMVs can reveal complementary patterns and enable the use of powerful pre-trained large models, such as large vision models (LVMs), for long-term time series forecasting (LTSF). However, as we identified in this work, the state-of-the-art (SOTA) LVM-based forecaster poses an inductive bias towards "forecasting periods". To harness this bias, we propose DMMV, a novel decomposition-based multi-modal view framework that leverages trend-seasonal decomposition and a novel backcast-residual based adaptive decomposition to integrate MMVs for LTSF. Comparative evaluations against 14 SOTA models across diverse datasets show that DMMV outperforms single-view and existing multi-modal baselines, achieving the best mean squared error (MSE) on 6 out of 8 benchmark datasets. The code for this paper is available at: https://github.com/D2I-Group/dmmv.

## 1 Introduction

Long-term time series forecasting (LTSF) is vital across domains such as geoscience [1], neuroscience [3], energy [18], healthcare [28], and smart city [27]. Inspired by the success of Transformers and Large Language Models (LLMs) in the language domain, recent research has explored similar architectures for time series [40, 14, 51]. Meanwhile, Large Vision Models (LVMs) like ViT [7], BEiT [2] and MAE [12], have achieved comparable breakthroughs in the vision domain, prompting interest in their application to LTSF [4]. These approaches transform time series into image-like representations, enabling LVMs to extract embeddings for forecasting [29]. The rationale is that LVMs, pre-trained on large-scale image datasets, can transfer useful knowledge to LTSF due to a structural similarity: each image channel contains sequences of *continuous* pixel values analogous to univariate time series (UTS). This alignment suggests LVMs may be better suited to time series than LLMs, which process *discrete* tokens.

This hypothesis is partially validated by the SOTA VisionTS model [4], which applies MAE [12] to imaged time series and achieves impressive forecasting performance. This progress has spurred interests in combining imaged time series with other representations. In the past, time series have been studied through various forms: (1) raw numerical sequences [49, 30], (2) imaged representations [43, 4], and (3) verbalized (textual) descriptions [46, 10]. While they differ in modality, they represent alternative views of the same underlying data – unlike typical multi-modal data, where modalities originate from distinct sources [23]. However, these *multi-modal views (MMVs)* enable the application of large pre-trained models, such as LLMs, LVMs, and vision-language models (VLMs) [16, 33], to time series analysis, specializing them from those in conventional multi-view learning [39], where multi-view is a broader notion including both MMVs and views of the same modality (*e.g.*, augmented image views [39]). To distinguish, we use MMVs for time series throughout this paper.

Leveraging MMVs offers two key advantages: (1) augmenting time series with alternative views can reveal patterns not evident in the original numerical data, and (2) pre-trained large models can extract complex patterns specific to certain views, such as visual representations. Motivated by these benefits and the recent success of LVMs, this work investigates the synergy of MMVs for LTSF, with a focus on incorporating LVMs. To our knowledge, integrating the visual view of time series via LVMs alongside other modalities remains underexplored. The most related effort, Time-VLM [52], uses a VLM (ViLT [16]) to encode visual view and contextual texts of time series, augmented with a Transformer for the numerical view. All embeddings are combined through a fusion layer. However, this simple combination strategy overlooks the unique inductive biases of individual views, leading to suboptimal performance (see §4.1). Moreover, its use of textual inputs provides only marginal improvements while introducing significant computational overhead due to the large language encoder.

We propose DMMV, a Decomposition-based Multi-Modal View Framework for LTSF, which integrates numerical and visual views in a compact architecture. We exclude the textual view due to its marginal gains in Time-VLM [52] and recent doubts about the effectiveness and cost-efficiency of LLMs for LTSF [36]. DMMV comprises two specialized forecasters: a numerical forecaster and a visual forecaster. The visual forecaster, inspired by VisionTS [4], uses MAE [12] - a self-supervised LVM capable of reconstructing masked images - leveraging its strong performance on continuous values (i.e., pixels). Time series are transformed into images using a period-based patching technique [43], which, although effective, imposes an inductive bias on LVMs towards periodic signals. To address this, we design two DMMV variants as illustrated in Fig. 1: (a) DMMV-S (simple decomposition), which splits the time series into trend and seasonal components, assigning them to the numerical and visual forecasters, respectively; (b) DMMV-A (adaptive decomposition), which adaptively learns the decomposition via a backcast-residual mechanism aligned with the two forecasters. DMMV employs late fusion [17] via a gating mechanism, as intermediate fusion (e.g., embedding-level) underutilizes MAE's decoder, which plays a crucial role in pixel prediction. Extensive experiments show that DMMV significantly outperforms both SOTA single-view methods and Time-VLM, despite the latter incorporating an additional text encoder. To sum up, our contributions are as follows.

- We distinguish MMVs in time series analysis from the broader notion in conventional multi-view learning and study the emergent yet underexplored problem of MMV-based LTSF.
- We propose DMMV, a novel MMV framework that is carefully designed to harness an inductive bias we identified in SOTA LVM-based forecasters, complemented by the strength of a numerical forecaster, with two technical variants DMMV-S and DMMV-A.
- We conduct comprehensive experiments on benchmark datasets to evaluate DMMV, demonstrating
  its superior performance over 14 SOTA baselines and highlighting its potential as a new paradigm
  for MMV-based time series learning.

## 2 Related Work

To the best of our knowledge, this is the first work to explore LVMs in a decomposition-based MMV framework for LTSF. Our work relates to LVM-based time series forecasting (TSF), Multi-modal TSF, and Decomposition-based TSF, which are discussed below.

**LVM-based TSF.** Various vision models, such as ResNet [13], VGG-Net [34], and ViT [7], have been applied to TSF [50, 47], with some studies exploring image-pretrained CNNs like ResNet [13], Inception-v1 [35], and VGG-19 [34] for LTSF [19]. The use of LVMs in this area is still emerging, with most efforts focused on time series classification (e.g., AST [9] uses DeiT [37], ViTST [20] uses Swin [26]). In contrast, LVMs have seen limited use in TSF, likely due to their lower effectiveness on low-level (*i.e.*, numerical-level) tasks. The most notable method is VisionTS [4], which adapts MAE [12] for zero-shot and few-shot TSF. Another method, ViTime [47], trains ViT [7] from scratch on synthetic imaged time series but does not explore transferring knowledge from image-pretrained LVMs. Importantly, these approaches rely solely on vision models without incorporating other views or modalities.

Multi-modal TSF. Recently, large VLMs such as LLaVA [23], CLIP [33], and ViLT [16] have been explored for time series analysis [41, 32, 56, 52]. The most relevant is Time-VLM [52], which builds a forecaster on ViLT [16] to encode numerical and visual views, along with contextual texts. While integrating rich information with a large model, Time-VLM demonstrates promising results in TSF.

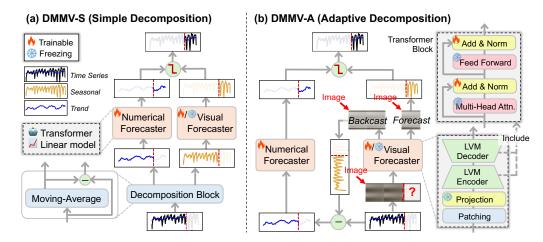


Figure 1: An overview of DMMV framework. (a) DMMV-S uses moving-average to extract trend and seasonal components. (b) DMMV-A uses a backcast-residual decomposition to automatically learn trend and seasonal components. In (b), the gray blocks are gray-scale images. "?" marks masks.

However, its fusion strategy closely follows the ViLT backbone and lacks time-series-specific design, leading to potentially suboptimal performance.

**Decomposition-based TSF.** Decomposition is a common technique in TSF, with seasonal-trend decomposition (STD) employed by models like Autoformer [44], FEDformer [54], and DLinear [49]. Recent work such as Leddam [48] replaces the traditional moving-average kernel in STD with a learnable one. Residual decomposition is another approach, used by N-BEATS [31] to reduce forecasting errors and later adopted by DEPTS [8] and CycleNet [21] for period-trend decomposition. While SparseTSF [22] predicts periods without explicit decomposition, SSCNN [6] introduces an attention-based method to extract long-term, short-term, seasonal, and spatial components. However, none of these methods incorporate LVMs. In contrast, our proposed DMMV shares insights with residual decomposition but is uniquely designed to exploit the inductive bias of LVMs for adaptive decomposition, setting it apart from prior works.

In summary, the proposed DMMV framework is distinct from existing approaches, yet integrates the strengths of pre-trained LVMs, the MMV framework, and decomposition techniques.

# 3 Decomposition-Based Multi-Modal View (DMMV) Framework

**Problem Statement**. Given a multivariate time series (MTS)  $\mathbf{X} = [\mathbf{x}^1,...,\mathbf{x}^D]^\top \in \mathbb{R}^{D \times T}$  within a look-back window of length T, where  $\mathbf{x}^i \in \mathbb{R}^T$   $(1 \leq i \leq D)$  is a UTS of the i-th variate, the goal of LTSF is to estimate the most likely values of the MTS at future H time steps, i.e.,  $\hat{\mathbf{Y}} \in \mathbb{R}^{D \times H}$ , such that the difference between the estimation and the ground truth  $\mathbf{Y} = \mathbf{X}_{T+1:T+H} \in \mathbb{R}^{D \times H}$  is minimized in terms of mean squared error (MSE), i.e.,  $\frac{1}{D:H} \sum_{i=1}^{D} \sum_{t=1}^{H} \|\hat{\mathbf{Y}}_{it} - \mathbf{Y}_{it}\|_2^2$ .

**Preliminaries.** Masked autoencoder (MAE) [12] is pre-trained self-supervisedly by reconstructing masked image patches using ImageNet dataset [5]. To adapt it to LTSF, VisionTS [4] adopts a period-based imaging technique introduced by TimesNet [43]. Specifically, each length-T UTS  $\mathbf{x}^i$  is segmented into  $\lfloor T/P \rfloor$  subsequences of length P, where P is set to be the period of  $\mathbf{x}^i$ , which can be obtained using Fast Fourier Transform (FFT) on  $\mathbf{x}^i$  [43] or from prior knowledge on sampling frequency. The

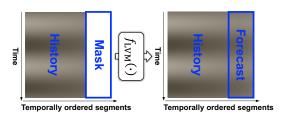


Figure 2: An illustration of an LVM forecaster

subsequences are stacked to form a 2D image  $\mathbf{I}^i \in \mathbb{R}^{P \times \lfloor T/P \rfloor}$ . After standard-deviation normalization,  $\mathbf{I}^i$  is duplicated 3 times to form a gray image of size  $P \times \lfloor T/P \rfloor \times 3$ , followed by a bilinear interpolation to resize it to an image  $\tilde{\mathbf{I}}^i$  of size  $224 \times 224 \times 3$  to fit the input requirement of MAE. As Fig.

2 shows, the forecast is achieved by reconstructing a right-appended masked area of  $\tilde{\mathbf{I}}^i$ , corresponding to the future horizon of  $\mathbf{x}^i$ . The forecast  $\hat{\mathbf{y}}^i \in \mathbb{R}^H$  can be recovered from the reconstructed area by de-normalization and reverse transformation. The forecast of MTS  $\mathbf{X}$  is achieved by forecasting over  $\mathbf{x}^1, ..., \mathbf{x}^D$  in parallel, following the channel-independence assumption [30].

An Inductive Bias. Due to period-based imaging and the spatial consistency enforced during MAE's pixel inference, VitionTS exhibits a strong bias toward *inter-period consistency*, often overshadowing the global trend. Fig. 3 illustrates VisionTS's forecasts on a synthetic sinusoidal time series with a period of 24. As shown in Fig. 3(a)-(d), where the segment length P varies from 24 to 48, forecasts alternate between accurate and inaccurate as P shifts from 1×period to 2×period, highlighting a strong inductive bias toward periodicity. Notably, the forecasts aren't mere repetitions – the decreasing intra-period amplitude indicates that LVMs can still capture local trends within each period. More quantitative results are deferred to Appendix D.2.

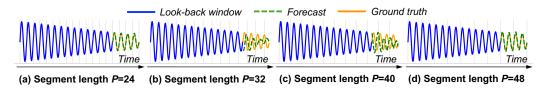


Figure 3: An illustration of LVM forecaster's inductive bias. The time series has a period of 24. The vertical dashed lines mark the segment points. The example indicates a bias towards segment lengths that are multiples of the period in (a)(d) over other segment lengths in (b)(c).

Motivated by this observation, we design the DMMV framework to leverage the inductive bias of LVMs while addressing their limitations. Specifically, the *visual forecaster*  $f_{vis}(\cdot)$  (*i.e.*, LVM) focuses on capturing periodic patterns from the *visual view*, while the *numerical forecaster*  $f_{num}(\cdot)$  models global trends from the *numerical view*, resulting in more balanced forecasting. Fig. 1 presents the two DMMV variants – DMMV-S and DMMV-A – within a decomposition-based architecture. Unlike prior approaches [44, 22, 6, 21], DMMV is explicitly designed to align with the inductive bias of LVMs.

## 3.1 DMMV with Simple Decomposition (DMMV-S)

DMMV-S adopts a simple moving-average (MOV) decomposition [44], which explicitly decomposes an input time series  $\mathbf{x}^i$  into a trend part and a seasonal (or periodic) part, reflecting the long-term progression and the seasonality of  $\mathbf{x}^i$ , respectively. Basically, MOV uses a kernel (*i.e.*, a sliding window) of length  $2\lfloor P/2\rfloor + 1$  to extract the component with frequency lower than  $\mathbf{x}^i$ 's sampling frequency (*i.e.*, 1/P), highlighting the global trend. The residual component is the seasonal part. This operation constitutes the decomposition block in Fig. 1(a).

$$\mathbf{x}_{\text{trend}}^i = \texttt{Moving-Average}(\texttt{Padding}(\mathbf{x}^i)), \quad \mathbf{x}_{\text{season}}^i = \mathbf{x}^i - \mathbf{x}_{\text{trend}}^i, \quad 1 \leq i \leq D \tag{1}$$

where  $Padding(\cdot)$  keeps the length of  $\mathbf{x}^i$  fixed.

The visual forecaster  $f_{\text{vis}}(\cdot)$  transforms the input  $\mathbf{x}_{\text{season}}^i$  into a  $224 \times 224 \times 3$  image  $\tilde{\mathbf{I}}_{\text{season}}^i$ , and outputs the forecast  $\hat{\mathbf{y}}_{\text{season}}^i \in \mathbb{R}^H$  for the seasonal component. For the numerical forecaster  $f_{\text{num}}(\cdot)$ , rather than imposing a specialized inductive bias, we adopt a general-purpose architecture capable of capturing long-term dependencies. We investigate the feasibility of two options and leave other explorations as a future work: (1) A simple linear model motivated by the proven effectiveness of linear methods in LTSF [49, 22, 21], i.e.,  $\hat{\mathbf{y}}_{\text{trend}}^i = f_{\text{num}}(\mathbf{x}_{\text{trend}}^i) = \mathbf{W}\mathbf{x}_{\text{trend}}^i + \mathbf{b}$ , where  $\mathbf{W} \in \mathbb{R}^{H \times T}$  and  $\mathbf{b} \in \mathbb{R}^H$  are weight and bias, respectively; and (2) A Transformer-based model inspired by PatchTST [30], which segments  $\mathbf{x}_{\text{trend}}^i$  into N length-L patches, where  $N = \lfloor T/L \rfloor + 1$ , to form the input  $\mathbf{X}_{\text{trend}}^i \in \mathbb{R}^{L \times N}$ , and performs

$$\mathbf{\tilde{X}}_{\text{trend}}^{i} = \mathbf{W}_{\text{pro}}\mathbf{X}_{\text{trend}}^{i} + \mathbf{W}_{\text{pos}} \rightarrow \mathbf{\hat{X}}_{\text{trend}}^{i} = \texttt{Transformer}(\mathbf{\tilde{X}}_{\text{trend}}^{i}) \rightarrow \mathbf{\hat{y}}_{\text{trend}}^{i} = \texttt{Linear}(\texttt{Flatten}(\mathbf{\hat{X}}_{\text{trend}}^{i})) \quad (2)$$

to achieve the forecast  $\hat{\mathbf{y}}_{\text{trend}}^i \in \mathbb{R}^H$  for the trend part, where  $\mathbf{W}_{\text{pro}} \in \mathbb{R}^{D' \times L}$  is the weight to project the patches to D'-dimensional embeddings,  $\mathbf{W}_{\text{pos}} \in \mathbb{R}^{D' \times N}$  is a learnable positional encoding,  $\text{Flatten}(\cdot)$  and  $\text{Linear}(\cdot)$  are flatten and linear operators.

Finally,  $\hat{\mathbf{y}}_{\text{season}}^i$  and  $\hat{\mathbf{y}}_{\text{trend}}^i$  are merged to produce the overall forecast  $\hat{\mathbf{y}}^i$  for the *i*-th variate. In particular, instead of using the regular summation-based merge, we design an adaptive merge function

with a light-weight gate  $g = \text{sigmoid}(w_g) \in [0,1]$ , where  $w_g$  is a learnable scalar parameter. To sum up, the overall process of DMMV-S is as follows.

$$\hat{\mathbf{y}}^i = g \circ \hat{\mathbf{y}}_{\text{season}}^i + (1 - g) \circ \hat{\mathbf{y}}_{\text{trend}}^i, \text{ where } \hat{\mathbf{y}}_{\text{season}}^i = f_{\text{vis}}(\tilde{\mathbf{I}}_{\text{season}}^i), \hat{\mathbf{y}}_{\text{trend}}^i = f_{\text{num}}(\mathbf{x}_{\text{trend}}^i)$$
 (3)

**Remark**. One limitation of DMMV-S is the explicit trend-seasonal decomposition placed on the input  $\mathbf{x}^i$ , which will enforce  $f_{\text{num}}(\cdot)$  and  $f_{\text{vis}}(\cdot)$  to fit **pre-defined components** extracted by a certain kernel size. This is not flexible and may not fully leverage LVMs' potential. To address it, we develop DMMV-A to have an adaptive decomposition in the next.

## 3.2 DMMV with Adaptive Decomposition (DMMV-A)

Unlike DMMV-S, DMMV-A *implicitly* decomposes the input  $\mathbf{x}^i$  into trend and seasonal components tailored to the strengths of the numerical and visual forecasters, respectively. This is achieved via a backcast-residual mechanism (Fig. 1(b)) that leverages LVMs' bias toward periodic patterns. The input  $\mathbf{x}^i$  is first transformed into an image  $\tilde{\mathbf{I}}^i$  using period-based imaging. Before forecasting,  $f_{\text{vis}}(\cdot)$ is used to backcast the look-back window by reconstructing masked segments of  $\tilde{\mathbf{I}}^i$ . An effective masking strategy must: (1) enable full-window reconstruction; (2) align with the forecasting setup (Fig. 2); and (3) minimize the usage of  $f_{vis}(\cdot)$  to avoid computational overhead. To meet these criteria, we propose an efficient *BackCast-Masking* (BCMASK) strategy (Fig. 4), which applies two passes: masking and reconstructing the left and right halves of  $\tilde{\mathbf{I}}^i$ , respectively.

$$\hat{\mathbf{I}}^{i} = [\hat{\mathbf{I}}_{\text{left}}^{i}, \hat{\mathbf{I}}_{\text{right}}^{i}], \quad \text{with} \quad \hat{\mathbf{I}}_{\text{left}}^{i} = f_{\text{vis}}(\tilde{\mathbf{I}}_{\text{right}}^{i}), \quad \hat{\mathbf{I}}_{\text{right}}^{i} = f_{\text{vis}}(\tilde{\mathbf{I}}_{\text{left}}^{i})$$
(4)

where  $\tilde{\mathbf{I}}_{\text{right}}^i$  ( $\tilde{\mathbf{I}}_{\text{left}}^i$ ) is the masked image with right (left) area unmasked,  $\hat{\mathbf{I}}_{\text{left}}^i$  ( $\hat{\mathbf{I}}_{\text{right}}^i$ ) is the reconstructed left (right) area,  $\hat{\mathbf{I}}^i$  is the reconstruction, or backcast, of  $\tilde{\mathbf{I}}^i$  by concatenating  $\hat{\mathbf{I}}^i_{\text{left}}$  and  $\hat{\mathbf{I}}^i_{\text{right}}$ .

BCMASK satisfies all three criteria: (1) it enables full reconstruction of  $\tilde{\mathbf{I}}^i$ ; (2) it uses contiguous segments in  $\tilde{\mathbf{I}}_{\text{right}}^i$  ( $\tilde{\mathbf{I}}_{\text{left}}^i$ ) to predict adjacent segments  $\hat{\mathbf{I}}_{\text{left}}^{i}$  ( $\hat{\mathbf{I}}_{\text{right}}^{i}$ ), mirroring the forecasting process in Fig. 2; and (3) it minimizes the use of  $f_{vis}(\cdot)$  – only two passes are needed, as some unmasked regions of  $\tilde{\mathbf{I}}^i$  are required for prediction and must later be masked to complete the full reconstruction.

Notably, the backcast in image  $\hat{\mathbf{I}}^i$  is biased toward the periodic patterns in  $\hat{\mathbf{I}}^i$ . After de-normalization and reverse transformation, a backcast time series

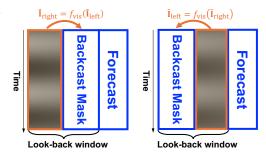


Figure 4: An illustration of BCMASK.

rigure 4. An inustration of Benjack.  $\hat{\mathbf{x}}^i \in \mathbb{R}^T$  is recovered, reflecting periodic component in  $\mathbf{x}^i$ . The residual  $\Delta \mathbf{x}^i = \mathbf{x}^i - \hat{\mathbf{x}}^i$  therefore emphasizes the trend. As shown in Fig. 1(b), we feed  $\Delta \mathbf{x}^i$  into  $f_{\text{num}}(\cdot)$  to produce  $\hat{\mathbf{y}}^i_{\text{trend}} \in \mathbb{R}^H$ , analogous to its role in DMMV-S. Meanwhile,  $f_{\text{vis}}(\cdot)$ predicts from  $\tilde{\mathbf{I}}^i$ , likely yielding the forecast of seasonal component  $\hat{\mathbf{y}}_{\text{season}}^i \in \mathbb{R}^H$ . Finally,  $\hat{\mathbf{y}}_{\text{trend}}^i$ and  $\hat{\mathbf{y}}_{\text{season}}^i$  are fused via the same gating mechanism as Eq. (3). In summary, this defines the overall process of DMMV-A as follows.

$$\hat{\mathbf{y}}^{i} = g \circ \hat{\mathbf{y}}_{\text{season}}^{i} + (1 - g) \circ \hat{\mathbf{y}}_{\text{trend}}^{i}, \quad \text{where} \quad \hat{\mathbf{y}}_{\text{season}}^{i} = f_{\text{vis}}(\tilde{\mathbf{I}}^{i}), \quad \hat{\mathbf{y}}_{\text{trend}}^{i} = f_{\text{num}}(\Delta \mathbf{x}^{i})$$
 (5)

**Remark**. Unlike DMMV-S, DMMV-A automatically learns a decomposition of  $\mathbf{x}^i$  that optimally aligns  $f_{\text{num}}(\cdot)$  and  $f_{\text{vis}}(\cdot)$  with the forecasting task. As shown in §4.3, this adaptive decomposition effectively separates seasonal and trend components, leveraging the inductive bias of  $f_{\text{vis}}(\cdot)$ . Unlike the backcast in N-BEATS [31] - designed merely to extract predictive errors - our approach is specifically tailored to exploit LVMs' bias toward periodic patterns, making it fundamentally different.

## 3.3 Model Optimization

After obtaining  $\hat{\mathbf{Y}} = [\hat{\mathbf{y}}^1,...,\hat{\mathbf{y}}^D]^{\top} \in \mathbb{R}^{D \times H}$ , DMMV is trained by minimizing the MSE between  $\hat{\mathbf{Y}}$  and  $\mathbf{Y}, i.e., \frac{1}{D \cdot H} \sum_{i=1}^D \sum_{t=1}^H \|\hat{\mathbf{Y}}_{it} - \mathbf{Y}_{it}\|_2^2$ . As shown in Fig. 1,  $f_{\text{num}}(\cdot)$  is trained from scratch, while  $f_{\text{vis}}(\cdot)$  uses pre-trained LVM weights with partial fine-tuning. We find that fine-tuning only

the normalization layers yields the best performance, consistent with the findings in [55]. For the choice of LVM, we tested MAE [12] and SimMIM [45], both are self-supervisedly pre-trained LVMs, in §4.2. MAE performs better and is set as the default. Training begins with  $f_{\text{vis}}(\cdot)$  frozen while  $f_{\text{num}}(\cdot)$  is trained for a number of epochs (e.g.,  $\sim$ 30). Then, the norm layers of  $f_{\text{vis}}(\cdot)$  are unfrozen and fine-tuned jointly with  $f_{\text{num}}(\cdot)$  until convergence or early stopping.

# 4 Experiments

In this section, we compare DMMV with the SOTA methods on LTSF benchmark datasets, and analyze its effectiveness with both quantitative and qualitative studies.

**Datasets**. We adopt 8 widely used MTS benchmarks: ETT (Electricity Transformer Temperature) [53], including ETTh1, ETTh2, ETTm1, ETTm2; Weather [44], Illness [44], Traffic [44], and Electricity [38]. Following standard protocols [44], we split the datasets chronologically into training/validation/test sets using a 60%/20%/20% ratio for ETT and 70%/10%/20% for the others. The prediction horizon H is set to {24, 36, 48, 60} for Illness, and {96, 192, 336, 720} for the remaining datasets. By default, look-back window T is 336. Full dataset details are provided in Appendix B.1.

The Compared Methods. We compare DMMV with the SOTA methods, including VLM-based multimodal model: (1) Time-VLM [52]; LVM-based model: (2) VisionTS [4]; LLM-based models: (3) Time-LLM [15], (4) GPT4TS [55], (5) CALF [24]; Transformer-based models: (6) PatchTST [30], (7) FEDformer [54], (8) Autoformer [44], (9) Stationary [25], (10) ETSformer [42], (11) Informer [53]; and non-Transformer models: (12) DLinear [49], (13) TimesNet [43], (14) CycleNet [21]. As we use the standard evaluation protocol, we collect results from prior works: Time-VLM ([52]), LLM-based models (reproduced by [36]), Transformer-based models, PatchTST and TimesNet ([4]), and CycleNet ([21]). Since VisionTS in [4] originally uses dynamic look-back windows such as T=1152 and T=2304 for different datasets, we re-run it with T=336. CycleNet's results on the Illness dataset is unavailable in [21]. Thus we run its official code on the Illness dataset.

We evaluate both DMMV variants – DMMV-s and DMMV-A. A linear forecaster (§3.1) and MAE are set as the default in  $f_{\text{num}}(\cdot)$  and  $f_{\text{vis}}(\cdot)$ , respectively. Ablation studies (§4.2) include variants with a Transformer-based  $f_{\text{num}}(\cdot)$  and SimMIM as  $f_{\text{vis}}(\cdot)$ . Following [4], the imaging period P (§3) for both VisionTS and DMMV is set based on each dataset's sampling frequency (see Appendix B.1). Additional details on all compared methods are in Appendix B.2.

**Evaluation**. Following [30, 49, 36], we use Mean Squared Error (MSE) and Mean Absolute Error (MAE) to evaluate the LTSF performance of the compared methods.

## 4.1 Experimental Results

Table 1 summarizes the LTSF performance of 10 representative methods across four categories: MMV-based, visual-view-based, language-view-based, and numerical-view-based approaches, with full results for all 16 methods provided in Appendix D.1. Time-VLM's results on the Illness dataset are not reported in [52] and its code is unavailable at the time of this experiment, thus are marked by "-". For DMMV, the stronger variant, DMMV-A, is reported. In Table 1, several key insights emerge: (1) MMV and visual-view methods generally outperform language-view methods, underscoring the effectiveness of LVMs, particularly when integrated within MMV frameworks; (2) Numerical-view models such as PatchTST and CycleNet remain competitive, especially on datasets where VisionTS underperforms (e.g., ETTm2 and Electricity), highlighting their potential to complement visual models; (3) The strong results of CycleNet, a lightweight model with learnable decomposition, demonstrate the value of combining simplicity with structure in LTSF; (4) Notably, DMMV-A, which unifies visual and numerical views through a novel adaptive decomposition, outperforms the baselines in most cases, achieving 43 first-places and confirming its effectiveness; (5) Lastly, while VisionTS performs well on highly periodic datasets (e.g., ETTh1, ETTm1, Traffic) due to MAE's inductive bias toward periodicity, DMMV-A alleviates this bias, resulting in more generalizable forecasts.

Fig. 5 presents critical difference (CD) diagrams [11] showing the average rank of all 16 methods based on MSE and MAE across prediction lengths and all datasets. DMMV-s ranks 4.5/16 in MSE and 7.1/16 in MAE, underscoring the benefit of the adaptive decomposition used in DMMV-A (Fig. 1(b)). From Fig. 5, DMMV-s's comparable ranks to CycleNet indicate even with a simpler, fixed decomposition, DMMV-s exhibits an ability that a strong SOTA model with learnable decomposition has. In §4.3, a detailed comparison between DMMV-s and DMMV-A is provided.

Table 1: LTSF performance comparison on the benchmark datasets. Lower MSE and MAE indicate better performance. **Red** values indicate the best MSE and MAE per row. Time-VLM's results on the Illness dataset are unavailable in [52] and its code was unavailable at the time of this experiment.

Vie	w		Multi-	Modal		Vis	ual		Lang	uage						Nume	rical				
Мо	del	DMM	AV-A	Time	-VLM	Visio	nTS	GPT	'4TS	Time	LLM	Patch	nTST	Cycl	eNet	Time	esNet	DLi	near	FEDf	ormer
Me	tric	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
	96	0.354	0.389	0.361	0.386	0.355	0.386	0.370	0.389	0.376		0.370	0.399	0.374	0.396		0.402				
ľh1	192	0.393	0.405	0.397	0.415	0.395	0.407	0.412	0.413	0.407		0.413	0.421	0.406	0.415		0.429		0.416		0.448
ETTh1	336 720	0.387 0.445	<b>0.413</b> 0.450	0.420 0.441	0.421 0.458	0.419 0.458	0.421	0.448 <b>0.441</b>	0.431 <b>0.449</b>	0.430		0.422	0.436	0.431	0.430	l	0.469	1	0.416	0.459	
	Avg.	0.395	0.414	0.405	0.420	0.407	0.419	0.418	0.421	0.437		0.413	0.431	0.430	0.426		0.450		0.430		0.460
	96	0.294	0.349	0.267	0.335	0.288	0.334	0.280	0.335	0.286	0.346	0.274	0.336	0.279	0.341	0.340	0.374	0.289	0.353	0.358	0.397
ETTh2	192	0.339	0.395	0.326	0.373	0.349	0.380	0.348	0.380	0.361		0.339	0.379	0.342	0.385		0.414		0.418		0.439
Ħ	336	0.322	0.384	0.357	0.406	0.364	0.398	0.380	0.405	0.390		0.329	0.380	0.371	0.413		0.452	0.448		1	0.487
	720	0.392	0.425	0.412	0.449	0.403	0.431	0.406	0.436	0.405		0.379	0.422	0.426	0.451		0.468	0.605			0.474
	Avg. 96	0.337	0.388	0.341	0.391	0.351	0.386	0.354	0.389	0.361	0.396	0.330	0.379	0.355	0.398	0.414	0.427	0.431	0.447	0.437	0.449
=	192	0.279	0.357	0.332	0.346	0.284	0.362	0.343	0.340	0.291		0.230	0.342	0.233	0.348		0.373	0.235		0.426	
ETTm1	336	0.351	0.381	0.364	0.383	0.354	0.382	0.376	0.386	0.359		0.366	0.392	0.368	0.386		0.411		0.386		0.459
H	720	0.411	0.415	0.402	0.410	0.411	0.415	0.431	0.416	0.433		0.416	0.420	0.417	0.414		0.450	0.425			0.490
	Avg.	0.340	0.371	0.351	0.376	0.344	0.373	0.363	0.378	0.356	0.377	0.351	0.381	0.355	0.379	0.400	0.406	0.357	0.379	0.448	0.452
	96	0.172	0.260	0.160	0.250	0.174	0.262	0.163	0.249		0.248	0.165	0.255	0.159	0.247	0.187	0.267	0.167	0.260	0.203	0.287
m2	192	0.227	0.298	0.215	0.291	0.228	0.297	0.222	0.291	0.235		0.220	0.292	0.214	0.286		0.309	0.224			0.328
ETTm2	336	0.272	0.327	0.270	0.325	0.281	0.337	0.273	0.327	0.280		0.274	0.329	0.269	0.322		0.351		0.342		0.366
	720 Avg.	0.351	0.381	0.348 0.248	0.378	0.384	0.410 0.327	0.357	<b>0.376</b> 0.311	0.366		0.362 0.255	0.385	0.363	0.382 0.309	0.408	0.403	0.397	0.421		0.415
	24	1.409	0.317	0.240	0.311	1.613	0.327	1.869	0.823	1.792	0.807	1.319	0.313	2.255	1.017	2.317	0.934	2.215	1.081	3.228	1.260
S	36	1.290	0.745	_	_	1.316	0.750	1.853	0.854	1.833		1.430	0.834	2.121	0.950		0.920		0.963		1.080
Illness	48	1.499	0.810	l –	_	1.548	0.818	1.886	0.855	!	1.012	1.553	0.815	2.187	1.007	2.238	0.940	2.130	1.024	2.622	1.078
$\equiv$	60	1.428	0.773	-	_	1.450	0.783	1.877	0.877	2.177	0.925	1.470	0.788	2.185	0.997	2.027	0.928	2.368	1.096	2.857	1.157
	Avg.	1.407	0.771	_	-	1.482	0.796	1.871	0.852	2.018		1.443	0.798	2.187	0.992	2.139	0.931	2.169	1.041	2.847	
≥	96	0.126	0.213	0.142	0.245	0.127	0.217	0.141	0.239		0.233	0.129	0.222	0.128	0.223	0.168	0.272	0.140	0.237	0.193	0.308
:5:	192	0.145	0.237	0.157	0.260	0.148	0.237	0.158	0.253	0.152		0.157	0.240	0.144	0.237		0.289		0.249		0.315
Electricity	336 720	0.162 0.197	0.254 0.286	0.174	0.276 0.308	0.163	<b>0.253</b> 0.293	0.172 0.207	0.266	0.169		0.163 <b>0.197</b>	0.259	<b>0.160</b> 0.198	0.254 0.287		0.300	0.169			0.329
	Avg.	0.157	0.248	0.214	0.308	0.159	0.250	0.207	0.263	0.200		0.162	0.253		0.250		0.320		0.361		0.333
	96	0.143	0.195	0.148	0.200	0.146	0.191	0.148	0.188		0.199	0.149	0.198	0.167	0.221	0.172		0.176	0.237	0.217	0.296
ЭC	192	0.187	0.242	0.193	0.240	0.194	0.238	0.192	0.230	0.223	0.261	0.194	0.241	0.212	0.258	l	0.261	0.220	0.282	0.276	0.336
Weather	336	0.237	0.273	0.243	0.281	0.243	0.275	0.246	0.273	0.251	0.279	0.245	0.282	0.260			0.306	0.265	0.319	0.339	0.380
⋛	720	0.302	0.315	0.312	0.332	0.318	0.328	0.320	0.328	0.345		0.314	0.334		0.339		0.359	0.333			0.428
	Avg.	0.217	0.256	0.224	0.263	0.225	0.258	0.227	0.255	0.244		0.226	0.264	0.242	0.278		0.287		0.300		0.360
	96	0.344	0.237	0.393	0.290	0.346	0.232	0.396	0.264	0.392		0.360	0.249	0.397	0.278	l	0.321	0.410	0.282	0.587	0.366
Traffic	192 336	0.363 0.387	0.249 0.256	0.405	0.296 0.305	0.376	0.245 0.252	0.412	0.268 0.273	0.409		0.379 0.392	0.256 0.264	0.411	0.283		0.336	0.423			0.373
Ë	720	0.433	0.230	0.420	0.303	0.389		0.421	0.273	0.451		0.392	0.286		0.205		0.350	0.466			0.383
	Avg.	0.382	0.257	0.419	0.304	0.386	0.256	0.421	0.274	0.422		0.391	0.264	0.421	0.289	l	0.336		0.295		0.376
# W		4		9		9	)		7		1	9		1		(		_	)	_	0
	(a) MSE Ranking (b) MAE Ranking																				
		1	2 3 4	5 6	7 8	9 10	11 12 1	3 14 1	5 16 <b>L</b>			1	2 3	4 5	6 7 8	9 1	0 11 1	2 13 1	4 15 16	6	
1	DMM' Patch' ime-V Visior DMM' Cyclel	TST —— TLM —— TTS —— V-S ——						L	Aut Sta ETS FEI	ormer toforme tionary former Oformer	r Visi GP Pato Time	MV-A — onTS — T4TS — hTST — -VLM — eNet —								- Inform - Autofo - FEDfo - ETSfo - Statio - Times	ormer rmer rmer nary

Figure 5: Critical difference (CD) diagram on the average rank of all 16 compared methods in terms of (a) MSE and (b) MAE over all benchmark datasets. The lower rank (left of the scale) is better.

## 4.2 Ablation Analysis

We validate the design of DMMV-A through ablation studies on four datasets; DMMV-S results are deferred to Appendix D.3 for brevity. Table 2 summarizes the analysis: (a) replaces the linear model in  $f_{\text{num}}(\cdot)$  with a PatchTST-style Transformer (see §3.1); (b) swaps MAE with SimMIM [45] as  $f_{\text{vis}}(\cdot)$ ; (c) replaces the gating fusion with a simple sum; (d) removes BCMASK, performing backcasting and forecasting on the full, unmasked look-back window; (e) substitutes BCMASK with random masking; (f) freezes the entire  $f_{\text{vis}}(\cdot)$  instead of fine-tuning norm layers; and (g) removes the backcast-residual mechanism, feeding both  $f_{\text{num}}(\cdot)$  and  $f_{\text{vis}}(\cdot)$  the same input  $\mathbf{x}^i$  and merging their outputs via gating.

Table 2 reveals key insights into DMMV-A's design. In (a), replacing the linear numerical forecaster with a Transformer slightly degrades performance, likely due to the increased difficulty of jointly training Transformer with LVMs. In (b), MAE outperforms SimMIM as  $f_{\rm vis}(\cdot)$ , likely due to its ViT-based reconstruction decoder being better suited for pixel-level tasks like LTSF than SimMIM's linear decoder, while both models share similar encoder architectures. In (c), gate-based fusion

Table 2: Ablation analysis of DMMV-A. MSE and MAE are averaged over different prediction lengths. Lower MSE and MAE are better. "Improvement" of each ablation is relative to DMMV-A.

Dataset $(\rightarrow)$	ET	Γh1	ETT	Γm1	Illr	iess	Wea	ther
Method $(\downarrow)$ , Metric $(\rightarrow)$	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
DMMV-A	0.395	0.414	0.340	0.371	1.407	0.771	0.217	0.256
(a) $f_{ ext{num}}(\cdot)  o  ext{Transformer}$	0.407	0.421	0.339	0.372	1.442	0.786	0.219	0.260
Improvement	-3.04%	-1.69%	+0.29%	-0.27%	-2.49%	-1.95%	-0.92%	-1.56%
(b) $f_{ ext{vis}}(\cdot)  o  ext{SimMIM}$	0.407	0.415	0.345	0.377	1.649	0.814	0.227	0.261
Improvement	-3.04%	-0.24%	-1.47%	-1.62%	-17.20%	-5.58%	-4.61%	-1.95%
(c) Gate $\rightarrow$ Sum	0.414	0.427	0.352	0.383	1.606	0.863	0.233	0.278
Improvement	-4.81%	-3.14%	-3.53%	-3.23%	-14.14%	-11.93%	-7.37%	-8.59%
(d) BCMASK→ No mask	0.426	0.441	0.349	0.377	1.493	0.828	0.221	0.267
Improvement	-7.85%	-6.52%	-2.65%	-1.62%	-6.11%	-7.39%	-1.84%	-4.30%
(e) BCMASK→ Random mask	0.394	0.414	0.340	0.372	1.472	0.829	0.223	0.262
Improvement	0.25%	0.00%	0.00%	-0.27%	-4.62%	-7.52%	-2.76%	-2.34%
(f) Freeze $f_{\text{vis}}(\cdot)$	0.431	0.428	0.358	0.380	1.442	0.773	0.246	0.288
Improvement	-9.11%	-3.38%	-5.29%	-2.43%	-2.49%	-0.26%	-13.36%	-12.50%
(g) W/o decomposition	0.408	0.424	0.338	0.373	1.712	0.903	0.219	0.268
Improvement	-3.29%	-2.42%	0.59%	-0.54%	-21.68%	-17.12%	-0.92%	-4.69%
0.84 0.83		0.84	0	.83	0.89		0.89	
						0.65	0.0	0.80
0.54 0.51 0.49	0.53		0.61	0.53		0.05	0.1	,,
0.46	0.47		0.39	0	.47	0.35	0.35	0
						0.33	0.33	0.20
6 0.16	0.17	0.16		0.17	0.1	1	0.11	0.20

Figure 6: Comparing DMMV-S and DMMV-A w.r.t. gate weights on visual and numerical forecasters.

weather

DMMV-S Visual

ETTm2

DMMV-A Numerical

electricity

DMMV-S Numerica

illness

traffic

outperforms simple summation, highlighting its adaptability to the distinct outputs of  $f_{\text{num}}(\cdot)$  and  $f_{\text{vis}}(\cdot)$ . (d) and (e) underscore the importance of BCMASK: removing it (i.e., (d) "No mask") recovers the full look-back window as the backcasted seasonal component, diminishing the trend signal and weakening  $f_{\text{num}}(\cdot)$ , while "Random mask" (i.e., (e)) performs slightly worse due to poorer periodic pattern extraction, which leads to many fluctuations. In §4.3, we provide visual examples to compare these masking strategies. In (f), fine-tuning only the norm layers significantly improves performance over freezing, confirming the benefit of coordinated learning between forecasters, as described in §3.3. Finally, (g) shows that removing the backcast-residual mechanism causes a major performance drop, affirming its role in effective decomposition. Overall, the LVM decoder, fusion strategy, masking method, training approach, and decomposition mechanism are crucial to DMMV-A's success.

## 4.3 Performance Analysis

ETTh1

ETTm1

DMMV-A Visual

ETTh2

In this section, we perform an in-depth analysis of DMMV using the same four datasets as in §4.2.

The Difference between DMMV-s and DMMV-A. Fig. 5 highlights DMMV-A's superiority over DMMV-s, largely due to its adaptive decomposition mechanism. A key advantage of the gate-based fusion is its interpretability. As shown in Fig. 6, which presents average gate weights across datasets, DMMV-A consistently places more weight on  $f_{\text{vis}}(\cdot)$ , while DMMV-S tends to balance both  $f_{\text{num}}(\cdot)$  and  $f_{\text{vis}}(\cdot)$  but leans toward  $f_{\text{num}}(\cdot)$ . In DMMV-A, these weights are learned based on forecasting performance, emphasizing  $f_{\text{vis}}(\cdot)$ 's importance. Notably, although  $f_{\text{num}}(\cdot)$  receives less weight, it remains essential — as evidenced by DMMV-A outperforming the visual-only baseline VisionTS in Table 1. In contrast, DMMV-S's weights are limited by its fixed moving-average decomposition, leading to a non-adaptive and suboptimal allocation of forecasting roles.

Fig. 7 provides example decompositions by DMMV-S and DMMV-A (additional cases in Appendix D.4). DMMV-A produces a smooth, clearly periodic component – consistent with expectations, and a trend component with some noises. In contrast, DMMV-S's moving-average yields a smoother trend by absorbing fluctuations, pushing noise into the seasonal component. This makes forecasting harder for  $f_{\text{vis}}(\cdot)$ , which is more sensitive to fluctuations than  $f_{\text{num}}(\cdot)$ , resulting in lower weights of  $f_{\text{vis}}(\cdot)$  in Fig. 6. Since periodic patterns are crucial for long-term forecasting, as identified by [22, 21], the clearer period separation in DMMV-A leads to forecasts that better match the ground truth.

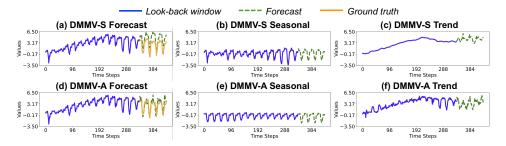


Figure 7: The decompositions of DMMV-S and DMMV-A on the same example in ETTh1: (a)(d) input time series and forecasts, (b)(e) seasonal component, and (c)(f) trend component.

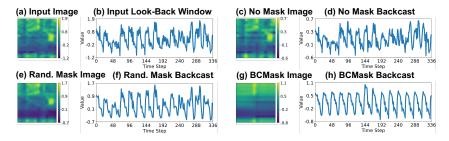


Figure 8: Comparison of different masking methods on the same example in ETTh1. (a) image of input look-back window; (c)(e)(g) are images of backcast output by DMMV-A: (c) uses "No mask"; (e) uses "Random mask"; (g) uses BCMASK. (b)(d)(f)(h) are their recovered time series, respectively.

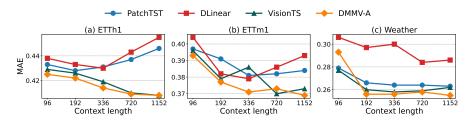


Figure 9: Average MAE comparison with varying look-back window (or context) lengths.

**The Effectiveness of BCMASK.** Fig. 8 compares the backcast results of DMMV-A using BCMASK, "No mask", and "Random mask" (as in Table 2) on a sample case; more examples are in Appendix D.5. BCMASK produces a smooth image along the temporal (*x*-axis) segments, effectively capturing clean periodic patterns. In contrast, "No mask" closely replicates the input, offering no meaningful decomposition. "Random mask" performs moderately well, resembling BCMASK but with less temporal smoothness, indicating a less optimal decomposition.

Impact of Look-Back Window. Fig. 9 compares DMMV-A with a visual forecaster (VisionTS) and two numerical forecasters (PatchTST, DLinear), which can serve as its single-view ablations. Illness dataset is excluded due to its short time series (966 time steps). Using MAE metric (MSE results in Appendix D.6), we observe that DMMV-A and VisionTS benefit from longer look-back windows, while PatchTST and DLinear degrade beyond a length of 336. Notably, DMMV-A outperforms VisionTS at length 1152, highlighting the advantage of explicitly modeling global trends.

## 5 Conclusion

This paper introduces DMMV, a novel MMV framework that leverages LVMs and adaptive decomposition to enhance LTSF. By addressing the inductive bias of LVMs toward periodicity through a tailored backcast-residual decomposition, DMMV effectively integrates numerical and visual perspectives. Extensive experiments on benchmark datasets demonstrate that DMMV outperforms both single-view and SOTA multi-modal baselines, validating its effectiveness. This work highlights the potential of MMVs and LVMs in advancing LTSF, offering a new direction for future research in this domain.

## References

- [1] A. Ardid, D. Dempsey, C. Caudron, S. Cronin, B. Kennedy, T. Girona, D. Roman, C. Miller, S. Potter, O. D. Lamb, et al. Ergodic seismic precursors and transfer learning for short term eruption forecasting at data scarce volcanoes. *Nat. Commun.*, 16(1):1758, 2025.
- [2] H. Bao, L. Dong, S. Piao, et al. Beit: Bert pre-training of image transformers. In *ICLR*, 2022.
- [3] J. O. Caro, A. H. de Oliveira Fonseca, S. A. Rizvi, M. Rosati, C. Averill, J. L. Cross, P. Mittal, E. Zappala, R. M. Dhodapkar, C. Abdallah, et al. Brainlm: A foundation model for brain activity recordings. In *ICLR*, 2024.
- [4] M. Chen, L. Shen, Z. Li, X. J. Wang, J. Sun, and C. Liu. VisionTS: Visual masked autoencoders are free-lunch zero-shot time series forecasters. In *ICML*, 2025.
- [5] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009.
- [6] J. Deng, F. Ye, D. Yin, X. Song, I. Tsang, and H. Xiong. Parsimony or capability? decomposition delivers both in long-term time series forecasting. In *NeurIPS*, 2024.
- [7] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, et al. An image is worth 16 x 16 words: Transformers for image recognition at scale. In *ICLR*, 2021.
- [8] W. Fan, S. Zheng, X. Yi, W. Cao, Y. Fu, J. Bian, and T.-Y. Liu. Depts: Deep expansion learning for periodic time series forecasting. In *ICLR*, 2022.
- [9] Y. Gong, Y.-A. Chung, and J. Glass. Ast: Audio spectrogram transformer. In *INTERSPEECH*, 2021.
- [10] N. Gruver, M. Finzi, S. Qiu, and A. G. Wilson. Large language models are zero-shot time series forecasters. In *NeurIPS*, 2023.
- [11] S. Han, X. Hu, H. Huang, M. Jiang, and Y. Zhao. Adbench: Anomaly detection benchmark. In *NeurIPS*, 2022.
- [12] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick. Masked autoencoders are scalable vision learners. In CVPR, 2022.
- [13] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [14] Y. Jiang, Z. Pan, X. Zhang, S. Garg, A. Schneider, Y. Nevmyvaka, and D. Song. Empowering time series analysis with large language models: a survey. In *IJCAI*, 2024.
- [15] M. Jin, S. Wang, L. Ma, Z. Chu, J. Y. Zhang, X. Shi, P.-Y. Chen, Y. Liang, Y.-F. Li, S. Pan, and Q. Wen. Time-LLM: Time series forecasting by reprogramming large language models. In *ICLR*, 2024.
- [16] W. Kim, B. Son, and I. Kim. Vilt: Vision-and-language transformer without convolution or region supervision. In *ICML*, 2021.
- [17] A. Kline, H. Wang, Y. Li, S. Dennis, M. Hutch, Z. Xu, F. Wang, F. Cheng, and Y. Luo. Multimodal machine learning in precision health: A scoping review. *npj Digit. Med.*, 5(1):171, 2022.
- [18] I. Koprinska, D. Wu, and Z. Wang. Convolutional neural networks for energy time series forecasting. In *IJCNN*, 2018.
- [19] X. Li, Y. Kang, and F. Li. Forecasting with time series imaging. Expert Syst. Appl., 160:113680, 2020.
- [20] Z. Li, S. Li, and X. Yan. Time series as images: Vision transformer for irregularly sampled time series. In *NeurIPS*, 2023.

- [21] S. Lin, W. Lin, X. Hu, W. Wu, R. Mo, and H. Zhong. Cyclenet: enhancing time series forecasting through modeling periodic patterns. In *NeurIPS*, 2024.
- [22] S. Lin, W. Lin, W. Wu, H. Chen, and J. Yang. Sparsetsf: Modeling long-term time series forecasting with 1k parameters. In *ICML*, 2024.
- [23] H. Liu, C. Li, Q. Wu, and Y. J. Lee. Visual instruction tuning. In NeurIPS, 2023.
- [24] P. Liu, H. Guo, T. Dai, N. Li, J. Bao, X. Ren, Y. Jiang, and S.-T. Xia. Calf: Aligning Ilms for time series forecasting via cross-modal fine-tuning. In AAAI, 2025.
- [25] Y. Liu, H. Wu, J. Wang, and M. Long. Non-stationary transformers: Exploring the stationarity in time series forecasting. In *NeurIPS*, 2022.
- [26] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, 2021.
- [27] X. Ma, Z. Dai, Z. He, J. Ma, Y. Wang, and Y. Wang. Learning traffic as images: A deep convolutional neural network for large-scale transportation network speed prediction. *Sensors*, 17(4):818, 2017.
- [28] M. A. Morid, O. R. L. Sheng, and J. Dunbar. Time series prediction using deep learning methods in healthcare. *Trans. Manag. Inf. Syst.*, 14(1):1–29, 2023.
- [29] J. Ni, Z. Zhao, C. Shen, H. Tong, D. Song, W. Cheng, D. Luo, and H. Chen. Harnessing vision models for time series analysis: A survey. In *IJCAI*, 2025.
- [30] Y. Nie, N. H. Nguyen, P. Sinthong, and J. Kalagnanam. A time series is worth 64 words: Long-term forecasting with transformers. In *ICLR*, 2023.
- [31] B. N. Oreshkin, D. Carpov, N. Chapados, and Y. Bengio. N-beats: Neural basis expansion analysis for interpretable time series forecasting. In *ICLR*, 2020.
- [32] V. Prithyani, M. Mohammed, R. Gadgil, R. Buitrago, V. Jain, and A. Chadha. On the feasibility of vision-language models for time-series classification. *arXiv*:2412.17304, 2024.
- [33] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021.
- [34] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.
- [35] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *CVPR*, 2015.
- [36] M. Tan, M. Merrill, V. Gupta, T. Althoff, and T. Hartvigsen. Are language models actually useful for time series forecasting? In *NeurIPS*, 2024.
- [37] H. Touvron, M. Cord, M. Douze, F. Massa, A. Sablayrolles, and H. Jégou. Training data-efficient image transformers & distillation through attention. In *ICML*, 2021.
- [38] A. Trindade. Electricity Load Diagrams. UCI Machine Learning Repository, 2015.
- [39] W. Wang, R. Arora, K. Livescu, and J. Bilmes. On deep multi-view representation learning. In ICML, 2015.
- [40] Q. Wen, T. Zhou, C. Zhang, W. Chen, Z. Ma, J. Yan, and L. Sun. Transformers in time series: a survey. In *IJCAI*, 2023.
- [41] C. Wimmer and N. Rekabsaz. Leveraging vision-language models for granular market change prediction. *arXiv:2301.10166*, 2023.
- [42] G. Woo, C. Liu, D. Sahoo, A. Kumar, and S. Hoi. Etsformer: Exponential smoothing transformers for time-series forecasting. *arXiv preprint arXiv:2202.01381*, 2022.

- [43] H. Wu, T. Hu, Y. Liu, H. Zhou, J. Wang, and M. Long. Timesnet: Temporal 2d-variation modeling for general time series analysis. In *ICLR*, 2023.
- [44] H. Wu, J. Xu, J. Wang, and M. Long. Autoformer: Decomposition transformers with Auto-Correlation for long-term series forecasting. In *NeurIPS*, 2021.
- [45] Z. Xie, Z. Zhang, Y. Cao, Y. Lin, J. Bao, Z. Yao, Q. Dai, and H. Hu. Simmim: A simple framework for masked image modeling. In *CVPR*, 2022.
- [46] H. Xue and F. D. Salim. Promptcast: A new prompt-based learning paradigm for time series forecasting. *Trans. Knowl. Data Eng.*, 36(11):6851–6864, 2023.
- [47] L. Yang, Y. Wang, X. Fan, I. Cohen, J. Chen, Y. Zhao, and Z. Zhang. Vitime: A visual intelligence-based foundation model for time series forecasting. arXiv preprint arXiv:2407.07311, 2024.
- [48] G. Yu, J. Zou, X. Hu, A. I. Aviles-Rivero, J. Qin, and S. Wang. Revitalizing multivariate time series forecasting: Learnable decomposition with inter-series dependencies and intra-series variations modeling. In *ICML*, 2024.
- [49] A. Zeng, M. Chen, L. Zhang, and Q. Xu. Are transformers effective for time series forecasting? In AAAI, 2023.
- [50] Z. Zeng, R. Kaur, S. Siddagangappa, et al. From pixels to predictions: Spectrogram and vision transformer for better time series forecasting. In *ICAIF*, 2023.
- [51] X. Zhang, R. R. Chowdhury, R. K. Gupta, and J. Shang. Large language models for time series: A survey. In *IJCAI*, 2024.
- [52] S. Zhong, W. Ruan, M. Jin, H. Li, Q. Wen, and Y. Liang. Time-VLM: Exploring multimodal vision-language models for augmented time series forecasting. arXiv preprint arXiv:2502.04395, 2025.
- [53] H. Zhou, S. Zhang, J. Peng, S. Zhang, J. Li, H. Xiong, and W. Zhang. Informer: Beyond efficient transformer for long sequence time-series forecasting. In *AAAI*, 2021.
- [54] T. Zhou, Z. Ma, Q. Wen, X. Wang, L. Sun, and R. Jin. Fedformer: Frequency enhanced decomposed transformer for long-term series forecasting. In *ICML*, 2022.
- [55] T. Zhou, P. Niu, X. Wang, L. Sun, and R. Jin. One fits all: Power general time series analysis by pretrained LM. In *NeurIPS*, 2023.
- [56] J. Zhuang, L. Yan, Z. Zhang, R. Wang, J. Zhang, and Y. Gu. See it, think it, sorted: Large multi-modal models are few-shot time series anomaly analyzers. arXiv preprint arXiv:2411.02465, 2024.

# **NeurIPS Paper Checklist**

## 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: The main claims in the abstract accurately reflect our contributions.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

#### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: We discuss the limitations of this work in Appendix A.1.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

## 3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: There is no theoretical result in this paper.

#### Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

# 4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We provide complete experimental details in Appendix C.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
- (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

## 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We provide a single zip file of code with the additional supplementary materials. We describe the code usage in the README file.

#### Guidelines

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how
  to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We describe the complete training details and hyperparameter choices in Appendix C.2.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental
  material.

## 7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: We discuses about the stantard deviations in Appendix D.7.

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).

- It should be clear whether the error bar is the standard deviation or the standard error
  of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: We report the computational resource requirements of our proposed method in Appendix C.3.

## Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

## 9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: Our research aligns with the NeurIPS Code of Ethics.

#### Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a
  deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

## 10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: We discuss the broader impacts of this work in Appendix A.2.

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: This paper poses no such risks.

#### Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do
  not require this, but we encourage authors to take this into account and make a best
  faith effort.

## 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We use only publicly released assets, each of which is properly cited with clear attribution and compliance with the associated licenses.

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the
  package should be provided. For popular datasets, paperswithcode.com/datasets
  has curated licenses for some datasets. Their licensing guide can help determine the
  license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

• If this information is not available online, the authors are encouraged to reach out to the asset's creators.

#### 13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: The code will be made publicly available after acceptance, along with comprehensive documentation and scripts to reproduce the main experimental results.

#### Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

### 14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: This work does not involve crowdsourcing nor research with human subjects. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

# 15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: This work does not involve crowdsourcing nor research with human subjects. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

## 16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: LLMs were not used in the development of the core methods of this research. Any usage was limited to non-substantive tasks such as proofreading or formatting.

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.

# A Limitation and Broader Impact

#### A.1 Discussion of Limitations

Adapting LVMs to LTSF is an emerging area of active research. This work serves a pioneering effort in investigating LVMs within an MMV framework for LTSF. As an initial exploration, we acknowledge some limitations in this work. First, one limitation of the current best LVM forecaster (e.g., VisionTS) is its sensitivity to segment length used in image construction due to its inductive bias, as discussed in §3 (Fig. 3). By incorporating an additional numerical view for modeling the global trend, the proposed DMMV is expected to alleviate this sensitivity. In our further analysis in Appendix D.2, we observe DMMV-A is less sensitive to the change of segment lengths than VisionTS on some datasets, but cannot consistently enhance the robustness over different datasets, despite the improved overall forecasting performance. This may be caused by the higher weights automatically allocated to  $f_{vis}(\cdot)$  than  $f_{num}(\cdot)$  by the gate fusion mechanism (Fig. 6), which could make the model prone to inherit the behavior of the LVM used in  $f_{vis}(\cdot)$  to some extent, including its sensitivity to segment length, but with a less extent than a sole LVM. As such, a future work to improve DMMV is to reduce such sensitivity to an unnoticeable effect. Second, under our proposed  $\hat{BCMASK}$  strategy, the vision backbone in  $f_{vis}(\cdot)$  is reused three times during training and inference – once for forecasting and twice for reconstructing different masked parts of the look-back window. The triple use of  $f_{\text{vis}}(\cdot)$  could lead to a non-trivial computational overhead. In this work, considering the remarkable performance improvement of BCMASK over other masking strategies (Table 2) and its possibly minimum use of  $f_{vis}(\cdot)$ , as analyzed in §3.2, we take it as the current solution. However, as a future work, we expect to further reduce the use of  $f_{vis}(\cdot)$  and improve the efficiency, by methods such as joint backcasting and forecasting within a single forward pass or amortizing the use of LVMs across multiple time series samples. Finally, the proposed method shares a similar limitation of the existing LVM forecasters. The imaging process transforms an input time series into a 2D image-like representation to fit pre-trained LVMs, which typically expect a high-resolution input of size such as  $224 \times 224 \times 3$ . Therefore, upsampling is performed on a smaller image obtained from the patched time series to fit the input requirements of LVMs. While critical to compatibility, resizing may introduce subtle changes that may distort the original temporal structures to some extent. Thus an imaging process that can reflect the temporal patterns more accurately is in demand in future works.

## A.2 Broader Impact

LTSF plays a vital role across various domains, including geoscience [1], neuroscience [3], energy [18], healthcare [28], and smart city [27]. This work proposes a novel MMV framework DMMV that integrates LVMs and a numerical forecaster, which could serve as a groundwork in the emerging area of LVM-based time series analysis and shed some lights on broader areas that integrate LLMs, VLMs, and large multi-modal models (LMMs) for future research on multi-modal and agentic time series analysis. This work does not involve sensitive data, legal risks, or ethical concerns. To the best of our knowledge, it does not adversely affect any specific population. The proposed method could serve as a general-purpose time series forecasting technique with a relatively broad applicability and social acceptability.

## **B** Benchmark and Baseline

#### **B.1** Benchmark Datasets

Following [53, 44, 30, 49, 36, 4], our experiments are conducted on 8 widely used LTSF benchmark datasets that cover a wide range of sampling frequencies, number of variates, levels of periodicity, and real-world domains. The four ETT datasets (ETTh1, ETTh2, ETTm1, ETTm2) record oil temperature from two electric transformers, sampled at 15-minute and hourly intervals. The Weather dataset collects measurements of meteorological indicators in Germany every 10 minutes. The Illness dataset keeps weekly counts of patients and the influenza-like illness ratio from the United States. The Traffic dataset measures hourly road occupancy rates from sensors on San Francisco freeways. The Electricity dataset records hourly electricity consumption of Portuguese clients. Table 3 summarizes the statistics of the datasets.

Table 3: Statistics of the benchmark datasets. "Dataset Size" is organized in (Train, Validation, Test).

Dataset	# Variates	Series Length	Dataset Size	Frequency
ETTh1	7	17420	(8545, 2881, 2881)	Hourly
ETTh2	7	17420	(8545, 2881, 2881)	Hourly
ETTm1	7	69680	(34465, 11521, 11521)	15 mins
ETTm2	7	69680	(34465, 11521, 11521)	15 mins
Weather	321	52696	(36792, 5271, 10540)	10 mins
Illness	7	966	(617, 74, 170)	Weekly
Traffic	862	17544	(12185, 1757, 3509)	Hourly
Electricity	21	26304	(18317, 2633, 5261)	Hourly

#### **B.2** Baselines

In the following, we provide a brief description for each baseline method involved in our experiments.

- Time-VLM [52] integrates time series data with visual views and contextual texts using a pre-trained VLM, ViLT, to enhance forecasting performance.
- VisionTS [4] reformulates time series forecasting as an image reconstruction problem using an LVM, MAE, for zero/few/full-shot forecasting.
- Time-LLM [15] reprograms LLMs by aligning time series patches with text tokens, enabling time series forecasting without re-training LLMs.
- GPT4TS [55] demonstrates that frozen pretrained LLMs, e.g., GPT, can be directly applied to a variety of time series tasks with strong performance.
- CALF [24] adapts LLMs to time series forecasting via cross-modal fine-tuning, bridging the distribution gap between textual and temporal data.
- CycleNet [21] enhances LTSF by explicitly modeling the periodic patterns in time series through a residual cycle forecasting technique.
- PatchTST [30] introduces a patching strategy and a channel-independence strategy for LTSF. It uses patches of time series as the input to a Transformer to capture the temporal dependency of semantically meaningful tokens (*i.e.*, patches).
- TimesNet [43] transforms an input time series into a 2D image-like representation and models temporal variations in the image using inception-like blocks for time series analysis.
- DLinear [49] decomposes an input time series into trend and seasonal components, each of which is modeled by linear layers for time series forecasting.
- FEDformer [54] incorporates frequency-enhanced attention mechanisms by combining Fourier transforms with seasonal-trend decomposition in a Transformer framework.
- Autoformer [44] introduces an auto-correlation mechanism within a Transformer architecture to capture long-term dependencies in time series data.
- Stationary [25] combines series stationarization and de-stationary attention mechanisms to solve the over-stationarization problem in time series forecasting.
- ETSformer [42] decomposes an input time series into interpretable components with exponential smoothing attention and frequency attention for time series forecasting.
- Informer [53] proposes a ProbSparse self-attention mechanism to reduce the computational complexity of LTSF with Transformer models.

## C Implementation Details

## C.1 Pre-trained LVM Checkpoints

As described in §3.3,  $f_{vis}(\cdot)$  uses pre-trained LVMs. For MAE, we use the checkpoint released by  $Meta\ Research^{-1}$ , which was pretrained on  $224 \times 224 \times 3$  sized images from ImageNet-1K [5] with

<sup>1</sup>https://github.com/facebookresearch/mae

# Algorithm 1: The Training Algorithm of DMMV-A

```
Input: training dataset \mathcal{D}_{\text{train}} = \{\mathbf{X}_i, \mathbf{Y}_i\}_{i=1}^n, where \mathbf{X}_i \in \mathbb{R}^{D \times T} is an MTS, \mathbf{Y}_i \in \mathbb{R}^{D \times H} is
             the ground truth of forecast
   Output: model parameters of DMMV-A
 1 Load pre-trained f_{vis}(\cdot) and freeze its weights
2 Randomly initialize f_{\text{num}}(\cdot) and the gating parameter g
    /* stage 1: numerical forecaster training */
3 for i \leftarrow 1 to MaxEpoch do
        for (X, Y) in Dataloader (\mathcal{D}_{train}) do
              /* channel-independence strategy is applied in the following */
              \mathbf{\hat{X}}_{\text{season}}, \mathbf{\hat{Y}}_{\text{season}} \leftarrow f_{\text{vis}}(\mathbf{X}, \text{BCMASK}) // backcast/forecast seasonal part
 5
              \mathbf{X}_{\text{trend}} \leftarrow \mathbf{X} - \mathbf{\hat{X}}_{\text{season}}
                                                                                     // extract trend component
              \hat{\mathbf{Y}}_{\text{trend}} \leftarrow f_{\text{num}}(\mathbf{X}_{\text{trend}})
                                                                                  // forecast trend with f_{	exttt{num}}(\cdot)
 7
              \mathbf{\hat{Y}} \leftarrow g \circ \mathbf{\hat{Y}}_{\text{season}} + (1 - g) \circ \mathbf{\hat{Y}}_{\text{trend}}
                                                                                                           // gate fusion
              Calculate \ell_{\text{MSE}}(\hat{\mathbf{Y}}, \mathbf{Y})
                                               // calculate MSE loss as specified in §3.3
              Update model parameters of f_{\text{num}}(\cdot) and g
10
             if Early stopping condition is TRUE then
11
                   Break
12
             end
13
14
        end
15 end
    /* stage 2: joint training */
16 Unfreeze the norm layers in f_{\text{vis}}(\cdot)
17 for i \leftarrow 1 to MaxEpoch do
        for (X, Y) in Dataloader (\mathcal{D}_{train}) do
              /* Repeat lines 5-9 */
              Update model parameters of f_{\text{num}}(\cdot), norm layers in f_{\text{vis}}(\cdot), and parameter g
19
             if Early stopping condition is TRUE then
20
                   Break
21
             end
22
23
        end
24 end
```

ViT-Base Backbone. For SimMIM, we adopt the checkpoint released by *Microsoft*<sup>2</sup>, which has the same pretraining setting as aforementioned for MAE. For these two LVM backbones, the base versions are adopted to balance the performance and computational costs.

## **C.2** Training Details

For training the proposed DMMV-S and DMMV-A models, we adopt AdamW optimizer throughout the experiments. The batch size is set to 64 for the ETT datasets and Illness dataset, and set to 8 for the other three datasets to balance training stability and memory consumption.

For both DMMV-S and DMMV-A, we propose a two-stage training scheme to facilitate effective integration of numerical and visual features:

• Stage 1 (Numerical forecaster training). In this stage, we freeze all parameters of  $f_{\text{vis}}(\cdot)$  and train  $f_{\text{num}}(\cdot)$  only. This warm-up step prevents  $f_{\text{vis}}(\cdot)$  from updating with unstable gradients caused by the random representations from the under-trained  $f_{\text{num}}(\cdot)$ . In this stage, the learning rate is set to 0.01. The training runs up to a maximum of 50 epochs on the training set. Early stopping is applied with a patience of 10 epochs.

<sup>&</sup>lt;sup>2</sup>https://github.com/microsoft/SimMIM

• Stage 2 (Joint training). In this stage, we unfreeze the layer normalization parameter in  $f_{\text{vis}}(\cdot)$  and jointly train them with  $f_{\text{num}}(\cdot)$  to enable deep fusion of visual and numerical views. The learning rate is reduced to 0.005 to preserve learned features and stabilize training. The training at this stage runs up to 5 epochs. Early stopping is applied with a patience of 2 epochs.

The detailed training algorithm of DMMV-A is summarized in Algorithm 1.

## **C.3** Running Environment

The experiments are conducted on a Linux server (kernel 5.15.0-139) with 8x NVIDIA RTX 6000 Ada GPUs (48 GB each). The environment uses Python 3.12.8, PyTorch 2.5.1 with CUDA 12.4 and cuDNN 9.1. The key libraries include NumPy 2.1.3, Pandas 2.2.3, Matplotlib 3.10.0, SciPy 1.15.1, scikit-learn 1.6.1, and torchvision 0.20.1.

# D More Experimental Results

## D.1 Comparison with All Baselines

Table 4 provides the full results of comparing DMMV-A and DMMV-S with all of the 14 baseline methods, which complements Table 1 in the paper. In Table 4, Time-VLM's results on Illness dataset is marked by "-" since its paper doesn't report the results and its code is not publicly available at the time of this experiment. CycleNet's paper doesn't report its results on Illness dataset, so we run its code and reproduce its results on Illness dataset in Table 4.

From Table 4, we can observe that DMMV-A maintains a clear advantage when compared against all of the baseline methods. It achieves 41 first-place results, significantly surpassing the second-best method. Additionally, taking a closer look at all compared methods, MMV-based methods LVM-based methods, and decomposition-based methods demonstrate superiority over other baseline methods. This suggests the synergy of MMV framework, LVMs, and decomposition strategy, which are explored by the proposed DMMV model.

#### **D.2** Further Analysis of The Inductive Bias

A contribution of our work lies in the in-depth analysis of an inductive bias of the current best LVM forecasters. In §3, we have discussed the impact of the alignment of the segment length and the period of time series on model performance. We find that the LVM exhibits a strong *inter-period consistency* when applied to synthetic data. The function of the synthetic time series is  $x(t) = A(t) \cdot \sin\left(\frac{2\pi t}{P}\right)$ , where the period P is set to 24 and the amplitude function A(t) decreases linearly over time. The forecasts are more accurate when the segment length is a multiple of the period (e.g., 24, 48) than other values. This section provides detailed quantitative results on the synthetic data in Table 5. From Table 5, the fluctuations in MSEs and MAEs across different segment lengths other than 24 and 48 support the findings of the inductive bias toward "forecasting periods".

In addition, we evaluate the performance of the proposed method DMMV-A and VisionTS *w.r.t.* varying segment lengths to compare their robustness to the change of segment length. Fig. 10 summarizes the results in terms of MSE on four benchmark datsets, where the segment length varies from  $\frac{P}{6}$  to  $\frac{6P}{6}$  and P is a period of the input time series. From Fig. 10, we have several observations. First, DMMV-A consistently outperforms VisionTS, validating the effectiveness of the proposed MMV framework. Second, in contrast to VisionTS, DMMV-A exhibits a better robustness to the change of segment length on ETTh1 and Weather datasets, but has a similar sensitivity to the change of segment length as VisionTS on ETTm1 and Illness datasets. This implies that by incorporating  $f_{\text{num}}(\cdot)$ , DMMV-A can alleviate  $f_{\text{vis}}(\cdot)$ 's sensitivity to the inductive bias to some extent. However, the current DMMV-A does not fully mitigate this limitation, suggesting a future work for method development as discussed in Appendix A.

Table 4: Full LTSF performance of the compared methods on the benchmark datasets. Lower MSE and MAE indicate better performance. The best performance is highlighted in red. Time-VLM results on the Illness dataset are unavailable in [52]. Its code was not publicly available at the time of this experiment. As such, its results on Illness dataset are marked by "-".

er	MAE 0.713 0.792 0.809 0.865	1,525 1,931 1,835 1,625 1,729	0.571 0.669 0.871 0.823 0.734	0.453 0.563 0.887 0.810	1.677 1.467 1.469 1.564 1.544	)368 )386 )394 )439 )397	0.384 0.544 0.523 0.741 0.548	),391 ),379 ),420 ),472 ),416	
Informer	MSE 1.008 (1.107 (1.181 (1.040 (1.181 (1.040 (1.181 (1.040 (1.181 (1.040 (1.181 (1.040 (1.181 (1.040	3.755 5.602 1.721 8.647 1.431	.672 .795 .212 .166 .961	).365 ).533 (1.363 (3.379 1.410	5.764 4.755 4.763 5.264 5.137	) 274 ) 296 ) 300 ) 373 ) 311	300 598 578 1.059 1.634	).719 ).696 ).777 ).864 ).764	0
ner	H 0 4 - 0 0	391 3 3,439 5 3,479 4 3,497 3 4,452 4	0.398 C 0.410 C 0.428 1 0.462 1 0.425 C		1.020 1.007 2.972 4 1.016 5	304 (0.315 (0.329 (0.345 (0.323 (0.32		0.392 0.399 0.396 0.396 0.396	
ETSformer	MSE 19.494 (19.538 (19.574 (19.552 (19	),340 (),430 (),485 (),500 (),439 ()	0.375 ( 0.408 ( 0.435 ( 0.499 (	).189 () ).253 () ).314 () ).414 ()	2.527 1 2.615 1 2.359 ( 2.487 1	0.187 (0.1199 (0.1212 (0.1233 (0.1208 (0.1212	),197 (),237 (),298 (),352 (),271	).607 (0.621 (0.622 (0.632 (0.	0
ary	MAE 10.2504 (0.535 (0.5	),458 (),493 (),551 (),560 (),516 ()	0.398 ( 0.444 ( 0.464 ( 0.516 (	0.274 (0.339 (0.361 (0.413 (0.347 (0.	).945 ).848 ).900 ).963	) 273 ) 286 ) 304 ) 321 ) 296	0.223 (0.285 (0.338 (0.410 (0.314 (	0.338 (0.340 (0.328 (0.355 (0.340 (0.340 (0.355 (0.340 (0.	
Stationary	MSE 0.513 0.534 0.588 0.643 0.570	0.476 0.512 0.552 0.562 0.526	0.386 0.459 0.495 0.585 0.481	0.192 0.280 0.334 0.417 0.306	2.294 1.825 2.010 2.178 2.077	0.169 0.182 0.200 0.222 0.193	0.173 0.245 0.321 0.414 0.288	0.612 0.613 0.618 0.653 0.624	0
rmer	MAE 0.489 0.496 0.512 0.487	0.388 0.452 0.486 0.511 0.459	0.475 0.496 0.537 0.561 0.517	0.339 0.340 0.372 0.432 0.371	1.287 1.148 1.085 1.125 1.161	0.317 0.334 0.338 0.361 0.338	0.336 0.367 0.395 0.428 0.382	0.388 0.382 0.337 0.408 0.379	
Autoformer	MSE 0.449 0.500 0.521 0.514 0.496	0.346 0.456 0.482 0.515 0.450	0.505 0.553 0.621 0.671 0.588	0.255 0.281 0.339 0.433 0.327	3.483 3.103 2.669 2.770 3.006	0.201 0.222 0.231 0.254 0.257	0.266 0.307 0.359 0.419 0.338	0.613 0.616 0.622 0.660 0.628	0
rmer	MAE 0.419 0.448 0.465 0.507 0.460	0.397 0.439 0.487 0.474 0.449	0.419 0.441 0.459 0.490 0.452	0.287 0.328 0.366 0.415 0.349	1.260 1.080 1.078 1.157 1.144	0.308 0.315 0.329 0.355 0.327	0.296 0.336 0.380 0.428 0.360	0.366 0.373 0.383 0.382 0.376	
FEDformer	MSE 0.376 0.420 0.459 0.506 0.440	0.358 0.429 0.496 0.463 0.437	0.379 0.426 0.445 0.543 0.448	0.203 0.269 0.325 0.421 0.305	3.228 2.679 2.622 2.857 2.847	0.193 0.201 0.214 0.246 0.214	0.217 0.276 0.339 0.403 0.309	0.587 0.604 0.621 0.626 0.610	0
DLinear	MAE 0.399 0.416 0.490 0.430	0.353 0.418 0.465 0.551 0.447	0.343 0.365 0.386 0.421 0.379	0.260 0.303 0.342 0.421 0.332	1.081 0.963 1.024 1.096 1.041	0.237 0.249 0.267 0.301 0.264	0.237 0.282 0.319 0.362 0.300	0.282 0.287 0.296 0.315 0.295	
DLi	MSE 0.375 0.405 0.439 0.472 0.423	0.289 0.383 0.448 0.605 0.431	0.299 0.335 0.369 0.425 0.357	0.167 0.224 0.281 0.397 0.267	2.215 1.963 2.130 2.368 2.169	0.140 0.153 0.169 0.203 0.166	0.176 0.220 0.265 0.333 0.249	0.410 0.423 0.436 0.466 0.466	
TimesNet	MAE 0.402 0.429 0.469 0.500 0.450	0.374 0.414 0.452 0.468 0.427	0.375 0.387 0.411 0.450 0.406	0.267 0.309 0.351 0.403 0.333	0.934 0.920 0.940 0.928 0.931	0.272 0.289 0.300 0.320 0.295	0.220 0.261 0.306 0.359 0.359	0.321 0.336 0.336 0.350 0.350	0
ŢŢ	MSE 0.384 0.436 0.491 0.521 0.458	0.340 0.402 0.452 0.462 0.414	0.338 0.374 0.410 0.478 0.400	0.187 0.249 0.321 0.408 0.291	2.317 1.972 2.238 2.027 2.139	0.168 0.184 0.198 0.220 0.193	0.172 0.219 0.280 0.365 0.259	0.593 0.617 0.629 0.640 0.620	
PatchTST	MAE 0.399 0.421 0.436 0.466 0.461	0.336 0.379 0.380 <b>0.422</b> <b>0.379</b>	0.342 0.369 0.392 0.420 0.381	0.255 0.292 0.329 0.385 0.315	0.754 0.834 0.815 0.788 0.798	0.222 0.240 0.259 0.290 0.253	0.198 0.241 0.282 0.334 0.264	0.249 0.256 0.264 0.286 0.264	_
Patch	MSE 0.370 0.413 0.447 0.447	0.274 0.339 0.329 <b>0.379</b> <b>0.330</b>	0.290 0.332 0.366 0.416 0.351	0.165 0.220 0.274 0.362 0.255	1.319 1.430 1.553 1.470 1.443	0.129 0.157 0.163 <b>0.197</b> 0.162	0.149 0.194 0.245 0.314 0.226	0.360 0.379 0.392 <b>0.432</b> 0.391	
Net	MAE 0.396 0.415 0.430 0.464 0.426	0.341 0.385 0.413 0.451 0.398	0.348 0.367 0.386 0.414 0.379	0.247 0.286 0.322 0.382 0.309	1.017 0.950 1.007 0.997 0.992	0.223 0.237 0.254 0.287 0.287	0.221 0.258 0.293 0.339 0.278	0.278 0.283 0.289 0.305 0.289	
CycleNet	MSE 0.374 0.406 0.431 0.450 0.415	0.279 0.342 0.371 0.426 0.355	0.299 0.334 0.368 0.417 0.355	0.159 0.214 0.269 0.363 0.251	2.255 2.121 2.187 2.187 2.185	0.128 0.144 0.160 0.198 0.158	0.167 0.212 0.260 0.328 0.242	0.397 0.411 0.424 0.450 0.421	1
ц	MAE 0.393 0.426 0.440 0.466	0.336 0.378 0.394 0.428 0.384	0.350 0.376 0.401 0.438 0.391	0.255 0.300 0.341 0.395 0.323	0.788 0.837 0.890 0.962 0.869	0.240 0.254 0.270 0.300 0.266	0.207 0.251 0.292 0.345 0.274	0.274 0.276 0.286 0.301 0.284	
CALF	MSE 0.370 0.429 0.451 0.476	0.284 0.353 0.361 0.406 0.351	0.323 0.375 0.411 0.476 0.396	0.177 0.245 0.309 0.402 0.283	1.460 1.573 1.784 1.982 1.700	0.147 0.163 0.178 0.215 0.176	0.168 0.216 0.271 0.350 0.251	0.416 0.430 0.451 0.478 0.444	0
ITS	MAE 0.389 0.413 0.431 0.449	0.335 0.380 0.405 0.436 0.389	0.340 0.368 0.386 0.416 0.378	0.249 0.291 0.327 <b>0.376</b> 0.311	0.823 0.854 0.855 0.877 0.852	0.239 0.253 0.266 0.293 0.263	0.188 0.230 0.273 0.328 0.255	0.264 0.268 0.273 0.291 0.274	
GPT4TS	MSE 0.370 0.412 0.448 0.441	0.280 0.348 0.380 0.406 0.354	0.300 0.343 0.376 0.431 0.363	0.163 0.222 0.273 0.357 0.254	1.869 1.853 1.886 1.877 1.871	0.141 0.158 0.172 0.207 0.170	0.148 0.192 0.246 0.320 0.227	0.396 0.412 0.421 0.455 0.451	7
TW	MAE 0.402 0.421 0.438 0.468	0.346 0.391 0.414 0.434 0.396	0.341 0.369 0.379 0.419 0.377	0.248 0.304 0.329 0.382 0.316	0.807 0.833 1.012 0.925 0.894	0.233 0.247 0.267 0.290 0.259	0.199 0.261 0.279 0.342 0.270	0.267 0.271 0.296 0.291 0.281	
Time-I	MAE MSE MAE  0.386 0.376 0.402  0.407 0.407 0.421  0.450 0.457 0.468  0.419 0.418 0.432	0.286 0.361 0.390 0.405 0.361	0.291 0.341 0.359 0.433 0.356	0.162 0.235 0.280 0.366 0.261	1.792 1.833 2.269 2.177 2.018	0.137 0.152 0.169 0.200 0.165	0.155 0.223 0.251 0.345 0.244	0.392 0.409 0.434 0.451 0.422	1
LS	MAE 0.386 0.407 0.421 0.460	0.334 0.380 0.398 0.431 0.386	0.332 0.362 0.382 0.415 0.373	0.262 0.297 0.337 0.410 0.327	0.834 0.750 0.818 0.783 0.796	0.217 0.237 0.253 0.293 0.250	0.191 0.238 0.275 0.328 0.258	0.232 0.245 0.252 0.293 0.256	
VisionTS	MSE 0.355 0.395 0.419 0.458	0.288 0.349 0.364 0.403 0.351	0.284 0.327 0.354 0.411 0.344	0.174 0.228 0.281 0.384 0.267	1.613 1.316 1.548 1.450 1.482	0.127 0.148 0.163 0.199 0.159	0.146 0.194 0.243 0.318	0.346 0.376 0.389 <b>0.432</b> 0.386	6
LM	MAE 0.386 0.415 0.421 0.458	0.335 0.373 0.406 0.449 0.391	0.346 0.366 0.383 <b>0.410</b> 0.376	0.250 0.291 0.325 0.378 0.311		0.245 0.260 0.276 0.308 0.272	0.200 0.240 0.281 0.332 0.263	0.290 0.296 0.305 0.323 0.304	
Time-VLM	MSE 0.361 0.397 0.420 0.441	0.267 0.326 0.357 0.412 0.341	0.304 0.332 0.364 0.402 0.351	0.160 0.215 0.270 0.348 0.248		0.142 0.157 0.174 0.214 0.172	0.148 0.193 0.243 0.312 0.224	0.393 0.405 0.420 0.459 0.419	∞
- s	MAE 0.388 0.420 0.415 0.479 0.426	0.360   0.387   0.462   0.397	0.349   0.370   0.393   0.414   0.382   0	0.254   0.293   0.332   0.393   0.318	0.838   0.753   0.851   0.810   0.813	0.267   0.296   0.344   0.296	0.218   0.259   0.304   0.343   0.281	0.253   0.262   0.262   0.262   0.268	-
DMMV-S	MSE N 0.350 (0.399 (0.399 (0.472 (0.472 (0.405 (0.4	0.286 (0.331 (0.309 (0.430 (0.339 (0.	0.296 (0.328 (0.369 (0.349 (0.349 (	0.164 ( 0.217 ( 0.273 ( 0.362 ( 0.254 (	1.638 (1.323 (1.644 (1.473 (1.520 (1.	0.165 0.172 0.190 0.242 0.192	0.168 ( 0.220 ( 0.267 ( 0.322 ( 0.244 (	0.362 (0.385 (0.396 (0.436 (0.395 (0.	4
	MAE N. 2389 0.0405 0.413 0.0450 0.450 0.414 0.0414	0.349 (0.395 (0.384 (0.425 (0.3888 (0.388 (0.388 (0.388 (0.388 (0.388 (0.388 (0.388 (0.388 (0.3888 (0.388 (0.388 (0.388 (0.388 (0.3888 (0.388 (0.388 (0.388 (0.388 (0.388 (0.388 (0.388 (0.388 (0.388 (0.388 (0.388 (0.388 (0.3888	0.329 (0.381 (0.	0.260 0.298 0.327 0.381 0.317	0.754 1 0.745 1 0.810 1 0.773 1 0.771 1	0.213 (0.254 (0.254 (0.286 (0.248 (0.	0.195 (0.242 (0.273 (0.315 (0.256 (0.	0.237 (0.249 (0.256 (0.284 (0.257 (0.	
DMMV-A	MSE N 0.354 0 0.393 0 0.387 0 0.445 0	0.294 0 0.339 0 0.322 0 0.392 0 0.337 0	0.279 0 0.317 0 0.351 0 0.411 0	0.172 0 0.227 0 0.272 0 0.351 0	1.409 0 1.290 0 1.499 0 1.407 0	0.126 0 0.145 0 0.162 0 0.197 0 0.158 0	0.143 0 0.187 0 0.237 0 0.302 0	0.344 0 0.363 0 0.387 0 0.433 0	41
Model	Metric N 96 0 96 192 0.336 0.720 0.44 Avg.	96 0 192 0 336 0 720 0 Avg. 0	96 0 192 0 336 0 720 0 Avg. 0	96 0 192 0 336 0 720 0 Avg. 0	224 1 48 1 1 48 1 1 4 48 1 1 1 4 48 1 1 1 1	96 0 192 0 336 0 720 0	96 0 192 0 336 0 720 0 Avg. 0	96 0 192 0 336 0 720 0 Avg. 0	su
×	ETThi	ELLPS	ETTm1	ETTm2	Illness	Electricity	Weather	Traffic	# Wins

Table 5: Forecasting performance of an LVM w.r.t. varying segment length on a synthetic dataset. The function of the synthetic time series is  $x(t) = A(t) \cdot \sin\left(\frac{2\pi t}{P}\right)$ , where the period P = 24 and the amplitude function A(t) decreases linearly over time.

Segment Length	16	20	24	28	32	36	40	44	48
MSE	0.043	0.099	0.001	0.147	0.154	0.143	0.221	0.114	0.002
MAE	0.177	0.257	0.024	0.342	0.347	0.315	0.408	0.289	0.045

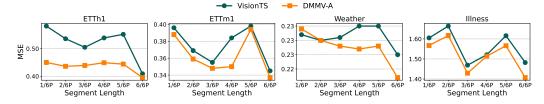


Figure 10: MSE Performance of DMMV-A and VisionTS w.r.t. varying segment length that is used in image construction. The x-axis indicates the segment length varies from  $\frac{1}{6}$  period to  $\frac{6}{6}$  period.

## **D.3** Ablation Study

In Table 2 (§4.2), we provide ablation analyses for DMMV-A. Table 6 provides the ablation analysis for DMMV-S, where MSE and MAE are averaged over different prediction lengths. In addition, Tables 7 (Table 8) includes the full results for Table 2 (Table 6) with all prediction lengths.

In Table 6, from (a), replacing the linear numerical forecaster with PatchTST can slightly improve the performance of DMMV-s, likely because DMMV-s relies more on the predictions from the numerical view than visual view (Fig. 6). Therefore, in this case, increasing the complexity of the numerical model can improve the ability of  $f_{\text{num}}(\cdot)$  and finally improve the overall performance. From (b), replacing MAE with SimMiM reduces the overall performance, this is the same as the findings in Table 2 for DMMV-A. From (c), gate-based fusion outperforms simple summation for DMMV-s, highlighting the effectiveness of gate fusion. From (d), fine-tuning the norm layers of  $f_{\text{vis}}(\cdot)$  improves the performance for DMMV-s, suggesting the used fine-tuning strategy.

Table 6: Ablation analysis of DMMV-S. MSE and MAE are averaged over different prediction lengths. Lower MSE and MAE are better. "Improvement" of each ablation is relative to DMMV-S.

Dataset $(\rightarrow)$	ETTh1		ETTm1		Illness		Wea	ther
Method $(\downarrow)$ , Metric $(\rightarrow)$	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
Dmmv-s	0.405	0.426	0.349	0.382	1.520	0.813	0.244	0.281
(a) $f_{\text{num}}(\cdot) \to \texttt{Transformer}$	0.402	0.423	0.342	0.376	1.544	0.841	0.229	0.264
Improvement	0.74%	0.47%	2.01%	1.31%	-1.65%	-3.44%	6.15%	6.41%
$(b) f_{\mathrm{vis}}(\cdot) \to \mathtt{SimMIM}$	0.415	0.423	0.355	0.382	1.810	0.875	0.233	0.272
Improvement	-2.47%	0.47%	-1.72%	0.00%	-19.16%	-7.63%	4.92%	3.20%
(c) Gate → Sum	0.419	0.435	0.355	0.379	1.453	0.790	0.256	0.297
Improvement	-3.46%	-2.24%	-2.01%	0.52%	4.34%	2.95%	-4.51%	-5.69%
(d) Freeze $f_{\text{vis}}(\cdot)$	0.436	0.442	0.368	0.386	2.125	0.969	0.251	0.288
Improvement	-7.41%	-3.76%	-5.75%	-1.05%	-39.83%	-19.07%	-2.46%	-2.14%

## D.4 Additional Visualizations on Decomposition

Fig. 11 and Fig. 12 several more examples the decomposed time series of DMMV-S and DMMV-A. Fig. 11 illustrates a case where the series has a localized periodic anomaly at time step around 192, which poses a challenge for detecting periodic patterns. In this case, DMMV-A effectively suppresses the influence of the anomaly and extracts a clear periodic pattern from the time series series. In contrast, DMMV-S is affected by the anomaly and fails to capture a smooth periodic pattern. Fig. 12 is an example with weak periodicity, where the periodic signal is either faint or overwhelmed by trend. In this case, DMMV-A is able to extract and utilize the underlying periodicity to produce reasonable forecasts, which is better than DMMV-S, suggesting the importance of the

Table 7: Full results of the ablation analysis of DMMV-A. Lower MSE and MAE are better. The Illness dataset uses prediction lengths of  $\{24, 36, 48, 60\}$  due to its short time series (in total 966 time steps), which is different from the prediction lengths of other datasets.

$\begin{array}{c} \textbf{Dataset}(\rightarrow) \\ \textbf{Method}(\downarrow), \ \ \textbf{Metric}(\rightarrow) \end{array}$	Length	ET MSE	Th1 MAE	ET MSE	Γm1 MAE	Illr MSE	ness MAE	Wea MSE	ther MAE
Method $(\downarrow)$ , Metric $(\rightarrow)$	Length	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
	96	0.354	0.389	0.279	0.329	1.409	0.754	0.143	0.195
	192	0.393	0.405	0.317	0.357	1.290	0.745	0.187	0.242
DMMV-A	336	0.387	0.413	0.351	0.381	1.499	0.810	0.237	0.273
	720	0.445	0.450	0.411	0.415	1.428	0.773	0.302	0.315
	Avg.	0.395	0.414	0.340	0.371	1.407	0.771	0.217	0.256
	96	0.357	0.389	0.279	0.329	1.604	0.823	0.145	0.193
	192	0.407	0.420	0.318	0.359	1.250	0.742	0.187	0.239
(a) $f_{ ext{num}}(\cdot)  o  ext{Transformer}$	336	0.389	0.411	0.352	0.382	1.555	0.803	0.241	0.283
	720	0.474	0.462	0.407	0.416	1.359	0.774	0.301	0.326
	Avg.	0.407	0.421	0.339	0.372	1.442	0.786	0.219	0.260
	96	0.358	0.383	0.301	0.348	1.729	0.832	0.145	0.194
	192	0.405	0.41	0.325	0.363	1.643	0.734	0.192	0.242
(b) $f_{ ext{vis}}(\cdot)  o  ext{ t SimMiM}$	336	0.412	0.414	0.354	0.383	1.689	0.845	0.241	0.275
	720	0.453	0.452	0.398	0.412	1.534	0.845	0.328	0.332
	Avg.	0.407	0.415	0.345	0.377	1.649	0.814	0.227	0.261
	96	0.373	0.400	0.286	0.339	1.728	0.845	0.156	0.214
	192	0.414	0.424	0.329	0.369	1.423	0.795	0.204	0.261
(c) Gate $\rightarrow$ Sum	336	0.411	0.422	0.364	0.392	1.693	0.920	0.258	0.302
	720	0.457	0.461	0.427	0.430	1.580	0.890	0.315	0.335
	Avg.	0.414	0.427	0.352	0.383	1.606	0.863	0.233	0.278
	96	0.384	0.402	0.288	0.342	1.628	0.840	0.145	0.198
	192	0.413	0.440	0.325	0.363	1.325	0.796	0.191	0.244
$(d)BCMASK \rightarrow No mask$	336	0.434	0.448	0.361	0.384	1.606	0.865	0.241	0.285
	720	0.474	0.473	0.421	0.419	1.414	0.811	0.308	0.340
	Avg.	0.426	0.441	0.349	0.377	1.493	0.828	0.221	0.267
	96	0.348	0.384	0.279	0.329	1.618	0.859	0.146	0.197
	192	0.388	0.405	0.318	0.360	1.318	0.798	0.189	0.240
(e)BCMASK $\rightarrow$ Random mask	336	0.383	0.404	0.350	0.381	1.560	0.858	0.243	0.282
	720	0.458	0.462	0.414	0.418	1.392	0.800	0.312	0.328
	Avg.	0.394	0.414	0.340	0.372	1.472	0.829	0.223	0.262
	96	0.389	0.402	0.293	0.342	1.482	0.761	0.161	0.224
	192	0.434	0.425	0.335	0.367	1.218	0.694	0.203	0.287
(f) Freeze $f_{\text{vis}}(\cdot)$	336	0.431	0.428	0.372	0.389	1.58	0.82	0.285	0.302
	720	0.468	0.457	0.431	0.422	1.489	0.815	0.335	0.338
	Avg.	0.431	0.428	0.358	0.380	1.442	0.773	0.246	0.288
	96	0.352	0.387	0.274	0.329	1.728	0.938	0.143	0.195
	192	0.402	0.414	0.315	0.358	1.841	0.940	0.187	0.242
(g) W/o decomposition	336	0.391	0.410	0.347	0.382	1.672	0.886	0.237	0.284
	720	0.487	0.486	0.417	0.422	1.606	0.846	0.309	0.350
	Avg.	0.408	0.424	0.338	0.373	1.712	0.903	0.219	0.268

proposed adaptive decomposition method. In summary, the results demonstrate that DMMV-A has a strong modeling ability of temporal structures and robustness to fluctuations even when dealing with anomalous or weakly periodic time series, validating its reliability and applicability across a broad range of scenarios.

## D.5 Additional Visualizations on Masking Strategies

Fig. 13 and Fig. 14 present additional examples of BCMASK in DMMV-A. Similar to §4.3, Fig. 13 and Fig. 14 compare different masking methods. From both figures, we observe that BCMASK produces smooth patterns along the temporal (*x*-axis) dimension, effectively capturing periodic structures. Notably, when the input time series contains an anomaly (*e.g.*, Fig. 13, time steps 288-336), BCMASK can effectively extract the periodic patterns.

### **D.6** Impact of Look-Back Window

Fig. 15 provides the MSE results that compare DMMV-A with the other three models. Fig. 15 demonstrate a similar trend as that of the MAE results in Fig. 9.

Table 8: Full results of the ablation analysis of DMMV-s. Lower MSE and MAE are better. The Illness dataset uses prediction lengths of  $\{24, 36, 48, 60\}$  due to its short time series (in total 966 time steps), which is different from the prediction lengths of other datasets.

$Dataset(\rightarrow)$		ET	Th1	ET	Γm1	Illr	iess	Wea	ther
$\mathbf{Method}(\downarrow), \ \ \mathbf{Metric}(\rightarrow)$	Length	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
	96	0.350	0.388	0.296	0.349	1.638	0.838	0.168	0.218
	192	0.399	0.420	0.328	0.370	1.323	0.753	0.220	0.259
DMMV-S	336	0.399	0.415	0.369	0.393	1.644	0.851	0.267	0.304
	720	0.472	0.479	0.401	0.414	1.473	0.810	0.322	0.343
	Avg.	0.405	0.426	0.349	0.382	1.520	0.813	0.244	0.281
	96	0.352	0.387	0.286	0.339	1.613	0.829	0.148	0.194
	192	0.401	0.420	0.325	0.364	1.417	0.825	0.193	0.240
(a) $f_{\text{num}}(\cdot) \to \texttt{Transformer}$	336	0.395	0.415	0.354	0.387	1.610	0.853	0.246	0.280
	720	0.460	0.471	0.401	0.414	1.536	0.858	0.330	0.341
	Avg.	0.402	0.423	0.342	0.376	1.544	0.841	0.229	0.264
	96	0.366	0.391	0.323	0.360	1.923	0.901	0.153	0.210
	192	0.412	0.420	0.331	0.364	1.812	0.863	0.194	0.248
(b) $f_{\mathrm{vis}}(\cdot)  o \mathtt{SimMiM}$	336	0.419	0.420	0.361	0.386	1.793	0.854	0.245	0.279
	720	0.464	0.461	0.404	0.416	1.712	0.883	0.339	0.352
	Avg.	0.415	0.423	0.355	0.382	1.810	0.875	0.233	0.272
	96	0.356	0.389	0.300	0.346	1.503	0.763	0.183	0.234
	192	0.403	0.417	0.334	0.365	1.350	0.746	0.236	0.277
(c) Gate $\rightarrow$ Sum	336	0.414	0.426	0.362	0.385	1.530	0.820	0.271	0.308
	720	0.504	0.506	0.424	0.420	1.429	0.830	0.333	0.369
	Avg.	0.419	0.435	0.355	0.379	1.453	0.790	0.256	0.297
	96	0.386	0.404	0.306	0.352	1.966	0.921	0.156	0.225
	192	0.436	0.434	0.347	0.375	2.050	0.945	0.240	0.261
(d) Freeze $f_{\text{vis}}(\cdot)$	336	0.436	0.440	0.377	0.392	2.223	0.999	0.271	0.312
* *	720	0.484	0.488	0.443	0.424	2.259	1.009	0.335	0.353
	Avg.	0.436	0.442	0.368	0.386	2.125	0.969	0.251	0.288

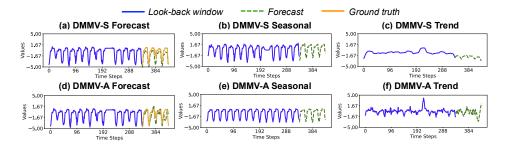


Figure 11: The decompositions of DMMV-S and DMMV-A on the same example in ETTh1: (a)(d) input time series and forecasts, (b)(e) seasonal component, and (c)(f) trend component.

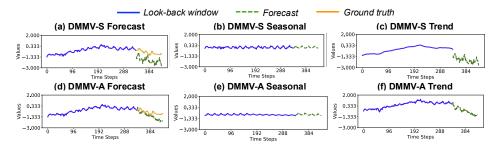


Figure 12: The decompositions of DMMV-S and DMMV-A on the same example in ETTh2: (a)(d) input time series and forecasts, (b)(e) seasonal component, and (c)(f) trend component.

#### **D.7** Standard Deviations

To assess the uncertainty and stability of the forecasting performance, we report the standard deviations of DMMV-S and DMMV-A on the four benchmark datasets used in  $\S4.2$  and  $\S4.3$  in Table 9. From Table 9, the relative standard deviations of the proposed models, which are calculated as

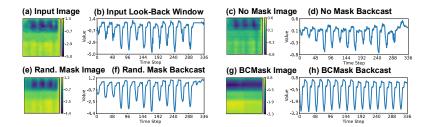


Figure 13: Comparison of different masking methods on the same example in ETTh1. (a) image of input look-back window; (c)(e)(g) are images of backcast output by DMMV-A: (c) uses "No mask"; (e) uses "Random mask"; (g) uses BCMASK. (b)(d)(f)(h) are their recovered time series, respectively.

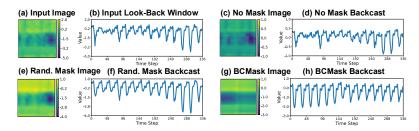


Figure 14: Comparison of different masking methods on the same example in ETTh2. (a) image of input look-back window; (c)(e)(g) are images of backcast output by DMMV-A: (c) uses "No mask"; (e) uses "Random mask"; (g) uses BCMASK. (b)(d)(f)(h) are their recovered time series, respectively.

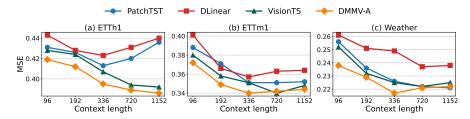


Figure 15: Average MSE comparison with varying look-back window (or context) lengths.

Table 9: Standard Deviations of DMMV-s and DMMV-A in terms of MSE and MAE on four LTSF benchmark datasets.

Mo	del	DMN	/IV-A	DMN	MV-S
Me	tric	MSE	MAE	MSE	MAE
ETTh1	96 192 336 720	$ \begin{array}{c} 0.354 \! \pm 0.001 \\ 0.393 \! \pm 0.001 \\ 0.387 \! \pm 0.001 \\ 0.447 \! \pm 0.002 \end{array} $	$\begin{array}{c} 0.390 \!\pm 0.001 \\ 0.405 \!\pm 0.001 \\ 0.413 \!\pm 0.001 \\ 0.451 \!\pm 0.001 \end{array}$	$ \begin{vmatrix} 0.350 \pm 0.001 \\ 0.399 \pm 0.002 \\ 0.401 \pm 0.002 \\ 0.472 \pm 0.001 \end{vmatrix} $	$\begin{array}{c} 0.388 \pm 0.002 \\ 0.420 \pm 0.001 \\ 0.415 \pm 0.001 \\ 0.480 \pm 0.002 \end{array}$
ETTm1	96 192 336 720	$ \begin{array}{c} 0.278 \pm 0.001 \\ 0.317 \pm 0.001 \\ 0.351 \pm 0.001 \\ 0.411 \pm 0.000 \end{array} $	$\begin{array}{c} 0.329 \!\pm 0.000 \\ 0.358 \!\pm 0.001 \\ 0.381 \!\pm 0.000 \\ 0.415 \!\pm 0.000 \end{array}$	$ \begin{vmatrix} 0.296 \pm 0.001 \\ 0.328 \pm 0.001 \\ 0.367 \pm 0.002 \\ 0.401 \pm 0.002 \end{vmatrix} $	$\begin{array}{c} 0.348 \pm 0.002 \\ 0.368 \pm 0.002 \\ 0.393 \pm 0.002 \\ 0.415 \pm 0.003 \end{array}$
Illness	24 36 48 60	1.409± 0.001 1.291± 0.002 1.499± 0.002 1.430± 0.003	$\begin{array}{c} 0.754 {\pm} \ 0.001 \\ 0.742 {\pm} \ 0.003 \\ 0.810 {\pm} \ 0.011 \\ 0.774 {\pm} \ 0.001 \end{array}$	$ \begin{vmatrix} 1.638 \pm 0.003 \\ 1.329 \pm 0.012 \\ 1.643 \pm 0.002 \\ 1.473 \pm 0.002 \end{vmatrix} $	$\begin{array}{c} 0.842 \pm 0.005 \\ 0.751 \pm 0.002 \\ 0.853 \pm 0.005 \\ 0.810 \pm 0.002 \end{array}$
Weather	96 192 336 720	$ \begin{array}{c} 0.143 \pm 0.001 \\ 0.187 \pm 0.001 \\ 0.237 \pm 0.001 \\ 0.300 \pm 0.002 \end{array} $	$\begin{array}{c} 0.196 \!\pm 0.002 \\ 0.245 \!\pm 0.003 \\ 0.272 \!\pm 0.003 \\ 0.318 \!\pm 0.003 \end{array}$	$ \begin{vmatrix} 0.168 \pm 0.001 \\ 0.221 \pm 0.002 \\ 0.267 \pm 0.002 \\ 0.323 \pm 0.001 \end{vmatrix} $	$\begin{array}{c} 0.218 \pm 0.002 \\ 0.259 \pm 0.002 \\ 0.305 \pm 0.001 \\ 0.341 \pm 0.003 \end{array}$

the ratio between standard deviation and mean, are all below 1.30% across different datasets and evaluation metrics, demonstrating their stability and robustness over different runs.