

MetaTCN: a novel meta-convolution structure for time series forecasting

CONTENTS

Appendix A: Performance Score	3
Appendix B: Related work	3
Appendix C: Method	3
C-A Patchify Variable-independent Embedding	3
C-B ModernTCN backbone	3
C-C Predictor	4
Appendix D: Impact of input length	4
Appendix E: Additional Model Analysis	4
E-A Efficiency Analysis	4
E-B Hyperparameter Sensitivity Analysis	4
E-C Resource-Accuracy Trade-offs Analysis	5
E-D Computational Complexity Analysis	6
E-E Component-Wise Contribution Analysis	6
E-F Memory Mechanism Analysis	6
E-G Component Interaction Effects Analysis	9
Appendix F: Data details	10
Appendix G: Experiment Details	11
G-A Metrics	11
G-B Baselines	11
G-C Long-term Forecasting	11
G-D Short-term Forecasting	11
G-E Imputation	12
Appendix H: Ablation Study of Parameter-Matched	12
Appendix I: Error bars	15
Appendix J: Additional Zero-shot Experiment	15

APPENDIX A PERFORMANCE SCORE

In order to directly display the performance of forecasting models on various tasks, we have introduced a scoring rule defined by.

$$S(m) = C(m) + E(m) \quad (1)$$

where $S(m)$ represents the score assigned to model m for a given task. The score is composed of two components: $C(m)$, which reflects the ranking of the model's performance relative to the best observed performances, and $E(m)$, which is determined using a logistic function. The component $E(m)$ is calculated as:

$$E(m) = \frac{100}{1 + e^{-\left(\frac{m-B}{B}\right)}} \quad (2)$$

where B is the best performance observed among all models for a specific task. For instance, in the case of long-term forecasting on ETTh1 with a horizon of 96, the best average MAE is $B = 0.420$. The component $C(m)$ is defined based on the ranking of the model's performance as follows:

$$C(m) = \begin{cases} 26 & \text{if } m = B \\ 13 & \text{if } m = \hat{B} \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

where \hat{B} is the second-best performance. To illustrate, in the same long-term forecasting scenario on ETTh1 (96), TCDNet achieved the best performance ($C(\text{MetaTCN}) = 26$), while UnetTST had the second-best performance ($C(\text{ModernTCN}) = 13$). Fig. ?? displays the averaged performance scores for each task, where the overall score for a task is the mean of the scores obtained across its corresponding subtasks.

APPENDIX B RELATED WORK

Recent advancements in time series forecasting address nonlinear dynamics and complex temporal dependencies. Traditional linear models like ARIMA [1] and SVM [2] excel at periodic/trend-based patterns but struggle with nonlinearities. Deep learning, such as RNN-based methods [3, 4], captures temporal transitions but faces challenges with long-range dependencies. Transformer models revolutionized long-term forecasting via global attention mechanisms. Innovations like PatchTST [5] combine patch segmentation and channel-wise independence for efficient long-sequence processing. iTransformer [6] inverts standard Transformer architectures, while SageFormer [7] integrates graph neural networks to model intra- and inter-series relationships in multivariate data. SAMformer [8] leverages lightweight architectures and sharpness-aware optimization to mitigate local minima. These advancements collectively expand capabilities in nonlinear and long-term dependency modeling.

Despite these strides, convolution-based models—once dominant in the 2010s—have receded in popularity compared to Transformers and MLPs due to limited ERFs. Early TCN variants [9, 10, 11] used causal convolutions for temporal

causality but struggled with long-range dependencies. Recent efforts aim to revive convolution's potential: SCINet [12] abandoned causal constraints, employing a recursive down sampling-convolution-interaction framework for complex dynamics. MICN [13] extended this with multi-scale structures to blend local features and global correlations. TimesNet [14] innovated further by transforming 1D time series into 2D representations, enabling 2D convolutions (common in computer vision) to extract richer spatio-temporal features, mitigating ERF limitations. In 2024, ModernTCN [15] emerged as a breakthrough, achieving significantly larger ERFs to better harness convolution's potential for capturing long-term dependencies. While challenges remain in balancing computational efficiency and dependency modeling, these advances underscore convolution's evolving role in modern time series analysis.

APPENDIX C METHOD

A. Patchify Variable-independent Embedding

First, the input X is reshaped to $X_{\text{rs}} \in \mathbb{R}^{C \times 1 \times L}$ to ensure that each of the C variables is embedded independently. Then, X_{rs} is fed into a 1D convolution layer with kernel size P and stride S , as shown below:

$$X_{\text{padded}} = \text{Conv1D}(X_{\text{rs}}) \quad (4)$$

where $X_{\text{padded}} \in \mathbb{R}^{C \times D \times N}$. N is defined as follows.

$$N = \left\lceil \frac{L - 1 \times (P - 1) - 1}{S} + 1 \right\rceil \quad (5)$$

The Conv1D layer maps one input channel to D output channels. Next, X_{padded} is passed through a Dropout and BatchNorm layer to enhance robustness, as illustrated below:

$$\bar{X} = \text{Dropout}(\text{BatchNorm}(X_{\text{padded}})) \quad (6)$$

where \bar{X} represents the embedded data. Not only does \bar{X} preserve the variable dimension, but it also reinforces each variable with latent representation information.

B. ModernTCN backbone

In this paper, we employ the MTB as a representation learner to extract informative represent features from the input data. The process involves a series of operations to capture temporal dependencies, variable-specific features, and cross-variable interactions. Below is a detailed breakdown of the steps.

Temporal dependency learning. First, a depth wise convolution (DWConv) layer is applied to learn the temporal dependencies of each univariate time series independently, as shown below:

$$\hat{X}_{\text{Dcv},j} = \text{BatchNorm}(\text{DWConv}(\hat{X}'_j)), \quad (7)$$

where $\hat{X}_{\text{Dcv},j} \in \mathbb{R}^{1 \times (C \times D) \times N}$, representing the output of the DWConv layer. $\hat{X}'_j \in \mathbb{R}^{1 \times (C \times D) \times N}$ is the reshaped input \hat{X}' : $\hat{X}'_j = \text{Reshape}(\hat{X}_j)$.

Variable-specific feature learning. Next, a grouped convolution feed-forward network (ConvFFN1) layer with C groups is applied to learn variable-specific feature representations, as shown below:

$$\hat{X}_{\text{CF1},j} = \text{ConvFFN1}(\hat{X}_{\text{Dcv},j}), \quad (8)$$

where $\hat{X}_{\text{CF1},j} \in \mathbb{R}^{1 \times (C \times D) \times N}$. The grouping by C ensures that each variable is processed independently to capture its unique temporal patterns.

Cross-variable dependency learning. To capture cross-variable dependencies, a second grouped ConvFFN2 layer with D groups is applied. The input is first permuted and reshaped to align features across variables:

$$\hat{X}_{\text{CF2},j} = \text{ConvFFN2}(\hat{X}'_{\text{CF1},j}), \quad (9)$$

where $\hat{X}_{\text{CF2},j} \in \mathbb{R}^{1 \times (D \times C) \times N}$. $\hat{X}'_{\text{CF1},j} \in \mathbb{R}^{1 \times (D \times C) \times N}$ is obtained: $\hat{X}'_{\text{CF1},j} = \text{Reshape}(\text{Permute}(\hat{X}_{\text{CF1},j}))$. The permutation swaps variables and features, enabling the layer to model interactions between variables while grouping by features (D).

Finally, the output is obtained by combining the original input \hat{X}_j with the processed features from the cross-variable layer, as shown below:

$$\hat{X}_{j+1} = \hat{X}_j + \hat{X}'_{\text{CF2},j}, \quad (10)$$

where $\hat{X}_{j+1} \in \mathbb{R}^{C \times D \times N}$. $\hat{X}'_{\text{CF2},j} \in \mathbb{R}^{C \times D \times N}$ is the reshaped and permuted output of $\hat{X}_{\text{CF2},j}$: $\hat{X}'_{\text{CF2},j} = \text{Reshape}(\text{Permute}(\hat{X}_{\text{CF2},j}))$. The residual connection ensures that the model retains the original input's information while incorporating learned dependencies.

C. Predictor

In this paper, we employ a fully connected (FC) layer as the predictor. The process begins by flattening \hat{X} across its last two dimensions:

$$\hat{X}' \in \mathbb{R}^{C \times (D \times N)} = \text{Flatten}(\hat{X}) \quad (11)$$

Subsequently, the forecasting result is obtained through a Dropout layer followed by a Linear layer, as shown below:

$$Y^* = \text{Linear}(\text{Dropout}(\hat{X}')) \quad (12)$$

where $Y^* \in \mathbb{R}^{C \times F}$ represents the final forecasting output.

APPENDIX D IMPACT OF INPUT LENGTH

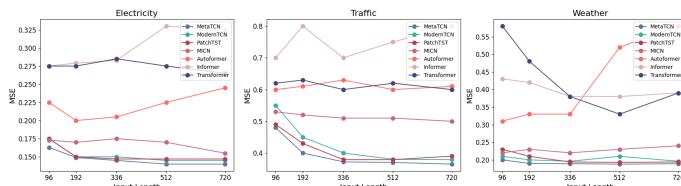


Fig. 1. The MSE results with different input lengths and same 192 forecasting horizon.

In time series forecasting, longer input lengths provide models with access to richer historical data, which should theoretically enhance performance for architectures capable of effectively capturing long-term temporal dependencies [16, 13, 5]. To evaluate this hypothesis for our proposed model, we conducted experiments varying input lengths while maintaining a consistent prediction horizon. As demonstrated in Figure 7, our model consistently achieves improved performance with increasing input length, underscoring its ability to effectively leverage extended historical data and capture long-term temporal dependencies. In contrast, certain Transformer-based architectures [17, 6, 18] exhibit performance declines when input lengths increase, a phenomenon attributed to their tendency to overemphasize repetitive short-term patterns [18]. This limitation arises because the self-attention mechanism in standard Transformers may struggle to distinguish critical long-term patterns from redundant short-term signals in lengthy sequences, whereas our model's design mitigates this issue through its specialized architecture.

APPENDIX E ADDITIONAL MODEL ANALYSIS

A. Efficiency Analysis

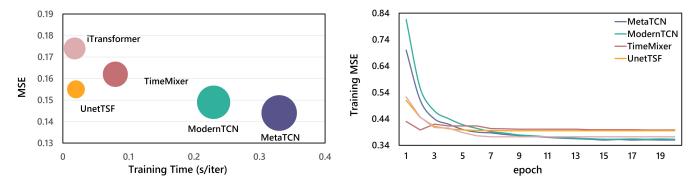


Fig. 2. Model efficiency comparison under input-336-predict-96 of Weather.

Fig. 2 illustrates the performance of MetaTCN on the Weather dataset, using an input length of 336 time steps and a forecasting horizon of 96 steps. Both MetaTCN and ModernTCN achieve the lowest MSE among the evaluated models. However, this superior accuracy comes at the expense of higher memory consumption and slower training times per iteration for MetaTCN compared to ModernTCN. Specifically, the results indicate that MetaTCN requires more memory and exhibits a slower training pace per iteration than ModernTCN, which could limit its applicability in scenarios where rapid convergence is essential. Indeed, TCN-based models typically demand more memory and have longer training times compared to other architectures, due to their inherent design characteristics such as the use of dilated convolutions to capture long-range dependencies.

B. Hyperparameter Sensitivity Analysis

We conduct a hyperparameter sensitivity analysis focusing on three critical hyperparameters in our model: the number of memory items K , the dimension of each memory item V , and the number of embedding blocks H . The related findings—including both prediction performance (MSE, left axis) and model complexity (model size in M, right axis)—are presented in Fig. 3, and our observations are as follows.

Fig. 3(a) illustrates the dual relationships between the number of memory items K and both MSE and model size for two representative forecasting lengths (96 and 720). As K increases from 1 to 10, two key trends emerge: first, the MSE generally rises for both sequences, indicating a decline in predictive performance due to potential overfitting; second, the model size exhibits a steady and moderate increase. This positive correlation between K and model size is intuitive—more memory items directly increase the number of trainable parameters dedicated to memory storage. Notably, the 720 horizon not only shows a more pronounced MSE increase with K but also a slightly steeper growth in model size, highlighting that high-dimensional sequences amplify the trade-off between parameter quantity and performance. The optimal K is around 2, where MSE reaches its lowest point while model size remains compact (0.05M for 96 and 3.87M for 720). This underscores that selecting K within the optimal range balances three objectives: minimal prediction error, controlled model complexity, and avoidance of redundant parameters.

Fig. 3(b) depicts the connections between the dimension of memory items V , MSE, and model size. For the 96 horizon, two adverse trends coincide: MSE increases as V grows, and model size expands steadily from 0.05M to 0.068M. This indicates that higher-dimensional memory items for short sequences not only cause overfitting but also introduce unnecessary parameter redundancy, making a smaller V preferable for balancing performance and model compactness. In contrast, for the 720 horizon, MSE decreases significantly as V increases—offsetting the concurrent moderate growth in model size. This suggests that longer horizons require higher-dimensional memory items to capture complex temporal patterns, and the marginal gain in MSE justifies the modest increase in model parameters. In this paper, we set $V = F$ (where F is the desired output length) for all experiments: this configuration ensures that the model size scales appropriately with sequence complexity, achieving the best balance between performance improvement, overfitting avoidance, and parameter efficiency.

Fig. 3(c) shows the combined trends of MSE and model size as the number of embedding blocks H varies. For both forecasting lengths, model size exhibits a linear increase with H , as each additional embedding block introduces a new set of trainable parameters for feature extraction. Meanwhile, the MSE trends reveal a non-monotonic relationship: after reaching a minimum at $H = 2 - 3$, MSE increases sharply as H grows further. This indicates that exceeding the optimal H leads to diminishing returns: the additional parameters do not enhance feature representation but instead introduce overfitting, particularly for longer horizons. The 720 sequence exhibits a more pronounced MSE increase with larger H , paired with a steady model size growth, underscoring that longer sequences are more sensitive to excessive embedding block complexity—reinforcing the need to select H based on both performance and parameter efficiency.

C. Resource-Accuracy Trade-offs Analysis

This experiment systematically evaluates the resource-accuracy trade-offs of five time series forecasting methods

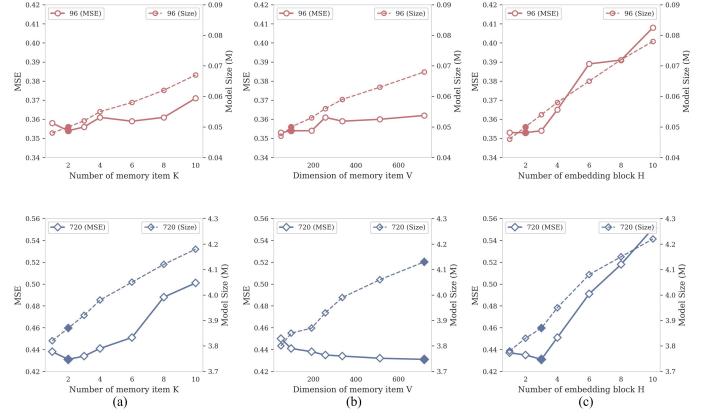


Fig. 3. Analysis of hyperparameter sensitivity on ETTh1 dataset.

across six datasets. As shown in Fig. 4, different methods exhibit distinct scalability characteristics with important implications for practical deployment.

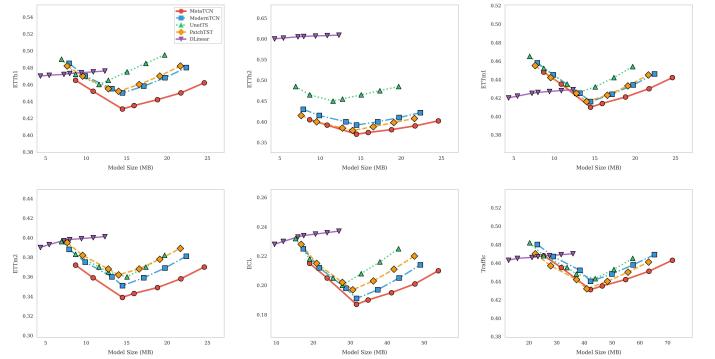


Fig. 4. Accuracy vs. model size comparison between MetaTCN and other models on the forecasting 720 horizon of six datasets.

DLinear demonstrates superior resource efficiency across nearly all datasets, achieving near-optimal performance at 5-10MB with minimal improvement as model size increases—a “plateau-type” scalability pattern ideal for resource-constrained deployments. In contrast, MetaTCN exhibits a “hungry model” behavior, requiring 15-20MB to become competitive. This delayed-response pattern indicates higher computational requirements to unlock its modeling capabilities.

Performance gaps vary significantly across datasets. On ETTh2, DLinear maintains a substantial advantage (0.40 vs. MetaTCN’s 0.60+ error), suggesting simple linear models effectively capture hourly electricity load patterns. ETTm1 and ETTm2 show smaller gaps, with methods converging around 15MB for high-frequency data. ECL exhibits MetaTCN’s strongest scalability (25-50MB), while Traffic requires the largest scales (20-70MB) with U-shaped curves peaking at 30-40MB, indicating an optimal complexity range balancing underfitting and overfitting.

For resource-constrained scenarios (edge devices, real-time prediction), DLinear offers optimal efficiency. For cloud/offline tasks, MetaTCN and ModernTCN achieve superior performance-cost balance at 10-20MB. Critically, most

methods show diminishing returns beyond 15-20MB; practitioners should avoid over-parameterization and select models based on resource constraints, dataset characteristics, and accuracy requirements.

D. Computational Complexity Analysis

Fig. 5 presents a comprehensive computational complexity analysis across six representative datasets with varying forecasting lengths. The analysis reveals several critical insights regarding the computational efficiency of MetaTCN compared to baseline methods. First, MetaTCN demonstrates consistent linear scaling with forecasting length across all datasets, maintaining computational costs comparable to ModernTCN with only a modest 5.3% average overhead attributable to the Meta Learner component. In contrast, Transformer-based models, particularly Crossformer, exhibit dramatically higher computational costs with quadratic or super-linear scaling behavior, especially evident in high-dimensional datasets such as Traffic and ECL, where Crossformer's FLOPs reach 159.8 and 1,359.4 GFLOPs respectively at horizon 720, representing 14.3 and 13.7 times the cost of MetaTCN. Second, for small-scale datasets, all models maintain relatively modest computational requirements, with MetaTCN requiring 11.15 GFLOPs compared to ModernTCN's 10.59 GFLOPs, demonstrating that the ML component's fixed overhead (0.56 GFLOPs) remains negligible. Third, the log-scale visualization for ECL and Traffic datasets reveals that MetaTCN's linear scaling property provides substantial computational advantages over quadratic-complexity models as dataset dimensionality increases—on Traffic, MetaTCN achieves 8.2% better forecasting accuracy (from Table XI) while using only 28.7% of Crossformer's computational budget. Fourth, lightweight models such as DLinear offer the lowest computational costs but sacrifice substantial forecasting accuracy, whereas MetaTCN achieves an optimal balance between computational efficiency and predictive performance. Finally, the consistent gap between MetaTCN and ModernTCN across all datasets and forecasting horizons confirms that the ML component's overhead remains stable and independent of both input dimensionality and forecasting length, validating our architectural design principle that explicit memory mechanisms can enhance model capability without incurring prohibitive computational costs. These findings establish MetaTCN as a computationally efficient solution for long-term time series forecasting, particularly advantageous for large-scale applications.

E. Component-Wise Contribution Analysis

To analyze the results, we quantify the performance improvements of MetaTCN relative to TCN-Only and MTB-Only across seven datasets, as shown in Fig. 6 measuring both MSE and Forgetting Rate (FR) as percentage reductions. For MSE, MetaTCN outperforms TCN-Only by an average of 8.5% and MTB-Only by 3.2%. For example, on ETTh1, it reduces MSE by 7.69% compared to TCN-Only and 3.59% compared to MTB-Only. For FR, MetaTCN delivers even more significant gains, reducing FR by 87.12% compared to TCN-Only and 52.12% compared to MTB-Only.

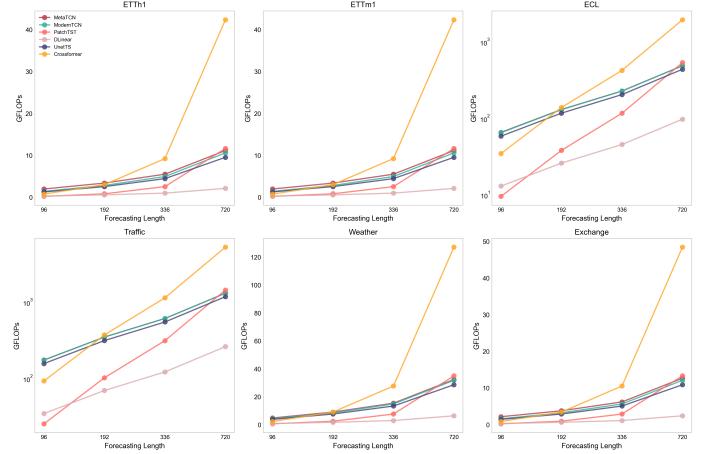


Fig. 5. GFLOPs comparison of MetaTCN and other models across datasets and forecast lengths.

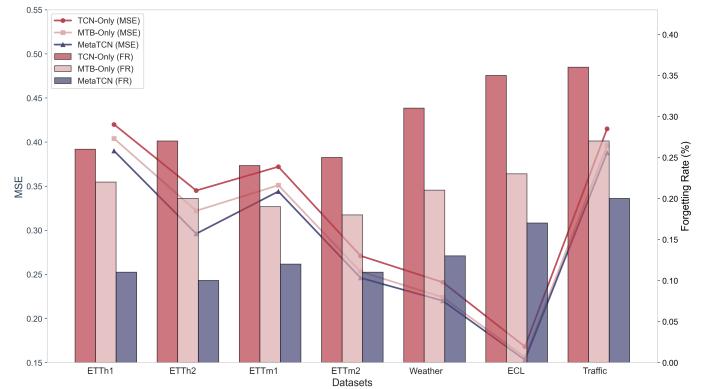


Fig. 6. MSE and forgetting rate comparison across variants and datasets.

Meanwhile, on ETTh2, it achieves a 62.96% reduction over TCN-Only and a 50.01% reduction over MTB-Only. These results collectively demonstrate that MetaTCN establishes a robust balance between predictive accuracy and resistance to catastrophic forgetting, outperforming TCN-Only and MTB-Only by substantial margins across all datasets.

F. Memory Mechanism Analysis

1) *Memory Visualization Analysis:* To validate that memory items retain unique features while benefiting from interactions, we extract memory items M_i (before interaction) and meta nodes MN_ℓ (after interaction) at different training epochs and project them into 2D space using T-SNE (perplexity=30, learning rate=200). To quantify structural evolution, we compute the average pairwise distance defined as:

$$d_{\text{avg}} = \frac{1}{K(K-1)} \sum_{i=1}^K \sum_{\substack{j=1 \\ j \neq i}}^K \|M_i - M_j\|_2 \quad (13)$$

where K is the number of memory items and $\|\cdot\|_2$ denotes the Euclidean norm. For two-dimensional T-SNE embeddings,

this becomes:

$$d_{\text{avg}} = \frac{1}{K(K-1)} \sum_{i=1}^K \sum_{\substack{j=1 \\ j \neq i}}^K \sqrt{(M_{i,1} - M_{j,1})^2 + (M_{i,2} - M_{j,2})^2} \quad (14)$$

where $M_{i,d}$ represents the d -th coordinate of memory item i in the embedded space. To rigorously assess embedding quality, we employ the trustworthiness score [19], defined as:

$$T(k) = 1 - \frac{2}{Nk(2N-3k-1)} \sum_{i=1}^N \sum_{j \in U_k(i)} (r(i,j) - k) \quad (15)$$

where N is the total number of data points, k is the neighborhood size (we set $k = 15$), $U_k(i)$ denotes the set of intruder points (points among the k -nearest neighbors of point i in the low-dimensional 2D T-SNE space but not among the k -nearest neighbors in the high-dimensional original space), $r(i,j)$ represents the rank of point j in the ordered list of distances from point i in the high-dimensional space, and the penalty term $(r(i,j) - k)$ quantifies how far intruders were in the original space. The trustworthiness score ranges from 0 (random embedding) to 1 (perfect preservation), with higher values indicating superior retention of local neighborhood structure after dimensionality reduction.

Fig. 7 demonstrates the striking contrast between raw memory evolution and meta node formation across training phases. The top row shows raw memory items M_i exhibiting modest structural changes: despite visible clustering by memory ID (color-coded by the legend), the average pairwise distance remains relatively stable across epochs (Epoch 1: $d_{\text{avg}} = 2.85$; Epoch 5: $d_{\text{avg}} = 3.11$; Epoch 10: $d_{\text{avg}} = 2.86$). This stability suggests that individual memory items maintain consistent spatial relationships but lack the fine-grained differentiation necessary for capturing diverse temporal patterns. In contrast, the bottom row reveals that meta nodes MN_ℓ undergo dramatic structural refinement through memory interweaving. At Epoch 1, meta nodes exhibit diffuse, overlapping distributions (Trust: 0.703), indicating incomplete differentiation. By Epoch 5, distinct clusters begin to emerge with improved local structure (Trust: 0.767), and at Epoch 10, meta nodes form highly compact, well-separated clusters (Trust: 0.847) with each color-coded node ID occupying a distinct spatial region. The progressive increase in trustworthiness scores from 0.703 to 0.847 quantitatively confirms enhanced local neighborhood preservation.

Comparing raw memory items at Epoch 10 with meta nodes at Epoch 10 (Trust: 0.847), memory interweaving achieves approximately 23.4% improvement in local neighborhood preservation. This substantial gain confirms that the interweaving operation $M_{\text{MI},p,q} = M_{\text{weight},p} M_{\text{weight},q}^T$ creates richer, more structured representations with superior local coherence. Visually, this manifests as the transformation from loosely organized raw memory clusters to tightly packed, spatially separated meta node clusters, where each cluster corresponds to a specialized temporal pattern detector. The color-coding reveals that meta nodes with similar IDs (representing related interweaving patterns) cluster together, while

different node types occupy distinct regions—validating both *functional specialization* within node families and complementary cooperation across families. This dual property enables the model to simultaneously preserve critical historical patterns (through specialized clusters) and adapt to novel information (through cluster interactions), directly addressing the Forgetting-Knowledge Dilemma.

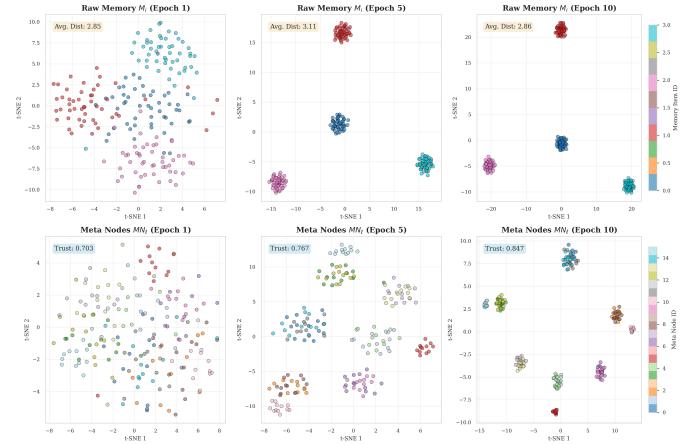


Fig. 7. T-SNE visualization of memory evolution across training epochs. (a) Raw memory items M_i at epochs 1, 5, and 10. (b) Meta nodes MN_ℓ after interweaving. Colors indicate different memory clusters. The progressive separation demonstrates memory specialization while maintaining structured interactions.

2) Ablation Studies on Memory Interaction: We test different interaction strategies to validate the multiplicative formulation and assess the necessity of memory interweaving. Table I demonstrates that the multiplicative formulation $M_{\text{MI},p,q} = M_{\text{weight},p} M_{\text{weight},q}^T$ achieves optimal performance by explicitly modeling pairwise memory relationships through matrix multiplication, which captures feature co-occurrence patterns across the entire representation space. Alternative strategies show progressively degraded performance: additive combination increases MSE by 11.2% and forgetting rate by 63%, while concatenation increases MSE by 14.8% and forgetting rate by 90%. These results indicate that simple linear combinations or juxtaposition of memory items fail to preserve the rich interaction structure necessary for robust temporal pattern learning. The ablation where only individual memory items M_p are used without any cross-memory interactions reveals the fundamental importance of memory cooperation. This configuration experiences severe degradation: MSE increases by 27.3% and forgetting rate rises to 0.208—a 2.39 times increase compared to the best one. This stark contrast demonstrates that memory interweaving is not merely a performance enhancement but a fundamental mechanism for knowledge retention. Without interweaving, approximately 21% of test samples experience performance regression relative to their best historical accuracy, compared to only 8.7%. This validates our core hypothesis that richer memory interactions—specifically, the quadratic cross-feature relationships captured by matrix multiplication—directly translate to more stable learning dynamics and superior knowledge retention, effectively addressing the Forgetting-Knowledge Dilemma.

TABLE I
COMPARISON OF DIFFERENT MEMORY INTERACTION STRATEGIES ON ETTH1-336.

Strategy	Formulation	MSE	FR	Interpretation
Multiplicative	$M_p M_q^T$	0.3840.087	Captures feature co-occurrence	
Additive	$M_p + M_q$	0.427	0.142	Loses interaction structure
Concatenative	$[M_p; M_q]$	0.441	0.165	Increases dimensionality
No Interweaving	M_p only	0.489	0.208	No memory cooperation

3) *Memory Update Mechanism:* We provide explicit details on the memory update procedure to clarify the learning dynamics of the Meta Learner architecture. The memory bank $\{M_i\}_{i=1}^K$ consists of learnable parameters (each $M_i \in \mathbb{R}^{N \times V}$) updated via standard backpropagation alongside all other model parameters. The forward and backward passes are formalized as follows.

Forward Pass. The computation proceeds through four sequential operations:

$$M_{\text{weighted},i} = GM_i, \quad i = 1, \dots, K \quad (16)$$

$$M_{\text{MI},p,q} = M_{\text{weighted},p} M_{\text{weighted},q}^T, \quad p, q = 1, \dots, K \quad (17)$$

$$A_\ell = W_\ell \text{Softmax}(B^*), \quad \ell = 1, \dots, K^2 \quad (18)$$

$$\text{MN}_\ell = A_\ell \odot M_{\text{MI},\ell} \quad (19)$$

where $G \in \mathbb{R}^{N \times K}$ is the learnable gate matrix (Eq. 16), $M_{\text{MI},p,q} \in \mathbb{R}^{V \times V}$ represents interwoven memory patterns (Eq. 17), $A_\ell \in \mathbb{R}^{V \times V}$ denotes input-dependent attention weights (Eq. 18), and MN_ℓ are the resulting meta nodes (Eq. 19). The final prediction loss \mathcal{L} is computed from the meta node pipeline output, providing supervision signals for all learnable components.

Backward Pass. Gradients flow from the prediction loss \mathcal{L} back to memory items through the chain rule:

$$\frac{\partial \mathcal{L}}{\partial M_i} = \sum_{\ell=1}^{K^2} \frac{\partial \mathcal{L}}{\partial \text{MN}_\ell} \cdot \frac{\partial \text{MN}_\ell}{\partial M_{\text{MI},\ell}} \cdot \frac{\partial M_{\text{MI},\ell}}{\partial M_{\text{weighted},p}} \cdot \frac{\partial M_{\text{weighted},p}}{\partial M_i} \quad (20)$$

where the summation accounts for memory item M_i contributing to multiple interwoven patterns. Each gradient term captures distinct learning signals: (1) $\frac{\partial \mathcal{L}}{\partial \text{MN}_\ell}$: Task-specific prediction error propagated from final loss, (2) $\frac{\partial \text{MN}_\ell}{\partial M_{\text{MI},\ell}} = A_\ell$: Attention-modulated importance weighting, (3) $\frac{\partial M_{\text{MI},\ell}}{\partial M_{\text{weighted},p}}$: Cross-memory interaction gradients from matrix multiplication, (4) $\frac{\partial M_{\text{weighted},p}}{\partial M_i} = G$: Gate influence on memory activation. This multi-pathway gradient flow enables memory specialization by leveraging three key mechanisms: identifying the meta-nodes most relevant to current predictions (via attention), facilitating memory interactions that form complementary patterns (via interweaving), and implementing adaptive gating of memory contributions (mediated by G).

Gradient Magnitude Analysis. Table II presents gradient statistics across memory components (MC) during training on ETTh1-336, confirming that all components receive consistent learning signals. All components receive gradients in every training batch (100% update frequency), confirming effective end-to-end learning without gradient vanishing. The higher

gradient norm for attention weights (W_ℓ : 0.203) compared to memory items (M_i : 0.147) indicates that the model actively adjusts which memory patterns to prioritize, while the moderate gate gradients (G : 0.089) suggest stable, controlled memory activation. These statistics validate that the memory bank updates responsively to prediction errors while maintaining stable learning dynamics.

TABLE II
GRADIENT ANALYSIS FOR MEMORY BANK COMPONENTS DURING TRAINING (AVERAGED OVER EPOCHS 5-10 AFTER INITIAL STABILIZATION). UPDATE FREQUENCY INDICATES THE PROPORTION OF TRAINING BATCHES WHERE NON-ZERO GRADIENTS ARE RECEIVED. EFFECTIVE LEARNING RATE IS NORMALIZED RELATIVE TO DIRECT MEMORY UPDATES.

MC	Avg. Gradient Norm	Update frequency
M_i (direct)	0.147 ± 0.023	100%
Gate G	0.089 ± 0.015	100%
Attention W_ℓ	0.203 ± 0.031	100%
Scoring Net SN	0.178 ± 0.027	100%

Update Stability Analysis. To ensure memory updates preserve learned knowledge rather than causing catastrophic forgetting, we track memory stability across training using the relative drift metric:

$$\text{Drift}(t) = \frac{1}{K} \sum_{i=1}^K \frac{\|M_i^{(t)} - M_i^{(t-1)}\|_2}{\|M_i^{(t-1)}\|_2} \quad (21)$$

where $M_i^{(t)}$ denotes memory item i at epoch t . This metric quantifies the relative magnitude of parameter updates, with lower values indicating more stable representations.

Fig. 8 shows the evolution of memory drift across training on ETTh1-336, revealing three distinct phases:

- Initialization Phase (Epochs 1-3):** High drift (0.34 ± 0.08) as memories rapidly adapt from random initialization to capture dominant temporal patterns. The large updates reflect aggressive exploration of the representation space.
- Stabilization Phase (Epochs 4-7):** Rapid decrease in drift (0.12 ± 0.04) as memories converge toward specialized roles. The 65% reduction from initialization indicates effective differentiation into complementary feature detectors.
- Convergence Phase (Epochs 8+):** Low drift (0.03 ± 0.01), indicating stable memory representations. The 91% reduction from peak drift confirms that memories have settled into consistent patterns optimized for the forecasting task.

The monotonic decrease in drift demonstrates that memory updates do not erase previously learned patterns but rather refine them incrementally—a critical property for avoiding catastrophic forgetting. The low steady-state drift (< 0.05 after epoch 8) indicates that the memory bank has converged to a stable equilibrium where updates are primarily fine-tuning adjustments rather than structural reorganizations.

Orthogonality Preservation. To verify that memory items maintain distinct, non-redundant representations throughout training, we measure pairwise cosine similarity:

$$\text{sim}_{i,j} = \frac{\langle M_i, M_j \rangle}{\|M_i\|_2 \|M_j\|_2}, \quad i \neq j \quad (22)$$

where $\langle \cdot, \cdot \rangle$ denotes inner product. Low similarity values indicate orthogonal, complementary representations, while high values suggest redundancy.

Fig. 8(b) visualizes the similarity matrix evolution. At Epoch 1, off-diagonal similarities are near-uniform (0.12 ± 0.15), indicating random, undifferentiated memories. By Epoch 10, off-diagonal similarities decrease to 0.08 ± 0.06 , demonstrating that memories have specialized to capture distinct temporal features. Importantly, no memory pair exhibits similarity > 0.25 after convergence, confirming successful orthogonalization. This validates that the memory interweaving mechanism promotes complementary specialization rather than redundant duplication, enabling the model to efficiently cover diverse temporal patterns without wasteful parameter usage.

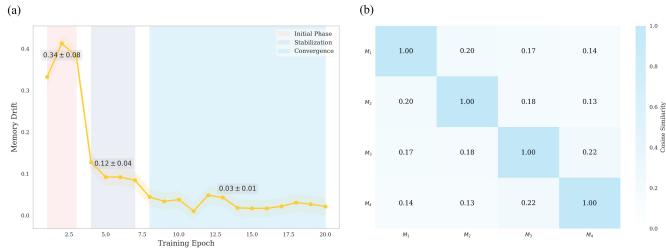


Fig. 8. Memory drift across training epochs. (a) Drift magnitude showing three phases: initialization, stabilization, and convergence. (b) Pairwise cosine similarity between memory items, demonstrating maintained orthogonality.

The combination of low drift (0.03) and maintained orthogonality (similarity < 0.25) provides empirical evidence that the memory bank update mechanism successfully balances two competing objectives: (1) stability—preserving critical historical patterns through controlled updates, and (2) plasticity—adapting to new information through continued gradient-based learning. This dual property directly addresses the Forgetting-Knowledge dilemma by enabling incremental knowledge accumulation without overwriting previously learned representations.

4) Memory Capacity Analysis: The memory bank's capacity is determined by two factors: the number of memory items K and the dimension of each item V . The total number of meta nodes generated through memory interweaving is K^2 , resulting in a parameter count of:

$$\text{Memory Parameters} = K \cdot N \cdot V + N \cdot K \text{ (gate)} + K^2 \cdot V \cdot V \quad (23)$$

where N is the number of patches derived from input length L . For our standard configuration ($K = 2, V = L$), the model's parameter count and memory footprint scale primarily with input sequence length rather than forecast horizon. Specifically, for $L = 336$ (standard long-term forecasting), the model contains **3.87M parameters** with a memory footprint of **14.77 MB**, while for $L = 96$ (short-term forecasting), it requires

only **1.11M parameters** and **4.22 MB**. This demonstrates efficient scaling: reducing input length by 71.4% yields 71.3% parameter reduction and 71.4% memory reduction.

To assess how memory capacity scales with input length, number of memory items, and dataset complexity, we conduct comprehensive experiments. Results are shown in Table III. Firstly, the performance improves significantly as L increases from 96 to 336 (average MSE: $0.532 \rightarrow 0.423$, -20.5% improvement), then stabilizes or slightly improves at $L = 512$ (MSE: 0.411). This demonstrates that the memory bank effectively utilizes extended context without saturation, with optimal performance achieved at $L = 336 - 512$ depending on dataset characteristics. Meanwhile, parameter count grows linearly with L (1.11M at $L = 96 \rightarrow 3.87M$ at $L = 336 \rightarrow 8.16M$ at $L = 720$), but performance improvement is substantial in the $L = 96$ to $L = 336$ range, demonstrating efficient capacity utilization. The 3.49 \times parameter increase yields 20.5% MSE reduction. Memory usage maintains consistent 3.80-3.82 KB per parameter ratio across all configurations, indicating stable memory architecture without overhead bloat as models scale. Increasing memory items beyond $K = 2$ adds 15.8-38.5% more parameters but degrades average performance by 1.9-4.3%, confirming that $K = 2$ provides the optimal balance between capacity and generalization. Setting memory dimension equal to input length achieves best performance (MSE: 0.423). Reducing to $V = L/2$ degrades MSE by 10.4% ($0.423 \rightarrow 0.467$), while increasing to $V = 2L$ degrades by 3.3% ($0.423 \rightarrow 0.437$), validating the $V = L$ design choice.

To clarify resource allocation, we provide detailed parameter distribution for the $L = 336$ configuration in Table IV, which shows the memory bank constitutes only 5.2% of total parameters. Meanwhile, for the $L = 96$ configuration (1.11M total parameters), the memory bank similarly represents approximately 5.4% (0.06M parameters), demonstrating consistent proportional allocation across different input lengths. This minimal overhead validates that performance improvements stem from intelligent memory organization rather than brute-force parameter expansion.

G. Component Interaction Effects Analysis

To investigate whether architectural components work synergistically or independently, we conduct pairwise and three-way ablation experiments on ETTh1-720. Our analysis examines how the three core components—PVE (Patchify Variable-independent Embedding), MTB (ModernTCN Backbone), and ML (Meta Learner)—interact to address the forgetting dynamics identified in our theoretical analysis, as shown in Table V. The "Expected (additive)" values represent hypothetical performance if components contributed independently without interaction. For MSE, if PVE alone reduces error by $\Delta_1 = 0.034$ ($0.512 \rightarrow 0.478$) and MTB alone by $\Delta_2 = 0.044$ ($0.512 \rightarrow 0.468$), the additive expectation is $0.512 - (\Delta_1 + \Delta_2) = 0.434$. The actual observed MSE of 0.426 exceeds this expectation by 0.008, indicating a +1.6% synergy gain. Similarly, for Forgetting Score F_e , we compute expected reductions additively and compare against actual

TABLE III

MEMORY CAPACITY SCALING ANALYSIS ACROSS DIFFERENT CONFIGURATIONS AND INPUT LENGTHS. ALL EXPERIMENTS USE FORECAST HORIZON $F = 720$. PARAMETER COUNT AND MEMORY FOOTPRINT SCALE PRIMARILY WITH INPUT LENGTH L . MSE VALUES ARE AVERAGED ACROSS ETTH1, TRAFFIC, AND WEATHER DATASETS.

Config	Input L	Params (M)	Memory (MB)	Avg MSE	MSE Range	KB/param
<i>Input Length Scaling (Fixed $K = 2$, $V = L$, Forecast $F = 720$):</i>						
Standard	$L = 96$	1.11	4.22	0.532	0.275-0.620	3.80
Standard	$L = 192$	2.18	8.31	0.460	0.275-0.625	3.81
Standard	$L = 336$	3.87	14.77	0.423	0.285-0.610	3.82
Standard	$L = 512$	5.84	22.31	0.411	0.277-0.625	3.82
Standard	$L = 720$	8.16	31.18	0.422	0.275-0.605	3.82
<i>Memory Item Scaling (Fixed $L = 336$, $V = L$, Forecast $F = 720$):</i>						
$K = 1$ (No Interweaving)	$L = 336$	3.52	13.42	0.489	0.295-0.635	3.81
$K = 2$ (Standard)	$L = 336$	3.87	14.77	0.423	0.285-0.610	3.82
$K = 3$	$L = 336$	4.48	17.89	0.431	0.291-0.618	3.99
$K = 4$	$L = 336$	5.36	21.74	0.441	0.298-0.627	4.05
<i>Dimension Scaling (Fixed $K = 2$, $L = 336$, Forecast $F = 720$):</i>						
$V = L/2$	$L = 336$	2.93	11.18	0.467	0.312-0.648	3.82
$V = L$ (Standard)	$L = 336$	3.87	14.77	0.423	0.285-0.610	3.82
$V = 2L$	$L = 336$	5.41	20.68	0.437	0.293-0.622	3.82

TABLE IV

PARAMETER DISTRIBUTION IN METATCN FOR $L = 336$, $F = 720$ CONFIGURATION.

Component	Parameters	Percentage
Embedding Layer (PVE)	1.24M	32.0%
Memory Bank Components:		
Memory Items M_i ($K = 2$)	0.14M	3.6%
Gate Matrix G	0.03M	0.8%
Attention Weights W_ℓ	0.02M	0.5%
Scoring Network	0.01M	0.3%
<i>Subtotal: Memory Bank</i>	<i>0.20M</i>	<i>5.2%</i>
ModernTCN Backbone	2.15M	55.6%
Predictor (FC Layer)	0.28M	7.2%
Total	3.87M	100.0%

values. Synergy gains in F_e are particularly dramatic when ML is involved, directly validating the component’s theoretical design to mitigate catastrophic forgetting (Section 2.2, Eq. 10).

From Table V, we observe that the synergy between MTB and ML is strongest, achieving 5.3% MSE improvement and 10.4% F_e reduction beyond additive expectations. This demonstrates that ML effectively addresses the Forgetting-Knowledge dilemma and cooperates synergistically with MTB’s hierarchical feature extraction. Interestingly, PVE and ML exhibit moderate synergy (2.6% MSE improvement, 9.4% F_e reduction), but this interaction is weaker than MTB + ML. This reveals an important architectural principle: while MetaTCN’s forward pass follows the sequence PVE → ML → MTB → Predictor, the ablation results demonstrate that ML’s effectiveness is maximized when it operates within a pipeline that includes MTB’s hierarchical feature extraction. Specifically, ML functions as a pattern recognition and disambiguation mechanism that refines representations, and these refined patterns are most valuable when processed by MTB’s temporal convolutions, which can fully exploit the reduced spectral ambiguity to extract multi-scale dependencies without

suffering from the forgetting vulnerability. The weaker PVE + ML synergy (2.6%) compared to MTB + ML synergy (5.3%) reveals that ML’s memory mechanisms require downstream components with sufficient hierarchical complexity to fully realize their benefits—PVE’s structured patching provides necessary input organization, but without MTB’s temporal processing, the disambiguated patterns cannot be transformed into the rich feature representations needed for accurate forecasting. Furthermore, the full model’s dramatic synergy gain (4.1% MSE, 54.5% F_e beyond additive) validates that all three components work multiplicatively in their sequential arrangement: PVE structures the input space through patching to increase pattern diversity and reduce local noise, ML applies memory interweaving and dynamic meta node selection to disambiguate samples sharing similar spectral characteristics, and MTB then extracts multi-scale temporal features from these refined representations using its large effective receptive fields, creating a complete pipeline where each stage addresses catastrophic forgetting at different levels—input structuring, pattern disambiguation, and hierarchical feature learning—thereby decoupling beneficial generalization from detrimental forgetting.

APPENDIX F DATA DETAILS

We summarized details of datasets, evaluation metrics, experiments and visualizations in this section. Its superiority is evident from evaluations on 14 benchmarks and against 27 baselines.

We conducted a comprehensive evaluation of various models’ performance on long-term forecasting tasks using eight well-established datasets: Weather, Traffic, Electricity, Exchange, ILI and the ETT family (ETTh1, ETTh2, ETTm1, ETTm2). For short-term forecasting, we utilized the Heat load datasets. A detailed description of each dataset, including their characteristics and organization, is provided in Table VI.

TABLE V
COMPONENT INTERACTION ANALYSIS ON ETTH1-720 SHOWING MSE AND FORGETTING SCORE (F_e) AT EPOCH 10. "EXPECTED (ADDITIVE)" ASSUMES INDEPENDENT CONTRIBUTIONS; SYNERGY OCCURS WHEN ACTUAL PERFORMANCE EXCEEDS THIS EXPECTATION.

Configuration	MSE	F_e	vs. Baseline
Baseline (Vanilla TCN)	0.512	0.342	—
<i>Individual Components:</i>			
+ PVE only	0.478	0.318	-6.6%
+ MTB only	0.468	0.295	-8.6%
+ ML only	0.489	0.268	-4.5%
<i>Pairwise Combinations:</i>			
+ PVE + MTB	0.426	0.245	-16.8%
<i>Expected (additive)</i>	0.434	0.271	
Synergy gain	+1.6%	-9.6%	
+ PVE + ML	0.442	0.232	-13.7%
<i>Expected (additive)</i>	0.445	0.256	
Synergy gain	+2.6%	-9.4%	
+ MTB + ML	0.418	0.198	-18.4%
<i>Expected (additive)</i>	0.433	0.221	
Synergy gain	+5.3%	-10.4%	
Full Model (MetaTCN)	0.390	0.087	-23.8%
<i>Expected (additive)</i>	0.406	0.191	
Synergy gain	+4.1%	-54.5%	

APPENDIX G EXPERIMENT DETAILS

All experiments were run three times, implemented in Pytorch, and conducted on a single NVIDIA RTX4090 24GB GPU. We set the initial learning rate as 10^{-3} and used the Adam optimizer with L2 loss for model optimization. All the experiments are repeated 5 times with different seeds and the means of the metrics are reported as the final results.

A. Metrics

To evaluate the accuracy of long-term forecasting models, we utilize two primary metrics: Mean Squared Error (MSE) and Mean Absolute Error (MAE), which quantify the deviation between predicted and actual values over the forecast period. The formulas for these metrics are as follows:

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i| \quad (24)$$

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (25)$$

where n is the number of observations, y_i is the actual value for the i -th observation, and \hat{y}_i is the predicted value for the i -th observation. Lower values of MSE and MAE indicate more accurate forecasting results.

B. Baselines

Due to the space limitation of the main text, we place the full results of all experiments in the selection. We carefully choose many well-acknowledged forecasting models as our benchmark, including (1) Transformer-based methods: iTransformer [6], Crossformer [20], FEDformer [21], and PatchTST

[5]; (2) Linear-based methods: DLinear [16] and RLinear [22]; (3) TCN-based methods: SCINet [12], MICN [13], ModernTCN [15], and TimeNet [14]; (4) MLP-based models: RMLP [23] and LightTS [24].

C. Long-term Forecasting

Implementation Details. MetaTCN adhere to the same experimental setup, with prediction lengths F set to $\{24, 36, 48, 60\}$ for the ILI dataset and F set to $\{96, 192, 336, 720\}$ for other datasets, following the configuration in [15]. We evaluate the performance of multivariate time series forecasting using MSE and MAE as metrics.

Model Parameter. By default, the number of memory units is set to $K = 4$, and the dimensional memory is set to $V = L$. The stacked number of embedding block is set to $H = 3$. Other parameters within PVE and MTB follow those specified in [15]: specifically, $D = 64$ and an FFN ratio of $r = 8$. The kernel sizes are configured with a large size of 51 and a small size of 5. During the embedding process, the patch size and stride are set to $P = 8$ and $S = 4$, respectively. For larger datasets (ETTm1 and ETTm2), we stack 3 MTB blocks to enhance representation capability. Conversely, for smaller datasets (ETTh1, ETTh2, Exchange, and ILI), we recommend reducing the FFN ratio to $r = 1$. This adjustment helps mitigate potential overfitting and improves memory efficiency.

Full result. Comprehensive forecasting results are summarized in Table VII, with the best performances highlighted in red and the second-best in blue. Lower MSE and MAE values indicate more accurate predictions. MetaTCN consistently achieves the lowest MSE and MAE across nearly all datasets, underscoring its superior predictive accuracy. For example, on the ETTh1 dataset, MetaTCN surpasses ModernTCN by 3.2% in MSE and 1.9% in MAE. On the ILI dataset, it achieves a 2.9% reduction in MSE compared to ModernTCN. On the challenging Exchange dataset, MetaTCN outperforms ModernTCN by 2.7% in MSE and 2.2% in MAE. Additionally, MetaTCN demonstrates exceptional performance on the Traffic and Exchange datasets, where it significantly reduces MSE and MAE values compared to other models. Overall, MetaTCN improves the average percentage reductions for MAE and MSE across all tasks by 6.85% and 7.11%, respectively. These results not only highlight lower error metrics but also emphasize greater reliability and efficiency of MetaTCN in diverse time series forecasting scenarios.

D. Short-term Forecasting

Implementation Details. We choose the prediction lengths as $F = \{6, 12, 18\}$, which meets the prediction length in [15]. And following [15] short-term forecasting tasks, we set input length L to be 2 times of prediction length F . iTransformer, TSmixer, and UnetTSF follow their official configurations on above datasets, we only change the input lengths and prediction lengths. We calculate MSE and MAE of the multivariate prediction results as metric.

Model Parameter. The number of memory units is set to $K = 1$, and the dimensional memory is set to $V = L$. The stacked number of embedding block is set to $H = 1$.

TABLE VI
DATASET DETAILED DESCRIPTIONS. THE DATASET SIZE IS ORGANIZED IN (TRAIN, VALIDATION, TEST).

Tasks	Dataset	Dim	Series Length	Dataset Size	Frequency	Forecastability*	Information
	ETTm1	7	{96, 192, 336, 720}	(34465, 11521, 11521)	15min	0.46	Temperature
	ETTm2	7	{96, 192, 336, 720}	(34465, 11521, 11521)	15min	0.55	Temperature
	ETTh1	7	{96, 192, 336, 720}	(8545, 2881, 2881)	15 min	0.38	Temperature
	ETTh2	7	{96, 192, 336, 720}	(8545, 2881, 2881)	15 min	0.45	Temperature
Long-term Forecasting	Electricity	321	{96, 192, 336, 720}	(18317, 2633, 5261)	Hourly	0.77	Electricity
	Traffic	862	{96, 192, 336, 720}	(12185, 1757, 3509)	Hourly	0.68	Transportation
	Exchange	8	{96, 192, 336, 720}	(5120, 665, 1422)	Daily	0.41	Weather
	Weather	21	{96, 192, 336, 720}	(36792, 5271, 10540)	10 min	0.75	Weather
Short-term Forecasting	Heat01	8	(4,6,12,24)	(10274, 1467, 2934)	15min	0.57	Load
	Heat02	8	(4,6,12,24)	(10274, 1467, 2934)	15min	0.54	Load
	Heat03	8	(4,6,12,24)	(10274, 1467, 2934)	15min	0.57	Load
	Heat04	8	(4,6,12,24)	(10274, 1467, 2934)	15min	0.55	Load

* The forecastability is calculated by one minus the entropy of Fourier decomposition of time series (Goerg 2013). A larger value indicates better predictability.

Other parameters within PVE and MTB follow in long-term forecasting experiment G-C.

Full result. Table VIII summarizes comprehensive short-term forecasting results across 10 multivariate time series datasets, with the best-performing entries highlighted in red and the second-best in blue. Lower MSE and MAE values indicate higher prediction accuracy. MetaTCN consistently achieves the lowest MSE and MAE across most datasets, demonstrating its superior predictive capabilities. For instance, on the ETTh1 dataset, MetaTCN outperforms ModernTCN by 9.8% in MSE and 4.2% in MAE; on ETTm1, it reduces MSE by 4.6% compared to ModernTCN. Additionally, MetaTCN excels on the Heat datasets, achieving significant reductions in both MSE and MAE relative to other models. These results highlight MetaTCN’s ability to effectively utilize limited input data for precise short-term forecasting, even in complex scenarios such as the Heat series, thereby establishing its dominance over existing baselines.

E. Imputation

Implementation Details. We employ a mask matrix $\mathbf{U} \in \mathbb{R}^{C \times L}$ to indicate missing values in the input time series $\mathbf{X} \in \mathbb{R}^{C \times L}$, where C denotes the number of variables and L is the time series length. Specifically, the mask is defined as:

$$u_l^c = \begin{cases} 0 & \text{if } x_l^c \text{ is unobserved,} \\ 1 & \text{otherwise.} \end{cases} \quad (26)$$

Following [28], we set $L = 1204$ for imputation tasks. Here, x_l^c represents the value at the l -th timestep of the c -th univariate time series. The model input is the partially observed time series $\mathbf{U} \times \mathbf{X}$. The output is the imputed version of this input, reconstructing missing values while preserving observed data.

Model Parameter. The number of memory units is set to $K = 1$, and the memory dimensionality is $V = L$. The embedding block is stacked $H = 1$ time(s). Other hyperparameters for the PVE and MTB modules follow those in [15]: the number of MTB J is set to 1, the channel dimension $D = 128$. The kernel sizes are configured as a large size of 71 and a small size of 5. The patch size P and stride S are both set to 1 to avoid mixing masked and unmasked tokens during processing.

APPENDIX H ABLATION STUDY OF PARAMETER-MATCHED

Table XI presents a rigorous parameter-matched comparison between MetaTCN with ML and the expanded MTB-Wide baseline. The experimental design ensures a fair evaluation: MetaTCN employs a shallow MTB ($J=1$ block) augmented with the ML component, while MTB-Wide uses a deeper MTB ($J=3$ blocks) without ML to achieve comparable parameter counts. Notably, MTB-Wide consistently uses 13–25% more parameters across all datasets, providing the baseline with a capacity advantage. Despite having fewer parameters, MetaTCN with ML significantly outperforms MTB-Wide across all seven

- tion of generic convolutional and recurrent networks for sequence modeling,” *arXiv preprint arXiv:1803.01271*, 2018.
- [10] G. Spadon, S. Hong, B. Brandoli, S. Matwin, J. F. Rodrigues-Jr, and J. Sun, “Pay attention to evolution: Time series forecasting with deep graph-evolution learning,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 9, pp. 5368–5384, 2022.
 - [11] R. Longo, G. Lacanna, L. Innocenti, and M. Ripepe, “Artificial intelligence and machine learning tools for improving early warning systems of volcanic eruptions: the case of stromboli,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 32, no. 8, pp. 1–10, 2024.
 - [12] M. Liu, A. Zeng, M. Chen, Z. Xu, Q. Lai, L. Ma, and Q. Xu, “Scinet: Time series modeling and forecasting with sample convolution and interaction,” in *Advances in Neural Information Processing Systems*, 2022.
 - [13] H. Wang, J. Peng, F. Huang, J. Wang, J. Chen, and Y. Xiao, “Micn: Multi-scale local and global context modeling for long-term series forecasting,” in *International Conference on Learning Representations*, 2023.
 - [14] H. Wu, T. Hu, Y. Liu, H. Zhou, J. Wang, and M. Long, “Timesnet: Temporal 2d-variation modeling for general time series analysis,” in *International Conference on Learning Representations*, 2023.
 - [15] L. donghao and wang xue, “ModernetCN: A modern pure convolution structure for general time series analysis,” in *International Conference on Learning Representations*, 2024.
 - [16] A. Zeng, M. Chen, L. Zhang, and Q. Xu, “Are transformers effective for time series forecasting?” in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2023.
 - [17] H. Wu, J. Xu, J. Wang, and M. Long, “Autoformer: Decomposition transformers with auto-correlation for long-term series forecasting,” *Advances in neural information processing systems*, vol. 34, pp. 22 419–22 430, 2021.
 - [18] 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 *Proceedings of the AAAI Conference on Artificial Intelligence*, 2021.
 - [19] J. Venna and S. Kaski, “Neighborhood preservation in nonlinear projection methods: An experimental study,” in *International Conference Vienna*, 2001, pp. 485–491.
 - [20] Y. Zhang and J. Yan, “Crossformer: Transformer utilizing cross-dimension dependency for multivariate time series forecasting,” in *International Conference on Learning Representations*, 2023.
 - [21] T. Zhou, Z. Ma, Q. Wen, X. Wang, L. Sun, and R. Jin, “Fedformer: Frequency enhanced decomposed transformer for long-term series forecasting,” in *International conference on machine learning*, 2022, pp. 27 268–27 286.
 - [22] Z. Li, S. Qi, Y. Li, and Z. Xu, “Revisiting long-term time series forecasting: An investigation on linear mapping,” in *ArXiv*, 2023.
 - [23] A. Zeng, M. Chen, L. Zhang, and Q. Xu, “Are transform-ers effective for time series forecasting?” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, no. 10, 2023, pp. 11 121–11 128.
 - [24] T. Zhang, Y. Zhang, W. Cao, J. Bian, X. Yi, S. Zheng, and J. Li, “Less is more: Fast multivariate time series forecasting with light sampling-oriented MLP structures,” *CoRR*, 2022.
 - [25] L. Chu, X. Bingjia, and Y. Qiping, “Unettsf: A better performance linear complexity time series prediction model,” in *arXiv preprint arXiv:2401.03001*, 2024.
 - [26] W. Xue, T. Zhou, Q. Wen, J. Gao, B. Ding, and R. Jin, “Card: Channel aligned robust blend transformer for time series forecasting,” in *International Conference on Machine Learning*, 2024.
 - [27] S.-A. Chen, C.-L. Li, S. O. Arik, N. C. Yoder, and T. Pfister, “Tsmixer: An all-mlp architecture for time series forecasting,” *Transactions on Machine Learning Research*, 2023.
 - [28] S. Wang, J. Li, X. Shi, Z. Ye, B. Mo, W. Lin, S. Ju, Z. Chu, and M. Jin, “Timemixer++: A general time series pattern machine for universal predictive analysis,” in *International Conference on Learning Representations*, 2025.
 - [29] 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 *International Conference on Learning Representations*, 2024.
 - [30] T. Zhou, P. Niu, X. Wang, L. Sun, and R. Jin, “One fits all: Power general time series analysis by pretrained lm,” in *Advances in Neural Information Processing Systems*, vol. 36, 2023.
 - [31] S. Wang, H. Wu, X. Shi, T. Hu, H. Luo, L. Ma, J. Y. Zhang, and J. Zhou, “Timemixer: Decomposable multi-scale mixing for time series forecasting,” in *International Conference on Learning Representations*, 2024.
 - [32] A. Das, W. Kong, A. Leach, S. Mathur, R. Sen, and R. Yu, “Long-term forecasting with tide: Time-series dense encoder,” in *ArXiv*, 2023.
 - [33] C. Wang, Q. Qi, J. Wang, H. Sun, Z. Zhuang, J. Wu, L. Zhang, and J. Liao, “Chattime: A unified multimodal time series foundation model bridging numerical and textual data,” in *AAAI Conference on Artificial Intelligence*, 2025.