FITS: CONDITIONAL DIFFUSION MODEL FOR IRREGULAR TIME SERIES FORECASTING WITH PSEUDO-FUTURE EXOGENOUS COVARIATES

Anonymous authors

Paper under double-blind review

ABSTRACT

Irregular multivariate time series (IMTS) present unique challenges due to non-uniform intervals and different sampling rates. While existing methods struggle to capture both long-term dynamics and cross-channel dependencies under such irregularities, we tackle this by formulating time series forecasting as a conditional generation problem and introducing FITS, a conditional diffusion model for IMTS forecasting that leverages pseudo-future exogenous covariates. Our approach incorporates two key innovations. First, we propose a novel entropy-aware adaptive patching scheme that generates data-driven segments with dynamic boundaries determined by the information density. This scheme overcomes the limitations of traditional fixed-length or fixed-span segmentation in preserving continuous local semantics and modeling inter-time series correlations. Second, we develop a transformer-based prior knowledge extractor that captures forward-looking covariate dependencies via a novel cross-variate attention mechanism. The transformer structure is integrated into the conditional diffusion generative process as a unified framework, enabling precise distributional forecasting for IMTS. Extensive experiments on multiple datasets with four evaluation metrics validate the effectiveness of FITS.

1 Introduction

Time series forecasting (TSF) plays a crucial role in numerous real-world applications, facilitating data-driven decision-making across diverse fields. It is widely utilized in domains such as stock price prediction (Li et al., 2024a), weather prediction, transportation planning (Guo et al., 2022), and healthcare. Many approaches, such as autoregressive models (Salinas et al., 2020) and sequence-to-sequence modeling (Wen et al., 2017), frame forecasting as a conditional generative task. In particular, diffusion-based generative models have attracted considerable attention owing to their capabilities in image, video, and text generation (Ho et al., 2020a; Dhariwal & Nichol, 2021; Kong et al., 2021).

Most existing time series diffusion models are designed for regularly sampled time series, such as Li et al. (2024c); Shen et al. (2024); Wang et al. (2025), however, when dealing with sparse and irregularly observed data, there are several obstacles: (1) how to capture irregularities in intra-series dependencies and asynchronies in inter-series correlations amid varying time intervals between adjacent observations; (2) how to extract critical insights from all available historical data, which can then serve as prior knowledge to capture covariate dependencies in both forward and reverse processes within the diffusion model. While prior studies such as Li et al. (2024b) and Shen & Kwok (2023) have proposed effective conditional embeddings to guide the diffusion process, when the conditional inputs (e.g., historical observations) are highly sparse, models face challenges in extracting adequate contextual information as they are unable to capture the temporal dependencies, compromising the reliability of time series prediction.

To this end, we propose a conditional diffusion model for irregular time series forecasting with pseudo-future exogenous covariates (FITS), which integrates a transformer-enhanced modeling approach to capture the forward - backward covariate dynamics. It then leverages this model to generate pseudo forecasts of the target variable, which essentially serve as conditional guidance for generating the unobserved segments of sparse time series, supporting downstream prediction tasks (Fig. 1).

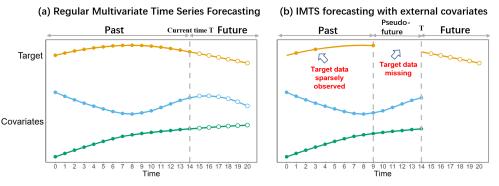


Figure 1: The comparison of regular multivariate time series forecasting and IMTS forecasting frmework considered in this work. (a) The goal of regular multivariate time series forecasting is to simultaneously forecast all the variables in the system. (b) The goal of this work is to forecast the sparse and irregularly observed target time series, given all the external covariates.

Our contributions are summarized as follows: We introduce FITS, a conditional diffusion model designed for forecasting sparse and irregularly observed time series. Specifically: (1) FITS incorporates an adaptive entropy-based patching approach tailored to irregular time series, which leverages local semantic granularities and enables more accurate modeling of inter-series correlations. (2) FITS also employs a transformer-based predictive model learned from external covariates, equipped with forward-looking cross-variate attention mechanisms. During the reverse diffusion process, this learned model is leveraged as a conditional representation to generate accurate probability distributions of future time series. (3) In our experiments, besides standard evaluation metrics such as mean squared error (MSE) and mean absolute error (MAE), we employed Prediction Interval Coverage Probability (PICP) (Yao et al., 2019) and Quantile Interval Coverage Error (QICE) (Han et al., 2022) as metrics for the probabilistic multivariate time series forecasting task. Extensive experiments demonstrate that FITS outperforms state-of-the-art time series diffusion models and performs better than or comparable to various advanced time series prediction models.

2 RELATED WORK

2.1 IRREGULAR MULTIVARIATE TIME SERIES FORECASTING

Existing works have primarily focused on IMTS classification (Yalavarthi et al., 2022; Horn et al., 2020; Tashiro et al., 2021), imputation (Shukla & Marlin, 2021; Yalavarthi et al., 2023) and forecasting (Zhang et al., 2023; Mercatali et al., 2024; Yalavarthi et al., 2024). To summarize the core mechanism of the IMTS forecasting methods in addressing the data irregularities, some authors proposed novel data preprocessing and representation methods, for example, in the patching-based approach, the input time series is represented as matrices with temporal and variable dimensions, and model components are designed to learn dependencies along both dimensions Zhang et al. (2024), however, in the case of sparsely observed time series, the number of observations within a patch may be scarce, resulting in an excessive number of unin-

 formative patches under a temporal resolution; there are also other non-patching approaches that use bipartite graphs (Yalavarthi et al., 2024), or hypergraphs (Li et al., 2025), but their model architectures restrict the ability to capture dependencies in high dimensional or highly sparse IMTS.

In addition to data representation methods for IMTS forecasting, some authors also proposed novel deep architectures and attention mechanisms. For example, *T-PATCHGNN* (Zhang et al., 2024) proposed a time-adaptive graph neural network to model the dynamic intra-patch and inter-patch dependencies. *Warpformer* (Zhang et al., 2023) proposed a doubly self-attention module within the transformer framework for representation learning on multiple sampling granularities. *ContiFormer* (Chen et al., 2023) adopted continuous-time Neural ordinary differential equations (ODEs) within the attention mechanism of Transformers to capture the temporal dynamics of the underlying IMTS system. These methods often presume a specific form of dependency, which introduces significant restrictiveness and fails to accommodate considerations of complex hierarchical, higher-order or multi-scale dependencies.

2.2 TIME SERIES DIFFUSION MODELS

The Denoising Diffusion Probabilistic Models proposed by Ho et al. (2020b) has become a powerful tool for time series modeling (Lin et al., 2024), due to their advantages in fine-grained temporal modeling. Many recent time-series diffusion models have focused on designing effective conditional embeddings to guide the reverse process (Li et al., 2024c; Tashiro et al., 2021; Rasul et al., 2021). For example, TimeGrad (Rasul et al., 2021) employs the hidden state from an RNN as the conditional embedding, Li et al. (2024c) utilized vanilla transformers to extract a representation from historical data, which is then used as a prior knowledge to recover the full distribution of future time series. In addition, Shen & Kwok (2023) further incorporated parts of the ground-truth future predictions for conditioning, which introduces additional inductive bias in the conditioning module for more accurate time series prediction. Shen et al. (2024) also considered other unique time series properties and proposed a multi-resolution diffusion model corresponding to a sequence of fine-to-coarse trend.

So far, the existing works on time series diffusion models have been focused on regularly sampled time series data, in the context of IMTS, representations extracted from historical data may fail to capture the underlying trends of the sequence, leading to a lack of reliable prior guidance, making it prone to generating sequences that are disconnected from historical patterns. Furthermore, in terms of model training during the reverse process, it is difficult to generate the desired series when there are limited fine granularity information (Coletta et al., 2023), which may provide unreliable underlying inputs for the multi-resolution framework and thus undermining the consistency of the overall trend.

3 Proposed Method

In this work, we assume that the total length of the observed time series is T, where the historical observed target time series $\mathbf{x}_{0:T-L}$ (0 < L < T) is sparse and irregularly sampled, with its last valid observation recorded at time T-L. Furthermore, we consider multiple exogenous covariates $\mathbf{z}_{0:T} \in \mathbb{R}^{T \times C}$, where C represents the dimensionality of the exogenous covariates; by definition, any time series that provides predictive value for the prediction target is classified as an exogenous covariate. The proposed diffusion-based forecasting framework aims to predict the future segment $\mathbf{x}_{T:T+H}$ using a model \mathcal{F}_{θ} that specifically captures all available information embedded in the historical observed time series $\mathbf{x}_{0:T-L}$ and exogenous covariates $\mathbf{z}_{0:T}$.

$$\widehat{\mathbf{x}}_{T:T+H} = \mathcal{F}_{\theta} \left(\mathbf{x}_{0:T-L}, \mathbf{z}_{0:T} \right). \tag{1}$$

Fig. 2 shows an overview of the proposed model.

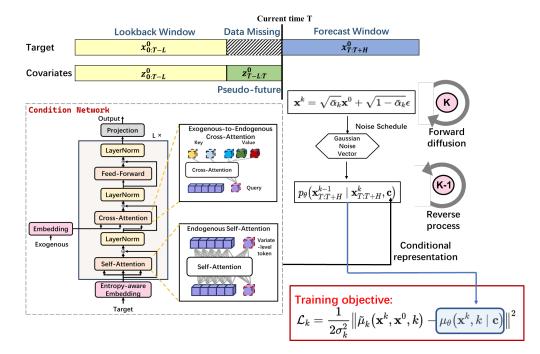


Figure 2: Overall framework of the proposed FITS framework.

3.1 FORWARD DIFFUSION PROCESS

During model training, the objective of the forward diffusion is to diffuse the "future" time steps $\mathbf{x}_{T:T+H}$ of the target time series. At the k-th step of the forward process, \mathbf{x}^k is parameterized by adding noise to the previous diffusion step k-1, scaled by $\sqrt{1-\beta_k}$:

$$q\left(\mathbf{x}^{k} \mid \mathbf{x}^{k-1}\right) = \mathcal{N}\left(\mathbf{x}^{k}; \sqrt{1-\beta_{k}}\mathbf{x}^{k-1}, \beta_{k}\mathbf{I}\right), \quad k = 1, \dots, K,$$
(2)

with $\beta_t \in (0,1)$ representing the noise variance following a predefined schedule. It can be shown that:

$$q\left(\mathbf{x}^{k} \mid \mathbf{x}^{0}\right) = \mathcal{N}\left(\mathbf{x}^{k}; \sqrt{\bar{\alpha}_{k}}\mathbf{x}^{0}, (1 - \bar{\alpha}_{k})\mathbf{I}\right), \tag{3}$$

where $\bar{\alpha}_k = \prod_{s=1}^k \alpha_s$, and $\alpha_k = 1 - \beta_k$. Then, \mathbf{x}^k is given as:

$$\mathbf{x}^k = \sqrt{\bar{\alpha}_k} \mathbf{x}^0 + \sqrt{1 - \bar{\alpha}_k} \epsilon, \quad \epsilon \in \mathcal{N}(0, \mathbf{I}).$$
 (4)

The subscript of $\mathbf{x}_{T:T+H}$ is omitted for notational simplicity.

3.2 CONDITIONING THE BACKWARD DENOISING PROCESS

Existing time series diffusion models typically incorporate either the original historical observation segment $\mathbf{x}_{0:T}$ (Tashiro et al., 2021) or a derived representation $\mathcal{F}(\cdot)$ from historical data (Li et al., 2024b) as input to their conditioning networks. In contrast, this study proposes to leverage the evolutionary dynamics embedded in external covariates, which capture the relationship from $\mathbf{z}_{0:T-L}$ to $\mathbf{z}_{T-L:T}$. This latent process characterize the potential variation patterns of the target variable from the historically observed part $\mathbf{x}_{0:T-L}$ to the "pseudo-future" segment $\mathbf{x}_{T-L:T}$, thereby facilitating predictive inference.

3.2.1 Entropy-aware patching and encoding for irregular time series

In this subsection, we propose a novel information density-based patching and encoding approach applied to all variables. For IMTS, it is difficult to capture the local dynamic granular scemantics due to discretionary segmentation of continuous observations, which hinders the effective extraction of low-dimensional latent factors and state evolution patterns. For example, a patient's sudden health deterioration may be segmented across two time windows, which fragments this critical pattern and prevents it from being fully captured.

Entropy-aware module to compute dynamic window boundaries. To fully leverage temporal information, we first enrich each raw observation by filling the missing points with zero. Motivated by Liu et al. (2025), assume the historical observation is initially divided into P patches with length $S_{\rm init} = T/P$. For each patch p, the initial reference center is $c_p = (p-0.5) \cdot (T/P)$, and the window boundaries can be computed as:

$$t_p^{\text{left}} = c_p - \frac{S_{\text{init}}}{2} + \delta_p, \quad t_p^{\text{right}} = t_p^{\text{left}} + \exp\left(\lambda_p\right).$$
 (5)

In this work, we propose a novel boundary network (BoundaryNet) based on a sample entropy (SampEn) measure to specifically learn the parameters δ_p and λ_p in Eq. (5). Specifically, the SampEn measure proposed by Richman & Moorman (2000) quantifies the information richness of a time series: a higher entropy value indicates a more complex series that harbors dense implicit information. For a given patch \mathbf{x}_p , using a light-weight MLP network, we can map the entropy SampEn $_p$ to the latent space:

$$\mathbf{e}_p = \text{Linear}\left(\text{ELU}\left(\text{Linear}(\text{SampEn}_p)\right)\right) \in \mathbb{R}^{D_e}.$$
 (6)

Then, \mathbf{e}_p is concatenated with the Rotary Position Embedding (RoPE) (Su et al., 2024) $\mathrm{PE}(c_p) \in \mathbb{R}^{D_{pe}}$ to arrive at a enhanced time-information representation: $\widetilde{\mathrm{PE}}(c_p) = \mathrm{Concat}(\mathrm{PE}(c_p), \mathbf{e}_p) \in \mathbb{R}^{D_{pe}+D_e}$.

Finally, the proposed BoundaryNet to calculate the two scalar boundary parameters is given as:

$$[\delta_{p}, \lambda_{p}] = \operatorname{Linear}_{\text{output}} \left(\operatorname{SiLU} \left(\operatorname{Linear}_{\text{hidden}} \left(\widetilde{\operatorname{PE}} \left(c_{p} \right) \right) \right) \right). \tag{7}$$

Substitute Eq. (7) into Eq. (5), we can effectively compute the dynamically adjusted window boundaries based on the information density.

Adaptive patch representations. After defining the dynamic temporal windows, using the method proposed by Liu et al. (2025), we calculate a relevance weight $\alpha_{i,p}$ using $[\delta_p, \lambda_p]$ for each observation i in patch p and arrive at the final representation:

$$\bar{h}_p = \frac{\sum_{i=1}^{L_p} \alpha_{i,p} \cdot \tilde{v}_i}{\sum_{i=1}^{L} \alpha_{i,p} + \epsilon} \in \mathbb{R}^{1 + D_{pe} + D_e + D_{te}}, \tag{8}$$

where L_p denotes the number of observations in patch p, and $\tilde{v}_i = \operatorname{Concat}(\mathbf{x}_p(t_i), \widetilde{\operatorname{PE}}(c_p), \operatorname{TE}(t_i))$, $\operatorname{TE}(t_i) \in \mathbb{R}^{D_{te}}$ denotes the learnable time embedding. Then, \bar{h}_p is projected into the model's uniform hidden space via a linear layer: $h_p = \operatorname{Linear}_D(\bar{h}_p) \in \mathbb{R}^D$. Therefore, we have for the whole sequence: $H = [h_1, \dots, h_P] \in \mathbb{R}^{P \times D}$.

3.2.2 Learning conditional representation through reconciliating target and exogenous information

In this work, transformer is utilized as a prior knowledge extractor, capturing covariate-dependence in the reverse process within the diffusion model. In addition to the patch representation H derived in the previous section, the entire target time series $\mathbf{x}_{0:T}$ is also embedded into one single series-level global token embedding \mathbf{G}_{tar} via the same trainable linear MLP projector.

Intra-series self-attention. In the patch-level attention, we apply multi-head attention with causal masking to all variables to capture their intra-variate cross-time dependency. Taking the target variable as an example, and dropping layer index for brevity, this can be formalized as:

$$\begin{split} \widetilde{\mathbf{H}}_{:}^{pat} &= \mathrm{LN}\left(\mathbf{H}_{:} + \mathrm{MHA}\left(\mathbf{H}_{:}, \mathbf{H}_{:}, \mathbf{H}_{:}\right)\right), \\ \mathbf{H}_{:} &= \mathrm{LN}\left(\widetilde{\mathbf{H}}_{:}^{pat} + \mathrm{FFN}\left(\widetilde{\mathbf{H}}_{:}^{pat}\right)\right), \end{split} \tag{9}$$

where $\mathbf{H}_{:}$ denotes the collective token embeddings of a variable at all patch steps, LN denotes layer normalization, MHA($\mathbf{Q}, \mathbf{K}, \mathbf{V}$) denotes the multi-head attention layer where \mathbf{Q}, \mathbf{K} , and \mathbf{V} serve as queries, keys and values, and FFN denotes a feed-forward network. In addition, we also employ a series-level global token embedding \mathbf{G}_{tar} , which serves as a bridge that connects the patches in the target variable and the exogenous variables (Wang et al., 2024). Accordingly, we also employ a variate-to-patch attention $\mathbf{H}_{:}^{\text{var-to-pat}}$ and a patch-to-variate attention $\mathbf{G}^{\text{pat-to-var}}$ (Wang et al., 2024), which offers a holistic perspective of the temporal dependencies inherent to the target variable, while also enabling enhanced interactions with exogenous variables that exhibit arbitrary irregularity.

Inter-series cross-attention. Assume the last observed data point of the target variable occurs at time T-L; $\mathbf{z}_{T-L:T}$ thus constitutes a relative future segment relative to $\mathbf{x}_{0:T-L}$. To this end, we redesign the cross-attention layer: the global token of the target variable, \mathbf{G}_{tar} , remains the query (Q), while exogenous variables are split into two segments for the key (K) and value (V), where the embedding of the historical segment $\mathbf{z}_{0:T-L}$ serves as K and the embedding of the pseudo-future segment $\mathbf{z}_{T-L:T}$ serves as V. The learned global token of the target acts as a bridge to integrate and filter exogenous information, ensuring that only relevant insights support the prediction of the target variable.

3.2.3 CONDITIONING NETWORK

Following Shen & Kwok (2023), using the transformer network $\mathcal{T}(\cdot)$ derived from the previous sections, we adopt the future mixup strategy which combines the past information's mapping $\hat{\mathbf{x}}_{T:T+H} = \mathcal{T}(\mathbf{x}_{0:T-L})$ with the future ground-truth $\mathbf{x}_{T:T+H}^0$, which is only available during training. At diffusion step k, it produces the conditioning signal \mathbf{c} as:

$$\mathbf{c} = \mathbf{m}^k \mathcal{T} \left(\mathbf{x}_{0:T-L} \right) + \left(1 - \mathbf{m}^k \right) \mathbf{x}_{T:T+H}^0. \tag{10}$$

Here, $\mathbf{m}^k \in [0,1)^{1 \times H}$ is a mixing coefficient randomly sampled from the uniform distribution on [0,1). During inference, $\mathbf{x}_{T:T+H}^0$ is no longer available, and the condition \mathbf{c} is set to $\mathcal{T}(\mathbf{x}_{0:T-L})$.

3.3 Denoising reverse process

The reverse denoising process is a markov chain. At the k-th denoising step, $\mathbf{x}_{T:T+H}^{k-1}$ is generated from $\mathbf{x}_{T:T+H}^k$ by sampling from the following normal distribution, subject to the conditional representation \mathbf{c} :

$$p_{\theta}\left(\mathbf{x}_{T:T+H}^{k-1} \mid \mathbf{x}_{T:T+H}^{k}, \mathbf{c}\right) = \mathcal{N}\left(\mathbf{x}_{T:T+H}^{k-1}; \mu_{\theta}\left(\mathbf{x}_{T:T+H}^{k}, k \mid \mathbf{c}\right), \Sigma_{\theta}\left(\mathbf{x}_{T:T+H}^{k}, k\right)\right), \tag{11}$$

where the variance $\Sigma_{\theta}\left(\mathbf{x}_{T:T+H}^{k},k\right)$ is fixed to $\sigma_{k}^{2}\mathbf{I}$. The goal of this reverse process is to learn this mean function $\mu_{\theta}(\mathbf{x}_{T:T+H}^{k},k)$, which effectively produces $\mathbf{x}_{T:T+H}^{k-1}$ close to the ground truth. Through iterative denoising steps, the prediction result $\hat{x}_{T:T+H}^{0}$ is ultimately recovered to match the distribution of the original time series. To train the diffusion model, considering Eqs. (2) and (11), one uniformly samples k from $\{1,2,\ldots,K\}$ and then minimizes the KL (Kullback-Leibler) divergence:

$$\mathcal{L}_k = D_{\mathrm{KL}}(q(\mathbf{x}^{k-1} \mid \mathbf{x}^k) || p_{\theta}(\mathbf{x}^{k-1} \mid \mathbf{x}^k, \mathbf{c})), \tag{12}$$

where $q(\mathbf{x}^{k-1} \mid \mathbf{x}^k)$ is the ground-truth conditional data distribution.

Then, the training objective in (12) is then formulated as:

$$\mathcal{L}_{k} = \frac{1}{2\sigma_{k}^{2}} \left\| \tilde{\mu}_{k} \left(\mathbf{x}^{k}, \mathbf{x}^{0}, k \right) - \mu_{\theta} \left(\mathbf{x}^{k}, k \mid \mathbf{c} \right) \right\|^{2}.$$
(13)

The estimation of $\mu_{\theta}(\mathbf{x}^k, k \mid \mathbf{c})$ can be computed via a noise prediction model $\epsilon_{\theta}(\mathbf{x}^k, k)$ following Benny & Wolf (2022).

During inference, a noise vector $\mathbf{x}_{1:H}^K \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ is generated, and through the reverse denoising process, we can obtain the final prediction result $\hat{\mathbf{x}}_{T:T+H}^0$.

4 EXPERIMENTS

In this section, we perform extensive experiments to compare the proposed FITS with recent 6 state-of-the-art (SOTA) time series prediction models on 7 commonly used real-world datasets.

4.1 SETUP

Benchmark datasets. Experiments were performed on 7 public benchmark datasets with different levels of multivariate correlations. The datasets include: (i) Electricity Price Forecasting Dataset (EPF) (5 subdatasets from different major power markets) (Lago et al., 2021). (ii) Exchange (Daily exchange rates of eight different countries) (Lai et al., 2018). (iii) Weather (21 meteorological variables from Germany) (Zhou et al., 2021). Due to space constraints, detailed descriptions of the datasets are deferred to Appendix B.1. Appendix B.1 also includes the process of downsampling the original data to arrive at the sparse and irregular time series used in this work, which in turn contextualizes the data preparation aligned with our research focus. The data are processed using two random missingness strategies: (1) Random Missing (RM), and (2) Block Missing (BM).

Baselines. To establish a comprehensive benchmark for our proposed FITS method, we select baselines from four methodological domains. Specifically, we include: (i) **Time series diffusion models:** CSDI (Tashiro et al., 2021); Transformer-Modulated Diffusion Model (TMDM) (Li et al., 2024b); Diffusion-TS (Yuan & Qiao, 2024). (ii) **Time series transformers:** PatchTST (Nie et al., 2023); Crossformer (Zhang & Yan, 2023). (iii) **Other time series forecasting methods.** TiDE (Das et al., 2023); DLinear (Zeng et al., 2023). See Appendix B.2 for more details about the baselines.

Implementation details. In our experiments, we employed a linear noise schedule with $\beta_1 = 10^{-4}$ and $\beta_K = 0.02$, setting the number of diffusion timesteps to K = 1000. We approximated the data distribution using 100 samples, and all experiments were repeated 5 times with seeds $\{1, 2, 3, 4, 5\}$. The model was trained using the Adam optimizer with a learning rate of 10^{-4} and a batch size of 64. Additional details are given in Appendix B.3.

4.2 Main results

4.2.1 PROBABILISTIC FORECASTING

To intuitively illustrate the probabilistic distribution forecasting capabilities of the models, we present the forecasting results of our proposed FITS model alongside three comparative baseline models in Figure 4. Specifically, we visualize the 50% and 90% prediction intervals (denoted by dark green and light green, respectively) and overlay the true observed values for direct reference. It is worth noting that certain baseline models were originally devised for generative tasks rather than dedicated probabilistic forecasting; however, their authors have asserted that these models are capable of yielding probabilistic forecasting results (Yuan & Qiao, 2024).

Table 1: Performance comparisons in terms of QICE and CRPS. The best results are boldfaced, and the suboptimal results are underlined. The table presents both the scenarios of no missing values and random missingness of 0.5.

	4t	1 337	41	I 171		1 8	ID.	l Di	D. f	l D	E	l T	TD.			
metrci			Weather		Exchange		NP		PJM		BE		FR		DE	
model		QICE	CRPS	QICE	CRPS	QICE	CRPS	QICE	CRPS	QICE	CRPS	QICE	CRPS	QICE	CRPS	
CSDI	No-Missing RM=0.5	13.91 14.90	0.735 0.763	10.91 10.76	0.178 0.189	$\frac{1.96}{2.394}$	0.279 <u>0.324</u>	16.75 14.37	0.603 0.624	14.63 14.33	0.393 0.452	15.60 15.32	0.377 0.386	13.34 13.56	0.683 0.769	
TMDM	No-Missing RM=0.5	$\frac{10.86}{12.20}$	$\frac{0.485}{0.554}$	$\frac{7.439}{7.305}$	0.516 0.488	5.609 5.453	0.593 0.503	$\frac{4.541}{2.511}$	0.209 0.178	$\frac{6.096}{9.001}$	$\frac{0.353}{0.532}$	$\frac{5.125}{8.795}$	$\frac{0.294}{0.388}$	$\frac{4.577}{4.488}$	$\frac{0.407}{0.410}$	
Diffusion-TS	No-Missing RM=0.5	12.13 13.690	0.532 0.543	15.90 15.27	1.310 1.040	9.692 10.25	0.612 0.635	15.36 15.20	0.218 0.256	8.236 8.569	0.401 <u>0.490</u>	9.486 10.50	0.376 0.358	13.96 12.30	0.785 0.813	
our	No-Missing RM=0.5	3.275 4.170	0.409 0.497	4.966 4.979	$\frac{0.354}{0.359}$	1.800 1.924	$\frac{0.287}{0.277}$	3.007 2.976	0.176 0.177	2.854 3.023	0.226 0.226	3.739 3.830	0.182 0.190	2.162 1.024	0.390 0.377	

As evidenced by the visualization, FITS demonstrates superior performance in probabilistic distribution forecasting. This advantage can be attributed to the design of our conditional estimation module, which enables more accurate mean estimation even when the input data suffers from temporal misalignment and high proportions of missing values. In particular, the inter-series cross-attention component embedded within this module facilitates the model's effective extraction and utilization of latent information in pseudo-future data, thereby enhancing forecasting reliability. Nevertheless, in scenarios where there exist unobserved gaps between the historical information window and the target forecasting window, all models encounter heightened challenges in capturing future trend dynamics, resulting in elevated predictive uncertainty. Further details and supplementary analyses related to these experiments are provided in Appendix C.2.

To quantitatively analyze the models' probabilistic forecasting capabilities, we adopted CRPS (Continuous Ranked Probability Score) and QICE (Quantile Interval Coverage Error) as evaluation metrics, following the approach of Li et al. (2024b). For both metrics, smaller values indicate better performance. Table 1 shows the CPRS and QICE on the time series. Notably, our model achieves the optimal performance on nearly all datasets, with its CRPS and QICE values consistently remaining at the lowest level among all compared models, fully demonstrating its superior probabilistic forecasting capability.

4.2.2 Non-Probabilistic forecasting

Table 2 presents the Mean Squared Error (MSE) results on time series. It can be observed that the performance improvement of the model is particularly significant on more complex datasets such as BE and FR. Among all datasets, the FITS model ranks first in 4 of them; overall, its average MSE is better than all other baseline models, including the latest diffusion models. In addition, it can be found that our model achieves a better average ranking when the random missing situation is better. It should be noted that the model does not achieve performance improvement on long-term forecasting datasets such as exchange rate. This may be because there are no complex interdependencies between variables in such datasets, leading to the introduction of noise by the inter-variable attention mechanism in the conditional estimation model. TiDE and DLinear, the two channel-independent models, achieved the optimal and suboptimal performance respectively, which also corroborates this point. In contrast, the covariates of the EPF dataset have been confirmed to indeed have a positive effect on target prediction, so our model has achieved better performance on this dataset.

Furthermore, we also found that all diffusion models in the baselines perform poorly, which indicates that the current diffusion models are weaker than general models in terms of mean prediction ability.

4.2.3 ABLATION STUDY

Table 2: Performance comparisons in terms of MAE and MSE. The best results are boldfaced, and the suboptimal results are underlined. The table presents both the scenarios of no missing values and random missingness of 0.5.

metric		Wea	ther	Excl	nange	N	IP	PJ	M	B	E	F	R	I	DΕ	AV	G
model		MSE	MAE														
TIDE	No-Missing RM=0.5	$\frac{0.520}{0.612}$	$\frac{0.519}{0.562}$	0.073 0.462	0.215 0.494	0.515 0.608	0.458 0.520	0.159 0.193	0.261 0.283	0.595	0.377 0.417	$\frac{0.506}{0.481}$	0.302 0.339	0.869 0.931	0.602 0.604	0.462 0.573	$\frac{0.390}{0.454}$
DLinear	No-Missing RM=0.5	0.849 0.855	0.653 0.656	$\frac{0.260}{0.265}$	0.412 0.416	0.507 0.528	0.448 0.45	0.176 0.178	0.276 0.28	0.651 0.628	0.399 0.397	0.587 0.577	0.344 0.344	0.867 0.881	0.601 0.602	0.556 0.558	0.447 0.449
Crossformer	No-Missing RM=0.5	0.486 0.486	0.499 0.501	0.543 0.557	0.602 0.598	0.327 0.340	0.308 0.341	0.231 0.242	$\frac{0.234}{0.238}$	$\frac{0.524}{0.529}$	$\frac{0.373}{0.379}$	0.537 0.499	0.303 0.302	0.578 0.577	0.483 0.491	$\frac{0.460}{0.481}$	0.400 0.407
CSDI	No-Missing RM=0.5	0.813 0.830	0.632 0.652	0.507 0.542	0.701 0.781	0.526 0.731	0.441 0.532	0.370 0.453	0.336 0.312	0.590 0.683	0.384 0.543	0.556 0.506	0.332 0.352	0.838 0.871	0.609 0.632	0.600 0.659	0.490 0.543
TMDM	No-Missing RM=0.5	0.568 0.671	0.551 0.584	0.571 0.532	0.641 0.610	1.068 0.779	0.718 0.627	0.170 0.155	0.286 0.285	0.651 0.728	0.467 0.541	0.640 0.506	$\frac{0.280}{0.357}$	0.687 0.686	0.514 0.513	0.622 0.579	0.493 0.502
Diffusion-TS	No-Missing RM=0.5	0.713 0.766	0.621 0.649	2.507 1.960	1.280 1.240	1.260 1.380	0.869 0.869	0.181 0.186	0.292 0.290	0.709 0.750	0.501 0.576	0.583 0.455	0.363 0.372	0.832 0.817	0.5964 0.571	0.969 0.902	0.646 0.652
our	No-Missing RM=0.5	0.554 0.559	0.541 0.547	0.394 0.395	0.551 0.474	$\frac{0.364}{0.362}$	$\frac{0.368}{0.356}$	$\frac{0.175}{0.161}$	0.232 0.238	0.443 0.398	0.292 0.296	0.424 0.384	0.233 0.244	$\frac{0.659}{0.581}$	$\frac{0.501}{0.488}$	0.430 0.405	0.388 0.377

We compared prediction results of one full model and three ablation variants in Table 3. The **rp-atten** variant replaces the proposed inter-variable attention with standard cross-attention, leading to performance degradation; **w/o-covar** removes the inter-variable attention module for univariate prediction, causing significant performance decline—the most severe among all variants; **rp-patch** uses standard instead of attention-driven patch partitioning. The experimental results show that the inter-series attention plays an important role in the EPF dataset, and the other components also have a positive impact on the experimental results.

Table 3: Ablation experiment results in terms of MSE

	Wea.	FR	BE
rp-atten	0.568	0.388	0.414
w/o-covar	0.607	0.423	0.444
rp-patch	0.569	0.391	0.401
FITS	0.559	0.384	0.398

5 CONCLUSION

In this work, we propose FITS, an innovative framework that integrates a diffusion generative process with a newly designed transformer-based conditional representation learning framework. In particular, our approach introduces two key innovations: first, we propose an entropy-based adaptive patching method that leverages the sample entropy measure to effectively capture granular local semantics, which avoids information fragmentation caused by discretionary segmentation. Second, we propose a novel cross-variate attention module to effectively capture the evolutionary dynamics of covariates. By using this transformer-based representation module as a conditional guidance for generating future target variables, the diffusion model can be more effectively guided toward the true values. Extensive experiments demonstrate that FITS achieves superior performance in both point forecasting and probabilistic forecasting quality.

REFERENCES

Yaniv Benny and Lior Wolf. Dynamic dual-output diffusion models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022.

Yuqi Chen, Kan Ren, Yansen Wang, Yuchen Fang, Weiwei Sun, and Dongsheng Li. Contiformer: Continuous-time transformer for irregular time series modeling. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neu*ral Information Processing Systems, volume 36, pp. 47143–47175. Curran Associates, Inc.,

2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/ 9328208f88ec69420031647e6ff97727-Paper-Conference.pdf.

- Andrea Coletta, Sriram Gopalakrishnan, Daniel Borrajo, and Svitlana Vyetrenko. On the constrained timeseries generation problem. *Advances in Neural Information Processing Systems*, 36:61048–61059, 2023.
- Abhimanyu Das, Weihao Kong, Andrew Leach, Shaan K Mathur, Rajat Sen, and Rose Yu. Long-term forecasting with tiDE: Time-series dense encoder. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL https://openreview.net/forum?id=pCbC3aQB5W.
- Prafulla Dhariwal and Alex Nichol. Diffusion models beat gans on image synthesis. In NIPS'21: Proceedings of the 35th International Conference on Neural Information Processing Systems, pp. 8780 8794. NeurIPS, 12 2021.
- Kan Guo, Yongli Hu, Zhen Qian, Yanfeng Sun, Junbin Gao, and Baocai Yin. Dynamic graph convolution network for traffic forecasting based on latent network of laplace matrix estimation. *IEEE Transactions on Intelligent Transportation Systems*, 23(2):1009–1018, 2022. doi: 10.1109/TITS.2020.3019497.
- Xizewen Han, Huangjie Zheng, and Mingyuan Zhou. Card: Classification and regression diffusion models. In *Thirty-Sixth Conference on Neural Information Processing Systems*, 2022.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 6840-6851. Curran Associates, Inc., 2020a. URL https://proceedings.neurips.cc/paper_files/paper/2020/file/4c5bcfec8584af0d967flab10179ca4b-Paper.pdf.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*, volume 33, pp. 6840–6851, 2020b.
- Max Horn, Michael Moor, Christian Bock, Bastian Rieck, and Karsten Borgwardt. Set functions for time series. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 4353–4363. PMLR, Jul 13–18 2020. URL https://proceedings.mlr.press/v119/horn20a.html.
- Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. Diffwave: A versatile diffusion model for audio synthesis. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=a-xFK8Ymz5J.
- Jesus Lago, Grzegorz Marcjasz, Bart De Schutter, and Rafał Weron. Forecasting day-ahead electricity prices: A review of state-of-the-art algorithms, best practices and an open-access benchmark. *Applied Energy*, 293:116983, 2021. ISSN 0306-2619. doi: https://doi.org/10.1016/j.apenergy.2021.116983. URL https://www.sciencedirect.com/science/article/pii/S0306261921004529.
- Guokun Lai, Wei-Cheng Chang, Yiming Yang, and Hanxiao Liu. Modeling long-and short-term temporal patterns with deep neural networks. In *The 41st international ACM SIGIR conference on research & development in information retrieval*, pp. 95–104, 2018.
- Boyuan Li, Yicheng Luo, Zhen Liu, Junhao Zheng, Jianming Lv, and Qianli Ma. Hyperimts: Hypergraph neural network for irregular multivariate time series forecasting. In *Proceedings of the 42nd International Conference on Machine Learning (ICML)*, 2025.

- 473
- 474 475
- 476 477 478
- 479 480
- 481 482 483
- 484 485 486
- 487 488
- 489 490
- 491 492 493
- 494 495 496
- 497 498 499
- 500 501 502
- 503 504 505
- 506 507 508
- 509 510 511
- 512 513
- 514 515 516

- Tong Li, Zhaoyang Liu, Yanyan Shen, Xue Wang, Haokun Chen, and Sen Huang. Master: Market-guided stock transformer for stock price forecasting. Proceedings of the AAAI Conference on Artificial Intelligence, 38(1):162-170, Mar. 2024a. doi: 10.1609/aaai.v38i1.27767. URL https://ojs.aaai.org/ index.php/AAAI/article/view/27767.
- Yuxin Li, Wenchao Chen, Xinyue Hu, Bo Chen, baolin sun, and Mingyuan Zhou. Transformer-modulated diffusion models for probabilistic multivariate time series forecasting. In The Twelfth International Conference on Learning Representations, 2024b. URL https://openreview.net/forum?id= qae04YACHs.
- Yuxin Li, Wenchao Chen, Xinyue Hu, Bo Chen, Mingyuan Zhou, et al. Transformer-modulated diffusion models for probabilistic multivariate time series forecasting. In The Twelfth International Conference on Learning Representations, 2024c.
- L. Lin, Z. Li, R. Li, X. Li, and J. Gao. Diffusion models for time-series applications: A survey. Frontiers of *Information Technology & Electronic Engineering*, 25(1):19–41, 2024.
- Xvyuan Liu, Xiangfei Qiu, Xingjian Wu, Zhengyu Li, Chenjuan Guo, Jilin Hu, and Bin Yang. Rethinking irregular time series forecasting: A simple yet effective baseline, 2025.
- G. Mercatali, A. Freitas, and J. Chen. Graph neural flows for unveiling systemic interactions among irregularly sampled time series. In The Thirty-eighth Annual Conference on Neural Information Processing *Systems*, 2024.
- Yuqi Nie, Nam H Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth 64 words: Long-term forecasting with transformers. In The Eleventh International Conference on Learning Representations, 2023. URL https://openreview.net/forum?id=Jbdc0vT0col.
- Kashif Rasul, Calvin Seward, Ingmar Schuster, and Roland Vollgraf. Autoregressive denoising diffusion models for multivariate probabilistic time series forecasting. In International Conference on Machine Learning, pp. 8857-8868, 2021.
- Joshua S. Richman and J. Randall Moorman. Physiological time-series analysis using approximate entropy and sample entropy. American Journal of Physiology - Heart and Circulatory Physiology, 278(6):H2039-H2049, 2000. doi: 10.1152/ajpheart.2000.278.6.H2039.
- David Salinas, Valentin Flunkert, Jan Gasthaus, and Tim Januschowski. Deepar: Probabilistic forecasting with autoregressive recurrent networks. *International Journal of Forecasting*, 36(3):1181–1191, 2020.
- Lifeng Shen and James Kwok. Non-autoregressive conditional diffusion models for time series prediction. In International Conference on Machine Learning, 2023.
- Lifeng Shen, Weiyu Chen, and James Kwok. Multi-resolution diffusion models for time series forecasting. In The Twelfth International Conference on Learning Representations, 2024. URL https: //openreview.net/forum?id=mmjnr0G8ZY.
- S. N. Shukla and B. Marlin. Multi-time attention networks for irregularly sampled time series. In International Conference on Learning Representations, 2021.
- Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomput.*, 568(C), February 2024.
- Yusuke Tashiro, Jiaming Song, Yang Song, and Stefano Ermon. Csdi: Conditional score-based diffusion models for probabilistic time series imputation. In Advances in Neural Information Processing Systems 34, pp. 24804–24816, 2021.

Daoyu Wang, Mingyue Cheng, Zhiding Liu, and Qi Liu. TimeDART: A diffusion autoregressive transformer for self-supervised time series representation. In Forty-second International Conference on Machine Learning, 2025. URL https://openreview.net/forum?id=v2G9HML7ep.

- Yuxuan Wang, Haixu Wu, Jiaxiang Dong, Guo Qin, Haoran Zhang, Yong Liu, Yunzhong Qiu, Jianmin Wang, and Mingsheng Long. Timexer: Empowering transformers for time series forecasting with exogenous variables. *Advances in Neural Information Processing Systems*, 37:469–498, 2024.
- Ruofeng Wen, Kari Torkkola, Balakrishnan Narayanaswamy, and Dhruv Madeka. A multi-horizon quantile recurrent forecaster. *arXiv preprint arXiv:1711.11053*, 2017.
- V. K. Yalavarthi, K. Madhusudhanan, R. Scholz, N. Ahmed, J. Burchert, S. Jawed, S. Born, and L. Schmidt-Thieme. Grafiti: Graphs for forecasting irregularly sampled time series. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(15):16255–16263, 2024. ISSN 2374-3468. doi: 10.1609/aaai.v38i15. 29560.
- Vijaya Krishna Yalavarthi, Johannes Burchert, and Lars Schmidt-Thieme. Dcsf: Deep convolutional set functions for classification of asynchronous time series. In 2022 IEEE 9th International Conference on Data Science and Advanced Analytics (DSAA), pp. 1–10. IEEE, 2022.
- Vijaya Krishna Yalavarthi, Johannes Burchert, and Lars Schmidt-Thieme. Tripletformer for probabilistic interpolation of irregularly sampled time series. In 2023 IEEE International Conference on Big Data (BigData), pp. 986–995. IEEE, 2023.
- Jiayu Yao, Weiwei Pan, Soumya Ghosh, and Finale Doshi-Velez. Quality of uncertainty quantification for bayesian neural network inference, 2019.
- Xinyu Yuan and Yan Qiao. Diffusion-TS: Interpretable diffusion for general time series generation. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=4hlapFj099.
- Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are transformers effective for time series forecasting? *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(9):11121–11128, Jun. 2023. doi: 10. 1609/aaai.v37i9.26317. URL https://ojs.aaai.org/index.php/AAAI/article/view/26317.
- J. Zhang, S. Zheng, W. Cao, J. Bian, and J. Li. Warpformer: A multi-scale modeling approach for irregular clinical time series. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 3273–3285, 2023. doi: 10.1145/3580305.3599543.
- Weijia Zhang, Chenlong Yin, Hao Liu, Xiaofang Zhou, and Hui Xiong. Irregular multivariate time series forecasting: A transformable patching graph neural networks approach. In *Forty-first International Conference on Machine Learning*, 2024.
- Yunhao Zhang and Junchi Yan. Crossformer: Transformer utilizing cross-dimension dependency for multivariate time series forecasting. In *International Conference on Learning Representations*, 2023.
- Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang. Informer: Beyond efficient transformer for long sequence time-series forecasting. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(12):11106–11115, May 2021. doi: 10.1609/aaai.v35i12.17325. URL https://ojs.aaai.org/index.php/AAAI/article/view/17325.

564 Algorithm 1 Training 565 **Require:** Number of diffusion steps K. 566 1: repeat 567 Sample $\mathbf{x}_{T:T+H}^0$ from the training set; 2: 568 3: $k \sim \text{Uniform}(\{1, 2, \dots, K\}), \epsilon \sim \mathcal{N}(0, \mathbf{I});$ 569 Compute $\mathbf{x}_{T:T+H}^k$ following Eq. (4); 4: 570 5: Using the transformer network given in Section ??, obtain condition c based on Eq. (10); 571 6: Compute the mean function $\mu_{\theta}(\cdot)$ using the proposed GCN framework in Section ??; 572 Use the reverse denoising process to generate denoised sample $\mathbf{x}_{T:T+H}^{k-1}$ by Eq. (??); 7: 573 8: Calculate the loss $\mathcal{L}_k(\theta)$ in (13); 574 9: Take gradient descent step on $\nabla_{\theta} \mathcal{L}_k(\theta)$; 575 10: until converged 576 577 578 TRAINING ALGORITHM 579

The training procedure is provided in Algorithm 1 below.

B DATASETS AND BASELINES

B.1 DATASETS

580

581 582

583 584 585

586

587

588

590 591

592

593

595 596

597

598

599

600

601

602

603

605

606

607

608 609

610

We assessed the effectiveness of the proposed FITS model through extensive experiments on 7 time series forecasting datasets. As our focus is on sparse and irregularly sampled time series, we modified the originally regular datasets by applying a subsampling procedure with different filtering rates to induce sparsity.

First, detailed descriptions of the original datasets are provided below:

- (1) The **EPF** is an electricity price forecasting dataset, which contains five datasets representing five different day-ahead electricity markets spanning six years each (Lago et al., 2021).
 - **NP** represents the Nord Pool electricity market, recording the hourly electricity price, and corresponding grid load and wind power forecast from 2013-01-01 to 2018-12-24.
 - **PJM** corresponds to the Pennsylvania New Jersey Maryland (PJM) market. It contains the zonal electricity price in the Commonwealth Edison (COMED) area, along with the corresponding system load and COMED load forecast data, spanning from 2013-01-01 to 2018-12-24.
 - **BE** stands for Belgium's electricity market. It documents the hourly electricity prices, load forecast in Belgium, and generation forecast in France, covering the period from 2011-01-09 to 2016-12-31.
 - **FR** represents the electricity market in France. It records the hourly electricity prices and the corresponding load and generation forecast data, with the time range from 2012-01-09 to 2017-12-31.
 - **DE** corresponds to the German electricity market. It keeps track of the hourly electricity prices, the zonal load forecast in the TSO Amprion zone, and the wind and solar generation forecasts, spanning from 2012-01-09 to 2017-12-31.
- (2) The **Exchange** (Lai et al., 2018) dataset comprises of daily closing exchange rates of eight currencies against the USD from 1990 to 2016.

614

615 616

617 618 619

626 627 628

625

629 630 631

632

633

635

646 647 648

649

650

645

651 652 653

654 655 656

(3) The **Weather**(Zhou et al., 2021) dataset contains 21 meteorological variables recorded every 10 minutes at a weather station in Germany during 2020. In this work, we use the Wet Bulb factor as the target variable to be predicted and the other indicators as exogenous variables

Table 4 provides a summary of the data statistics.

Table 4: Full dataset descriptions. Training/Validation/Test dataset is split as 60%/10%/20%.

Datasets	Look-back period		Target variable	No. of Exogenous Variables	Sampling frequency	
EPF - NP	192	24	Nord Pool Electricity Price	2	1h	
EPF - PJM	192	24	Pennsylvania-New Jersey- Maryland Electricity Price	2	1h	
EPF - BE	192	24	Belgium's Electricity Price	2	1h	
EPF - FR	192	24	France's Electricity Price	2	1h	
EPF - DE	192	24	German's Electricity Price	2	1h	
Exchange	96	{96,192}	Exchange rates	7	1d	
Weather	96	{96,192}	CO ₂ -Concentration	20	10m	

This study employs two distinct downsampling procedures to generate sparse datasets for subsequent model training and inference. The first is a random missing (RM) approach, wherein a fraction α of data points is randomly removed from the original target time series, where α is set to range from 30% to 70%. The forecasting accuracy under different sparsity levels is evaluated in subsequent sections. The second is a block missing (BM) approach. For each sliding window, this method removes a continuous segment of length s from a random position within the window. For instance, from a time series segment of length 96, a contiguous segment of 24 points is removed. Figure 3 illustrates examples of the original time series and the sparsified series resulting from these two methods.

Original TS RM ВМ

Figure 3: Schematic diagrams of RM and BM, where the gray shaded areas represent the missing regions.

B.2 BASELINES

To comprehensively assess the capabilities of FITS, we benchmark it against state-of-the-art approaches, including time series diffusion models, and other leading methods. This diverse set of baselines ensures a rigorous and well-rounded comparison, highlighting FITS's performance across different learning paradigms and demonstrating its effectiveness in a wide range of scenarios.

- (1) Time series diffusion models:
 - CSDI: https://github.com/ermongroup/CSDI. CSDI proposes a novel time series imputation method that leverages score-based diffusion models conditioned on observed data.
 - TMDM: https://github.com/LiYuxin321/TMDM. TMDM introduces a Transformer-Modulated Diffusion Model, uniting conditional diffusion generative process with

transformers into a unified framework to enable precise distribution forecasting for MTS.

Diffusion-TS: https://github.com/Y-debug-sys/Diffusion-TS. DiffusionTS is a diffusion
model-based framework that decomposes time series into trend, seasonality, and residual components, integrates Transformer architectures to capture temporal dependencies, and aims to
produce interpretable and multimodal time series data.

(2) Long time series Forecasting models:

- TiDE:https://github.com/google-research/google-research/blob/master/tide/. TiDE proposes an MLP-based encoder-decoder model for long-term time-series forecasting, which handles covariates and non-linear dependencies.
- DLinear: https://github.com/ioannislivieris/DLinear. DLinear introduces simple one-layer linear models that bypass the temporal information loss inherent in Transformer-based selfattention, achieving superior performance in long-term time series forecasting across diverse datasets.
- Crossformer: https://github.com/Thinklab-SJTU/Crossformer. Crossformer proposes a novel transformer-based model for long-sequence time series forecasting (LSTF), which segments the input into smaller chunks and leveraging cross-attention mechanisms to effectively capture long-range temporal dependencies, thereby enhancing prediction accuracy for extended time horizons.

B.3 IMPLEMENTATION DETAILS

For all datasets, the pseudo length was fixed at 24. The input sequence length for the EPF dataset was set to 192, yielding a target length of 168, while for the other datasets, the input length was fixed at 96. The forecasting horizon was 24 for EPF and $\{96,192\}$ for the remaining datasets. For Quantile Interval Coverage Error(QICE), we divided samples into 10 quantile intervals. For models not inherently designed to handle mismatched lengths between covariates and targets, the target sequences were zero-padded to align dimensions. All implementations were based on PyTorch and executed on an NVIDIA RTX 5090D GPU with 32 GB of memory.

C More Result

C.1 QUALITATIVE ANYLIST AND VISUALIZATION

We visualized the prediction probabilities of the 0th, 200th, 400th, 600th, and 800th samples in the EPF dataset and compared them with the TMDM model, which performed well in the previous results. It can be clearly seen from the visualization results that our model is more advantageous in terms of the accuracy, concentration of the probability distribution and the fit with the ground truth, which fully demonstrates that our model has a strong ability in predicting probabilities.

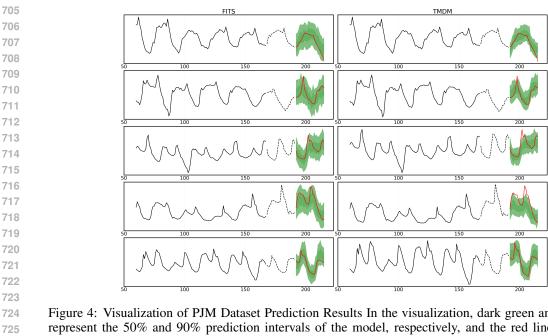


Figure 4: Visualization of PJM Dataset Prediction Results In the visualization, dark green and light green represent the 50% and 90% prediction intervals of the model, respectively, and the red line denotes the ground truth.

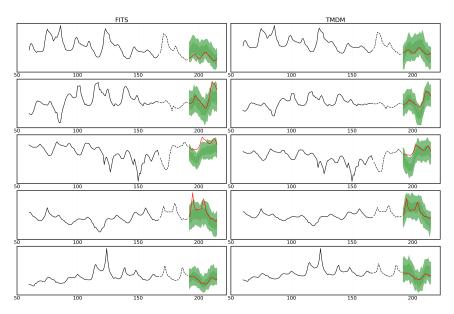


Figure 5: Visualization of NP Dataset Prediction Results In the visualization, dark green and light green represent the 50% and 90% prediction intervals of the model, respectively, and the red line denotes the ground truth.

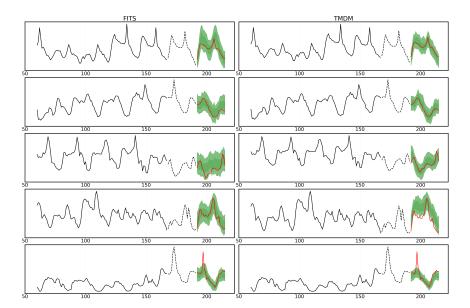


Figure 6: Visualization of FR Dataset Prediction Results In the visualization, dark green and light green represent the 50% and 90% prediction intervals of the model, respectively, and the red line denotes the ground truth.

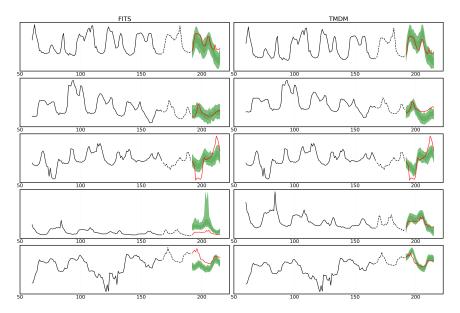


Figure 7: Visualization of DE Dataset Prediction Results In the visualization, dark green and light green represent the 50% and 90% prediction intervals of the model, respectively, and the red line denotes the ground truth.

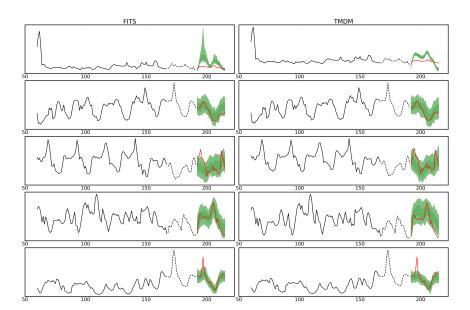


Figure 8: Visualization of BE Dataset Prediction Results In the visualization, dark green and light green represent the 50% and 90% prediction intervals of the model, respectively, and the red line denotes the ground truth.

C.2 PSEUDO LENGTH

We investigated the impact of different pseudo-future lengths (p) on prediction performance. Here, we selected DLinear and Crossformer, which exhibited the best prediction performance, and used two subsets of the EPF dataset as well as two long-term forecasting datasets, respectively. The results indicate that the pseudo-future length has a significant influence on prediction performance. However, a closer examination reveals that our model shows relative insensitivity on the EPF dataset, demonstrating that the attention mechanism we designed effectively addresses the forecasting problem under such circumstances.

Table 5: Performance comparisons in terms of MAE and MSE. The table presents the scenarios random missingness of 0.5.

metric		DLinear					Crossi	former		FITS			
		p=0	p=12	p=24	p=48	p=0	p=12	p=24	p=48	p=0	p=12	p=24	p=48
Weather	MSE	0.584	0.725	0.855	1.026	0.403	0.419	0.486	0.564	0.428	0.457	0.559	0.569
	MAE	0.533	0.604	0.656	0.704	0.446	0.456	0.501	0.542	0.472	0.495	0.547	0.532
Exchange	MSE	0.186	0.215	0.265	0.405	0.235	0.263	0.557	0.621	0.363	0.354	0.395	0.435
	MAE	0.348	0.378	0.416	0.520	0.401	0.429	0.598	0.655	0.435	0.456	0.474	0.513
FR	MSE	0.5150	0.551	0.577	0.605	0.472	0.496	0.529	0.521	0.386	0.392	0.398	0.463
	MAE	0.297	0.323	0.344	0.360	0.234	0.276	0.379	0.321	0.246	0.275	0.296	0.356
BE	MSE MAE	0.5363 0.349	0.599 0.3781	0.628 0.397	0.631 0.412	0.450 0.311	0.492 0.331	0.499 0.302	0.525 0.427	0.375 0.245	0.376 0.238	0.384 0.244	0.412 0.309

D THE USE OF LARGE LANGUAGE MODELS (LLM)

This work used LLMs to fix grammar mistakes and spelling errors.