

MULTI-RESOLUTION DIFFUSION MODELS FOR TIME SERIES FORECASTING

Lifeng Shen¹, Weiyu Chen², James T. Kwok²

¹ Division of Emerging Interdisciplinary Areas, Hong Kong University of Science and Technology

² Department of Computer Science and Engineering, Hong Kong University of Science and Technology
 lshenae@connect.ust.hk, wchenbx@cse.ust.hk, jamesk@cse.ust.hk

ABSTRACT

The diffusion model has been successfully used in many computer vision applications, such as text-guided image generation and image-to-image translation. Recently, there have been attempts on extending the diffusion model for time series data. However, these extensions are fairly straightforward and do not utilize the unique properties of time series data. As different patterns are usually exhibited at multiple scales of a time series, we in this paper leverage this multi-resolution temporal structure and propose the multi-resolution diffusion model (`mr-Diff`). By using the seasonal-trend decomposition, we sequentially extract fine-to-coarse trends from the time series for forward diffusion. The denoising process then proceeds in an easy-to-hard non-autoregressive manner. The coarsest trend is generated first. Finer details are progressively added, using the predicted coarser trends as condition variables. Experimental results on nine real-world time series datasets demonstrate that `mr-Diff` outperforms state-of-the-art time series diffusion models. It is also better than or comparable across a wide variety of advanced time series prediction models.

1 INTRODUCTION

Time series data are prevalent in many real-world applications. In particular, time series forecasting facilitates users to identify patterns and make predictions based on historical data. Examples include stock price prediction in finance, patient health monitoring in healthcare, machine monitoring in manufacturing, and traffic flow optimization in transportation. Over the years, significant advancements in time series analysis have been made through the development of various deep neural networks, including recurrent neural networks (Hewamalage et al., 2021), convolutional neural networks (Yue et al., 2022), and transformers (Vaswani et al., 2017).

Besides these prominent deep neural networks, the diffusion model has recently emerged as a strong generative modeling tool. It has outperformed many other generative models in areas such as image synthesis (Ho et al., 2020; Dhariwal & Nichol, 2021), video generation (Harvey et al., 2022; Blattmann et al., 2023), and multi-modal applications (Rombach et al., 2022; Saharia et al., 2022). Very recently, researchers have sought to leverage its strong generative capacity in the time-series domain. A number of time-series diffusion models have been developed (Rasul et al., 2021; Tashiro et al., 2021; Alcaraz & Strodthoff, 2022; Shen & Kwok, 2023). For example, TimeGrad (Rasul et al., 2021) integrates the standard diffusion model with recurrent neural network’s hidden states. CSDI (Tashiro et al., 2021) uses self-supervised masking to guide a non-autoregressive denoising process. While these time series diffusion models have demonstrated their efficacy, they do not fully utilize the unique structural properties in time series data and are still constrained to generate the time series directly from random vectors (Figure 1(a)). This can pose a significant challenge when working with real-world time series that are non-stationary and noisy.

As time series usually exhibit complex patterns over multiple scales, using the underlying multi-resolution temporal structure has been a cornerstone in traditional time series analysis. In particular, the seasonal-trend decomposition (Robert et al., 1990) can extract the seasonal and trend components, and the coarser temporal patterns can be used to help modeling the finer patterns. Recently, some deep time series prediction models (Oreshkin et al., 2019; Wu et al., 2021; Zeng et al., 2023) have also

incorporated multi-resolution analysis techniques. For example, NBeats (Oreshkin et al., 2019) uses the Fourier and polynomial basis to approximate the seasonal and trend components in multiple layers, respectively. Autoformer (Wu et al., 2021) uses average pooling to extract seasonal components in various transformer layers. DLinear (Zeng et al., 2023) introduces an MLP with seasonal-trend branches. N-Hits (Challu et al., 2023) uses hierarchical interpolation to better leverage the multiscale temporal patterns. Fedformer (Zhou et al., 2022b) improves Autoformer by using mixture-of-experts decomposition in the frequency domain. While these recent transformer and MLP models demonstrate the effectiveness of multi-resolution analysis in deep time series modeling, the use of multi-resolution analysis in time series diffusion models has yet to be explored.

In this paper, we bridge this gap by proposing the multi-resolution diffusion (mr-Diff) model for time series forecasting. Unlike existing time series diffusion models that directly denoises from random vectors (Figure 1(a)), mr-Diff decomposes the denoising objective into several sub-objectives (Figure 1(b)), each of which corresponds to a trend extracted from a sequence of fine-to-coarse seasonal-trend decompositions. This encourages the denoising process to proceed in an easy-to-hard manner. The coarser trends are generated first and the finer details are then progressively added. By better exploiting the seasonal-trend structure and different temporal resolutions, this leads to a more accurate generation of the time series.

The contributions of this paper can be summarized as follows: (i) We propose the multi-resolution diffusion (mr-Diff) model, which is the first to integrate the seasonal-trend decomposition-based multi-resolution analysis into time series diffusion models. (ii) We perform progressive denoising in an easy-to-hard manner, generating coarser signals first and then finer details. This allows a more accurate prediction of time series. (iii) Extensive experiments show that mr-Diff outperforms state-of-the-art time series diffusion models. It is also better than or comparable across a wide variety of advanced time series prediction models.

2 RELATED WORKS: DEEP TIME SERIES MODELS

Recently, a number of deep time series models have been proposed that use the transformer (Vaswani et al., 2017) to capture temporal dependencies. Informer (Zhou et al., 2021) improves the vanilla transformer by avoiding its quadratic time complexity with sparse attention, and improves the inference speed by decoding in a non-autoregressive manner. Autoformer (Wu et al., 2021) replaces the transformer’s self-attention block with an auto-correlation layer. Fedformer (Zhou et al., 2022b) uses a frequency-enhanced module to capture important temporal structures by frequency domain mapping. Pyraformer (Liu et al., 2021) uses a pyramidal attention module for multi-resolution representation of the time series. Scaleformer (Shabani et al., 2023) generates the forecast progressively, starting from the coarser level and then progressing to the finer levels. PatchTST (Nie et al., 2022) is similar to the vision transformer (ViT) (Dosovitskiy et al., 2020), and performs time series prediction by patching the time series and self-supervised pre-training to extract local semantic information. It also replaces the transformer’s decoder with a linear mapping and uses a channel-independence strategy to achieve good performance in multivariate time series prediction.

Besides the transformer-based models, some recent models leverage basis expansion to decompose the time series. FiLM (Zhou et al., 2022a) uses Legendre Polynomials projections to approximate historical information and Fourier projections to remove noise. NBeats (Oreshkin et al., 2019) represents the trends in the time series by polynomial coefficients and the seasonal patterns by Fourier coefficients. Depts (Fan et al., 2022) improves NBeats by using a periodicity module to model periodic time series. N-Hits (Challu et al., 2023) uses multi-scale hierarchical interpolations to further

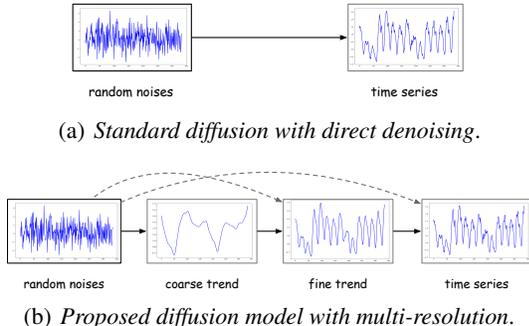


Figure 1: Direct denoising versus proposed multi-resolution denoising.

improve NBeats. In general, these models are easier to train than transformer-based models, though their performance can vary depending on the choice of the basis.

Besides these two types of deep learning models, other recent models are also very competitive. SCINet (Liu et al., 2022) uses a recursive downsample-convolve-interact architecture to extract temporal features from the downsampled sub-sequences or features. NLinear (Zeng et al., 2023) normalizes the time series and uses a linear layer for prediction. DLinear (Zeng et al., 2023) follows the Autoformer and uses seasonal-trend decomposition.

Very recently, diffusion models have also been developed for time series data. TimeGrad (Rasul et al., 2021) is a conditional diffusion model which predicts in an autoregressive manner, with the denoising process guided by the hidden state of a recurrent neural network. However, it suffers from slow inference on long time series because of the use of autoregressive decoding. To alleviate this problem, CSDI (Tashiro et al., 2021) uses non-autoregressive generation, and uses self-supervised masking to guide the denoising process. However, it needs two transformers to capture dependencies in the channel and time dimensions. Moreover, its complexity is quadratic in the number of variables and length of time series as in other transformer models. Besides, masking-based conditioning is similar to the task of image inpainting, and can cause disharmony at the boundaries between the masked and observed regions (Lugmayr et al., 2022; Shen & Kwok, 2023). SSSD (Alcaraz & Strodthoff, 2022) reduces the computational complexity of CSDI by replacing the transformers with a structured state space model. However, it uses the same masking-based conditioning as in CSDI, and thus still suffers from the problem of boundary disharmony. To alleviate this problem, the non-autoregressive diffusion model TimeDiff (Shen & Kwok, 2023) uses future mixup and autoregressive initialization for conditioning. However, all these time series diffusion models do not leverage the multi-resolution temporal structures and denoise directly from random vectors as in standard diffusion models.

In this paper, we propose to decompose the time series into multiple resolutions using seasonal-trend decomposition, and use the fine-to-coarse trends as intermediate latent variables to guide the denoising process. Recently, other multiresolution analysis techniques besides using seasonal-trend decomposition have also been used for time series modeling. For example, Yu et al. (2021) propose a U-Net (Ronneberger et al., 2015) for graph-structured time series, and leverage temporal information from different resolutions by pooling and unpooling. Mu2ReST (Niu et al., 2022) works on spatio-temporal data and recursively outputs predictions from coarser to finer resolutions. Yformer (Madhusudhanan et al., 2021) captures temporal dependencies by combining downscaling/upsampling with sparse attention. PSA-GAN (Jeha et al., 2022) trains a growing U-Net, and captures multi-resolution patterns by progressively adding trainable modules at different levels. However, all these methods need to design very specific U-Net structures.

3 BACKGROUND

3.1 DENOISING DIFFUSION PROBABILISTIC MODELS

A well-known diffusion model is the denoising diffusion probabilistic model (DDPM) (Ho et al., 2020). It is a latent variable model with forward diffusion and backward denoising processes. During forward diffusion, an input \mathbf{x}^0 is gradually corrupted to a Gaussian noise vector. Specifically, at the k th step, \mathbf{x}^k is generated by corrupting the previous iterate \mathbf{x}^{k-1} (scaled by $\sqrt{1 - \beta_k}$) with zero-mean Gaussian noise (with variance $\beta_k \in [0, 1]$):

$$q(\mathbf{x}^k | \mathbf{x}^{k-1}) = \mathcal{N}(\mathbf{x}^k; \sqrt{1 - \beta_k} \mathbf{x}^{k-1}, \beta_k \mathbf{I}), \quad k = 1, \dots, K.$$

It can be shown that this can also be rewritten as $q(\mathbf{x}^k | \mathbf{x}^0) = \mathcal{N}(\mathbf{x}^k; \sqrt{\bar{\alpha}_k} \mathbf{x}^0, (1 - \bar{\alpha}_k) \mathbf{I})$, where $\bar{\alpha}_k = \prod_{s=1}^k \alpha_s$, and $\alpha_k = 1 - \beta_k$. Thus, \mathbf{x}^k can be simply obtained as

$$\mathbf{x}^k = \sqrt{\bar{\alpha}_k} \mathbf{x}^0 + \sqrt{1 - \bar{\alpha}_k} \epsilon, \quad (1)$$

where ϵ is a noise from $\mathcal{N}(\mathbf{0}, \mathbf{I})$. This equation also allows \mathbf{x}^0 to be easily recovered from \mathbf{x}^k .

In DDPM, backward denoising is defined as a Markovian process. Specifically, at the k th denoising step, \mathbf{x}^{k-1} is generated from \mathbf{x}^k by sampling from the following normal distribution:

$$p_\theta(\mathbf{x}^{k-1} | \mathbf{x}^k) = \mathcal{N}(\mathbf{x}^{k-1}; \mu_\theta(\mathbf{x}^k, k), \Sigma_\theta(\mathbf{x}^k, k)). \quad (2)$$

Here, the variance $\Sigma_\theta(\mathbf{x}^k, k)$ is usually fixed as $\sigma_k^2 \mathbf{I}$, while the mean $\mu_\theta(\mathbf{x}^k, k)$ is defined by a neural network (parameterized by θ). This is usually formulated as a noise estimation or data prediction problem (Benny & Wolf, 2022). For noise estimation, a network ϵ_θ predicts the noise of the diffused input \mathbf{x}^k , and then obtains $\mu_\theta(\mathbf{x}^k, k)$ as $\frac{1}{\sqrt{\alpha_k}} \mathbf{x}^k - \frac{1-\alpha_k}{\sqrt{1-\bar{\alpha}_k} \sqrt{\alpha_k}} \epsilon_\theta(\mathbf{x}^k, k)$. Parameter θ is learned by minimizing the loss $\mathcal{L}_\epsilon = \mathbb{E}_{k, \mathbf{x}^0, \epsilon} [\|\epsilon - \epsilon_\theta(\mathbf{x}^k, k)\|^2]$. Alternatively, the data prediction strategy uses a denoising network \mathbf{x}_θ to obtain an estimate $\mathbf{x}_\theta(\mathbf{x}^k, k)$ of the clean data \mathbf{x}^0 given \mathbf{x}^k , and then set

$$\mu_\theta(\mathbf{x}^k, k) = \frac{\sqrt{\alpha_k}(1 - \bar{\alpha}_{k-1})}{1 - \bar{\alpha}_k} \mathbf{x}^k + \frac{\sqrt{\bar{\alpha}_{k-1}} \beta_k}{1 - \bar{\alpha}_k} \mathbf{x}_\theta(\mathbf{x}^k, k). \quad (3)$$

Parameter θ is learned by minimizing the loss

$$\mathcal{L}_\mathbf{x} = \mathbb{E}_{\mathbf{x}^0, \epsilon, k} \|\mathbf{x}^0 - \mathbf{x}_\theta(\mathbf{x}^k, k)\|^2. \quad (4)$$

3.2 CONDITIONAL DIFFUSION MODELS FOR TIME SERIES PREDICTION

In time series forecasting, one aims to predict the future values $\mathbf{x}_{1:H}^0 \in \mathbb{R}^{d \times H}$ given the past observations $\mathbf{x}_{-L+1:0}^0 \in \mathbb{R}^{d \times L}$ of the time series. Here, d is the number of variables, H is the length of the forecast window, and L is the length of the lookback window. When using conditional diffusion models for time series prediction, the following distribution is considered (Rasul et al., 2021; Tashiro et al., 2021; Shen & Kwok, 2023)

$$p_\theta(\mathbf{x}_{1:H}^{0:K} | \mathbf{c}) = p_\theta(\mathbf{x}_{1:H}^K | \mathbf{c}) \prod_{k=1}^K p_\theta(\mathbf{x}_{1:H}^{k-1} | \mathbf{x}_{1:H}^k, \mathbf{c}), \quad \mathbf{c} = \mathcal{F}(\mathbf{x}_{-L+1:0}^0), \quad (5)$$

where $\mathbf{x}_{1:H}^K \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, \mathbf{c} is the condition, and \mathcal{F} is a conditioning network that takes the past observations $\mathbf{x}_{-L+1:0}^0$ as input. Correspondingly, the denoising process at step k is given by

$$p_\theta(\mathbf{x}_{1:H}^{k-1} | \mathbf{x}_{1:H}^k, \mathbf{c}) = \mathcal{N}(\mathbf{x}_{1:H}^{k-1}; \mu_\theta(\mathbf{x}_{1:H}^k, k | \mathbf{c}), \sigma_k^2 \mathbf{I}), \quad k = K, K-1, \dots, 1. \quad (6)$$

During inference, we denote the generated sample corresponding to $\mathbf{x}_{1:H}^k$ by $\hat{\mathbf{x}}_{1:H}^k$. We first initialize $\hat{\mathbf{x}}_{1:H}^K$ as a noise vector from $\mathcal{N}(\mathbf{0}, \mathbf{I})$. By repeatedly running the denoising step in (6) till $k = 1$, the final generated sample is $\hat{\mathbf{x}}_{1:H}^0$.

4 MR-DIFF: MULTI-RESOLUTION DIFFUSION MODEL

As discussed in Section 1, recent transformer and MLP models demonstrate the effectiveness of seasonal-trend decomposition-based multi-resolution analysis in deep time series modeling. However, the use of multi-resolution temporal patterns in the diffusion model has yet to be explored. In this paper, we address this gap by proposing the multi-resolution diffusion (`mr-Diff`) model. An overview of the proposed model is shown in Figure 2.

The proposed `mr-Diff` can be viewed as a cascaded diffusion model (Ho et al., 2022), and proceeds in S stages, with the resolution getting coarser as the stage proceeds (Section 4.1). This allows capturing the temporal dynamics at multiple temporal resolutions. In each stage, the diffusion process is interleaved with seasonal-trend decomposition. For simplicity of notations, we use $\mathbf{X} = \mathbf{x}_{-L+1:0}$ and $\mathbf{Y} = \mathbf{x}_{1:H}$ for the time series segments in the lookback and forecast windows, respectively. Let the trend component of the lookback (resp. forecast) segment at stage $s+1$ be \mathbf{X}_s (resp. \mathbf{Y}_s). The trend gets coarser as s increases, and with $\mathbf{X}_0 = \mathbf{X}$ and $\mathbf{Y}_0 = \mathbf{Y}$. In each stage $s+1$, a conditional diffusion model is learned to reconstruct the trend component \mathbf{Y}_s extracted from the forecast window (Section 4.2). The reconstruction at stage 1 then corresponds to the target time series forecast.

While the forward diffusion process in this diffusion model is straightforward and similar to existing diffusion models, design of the denoising process, particularly on the denoising conditions and denoising network, are less trivial. During training, to guide the reconstruction of \mathbf{Y}_s , the proposed model takes the lookback segment \mathbf{X}_s (which has the same resolution as \mathbf{Y}_s) and the coarser trend \mathbf{Y}_{s+1} (which provides an overall picture of the finer \mathbf{Y}_s) as denoising condition. On inference, the ground-truth \mathbf{Y}_{s+1} is not available, and is replaced by its estimate $\hat{\mathbf{Y}}_{s+1}^0$ produced by the denoising process at stage $s+1$. By combining the diffusion model and seasonal-trend decomposition in a multi-resolution manner, the proposed model encourages a better modeling of real-world time series.

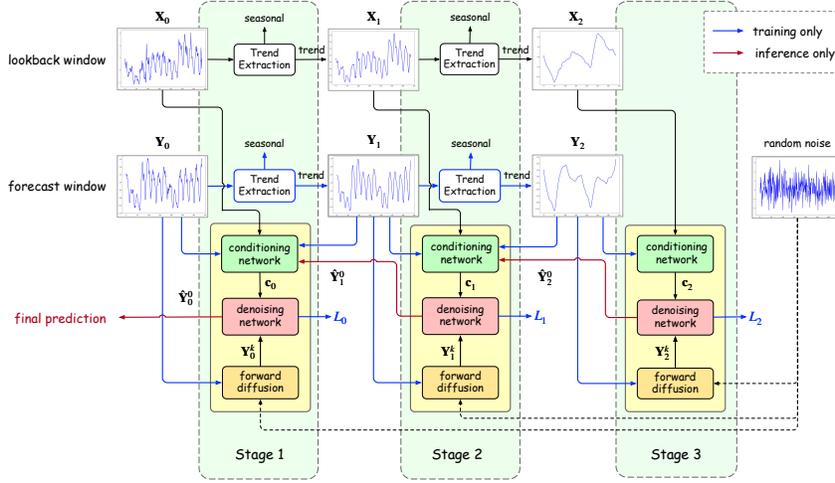


Figure 2: The proposed multi-resolution diffusion model $mr\text{-Diff}$. For simplicity of illustration, we use $S = 3$ stages. At stage $s + 1$, \mathbf{Y}_s denotes the corresponding trend component extracted from the segment in the forecast window, \mathbf{Y}_s^k is the diffusion sample at diffusion step k , and $\hat{\mathbf{Y}}_s^0$ is the denoised output.

4.1 EXTRACTING FINE-TO-COARSE TRENDS

For the given time series segment \mathbf{X}_0 in the lookback window, its trend components are successively extracted by the TrendExtraction module as:

$$\mathbf{X}_s = \text{AvgPool}(\text{Padding}(\mathbf{X}_{s-1}), \tau_s), \quad s = 1, \dots, S - 1,$$

where AvgPool is the average pooling operation (Wu et al., 2021), Padding keeps the lengths of \mathbf{X}_{s-1} and \mathbf{X}_s the same, and τ_s is the smoothing kernel size which increases with s so as to generate fine-to-coarse trends. Processing for the segment \mathbf{Y}_0 in the forecast window is analogous, and its trend components $\{\mathbf{Y}_s\}_{s=1, \dots, S-1}$ are extracted.

Note that while the seasonal-trend decomposition obtains both the seasonal and trend components, the focus here is on the trend. On the other hand, models such as the Autoformer (Wu et al., 2021) and Fedformer (Zhou et al., 2022b) focus on progressively decomposing the seasonal component. As we use the diffusion model for time series reconstruction at various stages/resolutions (Section 4.2), intuitively, it is easier to predict a finer trend from a coarser trend. On the other hand, reconstruction of a finer seasonal component from a coarser seasonal component may be difficult, especially as the seasonal component may not present clear patterns.

4.2 TEMPORAL MULTI-RESOLUTION RECONSTRUCTION

In each stage $s + 1$, we use a conditional diffusion model to reconstruct the future trend \mathbf{Y}_s extracted in Section 4.1. As in the standard diffusion model, it consists of a forward diffusion process and a backward denoising process. Forward diffusion does not involve learnable parameters, while the training procedure of the backward denoising process needs optimization as shown in Algorithm 1.

An d' -dimensional embedding \mathbf{p}^k of the diffusion step k is used in both the forward and backward denoising processes. As in (Rasul et al., 2021; Tashiro et al., 2021; Kong et al., 2020), this is obtained by first taking the sinusoidal position embedding (Vaswani et al., 2017): $k_{\text{embedding}} = \left[\sin(10^{\frac{0 \times 4}{w-1} t}), \dots, \sin(10^{\frac{w \times 4}{w-1} t}), \cos(10^{\frac{0 \times 4}{w-1} t}), \dots, \cos(10^{\frac{w \times 4}{w-1} t}) \right]$ where $w = \frac{d'}{2}$, and then passing it through two fully-connected (FC) layers to obtain

$$\mathbf{p}^k = \text{SiLU}(\text{FC}(\text{SiLU}(\text{FC}(k_{\text{embedding}})))), \quad (7)$$

where SiLU is the sigmoid-weighted linear unit (Elfwing et al., 2017). By default, d' is set to 128.

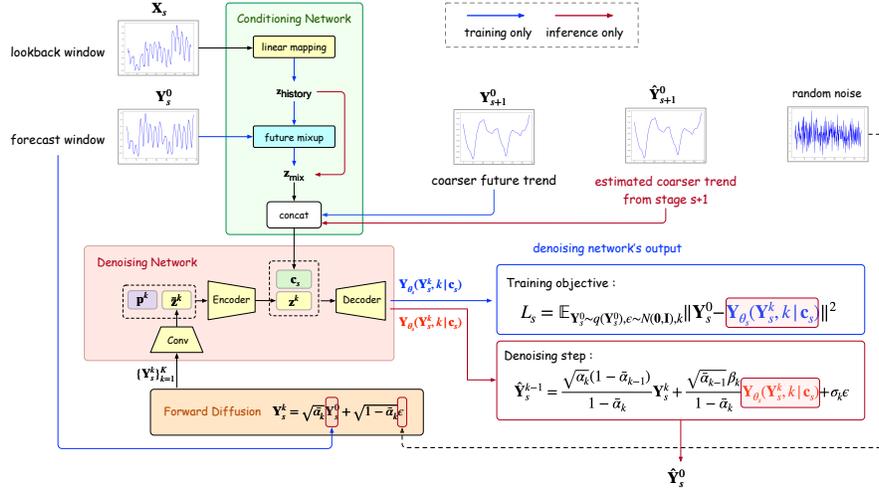


Figure 3: The conditioning network and denoising network.

4.2.1 FORWARD DIFFUSION

Forward diffusion is straightforward. Analogous to (1), with $\mathbf{Y}_s^0 = \mathbf{Y}_s$, we obtain at step k ,

$$\mathbf{Y}_s^k = \sqrt{\alpha_k} \mathbf{Y}_s^0 + \sqrt{1 - \alpha_k} \epsilon, \quad k = 1, \dots, K, \quad (8)$$

where the noise matrix ϵ is sampled from $\mathcal{N}(\mathbf{0}, \mathbf{I})$ with the same size as \mathbf{Y}_s .

4.2.2 BACKWARD DENOISING

Standard diffusion models perform one-stage denoising directly from random vectors. In this section, we decompose the denoising objective of `mr-Diff`, into S sub-objectives, one for each stage. As will be seen, this encourages the denoising process to proceed in an easy-to-hard manner, such that the coarser trends are generated first and the finer details are then progressively added.

Conditioning network. The conditioning network (Figure 3, top) constructs a condition to guide the denoising network (to be discussed). Existing time series diffusion models (Rasul et al., 2021; Tashiro et al., 2021) simply use the original time series’ lookback segment \mathbf{X}_0 as condition \mathbf{c} in (5). With the proposed multi-resolution seasonal-trend decomposition, we instead use the lookback segment \mathbf{X}_s at the same decomposition stage s . This allows better and easier reconstruction, as \mathbf{X}_s has the same resolution as \mathbf{Y}_s to be reconstructed. On the other hand, when \mathbf{X}_0 is used as in existing time series diffusion models, the denoising network may overfit temporal details at the finer level.

A linear mapping is applied on \mathbf{X}_s to produce a tensor $\mathbf{z}_{\text{history}} \in \mathbb{R}^{d \times H}$. For better denoising performance, during training we use future-mixup (Shen & Kwok, 2023) to enhance $\mathbf{z}_{\text{history}}$. It combines $\mathbf{z}_{\text{history}}$ and the (ground-truth) future observation \mathbf{Y}_s^0 with mixup (Zhang et al., 2018):

$$\mathbf{z}_{\text{mix}} = \mathbf{m} \odot \mathbf{z}_{\text{history}} + (1 - \mathbf{m}) \odot \mathbf{Y}_s^0, \quad (9)$$

where \odot denotes the Hadamard product, and $\mathbf{m} \in [0, 1]^{d \times H}$ is a mixing matrix in which each element is randomly sampled from the uniform distribution on $[0, 1]$. Future-mixup is similar to teacher forcing (Williams & Zipser, 1989), which mixes the ground truth with previous prediction output at each decoding step of autoregressive decoding.

Besides using \mathbf{X}_s on training, the coarser trend \mathbf{Y}_{s+1} ($= \mathbf{Y}_{s+1}^0$) can provide an overall picture of the finer trend \mathbf{Y}_s and is thus also useful for conditioning. Hence, \mathbf{z}_{mix} in (9) is concatenated with \mathbf{Y}_{s+1}^0 along the channel dimension to produce the condition \mathbf{c}_s (a $2d \times H$ tensor). For $s = S$ (the last stage), there is no coarser trend and \mathbf{c}_s is simply equal to \mathbf{z}_{mix} .

On inference, the ground-truth \mathbf{Y}_s^0 is no longer available. Hence, future-mixup in (9) is not used, and we simply set $\mathbf{z}_{\text{mix}} = \mathbf{z}_{\text{history}}$. Moreover, the coarser trend \mathbf{Y}_{s+1} is also not available, and we concatenate \mathbf{z}_{mix} with the estimate $\hat{\mathbf{Y}}_{s+1}^0$ generated from stage $s + 2$ instead.

Denoising network. Analogous to (6), the denoising process at step k of stage $s + 1$ is given by

$$p_{\theta_s}(\mathbf{Y}_s^{k-1} | \mathbf{Y}_s^k, \mathbf{c}_s) = \mathcal{N}(\mathbf{Y}_s^{k-1}; \mu_{\theta_s}(\mathbf{Y}_s^k, k | \mathbf{c}_s, \sigma_k^2 \mathbf{I})), \quad k = K, \dots, 1, \quad (10)$$

where θ_s includes all parameters in the conditional network and denoising network at stage $s + 1$, and

$$\mu_{\theta_s}(\mathbf{Y}_s^k, k | \mathbf{c}_s, \sigma_k^2 \mathbf{I}) = \frac{\sqrt{\alpha_k}(1 - \bar{\alpha}_{k-1})}{1 - \bar{\alpha}_k} \mathbf{Y}_s^k + \frac{\sqrt{\bar{\alpha}_{k-1}}\beta_k}{1 - \bar{\alpha}_k} \mathbf{Y}_{\theta_s}(\mathbf{Y}_s^k, k | \mathbf{c}_s). \quad (11)$$

As in (3), $\mathbf{Y}_{\theta_s}(\mathbf{Y}_s^k, k | \mathbf{c}_s)$ is an estimate of \mathbf{Y}_s^0 .

The denoising network (also shown in Figure 3) outputs $\mathbf{Y}_{\theta_s}(\mathbf{Y}_s^k, k | \mathbf{c}_s)$ with guidance from the condition \mathbf{c}_s (output by the conditioning network). Specifically, it first maps \mathbf{Y}_s^k to the embedding $\bar{\mathbf{z}}^k \in \mathbb{R}^{d' \times H}$ by an input projection block consisting of several convolutional layers. This $\bar{\mathbf{z}}^k$, together with the diffusion-step k 's embedding $\mathbf{p}^k \in \mathbb{R}^{d'}$ in (7), is fed to an encoder (which is a convolution network) to obtain the representation $\mathbf{z}^k \in \mathbb{R}^{d'' \times H}$. Next, we concatenate \mathbf{z}^k and \mathbf{c}_s along the variable dimension to form a tensor of size $(2d + d'') \times H$. This is then fed to a decoder, which is also a convolution network, and outputs $\mathbf{Y}_{\theta_s}(\mathbf{Y}_s^k, k | \mathbf{c}_s)$. Finally, analogous to (4), θ_s is obtained by minimizing the following denoising objective

$$\min_{\theta_s} \mathcal{L}_s(\theta_s) = \min_{\theta_s} \mathbb{E}_{\mathbf{Y}_s^0 \sim q(\mathbf{Y}_s), \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), k} \|\mathbf{Y}_s^0 - \mathbf{Y}_{\theta_s}(\mathbf{Y}_s^k, k | \mathbf{c}_s)\|^2. \quad (12)$$

On inference, for each $s = S, \dots, 1$, we start from $\hat{\mathbf{Y}}_s^K \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ (step 9 in Algorithm 2). Based on the data prediction strategy in (11), each denoising step from $\hat{\mathbf{Y}}_s^k$ (an estimate of \mathbf{Y}_s^k) to $\hat{\mathbf{Y}}_s^{k-1}$ is:

$$\hat{\mathbf{Y}}_s^{k-1} = \frac{\sqrt{\alpha_k}(1 - \bar{\alpha}_{k-1})}{1 - \bar{\alpha}_k} \hat{\mathbf{Y}}_s^k + \frac{\sqrt{\bar{\alpha}_{k-1}}\beta_k}{1 - \bar{\alpha}_k} \mathbf{Y}_{\theta_s}(\hat{\mathbf{Y}}_s^k, k | \mathbf{c}_s) + \sigma_k \epsilon, \quad (13)$$

where $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ when $k > 1$, and $\epsilon = 0$ otherwise.

The pseudocode for the training and inference procedures of the backward denoising process can be found in Appendix A.

Discussion. The proposed `mr-Diff` is related to the Scaleformer (Shabani et al., 2023), which also generates the forecast progressively from the coarser level to the finer ones. However, besides the obvious difference that Scaleformer is based on the transformer while `mr-Diff` is based on the diffusion model, Scaleformer uses one single forecasting module (a transformer) at all resolutions. However, as the time series patterns at different resolutions may be different, in `mr-Diff`, each resolution has its own denoising network and is thus more flexible. Moreover, in moving from a lower resolution to a higher resolution during both training and inference, Scaleformer needs linear interpolation and an additional cross-scale normalization mechanism. On the other hand, in `mr-Diff`, the lower-resolution prediction is fed into the conditioning network as a condition. The conditioning network can learn a more flexible nonlinear mapping to fuse the cross-scale trends. Coherent probabilistic forecasting (Rangapuram et al., 2023) also uses a temporal hierarchy to produce forecasts at multiple resolutions. However, it is more stringent than `mr-Diff` and requires the model to produce consistent forecasts at all resolutions. As will be shown in Section 5, this enables `mr-Diff` to have better empirical performance.

5 EXPERIMENTS

In this section, we conduct time series prediction experiments by comparing 22 recent strong prediction models on 9 popular real-world time series datasets. Due to space constraints, detailed introductions about the datasets and selected baselines are in Appendix B and Appendix C, respectively. For performance evaluation, we use the mean absolute error (MAE) and mean squared error (MSE) averaged over ten randomly sampled trajectories. Because of the lack of space, results on MSE are in Appendix D. Furthermore, ablation studies are provided in Appendix E.

Implementation Details. We train the proposed model using Adam with a learning rate of 10^{-3} . The batch size is 64, and training with early stopping for a maximum of 100 epochs. $K = 100$ diffusion steps are used, with a linear variance schedule (Rasul et al., 2021) starting from $\beta_1 = 10^{-4}$ to $\beta_K = 10^{-1}$. $S = 5$ stages are used. The history length (in {96, 192, 336, 720, 1440}) is selected by using the validation set. All experiments are run on an Nvidia RTX A6000 GPU with 48GB memory. More details can be found in Appendix F.

Table 1: Univariate prediction MAEs on the real-world time series datasets (subscript is the rank). Results of all baselines (except NLinear, SCINet, and N-Hits) are from (Shen & Kwok, 2023).

	<i>NorPool</i>	<i>Caiso</i>	<i>Traffic</i>	<i>Electricity</i>	<i>Weather</i>	<i>Exchange</i>	<i>ETTh1</i>	<i>ETTm1</i>	<i>Wind</i>	avg rank
mr-Diff	0.609 ₍₂₎	0.212 ₍₄₎	0.197 ₍₂₎	0.332 ₍₁₎	0.032 ₍₁₎	0.094 ₍₁₎	0.196 ₍₁₎	0.149 ₍₁₎	1.168 ₍₂₎	1.7
TimeDiff	0.613 ₍₃₎	0.209 ₍₃₎	0.207 ₍₃₎	0.341 ₍₃₎	0.035 ₍₄₎	0.102 ₍₇₎	0.202 ₍₂₎	0.154 ₍₆₎	1.209 ₍₅₎	4.0
TimeGrad	0.841 ₍₂₃₎	0.386 ₍₂₂₎	0.894 ₍₂₃₎	0.898 ₍₂₃₎	0.036 ₍₆₎	0.155 ₍₂₁₎	0.212 ₍₈₎	0.167 ₍₁₂₎	1.239 ₍₁₁₎	16.6
CSDI	0.763 ₍₂₀₎	0.282 ₍₁₄₎	0.468 ₍₂₀₎	0.540 ₍₁₉₎	0.037 ₍₇₎	0.200 ₍₂₃₎	0.221 ₍₁₂₎	0.170 ₍₁₄₎	1.218 ₍₇₎	15.1
SSSD	0.770 ₍₂₁₎	0.263 ₍₁₂₎	0.226 ₍₆₎	0.403 ₍₈₎	0.041 ₍₁₁₎	0.118 ₍₁₆₎	0.250 ₍₂₀₎	0.169 ₍₁₃₎	1.356 ₍₂₂₎	14.3
D ³ VAE	0.774 ₍₂₂₎	0.613 ₍₂₃₎	0.237 ₍₉₎	0.539 ₍₁₈₎	0.039 ₍₉₎	0.107 ₍₁₃₎	0.221 ₍₁₂₎	0.160 ₍₉₎	1.321 ₍₁₉₎	14.9
CPF	0.710 ₍₁₅₎	0.338 ₍₁₈₎	0.385 ₍₁₉₎	0.592 ₍₂₁₎	0.035 ₍₄₎	0.094 ₍₁₎	0.221 ₍₁₂₎	0.153 ₍₅₎	1.256 ₍₁₂₎	11.9
PSA-GAN	0.623 ₍₅₎	0.250 ₍₈₎	0.355 ₍₁₇₎	0.373 ₍₆₎	0.139 ₍₂₂₎	0.109 ₍₁₄₎	0.225 ₍₁₆₎	0.174 ₍₁₆₎	1.287 ₍₁₆₎	13.3
N-Hits	0.646 ₍₇₎	0.276 ₍₁₃₎	0.232 ₍₇₎	0.419 ₍₉₎	0.033 ₍₂₎	0.100 ₍₅₎	0.228 ₍₁₇₎	0.157 ₍₈₎	1.256 ₍₁₂₎	8.9
FiLM	0.654 ₍₉₎	0.290 ₍₁₅₎	0.315 ₍₁₄₎	0.362 ₍₅₎	0.069 ₍₁₄₎	0.104 ₍₁₀₎	0.210 ₍₆₎	0.149 ₍₁₎	1.189 ₍₃₎	8.6
Depts	0.616 ₍₄₎	0.205 ₍₁₎	0.241 ₍₁₀₎	0.434 ₍₁₂₎	0.102 ₍₁₉₎	0.106 ₍₁₂₎	0.202 ₍₂₎	0.165 ₍₁₀₎	1.472 ₍₂₃₎	10.3
NBeats	0.671 ₍₁₀₎	0.228 ₍₅₎	0.225 ₍₅₎	0.439 ₍₁₃₎	0.130 ₍₂₁₎	0.096 ₍₃₎	0.242 ₍₁₈₎	0.165 ₍₁₀₎	1.236 ₍₉₎	10.4
Scaleformer	0.687 ₍₁₂₎	0.320 ₍₁₆₎	0.375 ₍₁₈₎	0.430 ₍₁₀₎	0.083 ₍₁₇₎	0.148 ₍₁₉₎	0.302 ₍₂₂₎	0.210 ₍₂₂₎	1.348 ₍₂₁₎	17.4
PatchTST	0.590 ₍₁₎	0.260 ₍₁₁₎	0.269 ₍₁₁₎	0.478 ₍₁₇₎	0.098 ₍₁₈₎	0.111 ₍₁₅₎	0.260 ₍₂₁₎	0.174 ₍₁₆₎	1.338 ₍₂₀₎	14.4
FedFormer	0.725 ₍₁₇₎	0.254 ₍₉₎	0.278 ₍₁₂₎	0.453 ₍₁₄₎	0.057 ₍₁₃₎	0.168 ₍₂₂₎	0.212 ₍₈₎	0.195 ₍₂₀₎	1.271 ₍₁₄₎	14.3
Autoformer	0.755 ₍₁₉₎	0.339 ₍₁₉₎	0.495 ₍₂₁₎	0.623 ₍₂₂₎	0.040 ₍₁₀₎	0.152 ₍₂₀₎	0.220 ₍₁₁₎	0.174 ₍₁₆₎	1.319 ₍₁₈₎	17.3
Pyraformer	0.747 ₍₁₈₎	0.257 ₍₁₀₎	0.215 ₍₄₎	0.455 ₍₁₅₎	0.107 ₍₂₀₎	0.104 ₍₁₀₎	0.211 ₍₇₎	0.179 ₍₁₉₎	1.284 ₍₁₅₎	13.1
Informer	0.698 ₍₁₃₎	0.345 ₍₂₀₎	0.308 ₍₁₃₎	0.433 ₍₁₁₎	0.069 ₍₁₄₎	0.118 ₍₁₆₎	0.212 ₍₈₎	0.172 ₍₁₅₎	1.236 ₍₉₎	13.2
Transformer	0.723 ₍₁₆₎	0.345 ₍₂₀₎	0.336 ₍₁₆₎	0.469 ₍₁₆₎	0.071 ₍₁₆₎	0.103 ₍₉₎	0.247 ₍₁₉₎	0.196 ₍₂₁₎	1.212 ₍₆₎	15.4
SCINet	0.653 ₍₈₎	0.244 ₍₇₎	0.322 ₍₁₅₎	0.377 ₍₇₎	0.037 ₍₇₎	0.101 ₍₆₎	0.205 ₍₅₎	0.150 ₍₄₎	1.167 ₍₁₎	6.7
NLinear	0.637 ₍₆₎	0.238 ₍₆₎	0.192 ₍₁₎	0.334 ₍₂₎	0.033 ₍₂₎	0.097 ₍₄₎	0.203 ₍₄₎	0.149 ₍₁₎	1.197 ₍₄₎	3.3
DLinear	0.671 ₍₁₀₎	0.206 ₍₂₎	0.236 ₍₈₎	0.348 ₍₄₎	0.310 ₍₂₃₎	0.102 ₍₇₎	0.222 ₍₁₅₎	0.155 ₍₇₎	1.221 ₍₈₎	9.3
LSTMa	0.707 ₍₁₄₎	0.333 ₍₁₇₎	0.757 ₍₂₂₎	0.557 ₍₂₀₎	0.053 ₍₁₂₎	0.136 ₍₁₈₎	0.332 ₍₂₃₎	0.239 ₍₂₃₎	1.298 ₍₁₇₎	18.4

5.1 MAIN RESULTS

Table 1 shows the MAEs on the univariate time series. As can be seen, the proposed `mr-Diff` is the best in 5 of the 9 datasets. The improvement is particularly significant on the more complicated datasets, such as *Exchange* and *ETTh1*. On the remaining 4 datasets, `mr-Diff` ranks second in 3 of them. Overall, its average ranking is better than all other baselines (which include the most recent diffusion models). Note that there is no improvement on datasets such as *Caiso*. This is because there is no complex trend information on this dataset that can be leveraged (as can be seen in Figure 5(b) Appendix B).

Figure 4 shows example prediction results on *ETTh1* by `mr-Diff` and three competitive models: SCINet, NLinear and the recent diffusion model TimeDiff. As can be seen, `mr-Diff` produces higher-quality predictions than the others.

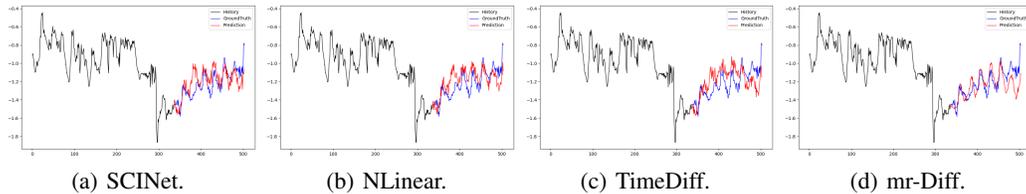


Figure 4: Visualizations on *ETTh1* by (a) SCINet (Liu et al., 2022), (b) NLinear (Zeng et al., 2023), (c) TimeDiff (Shen & Kwok, 2023) and (d) the proposed `mr-Diff`.

Table 2 shows the MAEs on the multivariate time series. Baselines N-Hits, FiLM and SCINet have strong performance because they also leverage multiresolution information in prediction. However, the proposed `mr-Diff` is still very competitive compared to these strong baselines. It ranks among the top-2 methods on 5 of the 9 datasets, and is the best overall. This is then followed by the recent diffusion-based model TimeDiff.

5.2 INFERENCE EFFICIENCY

In this experiment, we compare the inference efficiency of `mr-Diff` (with 1-4 stages) with the other time series diffusion model baselines (TimeDiff, TimeGrad, CSDI, SSSD) on the univariate *ETTh1* dataset. Due to space constraints, we compare the training efficiency in Appendix G. Besides using a

Table 2: Multivariate prediction MAEs on the real-world time series datasets (subscript is the rank). CSDI runs out of memory on *Traffic* and *Electricity*. Results of all baselines (except for NLinear, SCINet, and N-Hits) are from (Shen & Kwok, 2023).

	<i>NorPool</i>	<i>Caiso</i>	<i>Traffic</i>	<i>Electricity</i>	<i>Weather</i>	<i>Exchange</i>	<i>ETTh1</i>	<i>ETTm1</i>	<i>Wind</i>	avg rank
mr-Diff	0.604 ₍₂₎	0.219 ₍₃₎	0.320 ₍₅₎	0.252 ₍₃₎	0.324 ₍₂₎	0.082 ₍₃₎	0.422 ₍₂₎	0.373 ₍₂₎	0.675 ₍₁₎	2.6
TimeDiff	0.611 ₍₃₎	0.234 ₍₆₎	0.384 ₍₈₎	0.305 ₍₆₎	0.312 ₍₁₎	0.091 ₍₇₎	0.430 ₍₃₎	0.372 ₍₁₎	0.687 ₍₃₎	4.2
TimeGrad	0.821 ₍₁₉₎	0.339 ₍₁₈₎	0.849 ₍₂₂₎	0.630 ₍₂₁₎	0.381 ₍₁₄₎	0.193 ₍₁₉₎	0.719 ₍₂₂₎	0.605 ₍₂₂₎	0.793 ₍₂₁₎	19.8
CSDI	0.777 ₍₁₇₎	0.345 ₍₁₉₎	-	-	0.374 ₍₁₂₎	0.194 ₍₂₀₎	0.438 ₍₅₎	0.442 ₍₁₅₎	0.741 ₍₁₀₎	14.0
SSSD	0.753 ₍₁₃₎	0.295 ₍₁₀₎	0.398 ₍₁₃₎	0.363 ₍₁₂₎	0.350 ₍₈₎	0.127 ₍₁₂₎	0.561 ₍₁₇₎	0.406 ₍₁₀₎	0.778 ₍₁₈₎	12.6
D ³ VAE	0.692 ₍₉₎	0.331 ₍₁₆₎	0.483 ₍₁₇₎	0.372 ₍₁₄₎	0.380 ₍₁₃₎	0.301 ₍₂₂₎	0.502 ₍₁₄₎	0.391 ₍₈₎	0.779 ₍₁₉₎	14.7
CPF	0.889 ₍₂₁₎	0.424 ₍₂₁₎	0.714 ₍₂₁₎	0.643 ₍₂₂₎	0.781 ₍₂₃₎	0.082 ₍₃₎	0.597 ₍₂₀₎	0.472 ₍₁₆₎	0.757 ₍₁₅₎	18.0
PSA-GAN	0.890 ₍₂₂₎	0.477 ₍₂₂₎	0.697 ₍₂₀₎	0.533 ₍₂₀₎	0.578 ₍₂₂₎	0.087 ₍₆₎	0.546 ₍₁₆₎	0.488 ₍₁₈₎	0.756 ₍₁₃₎	17.7
N-Hits	0.643 ₍₇₎	0.221 ₍₄₎	0.268 ₍₂₎	0.245 ₍₂₎	0.335 ₍₄₎	0.085 ₍₅₎	0.480 ₍₈₎	0.388 ₍₆₎	0.734 ₍₈₎	5.1
FiLM	0.646 ₍₈₎	0.278 ₍₈₎	0.398 ₍₁₃₎	0.320 ₍₈₎	0.336 ₍₅₎	0.079 ₍₁₎	0.436 ₍₄₎	0.374 ₍₃₎	0.717 ₍₅₎	6.1
Depts	0.611 ₍₃₎	0.204 ₍₂₎	0.568 ₍₁₉₎	0.401 ₍₁₇₎	0.394 ₍₁₆₎	0.100 ₍₉₎	0.491 ₍₁₂₎	0.412 ₍₁₂₎	0.751 ₍₁₂₎	11.3
NBeats	0.832 ₍₂₀₎	0.235 ₍₇₎	0.265 ₍₁₎	0.370 ₍₁₃₎	0.420 ₍₁₇₎	0.081 ₍₂₎	0.521 ₍₁₅₎	0.409 ₍₁₁₎	0.741 _(10.5)	10.7
ScaleFormer	0.769 ₍₁₆₎	0.310 ₍₁₂₎	0.379 ₍₇₎	0.304 ₍₅₎	0.438 ₍₁₈₎	0.138 ₍₁₄₎	0.579 ₍₁₉₎	0.475 ₍₁₇₎	0.864 ₍₂₂₎	14.4
PatchTST	0.710 ₍₁₀₎	0.293 ₍₉₎	0.411 ₍₁₆₎	0.348 ₍₁₁₎	0.555 ₍₂₁₎	0.147 ₍₁₅₎	0.489 ₍₁₁₎	0.392 ₍₉₎	0.720 ₍₆₎	12.0
FedFormer	0.744 ₍₁₁₎	0.317 ₍₁₃₎	0.385 ₍₉₎	0.341 ₍₁₀₎	0.347 ₍₇₎	0.233 ₍₂₁₎	0.484 ₍₉₎	0.413 ₍₁₃₎	0.762 ₍₁₆₎	12.1
Autoformer	0.751 ₍₁₂₎	0.321 ₍₁₄₎	0.392 ₍₁₂₎	0.313 ₍₇₎	0.354 ₍₉₎	0.167 ₍₁₆₎	0.484 ₍₉₎	0.496 ₍₁₉₎	0.756 ₍₁₃₎	12.3
Pyraformer	0.781 ₍₁₈₎	0.371 ₍₂₀₎	0.390 ₍₁₀₎	0.379 ₍₁₅₎	0.385 ₍₁₅₎	0.112 ₍₁₁₎	0.493 ₍₁₃₎	0.435 ₍₁₄₎	0.735 ₍₉₎	13.9
Informer	0.757 ₍₁₄₎	0.336 ₍₁₇₎	0.391 ₍₁₁₎	0.383 ₍₁₆₎	0.364 ₍₁₀₎	0.192 ₍₁₈₎	0.605 ₍₂₁₎	0.542 ₍₂₀₎	0.772 ₍₁₇₎	16.0
Transformer	0.765 ₍₁₅₎	0.321 ₍₁₄₎	0.410 ₍₁₅₎	0.405 ₍₁₈₎	0.370 ₍₁₁₎	0.178 ₍₁₇₎	0.567 ₍₁₈₎	0.592 ₍₂₁₎	0.785 ₍₂₀₎	16.6
SCINet	0.601 ₍₁₎	0.193 ₍₁₎	0.335 ₍₆₎	0.280 ₍₄₎	0.344 ₍₆₎	0.137 ₍₁₃₎	0.463 ₍₇₎	0.389 ₍₇₎	0.732 ₍₇₎	5.8
NLinear	0.636 ₍₅₎	0.223 ₍₅₎	0.293 ₍₄₎	0.239 ₍₁₎	0.328 ₍₃₎	0.091 ₍₇₎	0.418 ₍₁₎	0.375 ₍₄₎	0.706 ₍₄₎	3.8
DLinear	0.640 ₍₆₎	0.497 ₍₂₃₎	0.268 ₍₂₎	0.336 ₍₉₎	0.444 ₍₁₉₎	0.102 ₍₁₀₎	0.442 ₍₆₎	0.378 ₍₅₎	0.686 ₍₂₎	9.1
LSTMa	0.974 ₍₂₃₎	0.305 ₍₁₁₎	0.510 ₍₁₈₎	0.444 ₍₁₉₎	0.501 ₍₂₀₎	0.534 ₍₂₀₎	0.782 ₍₂₃₎	0.699 ₍₂₃₎	0.897 ₍₂₃₎	20.3

prediction horizon of $H = 168$ as in the previous experiment, we also use $H = 96, 192, 336$ and 720 as in (Wu et al., 2021; Zhou et al., 2022b).

Table 3 shows the number of trainable parameters and inference time. As can be seen, `mr-Diff` has a smaller number of trainable parameters. Moreover, even with $S = 4$ or 5 stages, its inference is more efficient than SSSD, CSDI, and TimeGrad. When $S = 3$ (the setting used in the previous experiments), `mr-Diff` is even faster than TimeDiff. The inference efficiency of `mr-Diff` over existing diffusion models is due to: (i) It is non-autoregressive, which avoids autoregressive decoding in each time step as in TimeGrad. (ii) It uses convolution layers. This avoids the scaling problem in CSDI, whose complexity is quadratic with the number of variables and time series length. (iii) Use of the acceleration technique DPM-Solver (Lu et al., 2022), which further reduces the number of denoising steps in each stage. Moreover, it is also faster than TimeDiff when $S = 2$ or 3 , thanks to its use of a linear mapping for future mixup instead of employing additional deep layers for this purpose.

Table 3: Number of trainable parameters and inference time (in ms) of various time series diffusion models with different prediction horizons (H) on the univariate *ETTh1*.

	# of trainable params	inference time (ms)				
		$H=96$	$H=168$	$H=192$	$H=336$	$H=720$
<code>mr-Diff</code> ($S=2$)	0.9M	8.3	9.5	9.8	11.9	21.6
<code>mr-Diff</code> ($S=3$)	1.4M	12.5	14.3	14.9	16.8	27.5
<code>mr-Diff</code> ($S=4$)	1.8M	16.7	19.1	19.7	28.5	36.4
<code>mr-Diff</code> ($S=5$)	2.3M	30.0	30.2	30.2	35.0	43.6
TimeDiff	1.7M	16.2	17.3	17.6	26.5	34.6
TimeGrad	3.1M	870.2	1620.9	1854.5	3119.7	6724.1
CSDI	10M	90.4	128.3	142.8	398.9	513.1
SSSD	32M	418.6	590.2	645.4	1054.2	2516.9

6 CONCLUSION

In this paper, we propose the multi-resolution diffusion model `mr-Diff`. Different from the existing time series diffusion model, `mr-Diff` incorporates seasonal-trend decomposition and uses multiple temporal resolutions in both the diffusion and denoising processes. By progressively denoising the time series in a coarse-to-fine manner, `mr-Diff` is able to produce more reliable predictions. Experiments on a number of real-world univariate and multivariate time series datasets show that `mr-Diff` outperforms the state-of-the-art time series diffusion models. It also achieves very competitive performance even when compared to various families of non-diffusion-based time series prediction models.

ACKNOWLEDGMENTS

This research was supported in part by the Research Grants Council of the Hong Kong Special Administrative Region (Grant 16200021), and the Science and Technology Planning Project of Guangdong Province (Grant 2023A0505050106).

REFERENCES

- Juan Miguel Lopez Alcaraz and Nils Strodthoff. Diffusion-based time series imputation and forecasting with structured state space models. Technical report, arXiv, 2022.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. In *International Conference on Learning Representations*, 2015.
- Yaniv Benny and Lior Wolf. Dynamic dual-output diffusion models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022.
- Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023.
- Cristian Challu, Kin G Olivares, Boris N Oreshkin, Federico Garza, Max Mergenthaler-Canseco, and Artur Dubrawski. N-HITS: Neural hierarchical interpolation for time series forecasting. In *AAAI Conference on Artificial Intelligence*, 2023.
- Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. In *Advances in Neural Information Processing Systems*, 2021.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsbys. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2020.
- Stefan Elfving, Eiji Uchibe, and Kenji Doya. Sigmoid-weighted linear units for neural network function approximation in reinforcement learning. Technical report, arXiv, 2017.
- Wei Fan, Shun Zheng, Xiaohan Yi, Wei Cao, Yanjie Fu, Jiang Bian, and Tie-Yan Liu. DEPTS: Deep expansion learning for periodic time series forecasting. In *International Conference on Learning Representations*, 2022.
- William Harvey, Saeid Naderiparizi, Vaden Masrani, Christian Weilbach, and Frank Wood. Flexible diffusion modeling of long videos. In *Advances in Neural Information Processing Systems*, 2022.
- Hansika Hewamalage, Christoph Bergmeir, and Kasun Bandara. Recurrent neural networks for time series forecasting: Current status and future directions. *International Journal of Forecasting*, 37(1):388–427, 2021.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*, 2020.
- Jonathan Ho, Chitwan Saharia, William Chan, David J. Fleet, Mohammad Norouzi, and Tim Salimans. Cascaded diffusion models for high fidelity image generation. *Journal of Machine Learning Research*, 23(47):1–33, 2022.
- Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997.
- Paul Jeha, Michael Bohlke-Schneider, Pedro Mercado, Shubham Kapoor, Rajbir Singh Nirwan, Valentin Flunkert, Jan Gasthaus, and Tim Januschowski. PSA-GAN: Progressive self attention GANs for synthetic time series. In *International Conference on Learning Representations*, 2022.
- Taesung Kim, Jinhee Kim, Yunwon Tae, Cheonbok Park, Jang-Ho Choi, and Jaegul Choo. Reversible instance normalization for accurate time-series forecasting against distribution shift. In *International Conference on Learning Representations*, 2021.

- 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*, 2020.
- Guokun Lai, Wei-Cheng Chang, Yiming Yang, and Hanxiao Liu. Modeling long-and short-term temporal patterns with deep neural networks. In *International ACM SIGIR Conference on Research & Development in Information Retrieval*, 2018.
- Yan Li, Xinjiang Lu, Yaqing Wang, and Dejing Dou. Generative time series forecasting with diffusion, denoise, and disentanglement. In *Advances in Neural Information Processing Systems*, 2022.
- Minhao Liu, Ailing Zeng, Muxi Chen, Zhijian Xu, Qiuxia Lai, Lingna Ma, and Qiang Xu. Scinet: Time series modeling and forecasting with sample convolution and interaction. In *Advances in Neural Information Processing Systems*, 2022.
- Shizhan Liu, Hang Yu, Cong Liao, Jianguo Li, Weiyao Lin, Alex X. Liu, and Schahram Dustdar. Pyraformer: Low-complexity pyramidal attention for long-range time series modeling and forecasting. In *International Conference on Learning Representations*, 2021.
- Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. DPM-solver: A fast ODE solver for diffusion probabilistic model sampling in around 10 steps. In *Advances in Neural Information Processing Systems*, 2022.
- Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022.
- Kiran Madhusudhanan, Johannes Burchert, Nghia Duong-Trung, Stefan Born, and Lars Schmidt-Thieme. Yformer: U-net inspired transformer architecture for far horizon time series forecasting. Technical report, arXiv, 2021.
- Yuqi Nie, Nam H Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth 64 words: Long-term forecasting with transformers. Technical report, arXiv, 2022.
- Hao Niu, Chuizheng Meng, Defu Cao, Guillaume Habault, Roberto Legaspi, Shinya Wada, Chihiro Ono, and Yan Liu. Mu2ReST: Multi-resolution recursive spatio-temporal transformer for long-term prediction. In *Advances in Knowledge Discovery and Data Mining*, 2022.
- Boris N Oreshkin, Dmitri Carpov, Nicolas Chapados, and Yoshua Bengio. N-BEATS: Neural basis expansion analysis for interpretable time series forecasting. In *International Conference on Learning Representations*, 2019.
- Syama Sundar Rangapuram, Shubham Kapoor, Rajbir Singh Nirwan, Pedro Mercado, Tim Januschowski, Yuyang Wang, and Michael Bohlke-Schneider. Coherent probabilistic forecasting of temporal hierarchies. In *International Conference on Artificial Intelligence and Statistics*, 2023.
- 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*, 2021.
- Cleveland Robert, C William, and Terpenning Irma. STL: A seasonal-trend decomposition procedure based on loess. *Journal of Official Statistics*, 1990.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022.
- Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-Net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention*, 2015.
- Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L. Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, Jonathan Ho, David J. Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding. In *Advances in Neural Information Processing Systems*, 2022.

- Mohammad Amin Shabani, Amir H. Abdi, Lili Meng, and Tristan Sylvain. Scaleformer: Iterative multi-scale refining transformers for time series forecasting. In *International Conference on Learning Representations*, 2023.
- Lifeng Shen and James Kwok. Non-autoregressive conditional diffusion models for time series prediction. In *International Conference on Machine Learning*, 2023.
- 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*, 2021.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, 2017.
- Ronald J. Williams and David Zipser. A learning algorithm for continually running fully recurrent neural networks. *Neural Computation*, 1(2):270–280, 1989.
- Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Autoformer: Decomposition transformers with auto-correlation for long-term series forecasting. In *Advances in Neural Information Processing Systems*, 2021.
- Bing Yu, Haoteng Yin, and Zhanxing Zhu. St-unet: A spatio-temporal U-network for graph-structured time series modeling. Technical report, arXiv, 2021.
- Zhihan Yue, Yujing Wang, Juanyong Duan, Tianmeng Yang, Congrui Huang, Yunhai Tong, and Bixiong Xu. TS2Vec: Towards universal representation of time series. In *AAAI Conference on Artificial Intelligence*, 2022.
- Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are transformers effective for time series forecasting? In *AAAI Conference on Artificial Intelligence*, 2023.
- Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *International Conference on Learning Representations*, 2018.
- Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang. Informer: Beyond efficient transformer for long sequence time-series forecasting. In *AAAI Conference on Artificial Intelligence*, 2021.
- Tian Zhou, Ziqing MA, Xue Wang, Qingsong Wen, Liang Sun, Tao Yao, Wotao Yin, and Rong Jin. FiLM: Frequency improved legendre memory model for long-term time series forecasting. In *Advances in Neural Information Processing Systems*, 2022a.
- Tian Zhou, Ziqing Ma, Qingsong Wen, Xue Wang, Liang Sun, and Rong Jin. FEDformer: Frequency enhanced decomposed transformer for long-term series forecasting. In *International Conference on Machine Learning*, 2022b.

A PSEUDOCODE

The training and inference procedures are shown in Algorithms 1 and 2, respectively.

Algorithm 1 Training procedure of the backward denoising process.

```

1: repeat
2:    $(\mathbf{X}_0, \mathbf{Y}_0) \sim q(\mathbf{X}, \mathbf{Y})$ 
3:    $\{\mathbf{X}_s\}_{s=1, \dots, S-1} = \text{TrendExtraction}(\mathbf{X}_0)$ 
4:    $\{\mathbf{Y}_s\}_{s=1, \dots, S-1} = \text{TrendExtraction}(\mathbf{Y}_0)$ 
5:   for  $s = S - 1, \dots, 0$  do
6:      $k \sim \text{Uniform}(\{1, 2, \dots, K\})$ 
7:      $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
8:     Generate diffused sample  $\mathbf{Y}_s^k$  by  $\mathbf{Y}_s^k = \sqrt{\bar{\alpha}_k} \mathbf{Y}_s^0 + \sqrt{1 - \bar{\alpha}_k} \epsilon$ 
9:     Obtain diffusion step  $k$ 's embedding  $\mathbf{p}^k$  using (7)
10:    Randomly generate a matrix  $\mathbf{m}$  in (9)
11:    Obtain  $\mathbf{z}_{\text{history}}$  using linear mapping on  $\mathbf{X}_s$ 
12:    Obtain  $\mathbf{z}_{\text{mix}}$  using (9)
13:    if  $s < S - 1$  then
14:      Obtain condition  $\mathbf{c}_s$  by concatenating  $\mathbf{z}_{\text{mix}}$  and  $\mathbf{Y}_{s+1}^0$ 
15:    else
16:      Obtain condition  $\mathbf{c}_s = \mathbf{z}_{\text{mix}}$ 
17:    end if
18:    Calculate the loss  $\mathcal{L}_s^k(\theta_s)$  in (12)
19:    Take a gradient descent step based on  $\nabla_{\theta_s} \mathcal{L}_k(\theta_s)$ 
20:  end for
21: until converged

```

Algorithm 2 Inference procedure.

```

1:  $\{\mathbf{X}_s\}_{s=1, \dots, S-1} = \text{TrendExtraction}(\mathbf{X}_0)$ 
2: for  $s = S - 1, \dots, 0$  do
3:   Obtain  $\mathbf{z}_{\text{history}}$  using linear mapping on  $\mathbf{X}_s$ 
4:   if  $s < S - 1$  then
5:     Obtain condition  $\mathbf{c}_s$  by concatenating  $\mathbf{z}_{\text{history}}$  and  $\mathbf{Y}_{s+1}^0$ 
6:   else
7:     Obtain condition  $\mathbf{c}_s = \mathbf{z}_{\text{history}}$ 
8:   end if
9:   Initialization:  $\hat{\mathbf{Y}}_s^K \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
10:  for  $k = K, \dots, 1$  do
11:     $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , if  $k > 1$ , else  $\epsilon = 0$ 
12:    Obtain diffusion step  $k$ 's embedding  $\mathbf{p}^k$  using (7)
13:    Obtain the denoised output  $\hat{\mathbf{Y}}_s^{k-1}$  by (13)
14:  end for
15: end for
16: return  $\hat{\mathbf{Y}}_0^0$ 

```

B DATASETS

Experiments are performed on nine commonly-used real-world time series datasets (Zhou et al., 2021; Wu et al., 2021; Fan et al., 2022): (i) *NorPool*¹, which includes eight years of hourly energy production volume series in multiple European countries; (ii) *Caiso*², which contains eight years of hourly actual electricity load series in different zones of California; (iii) *Traffic*³, which records the hourly road occupancy rates generated by sensors in the San Francisco Bay area freeways; (iv) *Electricity*⁴, which includes the hourly electricity consumption of 321 clients over two years; (v) *Weather*⁵, which records 21 meteorological indicators at 10-minute intervals from 2020 to 2021; (vi)

¹<https://www.nordpoolgroup.com/Market-data/Power-system-data>

²<http://www.energyonline.com/Data>

³<http://pems.dot.ca.gov>

⁴<https://archive.ics.uci.edu/ml/datasets/ElectricityLoadDiagrams20112014>

⁵<https://www.bgc-jena.mpg.de/wetter/>

*Exchange*⁶ (Lai et al., 2018), which describes the daily exchange rates of eight countries (Australia, British, Canada, Switzerland, China, Japan, New Zealand, and Singapore); (vii)-(viii) *ETTh1* and *ETTm1*,⁷, which contain two years of electricity transformer temperature data (Zhou et al., 2021) collected in China, at 1-hour and 15-minute intervals, respectively; (ix) *Wind*⁸, which contains wind power records from 2020-2021 at 15-minute intervals (Li et al., 2022). Table 4 shows the dataset information. As different datasets have different sampling interval lengths, we follow (Shen & Kwok, 2023) and consider prediction tasks with more reasonable prediction lengths.

Table 4: Summary of dataset statistics, including the dimension, total number of observations, sampling frequency, and prediction length H used in the experiments.

	dimension	#observations	frequency	steps (H)
<i>NorPool</i>	18	70,128	1 hour	720 (1 month)
<i>Caiso</i>	10	74,472	1 hour	720 (1 month)
<i>Traffic</i>	862	17,544	1 hour	168 (1 week)
<i>Electricity</i>	321	26,304	1 hour	168 (1 week)
<i>Weather</i>	21	52,696	10 mins	672 (1 week)
<i>Exchange</i>	8	7,588	1 day	14 (2 weeks)
<i>ETTh1</i>	7	17,420	1 hour	168 (1 week)
<i>ETTm1</i>	7	69,680	15 mins	192 (2 days)
<i>Wind</i>	7	48,673	15 mins	192 (2 days)

Besides running directly on these multivariate datasets, we also convert them to univariate time series for performance comparison. For *NorPool* and *Caiso*, the univariate time series are extracted from all variables as in (Fan et al., 2022). For the other datasets, we follow (Wu et al., 2021; Zhou et al., 2021) and extract the univariate time series by using the last variable only.

Figure 5 shows examples of the time series data used in the experiments. Since all of them are multivariate, we only show the last variate. As can be seen, these datasets are representative in exhibiting various temporal characteristics (different stationary states, periodicities, and sampling rates). For example, *Electricity* contains obvious periodic patterns, while *ETTh1* and *ETTm1* exhibit non-stationary trends.

Real-world time series often have irregular dynamics over time or across variables. As in many deep time series models (Kim et al., 2021; Zhou et al., 2022a; Liu et al., 2022; Zeng et al., 2023), it is often useful to normalize the scales of the time series in each window. In the proposed method, we use instance normalization, which has also been used in RevIN (Kim et al., 2021) (which is referred to as reversible instance normalization) and FiLM (Zhou et al., 2022a). Specifically, we subtract the time series value in each window by its lookback window mean, and then divide by the lookback window’s standard deviation. During inference, the mean and standard deviation are added back to the prediction.

C BASELINES IN MAIN EXPERIMENTS

We compare with several types of baselines, including: (i) recent time series diffusion models: non-autoregressive diffusion model TimeDiff (Shen & Kwok, 2023), TimeGrad (Rasul et al., 2021), conditional score-based diffusion model for imputation (CSDI) (Tashiro et al., 2021), structured state space model-based diffusion (SSSD) (Alcaraz & Strodthoff, 2022); (ii) recent generative models for time series prediction: variational autoencoder with diffusion, denoise and disentanglement (D³VAE) (Li et al., 2022), coherent probabilistic forecasting (CPF) (Rangapuram et al., 2023), and PSA-GAN (Jeha et al., 2022); (iii) recent prediction models based on basis expansion: N-Hits (Challu et al., 2023), frequency improved Legendre memory model (FiLM) (Zhou et al., 2022a), Depts (Fan et al., 2022) and NBeats (Oreshkin et al., 2019); (iv) time series transformers: Scaleformer (Shabani et al., 2023), PatchTST (Nie et al., 2022), Fedformer (Zhou et al., 2022b), Autoformer (Wu et al., 2021), Pyraformer (Liu et al., 2021), Informer (Zhou et al., 2021) and the

⁶<https://github.com/laiguokun/multivariate-time-series-data>

⁷<https://github.com/zhouhaoyi/ETDataset>

⁸<https://github.com/PaddlePaddle/PaddleSpatial/tree/main/paddlespatial/datasets/WindPower>

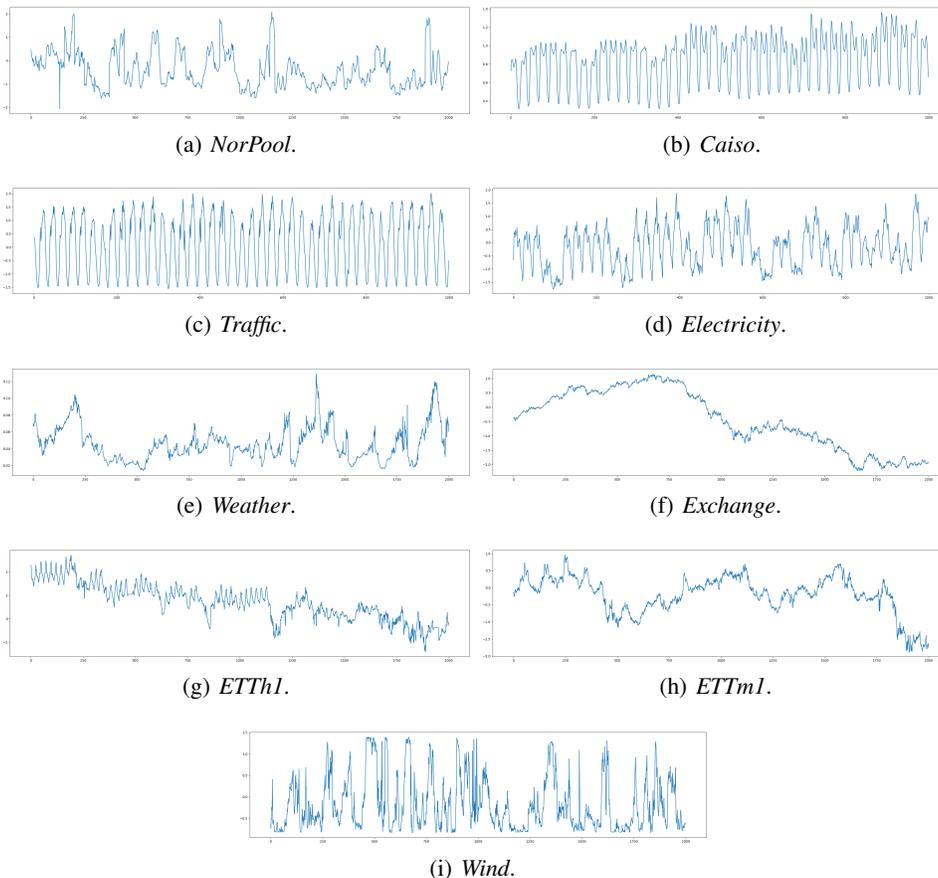


Figure 5: Visualization of the time series datasets.

standard Transformer (Vaswani et al., 2017); and (v) other competitive baselines: SCINet (Liu et al., 2022) that introduces sample convolution and interaction for time series prediction, NLinear (Zeng et al., 2023), DLinear (Zeng et al., 2023) LSTMa (Bahdanau et al., 2015), and an attention-based LSTM (Hochreiter & Schmidhuber, 1997). We do not compare with (Yu et al., 2021; Niu et al., 2022; Madhusudhanan et al., 2021) because their codes are not publicly available.

D RESULTS ON MSE

Tables 5 and 6 show the MSE results on the univariate and multivariate time series, respectively, for the experiment in Section 5.1. As can be seen, the proposed `mr-Diff` still achieves the best overall MSE performance in terms of average rank. Moreover, since deep models for time series forecasting may be influenced by different random initializations, Table 7 shows the prediction results on the univariate time series over 5 random runs. As can be seen, the largest standard deviation is 0.0042, indicating that the proposed model is robust towards different initializations.

E ABLATION STUDIES

In this section, we perform a number of ablation experiments on `mr-Diff` using the *Electricity*, *ETTh1* and *ETTM1* datasets. These datasets are representative in exhibiting various temporal characteristics (different stationary states, periodicities, and sampling rates). For example, *Electricity* contains obvious periodic patterns, while *ETTh1* and *ETTM1* exhibit non-stationary trends.

Table 5: Univariate prediction MSEs on the real-world time series datasets (subscript is the rank). Results of all baselines (except NLinear, SCINet, and N-Hits) are from (Shen & Kwok, 2023).

Method	<i>NorPool</i>	<i>Caiso</i>	<i>Traffic</i>	<i>Electricity</i>	<i>Weather</i>	<i>Exchange</i>	<i>ETTh1</i>	<i>ETTm1</i>	<i>Wind</i>	Rank
mr-Diff	0.667 ₍₄₎	0.122 ₍₃₎	0.119 ₍₁₎	0.234 ₍₃₎	0.002 ₍₁₎	0.016 ₍₁₎	0.066 ₍₁₎	<u>0.039</u> ₍₂₎	2.182 ₍₄₎	2.2
TimeDiff	<u>0.636</u> ₍₂₎	0.122 ₍₃₎	<u>0.121</u> ₍₂₎	<u>0.232</u> ₍₂₎	0.002 ₍₁₎	<u>0.017</u> ₍₄₎	0.066 ₍₁₎	0.040 ₍₅₎	2.407 ₍₁₃₎	3.7
TimeGrad	1.129 ₍₂₂₎	0.325 ₍₂₂₎	1.223 ₍₂₃₎	0.920 ₍₂₃₎	0.002 ₍₁₎	0.041 ₍₂₀₎	0.078 ₍₁₀₎	0.048 ₍₁₁₎	2.530 ₍₁₈₎	16.7
CSDI	0.967 ₍₂₁₎	0.192 ₍₁₄₎	0.393 ₍₂₀₎	0.520 ₍₁₈₎	0.002 ₍₁₎	0.071 ₍₂₃₎	0.083 ₍₁₅₎	0.050 ₍₁₅₎	2.434 ₍₁₅₎	15.8
SSSD	1.145 ₍₂₃₎	0.176 ₍₁₂₎	0.151 ₍₈₎	0.370 ₍₁₁₎	0.004 ₍₁₁₎	0.023 ₍₁₆₎	0.097 ₍₂₀₎	0.049 ₍₁₃₎	3.149 ₍₂₂₎	15.1
D ³ VAE	0.964 ₍₂₀₎	0.521 ₍₂₃₎	0.151 ₍₈₎	0.535 ₍₁₉₎	<u>0.003</u> ₍₉₎	0.019 ₍₁₂₎	0.078 ₍₁₀₎	0.044 ₍₉₎	2.679 ₍₂₀₎	14.4
CPF	0.855 ₍₁₅₎	0.260 ₍₂₁₎	0.279 ₍₁₈₎	0.609 ₍₂₁₎	0.002 ₍₁₎	0.016 ₍₁₎	0.080 ₍₁₃₎	0.041 ₍₆₎	2.430 ₍₁₄₎	12.2
PSA-GAN	0.658 ₍₃₎	0.150 ₍₇₎	0.250 ₍₁₇₎	0.273 ₍₇₎	0.035 ₍₂₁₎	0.020 ₍₁₃₎	0.084 ₍₁₆₎	0.051 ₍₁₆₎	2.510 ₍₁₇₎	13.0
N-Hits	0.739 ₍₉₎	0.170 ₍₁₁₎	0.147 ₍₇₎	0.346 ₍₈₎	0.002 ₍₁₎	<u>0.017</u> ₍₄₎	0.089 ₍₁₇₎	0.043 ₍₈₎	2.406 ₍₁₂₎	8.6
FiLM	0.707 ₍₇₎	0.185 ₍₁₃₎	0.198 ₍₁₃₎	0.260 ₍₅₎	0.007 ₍₁₄₎	0.018 ₍₉₎	<u>0.070</u> ₍₃₎	0.038 ₍₁₎	<u>2.143</u> ₍₂₎	7.4
Depts	0.668 ₍₅₎	0.107 ₍₁₎	0.151 ₍₈₎	0.380 ₍₁₄₎	0.024 ₍₁₉₎	0.020 ₍₁₃₎	<u>0.070</u> ₍₃₎	0.046 ₍₁₀₎	3.457 ₍₂₃₎	10.7
NBeats	0.768 ₍₁₁₎	0.125 ₍₅₎	0.142 ₍₆₎	0.378 ₍₁₃₎	0.137 ₍₂₂₎	0.016 ₍₁₎	0.095 ₍₁₉₎	0.048 ₍₁₁₎	2.434 ₍₁₅₎	11.4
Scaleformer	0.778 ₍₁₂₎	0.232 ₍₁₆₎	0.286 ₍₁₉₎	0.361 ₍₉₎	0.009 ₍₁₇₎	0.035 ₍₁₉₎	0.150 ₍₂₂₎	0.078 ₍₂₂₎	2.646 ₍₁₉₎	17.2
PatchTST	0.595 ₍₁₎	0.193 ₍₁₅₎	0.177 ₍₁₂₎	0.450 ₍₁₇₎	0.026 ₍₂₀₎	0.020 ₍₁₃₎	0.106 ₍₂₁₎	0.052 ₍₁₈₎	2.698 ₍₂₁₎	21
FedFormer	0.891 ₍₁₆₎	0.164 ₍₉₎	0.173 ₍₁₁₎	0.376 ₍₁₂₎	0.005 ₍₁₂₎	0.050 ₍₂₂₎	0.076 ₍₇₎	0.065 ₍₂₁₎	2.351 ₍₁₁₎	13.4
Autoformer	0.946 ₍₁₉₎	0.248 ₍₁₇₎	0.473 ₍₂₁₎	0.659 ₍₂₂₎	<u>0.003</u> ₍₉₎	0.041 ₍₂₀₎	0.081 ₍₁₄₎	0.051 ₍₁₆₎	2.349 ₍₁₀₎	16.4
Pyraformer	0.933 ₍₁₈₎	0.165 ₍₁₀₎	0.136 ₍₄₎	0.389 ₍₁₅₎	0.020 ₍₁₈₎	<u>0.017</u> ₍₇₎	0.076 ₍₇₎	0.054 ₍₁₉₎	2.279 ₍₆₎	11.2
Informer	0.804 ₍₁₃₎	0.250 ₍₁₈₎	0.213 ₍₁₅₎	0.363 ₍₁₀₎	0.007 ₍₁₄₎	0.023 ₍₁₆₎	0.076 ₍₇₎	0.049 ₍₁₃₎	2.297 ₍₇₎	12.6
Transformer	0.928 ₍₁₇₎	0.250 ₍₁₈₎	0.238 ₍₁₆₎	0.430 ₍₁₄₎	0.007 ₍₁₄₎	0.018 ₍₉₎	0.092 ₍₁₈₎	0.058 ₍₂₀₎	2.306 ₍₉₎	15.2
SCINet	0.746 ₍₁₀₎	0.154 ₍₈₎	0.212 ₍₁₄₎	0.272 ₍₆₎	0.002 ₍₁₎	0.018 ₍₉₎	0.071 ₍₆₎	<u>0.039</u> ₍₂₎	2.063 ₍₁₎	6.3
NLinear	0.708 ₍₈₎	0.147 ₍₆₎	0.124 ₍₃₎	0.231 ₍₁₎	0.002 ₍₁₎	0.017 ₍₄₎	<u>0.070</u> ₍₃₎	<u>0.039</u> ₍₂₎	2.193 ₍₅₎	3.7
DLinear	0.671 ₍₆₎	<u>0.118</u> ₍₂₎	0.139 ₍₅₎	0.244 ₍₄₎	0.168 ₍₂₃₎	<u>0.017</u> ₍₄₎	0.078 ₍₁₀₎	0.041 ₍₆₎	2.171 ₍₃₎	7.0
LSTMa	0.836 ₍₁₄₎	0.253 ₍₂₀₎	1.032 ₍₂₂₎	0.596 ₍₂₀₎	0.005 ₍₁₂₎	0.031 ₍₁₈₎	0.167 ₍₂₃₎	0.091 ₍₂₃₎	2.299 ₍₈₎	17.8

Table 6: Multivariate prediction MSEs on the real-world time series datasets (subscript is the rank). CSDI runs out of memory on *Traffic* and *Electricity*. Results of all baselines (except for NLinear, SCINet, and N-Hits) are from (Shen & Kwok, 2023).

Method	<i>NorPool</i>	<i>Caiso</i>	<i>Traffic</i>	<i>Electricity</i>	<i>Weather</i>	<i>Exchange</i>	<i>ETTh1</i>	<i>ETTm1</i>	<i>Wind</i>	Rank
mr-Diff	<u>0.645</u> ₍₂₎	0.127 ₍₃₎	0.474 ₍₆₎	0.155 ₍₃₎	0.296 ₍₁₎	0.016 ₍₁₎	0.411 ₍₃₎	<u>0.340</u> ₍₂₎	0.881 ₍₁₎	2.4
TimeDiff	0.665 ₍₄₎	0.136 ₍₆₎	0.564 ₍₇₎	0.193 ₍₅₎	<u>0.311</u> ₍₂₎	0.018 ₍₆₎	0.407 ₍₁₎	0.336 ₍₁₎	<u>0.896</u> ₍₂₎	3.8
TimeGrad	1.152 ₍₂₀₎	0.258 ₍₁₉₎	1.745 ₍₂₂₎	0.736 ₍₂₁₎	0.392 ₍₁₄₎	0.079 ₍₂₀₎	0.993 ₍₂₂₎	0.874 ₍₂₁₎	1.209 ₍₂₁₎	20.0
CSDI	1.011 ₍₁₉₎	0.253 ₍₁₈₎	-	-	0.356 ₍₉₎	0.077 ₍₁₉₎	0.497 ₍₇₎	0.529 ₍₁₇₎	1.066 ₍₉₎	14.0
SSSD	0.872 ₍₁₂₎	0.195 ₍₁₀₎	0.642 ₍₁₁₎	0.255 ₍₁₂₎	0.349 ₍₈₎	0.061 ₍₁₆₎	0.726 ₍₁₈₎	0.464 ₍₁₃₎	1.188 ₍₁₉₎	13.2
D ³ VAE	0.745 ₍₉₎	0.241 ₍₁₇₎	0.928 ₍₁₇₎	0.286 ₍₁₅₎	0.375 ₍₁₁₎	0.200 ₍₂₂₎	0.504 ₍₉₎	0.362 ₍₈₎	1.118 ₍₁₅₎	13.7
CPF	1.613 ₍₂₃₎	0.383 ₍₂₁₎	1.625 ₍₂₁₎	0.793 ₍₂₂₎	1.390 ₍₂₃₎	0.016 ₍₁₎	0.730 ₍₁₉₎	0.482 ₍₁₅₎	1.140 ₍₁₇₎	18.0
PSA-GAN	1.501 ₍₂₂₎	0.510 ₍₂₃₎	1.614 ₍₂₀₎	0.535 ₍₂₀₎	1.220 ₍₂₁₎	0.018 ₍₆₎	0.623 ₍₁₇₎	0.537 ₍₁₈₎	1.127 ₍₁₆₎	18.1
N-Hits	0.716 ₍₇₎	0.131 ₍₄₎	<u>0.386</u> ₍₂₎	<u>0.152</u> ₍₂₎	0.323 ₍₄₎	<u>0.017</u> ₍₅₎	0.498 ₍₈₎	0.353 ₍₆₎	1.033 ₍₆₎	4.9
FiLM	0.723 ₍₈₎	0.179 ₍₈₎	0.628 ₍₁₀₎	0.210 ₍₈₎	0.327 ₍₅₎	0.016 ₍₁₎	0.426 ₍₅₎	0.347 ₍₄₎	0.984 ₍₄₎	5.9
Depts	0.662 ₍₃₎	<u>0.106</u> ₍₂₎	1.019 ₍₁₉₎	0.319 ₍₁₇₎	0.761 ₍₁₉₎	0.020 ₍₉₎	0.579 ₍₁₃₎	0.380 ₍₁₀₎	1.082 ₍₁₂₎	11.6
NBeats	0.832 ₍₁₀₎	0.141 ₍₇₎	0.373 ₍₁₎	0.269 ₍₁₃₎	1.344 ₍₂₂₎	0.016 ₍₁₎	0.586 ₍₁₅₎	0.391 ₍₁₁₎	1.069 ₍₁₀₎	10.0
Scaleformer	0.983 ₍₁₅₎	0.207 ₍₁₃₎	0.618 ₍₉₎	0.195 ₍₆₎	0.462 ₍₁₆₎	0.036 ₍₁₂₎	0.613 ₍₁₆₎	0.481 ₍₁₄₎	1.359 ₍₂₂₎	13.7
PatchTST	0.851 ₍₁₁₎	0.193 ₍₉₎	0.831 ₍₁₆₎	0.225 ₍₁₀₎	0.782 ₍₂₀₎	0.047 ₍₁₄₎	0.526 ₍₁₁₎	0.372 ₍₉₎	1.070 ₍₁₁₎	12.3
FedFormer	0.873 ₍₁₃₎	0.205 ₍₁₁₎	0.591 ₍₈₎	0.238 ₍₁₁₎	0.342 ₍₇₎	0.133 ₍₂₁₎	0.541 ₍₁₂₎	0.426 ₍₁₂₎	1.113 ₍₁₄₎	12.1
Autoformer	0.940 ₍₁₄₎	0.226 ₍₁₅₎	0.688 ₍₁₅₎	0.201 ₍₇₎	0.360 ₍₁₀₎	0.056 ₍₁₅₎	0.516 ₍₁₀₎	0.565 ₍₁₉₎	1.083 ₍₁₃₎	13.1
Pyraformer	1.008 ₍₁₈₎	0.273 ₍₂₀₎	0.659 ₍₁₂₎	0.273 ₍₁₄₎	0.394 ₍₁₅₎	0.032 ₍₁₁₎	0.579 ₍₁₃₎	0.493 ₍₁₆₎	1.061 ₍₈₎	14.1
Informer	0.985 ₍₁₆₎	0.231 ₍₁₆₎	0.664 ₍₁₃₎	0.298 ₍₁₆₎	0.385 ₍₁₂₎	0.073 ₍₁₈₎	0.775 ₍₂₁₎	0.673 ₍₂₀₎	1.168 ₍₁₈₎	16.7
Transformer	1.005 ₍₁₇₎	0.206 ₍₁₂₎	0.671 ₍₁₄₎	0.328 ₍₁₈₎	0.388 ₍₁₃₎	0.062 ₍₁₇₎	0.759 ₍₂₀₎	0.992 ₍₂₂₎	1.201 ₍₂₀₎	17.0
SCINet	0.613 ₍₁₎	0.095 ₍₁₎	0.434 ₍₅₎	0.171 ₍₄₎	0.329 ₍₆₎	0.036 ₍₁₂₎	0.465 ₍₆₎	0.359 ₍₇₎	1.055 ₍₇₎	5.4
NLinear	0.707 ₍₆₎	0.135 ₍₅₎	0.430 ₍₄₎	0.147 ₍₁₎	0.313 ₍₃₎	0.019 ₍₈₎	<u>0.410</u> ₍₂₎	0.349 ₍₅₎	0.989 ₍₅₎	4.3
DLinear	0.670 ₍₅₎	0.461 ₍₂₂₎	0.389 ₍₃₎	0.215 ₍₉₎	0.488 ₍₁₇₎	0.022 ₍₁₀₎	0.415 ₍₄₎	<u>0.345</u> ₍₃₎	0.899 ₍₃₎	8.4
LSTMa	1.481 ₍₂₁₎	0.217 ₍₁₄₎	0.966 ₍₁₈₎	0.414 ₍₁₉₎	0.662 ₍₁₈₎	0.403 ₍₂₃₎	1.149 ₍₂₃₎	1.030 ₍₂₃₎	1.464 ₍₂₃₎	20.2

Length of lookback window L . Table 8 shows the prediction MAE of `mr-Diff` with different lengths (L) of the lookback window. We consider $L = \{96, 192, 336, 720, 1440\}$ as used in (Shen & Kwok, 2023; Zeng et al., 2023; Wu et al., 2021). As can be seen, on *Electricity*, *ETTh1* and

Table 7: Univariate prediction errors of $m_r\text{-Diff}$ obtained on five runs.

	<i>NorPool</i>		<i>Caiso</i>		<i>Traffic</i>		<i>Electricity</i>			
	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE
0	0.6096	0.6668	0.2129	0.1225	0.1973	0.1192	0.3322	0.2331		
1	0.6092	0.6664	0.2131	0.1226	0.1974	0.1195	0.3267	0.2293		
2	0.6087	0.6659	0.2091	0.1220	0.1970	0.1193	0.3327	0.2409		
3	0.6099	0.6672	0.2129	0.1225	0.1976	0.1192	0.3382	0.2347		
4	0.6091	0.6665	0.2120	0.1221	0.1973	0.1195	0.3321	0.2339		
mean	0.6093	0.6666	0.2120	0.1223	0.1973	0.1193	0.3324	0.2343		
std deviation	0.0005	0.0005	0.0017	0.0003	0.0002	0.0001	0.0041	0.0042		

	<i>Weather</i>		<i>Exchange</i>		<i>ETTh1</i>		<i>ETTm1</i>		<i>Wind</i>	
	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE
0	0.0322	0.0019	0.0936	0.0155	0.1964	0.0663	0.1486	0.0389	1.1704	2.1862
1	0.0321	0.0018	0.0935	0.0155	0.1968	0.0667	0.1489	0.0387	1.1682	2.1820
2	0.0323	0.0018	0.0947	0.0158	0.1956	0.0660	0.1483	0.0384	1.1690	2.1822
3	0.0322	0.0019	0.0954	0.0160	0.1959	0.0661	0.1486	0.0387	1.1676	2.1807
4	0.0322	0.0018	0.0946	0.0159	0.1965	0.0671	0.1489	0.0389	1.1670	2.1804
mean	0.0322	0.0018	0.0943	0.0157	0.1962	0.0664	0.1486	0.0387	1.1684	2.1823
std deviation	0.0001	0.0001	0.0008	0.0002	0.0005	0.0004	0.0003	0.0002	0.0013	0.0023

ETTm1 (corresponding to 168/168/192-step-ahead prediction, respectively), good performance can be obtained when L is 336 or above.

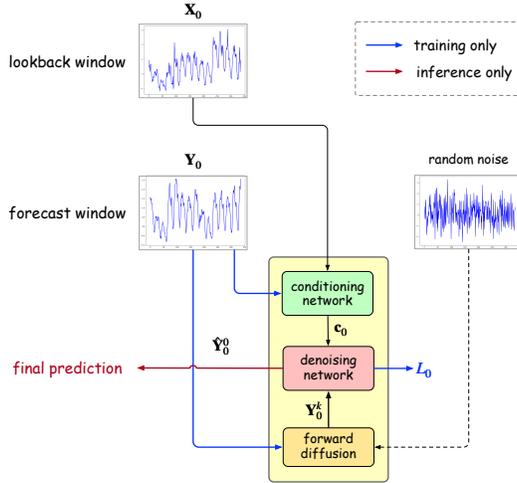


Figure 6: The baseline without seasonal-trend decomposition blocks.

Table 8: Prediction MAE versus lookback window length L .

L	<i>Electricity</i>	<i>ETTh1</i>	<i>ETTm1</i>
96	0.449	0.210	0.165
192	0.376	0.202	0.163
336	0.362	0.196	0.157
720	0.332	0.196	0.152
1,440	0.346	0.199	0.149

Number of stages S . In this experiment, we vary the number of stages S . The kernel size τ_s is set to increase with the stage number s , so as to generate fine-to-coarse trends. We also compare with a variant that does not perform seasonal-trend decomposition (illustrated in Figure 6). This variant directly generates the forecast window from a random noise vector without using the coarser

trend \mathbf{Y}_{s+1}^0 . Table 9 shows the prediction errors. As can be seen, using multiple stages improves performance, and not using the seasonal-trend decomposition leads to the worst performance.

Table 9: Prediction MAE versus number of stages S . L is set to the optimal setting in Table 8 (i.e., *Electricity*: 720, *ETTh1*: 336, *ETTm1*: 1440).

S	<i>Electricity</i>	<i>ETTh1</i>	<i>ETTm1</i>
1	0.403	0.208	0.1573
2	0.389	0.200	0.1529
3	0.363	0.196	0.1525
4	0.346	0.197	0.1509
5	0.332	0.197	0.1496

Number of diffusion steps K . Table 10 shows the prediction error of `mr-Diff` with K , the number of diffusion steps. As can be seen, setting $K = 100$ leads to stable performance across all four datasets. This also agrees with (Li et al., 2022) in that a small K may lead to incomplete diffusion, while a K too large may involve unnecessary computations.

Table 10: Prediction MAE versus number of diffusion steps K .

K	<i>Electricity</i>	<i>ETTh1</i>	<i>ETTm1</i>
50	0.348	0.199	0.151
100	0.332	0.196	0.149
500	0.335	0.198	0.151
1,000	0.337	0.199	0.151

Gaussian noise variance β_k . Recall from Section 5 that β_k is generated by a linear variance schedule (Rasul et al., 2021). In this experiment, we vary β_K in $\{0.001, 0.01, 0.1, 0.9\}$, with β_1 fixed to 10^{-4} . As can be seen from Table 11, setting $\beta_K = 0.1$ always has the best performance. Li et al. (2022) has a similar conclusion that a β_K too small can lead to a unsatisfactory diffusion, while a β_K too large can make the diffusion out of control.

Table 11: Prediction MAE versus Gaussian noise variance upper bound β_k .

β_K	<i>Electricity</i>	<i>ETTh1</i>	<i>ETTm1</i>
0.001	0.403	0.198	0.151
0.01	0.365	0.198	0.152
0.1	0.332	0.196	0.149
0.9	0.372	0.201	0.154

Future mixup. Recall that each element of the future mixup matrix \mathbf{m} in (9) is randomly sampled from the uniform distribution on $[0; 1)$. In this experiment, we study the effect of \mathbf{m} by sampling each of its elements from the Beta distribution $Beta(\gamma, \gamma)$ where $\gamma \in \{0.1, 1, 2\}$. Note that $Beta(1, 1)$ reduces to the uniform distribution. Also, we include a `mr-Diff` variant that does not use future mixup (by replacing (9) with $\mathbf{z}_{\text{mix}} = \mathbf{z}_{\text{history}}$). Table 12 shows the prediction results. As can be seen, the uniform distribution (with $\gamma = 1$) as used in `mr-Diff` has the best performance.

F MORE IMPLEMENTATION DETAILS

In the proposed `mr-Diff`, the conditioning network and the denoising network’s encoder/decoder are built by stacking a number of convolutional blocks. The default configuration of each convolutional block is shown in Table 13.

The number S is selected from $\{2, 3, 4, 5\}$ (the number of stages S is greater than the number of decompositions by 1). When $S = 2$, the kernel size τ_1 is selected from $\{5, 25\}$; When $S = 3$,

Table 12: Prediction MAE with future mixup matrix \mathbf{m} from different Beta distributions $Beta(\gamma, \gamma)$.

γ	<i>Electricity</i>	<i>ETTh1</i>	<i>ETTm1</i>
0.1	0.348	0.202	0.154
1	0.332	0.196	0.149
2	0.776	0.222	0.200
variant without future mixup	0.349	0.198	0.151

Table 13: Configuration of the convolutional block.

layer	operator	default parameters
1	Conv1d	in channel=256, out channel=256, kernel size=3, stride=1, padding=1
2	BatchNorm1d	number of features=256
3	LeakyReLU	negative slope=0.1
4	Dropout	dropout rate=0.1

the kernel size (τ_1, τ_2) is selected from $\{(5, 25), (25, 51), (51, 201)\}$; When $S = 4$, the kernel size (τ_1, τ_2, τ_3) is selected from $\{(5, 25, 51), (25, 51, 201)\}$. When $S = 5$, the kernel size $(\tau_1, \tau_2, \tau_3, \tau_4)$ is selected from $\{(5, 25, 51, 201)\}$. In practice, we run a grid search over S for each dataset. This is not expensive as the number of choices is small (as can be seen from above). Usually, $S = 2$ and the combinations (5,25) and (25,51) achieve promising performance.

G TRAINING EFFICIENCY

In Section 5.2, we discussed inference efficiency. In this section, we further compare training efficiency of `mr-Diff` with the other time series diffusion model baselines (TimeDiff, TimeGrad, CSDI, SSSD) on the univariate *ETTh1* dataset. As can be seen from Table 14, the proposed model requires less training time compared to TimeGrad, CSDI, and SSSD. This suggests that the proposed model is more efficient in terms of training. This is because, during training, `mr-Diff` uses convolution layers in its denoising networks, which eliminates the need for large modules like residual blocks (in TimeGrad and SSSD) and self-attention layers (in CSDI). Moreover, it is also faster than TimeDiff, thanks to its use of a linear mapping for future mixup instead of employing additional deep layers for this purpose.

Table 14: Training time (ms) of various time series diffusion models with different prediction horizons (H) on the univariate *ETTh1*.

	training time (ms)				
	$H = 96$	$H = 168$	$H = 192$	$H = 336$	$H = 720$
<code>mr-Diff</code> ($S=2$)	0.54	0.57	0.58	0.62	0.66
<code>mr-Diff</code> ($S=3$)	0.59	0.69	0.71	0.74	0.82
<code>mr-Diff</code> ($S=4$)	0.78	0.96	0.99	1.07	1.45
<code>mr-Diff</code> ($S=5$)	1.19	1.26	1.27	1.31	1.72
TimeDiff	0.71	0.75	0.77	0.82	0.85
TimeGrad	2.11	2.42	3.21	4.22	5.93
CSDI	5.72	7.09	7.59	10.59	17.21
SSSD	16.98	19.34	22.64	32.12	52.93