



**ICLR | 2024**

Twelfth International Conference  
on Learning Representations



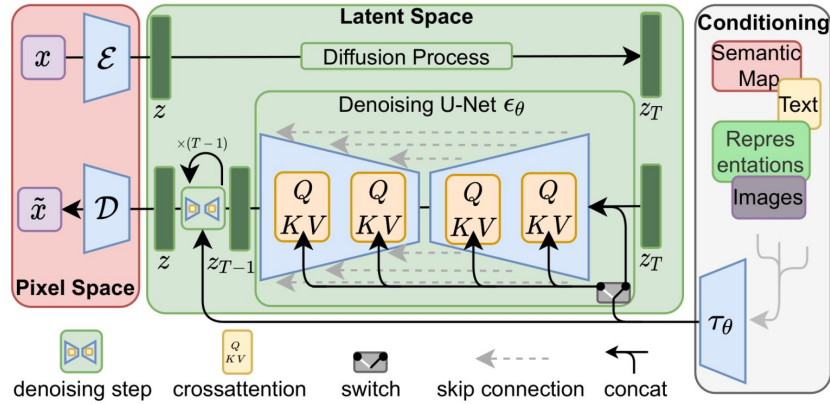
香港科技大學  
THE HONG KONG UNIVERSITY OF  
SCIENCE AND TECHNOLOGY

# Multi-Resolution Diffusion Models for Time Series Forecasting

Lifeng Shen, Weiyu Chen, James T. Kwok

Hong Kong University of Science and Technology

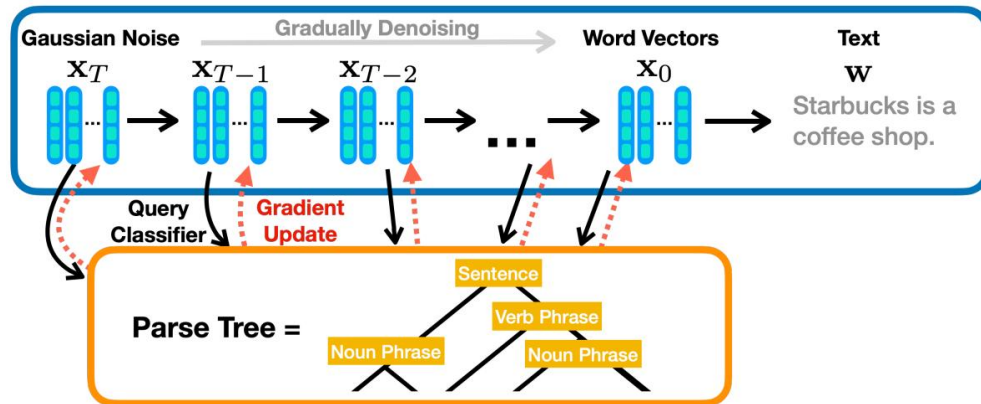
# Diffusion models



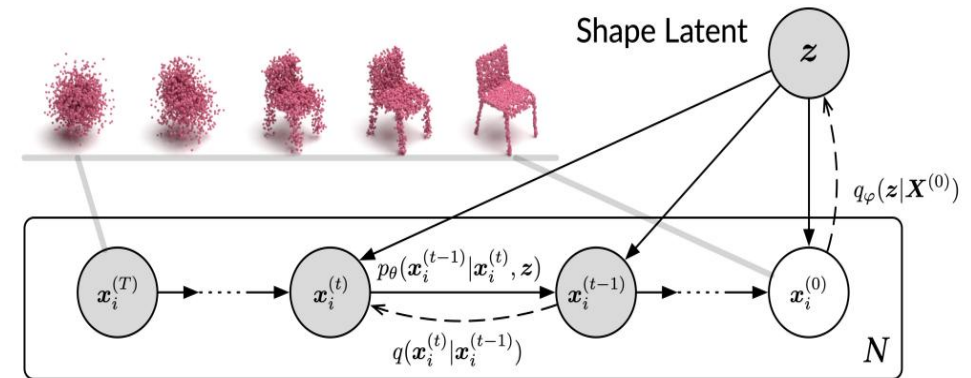
images: Stable Diffusion, Midjourney



video: Sora



text

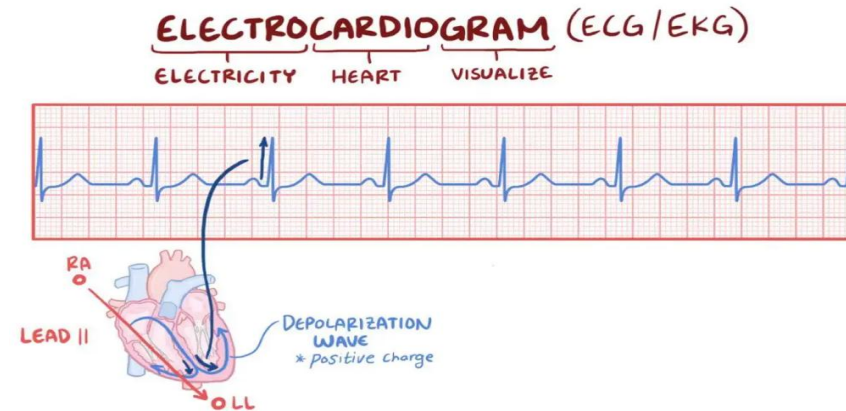


3D data

# Time series data are prevalent



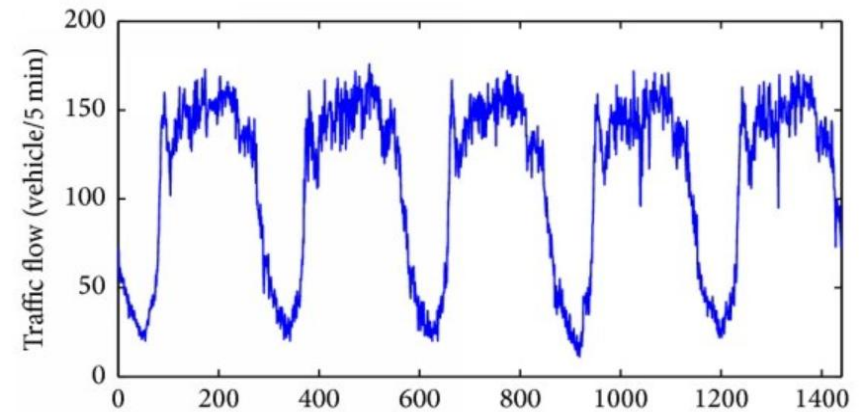
machine monitoring



patient health monitoring

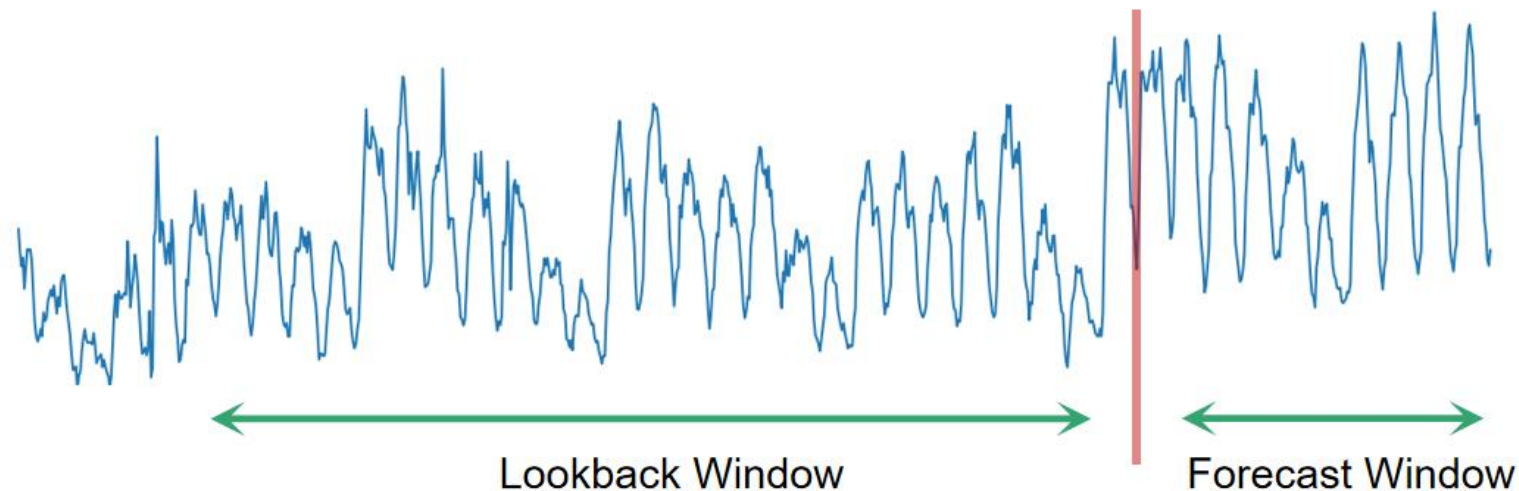


stock price prediction



traffic flow optimization

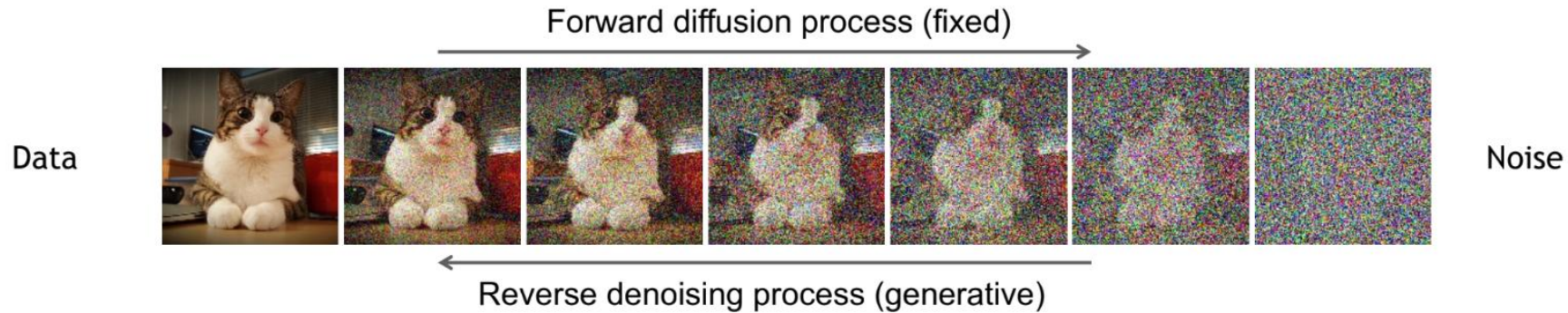
# Time series forecasting



predict forecast window  $\mathbf{x}_{1:H}^0 \in \mathbb{R}^{d \times H}$  given lookback window  $\mathbf{x}_{-L+1:0}^0 \in \mathbb{R}^{d \times H}$

- $d$ : number of variables
- $H$ : length of the forecast window
- $L$ : length of the lookback window

# Denoising diffusion probabilistic models (DPPM)



forward diffusion

- input  $\mathbf{x}^0$  is gradually corrupted to a Gaussian noise vector

$$q(\mathbf{x}^k | \mathbf{x}^{k-1}) = \mathbb{N}(\mathbf{x}^k; \sqrt{1 - \beta_k} \mathbf{x}^{k-1}, \beta_k \mathbf{I})$$

backward denoising

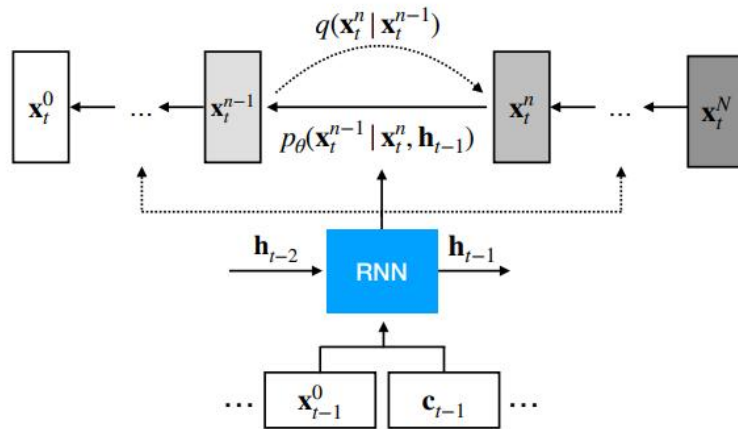
- learns to generate data by denoising

$$p_{\theta}(\mathbf{x}^{k-1} | \mathbf{x}^k) = \mathbb{N}(\mathbf{x}^{k-1}; \mu_{\theta}(\mathbf{x}^k, k), \sigma_k^2 \mathbf{I})$$

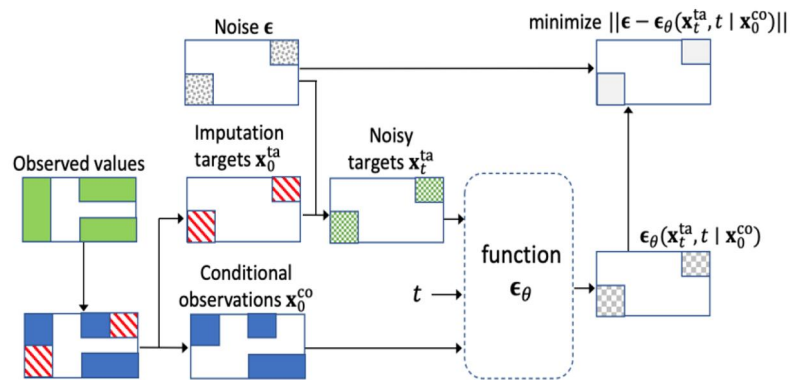
- $\mu_{\theta}(\mathbf{x}^k, k)$ : defined by a neural network (with parameter  $\theta$ ).

# Example diffusion models for time series

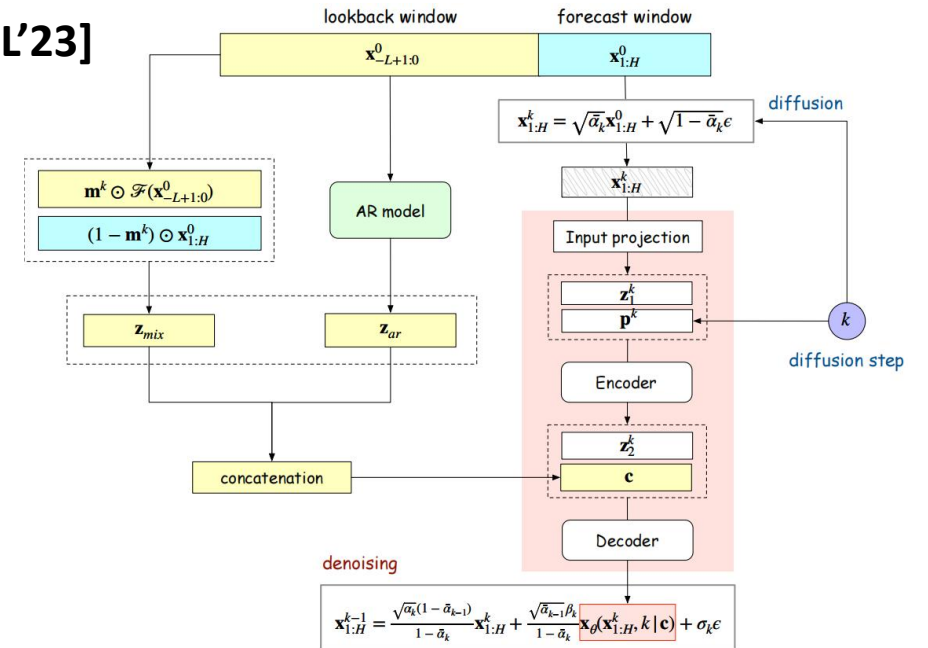
TimeGrad [ICML'21]



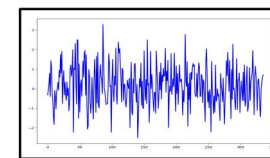
CSDI [NeurIPS'21]



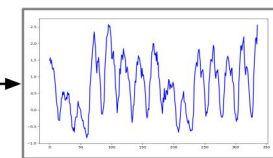
TimeDiff [ICML'23]



generate time series directly from random vectors



random noises



time series

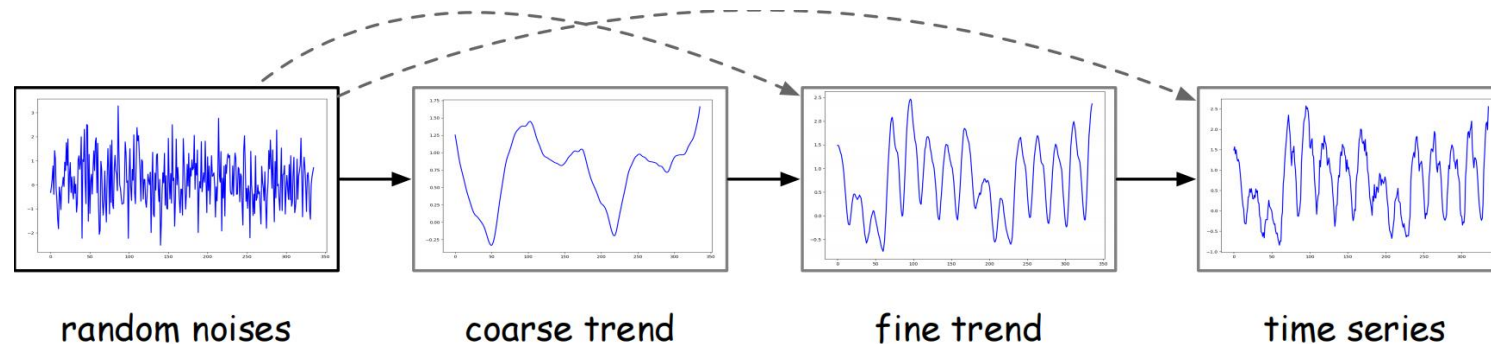
# Motivation

how to better utilize structural properties in time series?

multi-resolution temporal structure in time series

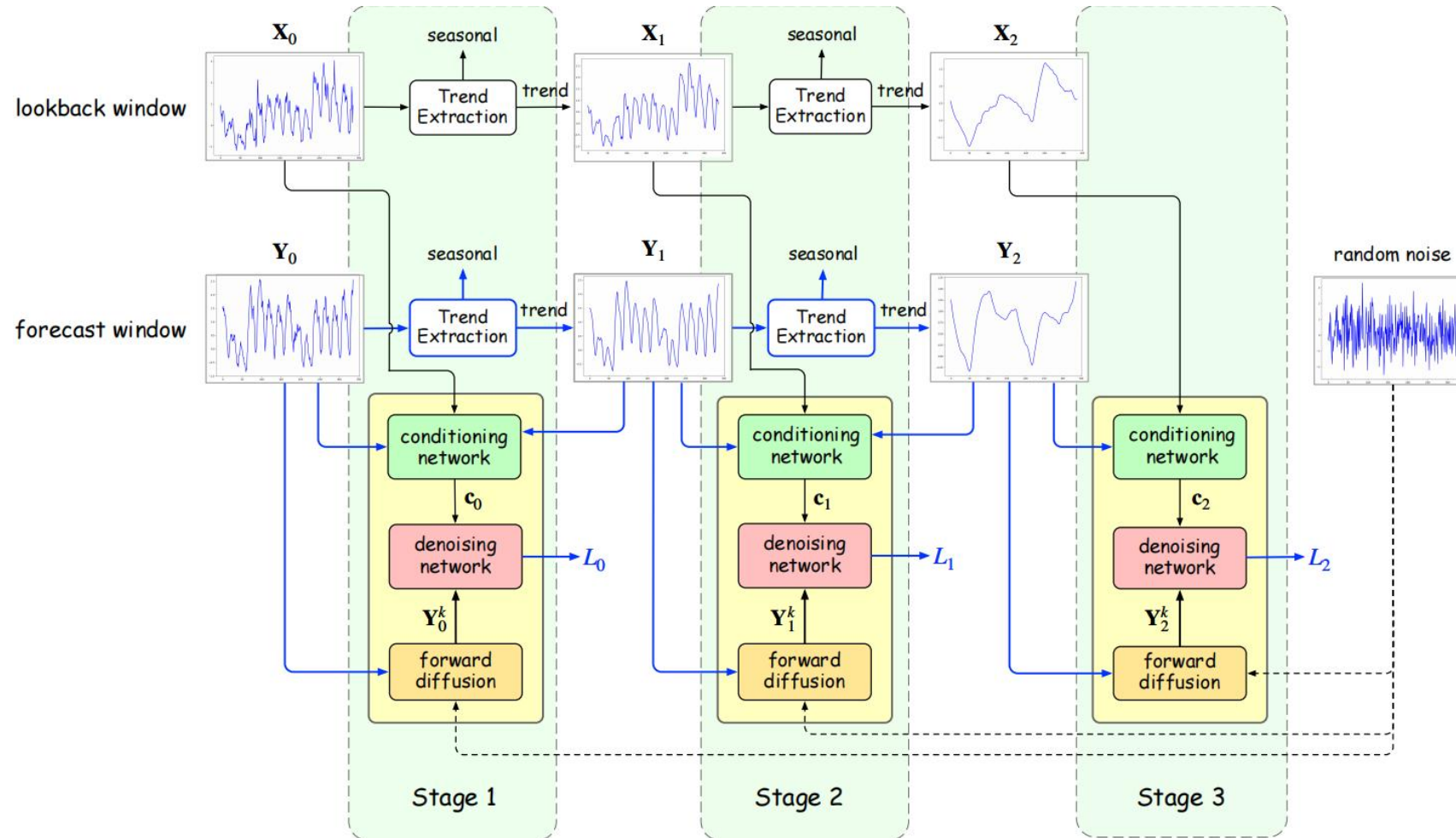
- **seasonal-trend decomposition:** extract the seasonal and trend components
- **the coarser temporal patterns can be used to help modeling the finer patterns**

how to use multi-resolution analysis in time series diffusion models?



- coarser trends are generated first and the finer details are then progressively added.
- decompose the denoising objective into several sub-objectives, each corresponds to a particular resolution.

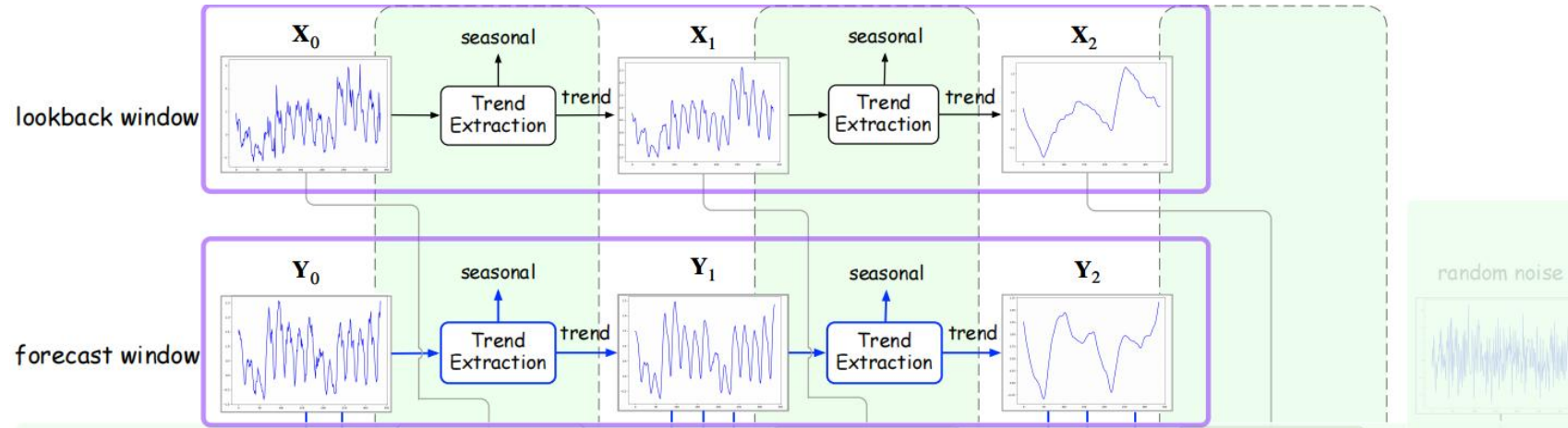
# Multi-resolution diffusion model (mr-Diff)



- in each stage, diffusion is interleaved with seasonal-trend decomposition



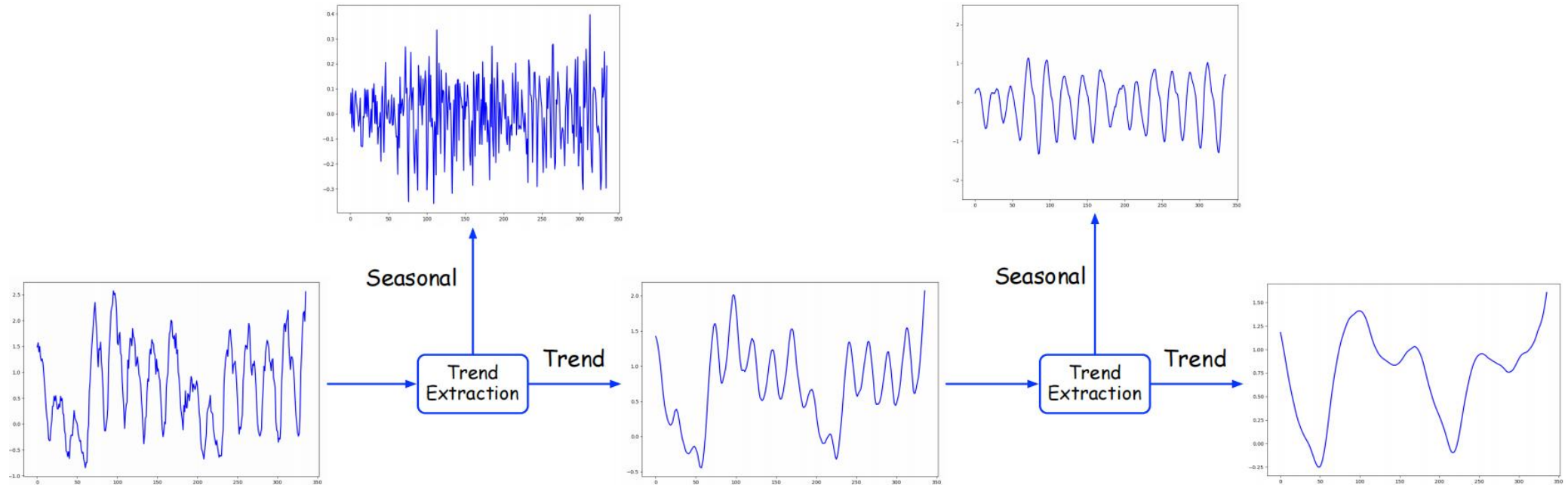
# Extract trend components successively



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

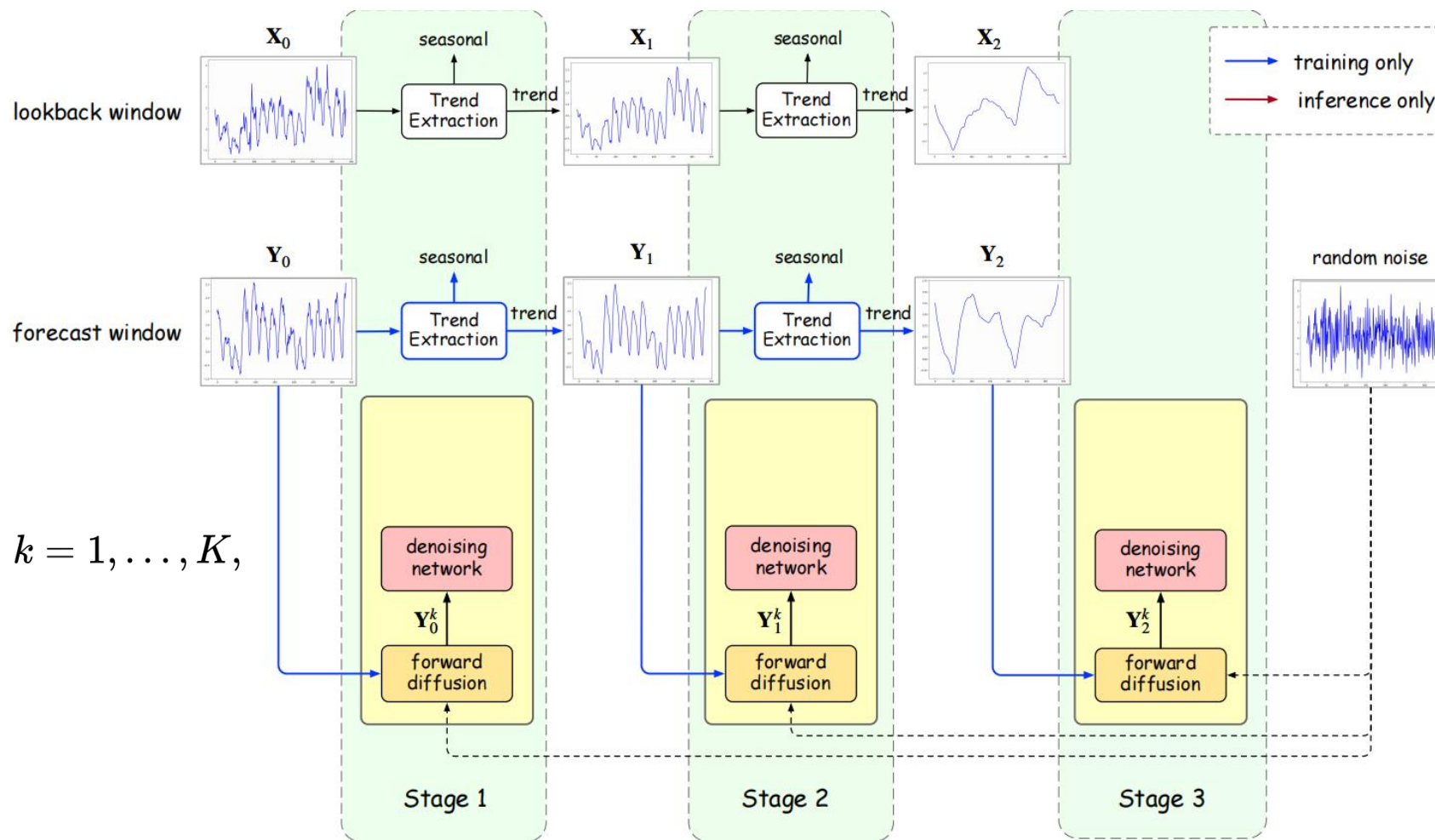
- $\mathbf{X}_s$  : trend component at stage  $s + 1$  ( $\mathbf{X}_0 = \mathbf{X}$ )
  - AvgPool: average pooling
  - Padding: keeps the lengths of  $\mathbf{X}_{s-1}$  and  $\mathbf{X}_s$  the same
  - $\tau_s$  : smoothing kernel size
    - increases with  $s$  ( trend gets coarser as  $s$  increases)
- similar processing for the segment  $\mathbf{Y}_0$  in the forecast window  
extract trend components  $\{\mathbf{Y}_s\}_{s=1, \dots, S-1}$

# Extract trend components successively



- focus here is on the trend component
- for time series, it is easier to predict a finer trend from a coarser trend
- stage  $s + 1$ : learns to reconstruct the trend component  $Y_s$  from  $Y_{s+1}$

# Forward diffusion

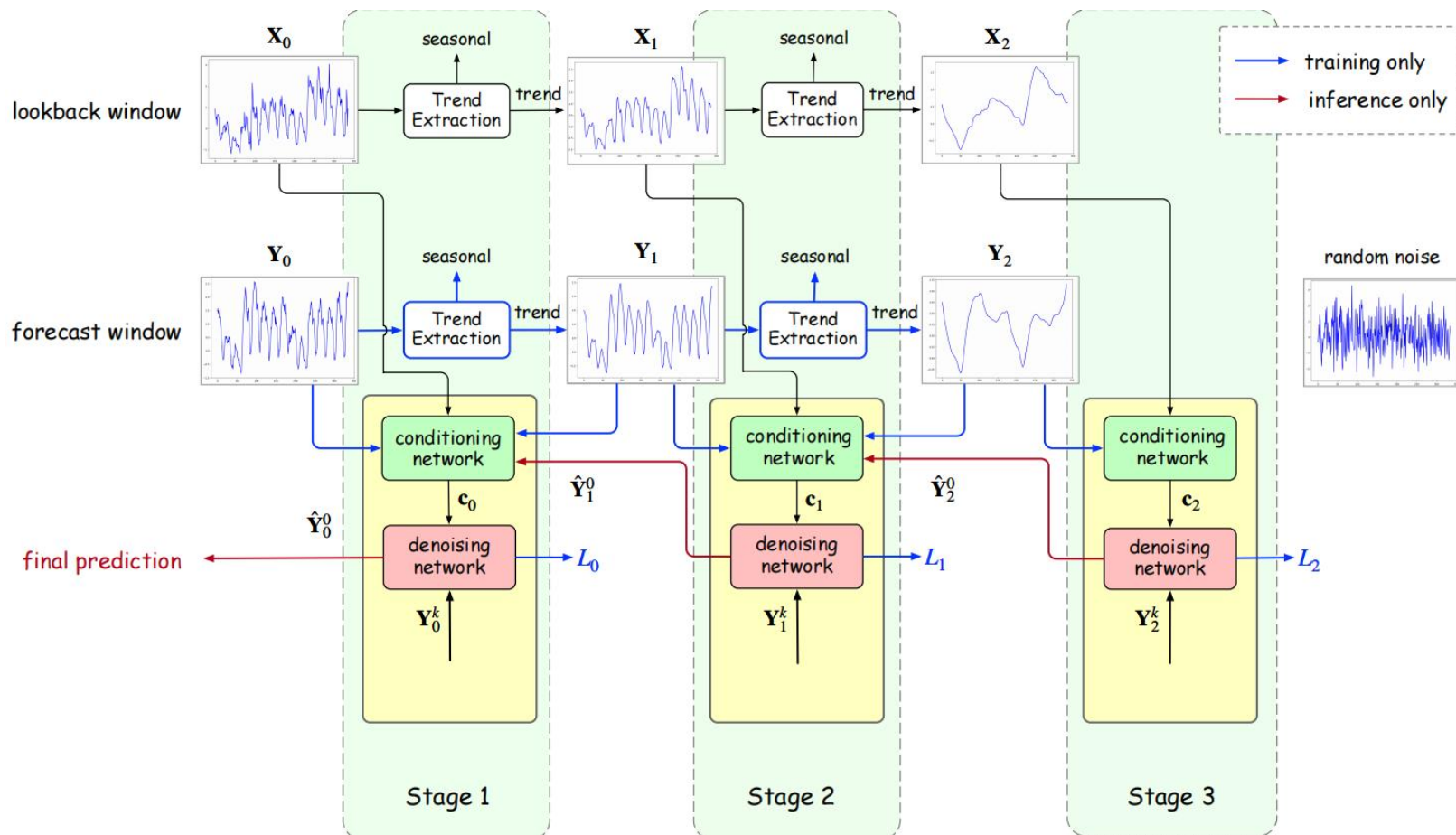


- same as in DDPM

$$\mathbf{Y}_s^k = \sqrt{\bar{\alpha}_k} \mathbf{Y}_s^0 + \sqrt{1 - \bar{\alpha}_k} \epsilon, \quad k = 1, \dots, K,$$

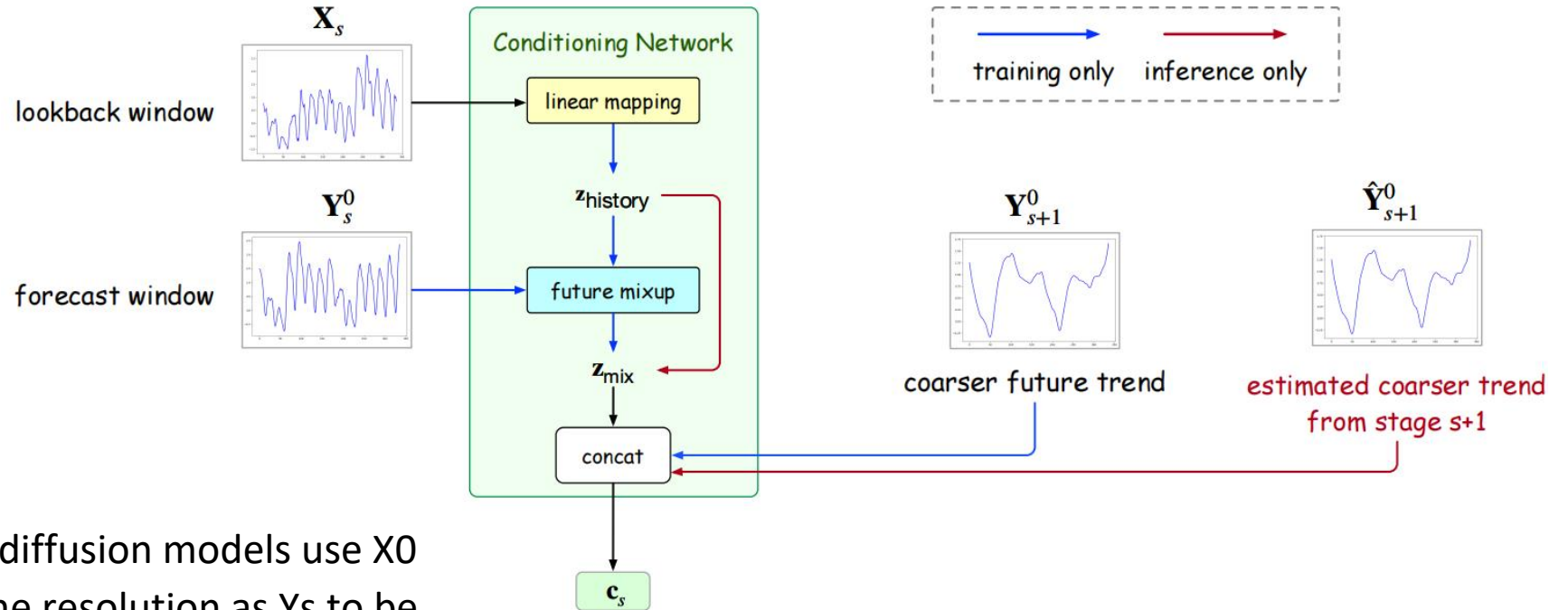
- no learnable parameters

# Backward denoising



- 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.

# Conditioning network

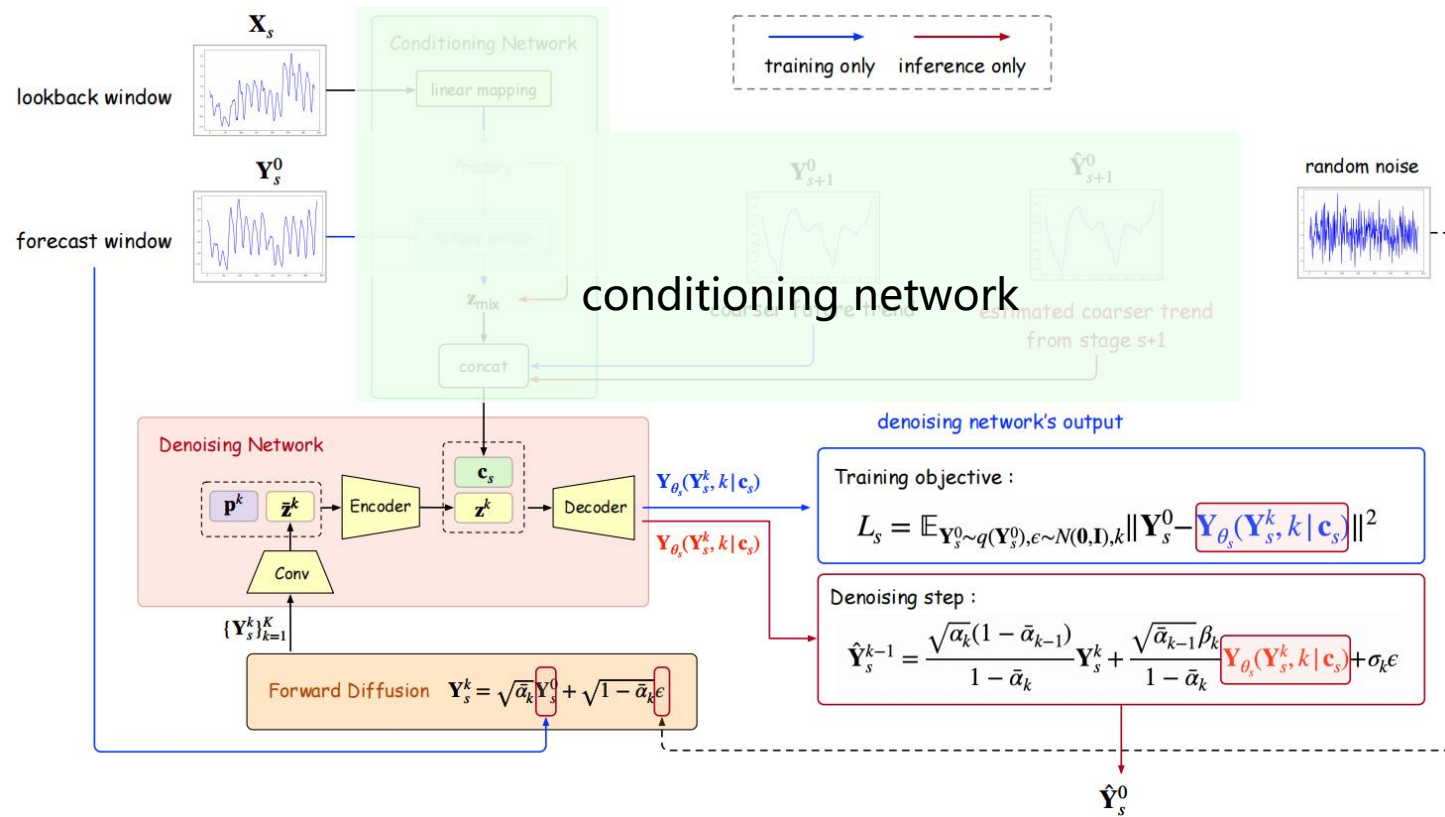


## Condition inputs:

- $\mathbf{X}_s$  (lookback segment)
  - existing time series diffusion models use  $\mathbf{X}_0$
  - ours:  $\mathbf{X}_s$  has the same resolution as  $\mathbf{Y}_s$  to be reconstructed
- $\mathbf{Y}_{s+1}$ : (ground-truth) coarser trend
  - provides an overall picture of the finer  $\mathbf{Y}_s$
- $\mathbf{Y}_s^0$ : (ground-truth) future observation
  - future-mixup with  $\mathbf{z}_{\text{history}}$  (a linear mapping on  $\mathbf{X}_s$ )
  - $s = S$  (last stage): no coarser trend and  $\mathbf{c}_s$  is simply  $\mathbf{z}_{\text{mix}}$

**Future mixup:** combine past and future time series information  $\mathbf{z}_{\text{mix}} = \mathbf{m} \odot \mathbf{z}_{\text{history}} + (1 - \mathbf{m}) \odot \mathbf{Y}_s^0$ , where  $\odot$  denotes the Hadamard product, and  $\mathbf{m} \in [0, 1)^{d \times H}$ .

# Denoising network

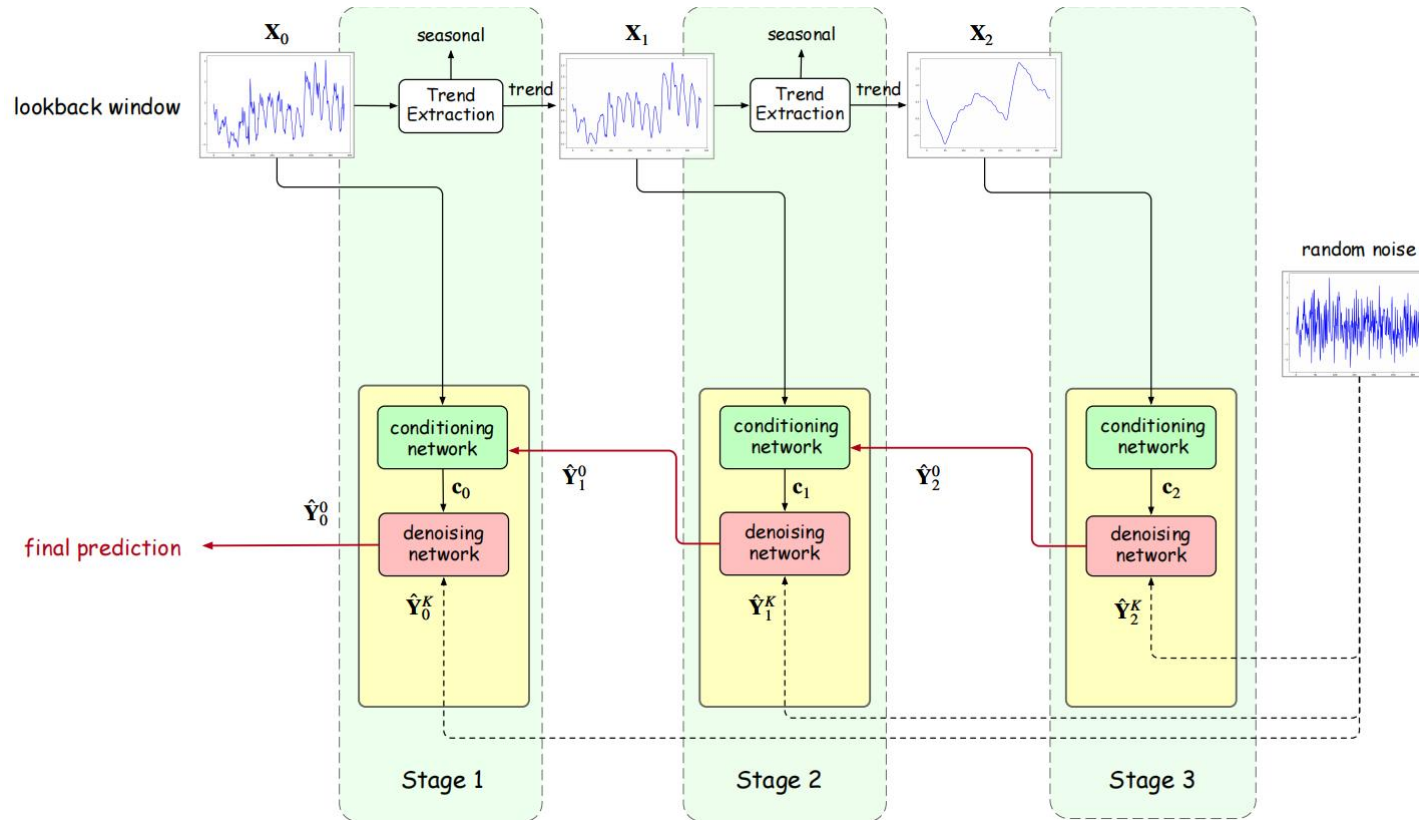


- map  $Y_s^k$  to the embedding by convolutional layers
- feed to an encoder
- concatenate with condition  $c_s$
- feed to a decoder, and output

- during training, one denoising objective one for each stage

$$\min_{\theta_s} \mathcal{L}_s(\theta_s) = \min_{\theta_s} \mathbb{E}_{Y_s^0 \sim q(Y_s^0), \epsilon \sim \mathcal{N}(0, I), k} \|Y_s^0 - Y_{\theta_s}(Y_s^k, k | c_s)\|^2$$

# Inference



for each  $s = S, \dots, 1$ , start from  $\hat{Y}_s^K \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

- encourages denoising to proceed in an easy-to-hard manner
- coarser trends are generated first, finer details are progressively added
- reconstruction at stage 1 corresponds to the target time series forecast

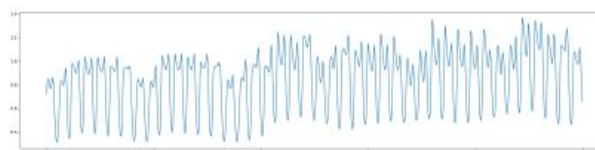
# Datasets

Summary of dataset statistics, including dimension, total observations, sampling frequency, and prediction length.

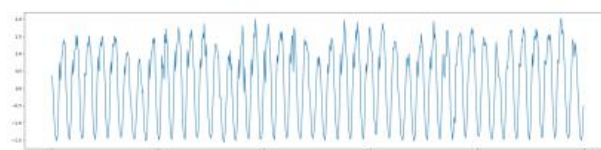
	dimension	#observations	frequency	steps ( $H$ )
<i>NorPool</i>	18	70,128	1 hour	720 (1 month)
<i>Caiso</i>	10	74,472	1 hour	720 (1 month)
<i>Traffic</i>	862	17,544	1 hour	168 (1 week)
<i>Electricity</i>	321	26,304	1 hour	168 (1 week)
<i>Weather</i>	21	52,696	10 mins	672 (1 week)
<i>Exchange</i>	8	7,588	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)



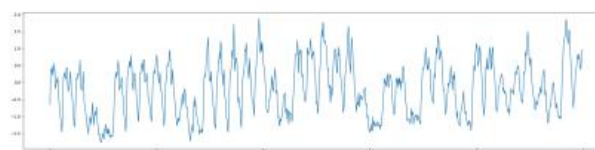
(a) *NorPool*.



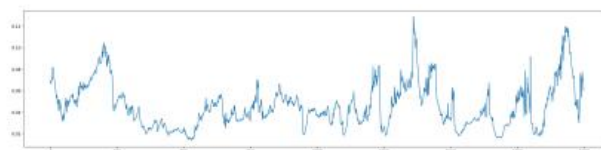
(b) *Caiso*.



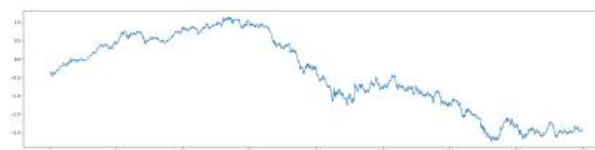
(c) *Traffic*.



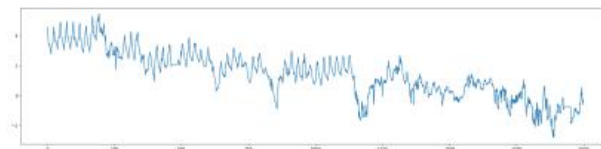
(d) *Electricity*.



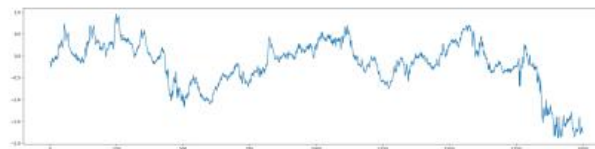
(e) *Weather*.



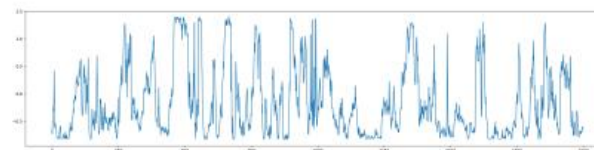
(f) *Exchange*.



(g) *ETTh1*.



(h) *ETTm1*.



(i) *Wind*.



# Baselines

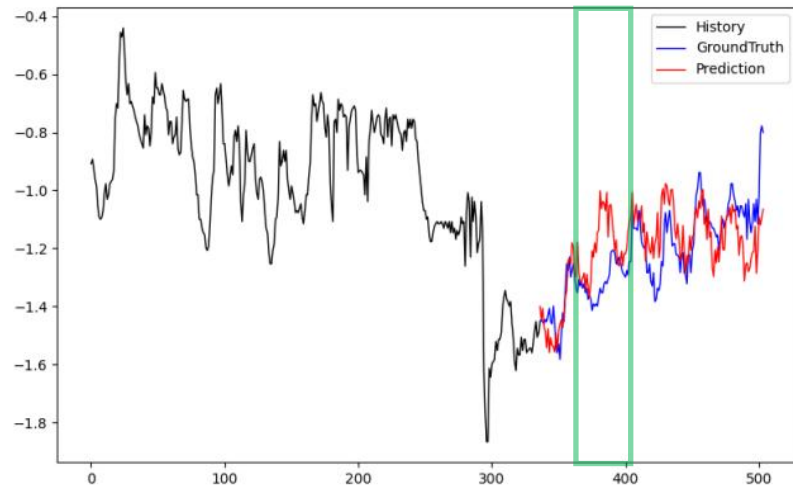
- (i) recent time series diffusion models: non-autoregressive diffusion model **TimeDiff** (Shen & Kwok, 2023), **TimeGrad** (Rasul et al., 2021), conditional score-based diffusion model for imputation (**CSDI**) (Tashiro et al., 2021), structured state space model-based diffusion (**SSSD**) (Alcaraz & 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: **NHits** (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 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).

# MAE on univariate time series

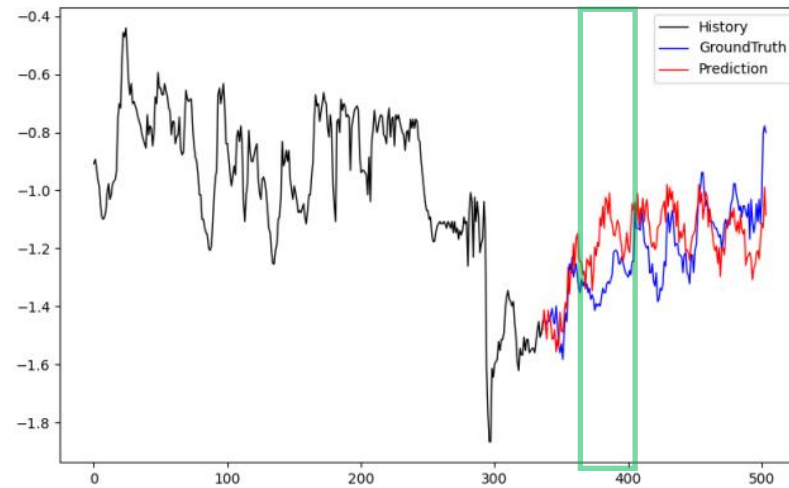
	<i>NorPool</i>	<i>Caiso</i>	<i>Traffic</i>	<i>Electricity</i>	<i>Weather</i>	<i>Exchange</i>	<i>ETTh1</i>	<i>ETTm1</i>	<i>Wind</i>	avg rank
<b>mr-Diff</b>	<b>0.609</b> <sub>(2)</sub>	0.212 <sub>(4)</sub>	<b>0.197</b> <sub>(2)</sub>	<b>0.332</b> <sub>(1)</sub>	<b>0.032</b> <sub>(1)</sub>	<b>0.094</b>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

- mr-Diff is the best in 5 of the 9 datasets on remaining 4 datasets,
- mr-Diff ranks second in 3 of them

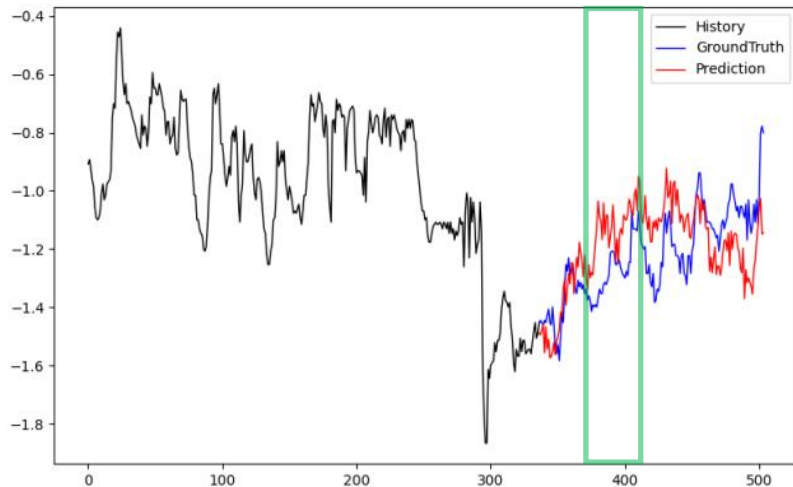
# Example prediction results on ETTh1



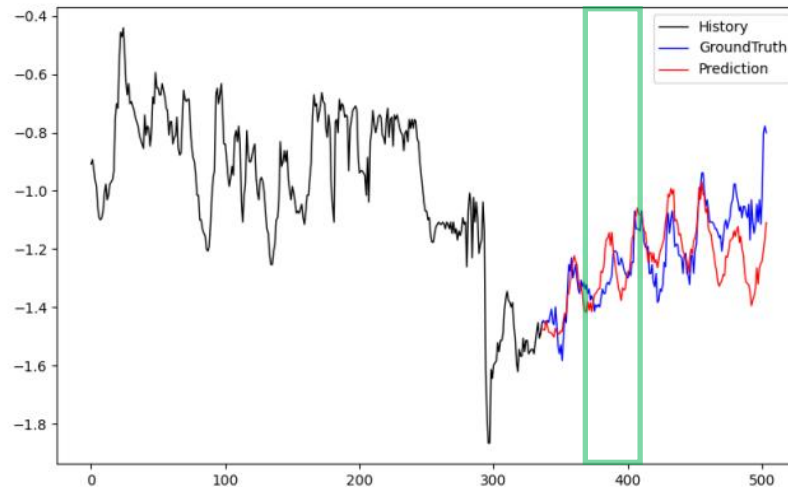
(a) SCINet.



(b) NLinear.



(c) TimeDiff.



(d) mr-Diff.

By progressively denoising the time series in a coarse-to-fine manner, mr-Diff produces higher-quality predictions than the others.



## Conclusions

- multi-resolution diffusion (**mr-Diff**), a new cascaded diffusion model for time series
- incorporates seasonal-trend decomposition and uses multiple temporal resolutions in both the diffusion and denoising processes
- progressive denoising the time series in a coarse-to-fine manner
  - more reliable predictions
- experiments show that mr-Diff outperforms the state-of-the-art time series diffusion models

