



IEEE ICDM 2023

23rd IEEE International Conference on Data Mining

December 1-4, 2023

Shanghai, China

## IEEE ICDM 2023 Tutorial



# Robust Time Series Analysis and Applications: An Interdisciplinary Approach



Qingsong Wen



Linxiao Yang



Tian Zhou



Weiqi Chen



Bingqing Peng



Liang Sun

Decision Intelligence Lab, DAMO Academy, Alibaba Group

- **Date and Time:** December 3<sup>rd</sup> (Sunday), 2:30-5:30
- **Location:** Room 6



# Who we are

- Alibaba DAMO Academy - Decision Intelligence Lab

- Research Focus:
  - AI for Time Series (AI4TS)
  - Explainable Artificial Intelligence (XAI)
  - Optimization
  - Others: ML, RL, ...
- Products and Applications:
  - Green Energy AI and Weather Forecasting:
    - Weather Forecasting,
    - Energy Forecasting and Scheduling
  - Optimization Solver:
    - MindOpt
  - AIOps:
    - Cloud Autoscaling AHPA
  - ...
- More Information:
  - <https://damo.alibaba.com/labs/decision-intelligence>

The screenshot shows the Alibaba DAMO Academy website. At the top, there is a navigation bar with links for 首页 (HOME), 实验室 (LABORATORIES), and 合作生态 (COLLABORATION). The LABORATORIES link is highlighted in red. Below the navigation, there are three main categories: Machine Intelligence, Data Computing, and Robotics. Under Machine Intelligence, there are Speech Lab, Vision Lab, Language Technology Lab, and Decision Intelligence Lab (which is also highlighted with a red border). Under Data Computing, there are Computing Technology Lab, Data Analytics and Intelligence Lab, Database and Storage Lab, and OS Lab. Under Robotics, there is Autonomous D. At the bottom of the page, there is a large image of a hand pointing at a digital interface displaying various charts and graphs, with the text "We are committed to the R&D of cutting-edge machine learning, deep learning, and AI technologies, and have made significant contributions to the field."



# Outline

- Introduction** (by Bingqing Peng)
- Preliminaries** (by Bingqing Peng)
- Robust Time Series Processing Blocks** (by Weiqi Chen)

- Time Series Periodicity Detection
- Time Series Trend Filtering
- Time Series Seasonal-Trend Decomposition
- Time Series Similarity

(10 minutes break)

- Robust Time Series Applications and Practices**

- Forecasting: Tree Models, Deep Ensemble, Transformers, etc. (by Tian Zhou)
- Autoscaling: Query Modeling, Scaling Decision, etc. (by Tian Zhou)
- Anomaly Detection: Decomposition Model, Deep State Space Model, Transformers, etc. (by Weiqi Chen)
- Fault Cause Localization: Rule Set Learning, Root Cause Analysis, etc. (by Linxiao Yang)
- XAI for Time Series Data (by Linxiao Yang)

(Q&A Session)



# Outline

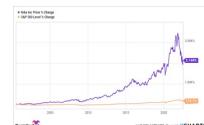
- *Introduction*
- Preliminaries
- Robust Time Series Processing Blocks
- Robust Time Series Applications and Practices



# Time Series Data is Ubiquitous

- A wide range of time series data
  - Retail
  - Finance
  - IoT
  - Healthcare
  - Energy
  - ...

stocks



sales



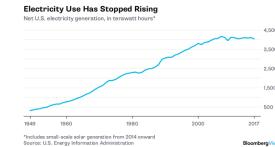
goods consumption



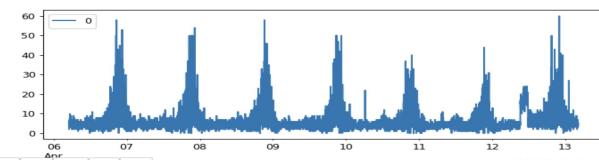
sensor



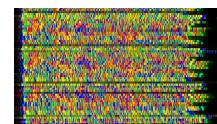
power demand



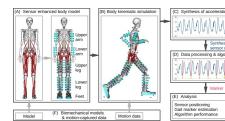
Cloud service monitoring



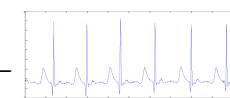
DNA sequence



motion detect

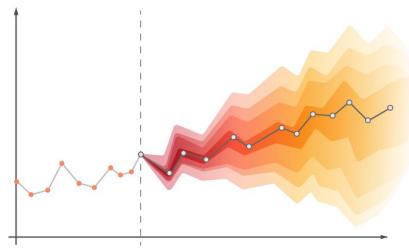


ECG

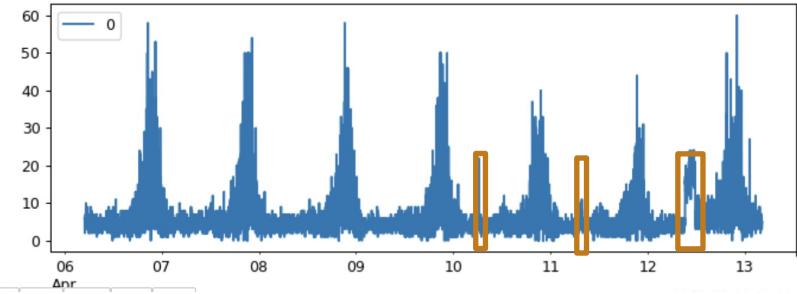


# Typical Applications of Time Series

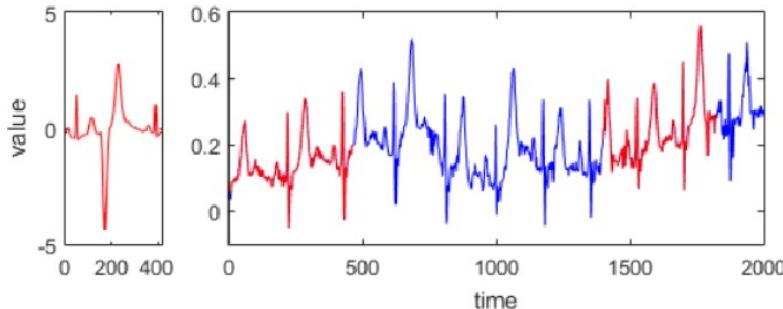
Time Series Forecasting



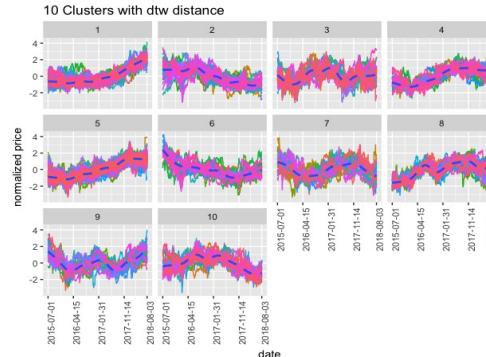
Time Series Anomaly Detection



Time Series Search/Query



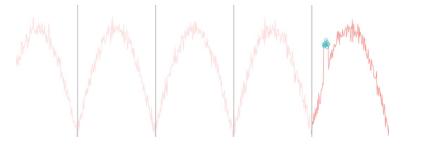
Time Series Classification/Clustering



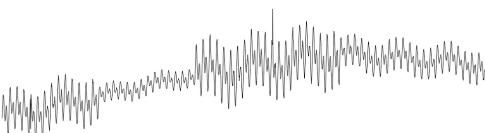
# AIOps: From Anomaly Detection to Root Cause Analysis

## Different Types of Anomalies

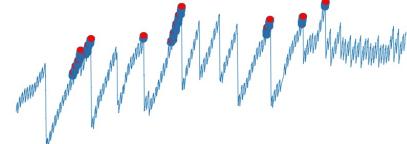
Spike & dip



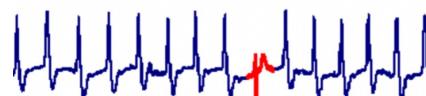
Change of mean



Change of variance



Long-term trend change



Subsequence anomaly

- Business service
- Software service
- Container and VM
- Server and components
- Network
- Data center infrastructure

In typical applications of Elastic Computing Service at Alibaba Cloud (2020):

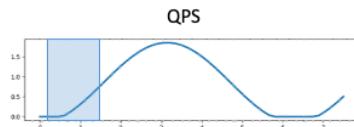
- 5M+ servers monitored
- 50M+ metrics monitored
- In 2020 some typical root identification time is ~ 1 person-month
- Some cases incur 1M+ USD loss

- ✓ *Huge amount of data calls for automatic and accurate anomaly detection algorithm*
- ✓ *Finding the root is the ultimate goal*



# From Forecasting to Decision-making: AutoScaling

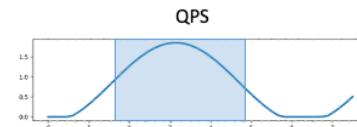
- Autoscaling in cloud computing is an effective method to improve the usage of computing resources
  - It automatically allocates resources for cloud-based applications while maintaining SLA (service level agreement)
  - Horizontal scaling (add/delete instances or VMs) vs vertical scaling (up/downgrade CPU, RAM, network, etc.)
  - Time series forecasting and decision-making on resources



load  
balancer

instance

instance



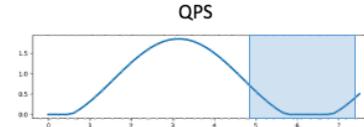
load  
balancer

instance

instance

instance

instance



load  
balancer

instance

instance



# Challenges of Time Series Analysis in Industry

- Huge amount of data
  - Hundreds of millions metrics or even more
  - **High frequency:** many time series data are sampled with higher frequency (e.g., every minute or several minutes)
  - **Very limited labels** in many applications
  - Many applications (e.g., anomaly detection) require fast or real-time processing
  - Low deployment cost
- Time series data with different complex types of characteristics
  - Noises and outliers
  - Periodicity/Seasonality: no/single/multiple periodic components, periodic component shift/change
  - Fully automatic solution is challenging but highly preferred
- Closely connected with other applications
  - Anomaly detection is closely related to root cause analysis
  - Forecasting is closely related to decision-making (e.g., autoscaling)

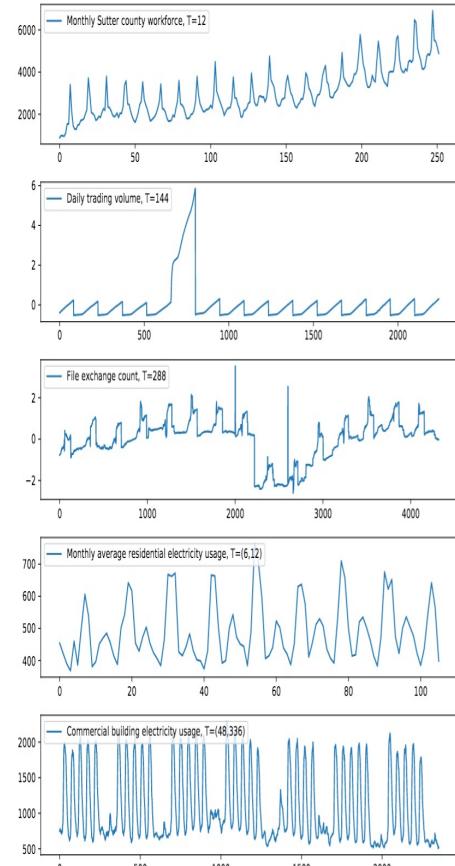
*Our solution: **robust** algorithm family for time series analysis*



# Why Periodicity so Important?

- Periodic time series dominate in many real-world applications, e.g.
  - In AIOps, 60%+ lower level monitoring data is periodic, and 90%+ higher level monitoring data is periodic
  - In electricity load data, 100% system load data is periodic, and 90%+ bus load data is periodic
- Specific processing required
  - Periodic and non-periodic time series
  - Different periodic patterns
- Challenges
  - Complicated periodic patterns, e.g.
    - multiple periodic components
    - periodic shift
    - periodic change
  - Change of trend
  - Noises/outliers

Different types of periodic time series





# Outline

## □ Introduction

### ➤ *Preliminaries*

- Robust Statistics: Robust Regression, M-estimators
- Optimization Algorithms: ADMM, MM Algorithm
- Signal Processing: Fourier, Wavelet
- Deep Learning: RNN, CNN, GNN, Transformer, Data Augmentation
- Explainable Artificial Intelligence (XAI): GAM, LIME, SHAP, Integrated Gradient

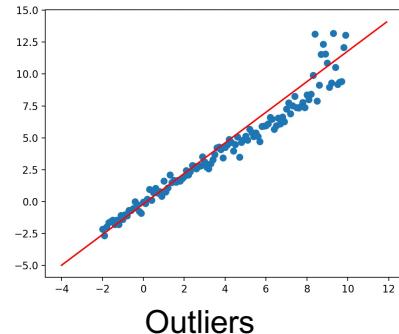
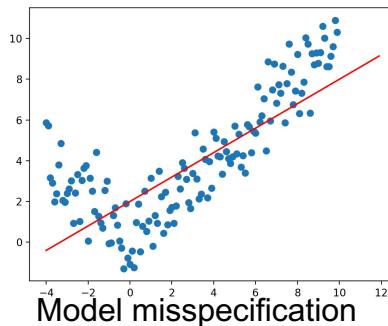
## □ Robust Time Series Processing Blocks

## □ Robust Time Series Applications and Practices



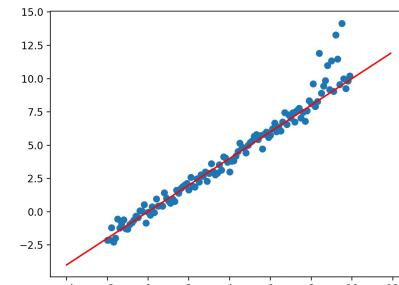
# Robust Regression

- Regression methods rely on assumptions
  - Assumptions on model, e.g. linear
  - Assumptions on noise, e.g. Gaussian
- When assumptions not hold



$$\mathbf{y} = \boxed{f(\mathbf{x}; \boldsymbol{\theta})} + \boxed{\epsilon}$$

linear      Gaussian



- Robust regression
  - Robust regression methods are designed to be not overly affected by violations of assumptions by the underlying data-generating process.



## M-estimator

- Generalized maximum likelihood estimation

$$\theta^* \leftarrow \arg \min_{\theta} \sum_{i=1}^N r^2(\theta)$$

Assuming Gaussian Error      residual

$\theta^* \leftarrow \arg \min_{\theta} \sum_{i=1}^N \rho(r(\theta))$

M-Estimator

- Properties of  $\rho$

- Continuous
- Symmetrical  $\rho(\xi) = \rho(-\xi)$
- $\rho(0) = 0$
- Positive:  $\rho(\xi) > 0$  for  $\xi \neq 0$
- Increasing monotonic  $\rho(\xi) > \rho(\xi')$  for  $|\xi| > |\xi'|$ , but it does not increase too much as  $\xi$  increases



# M-estimator

- Robust M-estimator

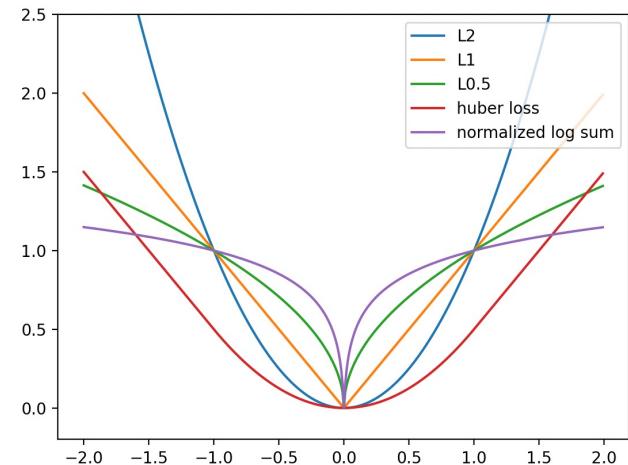
- L1-norm

$$\rho(\xi) = |\xi|$$

- Huber-loss

$$\rho_\epsilon(\xi) = \begin{cases} \frac{1}{2}\xi^2 & |\xi| \leq \epsilon \\ \epsilon(|\xi| - \frac{1}{2}\epsilon) & \text{otherwise} \end{cases}$$

- $L_p$  ( $p < 1$ ) pseudo-norm
  - Log-sum loss





# Alternating Direction Method of Multipliers

- Consider the problem

$$\min f(\mathbf{x}) + g(\mathbf{z}) \quad \text{s.t. } \mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z}$$

- The augmented Lagrangian of the problem is given as

$$L_\rho(\mathbf{x}, \mathbf{z}, \boldsymbol{\lambda}) = f(\mathbf{x}) + g(\mathbf{z}) + \boldsymbol{\lambda}^T(\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{B}\mathbf{z}) + \frac{\rho}{2}\|\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{B}\mathbf{z}\|^2$$

- ADMM process

$\rho > 0$  is called the penalty parameter

$$\mathbf{x}^{(k+1)} = \arg \min_{\mathbf{x}} L_\rho(\mathbf{x}, \mathbf{z}^{(k)}, \boldsymbol{\lambda}^{(k)})$$

$$\mathbf{z}^{(k+1)} = \arg \min_{\mathbf{z}} L_\rho(\mathbf{x}^{(k+1)}, \mathbf{z}, \boldsymbol{\lambda}^{(k)})$$

$$\boldsymbol{\lambda}^{(k+1)} = \boldsymbol{\lambda}^{(k)} + \rho(\mathbf{y} - \mathbf{A}\mathbf{x}^{(k+1)} - \mathbf{B}\mathbf{z}^{(k+1)})$$



# Alternating Direction Method of Multipliers

- Advantages of ADMM:
  - Not require gradient
  - Converges to modest accuracy within a few tens of iterations, sufficient for many applications in ML
- Apply to robust regression:

$\min_{\mathbf{x}}$   $f(\mathbf{x}) + g(\mathbf{x})$   
data fitting term  
regularization term,  
such as L1 norm, may  
be difficult to optimize  
**regression problem**

$\min_{\mathbf{x}} f(\mathbf{x}) + g(\mathbf{z}) \quad \text{s.t. } \mathbf{x} = \mathbf{z}$   
**variable split**

$$\begin{aligned}\mathbf{x}^{(k+1)} &= \text{prox}_{f, 1/\rho}(\mathbf{z}^{(k)} - \frac{1}{\rho} \boldsymbol{\lambda}^{(k)}) \\ \mathbf{z}^{(k+1)} &= \text{prox}_{g, 1/\rho}(\mathbf{x}^{(k+1)} + \frac{1}{\rho} \boldsymbol{\lambda}^{(k)}) \\ \boldsymbol{\lambda}^{(k+1)} &= \boldsymbol{\lambda}^{(k)} + \rho(\mathbf{x}^{(k+1)} - \mathbf{z}^{(k+1)})\end{aligned}$$
$$\text{prox}_f(v) = \arg \min_{x \in \mathcal{X}} \left( f(x) + \frac{1}{2} \|x - v\|_X^2 \right)$$

**update using ADMM**



# Majorization-Minimization (MM) Algorithms

- Consider the problem

$$\min_{\mathbf{x}} f(\mathbf{x}) \quad \text{difficult to optimize directly}$$

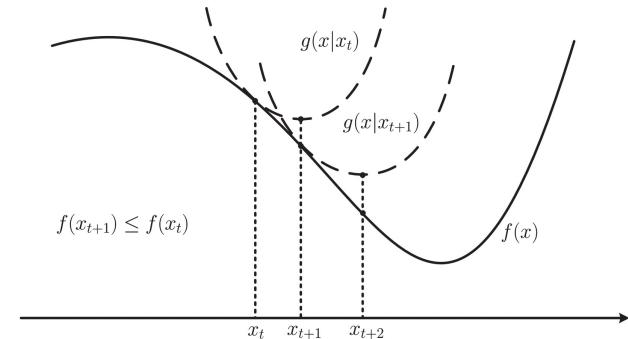
- Surrogate function

$$\underline{g(\mathbf{x}|\mathbf{x}_t)} \geq f(\mathbf{x}) + c_t \quad \text{where } c_t = g(\mathbf{x}_t|\mathbf{x}_t) - f(\mathbf{x}_t)$$

Easy to optimize

- Convergence guarantee

$$f(\mathbf{x}_{t+1}) \leq g(\mathbf{x}_{t+1}|\mathbf{x}_t) - c_t \leq g(\mathbf{x}_t|\mathbf{x}_t) - c_t = f(\mathbf{x}_t)$$





# Majorization-Minimization (MM) Algorithms

- Surrogate function construction

- Taylor Expansion

$$\text{Lp norm minimization: } |\mu|^{\frac{p}{2}} \leq \frac{p}{2} |\mu^{(t)}|^{(\frac{p}{2}-1)} |\mu| \xrightarrow{|\mu|=x^2} |x|^p \leq p |x^{(t)}|^{(p-2)} x^2 \quad 0 < p \leq 1$$

- Convexity Inequality

$$f_{\text{cvx}} \left( \sum_{i=1}^n w_i \mathbf{x}_i \right) \leq \sum_{i=1}^n w_i f_{\text{cvx}} (\mathbf{x}_i)$$

Examples: *Jensen's Inequality*

- Arithmetic-Geometric Mean Inequality

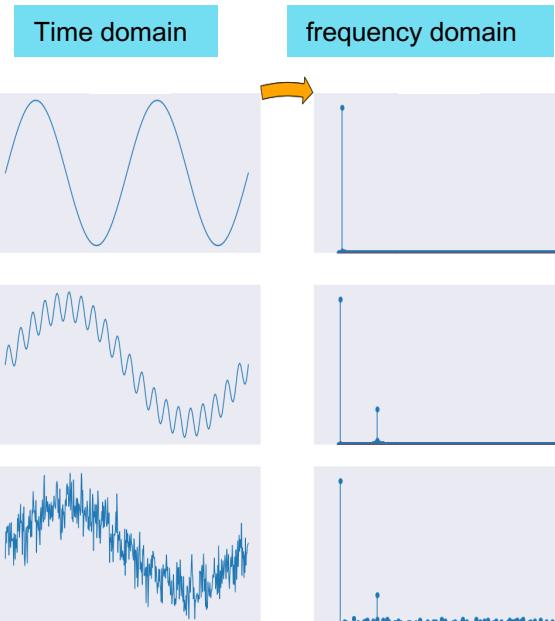
$$\text{L1 norm minimization: } x^2 + (x^{(t)})^2 \geq 2|x||x^{(t)}| \Rightarrow |x| \leq \frac{x^2 + (x^{(t)})^2}{2x^{(t)}}$$



# Signal Processing

- Time domain and frequency domain

- Time domain representation: how a signal changes over time
- Frequency domain representation: how much of a signal lies within each frequency over a band of frequencies



- Fourier transform

- Mathematical transform that decompose signal of time domain into signal of frequency domain using sine and cosine

$$\hat{f}(\xi) = \int_{-\infty}^{\infty} f(x) e^{-i2\pi\xi x} dx, \quad \forall \xi \in \mathbb{R}$$

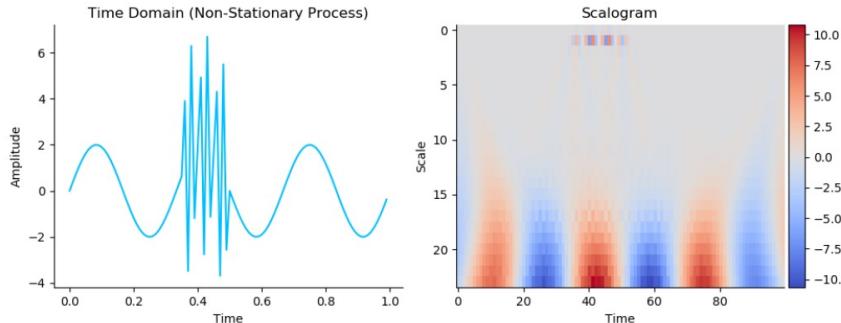


# Signal Processing

- Short-time Fourier transform (STFT)
  - Divide a long time signal into short segments of equal length and compute the Fourier transform separately in each segment
- Wavelet transform
  - Decompose the time signal by expanding by a family of orthonormal basis functions

$$X(a, b) = \frac{1}{\sqrt{a}} \int_{-\infty}^{\infty} \overline{\Psi\left(\frac{t-b}{a}\right)} x(t) dt$$

- Transform the time series from time domain to the time-scale (or time–frequency) domain





# Deep Learning: Model

- MLP (multiple layer perceptron)
  - Fully connected feedforward artificial neural network
- CNN (convolutional neural network)
  - Shared-weight architecture of filters that slide along input features
- RNN (recurrent neural network)
  - Connection between nodes form a directed or undirected graph along a temporal sequence
- Transformer
  - Deep learning model adopts the mechanism of self-attention, differentially weighting the significance of each part of input sequence

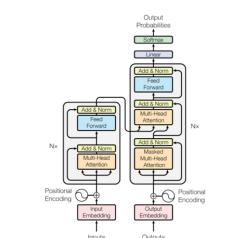
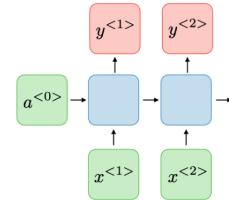
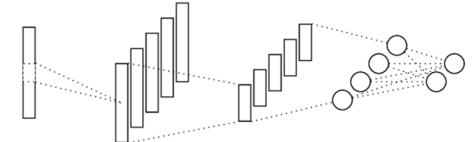
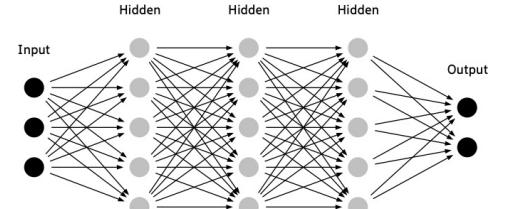
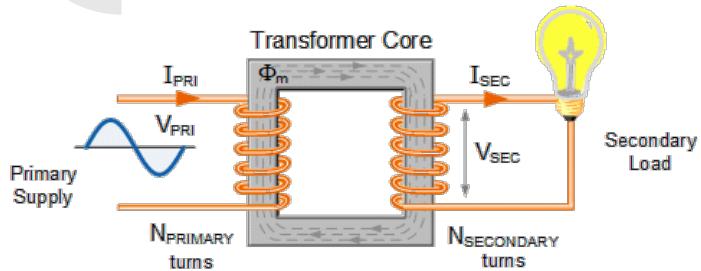


Figure 1: The Transformer - model architecture.

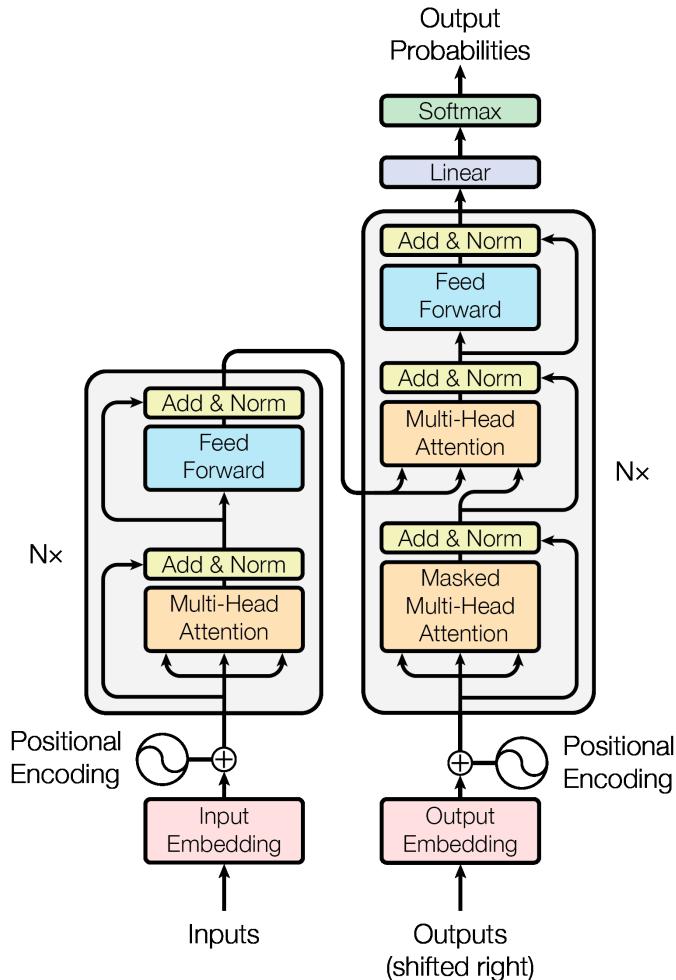


# Transformer



## Key Features:

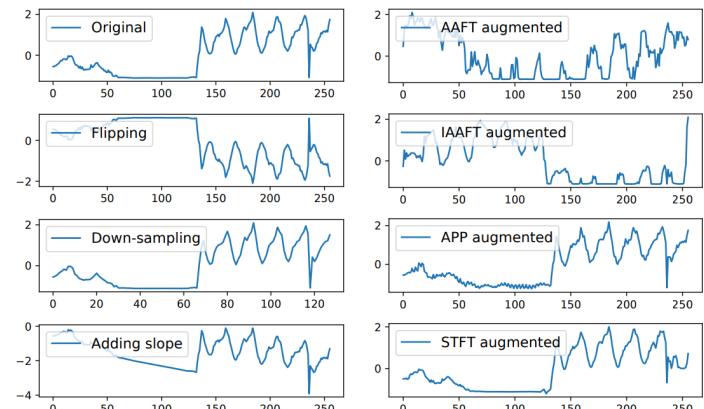
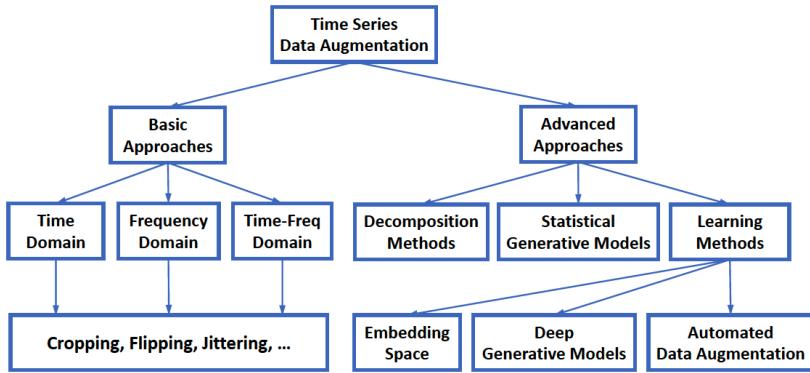
- **Attention Mechanism:** Allows the model to weigh the importance of different parts of the input sequence simultaneously.
- **Positional Encoding:** Preserves sequence order information of sequential data, thus enabling **parallelization**.





# Data Augmentation

- Taxonomy
  - Basic Approaches: time, freq, time-freq
  - Advanced Approaches: decomp, generative, learning
- Time domain
  - Flip & Down-sampling & Adding slope
  - Window warping & noise injection
- Frequency domain
  - Amplitude and phase perturbations (APP)
  - Amplitude adjusted Fourier transform (AAFT)
- Time-frequency domain
  - Short Fourier transform based augmentation (STFT)

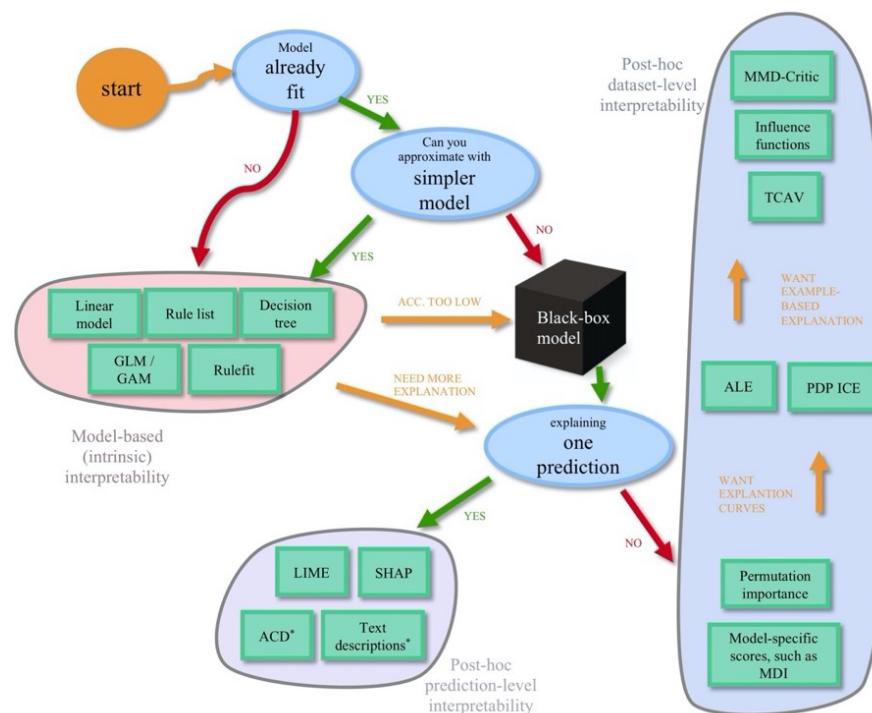


Data augmentations in time, frequency, time-frequency domains



# Explainable Artificial Intelligence (XAI)

- XAI aims to make AI systems more transparent, interpretable, and explainable to humans
  - Help users understand how AI systems work
  - Provide insights into the decision-making processes of AI models
  - Facilitates human-AI collaboration



# White Box Model: Generalized Additive Model (GAM)

Linear Models

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_m x_m + \varepsilon$$

Easy to interpret, less predictability

GAMs

$$g(E(y)) = \beta_0 + f_1(x_1) + f_2(x_2) + \dots + f_m(x_m)$$

*Link function*

*smooth nonlinear functions*

- Extensions

- Incorporate neural networks
- Allow second-order interactions between features

Black-box ML

Tree based models, DL models, etc



# How to Explain Black-box Models?

## Intuitively...

- Feature importance from tree based models
  - Correlated features
  - No depiction for complex interactions
  - Model specific computing methods
  - ...



## Better choice

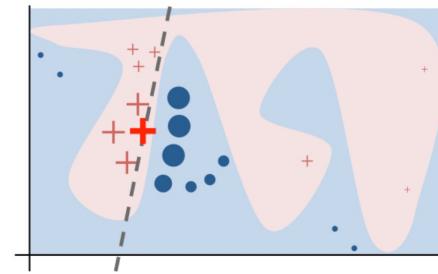
- Shapley additive explanations (SHAP)
  - Compute the contribution of each feature to the prediction fairly and robustly
$$\phi_i = \sum_{S \subseteq F \setminus \{i\}} \frac{|S|!(|F| - |S| - 1)!}{|F|!} [f_{S \cup \{i\}}(x_{S \cup \{i\}}) - f_S(x_S)]$$
  - ✓ Consistency and theoretical foundation
    - ✓ Game theory and Shapley values
  - ✓ Model Agnostic
  - ✓ Additivity
  - ✓ Local and global explanations
  - ✓ ...



# How to Explain Black-box Models?

- LIME (Local Interpretable Model-agnostic Explanations):
  - Approximate the black-box methods in a local area using interpretable models
  - Learn a local surrogate model to ensure both interpretability and local fidelity

$$\xi = \arg \min_{g \in \mathcal{G}} L(f, g, \pi_{x'}) + \Omega(g)$$



- Gradient based method
  - Vanilla gradient
  - Gradient X input
  - Smooth Gradient: average the gradients of the input with noise
  - Integrated Gradients: accumulate the gradients of samples along predefined path

$$h_g(x) = (x - \bar{x}) \odot \int_0^1 \frac{\partial g(\bar{x} + t(x - \bar{x}))}{\partial x} dt$$

- Layer-wise Relevance Propagation: propagates relevance backwards through the network



# Outline

## ❑ Introduction

## ❑ Preliminaries

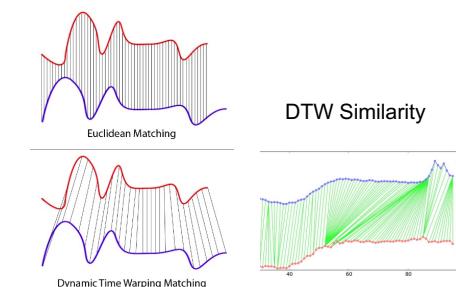
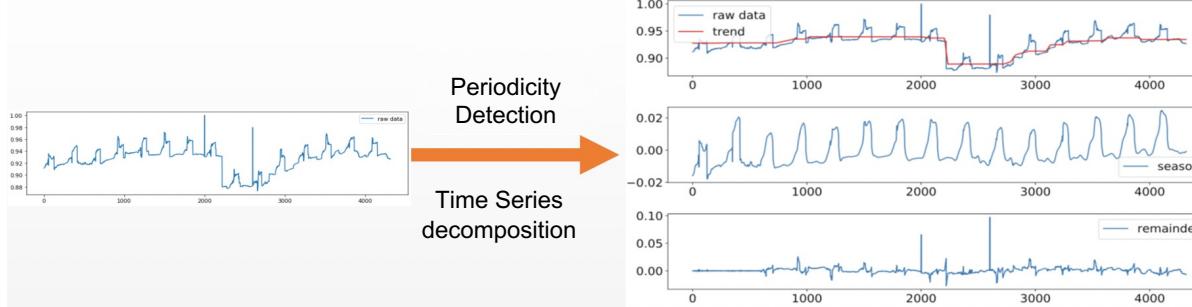
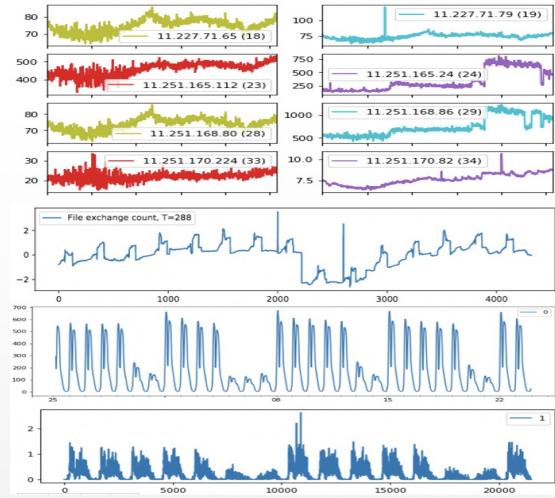
### ➤ ***Robust Time Series Processing Blocks***

- Time Series Periodicity Detection
- Time Series Trend Filtering
- Time Series Seasonal-Trend Decomposition
- Time Series Similarity

## ❑ Robust Time Series Applications and Practices

# Complex Periodic Time Series

- Time series
  - Periodicity/seasonality, trend, and noises are common
  - Periodicity detection, decomposition, similarity are crucial for downstream tasks (forecasting, anomaly detection, etc.)
- Challenges of robust time series processing blocks
  - Noise, outliers, missing data, abrupt trend changes
  - None/single/multiple periodicity
  - Need automatic and accurate processing for large-scale industrial time series



# Common Periodicity Detection Algorithms (ACF/FFT)

- Time domain: autocorrelation function (ACF)

- Calculate normalized ACF for time series  $x_t$

$$ACF(t) = \frac{1}{(N-t)\delta_x^2} \sum_{n=0}^{N-t-1} x_n x_{n+t}$$

- Calculate the (mean/median) distance of ACF peaks exceeding predefined threshold as period length

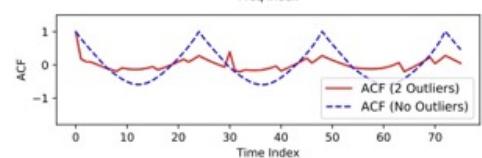
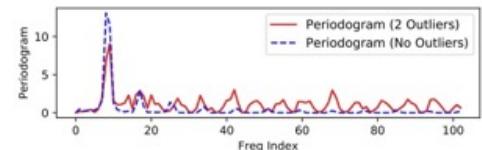
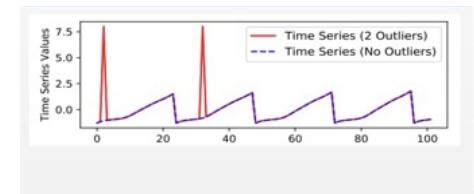
- Freq domain: Fisher's test by Periodogram (FFT)

- Fisher's test for dominant periodicity detection

$$g = \max_k P_k / (\sum_{j=1}^N P_j) \quad P_k = \|\text{DFT}\{x_t\}\|^2 = \frac{1}{N'} \sum_{t=0}^{N'-1} \|x_t e^{-i2\pi kt/N'}\|^2$$

- Distribution of the periodogram:  $P_k \sim (1/2)\sigma^2 \chi^2(2)$

- If passed Fisher's test, period length:  $N'/k$  where  $k = \arg \max_k P_k$



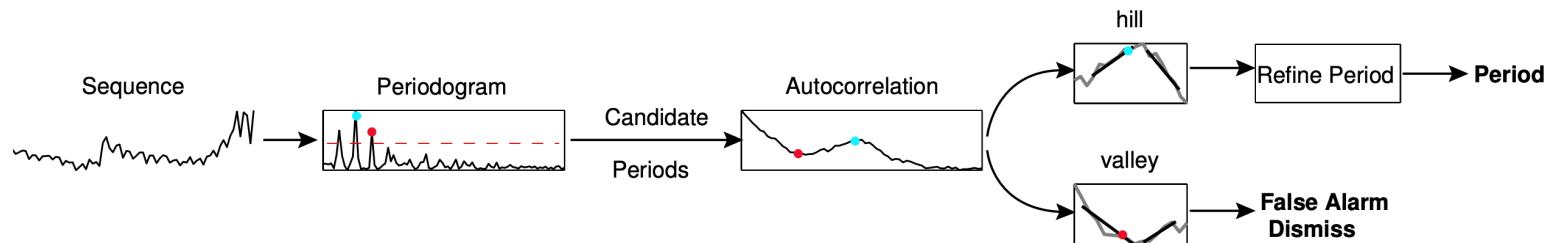
- Not robust to outliers
- Hard to handle multiple periodicities



# AUTOPERIOD

- Key idea: Combining time and frequency processing for periodicity detection
- **Pros:** more accurate than time/freq domain only methods, like ACF, FFT
- **Cons:** hard to handle multiple periodicities in complicated time series

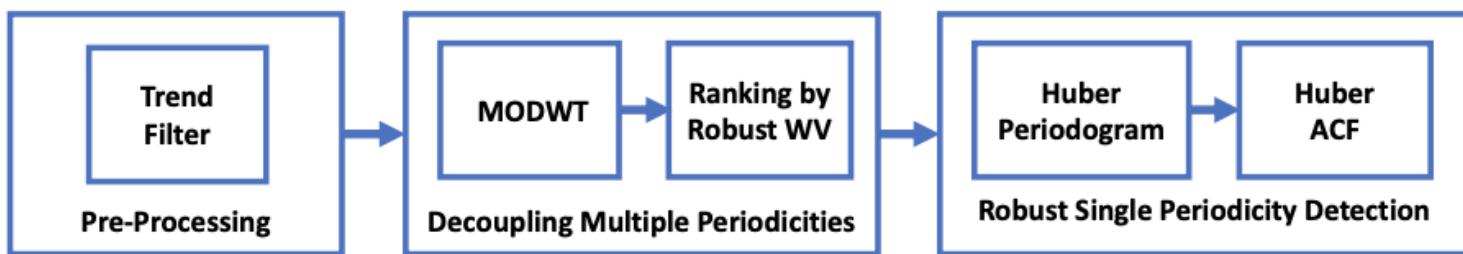
Method	Easy to threshold	Accurate short periods	Accurate large periods	Complexity
Periodogram	yes	yes	no	$O(N \log N)$
Autocorrelation	no	yes	yes	$O(N \log N)$
Combination	yes	yes	yes	$O(N \log N)$





# RobustPeriod Algorithm

- Robust and general periodicity detection: RobustPeriod
  - Key idea: *first isolate periodic components*, then detect each one robustly (“divide and conquer”)
  - Three blocks: Pre-processing, Decoupling Multiple Periodicities, Robust Single Periodicity Detection



$$y_t = \tau_t + \sum_{i=1}^m s_{i,t} + r_t$$

For periodic time series, period lengths are denoted as  $T_i, i = 1, \dots, m$ ,  $m$  is the number of periodic components

# Decouple Multiple Periodicities

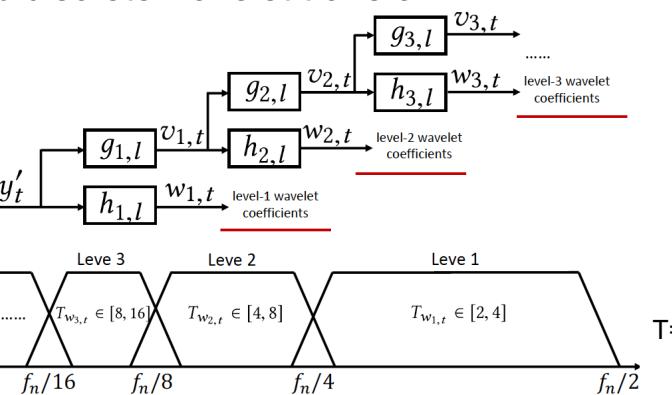
- MODWT for decoupling multiple periodicities
  - MODWT: maximal overlap discrete wavelet transform

*j*th level wavelet coefficients

$$w_{j,t} = \sum_{l=0}^{L_j-1} h_{j,l} y'_{t-l \bmod N}$$

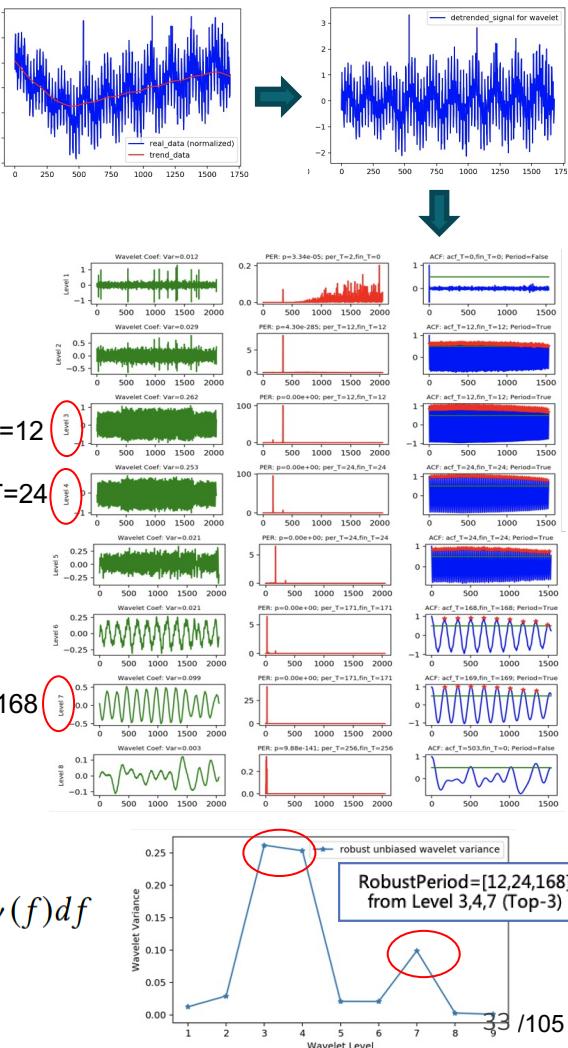
*j*th level scaling coefficients

$$v_{j,t} = \sum_{l=0}^{L_j-1} g_{j,l} y'_{t-l \bmod N}$$



- Wavelet variance ranking for speedup

- Wavelet variance decomposition:  $\hat{\sigma}_{y'}^2 = \sum_{j=1}^{J_0} \hat{\sigma}_{w_j}^2 + \hat{\sigma}_{v_{J_0}}^2$
- Relationship to power spectral density (PSD):  $\hat{\sigma}_{w_j}^2 \approx \int_{1/2^{j+1} \leq |f| \leq 1/2^j} S_{y'}(f) df$
- If there is a periodic component in the  $j$ th level wavelet coefficient, a large wavelet variance would be expected



# RobustPeriod: Robust Single Periodicity Detection

- Robust Huber-periodogram for Fisher's test

- Huber-periodogram: M-periodogram using Huber loss

$$P_k^M = \frac{N'}{4} \left| \hat{\beta}_M(k) \right|^2 = \frac{N'}{4} \left| \arg \min_{\beta \in R^2} \gamma(\phi\beta - x) \right|^2 \quad \gamma(x_t) = \gamma_{\zeta}^{hub}(x_t) = \begin{cases} \frac{1}{2}x_t^2, & |x_t| \leq \zeta \\ \zeta|x_t| - \frac{1}{2}\zeta^2, & |x_t| > \zeta \end{cases}$$

- Similar distribution as original periodogram

$$P_k^M \stackrel{A}{\sim} (1/2)m_2S_2\chi_2^2 \quad (\text{under practical mild conditions})$$

- Robust Huber-ACF for validating periodicity candidate

- Utilizing Wiener-Khinchin theorem based on Huber-periodogram

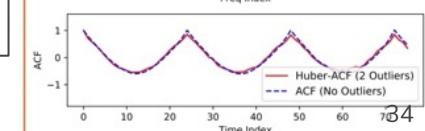
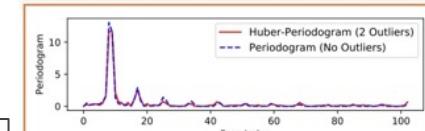
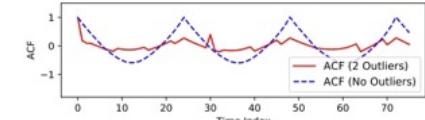
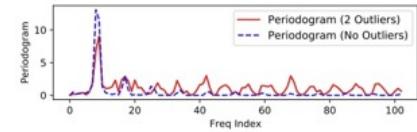
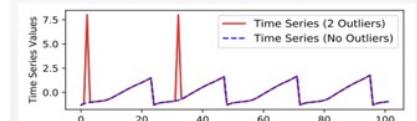
$$\text{HuberACF}(t) = \frac{p_t}{(N-t)p_0}$$

$$p_t = \text{IDFT}\{\bar{P}_k\} = \frac{1}{\sqrt{N'}} \sum_{k=0}^{N'-1} \bar{P}_k e^{i2\pi kt/N'}$$

$$\bar{P}_k = \begin{cases} P_k^M & k = 0, 1, \dots, N-1 \\ \left( \sum_{k=0}^{N-1} x_{2k} - x_{2k+1} \right)^2 / N' & k = N \\ P_{N'-k}^M & k = N+1, \dots, N'-1 \end{cases}$$

robust to outliers

Demo of Huber-Periodogram/ACF



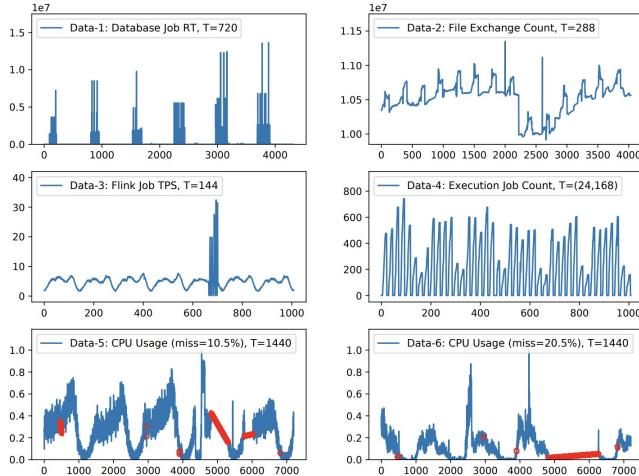


# Periodicity Detection Comparisons

- Real-world datasets:

- Alibaba Cloud monitoring data with outliers, noise, trend, missing data (**single/multiple** periodicity)

Monitoring datasets from Alibaba Cloud



Algorithms	Data-1, T=720 Database RT	Data-2, T=288 File Exchange	Data-3, T=144 Flink TPS
Siegel	(655,769,...)	(288,576,...)	(141,144)
AUTOPERIOD	(353,241,9)	(288,439,...)	(68,141)
Wavelet-Fisher	(372,745,...)	(282,585,...)	(73,146)
RobustPeriod	721	288	144

Algorithms	Data-4, T=(24,168) Job Count	Data-5, T=1440 CPU Usage	Data-6, T=1440 CPU Usage
Siegel	(24,168)	(1459,2597,...)	(1575,1063,...)
AUTOPERIOD	(24,26)	(1488,739,...)	(366,2880,...)
Wavelet-Fisher	(12,24,...)	(1489,712,...)	(1489,364,...)
RobustPeriod	(24,168)	1431	1426

Algorithms	Robust to outliers	Robust to amplitude changes	Robust to trend changes	Support multiple periods	Not need priors for number of periods
ACF	✗	✗	✗	✗	✗
FFT-Fisher	✗	✗	✗	✗	✗
FFT-Siegel	✗	✗	✗	✓	✓
Bandpass+ACF	✗	✓	✗	✓	✗
Wavelet-Fisher	✗	✗	✗	✓	✓
AUTOPERIOD	✗	✗	✗	✓	✓
<b>RobustPeriod</b>	✓	✓	✓	✓	✓



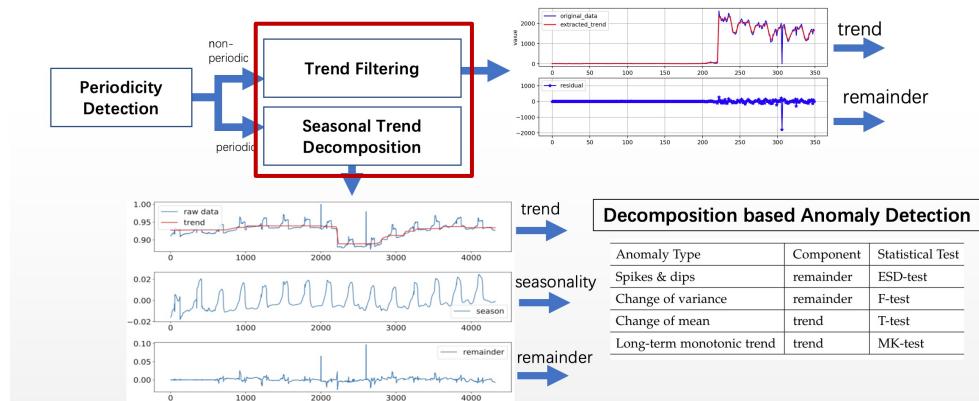
# Time Series Decomposition

- Trend filtering, seasonal-trend decomposition

$$y_t = \tau_t + r_t \quad y_t = \tau_t + \sum_{i=1}^m s_{i,t} + r_t$$

- Why need decomposition

- Insights and interpretability from different components
- Different utilization by components
  - e.g.: anomaly detection





# Non-periodic Time Series: Common Trend Filtering

$$y_t = \tau_t + r_t$$

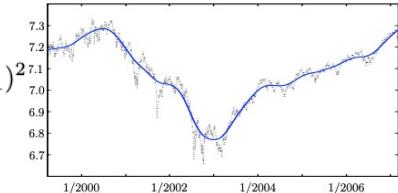
- Moving average: not accurate with delay

- Hodrick–Prescott filtering (H-P filter)

- Regularization with  $\ell_2$  norm of the second difference operator
- linear estimation, convex with closed-form solution

$$\frac{1}{2} \sum_{t=0}^{N-1} (y_t - \tau_t)^2 + \lambda \sum_{t=1}^{N-2} (\tau_{t-1} - 2\tau_t + \tau_{t+1})^2$$

$$\frac{1}{2} \|\mathbf{y} - \boldsymbol{\tau}\|_2^2 + \lambda \|\mathbf{D}^{(2)} \boldsymbol{\tau}\|_2$$

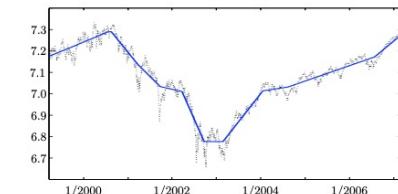


- $\ell_1$  trend filtering

- Regularization with  $\ell_1$  norm for piecewise linear trend estimation
- Nonlinear estimation, convex and linear computational complexity

$$\frac{1}{2} \|\mathbf{y} - \boldsymbol{\tau}\|_2^2 + \lambda \|\mathbf{D}^{(2)} \boldsymbol{\tau}\|_1$$

$$\mathbf{D}^{(2)} = \begin{bmatrix} 1 & -2 & 1 & & \\ & 1 & -2 & 1 & \\ & & \ddots & & \\ & & & 1 & -2 & 1 \end{bmatrix}$$



- Challenges

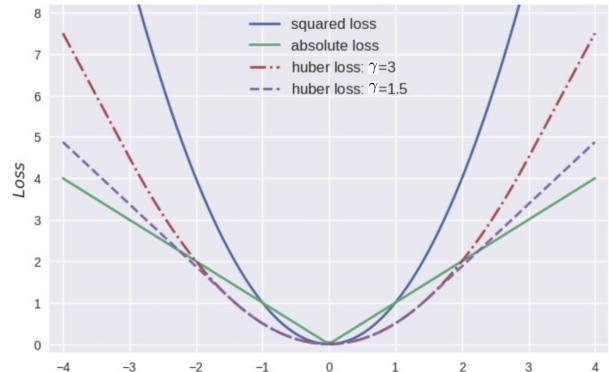
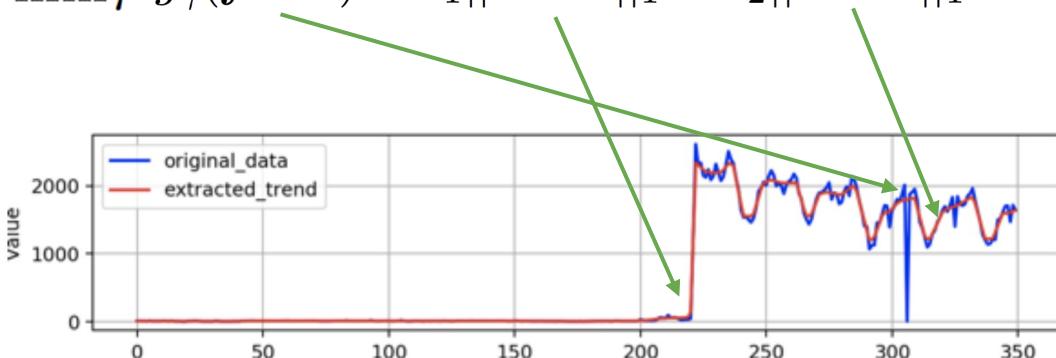
- Not robust to outliers and abrupt trend changes



# RobustTrend Filter

- Huber loss: robust to outliers
- 1st order L1 regularization: abrupt trend changes
- 2nd order L1 regularization: slow trend changes and can effectively reduce staircasing effect

$$\min_{\tau} g_\gamma(\mathbf{y} - \boldsymbol{\tau}) + \lambda_1 \|\mathbf{D}^{(1)} \boldsymbol{\tau}\|_1 + \lambda_2 \|\mathbf{D}^{(2)} \boldsymbol{\tau}\|_1 \quad \text{where}$$



$$g_\gamma(x_i) = \begin{cases} \frac{1}{2}x_i^2, & |x_i| \leq \gamma \\ \gamma|x_i| - \frac{1}{2}\gamma^2, & |x_i| > \gamma \end{cases}$$

$$\mathbf{D}^{(1)} = \begin{bmatrix} 1 & -1 & & \\ & 1 & -1 & \\ & & \ddots & \\ & & & 1 & -1 \end{bmatrix}, \quad \mathbf{D}^{(2)} = \begin{bmatrix} 1 & -2 & 1 & & \\ & 1 & 1 & -2 & 1 & \\ & & \ddots & & \ddots & \\ & & & 1 & -2 & 1 \end{bmatrix}$$



## ADMM for RobustTrend Filter

- Optimization formulation:  $\min_{\boldsymbol{\tau}} g_\gamma(\mathbf{y} - \boldsymbol{\tau}) + \lambda_1 \|\mathbf{D}^{(1)} \boldsymbol{\tau}\|_1 + \lambda_2 \|\mathbf{D}^{(2)} \boldsymbol{\tau}\|_1$

$$\begin{array}{ll} \min_{\boldsymbol{\tau}} & g_\gamma(\mathbf{y} - \boldsymbol{\tau}) + \|\mathbf{z}\|_1 \\ \text{s.t.} & \mathbf{D}\boldsymbol{\tau} = \mathbf{z} \end{array} \quad \text{where} \quad \mathbf{D} = \begin{bmatrix} \lambda_1 \mathbf{D}^{(1)} \\ \lambda_2 \mathbf{D}^{(2)} \end{bmatrix}$$

- Updating steps of ADMM: (Tau-minimization step is not efficient)

$$\boldsymbol{\tau}^{k+1} = \arg \min_{\boldsymbol{\tau}} \left( g_\gamma(\mathbf{y} - \boldsymbol{\tau}) + \frac{\rho}{2} \|\mathbf{D}\boldsymbol{\tau} - \mathbf{z}^k + \mathbf{u}^k\|_2^2 \right) \rightarrow \text{Not efficient, no closed form}$$

$$\begin{aligned} \mathbf{z}^{k+1} &= \arg \min_{\mathbf{z}} \left( \|\mathbf{z}\|_1 + \frac{\rho}{2} \|\mathbf{D}\boldsymbol{\tau}^{k+1} - \mathbf{z} + \mathbf{u}^k\|_2^2 \right) \rightarrow \text{Efficient by soft thresholding} \\ \mathbf{u}^{k+1} &= \mathbf{u}^k + \mathbf{D}\boldsymbol{\tau}^{k+1} - \mathbf{z}^{k+1} \end{aligned}$$

$$S_\rho(a) = \begin{cases} 0, & |a| \leq \rho \\ a - \rho \operatorname{sgn}(a), & |a| > \rho \end{cases}$$



## Efficient MM Algorithm for Tau–Minimization

- One-iteration *majorization minimization* (MM) for Tau–minimization step

$$g_\gamma(\mathbf{x}) \leq \eta_\gamma(\mathbf{x}|\mathbf{x}^k) = \sum_i \eta_\gamma(x_i|x_i^k) = \frac{1}{2}\mathbf{x}^T \mathbf{A}\mathbf{x} + C$$

where  $\mathbf{x} = \mathbf{y} - \boldsymbol{\tau}$  and  $\mathbf{A} = \text{diag}(g'_\gamma(x_i^k)) \text{diag}^{-1}(g_\gamma(x_i^k))$

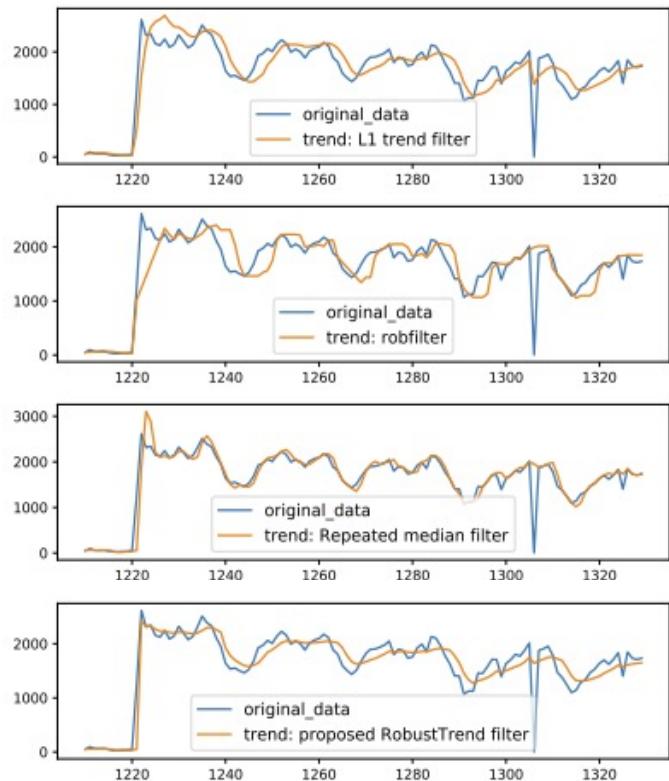
- Now the approximated Tau–minimization step

$$\boldsymbol{\tau}^{k+1} = \mathbf{y} - \rho(\mathbf{A} + \rho\mathbf{D}^T\mathbf{D})^{-1}\mathbf{D}^T(\mathbf{u}^k - \mathbf{z}^k + \mathbf{D}\mathbf{y}) \rightarrow \text{efficient with closed form}$$

- Theoretical motivation [Eckstein and Bertsekas, 1992]
  - ADMM can still converge when the updating steps are carried out approximately

# Experiments on Real-World Data

- Compare trend filters of SOTA models
- Performance highlights
  - L1 trend filter: sensitive to the outliers
  - robfilter: some delay when trend changes
  - Repeated median filter: overshoots trend estimation
  - **RobustTrend filter**: performs the best, captures the abrupt and slow trend change, and is robust to outliers

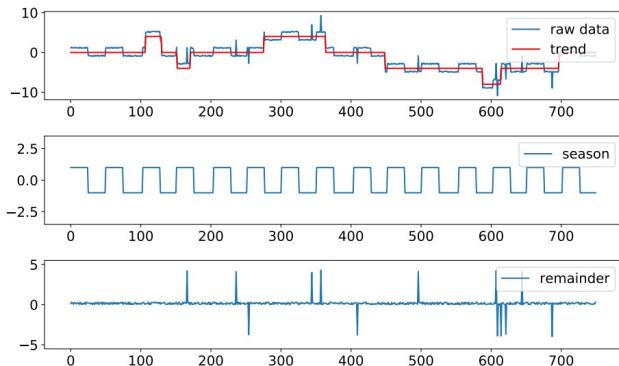




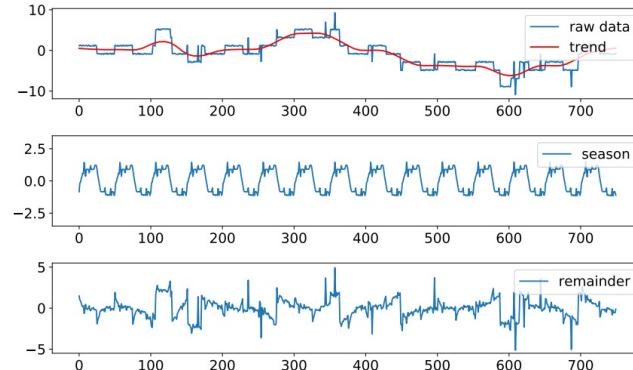
# Periodic Time Series: Common STL Algorithm

$$y_t = \tau_t + s_t + r_t$$

- Seasonal-trend decomposition using LOESS regression (STL)
  - A sequence of applications of the LOESS smoother
  - LOESS (locally estimated scatterplot smoothing): local regression
  - **Pros:** simple design, fast
  - **Cons:** not robust to trend changes and seasonality shift; cannot capture cross-period



complex synthetic data with trend changes and seasonality shift

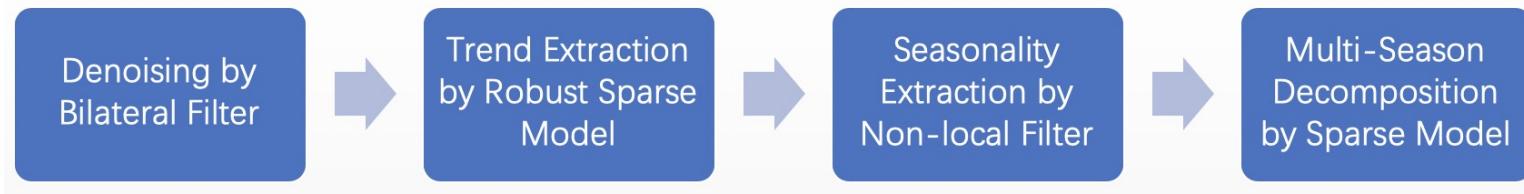


results of STL decomposition



# RobustSTL Algorithm

- Fast and robust seasonal-trend decomposition:
$$y_t = \tau_t + \sum_{i=1}^m s_{i,t} + r_t$$
  - Key idea: *sequentially and robustly* extract components in time series
  - Four blocks: noise removal, trend extraction, seasonality extraction, multi-season decomposition
  - *Efficient* GADMM for trend extraction and multi-season decomposition:  $O(N \log N)$



Qingsong Wen, Jingkun Gao, Xiaomin Song, Liang Sun, Huan Xu, Shenghuo Zhu, "RobustSTL: A Robust Seasonal-Trend Decomposition Algorithm for Long Time Series," in Proc. 33th AAAI Conference on Artificial Intelligence (AAAI 2019), Honolulu, Hawaii, Jan. 2019. ([single periodicity version](#))

Qingsong Wen, Zhe Zhang, Yan Li, Liang Sun, "Fast RobustSTL: Efficient and Robust Seasonal-Trend Decomposition for Time Series with Complex Patterns," in Proc. 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD 2020), San Diego, CA, Aug. 2020. ([multiple periodicities and high-speed version](#))

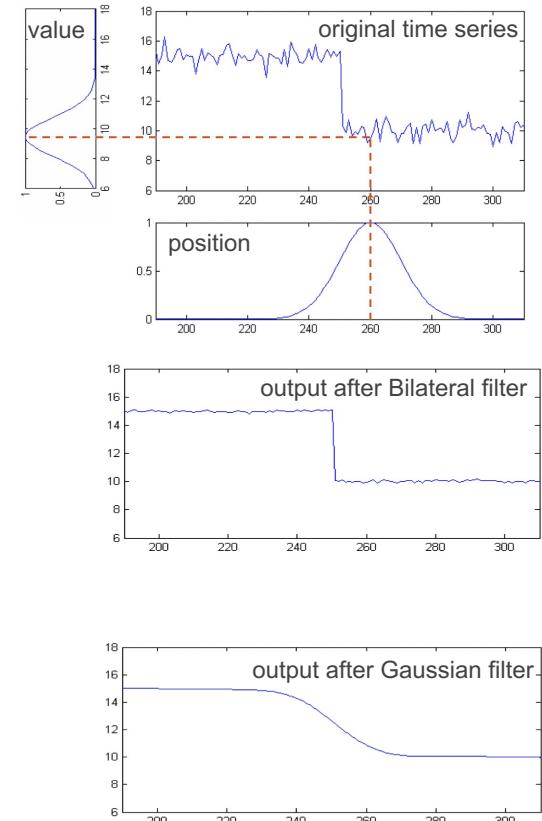
Linxiao Yang, Qingsong Wen, Bo Yang, Liang Sun, "A Robust and Efficient Multi-Scale Seasonal-Trend Decomposition," in Proc. IEEE 46th International Conference on Acoustics, Speech, and Signal Processing (ICASSP 2021), Toronto, Canada, Jun. 2021. ([multiple-scale version](#))

# Noise Removal

- Bilateral filtering
  - Consider both value and position difference in filtering
  - Smooth time series while preserving *abrupt changes*

$$y'_t = \sum_{j \in J} w_j^t y_j, \quad J = t, t \pm 1, \dots, t \pm H$$

$$w_j^t = \frac{1}{z} e^{-\frac{|j-t|^2}{2\delta_d^2}} e^{-\frac{|y_j - y_t|^2}{2\delta_i^2}}$$





# Trend Extraction

- Robust sparse model: LAD loss with two L1 regularizations

- Mitigate season effect by *seasonal difference*

$$g_t = \nabla_T y'_t = y'_t - y'_{t-T} = \nabla_T \tau_t + \nabla_T s_t + \nabla_T r'_t$$

$$T = \max T_i, i = 1, 2, \dots, m$$

- Estimate the *trend difference robustly*: LAD loss, two L1 regs

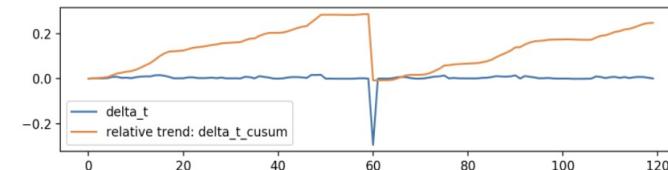
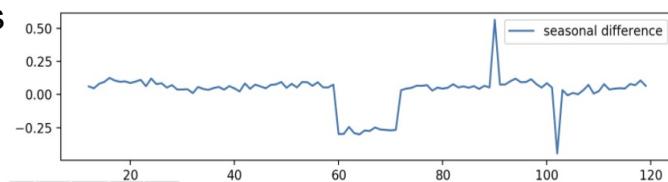
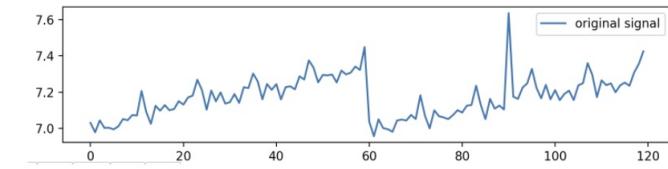
$$\sum_{t=T+1}^N |g_t - \sum_{i=0}^{T-1} \nabla \tau_{t-i}| + \lambda_1 \sum_{t=2}^N |\nabla \tau_t| + \lambda_2 \sum_{t=3}^N |\nabla^2 \tau_t|$$

To capture abrupt change:  $\nabla \tau_t = \tau_t - \tau_{t-1}$

To capture slow change:  $\nabla^2 \tau_t = \tau_t - 2\tau_{t-1} + \tau_{t-2}$

- Finally, recover trend by *cumulative sum*

$$\tilde{\tau}_t^r = \tilde{\tau}_t - \tau_1 = \tilde{\tau}_t - \tilde{\tau}_1 = \begin{cases} 0, & t = 1 \\ \sum_{i=2}^t \nabla \tilde{\tau}_i, & t \geq 2 \end{cases}$$



# Seasonality Extraction

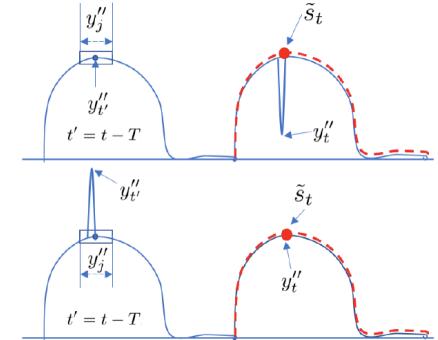
- Weighted non-local seasonal filtering
  - Consider a window of datapoints
  - Filter weights based on both value and position difference, similar to Bilateral filtering
  - Robust to *outlier/abrupt chang* and adaptive to *seasonal shift*

$$\tilde{s}_t = \frac{1}{z} \sum_{i=1}^m \alpha_i \sum_{(t'_i, j) \in \Omega} w_{(t'_i, j)}^t y''_j$$

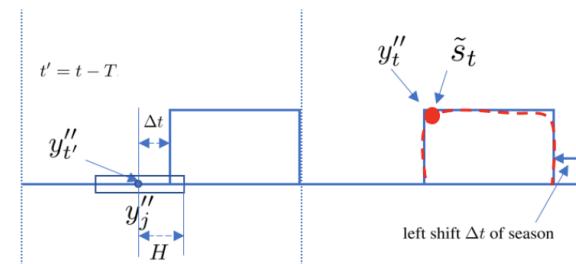
$$w_{(t'_i, j)}^t = e^{-\frac{|j-t'_i|^2}{2\delta_d^2} - \frac{|y''_j - y''_{t'_i}|^2}{2\delta_i^2}}$$

$$\Omega = \{(t'_i, j) | (t'_i = t \pm k \times T_i, j = t'_i \pm h)\}$$

$$k = 1, 2, \dots, K; h = 0, 1, \dots, H$$



(a) Outlier robustness



(b) Season shift adaptation

# Further Multi-Season Decomposition

- Sparse model: MSE loss with three L1 regularizations
  - *MSE loss*: Outliers have been removed by non-local seasonal filter
  - Regularizations:
    - First two regs: capture abrupt and slow season changes (similar to trend extraction)
    - Last *seasonal-wise second-order difference* reg: for seasonal-wise stability

$$\begin{aligned} \arg \min_{\{\mathbf{s}_i | i=1,2,\dots,m\}} \quad & \frac{1}{2} \|\tilde{\mathbf{s}} - \sum_{i=1}^m \mathbf{s}_i\|_2^2 + \sum_{i=1}^m \lambda_{1,i} \|\mathbf{D}\mathbf{s}_i\|_1 + \sum_{i=1}^m \lambda_{2,i} \|\mathbf{D}^2\mathbf{s}_i\|_1 \\ & + \sum_{i=1}^m \lambda_{3,i} \|\mathbf{D}_{T_i}^2 \mathbf{s}_i\|_1, \end{aligned}$$



# GADMM for Efficient Computation: $O(N \log N)$

- ADMM → GADMM for trend extraction

ADMM

$$\mathbf{x}_{t+1} = \arg \min_{\mathbf{x}} L_\rho(\mathbf{x}, \mathbf{y}, \mathbf{u}) = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T (\mathbf{b} + \mathbf{y}_t - \frac{1}{\rho} \mathbf{u}_t) \quad O(N^2)$$

$$\mathbf{y}_{t+1} = \arg \min_{\mathbf{y}} L_\rho(\mathbf{x}, \mathbf{y}, \mathbf{u}) = S_{1/\rho}(\mathbf{A}\mathbf{x}_{t+1} - \mathbf{b} + \frac{1}{\rho} \mathbf{u}_t) \quad O(N)$$

$$\mathbf{u}_{t+1} = \arg \min_{\mathbf{u}} L_\rho(\mathbf{x}, \mathbf{y}, \mathbf{u}) = \mathbf{u}_t + \rho(\mathbf{A}\mathbf{x}_{t+1} - \mathbf{y}_{t+1} - \mathbf{b}) \quad O(N)$$

$$\mathbf{A} = \begin{bmatrix} \mathbf{M}_{(N-T) \times (N-1)} \\ \lambda_1 \mathbf{I}_{(N-1) \times (N-1)} \\ \lambda_2 \mathbf{D}_{(N-2) \times (N-1)} \end{bmatrix}, \mathbf{b} = \begin{bmatrix} \mathbf{g}_{(N-T) \times 1} \\ \mathbf{0}_{(2N-3) \times 1} \end{bmatrix}$$

- Similar GADMM can be formulated for multi-season decomposition

GADMM  
X-update relaxation for trend extraction

$$\mathbf{x}_{t+1} = \arg \min_{\mathbf{x}} L_\rho(\mathbf{x}, \mathbf{y}, \mathbf{u}) = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T (\mathbf{b} + \mathbf{y}_t - \frac{1}{\rho} \mathbf{u}_t) \quad O(N^2)$$

$$\mathbf{x}_{t+1} \approx \arg \min_{\mathbf{x}} \left( L_\rho(\mathbf{x}, \mathbf{y}, \mathbf{u}) + \underbrace{\|\mathbf{x} - \mathbf{x}_t\|_{\hat{\mathbf{H}}}^2}_{\text{extra proximal term}} \right) \quad O(N^2)$$

$$\hat{\mathbf{H}} = \mathbf{H} - \mathbf{A}^T \mathbf{A} = \tilde{\mathbf{A}}^T \tilde{\mathbf{A}} - \mathbf{A}^T \mathbf{A} = \hat{\mathbf{M}}^T \hat{\mathbf{M}} + \lambda_2^2 \hat{\mathbf{D}}^T \hat{\mathbf{D}}$$

$$\mathbf{x}_{t+1} \approx \mathbf{H}^{-1} \mathbf{v} = \mathbf{F}^{-1} \text{diag}^{-1}(\mathbf{F}\mathbf{h}) \mathbf{F}\mathbf{v} \\ = \text{IDFT}(\text{DFT}(\mathbf{v}) ./ \text{DFT}(\mathbf{h})) \quad O(N \log N)$$

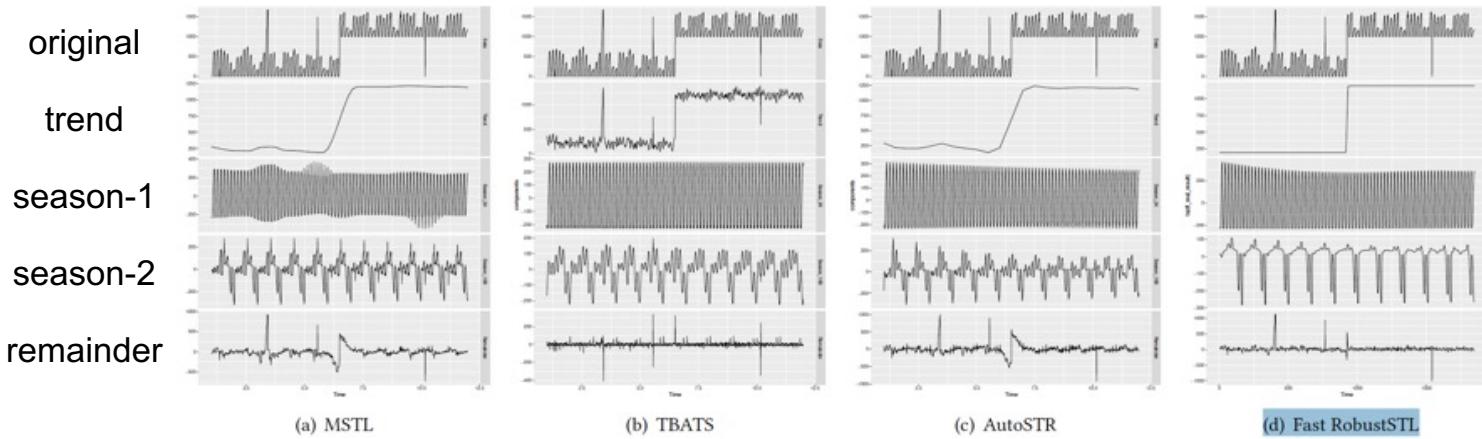
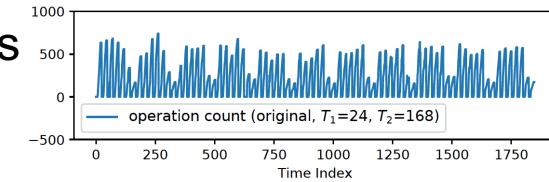
$$\hat{\mathbf{A}} = \begin{bmatrix} \hat{\mathbf{M}}_{(N-T) \times (N-1)} \\ \lambda_1 \mathbf{I}_{(N-1) \times (N-1)} \\ \lambda_2 \mathbf{D}_{(N-2) \times (N-1)} \end{bmatrix}, \hat{\mathbf{D}} = \begin{bmatrix} \mathbf{D} \\ \tilde{\mathbf{D}} \end{bmatrix} = \begin{bmatrix} 1 & -1 & & & \\ & 1 & -1 & & \\ & & \ddots & & \\ & & & 1 & -1 \\ -1 & & & & 1 \end{bmatrix} \mathbf{D}$$

construct circulant matrix

$$\hat{\mathbf{M}} = \begin{bmatrix} \mathbf{M} \\ \hat{\mathbf{M}} \end{bmatrix} = \begin{bmatrix} 1 & \cdots & 1 \\ & 1 & \cdots & 1 \\ & & \ddots & & \\ & & & 1 & \cdots & 1 \\ 1 & & & & 1 & \cdots \\ & \ddots & & & & \ddots \\ \cdots & & 1 & & & 1 \end{bmatrix} \mathbf{M}$$

# Decomposition on Real-World Data

- Real-world data with daily and weekly periodic components
  - adding abrupt trend change
  - adding 2 single-point outliers, 1 pattern outlier spanning 1 day



RobustSTL is robust to abrupt trend changes and outliers

# Comparisons

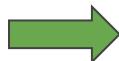
- Speed experiments

N	Metrics	RobustSTL-ADMM	RobustSTL-PDHG	<b>Fast RobustSTL</b>
N=1080	Iter	41	310	37
	Time(Sec)	0.142	0.0319	0.0109
	SpeedUp		4.45x	<b>13.0x</b>
N=2160	Iter	48	319	40
	Time(Sec)	1.11	0.0571	0.0295
	SpeedUp		19.4x	<b>37.6x</b>
N=4320	Iter	68	602	37
	Time(Sec)	5.98	0.191	0.0988
	SpeedUp		31.3x	<b>60.5x</b>
N=8640	Iter	92	1377	53
	Time(Sec)	36.7	0.62	0.254
	SpeedUp		59.1x	<b>144x</b>

Algorithm	Time (N=4320)
AutoSTR	4441 seconds
TBATS	60 seconds
(M)STL	0.2 second
<b>RobustSTL</b>	<b>0.1 second</b>

- Algorithm comparisons

Algorithm	Robust to outlier	Robust to seasonal shift	Robust to trend changes	Support multiple periods	Efficient computation
ARIMA/SEATS	✗	✗	✗	✗	✗
SSA	✗	✗	✗	✗	✗
TBATS	✗	✗	✓	✓	✗
STL	✗	✗	✗	✗	✓
MSTL	✗	✗	✗	✓	✗
STR	✓	✓	✗	✓	✗
<b>RobustSTL</b>	✓	✓	✓	✓	✓

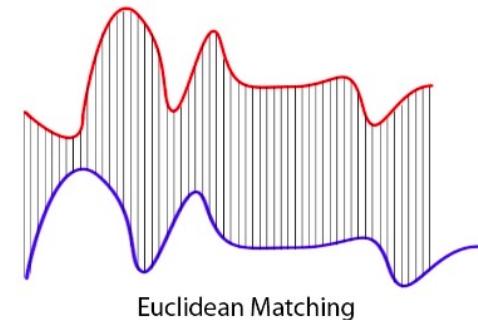


RobustSTL is more efficient than others

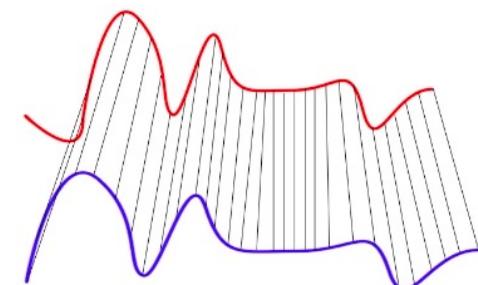


# Similarity Measure between Time Series

- Similarity measure is a core component in many tasks involving time series, e.g., time series classification, clustering and retrieval
- Different types of time series similarity
  - Euclidean distance
  - Feature-based measures (e.g., Fourier coefficients)
  - Model-based measures (e.g., auto-regressive)
  - Elastic measures (e.g., dynamic time warping and edit distance): introduce warping/aligning in the temporal domain
- Superior of DTW
  - A classical approach still among the top choices for time series similarity measure
  - Empirical studies on 38 datasets for 9 similarity measures in time series classification confirm the consistent superior of DTW
  - 1-NN with DTW beats most advanced time series classifiers



Euclidean Matching



Dynamic Time Warping Matching



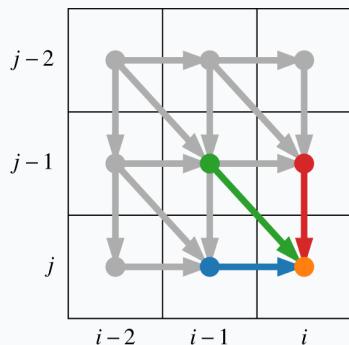
# Dynamic Time Warping (DTW)

- Dynamic Time Warping in the optimization form:

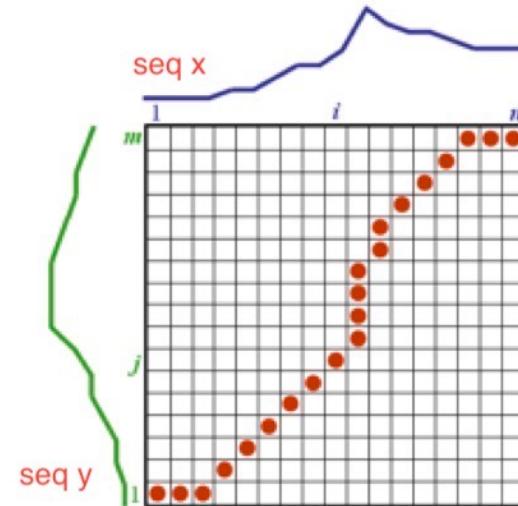
$$D(x, y) = \min \sum_t^T dist(x_t, y_{\phi(t)})$$

- Time warp function  $\phi(t)$  satisfies
  - Monotonicity
  - Continuity
  - Boundary

- How to compute DTW using dynamic programming?



$D(i, j) = Dist(i, j) + \min[D(i - 1, j), D(i, j - 1), D(i - 1, j - 1)]$   
where  $D(i, j)$  is the DTW distance between  $x_{1:i}$  and  $y_{1:j}$   
 $dist(i, j)$  is the distance between  $x_i$  and  $y_j$

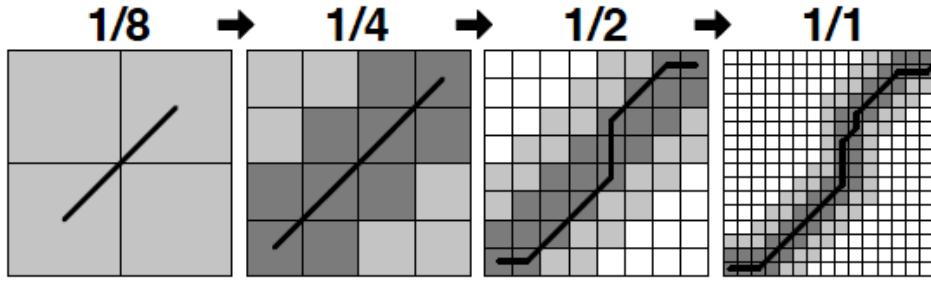


How time warp function changes  
the similarity computation



# FastDTW

- The time complexity of standard DTW is  $O(n^2)$
- How to reduce the computational cost significantly?
- FastDTW: linear time and space complexity  $O(nr)$
- Key ideas of FastDTW
  - It estimates the warp function in a multi-resolution framework
  - The warp function estimated at the current resolution is searched in a constrained space specified by the warping estimate from the low-resolution



It repeatedly uses low-resolution warping estimate as the constraint to generate the warping estimate in the current resolution



# RobustDTW: Key Ideas

- Challenges of DTW:
  - Noises and outliers bias the similarity
  - Block of missing values occurs in many real-world applications
  - How to perform DTW efficiently?
- Key ideas of RobustDTW:
  - It not only estimates the time warp function, also estimates the trend to handle noises and outliers
  - Trend estimation is achieved via graph detrending
  - A multi-resolution alternating framework is applied
- Advantages of RobustDTW:
  - Robust to noises and outliers
  - Efficient implementation thanks to the multi-resolution framework
  - Can handle (block) missing values effectively

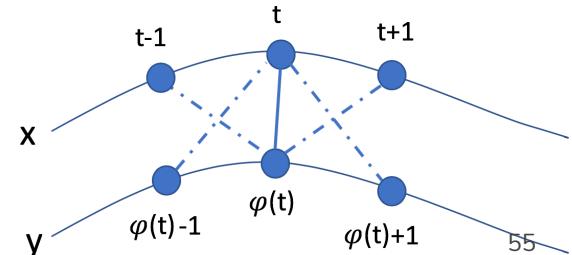
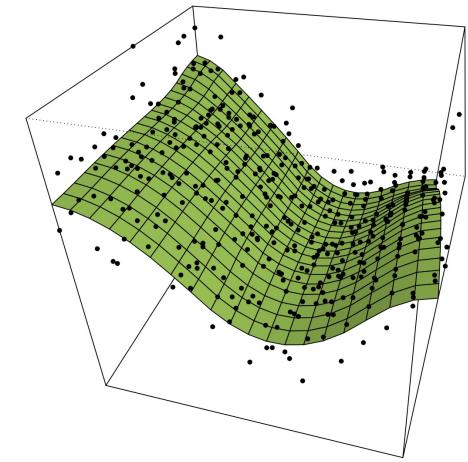
# Graph Detrending for Trend Estimation

- Graph smoothing: given a graph  $G = (V, E)$  with vertices denoted  $V = \{1, \dots, n\}$ , we assume  $y_i = \mu_i + \epsilon_i$ ,  $i = 1, \dots, n$ , where  $\epsilon_i$  assumed to have zero mean. Our goal is to estimate  $\mu_i$ , assumed to be smooth w.r.t. edges  $E$
- Graph trend filtering solves

$$\min_{\beta} \|y - \beta\|_2^2 + \lambda \|\Delta^{(k+1)} \beta\|_1$$

$\Delta^{(k+1)}$  is a graph difference operator of order  $k + 1$  over  $G$

- How to construct graph in DTW:
  - Each time point in time series corresponds to a vertex
  - In each time series, each time is connected to its previous and next time points
  - Cross time series connections are determined by the time warp function and extensions



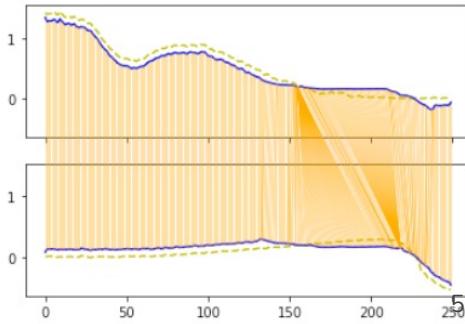
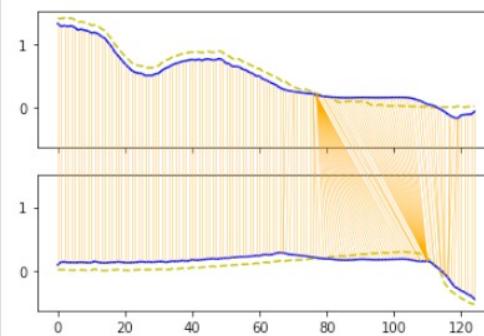
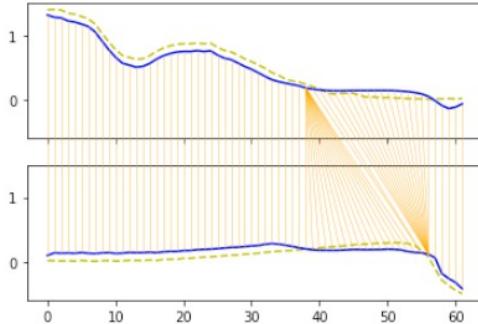


# Detailed Procedure of RobustDTW

- Step 1. Preprocess: detrend both ts separately
- Step 2. Learn the time warp function estimation
- Step 3. Construct a graph based on the learned time warp function, and perform graph detrending
- Step 4. Repeat Steps 2~3 until convergence

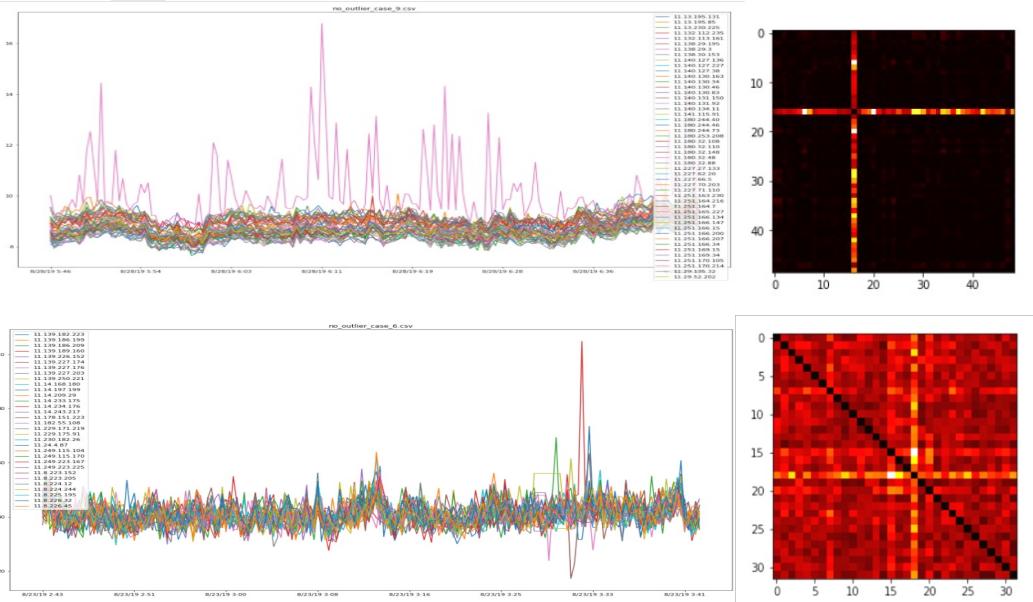


- Put the alternating procedure in a multi-resolution framework
- The time warp function is estimated based on that in the lower-resolution





# Application of RobustDTW: Time Series Outlier Detection



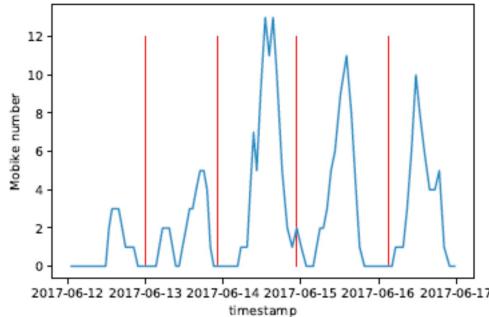
Comparison of AUC scores of outlier detection via local outlier factor (LOF) algorithm with different similarity measures and noise conditions

Dataset	Noise level	ED	DTW	FastDTW	<b>RobustDTW</b>
RT	raw data	0.887	0.986	0.949	<b>0.997</b>
	+ dips	0.811	0.965	0.918	<b>0.996</b>
	+ spikes	0.362	0.939	0.775	<b>0.972</b>
	+ spikes & dips	0.317	0.580	0.655	<b>0.969</b>
NetSpd	raw data	0.674	0.982	0.917	<b>0.988</b>
	+ dips	0.635	0.849	0.649	<b>0.982</b>
	+ spikes	0.642	0.915	0.874	<b>0.975</b>
	+ spikes & dips	0.508	0.627	0.616	<b>0.938</b>

- RobustDTW can effectively identify the outlier machine (represented by a time series) with noises and outliers in real-world

# Application of RobustDTW: Periodicity Detection

- RobustDTW can be applied in periodicity detection when the periodic length  $L$  is known
- Slice the input time series into segments with length  $L$ 
  - If the pairwise similarity with adjacent segments computed by RobustDTW is above a threshold, the input time series is predicted as periodic



Periodicity detection on 200 service monitor time series from AlibabaCloud, where 50% of ts are daily periodic

Methods	Precision	Recall	F1
ACF	0.960	0.701	0.810
AUTOPERIOD	<b>0.980</b>	0.715	0.827
RobustPeriod	0.920	0.902	0.911
Slicing	ED	0.919	0.800
	DTW	0.911	0.930
	FastDTW	0.873	0.970
	<b>RobustDTW</b>	0.951	<b>0.980</b>



# Outline

- ❑ **Introduction**
- ❑ **Preliminaries**
- ❑ **Robust Time Series Processing Blocks**

(10 minutes break)

## ➤ ***Robust Time Series Applications and Practices***

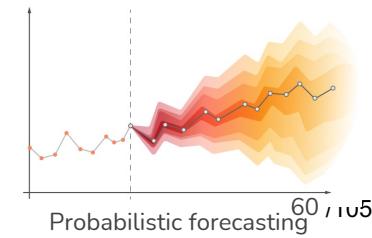
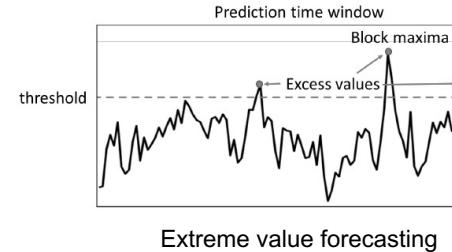
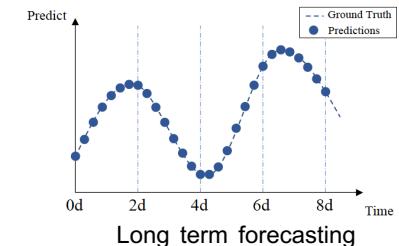
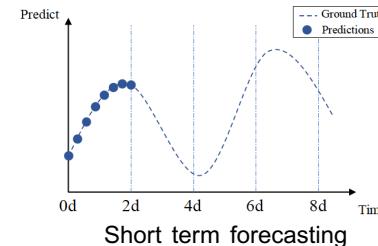
- **Forecasting:** Tree Models, Deep Ensemble, Transformers, etc.
- **Autoscaling (from Forecasting to Decision-Making):** Query Modeling, Scaling Decision, etc.
- **Anomaly Detection:** Decomposition Model, Deep State Space Model, Transformers, etc.
- **Fault Cause Localization (from Anomaly Detection to Localization):** Rule Set Learning, RCA, etc.
- **XAI for Time Series Data:** Add details Later

(Q&A Session)



# Forecasting: Background

- Different forecasting types
  - **Short-term** forecasting: predict the near future
  - **Long-term** forecasting: predict the future with an extended period
  - **Extreme value** forecasting: predict the extreme values
  - **Point or Probabilistic** forecasting: predict point value or interval/probability distribution
- Challenges:
  - Accuracy, robustness
- Models:
  - Traditional: Statistical (ARIMA, ETS, Prophet)
  - Ensemble: Tree, MLP
  - Deep Models: CNN, RNN, Transformers





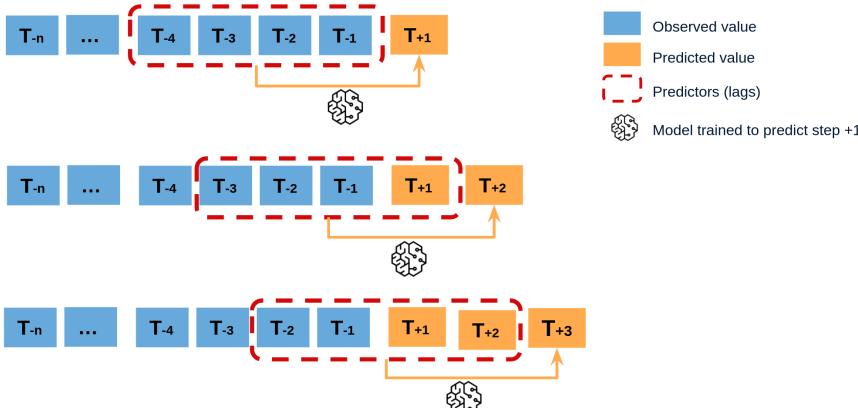
# Forecasting: Ensemble (Tree Models)

pros: robustness, effective, easy to deploy

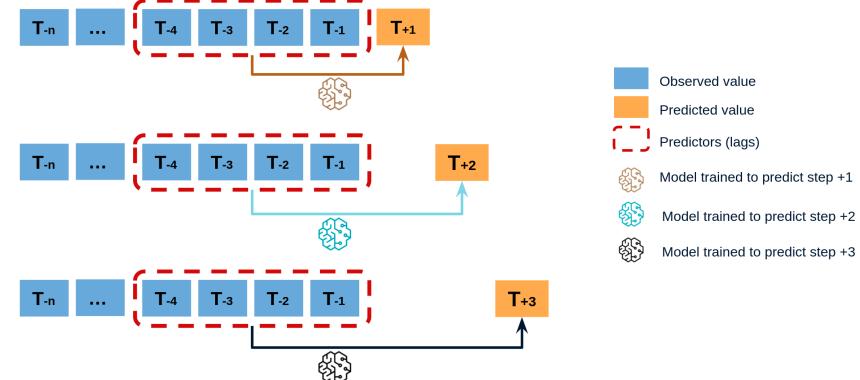
Cons: highly depend on feature engineering, global view

- Tree models: an ensemble of weak prediction models
  - XGBoost, LightGBM, CatBoost

Recursive forecasting with single tree based model



Direct forecasting with multiple tree based model for different steps

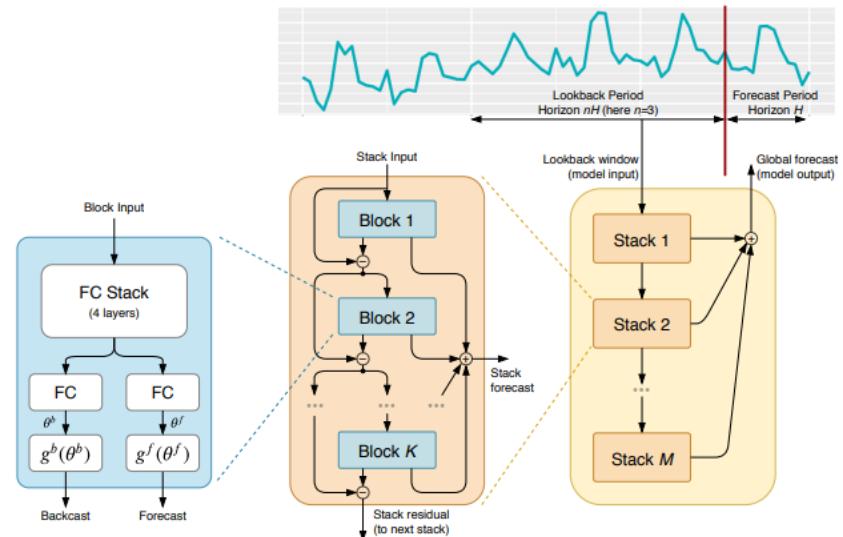




# Forecasting: Deep Ensemble (MLP based Models)

## ● N-BEATS

- Doubly residual stackings with forward and backward residual links
- Forecasts are aggregated in a hierarchical way
- Trend and seasonal models for interpretability
- Ensemble: e.g., 18 to 180 models
  - Fit on different metrics: sMAPE, MASE, MAPE
  - Train on input windows of different lengths
  - Train with different random initializations



# Forecasting: Dlinear

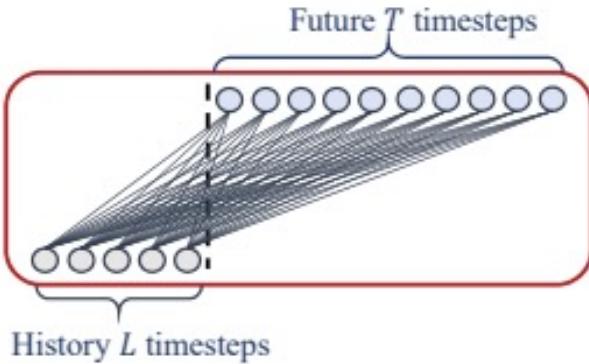


Figure 2. Illustration of the basic linear model.

**Key insight:**  
STL decomposition  
Channel independent

DLinear is a combination of a Decomposition scheme used in Autoformer and FEDformer with linear layers: 1. Decomposes a raw data input into a trend component by a moving average kernel and a remainder (seasonal) component. 2. Two one-layer linear layers are applied to each component, and we sum up the two features to get the final prediction.

NLinear first subtracts the input by the last value of the sequence. Then, the input goes through a linear layer, and the subtracted part is added back before making the final prediction.

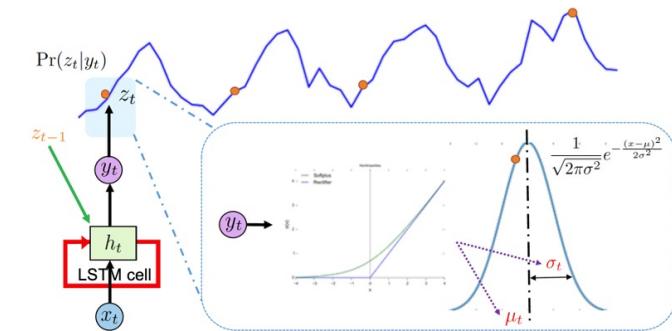
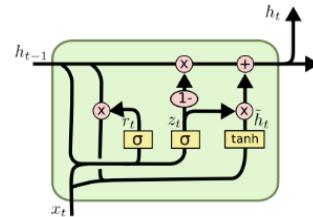
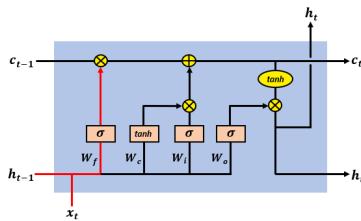
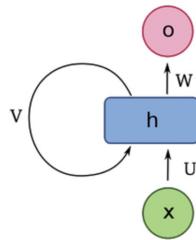
Methods	IMP	Linear*		NLinear*		DLinear*		FEDformer	Autoformer	Informer	Pyraformer*	LogTrans	Repeat*							
		MSE	MAE	MSE	MAE	MSE	MAE													
Electricity	96	27.40%	<b>0.140</b>	<b>0.237</b>	0.141	<b>0.237</b>	<b>0.237</b>	<b>0.193</b>	<b>0.308</b>	0.201	0.317	0.274	0.368	0.386	0.449	0.258	0.357	1.588	0.946	
	192	23.88%	<b>0.153</b>	0.250	0.154	<b>0.248</b>	<b>0.153</b>	0.249	<b>0.201</b>	<b>0.315</b>	0.222	0.334	0.296	0.386	0.386	0.443	0.266	0.368	1.595	0.950
	336	21.02%	<b>0.169</b>	0.268	0.171	<b>0.265</b>	<b>0.169</b>	0.267	<b>0.214</b>	<b>0.329</b>	0.231	0.338	0.300	0.394	0.378	0.443	0.280	0.380	1.617	0.961
	720	17.47%	<b>0.203</b>	0.301	0.210	<b>0.297</b>	<b>0.203</b>	0.301	<b>0.246</b>	<b>0.355</b>	0.254	0.361	0.373	0.439	0.376	0.445	0.283	0.376	1.647	0.975
Exchange	96	45.27%	0.082	0.207	0.089	0.208	<b>0.081</b>	0.203	<b>0.148</b>	<b>0.278</b>	0.197	0.323	0.847	0.752	0.376	1.105	0.968	0.812	<b>0.081</b>	<b>0.196</b>
	192	42.06%	0.167	0.304	0.180	0.300	<b>0.157</b>	0.293	<b>0.271</b>	<b>0.380</b>	0.300	0.369	1.204	0.895	1.748	1.151	1.040	0.851	0.167	<b>0.289</b>
	336	33.69%	0.328	0.432	0.331	0.415	<b>0.305</b>	0.414	<b>0.460</b>	<b>0.500</b>	0.509	0.524	1.672	1.036	1.874	1.172	1.659	1.081	<b>0.305</b>	<b>0.396</b>
	720	46.19%	0.964	0.750	1.033	0.780	<b>0.643</b>	<b>0.601</b>	<b>1.195</b>	<b>0.841</b>	1.447	0.941	2.478	1.310	1.943	1.206	1.941	1.127	0.823	0.681
Traffic	96	30.15%	<b>0.410</b>	0.282	<b>0.410</b>	<b>0.279</b>	<b>0.410</b>	0.282	<b>0.587</b>	<b>0.366</b>	0.613	0.388	0.719	0.391	2.085	0.468	0.684	0.384	2.723	0.179
	192	29.96%	<b>0.423</b>	0.287	<b>0.423</b>	<b>0.284</b>	<b>0.423</b>	0.287	<b>0.604</b>	<b>0.373</b>	0.616	0.382	0.696	0.379	0.867	0.467	0.685	0.390	2.756	1.087
	336	29.95%	0.436	0.295	<b>0.435</b>	<b>0.290</b>	0.436	0.296	<b>0.621</b>	0.383	0.622	<b>0.337</b>	0.777	0.420	0.869	0.469	0.734	0.408	2.791	1.095
	720	25.87%	0.466	0.315	<b>0.464</b>	<b>0.307</b>	0.466	0.315	<b>0.626</b>	<b>0.382</b>	0.660	0.408	0.864	0.472	0.881	0.473	0.717	0.396	2.811	1.097
Weather	96	18.89%	<b>0.176</b>	0.236	0.182	<b>0.232</b>	<b>0.176</b>	0.237	<b>0.217</b>	<b>0.296</b>	0.266	0.336	0.300	0.384	0.896	0.556	0.458	0.490	0.259	0.254
	192	21.01%	<b>0.218</b>	0.276	0.225	<b>0.269</b>	0.220	0.282	<b>0.276</b>	<b>0.336</b>	0.307	0.367	0.598	0.544	0.622	0.624	0.658	0.589	0.309	0.292
	336	22.71%	<b>0.262</b>	0.312	0.271	<b>0.301</b>	0.265	0.319	<b>0.339</b>	<b>0.380</b>	0.359	0.395	0.578	0.523	0.739	0.753	0.797	0.652	0.377	0.338
	720	19.85%	0.326	0.365	0.338	<b>0.348</b>	<b>0.323</b>	0.362	<b>0.403</b>	<b>0.428</b>	0.419	0.428	1.059	0.741	1.004	0.934	0.869	0.675	0.465	0.394



# Forecasting: RNN based Models

- Recurrent neural networks

- From RNN to LSTM/GRU: control information flow by gates, mitigating vanishing gradient problem
- DeepAR: time series probabilistic forecasting through autoregressive recurrent network



Hochreiter, Sepp, and Jürgen Schmidhuber. "Long short-term memory." Neural computation, 1997.

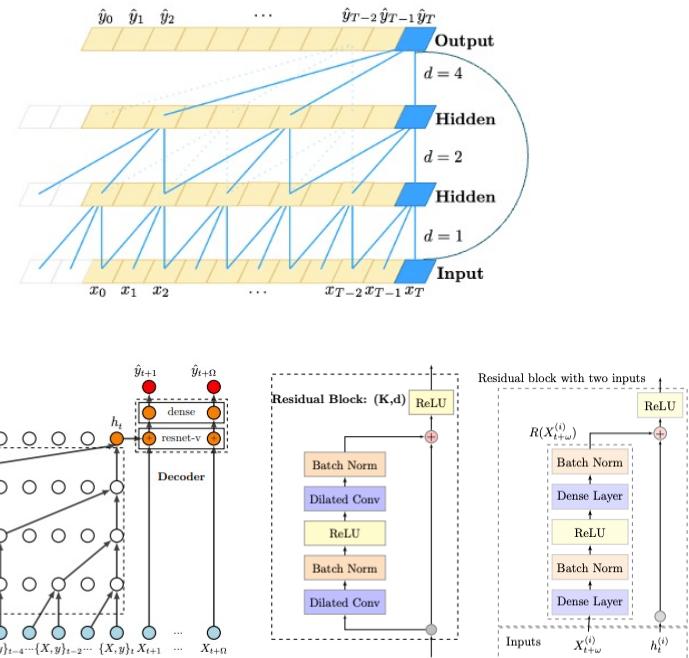
Chung, Junyoung, et al. "Empirical evaluation of gated recurrent neural networks on sequence modeling." arXiv preprint arXiv:1412.3555, 2014.

Salinas, D., Flunkert, V., Gasthaus, J., & Januschowski, T. DeepAR: Probabilistic forecasting with autoregressive recurrent networks. International Journal of Forecasting, 2020.



# Forecasting: CNN based Models

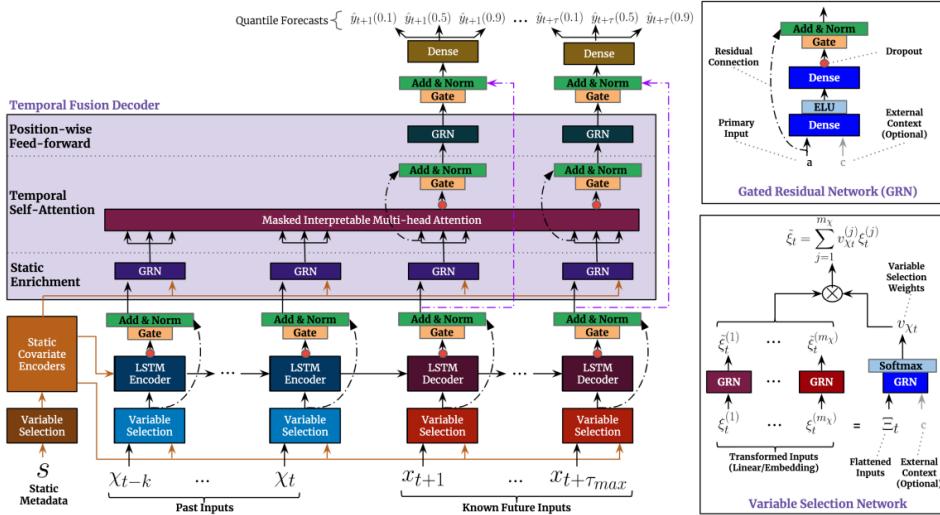
- TCN (temporal convolutional networks)
  - *1D dilated causal CNN*
  - Dilated convolutional layers for large receptive field
  - Skip connection for raw info flow like ResNet
  - Empirically good performance and fast training
- DeepTCN
  - Encoder and decoder with TCN
  - Probabilistic forecasting under both parametric (Gaussian) and non-parametric (quantile) settings





# Forecasting with Transformer: TFT

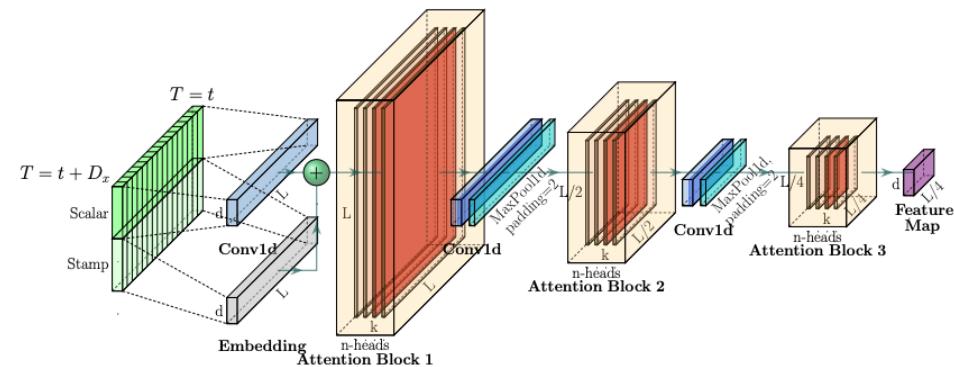
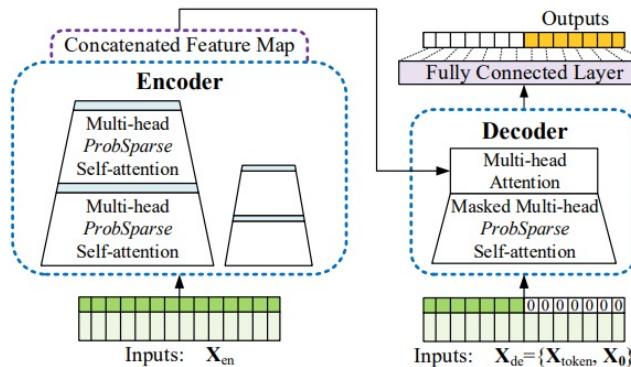
- TFT (temporal fusion Transformer)
  - Variable selection network: select most salient features
  - Gated residual network: efficient information flow with skip connections and gating layers
  - *Multi-scales network*: recurrent layers for local processing, interpretable self-attention layers for long-term dependencies



# Forecasting with Transformer: Informer

- **Informer**

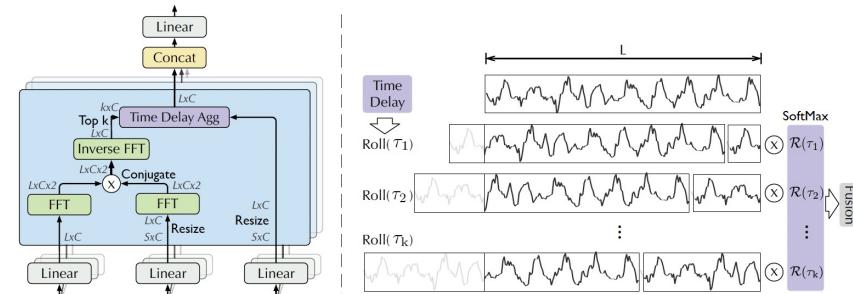
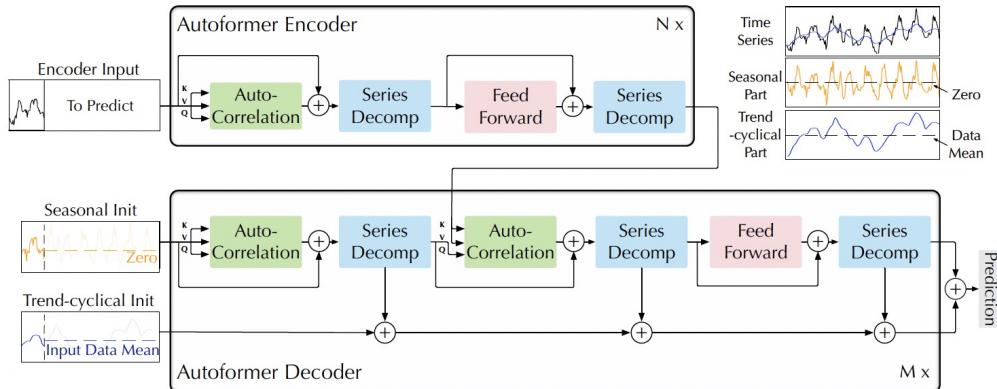
- ProbSparse self-attention for efficient and robust attention mechanism
- Self-attention distilling: extract dominating attention and reducing the network size
- *Generative style decoder*: produce long sequence forecasts with only one forward step, avoiding cumulative error spreading during inference



# Forecasting with Transformer: Autoformer

- Autoformer: Transformer with auto-correlation mechanism
  - *Decomposition* architecture to disentangle complex temporal patterns (seasonality, trend)
  - *Auto-correlation* instead of point-wise self-attention to utilize period-based dependencies and reduce complexity

Key insight: similarity measure using subsequence similarity; STL decomposition

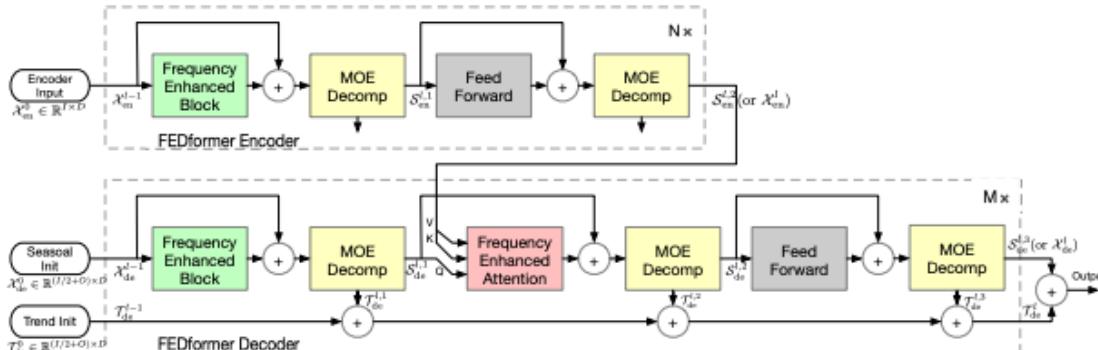


auto-correlation mechanism



# Forecasting with Transformer: FEDformer

- FEDformer: frequency enhanced decomposed Transformer
  - Efficient and robust *frequency domain processing*: to capture important structures in time series
    - Frequency enhanced block: substitute self-attention
    - Frequency enhanced attention: substitute cross-attention
  - Mixture of experts *seasonal-trend decomposition*: to better capture global properties in time series



Key insight: modeling in Frequency domain; moe stl decomposition

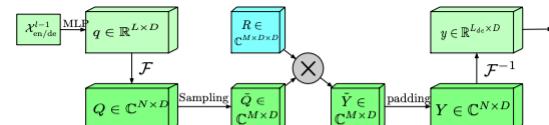


Figure 3. Frequency Enhanced Block with Fourier transform (FEB-f) structure.

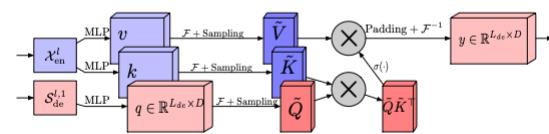


Figure 4. Frequency Enhanced Attention with Fourier transform (FEA-f) structure, \$\sigma(\cdot)\$ is the activation function.



# Forecasting with Transformer: FEDformer

## Empirical comparison of FEDformer on six benchmark datasets

Table 2. Multivariate long-term series forecasting results on six datasets with input length  $I = 96$  and prediction length  $O \in \{96, 192, 336, 720\}$  (For ILI dataset, we use input length  $I = 36$  and prediction length  $O \in \{24, 36, 48, 60\}$ ). A lower MSE indicates better performance, and the best results are highlighted in bold.

Methods	Metric	ETTm2				Electricity				Exchange				Traffic				Weather				ILI			
		96	192	336	720	96	192	336	720	96	192	336	720	96	192	336	720	96	192	336	720	24	36	48	60
FEDformer-f	MSE	<b>0.203</b>	<b>0.269</b>	<b>0.325</b>	<b>0.421</b>	0.193	0.201	0.214	0.246	0.148	0.271	0.460	1.195	0.587	0.604	0.621	0.626	<b>0.217</b>	<b>0.276</b>	<b>0.339</b>	<b>0.403</b>	3.228	2.679	2.622	2.857
	MAE	<b>0.287</b>	<b>0.328</b>	<b>0.366</b>	<b>0.415</b>	0.308	0.315	0.329	0.355	0.278	0.380	0.500	0.841	0.366	0.373	0.383	0.382	<b>0.296</b>	<b>0.336</b>	<b>0.380</b>	<b>0.428</b>	1.260	1.080	1.078	1.157
FEDformer-w	MSE	0.204	0.316	0.359	0.433	<b>0.183</b>	<b>0.195</b>	<b>0.212</b>	<b>0.231</b>	<b>0.139</b>	<b>0.256</b>	<b>0.426</b>	<b>1.090</b>	<b>0.562</b>	<b>0.562</b>	<b>0.570</b>	<b>0.596</b>	0.227	0.295	0.381	0.424	<b>2.203</b>	<b>2.272</b>	<b>2.209</b>	<b>2.545</b>
	MAE	0.288	0.363	0.387	0.432	<b>0.297</b>	<b>0.308</b>	<b>0.313</b>	<b>0.343</b>	<b>0.276</b>	<b>0.369</b>	<b>0.464</b>	<b>0.800</b>	<b>0.349</b>	<b>0.346</b>	<b>0.323</b>	<b>0.368</b>	0.304	0.363	0.416	0.434	<b>0.963</b>	<b>0.976</b>	<b>0.981</b>	<b>1.061</b>
Autoformer	MSE	0.255	0.281	0.339	0.422	0.201	0.222	0.231	0.254	0.197	0.300	0.509	1.447	0.613	0.616	0.622	0.660	0.266	0.307	0.359	0.419	3.483	3.103	2.669	2.770
	MAE	0.339	0.340	0.372	0.419	0.317	0.334	0.338	0.361	0.323	0.369	0.524	0.941	0.388	0.382	0.337	0.408	0.336	0.367	0.395	0.428	1.287	1.148	1.085	1.125
Informer	MSE	0.365	0.533	1.363	3.379	0.274	0.296	0.300	0.373	0.847	1.204	1.672	2.478	0.719	0.696	0.777	0.864	0.300	0.598	0.578	1.059	5.764	4.755	4.763	5.264
	MAE	0.453	0.563	0.887	1.338	0.368	0.386	0.394	0.439	0.752	0.895	1.036	1.310	0.391	0.379	0.420	0.472	0.384	0.544	0.523	0.741	1.677	1.467	1.469	1.564
LogTrans	MSE	0.768	0.989	1.334	3.048	0.258	0.266	0.280	0.283	0.968	1.040	1.659	1.941	0.684	0.685	0.7337	0.717	0.458	0.658	0.797	0.869	4.480	4.799	4.800	5.278
	MAE	0.642	0.757	0.872	1.328	0.357	0.368	0.380	0.376	0.812	0.851	1.081	1.127	0.384	0.390	0.408	0.396	0.490	0.589	0.652	0.675	1.444	1.467	1.468	1.560
Reformer	MSE	0.658	1.078	1.549	2.631	0.312	0.348	0.350	0.340	1.065	1.188	1.357	1.510	0.732	0.733	0.742	0.755	0.689	0.752	0.639	1.130	4.400	4.783	4.832	4.882
	MAE	0.619	0.827	0.972	1.242	0.402	0.433	0.433	0.420	0.829	0.906	1.016	0.423	0.420	0.420	0.420	0.423	0.596	0.638	0.596	0.792	1.382	1.448	1.465	1.483

## Linear complexity of FEDformer

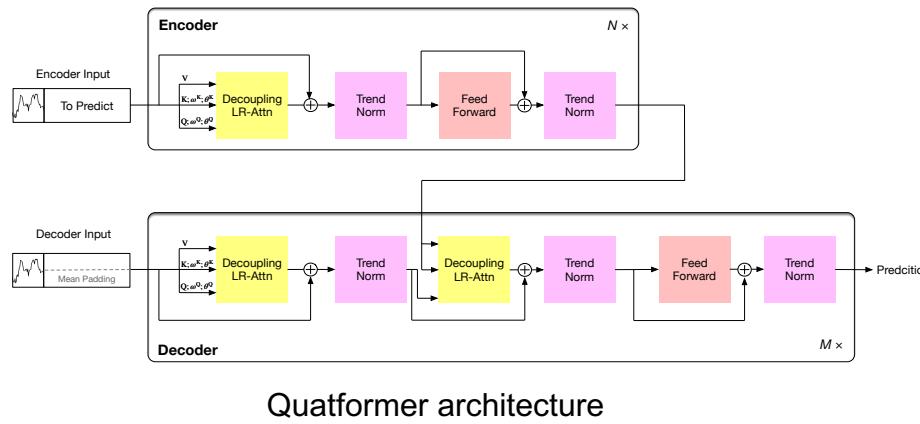
Table 1. Complexity analysis of different forecasting models.

Methods	Training		Testing
	Time	Memory	
FEDformer	$\mathcal{O}(L)$	$\mathcal{O}(L)$	1
Autoformer	$\mathcal{O}(L \log L)$	$\mathcal{O}(L \log L)$	1
Informer	$\mathcal{O}(L \log L)$	$\mathcal{O}(L \log L)$	1
Transformer	$\mathcal{O}(L^2)$	$\mathcal{O}(L^2)$	L
LogTrans	$\mathcal{O}(L \log L)$	$\mathcal{O}(L^2)$	1
Reformer	$\mathcal{O}(L \log L)$	$\mathcal{O}(L \log L)$	L
LSTM	$\mathcal{O}(L)$	$\mathcal{O}(L)$	L

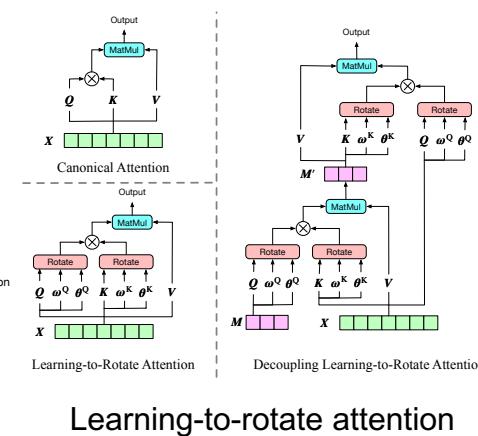


# Forecasting with Transformer: Quatformer

- Quatformer: Transformer with quaternions for periodic time series
  - Learning-to-rotate attention (LR-Attn): modeling multiple *periods and periodic changes* by quaternions
  - Trend normalization: modeling slowly varying trend



Quatformer architecture



Learning-to-rotate attention

$$\begin{aligned} & \frac{\gamma}{\sigma} \odot (\mathcal{X} - \text{MovingAvg}(\mathcal{X})) + \mathcal{T}, \\ & \sigma = \sqrt{\frac{1}{N} \sum_{i=1}^N (\mathcal{X}_i - \mu)^2}, \quad \mu = \frac{1}{N} \sum_{i=1}^N \mathcal{X}_i, \\ & \mathcal{T} = \sum_{i=0}^p \beta_i \text{pos}^i, \quad \text{pos} = [0, 1, 2, \dots, N-1]^\top / N. \end{aligned}$$

Trend normalization

Key insight: embedding in quaternion space; Trend norm

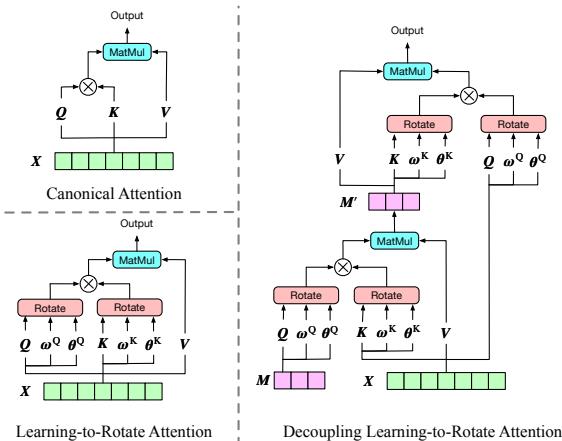


# Forecasting with Transformer: Quatformer

- Quatformer: Decoupling LR-Attn, Performance

- LR-Attn's complexity is  $O(N^2)$
- Decoupling LR-Attn introduces a momentum-updated  $c$ -length latent series  $\mathcal{M} \in \mathbb{R}^{c \times d}$ , and decouple  $\mathcal{H} = \text{LR-Attn}(\mathcal{X}, \mathcal{Y})$  into

$$\begin{aligned}\mathcal{H} &= \text{LR-Attn}(\mathcal{X}, \mathcal{M}') \in \mathbb{R}^{N \times d}, \\ \mathcal{M}' &= \text{LR-Attn}(\mathcal{M}, \mathcal{Y}) \in \mathbb{R}^{c \times d}.\end{aligned}$$



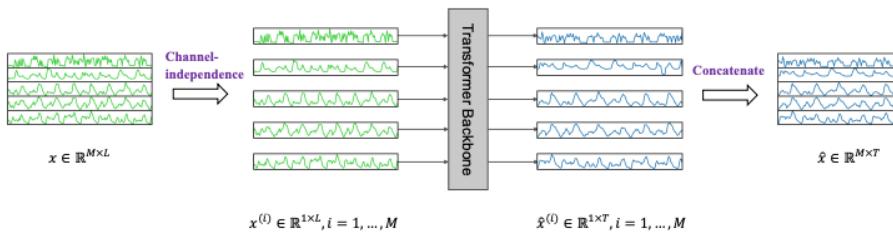
Complexity is decreased to  $O(2cN)$ .

Performance of Quatformer

Models	Quatformer		Quatformer <sup>†</sup>		Autoformer	
	MSE	MAE	MSE	MAE	MSE	MAE
ETTh	0.403	0.434	0.426	0.450	0.442	0.451
	0.444	0.463	0.453	0.466	0.500	0.482
	0.452	0.455	0.464	0.474	0.512	0.492
	0.477	0.490	0.474	0.487	0.514	0.512
ETTm	0.220	0.301	0.217	0.297	0.247	0.325
	0.279	0.333	0.269	0.329	0.278	0.335
	0.331	0.354	0.330	0.366	0.336	0.370
	0.422	0.413	0.433	0.428	0.439	0.435
Weather	0.211	0.279	0.213	0.287	0.259	0.332
	0.263	0.325	0.265	0.326	0.300	0.359
	0.310	0.344	0.315	0.354	0.364	0.401
	0.381	0.374	0.382	0.378	0.439	0.440
Exchange	0.147	0.274	0.148	0.276	0.154	0.284
	0.254	0.364	0.255	0.365	0.272	0.381
	0.427	0.481	0.425	0.480	0.461	0.509
	0.974	0.751	1.095	0.800	1.100	0.813
Traffic	0.618	0.384	0.617	0.387	0.636	0.397
	0.619	0.384	0.600	0.367	0.618	0.381
	0.622	0.384	0.618	0.385	0.626	0.388
	0.629	0.383	0.616	0.379	0.653	0.400
Electricity	0.197	0.308	0.200	0.311	0.203	0.318
	0.205	0.302	0.216	0.330	0.233	0.338
	0.220	0.329	0.228	0.343	0.259	0.359
	0.245	0.350	0.242	0.349	0.255	0.361

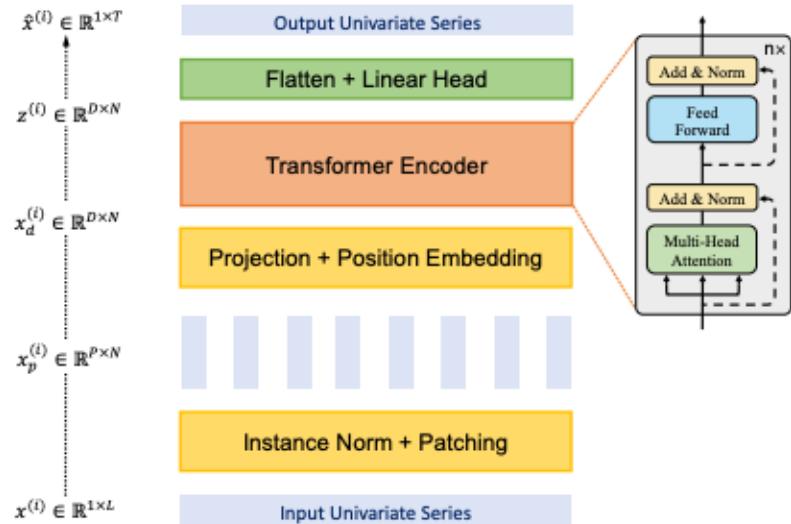


# PatchTST



Key insights: Channel independent; Patching

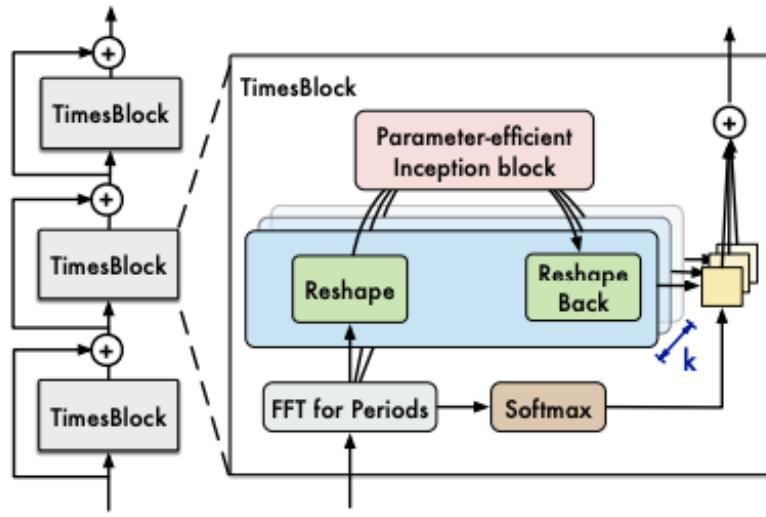
Models	PatchTST								FEDformer		
	P+CI		CI		P		Original				
Metric	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	
Weather	96	<b>0.152</b>	<b>0.199</b>	0.164	0.213	0.168	0.223	0.177	0.236	0.238	0.314
	192	<b>0.197</b>	<b>0.243</b>	0.205	0.250	0.213	0.262	0.221	0.270	0.275	0.329
	336	<b>0.249</b>	<b>0.283</b>	0.255	0.289	0.266	0.300	0.271	0.306	0.339	0.377
	720	<b>0.320</b>	<b>0.335</b>	0.327	0.343	0.351	0.359	0.340	0.353	0.389	0.409
Traffic	96	<b>0.367</b>	<b>0.251</b>	0.397	0.271	0.595	0.376	-	-	0.576	0.359
	192	<b>0.385</b>	<b>0.259</b>	0.411	0.276	0.612	0.387	-	-	0.610	0.380
	336	<b>0.398</b>	<b>0.265</b>	0.423	0.282	0.651	0.391	-	-	0.608	0.375
	720	<b>0.434</b>	<b>0.287</b>	0.457	0.309	-	-	-	-	0.621	0.375
Electricity	96	<b>0.130</b>	<b>0.222</b>	0.136	0.231	0.196	0.307	0.205	0.318	0.186	0.302
	192	<b>0.148</b>	<b>0.240</b>	0.164	0.263	0.215	0.323	-	-	0.197	0.311
	336	<b>0.167</b>	<b>0.261</b>	0.168	0.262	0.228	0.338	-	-	0.213	0.328
	720	<b>0.202</b>	<b>0.291</b>	0.219	0.312	0.244	0.345	-	-	0.233	0.344



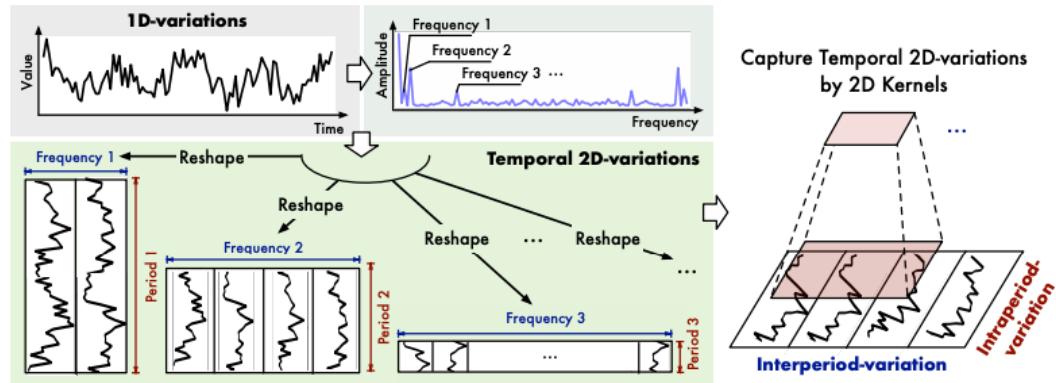
(b) Transformer Backbone (Supervised)

Models	PatchTST/64	PatchTST/42	DLinear	FEDformer		Autoformer		Informer		Pyraformer		LogTrans	
				MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
Weather	96	<b>0.149</b>	<b>0.198</b>	0.152	0.199	0.176	0.237	0.238	0.314	0.249	0.329	0.354	0.405
	192	<b>0.194</b>	<b>0.241</b>	0.197	0.243	0.220	0.282	0.275	0.329	0.325	0.370	0.419	0.434
	336	<b>0.245</b>	<b>0.282</b>	0.249	0.283	0.265	0.319	0.339	0.377	0.351	0.391	0.583	0.543
	720	<b>0.314</b>	<b>0.334</b>	0.320	0.335	0.323	0.362	0.389	0.409	0.415	0.426	0.916	0.705
Traffic	96	<b>0.360</b>	<b>0.249</b>	0.367	0.251	0.410	0.282	0.576	0.359	0.597	0.371	0.733	0.410
	192	<b>0.379</b>	<b>0.256</b>	0.385	0.259	0.423	0.287	0.610	0.380	0.607	0.382	0.777	0.435
	336	<b>0.392</b>	<b>0.264</b>	0.398	0.265	0.436	0.296	0.608	0.375	0.623	0.387	0.776	0.434
	720	<b>0.432</b>	<b>0.286</b>	0.434	0.287	0.466	0.315	0.621	0.375	0.639	0.395	0.827	0.466
Electricity	96	<b>0.129</b>	<b>0.222</b>	0.130	<b>0.222</b>	0.140	0.237	0.186	0.302	0.196	0.313	0.304	0.393
	192	<b>0.147</b>	<b>0.240</b>	0.148	<b>0.240</b>	0.153	0.249	0.197	0.311	0.211	0.324	0.327	0.417
	336	<b>0.163</b>	<b>0.259</b>	0.167	<b>0.261</b>	0.169	0.267	0.213	0.328	0.214	0.327	0.333	0.422
	720	<b>0.197</b>	<b>0.290</b>	0.202	0.291	0.203	0.301	0.233	0.344	0.236	0.342	0.351	0.427
ILI	24	<b>1.319</b>	<b>0.754</b>	1.522	0.814	2.215	1.081	2.624	1.095	2.906	1.182	4.657	1.449
	36	<b>1.579</b>	<b>0.870</b>	1.430	<b>0.834</b>	1.963	0.963	2.516	1.021	2.585	1.038	4.650	1.463
	48	<b>1.553</b>	<b>0.815</b>	1.673	<b>0.854</b>	2.130	1.024	2.505	1.041	3.024	1.145	5.004	1.542
	60	<b>1.470</b>	<b>0.788</b>	1.529	0.862	2.368	1.096	2.742	1.122	2.761	1.114	5.071	1.543

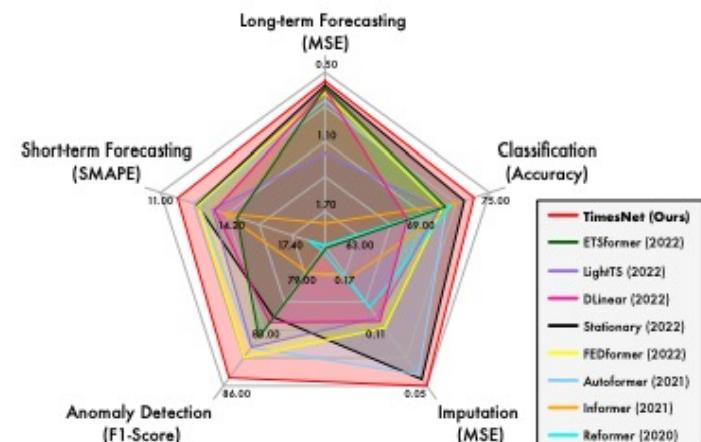
# TimesNet



Using an inception net to extract the features from 2D signal and make the final forecasting



Transform 1D-variations into 2D-variations using periodic signal extraction(FFT)



# One Fits All

“Foundation/Pretrain model” for TS

## Pretraining vs End2end

- Utilize cross modal pretrained model(LLM) for data sparsity problem(effectiveness)
- Understanding the why and how problem

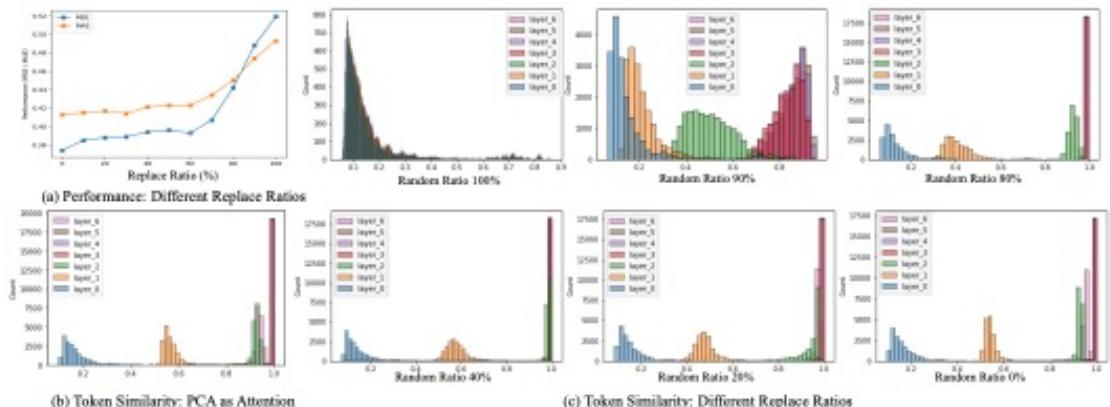
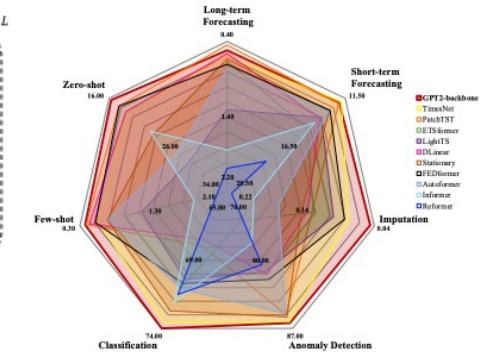
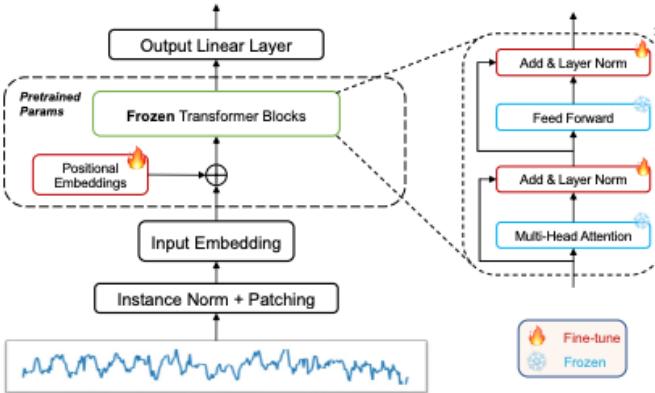
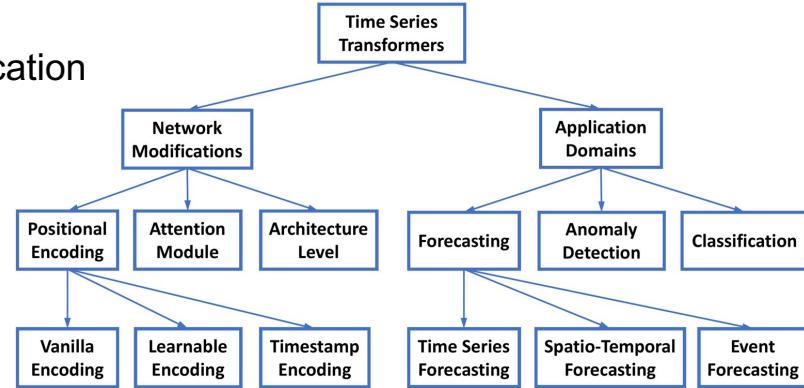


Figure 4: (a, c) The performance and token similarity within samples with respect to each layer with different random mixed ratio. Pre-trained parameters are mixed with random initial parameters according to certain proportions. (b) Token similarity within samples when replacing the attention with PCA.



# Forecasting with Transformer: Summary

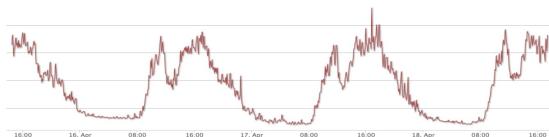
- Taxonomy of Transformers in time series
  - Network modifications
  - Applications: forecasting, anomaly detection, classification
- Evaluation and comparison
  - Robustness analysis
  - Model size analysis
  - Seasonal-trend decomposition analysis
- Future research opportunities
  - *Inductive bias* for time series Transformers (*seasonality, trend*)
  - Transformers and GNN for time series
  - Pre-trained Transformers for time series
  - Transformers with NAS/AutoML for time series



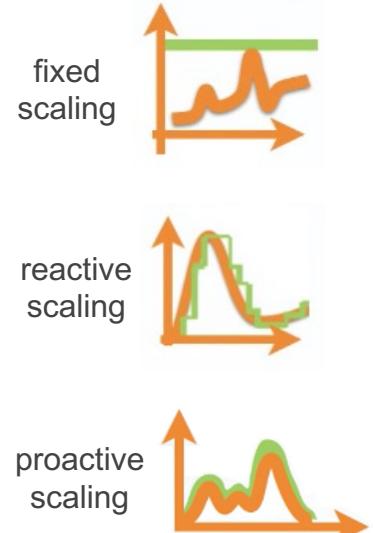
# From Forecasting to *Decision-making*: Autoscaling

- Autoscaling in cloud computing
  - Automatically add/delete resources to match the demand
- Strategies
  - Fixed, Reactive (k8s HPA), Proactive/Predictive
- Proactive/Predictive autoscaling
  - Many cloud applications with periodic pattern (often over 50%)
  - More potential in *periodic* scenario (reduce cost)

periodic scenario  
more potential for proactive autoscaling



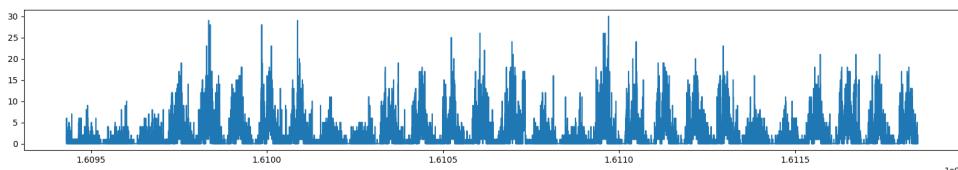
nonperiodic & random scenario  
less potential for proactive autoscaling



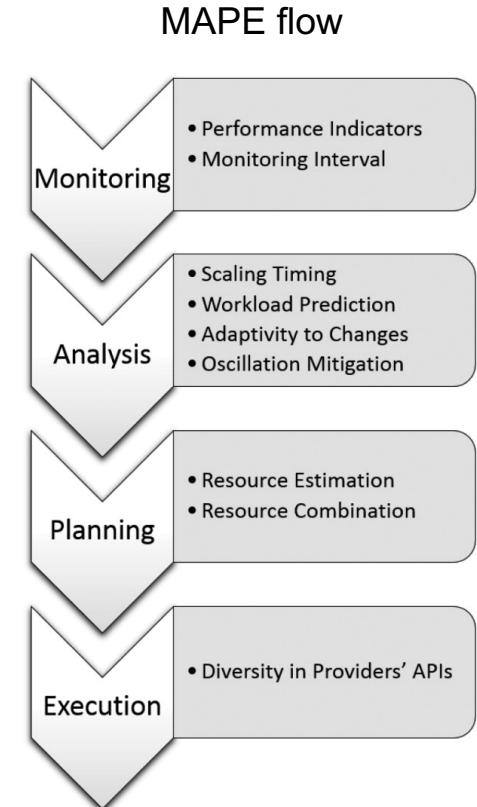


# Autoscaling

- Autoscaling procedure: continuously repeats MAPE
  - Monitoring
  - Analysis: time series forecasting/anomaly detection
  - Planning: decision with optimization
  - Execution
- Autoscaling challenges related to time series
  - Complex periodic patterns, data contamination
  - Uncertainty: query arrival time, workload amount



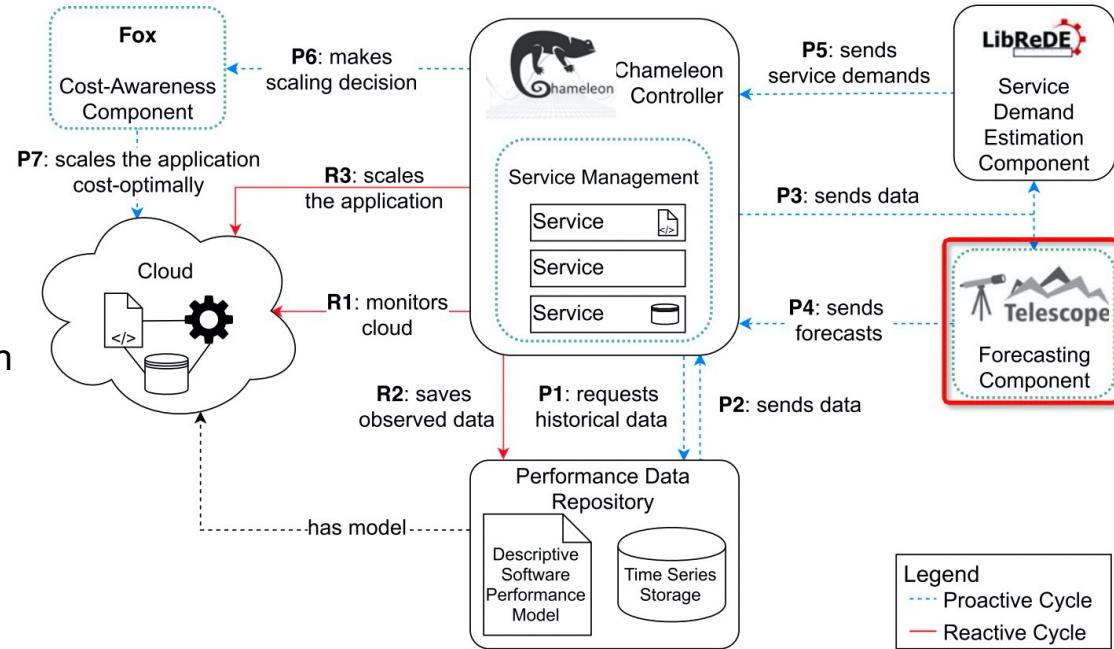
real-world QPS series from Alibaba container registry service



# Autoscaling: Chamalteon

- System:

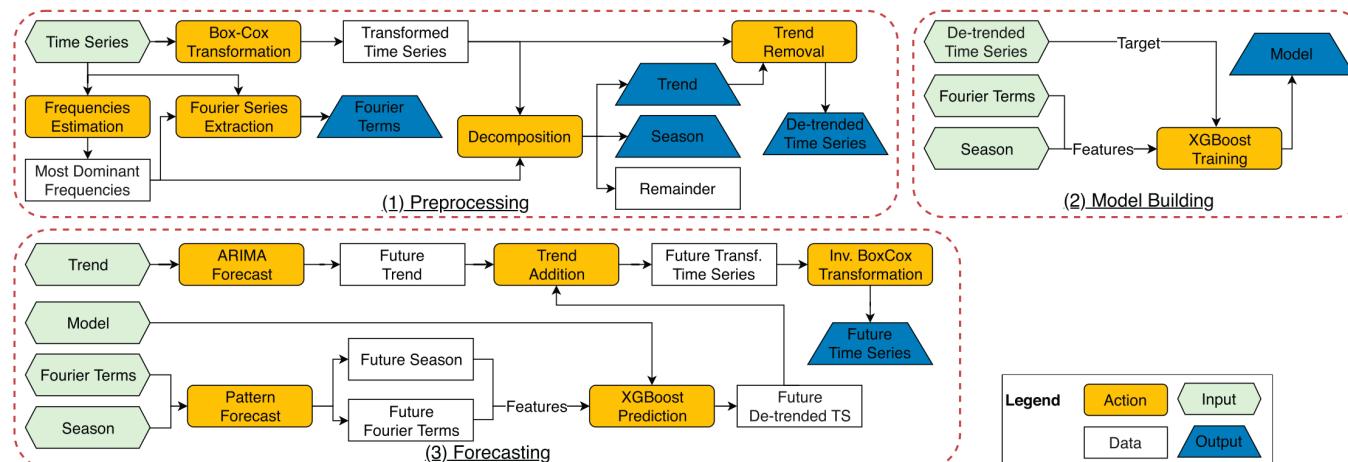
- Reactive
- Proactive
  - Monitor
  - Forecast
  - Demand estimation
  - Decision
  - Optimization





# Autoscaling: Chamalteon

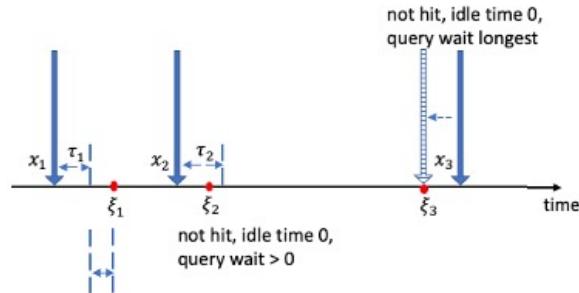
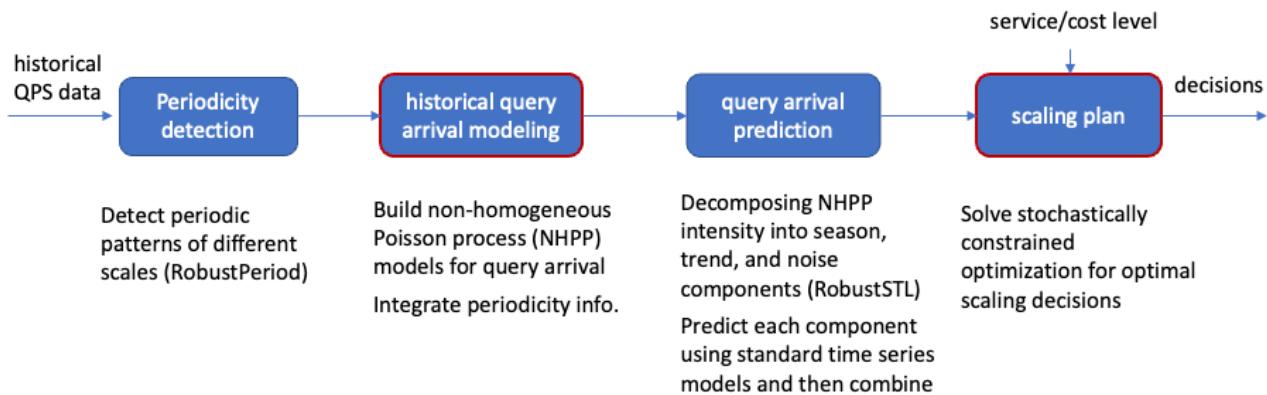
- Forecasting: Telescope
  - Designed for *seasonal* time series forecasting
    - Periodicity detection, seasonal-trend decomposition
    - Forecasting trend by ARIMA, seasonality by XGBoost





# Autoscaling: RobustScaler

- A robust and efficient proactive autoscaling scheme
  - Special non-homogeneous Poisson process for query arrival
  - Stochastically constrained optimization for uncertainty



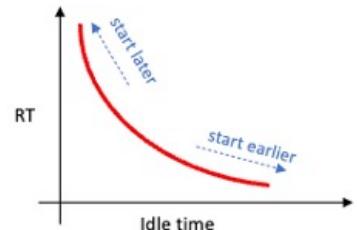
$x_1 < x_2 < \dots < x_i < \dots$  (Blue arrows): time of start of each instance  
 $\tau_1, \tau_2, \dots, \tau_i, \dots$ : startup delay of each instance  
 $\xi_1 < \xi_2 < \dots < \xi_i < \dots$  (Red dots): time of arrival of each query

#### QoS metrics:

- Response time (RT):  $\text{end of processing} - \text{time of arrival}$
- Hitting probability (HP): probability of a query being hit (i.e., warm instance available upon arrival)

#### Cost metrics:

- Idle time (of an instance):  $\text{time to start processing a query} - \text{time of finishing startup}$





# MagicScalar

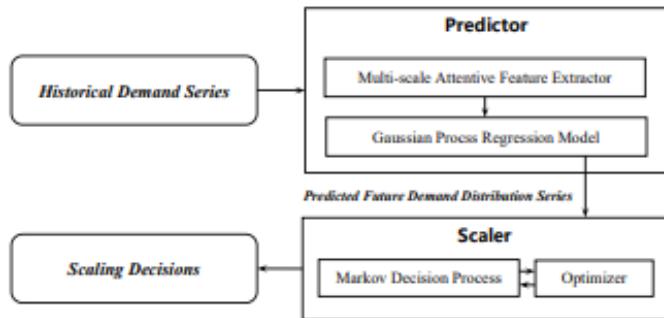


Figure 3: Framework overview of the proposed MagicScaler.

**Predictor:** The predictor first employs a multi-scale attentive feature extractor (MAFE) to capture multi-scale features from historical demand series. Then fed into a Gaussian process regression model to predict future demand distribution series.

**Scaler:** The scaler takes as input the predicted future demand distributions and models the scaling decisions into a Markov decision process. Finally, by considering the instance costs and QoS violation risks, the scaler returns the final scaling decisions, i.e., launching new instances or deleting existing instances.

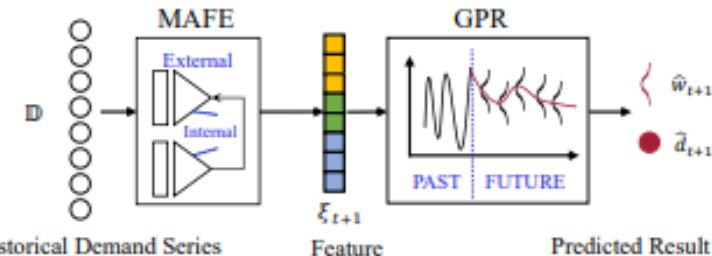


Figure 4: The overall workflow of the Predictor module.

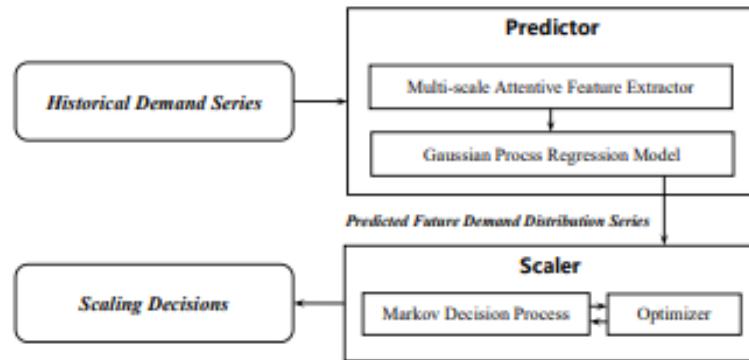


Figure 3: Framework overview of the proposed MagicScaler.



# Time Series Anomaly Detection: Background

- Time-series anomalies

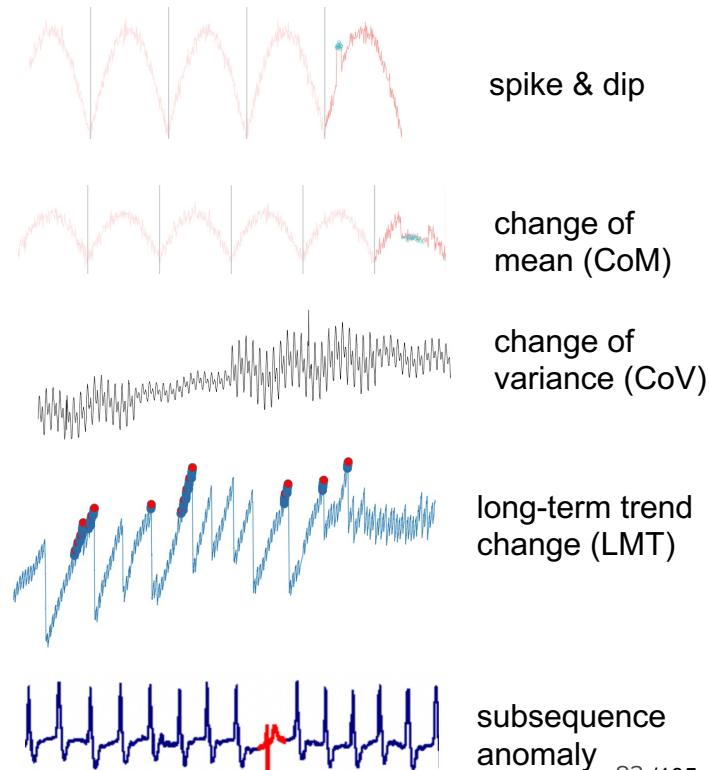
- Points/sequences that significantly deviate from the normal pattern
- **Differences** from static data anomaly detection:
  - Anomalies have temporal context
  - Noise is non-stationary (i.e., time-varying noise)
  - Lack of anomaly labels

- Types of anomalies

- Point: Spikes/Dips, CoM, CoV, LMT
- Subsequence anomaly: identifying anomalous subsequences (sequences of points)

- Models

- *Traditional methods, deep models*

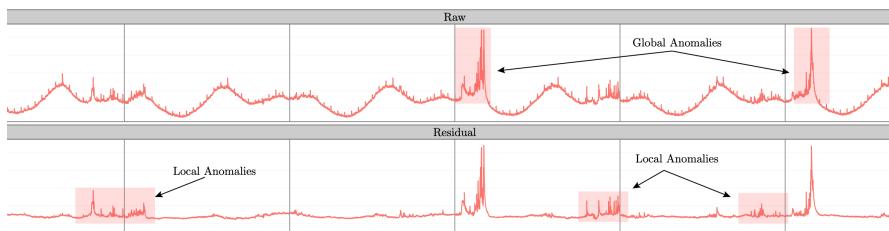




# Anomaly Detection: Decomposition based Model

- **Decomposition + Statistics**

- Seasonal ESD (S-ESD)
  - Seasonal-trend decomposition (STL) for residual
  - Extreme studentized deviate (ESD) test for anomalies
- Seasonal Hybrid ESD (S-H-ESD)
  - Robust statistics in ESD: median, median absolute deviation (MAD)
  - More robust to a higher percentage of anomalies
- Pros
  - Both global and local anomalies thanks to the STL decomposition
  - detect local anomalies that would otherwise be masked by seasonal data
- Cons
  - STL → not robustness to trend changes, seasonality changes



---

**Algorithm 1** S-ESD Algorithm**Input:**

$X$  = A time series

$n$  = number of observations in  $X$

$k$  = max anomalies (iterations in ESD)

**Output:**

$X_A$  = An anomaly vector wherein each element is a tuple  
(*timestamp, observed value*)

**Require:**

$k \leq (n \times .49)$

1. Extract seasonal component  $S_X$  using STL Variant
2. Compute median  $\tilde{X}$
- /\* Compute residual \*/
3.  $R_X = X - S_X - \tilde{X}$
- /\* Detect anomalies vector  $X_A$  using ESD \*/
4.  $X_A = \text{ESD}(R, k)$

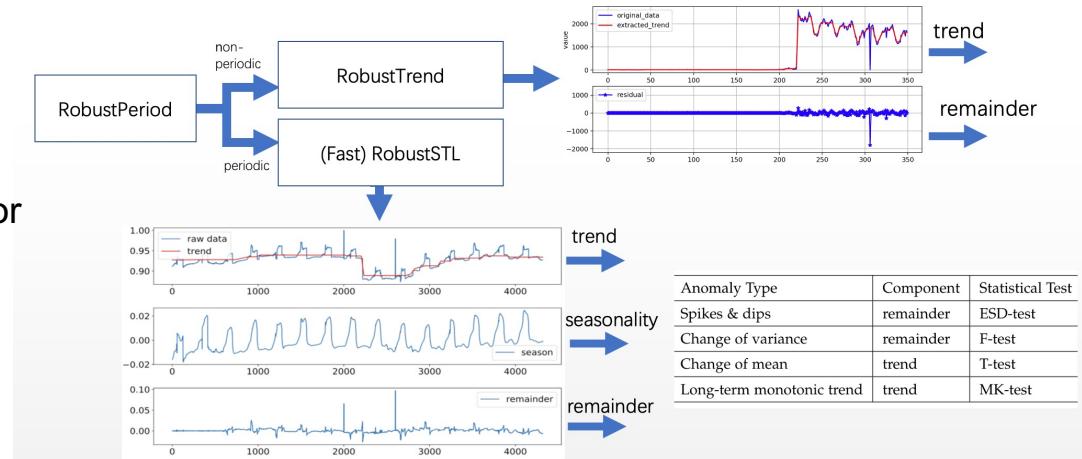
---

**return**  $X_A$

# Anomaly Detection: Decomposition based Model

- Robust Decomposition + Statistics

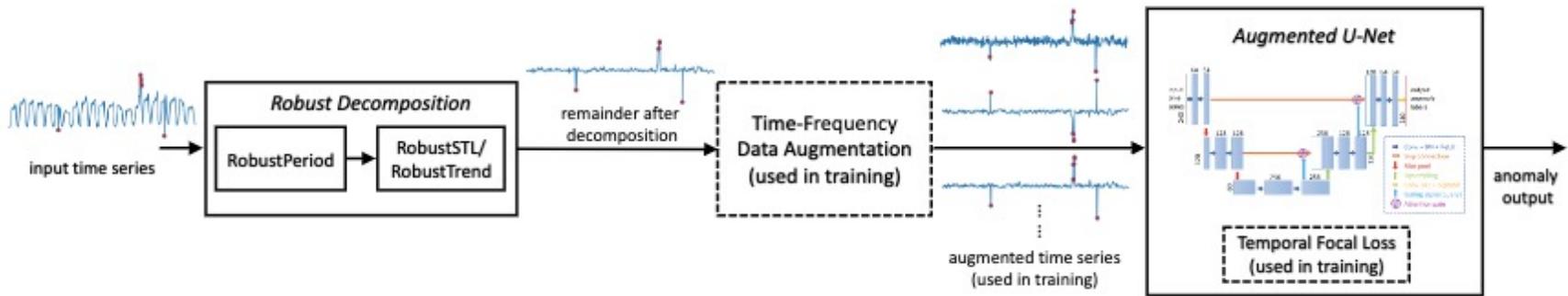
- Robust time series decompositions  
reduce complexity and bring  
explainability
- Robust statistical tests on each  
components lead to high accuracy for  
*different types of anomalies*





# Anomaly Detection: RobustTAD

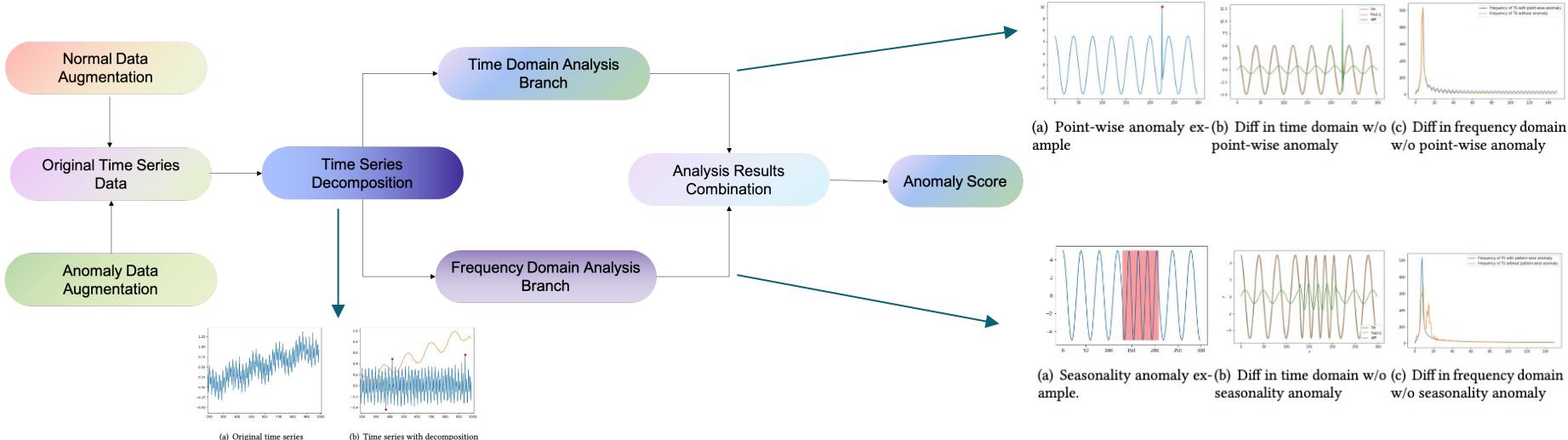
- Decomposition + Data Augmentation + U-Net (deep model)
  - Robust decomposition: handle complicated patterns, and simplify neural network
  - Data augmentation: mitigate the effects of limited labeled data
  - U-Net: capture multi-scale information of time series
  - Temporal focal loss: label-based weight and value-based weight for unbalanced label





# Anomaly Detection: TFAD

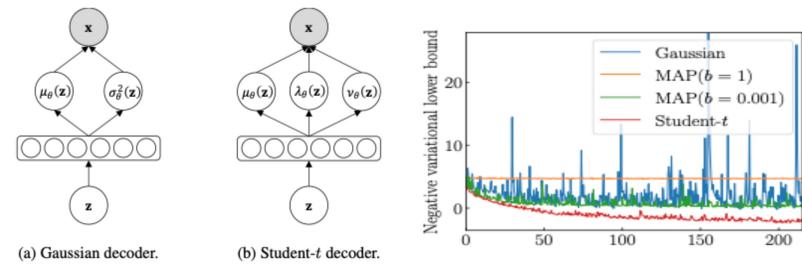
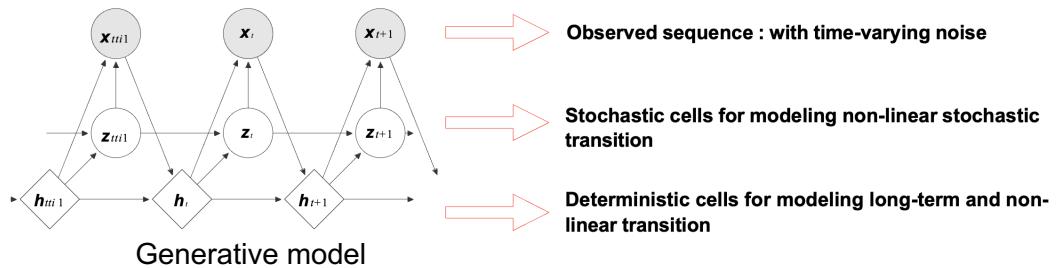
- Data Augmentation + Decomposition + TCN + *Time-Frequency Processing*
  - Anomalies types: seasonal anomaly, trend anomaly, global/context point anomaly
  - Point anomaly easier to detect in **Time**; seasonal anomaly easier to detect in **Frequency**



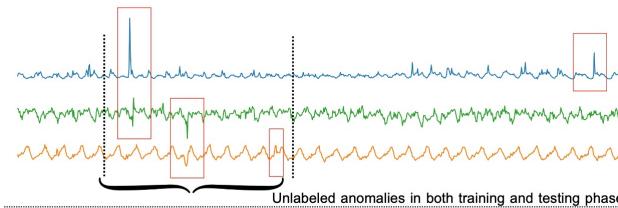


# Anomaly Detection: RDSSM

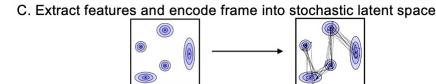
- Robust Deep State Space Model (*RDSSM*)
  - Robustness: training on contaminated data without labels
  - Model *temporal dependencies* with deep state space model
  - Adopt t-distribution in decoder for convergence on contaminated dataset



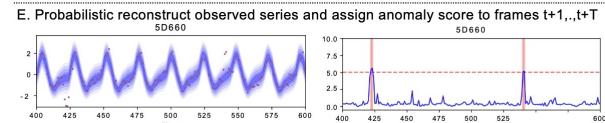
A. Slide window over noisy multi-variate time series contaminated with unlabeled anomalies



B. Get time series frames  $t+1, \dots, t+T$



C. Extract features and encode frame into stochastic latent space



D. Then random sample latent series from stochastic latent space

E. Probabilistic reconstruct observed series and assign anomaly score to frames  $t+1, \dots, t+T$



# Anomaly Detection: Anomaly Transformer

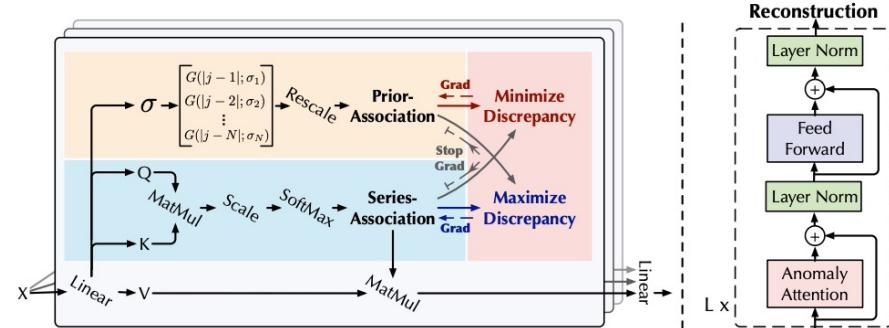
- Architecture: Anomaly-Attention

- Double branches model the prior-association and series-association simultaneously
- Association discrepancy: amplifying the difference between normal and abnormal points

$$\text{Prior-Association: } \mathcal{P}^l = \text{Rescale} \left( \left[ \frac{1}{\sqrt{2\pi}\sigma_i} \exp \left( -\frac{|j-i|^2}{2\sigma_i^2} \right) \right]_{i,j \in \{1, \dots, N\}} \right)$$

$$\text{Series-Association: } \mathcal{S}^l = \text{Softmax} \left( \frac{\mathcal{QK}^T}{\sqrt{d_{\text{model}}}} \right)$$

Adjacent-concentration inductive bias  
Find the most effective associations



$$\text{AssDis}(\mathcal{P}, \mathcal{S}; \mathcal{X}) = \left[ \frac{1}{L} \sum_{l=1}^L \left( \text{KL}(\mathcal{P}_{i,:}^l \| \mathcal{S}_{i,:}^l) + \text{KL}(\mathcal{S}_{i,:}^l \| \mathcal{P}_{i,:}^l) \right) \right]_{i=1, \dots, N}$$

**Symmetrized KL divergence** between **multi-level** prior- and series- associations  
(The adjacent-concentration property of series-association)

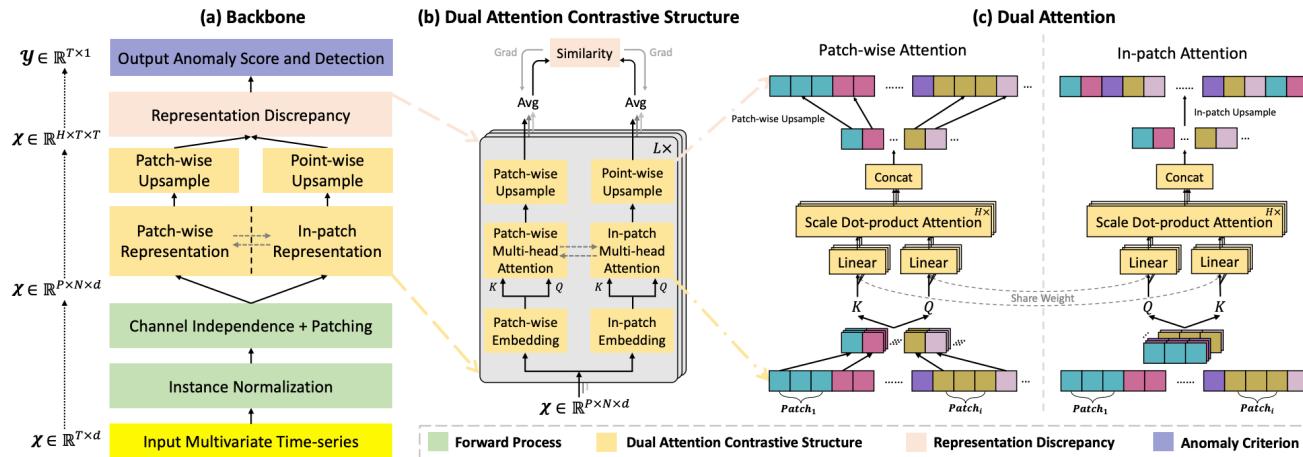
$$\mathcal{L}_{\text{Total}}(\hat{\mathcal{X}}, \mathcal{P}, \mathcal{S}, \lambda; \mathcal{X}) = \|\mathcal{X} - \hat{\mathcal{X}}\|_F^2 - \lambda \times \|\text{AssDis}(\mathcal{P}, \mathcal{S}; \mathcal{X})\|_1$$



# Anomaly Detection: DCdetector

- Dual Attention Contrastive Representation Learning

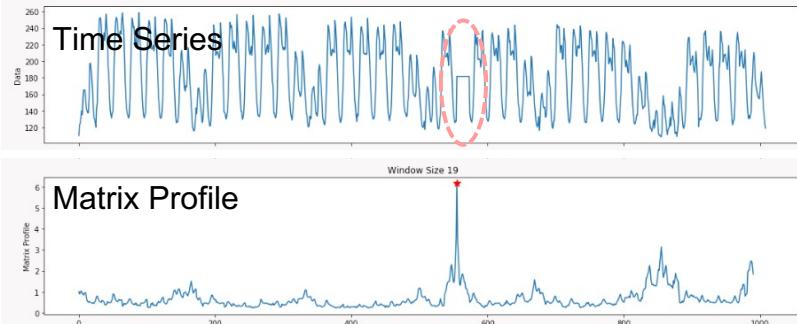
- Learning time series representations in different views: inner- and inter-patch
- Anomaly detection based on the discrepancy of the two representations



# Subsequence Anomaly Detection

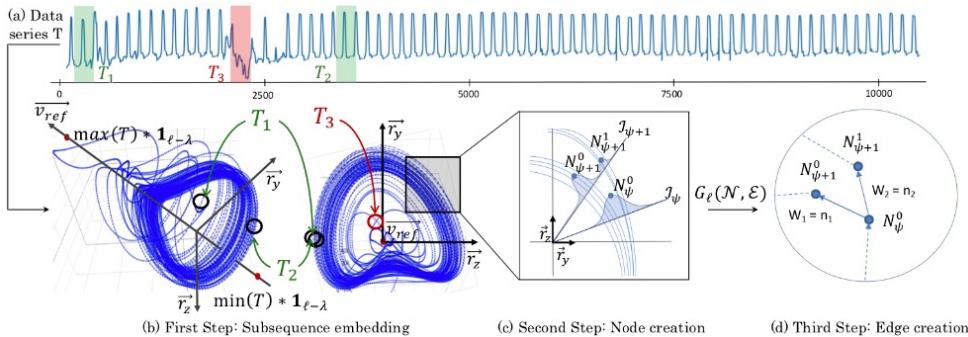
- Matrix Profile (MP)

- MP: records the distance of the subsequence to its nearest neighbor
- The higher the MP value, the greater the dissimilarity  
→ such areas are anomalies/discords



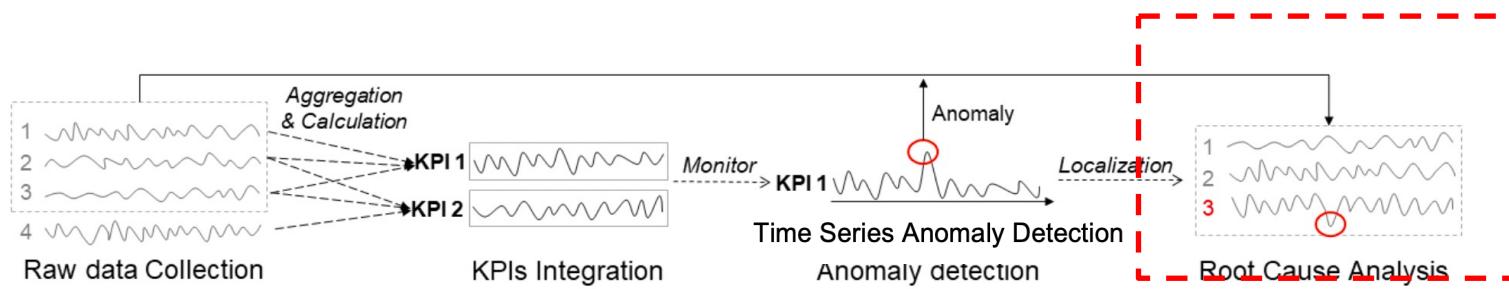
- Series2Graph

- Graph representation of subsequences
- Anomaly score based on the graph  
(edges/nodes and their weights/degrees)



# From Time Series Anomaly Detection to Localization

- Fault cause localization
  - Identify the cause for the fault by root cause analysis (RCA) that can best explain the observed system disorders



Challenge for RCA is how to accurately and quickly find a set of dimension-value combinations to resolve the anomaly detected given the huge search space

[1] CMMD: Cross-Metric Multi-Dimensional Root Cause Analysis, <https://arxiv.org/abs/2203.16280>

[2] Integrated Fault and Security Management, Ehab Al-Shaer, Yan Chen, in Information Assurance, 2008



# Multi-dimensional Fault Cause Localization

- Time series with multi-dimensional attributes are usually collected and monitored

Category	City	Online	Sales	Anomaly
shirt	Beijing	Y	3042	1
shirt	Beijing	N	5401	1
skirt	Beijing	Y	3103	0
shirt	Shanghai	N	7301	0
jacket	Shanghai	N	4133	0
skirt	Hangzhou	Y	2789	0



Root cause is  
Category = shirt & City = Beijing

- Fault cause localization aims to find the root cause (a combination of attribute-values) that contribute most to the total value of anomalous



# Fault Cause Localization

- Existing methods

- Bottom up
  - 1. Genetic algorithm: CMMD [1]
  - 2. Monte Carlo tree search: HotSpot [2], GMCTS [3]
- Top down
  - 1. Score based clustering: Squeeze [4], AutoRoot [5]



exponential complexity  
or sacrificing accuracy

[1] Yan, Shifu, et al. "CMMD: Cross-Metric Multi-Dimensional Root Cause Analysis." arXiv preprint arXiv:2203.16280 (2022).

[2] Sun, Yongqian, et al. "Hotspot: Anomaly localization for additive kpis with multi-dimensional attributes." IEEE Access 6 (2018): 10909-10923.

[3] Wang, Chunlin, et al. "Network Abnormality Location Algorithm Based on Greedy Monte Carlo Tree." 2022 14th International Conference on Machine Learning and Computing (ICMLC). 2022.

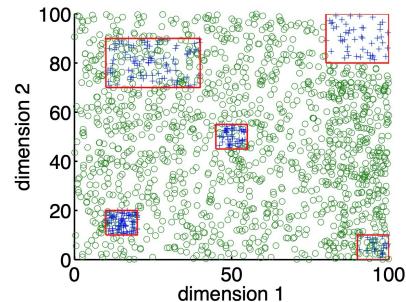
[4] Li, Zeyan, et al. "Generic and robust localization of multi-dimensional root causes." 2019 IEEE 30th International Symposium on Software Reliability Engineering (ISSRE). IEEE, 2019.

[5] Wang, Hanzhang, et al. "Groot: An event-graph-based approach for root cause analysis in industrial settings." 2021 36th IEEE/ACM International Conference on Automated Software Engineering (ASE). IEEE, 2021.



# Fault Cause Localization via Rule Set Learning

- Fault Cause Localization as a classification problem
  - Our goal is to find the discriminative attribute=value combinations best explain the outliers, but do not cover normal cases
- Our solutions: rule set learning



X	Y	Z	S	T	label
x1	y3	z1	s2	t1	1
x1	y2	z3	s1	t4	0
...	...	...	...	...	...

IF ( X=x1 AND Y=y2 )  
OR ( X=x2 AND Z=z3 )  
OR ( S=s1 AND t1<=T<=t3 )  
OR ( ... )  
THEN label=1

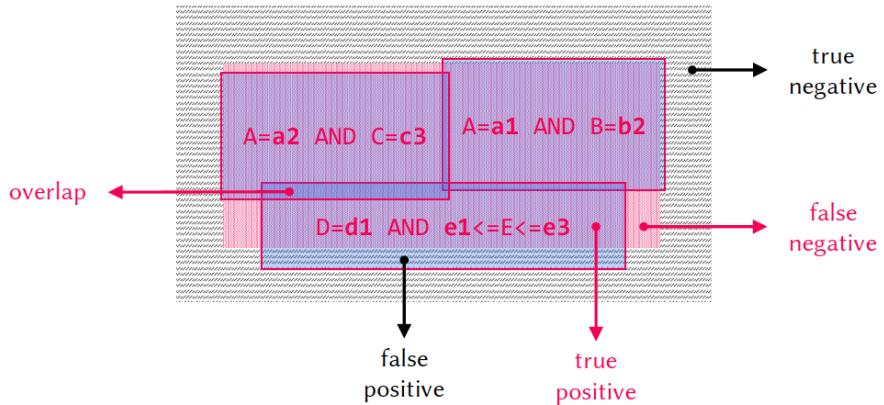
The rule set learning methods aim to find discriminative patterns that only covers the samples of interest class



# Interpretable Rule Set Learning

- High classification accuracy: small number of false positive and false negative
- Interpretable rules: less rules with shorter length and less overlap
- Good scalability

$$L(\mathcal{S}) = \beta_0 \sum_{\mathcal{R} \in \mathcal{S}} |\mathcal{X}_{\{\mathcal{R}\}}^-| + \beta_1 |\mathcal{X}^+ \setminus \mathcal{X}_{\mathcal{S}}^+| + \beta_2 \left( \sum_{\mathcal{R} \in \mathcal{S}} |\mathcal{X}_{\{\mathcal{R}\}}^+| - |\mathcal{X}_{\mathcal{S}}^+| \right) + \lambda \sum_{\mathcal{R} \in \mathcal{S}} |\mathcal{R}|$$



# Interpretable Rule Set Learning: A Submodular Optimization Approach

Minimize  $L(\mathcal{S}) = \beta_0 \sum_{\mathcal{R} \in \mathcal{S}} |\mathcal{X}_{\{\mathcal{R}\}}^-| + \beta_1 |\mathcal{X}^+ \setminus \mathcal{X}_{\mathcal{S}}^+| + \beta_2 \left( \sum_{\mathcal{R} \in \mathcal{S}} |\mathcal{X}_{\{\mathcal{R}\}}^+| - |\mathcal{X}_{\mathcal{S}}^+| \right) + \lambda \sum_{\mathcal{R} \in \mathcal{S}} |\mathcal{R}|$

FP	FN	Overlap	Complexity
----	----	---------	------------

Reorganize  $L(\mathcal{S}) = \beta_1 |\mathcal{X}^+| - (\beta_1 + \beta_2) |\mathcal{X}_{\mathcal{S}}^+| + \sum_{\mathcal{R} \in \mathcal{S}} \beta_0 |\mathcal{X}_{\{\mathcal{R}\}}^-| + \beta_2 |\mathcal{X}_{\{\mathcal{R}\}}^+| + \lambda |\mathcal{R}|$

Const	Submodular	Modular
-------	------------	---------

Maximize  $V(\mathcal{S}) = g(\mathcal{S}) - \sum_{\mathcal{R} \in \mathcal{S}} c(\mathcal{R})$

A cardinality constrained submodular maximization problem



# Interpretable Rule Set Learning: A Submodular Optimization Approach

- Distorted greedy algorithm

---

**Algorithm 1** Rule set learning

---

```
1 Input: Training data  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ , hyperparameters  $(\beta, \lambda)$ , cardinality  $K$ 
2 Initialize  $\mathcal{S} \leftarrow \emptyset$ 
3 for  $k = 1$  to  $K$  do
4   Define  $v_k(\mathcal{R}) = (1 - 1/K)^{K-k} g(\mathcal{R}|\mathcal{S}) - c(\mathcal{R})$            /*  $g(\mathcal{R}|\mathcal{S}) := g(\mathcal{S} \cup \{\mathcal{R}\}) - g(\mathcal{S})$  */
5   Solve  $\mathcal{R}^* \leftarrow \arg \max_{\mathcal{R} \subset [d]} v_k(\mathcal{R})$ 
6   if  $v_k(\mathcal{R}^*) > 0$  then  $\mathcal{S} \leftarrow \mathcal{S} \cup \{\mathcal{R}^*\}$  end if
7 end for
8 Output:  $\mathcal{S}$ 
```

---

Exhaustive  
enumeration  $O(2^d)$  ☹

- Approximation guarantee

$$V(\mathcal{S}) = g(\mathcal{S}) - \sum_{\mathcal{R} \in \mathcal{S}} c(\mathcal{R}) \geq (1 - 1/e)g(OPT) - \sum_{\mathcal{R} \in OPT} c(\mathcal{R})$$



# Interpretable Rule Set Learning: A Submodular Optimization Approach

- We rewrite the subobjective as a difference of submodular (DS) functions

**Maximize**  $v(\mathcal{R}) = \sum_{i=1}^n \omega_i \mathbb{1}_{\mathcal{R} \subseteq \mathbf{x}_i} - \lambda |\mathcal{R}|$

**Rewrite**  $v(\mathcal{R}) = \sum_{i=1}^n \omega_i + \sum_{i: \omega_i < 0} -\omega_i \mathbb{1}_{\mathcal{R} \not\subseteq \mathbf{x}_i} - \sum_{i: \omega_i > 0} \omega_i \mathbb{1}_{\mathcal{R} \not\subseteq \mathbf{x}_i} - \lambda |\mathcal{R}|$

**Submodular** **Submodular** **Modular**

**Minorize-Maximization**

Linearization

Modular lower bound      Modular upper bound

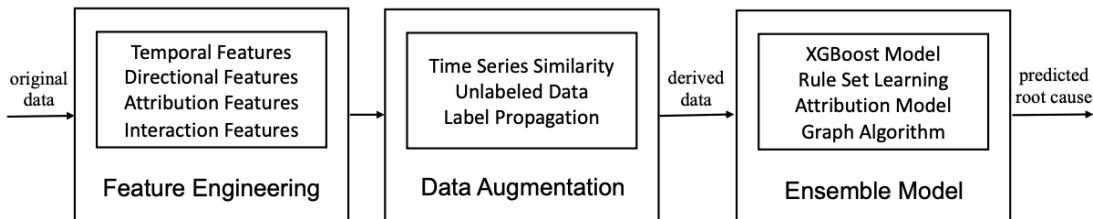
- If  $v(\mathcal{R})$  is nonnegative monotonic submodular function, it can be optimized using greedy algorithm
- We express it as the difference of submodular functions and apply MM algorithm



# NetRCA

- Effective and efficient root cause localization for network faults

- Key idea: 1. The root cause localization problem can be treated as a rule learning problem  
2. Multiple models ensemble can further improve the precision
- Three blocks: feature engineering, data augmentation, model ensemble





# XAI for Time Series

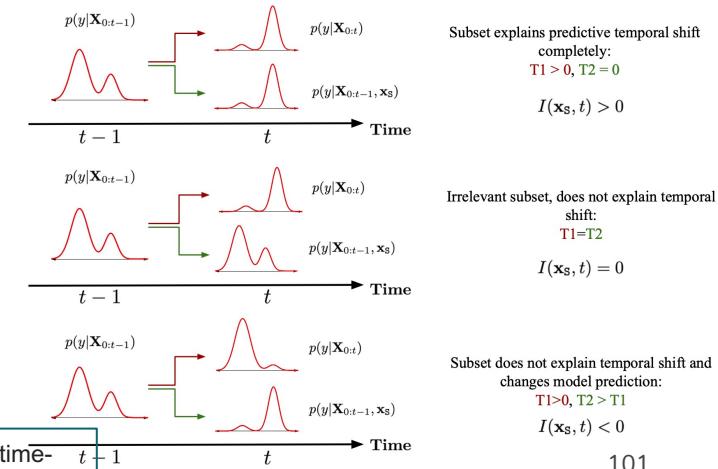
- Feature Importance in Time (FIT)

- Key idea: The importance of observations over time is related on their contribution to the temporal shift of the model output distribution.
- The importance of new set of observations  $S$  at time  $t$  is quantified by the shift in predictive distribution when only  $S$  are observed while accounting for the distributional shift over time

$$I(\mathbf{x}_S, t) = \underbrace{\text{KL}(p(y|\mathbf{X}_{0:t}) \parallel p(y|\mathbf{X}_{0:t-1}))}_{\text{T1: temporal distribution shift}} - \underbrace{\text{KL}(p(y|\mathbf{X}_{0:t}) \parallel p(y|\mathbf{X}_{0:t-1}, \mathbf{x}_S, t))}_{\text{T2: unexplained distribution shift}}$$

estimates the distributional shift between the predictive distribution at time  $t-1$  to  $t$

measures the residual shift after adding only  $\mathbf{x}_{S,t}$





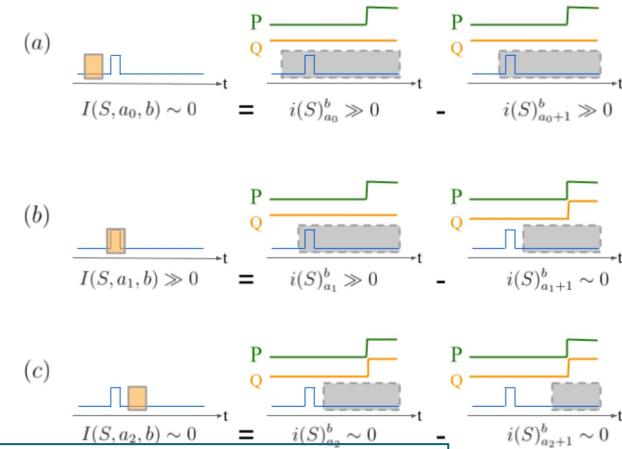
# XAI for Time Series

- Windowed Feature Importance in Time (WinIT)
  - Key idea: accounts for the temporal dependence between different observations of the same feature and computes the impact of a feature observation to predictions over a window of time steps.
  - Measure the importance score for a feature subset S at time t = a to the prediction at t = b by aggregate the importance scores across the forward prediction time window

$$\mathcal{I}(S, a) = \frac{1}{\hat{N} + 1} \sum_{b=a}^{a+\hat{N}} I(S, a, b), \quad I(S, a, b) = \begin{cases} i(S)_a^b - i(S)_{a+1}^b, & a < b < a + N \\ i(S)_a^a, & b = a \end{cases}$$

where  $i(S)_a^b$  measures the distance between  $P = p(y_b | \mathbf{X}_{1:b})$  and  $Q = p(y_b | \tilde{\mathbf{X}}_{a:b}^S)$ .

- Compared with FIT:
  - 1, better performance
  - 2, higher computational complexity

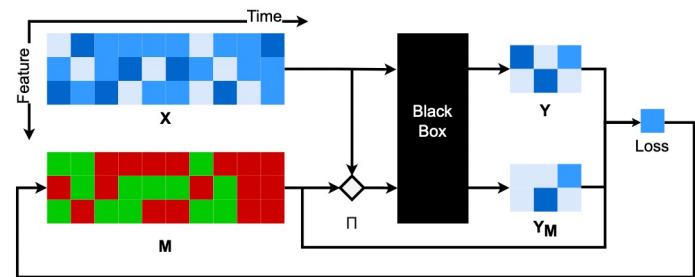




# XAI for Time Series

- Explaining Time Series Predictions with Dynamic Masks
  - Key idea: learn mask to highlight the regions of the image that are salient for the black-box model to issue its prediction
  - masking error:
$$\mathcal{L}_e(\mathbf{M}) = \sum_{t=t_y}^T \sum_{i=1}^{d_Y} \left( [(f \circ \Pi_{\mathbf{M}})(\mathbf{X})] \right)$$
  - Regularization:
    - 1, sparseness:  $\mathcal{L}_a(\mathbf{M}) = \|\text{vecsorth}(\mathbf{M}) - \mathbf{r}_a\|^2$
    - 2, smoothness:  $\mathcal{L}_c(\mathbf{M}) = \sum_{t=1}^{T-1} \sum_{i=1}^{d_X} |m_{t+1,i} - m_{t,i}|$
  - Overall objective:

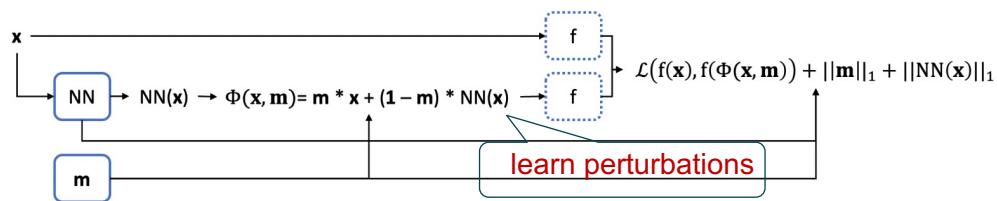
$$\mathbf{M}_a^* = \arg \min_{\mathbf{M} \in [0,1]^{T \times d_X}} \mathcal{L}_e(\mathbf{M}) + \lambda_a \cdot \mathcal{L}_a(\mathbf{M}) + \lambda_c \cdot \mathcal{L}_c(\mathbf{M}).$$





# XAI for Time Series

- Learning Perturbations to Explain Time Series Predictions
  - Key idea: simply masking part of time series will leads to out of distribution (OOD) problem, thus we can learn a perturbation for the masked timesteps.

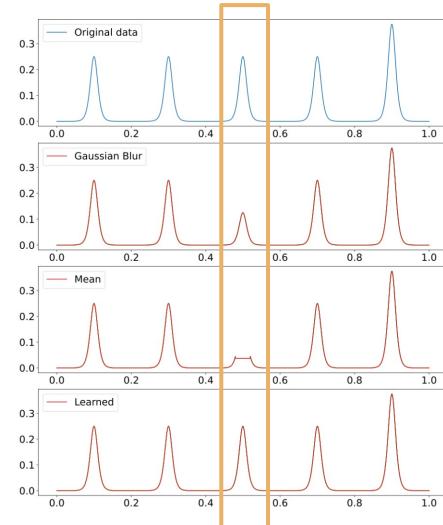


- Overall objective:

$$\arg \min_{\mathbf{m}, \Theta \in \text{NN}} \lambda_1 \|\mathbf{m}\|_1 + \lambda_2 \|\text{NN}(\mathbf{x})\|_1 + \mathcal{L}(f(\mathbf{x}), f(\Phi(\mathbf{x}, \mathbf{m})))$$



avoid the case that  $\mathbf{m} = \mathbf{0}$ , and  $\text{NN}(\mathbf{x}) \approx \mathbf{x}$





# Tutorial Summary

- ❑ **Introduction:** Real-world Challenges and Needs for Robust Time Series Processing
- ❑ **Preliminaries:** Multiple Disciplines: Statistics, Signal Processing, Optimization, Deep Learning, XAI
- ❑ **Robust Time Series Processing Blocks**
  - Time Series Periodicity Detection
  - Time Series Decomposition: Trend Filtering, Seasonal-Trend Decomposition
  - Time Series Similarity
- ❑ **Robust Time Series Applications and Practices**
  - Time Series Forecasting
  - From Forecasting to Decision-Making: Autoscaling
  - Time Series Anomaly Detection
  - From Anomaly Detection to Localization: Fault Cause Localization
  - XAI for Time Series Data