



KDD 2022 Tutorial



Robust Time Series Analysis and Applications: An Industrial Perspective



Qingsong Wen

Staff Engineer
Alibaba DAMO



Linxiao Yang

Senior Engineer
Alibaba DAMO



Tian Zhou

Senior Engineer
Alibaba DAMO



Liang Sun

Engineering
Director
Alibaba DAMO

- **Date and Time:** Aug. 14 (Sun.), 9:00 am - 12:00 pm
- **Location:** Room 207B, Washington DC Convention Center, USA



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:** 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 with a red border. To the right is the DAMO logo with the text ALIBABA DAMO ACADEMY. Below the navigation, there are three main categories: Machine Intelligence, Data Computing, and Robotics. Under Machine Intelligence, the Decision Intelligence Lab is highlighted with a red border. Other labs listed include Speech Lab, Vision Lab, Language Technology Lab, City Brain Lab, Computing Technology Lab, Data Analytics and Intelligence Lab, Database and Storage Lab, OS Lab, and Autonomous D. The background features a dark blue graphic with a hand pointing at a chart and text related to success and business strategy.



Outline

- Introduction** (by Liang Sun)
 - Preliminaries** (by Liang Sun)
 - Robust Time Series Processing Blocks**
 - Time Series Periodicity Detection (by Qingsong Wen)
 - Time Series Trend Filtering (by Qingsong Wen)
 - Time Series Seasonal-Trend Decomposition (by Qingsong Wen)
 - Time Series Similarity (by Liang Sun)
- (10 minutes break)
- Robust Time Series Applications and Practices**
 - Forecasting: Tree Models, Deep Ensemble, Transformers, etc. (by Qingsong Wen)
 - Autoscaling: Query Modeling, Scaling Decision, etc. (by Qingsong Wen)
 - Anomaly Detection: Decomposition Model, Deep State Space Model, Transformers, etc. (by Qingsong Wen)
 - Fault Cause Localization: Rule Set Learning, Root Cause Analysis, etc. (by Liang Sun)

(Q&A Session)



Outline

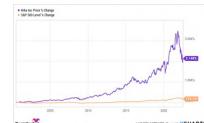
- *Introduction*
- Preliminaries
- Robust Time Series Processing Blocks
- Robust Time Series Applications and Practices



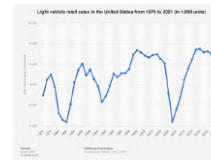
Time Series Data is Ubiquitous

- A wide range of time series data
 - AIOps
 - IoT
 - Business data, e.g., sales volume, stock price
 - Many others

stocks



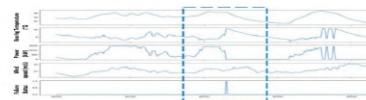
sales



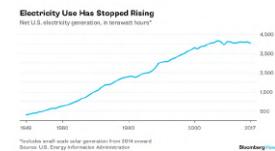
goods consumption



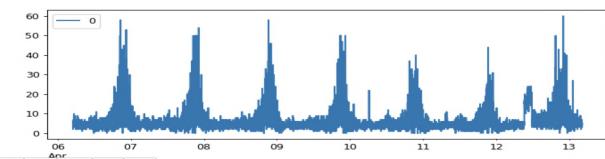
sensor



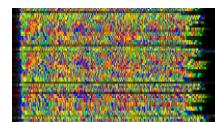
power demand



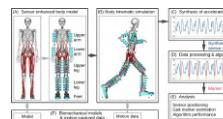
Cloud service monitoring



DNA sequence



motion detect

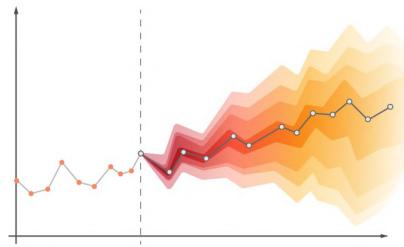


ECG

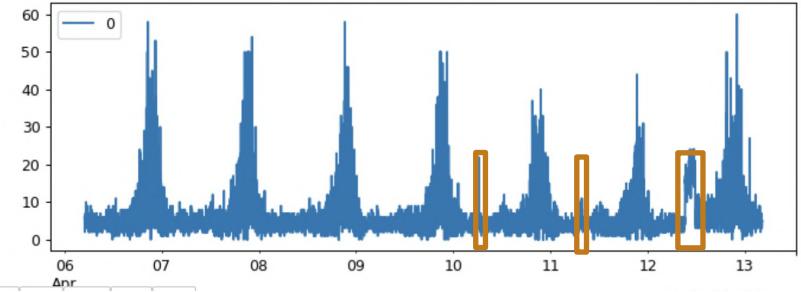


Typical Applications of Time Series

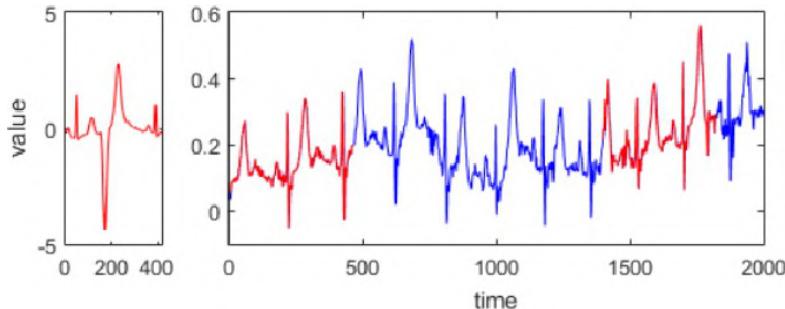
Time Series Forecasting



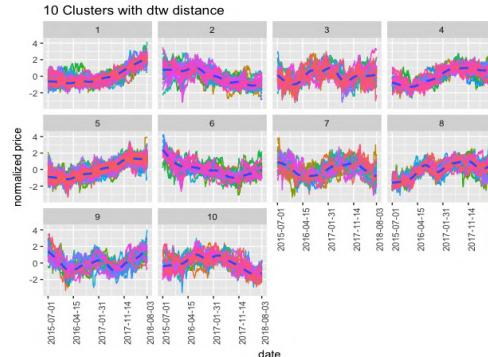
Time Series Anomaly Detection



Time Series Search/Query



Time Series Classification/Clustering

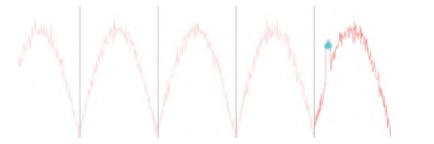




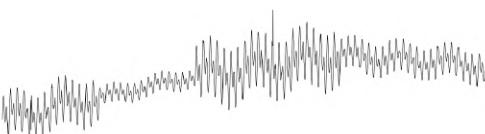
AIOps: From Anomaly Detection to Root Cause Analysis

Different Types of Anomalies

Spike & dip

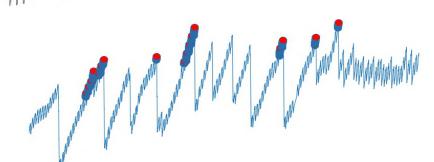


Change of mean

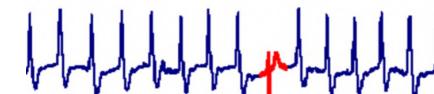


Change of variance

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):

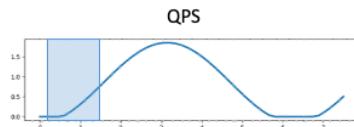
- 5M+ servers monitored
- 50M+ metrics monitored
- In 2020 some typical root identification time is ~ 1 person-month
- Some cases incur 1M+ USD loss

- ***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

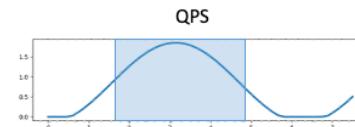
- Autoscaling in cloud computing is an effective method to improve the usage of computing resources
 - It automatically allocates resources for cloud-based applications while maintaining SLA (service level agreement)
 - Horizontal scaling (add/delete instances or VMs) vs vertical scaling (up/downgrade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



load
balancer

instance

instance

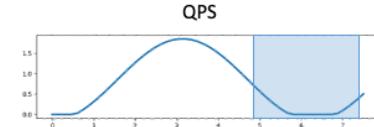


load
balancer

instance

instance

instance



load
balancer

instance

instance



Applications of Time Series in Industry

- Huge amount of data
 - Hundreds of millions metrics or even more
 - High frequency: unlike traditional time series, 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)

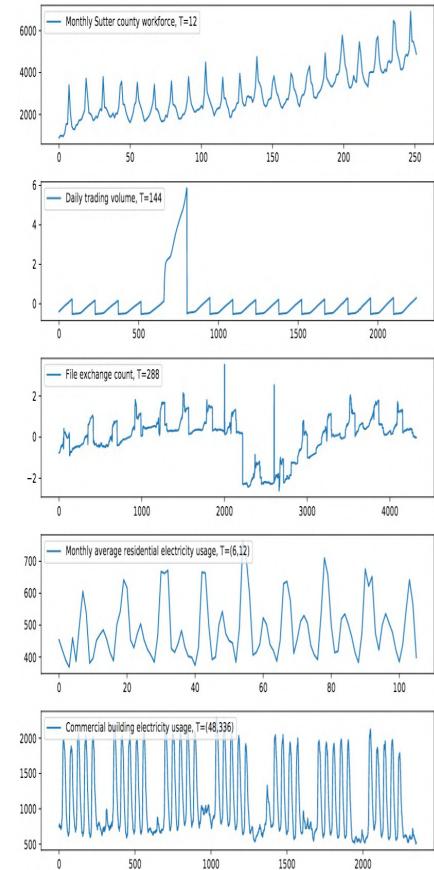
*Our solution: **robust** algorithm family for time series analysis*



Why Periodicity so Important?

- Periodic time series dominate in many real-world applications
 - 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
 - Many more periodic time series in a wide range of real-world applications
- Why periodicity so important?
 - Periodic and non-periodic time series requires quite different processing
 - Different periodic patterns requires different processing
- Challenges
 - Complicated periodic patterns: multiple periodic components, periodic shift and 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

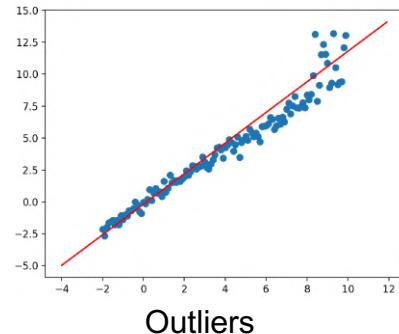
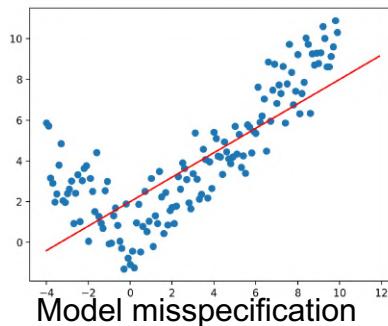
❑ Robust Time Series Processing Blocks

❑ Robust Time Series Applications and Practices



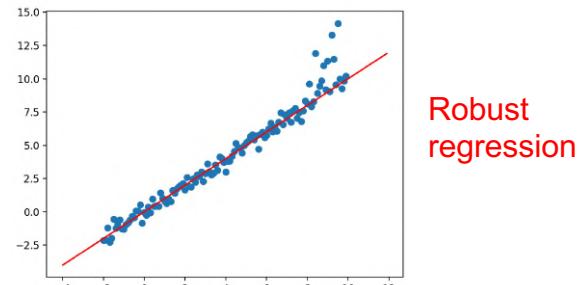
Robust Regression

- 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 ρ

- 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 ξ increases



M-estimator

- Robust M-estimator

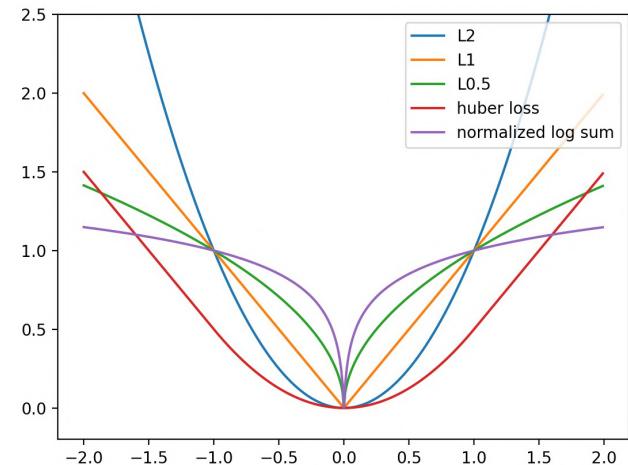
- L1-norm

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

- 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

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

- 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

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

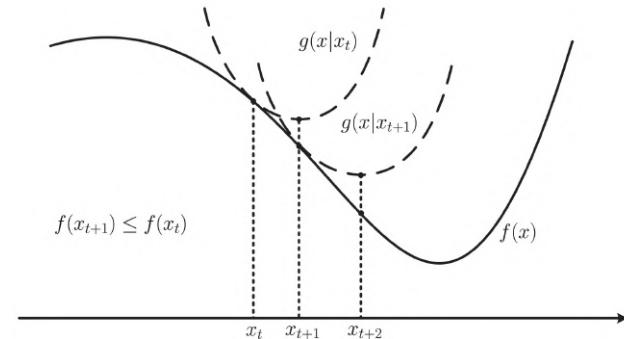
- Surrogate function

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

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

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

- 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

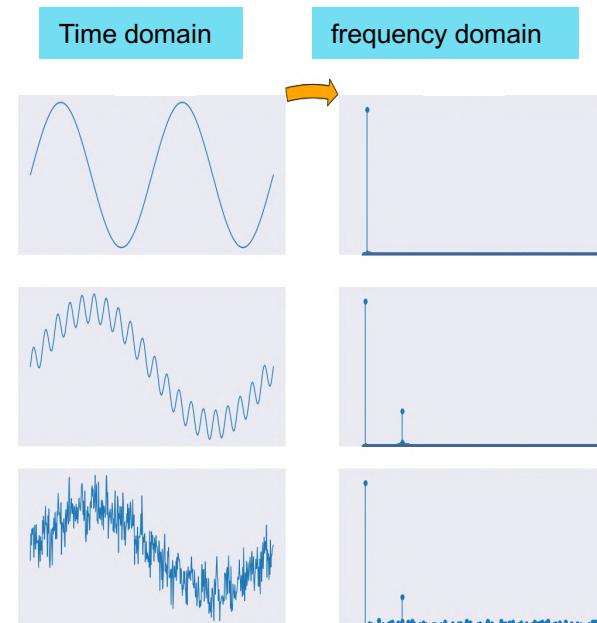
- Time domain and frequency domain

- 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

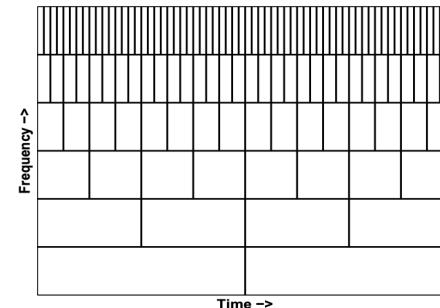
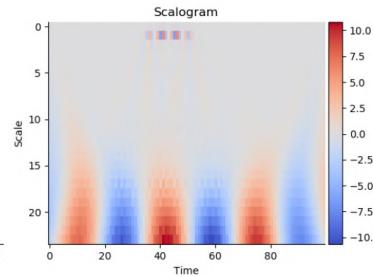
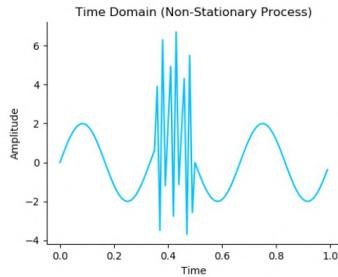
$$\hat{f}(\xi) = \int_{-\infty}^{\infty} f(x) e^{-i2\pi\xi x} dx, \quad \forall \xi \in \mathbb{R}$$





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

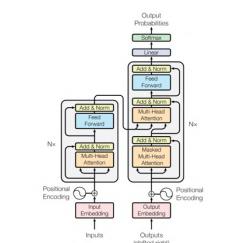
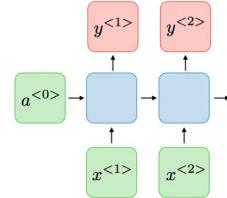
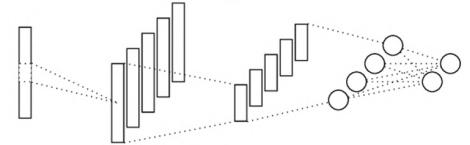
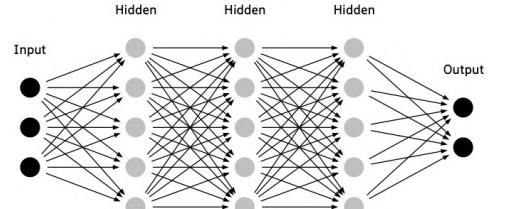
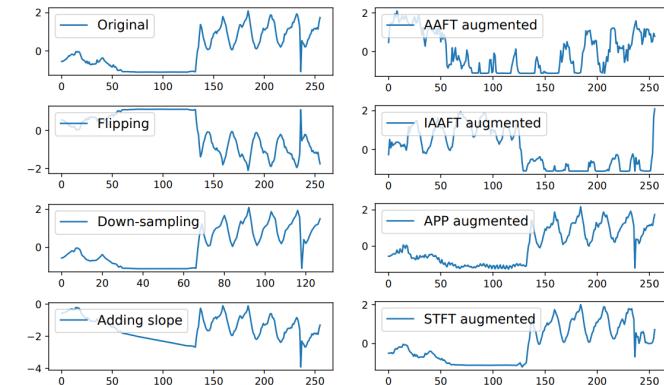
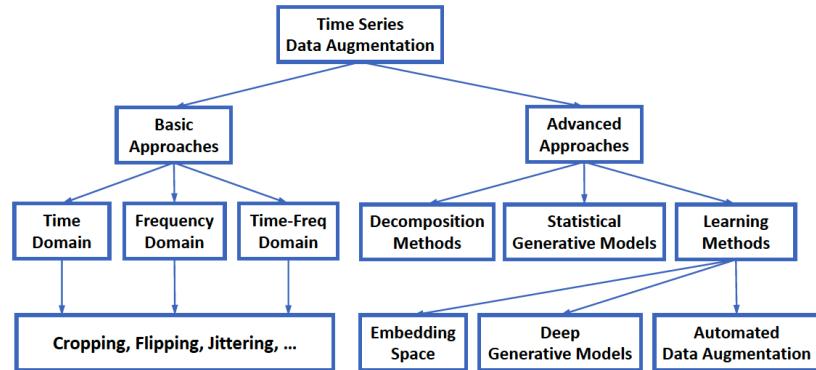


Figure 1: The Transformer - model architecture.



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)



(a) time domain

(b) (time-)frequency domain

Data augmentations in time, frequency, time-frequency domains



Outline

❑ Introduction

❑ Preliminaries

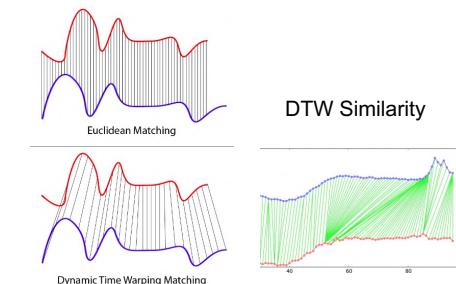
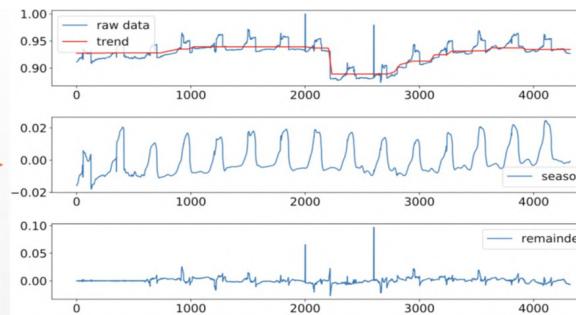
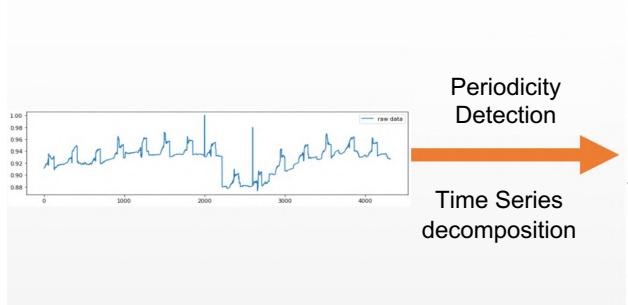
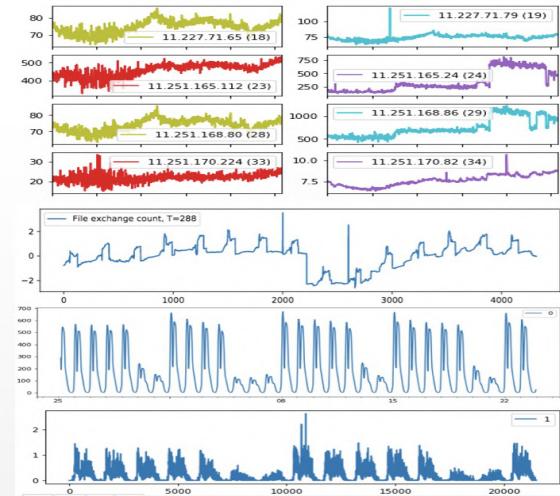
➤ ***Robust Time Series Processing Blocks***

- 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



DTW Similarity

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

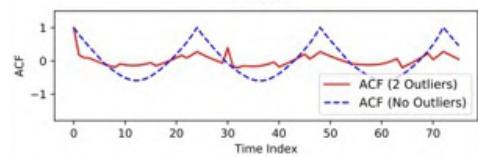
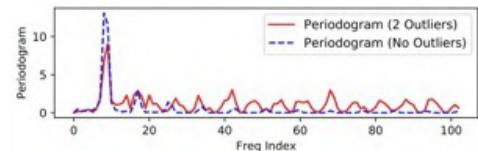
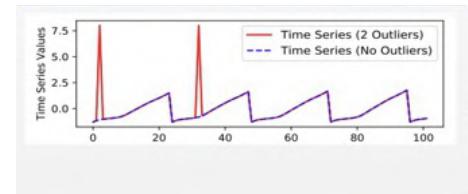
- Freq domain: Fisher's test by Periodogram (FFT)

- 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 $k = \arg \max_k P_k$



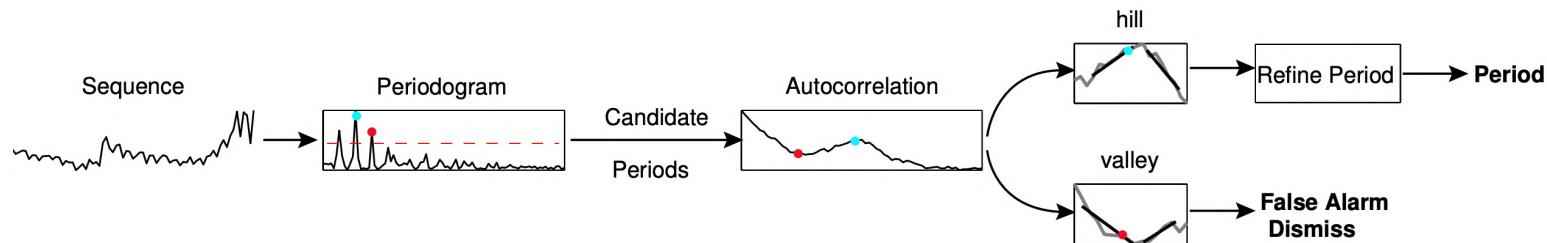
- Not robust to outliers
- Hard to handle multiple periodicities



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

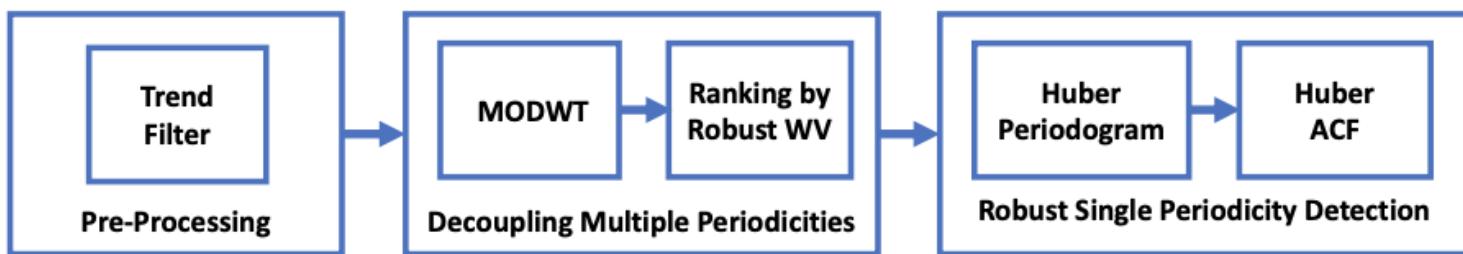
Method	Easy to threshold	Accurate short periods	Accurate large periods	Complexity
Periodogram	yes	yes	no	$O(N \log N)$
Autocorrelation	no	yes	yes	$O(N \log N)$
Combination	yes	yes	yes	$O(N \log N)$





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

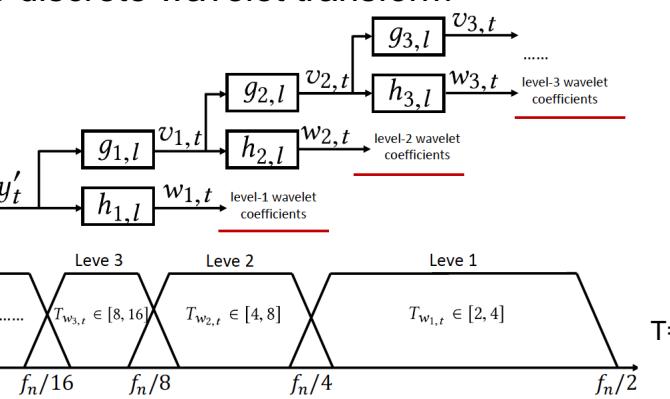
- MODWT for decoupling multiple periodicities
 - MODWT: maximal overlap discrete wavelet transform

*j*th level wavelet coefficients

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

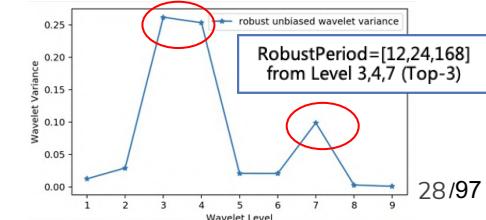
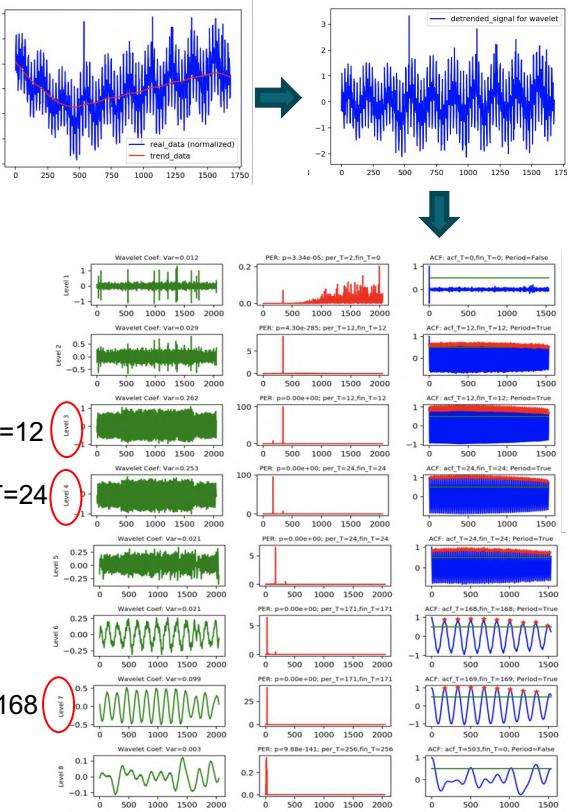
*j*th level scaling coefficients

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



- 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

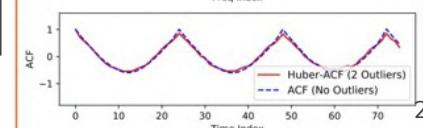
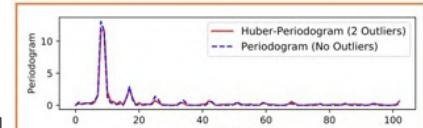
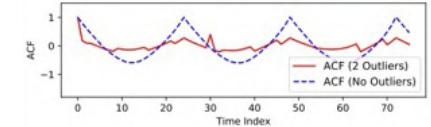
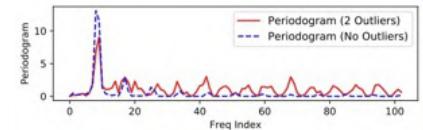
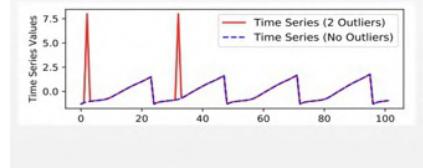
$$\text{HuberACF}(t) = \frac{p_t}{(N-t)p_0}$$

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

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

robust to outliers

Demo of Huber-Periodogram/ACF

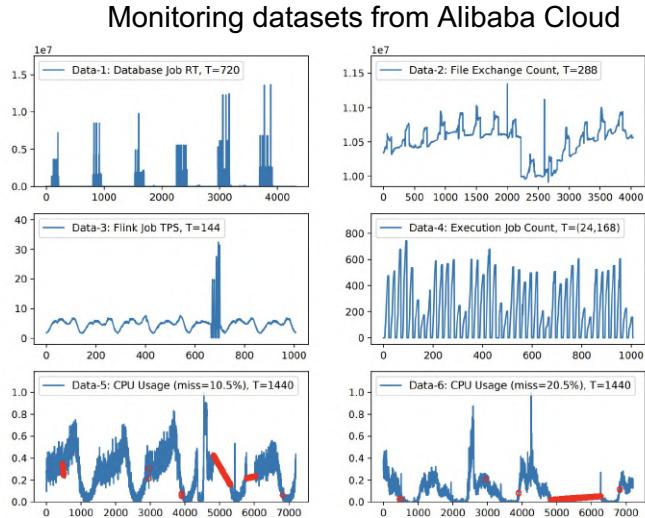




Periodicity Detection Comparisons

- Real-world datasets:

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



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	✗	✗	✗	✓	✓
RobustPeriod	✓	✓	✓	✓	✓



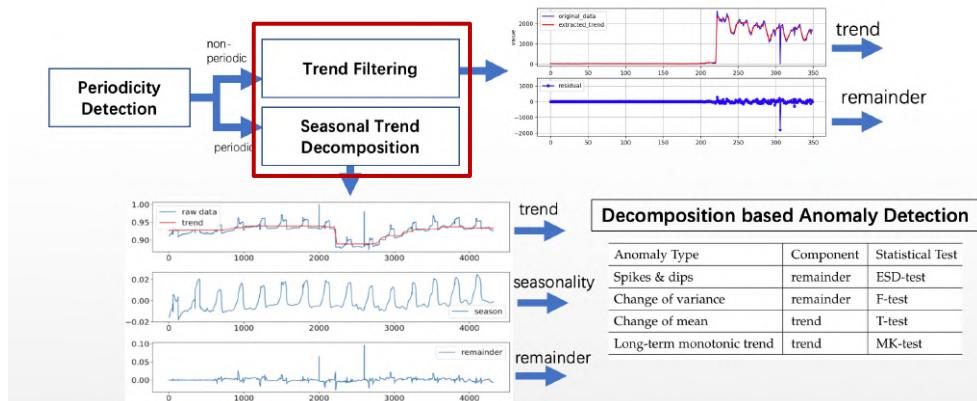
Time Series Decomposition

- 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

$$y_t = \tau_t + r_t$$

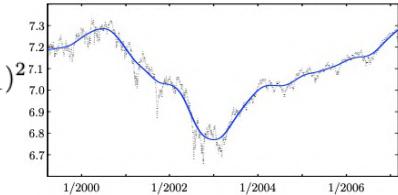
- Moving average: not accurate with delay

- Hodrick–Prescott filtering (H-P filter)

- Regularization with ℓ_2 norm of the second difference operator
- linear estimation, convex with closed-form solution

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

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

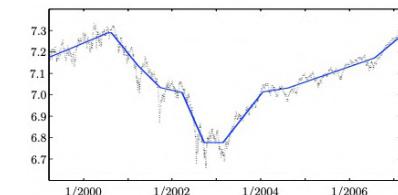


- ℓ_1 trend filtering

- Regularization with ℓ_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

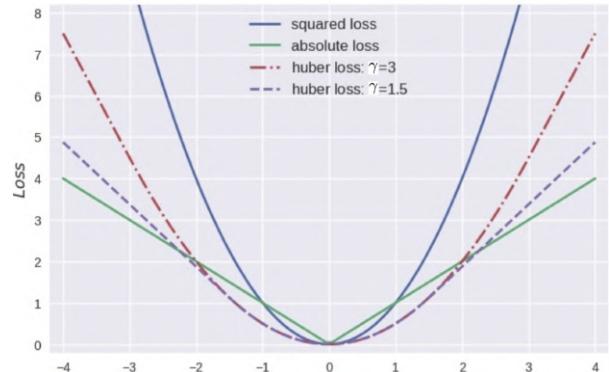
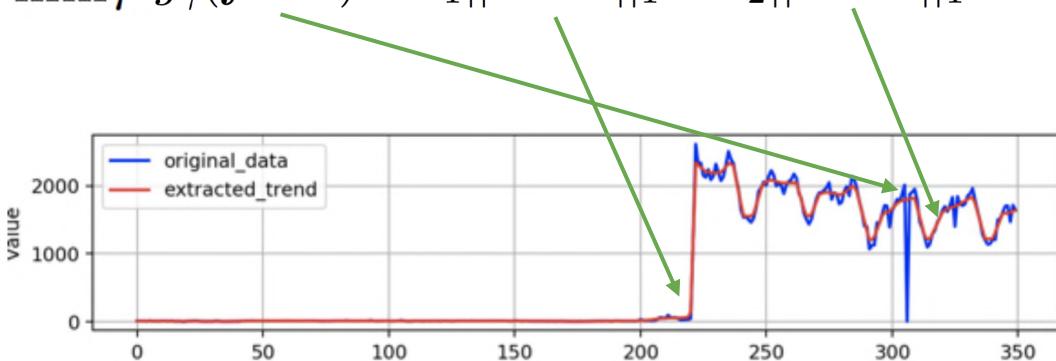
- Not robust to outliers and abrupt trend changes



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

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



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

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



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 $\mathbf{A} = \text{diag}(g'_\gamma(x_i^k)) \text{diag}^{-1}(g_\gamma(x_i^k))$

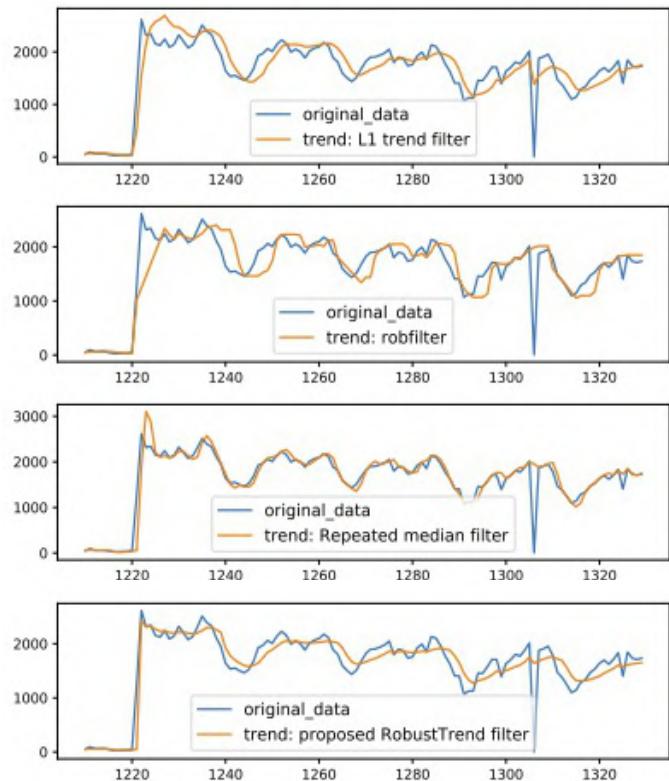
- Now the approximated Tau–minimization step

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

- 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**: best tradeoff under outliers and abrupt trend changes

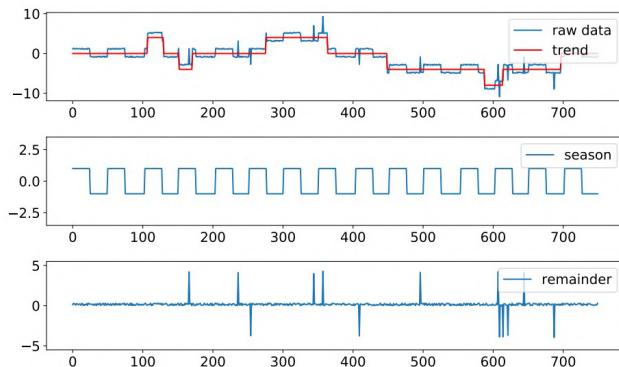




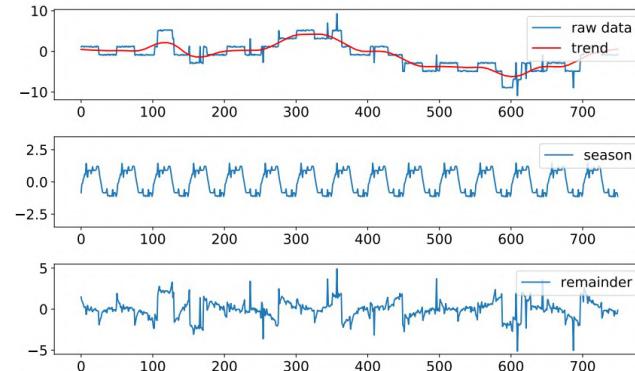
Periodic Time Series: Common STL Algorithm

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

- 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



complex synthetic data with trend changes and seasonality shift

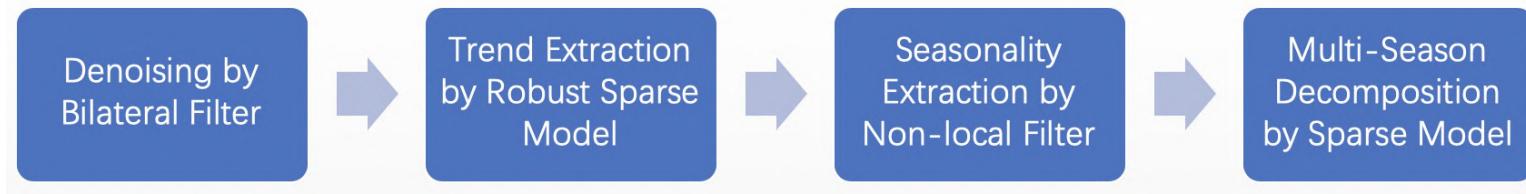


results of STL decomposition



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](#))

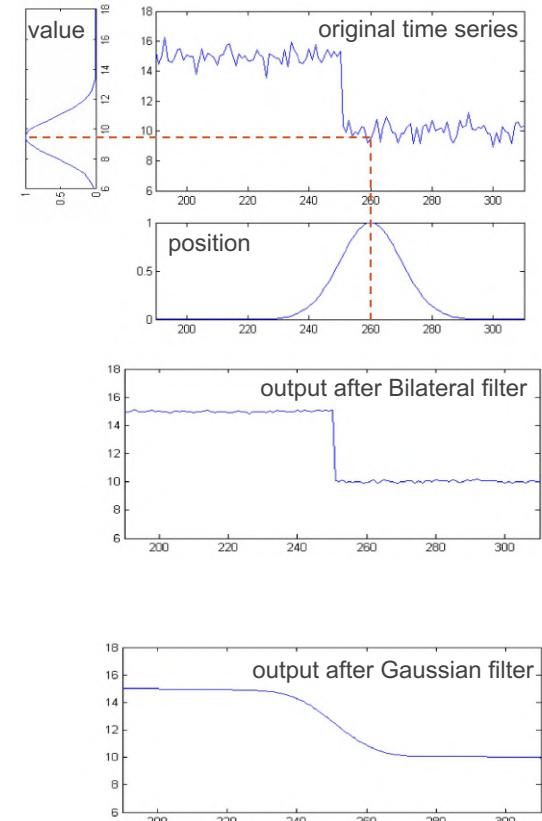
Linxiao Yang, Qingsong Wen, Bo Yang, Liang Sun, "A Robust and Efficient Multi-Scale Seasonal-Trend Decomposition," in Proc. IEEE 46th International Conference on Acoustics, Speech, and Signal Processing (ICASSP 2021), Toronto, Canada, Jun. 2021. ([multiple-scale version](#))

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

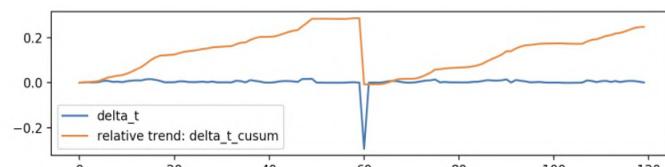
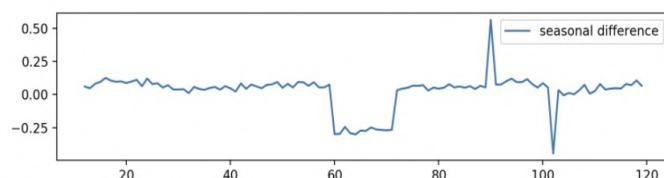
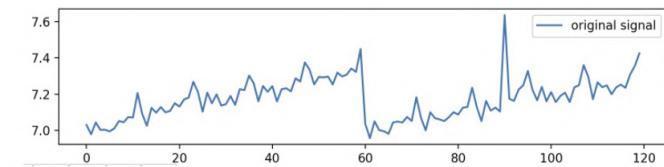
$$\sum_{t=T+1}^N |g_t - \sum_{i=0}^{T-1} \nabla \tau_{t-i}| + \lambda_1 \sum_{t=2}^N |\nabla \tau_t| + \lambda_2 \sum_{t=3}^N |\nabla^2 \tau_t|$$

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

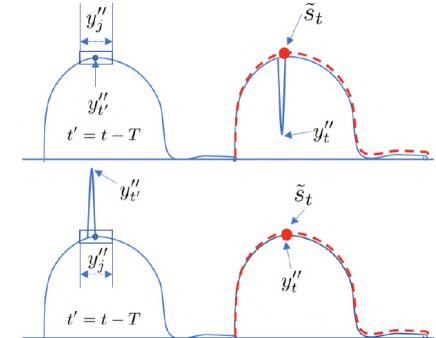
- Weighted non-local seasonal filtering
 - Consider K neighbourhoods for each season
 - Filter weights based on both value and position difference, similar to Bilateral filtering
 - Robust to *outlier* and adaptive to *seasonal shift*

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

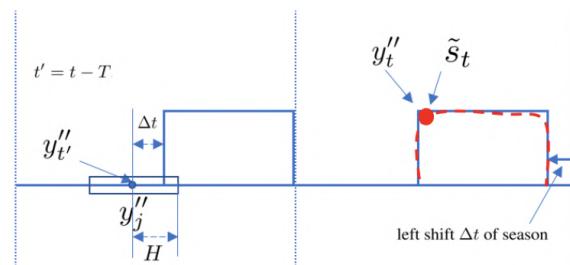
$$w_{(t'_i, j)}^t = e^{-\frac{|j-t'_i|^2}{2\delta_d^2} - \frac{|y''_j - y''_{t'_i}|^2}{2\delta_i^2}}$$

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

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



(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

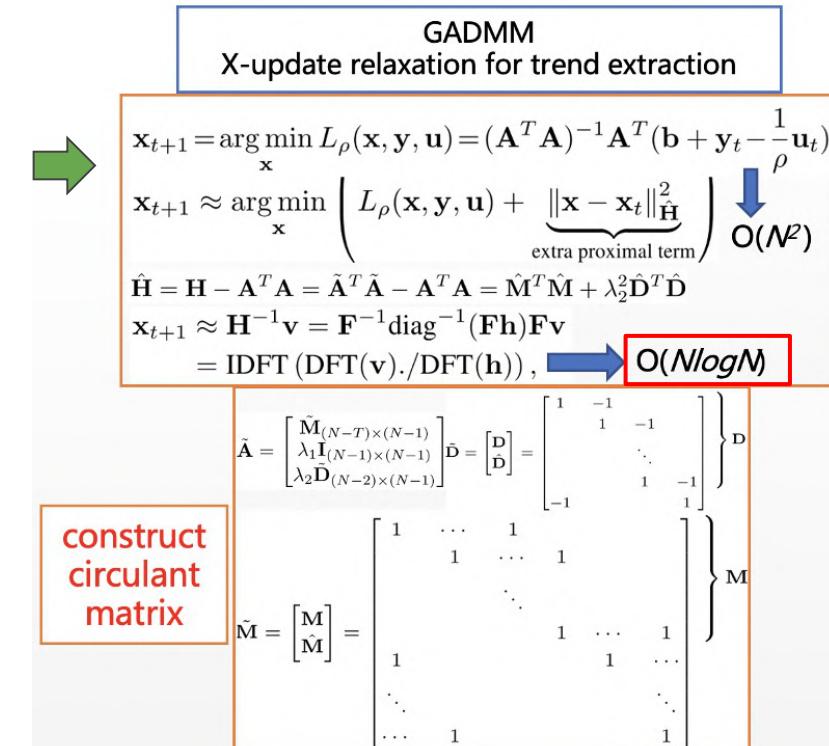
$$\mathbf{x}_{t+1} = \arg \min_{\mathbf{x}} L_\rho(\mathbf{x}, \mathbf{y}, \mathbf{u}) = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T (\mathbf{b} + \mathbf{y}_t - \frac{1}{\rho} \mathbf{u}_t) \quad O(N^2)$$

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

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

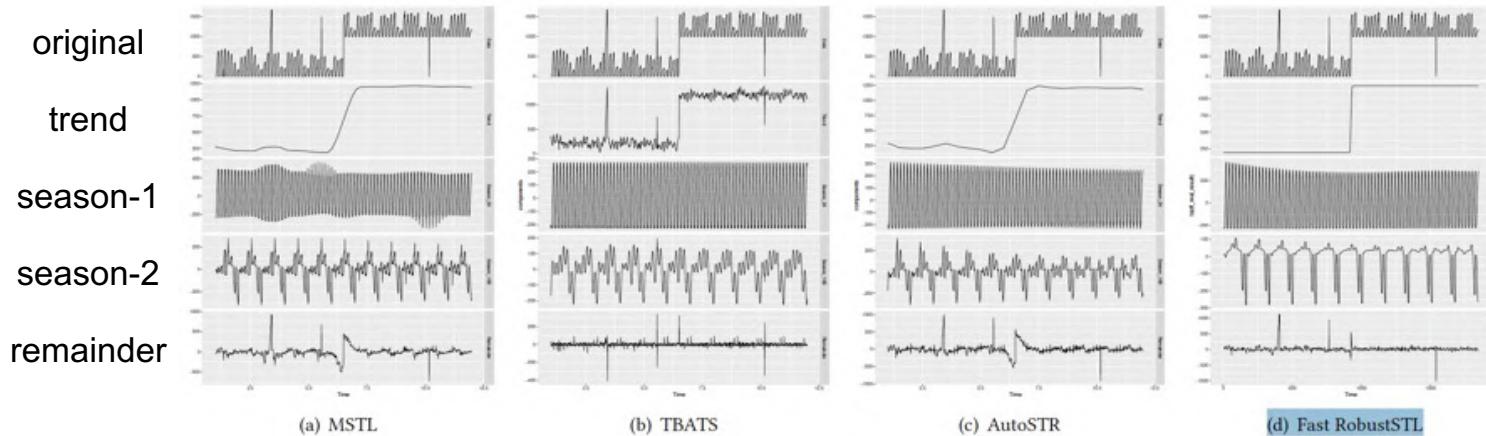
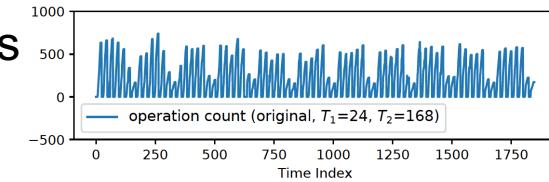
$$\mathbf{A} = \begin{bmatrix} \mathbf{M}_{(N-T) \times (N-1)} \\ \lambda_1 \mathbf{I}_{(N-1) \times (N-1)} \\ \lambda_2 \mathbf{D}_{(N-2) \times (N-1)} \end{bmatrix}, \mathbf{b} = \begin{bmatrix} \mathbf{g}_{(N-T) \times 1} \\ \mathbf{0}_{(2N-3) \times 1} \end{bmatrix}$$

- Similar GADMM can be formulated for multi-season decomposition



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

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

- Algorithm comparisons

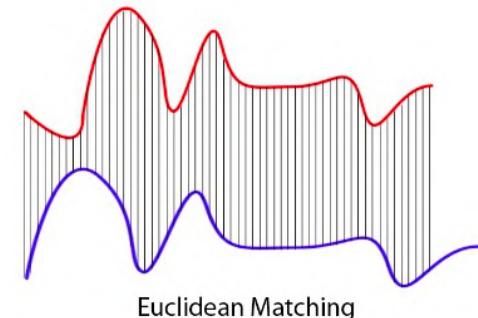
Algorithm	Robust to outlier	Robust to seasonal shift	Robust to trend changes	Support multiple periods	Efficient computation
ARIMA/SEATS	✗	✗	✗	✗	✗
SSA	✗	✗	✗	✗	✗
TBATS	✗	✗	✓	✓	✗
STL	✗	✗	✗	✗	✓
MSTL	✗	✗	✗	✓	✗
STR	✓	✓	✗	✓	✗
RobustSTL	✓	✓	✓	✓	✓

→ RobustSTL is more efficient than others

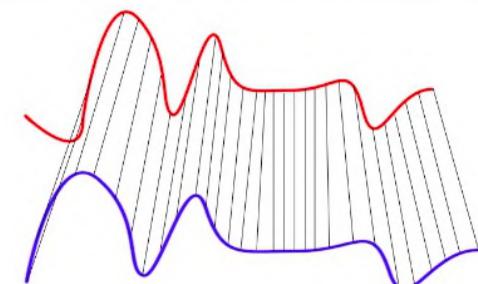


Similarity Measure between Time Series

- 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



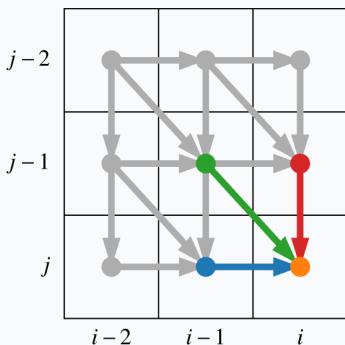
Dynamic Time Warping (DTW)

- Dynamic Time Warping in the optimization form:

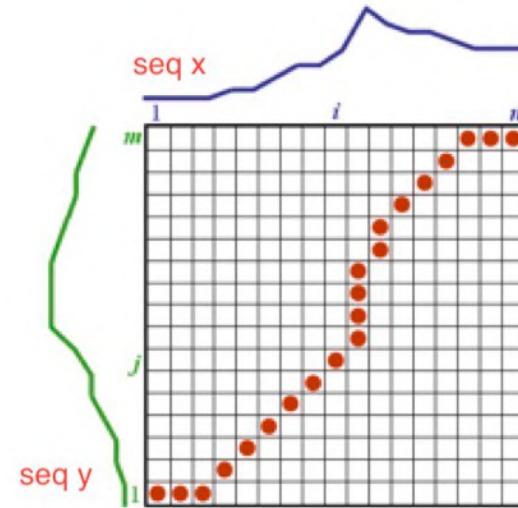
$$D(x, y) = \min \sum_{t=1}^T dist(x_t, y_{\phi(t)})$$

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

- How to compute DTW using dynamic programming?



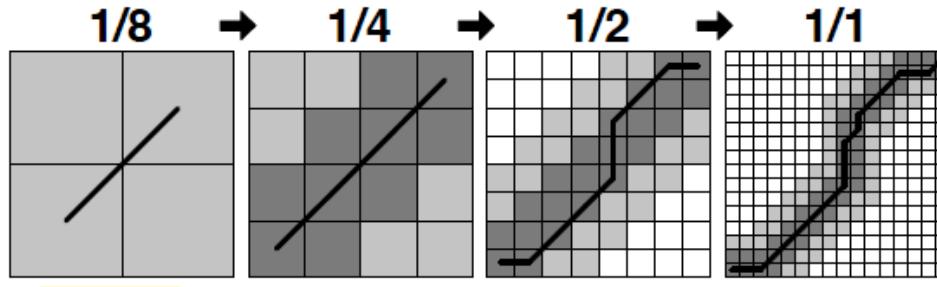
$D(i, j) = Dist(i, j) + \min[D(i - 1, j), D(i, j - 1), D(i - 1, j - 1)]$
where $D(i, j)$ is the DTW distance between $x_{1:i}$ and $y_{1:j}$
 $dist(i, j)$ is the distance between x_i and y_j



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

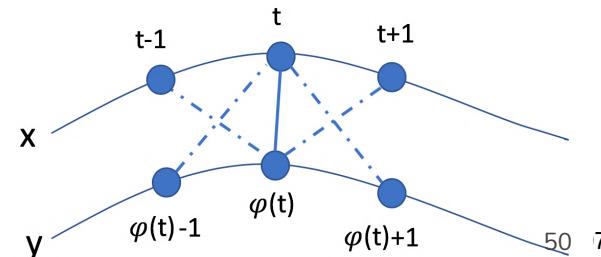
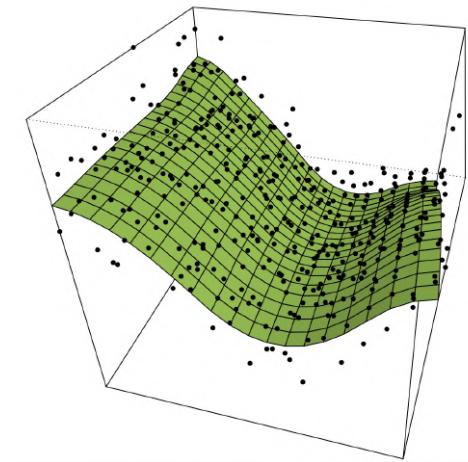
Graph Detrending for Trend Estimation

- 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 ϵ_i assumed to have zero mean. Our goal is to estimate μ_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



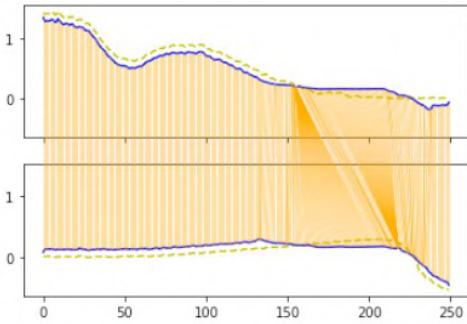
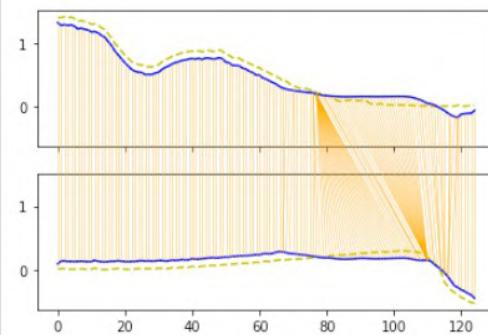
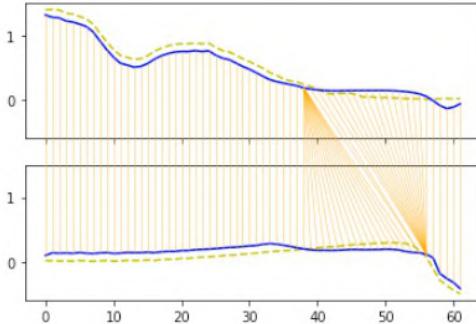


Detailed Procedure of RobustDTW

- 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



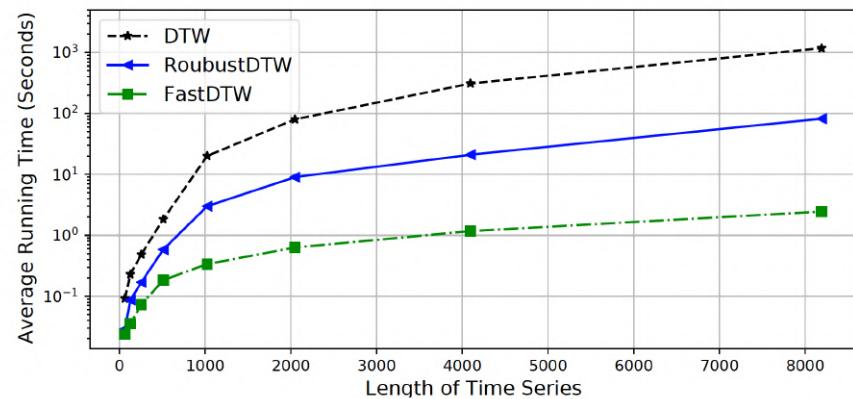


RobustDTW: Missing Value Processing and Efficiency

- Block missing data can be processed adequately in the low-resolution representation
- When mapping into the high-resolution, the corresponding segments are excluded in further computation

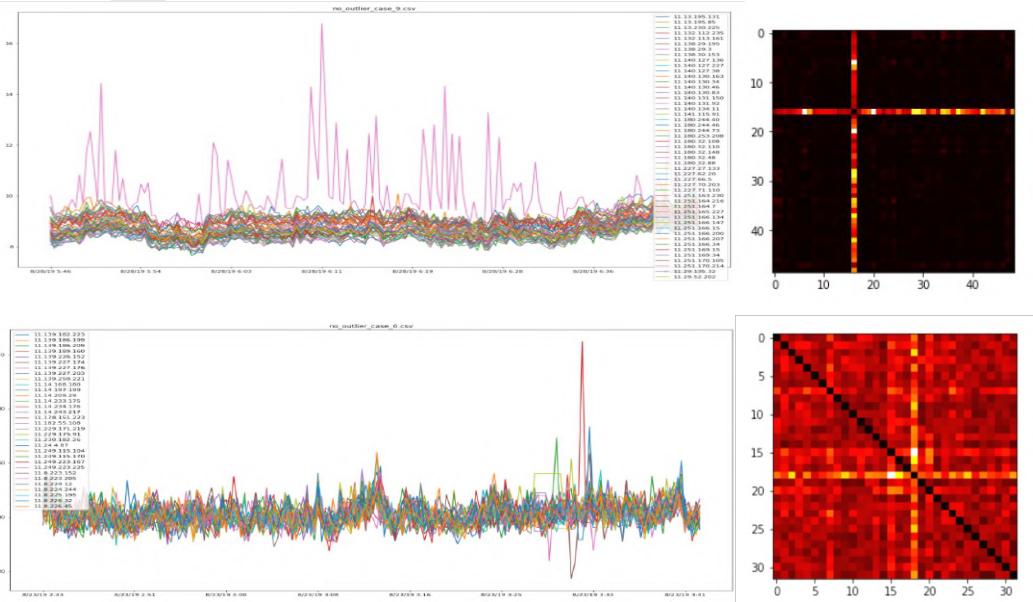


Running time comparison with different time series lengths





Application of RobustDTW: Time Series Outlier Detection



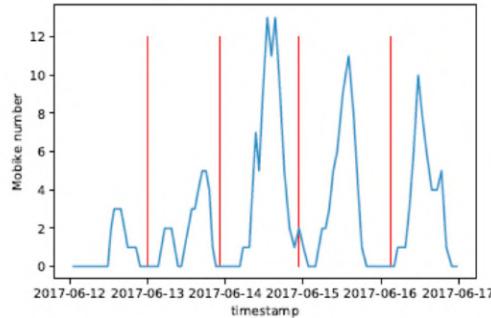
Comparison of AUC scores of outlier detection via local outlier factor (LOF) algorithm with different similarity measures and noise conditions

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

- 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	0.980	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
	RobustDTW	0.951	0.980



Outline

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

(10 minutes break)

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

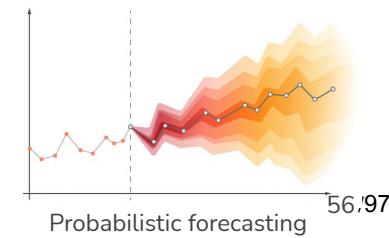
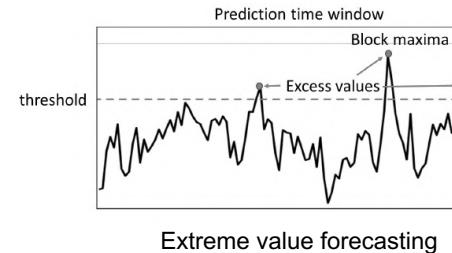
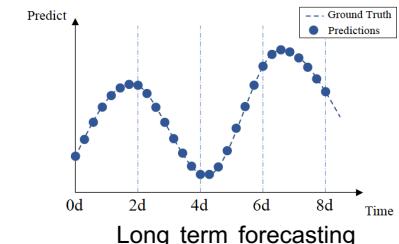
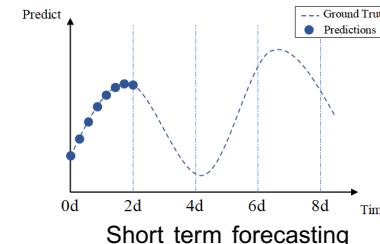
- **Forecasting:** Tree Models, Deep Ensemble, Transformers, etc.
- **Autoscaling (from Forecasting to Decision-Making):** Query Modeling, Scaling Decision, etc.
- **Anomaly Detection:** Decomposition Model, Deep State Space Model, Transformers, etc.
- **Fault Cause Localization (from Anomaly Detection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

(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

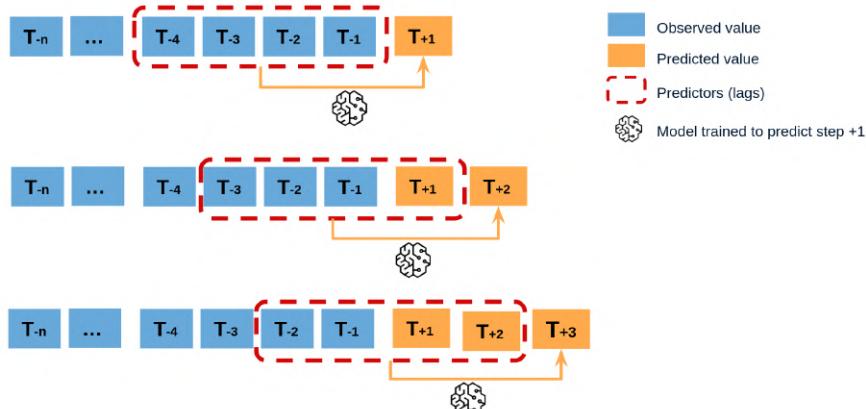




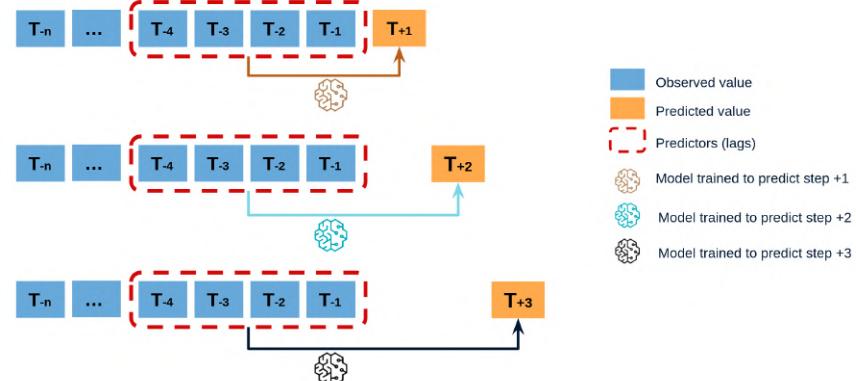
Forecasting: Ensemble (Tree Models)

- 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

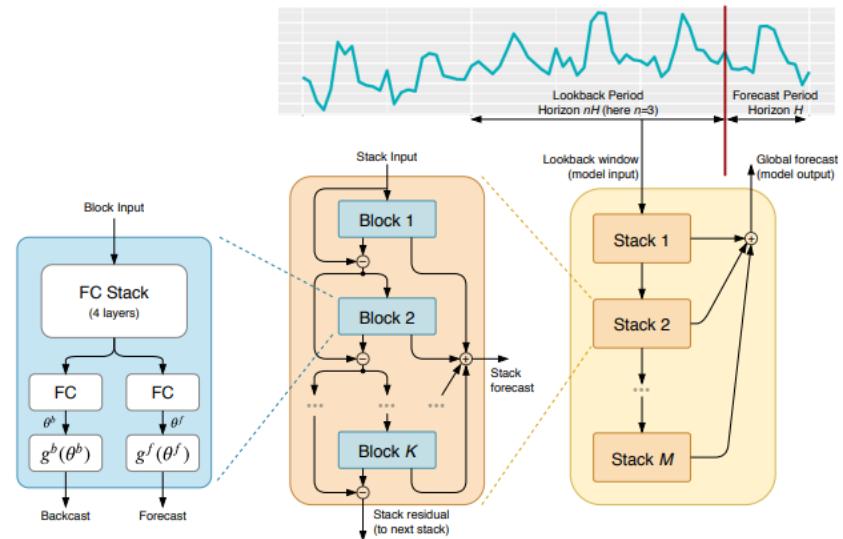




Forecasting: Deep Ensemble (MLP based Models)

● 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





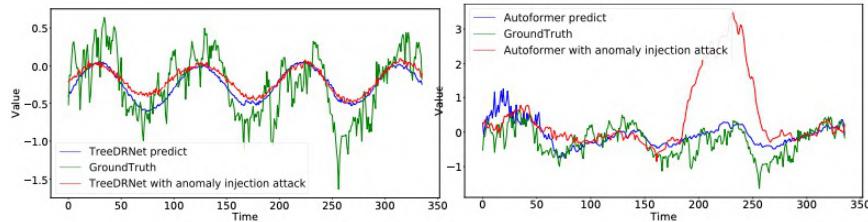
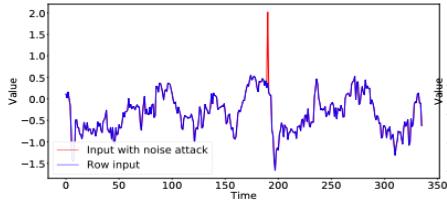
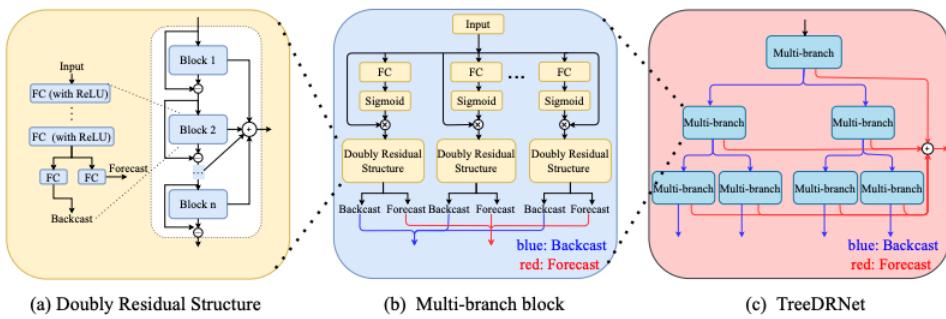
Forecasting: Deep Ensemble (MLP based Models)

- TreeDRNet

- Doubly residual link motivated by *robust* iterative reweighted algorithm
- Ensemble of models and *tree-based* aggregation motivated by Kolmogorov-Arnold representation theorem
- Robustness* against input anomalies/attacks

$$\delta y_k = y - \sum_{j=1}^{k-1} z_j, \quad \underbrace{x_k = f(x_{k-1}; \delta y_k) + x_{k-1}}_{\text{Backcast + Skip Link}}, \quad h(\mathbf{x}) = \sum_{k=0}^{2d} \Phi_k \left(\underbrace{\sum_{p=1}^m \phi_{k,p} \left(\underbrace{\mathbf{x} \circ \mathbf{m}_{k,p}}_{\text{Feature Selection}} \right)}_{\text{Model Ensemble}} \right),$$

$\Delta \beta_k = g(x_k, \delta y_k), \quad \underbrace{z_k = \langle \Delta \beta_k, x_k \rangle}_{\text{Forecast Link}}$

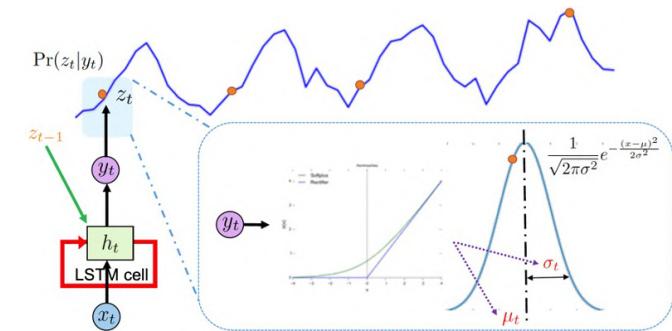
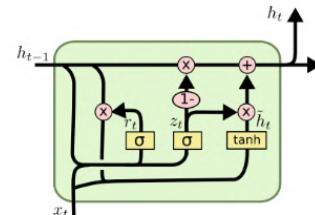
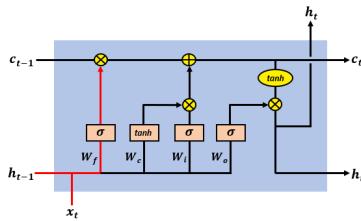
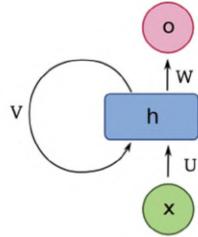




Forecasting: RNN based Models

- 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.

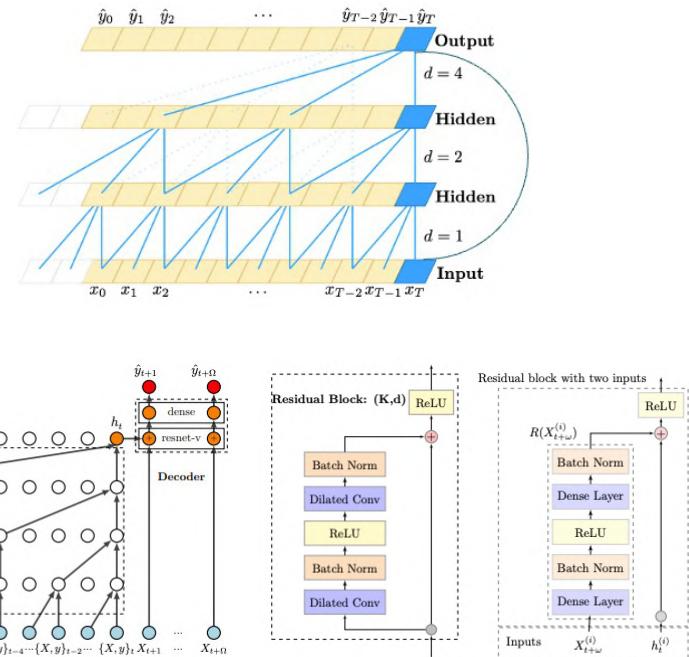
Chung, Junyoung, et al. "Empirical evaluation of gated recurrent neural networks on sequence modeling." arXiv preprint arXiv:1412.3555, 2014.

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



Forecasting: CNN based Models

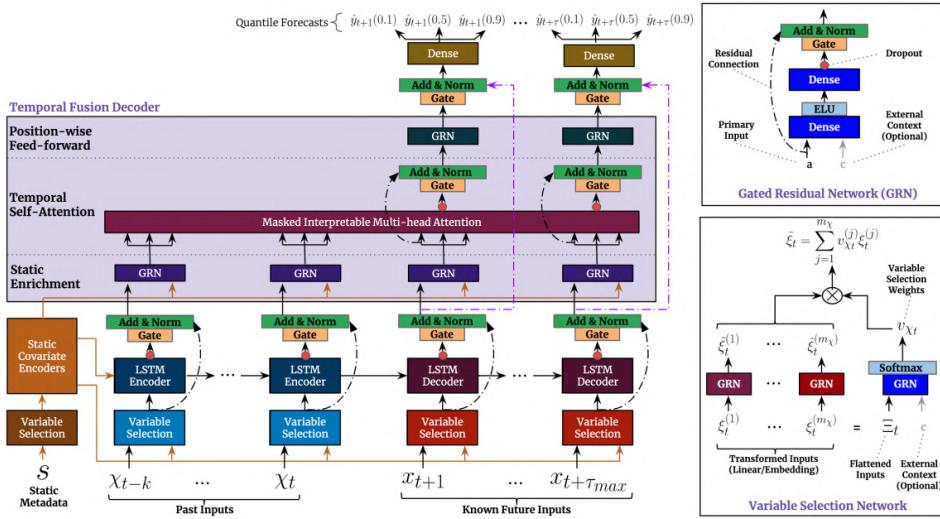
- TCN (temporal convolutional networks)
 - 1D dilated causal CNN
 - Dilated convolutional layers for large receptive field
 - Skip connection for raw info flow like ResNet
 - Empirically good performance and fast training
- DeepTCN
 - Encoder and decoder with TCN
 - Probabilistic forecasting under both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





Forecasting with Transformer: TFT

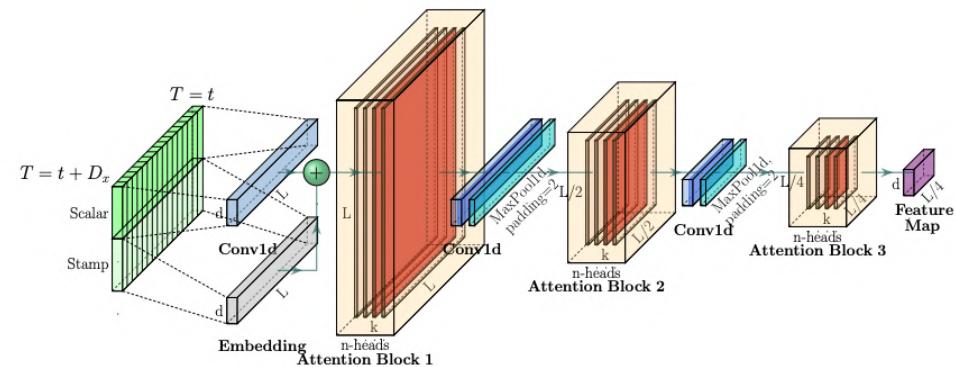
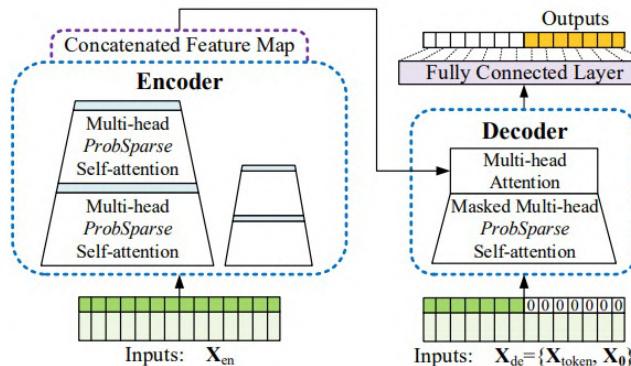
- TFT (temporal fusion Transformer)
 - Variable selection network: select most salient features
 - Gated residual network: efficient information flow with skip connections and gating layers
 - *Multi-scales network*: recurrent layers for local processing, interpretable self-attention layers for long-term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



Forecasting with Transformer: Informer

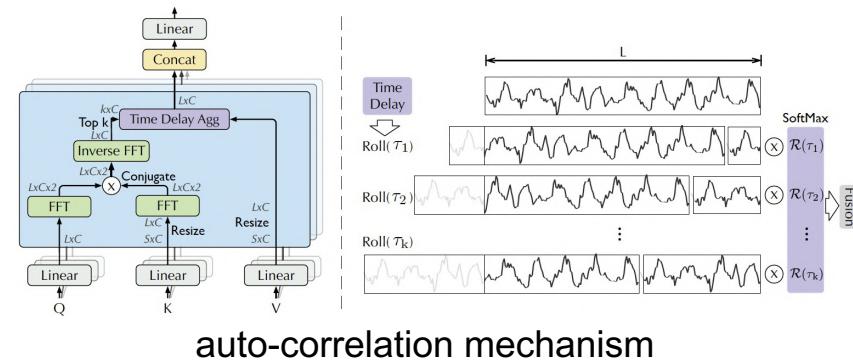
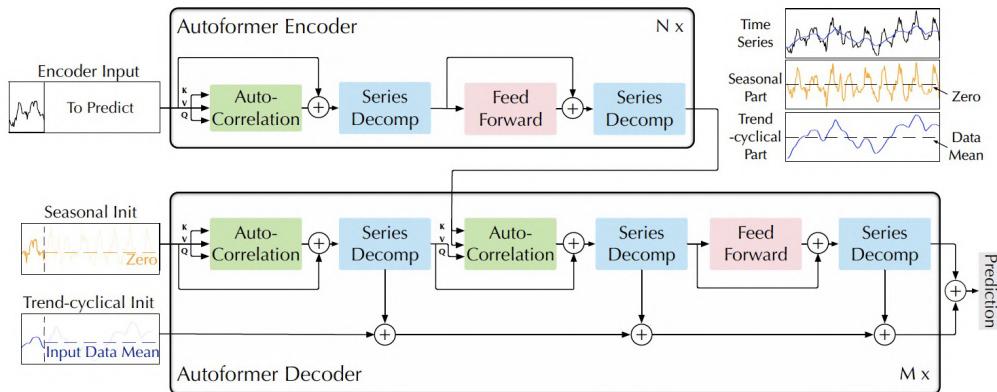
- **Informer**

- 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





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

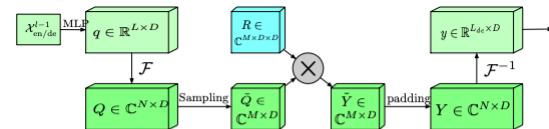
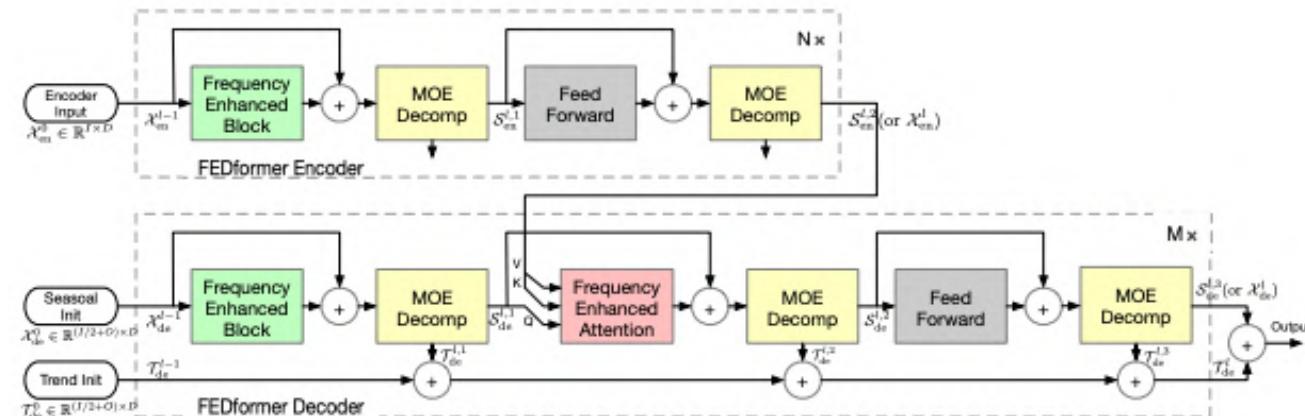


Figure 3. Frequency Enhanced Block with Fourier transform (FEB-f) structure.

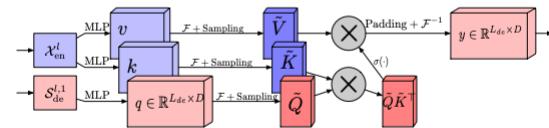


Figure 4. Frequency Enhanced Attention with Fourier transform (FEA-f) structure, $\sigma(\cdot)$ is the activation function.

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	0.203	0.269	0.325	0.421	0.193	0.201	0.214	0.246	0.148	0.271	0.460	1.195	0.587	0.604	0.621	0.626	0.217	0.276	0.339	0.403	3.228	2.679	2.622	2.857
	MAE	0.287	0.328	0.366	0.415	0.308	0.315	0.329	0.355	0.278	0.380	0.500	0.841	0.366	0.373	0.383	0.382	0.296	0.336	0.380	0.428	1.260	1.080	1.078	1.157
FEDformer-w	MSE	0.204	0.316	0.359	0.433	0.183	0.195	0.212	0.231	0.139	0.256	0.426	1.090	0.562	0.562	0.570	0.596	0.227	0.295	0.381	0.424	2.203	2.272	2.209	2.545
	MAE	0.288	0.363	0.387	0.432	0.297	0.308	0.313	0.343	0.276	0.369	0.464	0.800	0.349	0.346	0.323	0.368	0.304	0.363	0.416	0.434	0.963	0.976	0.981	1.061
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	1.023	0.423	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

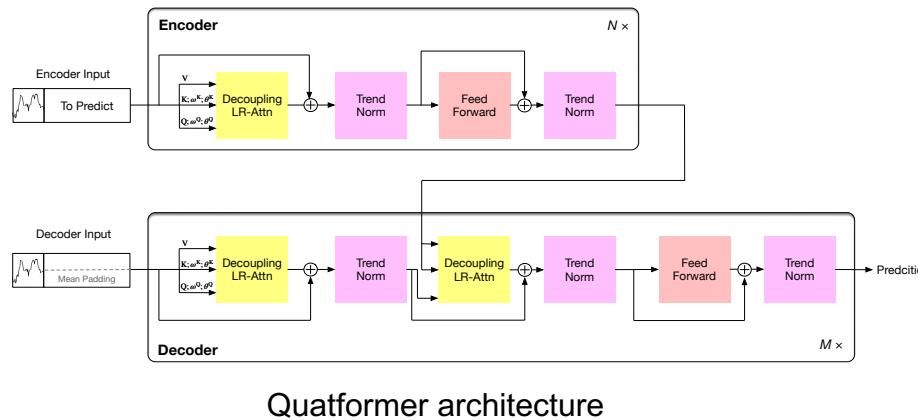
Table 1. Complexity analysis of different forecasting models.

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

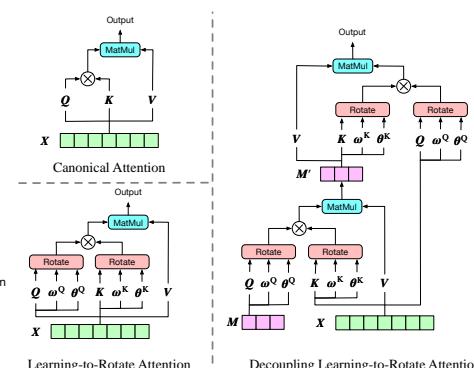


Forecasting with Transformer: Quatformer

- 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

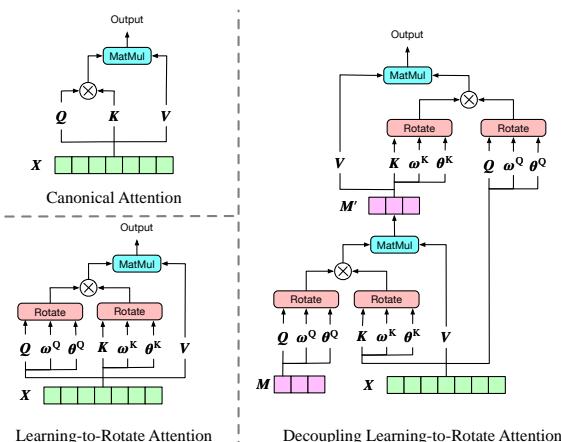
$$\begin{aligned} & \frac{\gamma}{\sigma} \odot (\mathcal{X} - \text{MovingAvg}(\mathcal{X})) + \mathcal{T}, \\ & \sigma = \sqrt{\frac{1}{N} \sum_{i=1}^N (\mathcal{X}_i - \mu)^2}, \quad \mu = \frac{1}{N} \sum_{i=1}^N \mathcal{X}_i, \\ & \mathcal{T} = \sum_{i=0}^p \beta_i \text{pos}^i, \quad \text{pos}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

Trend normalization



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 $O(2cN)$.

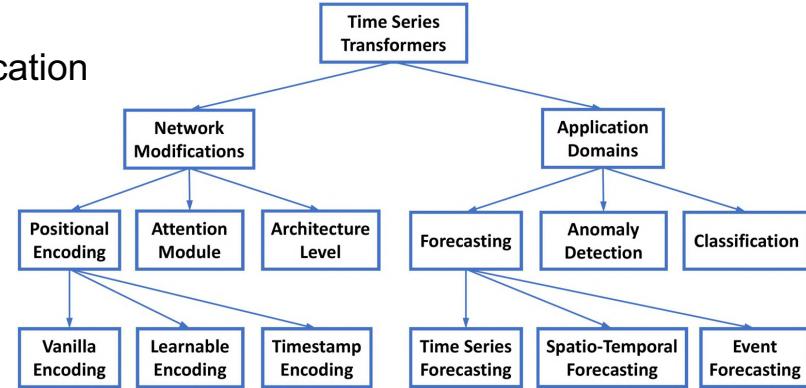
Performance of Quatformer

Models	Quatformer		Quatformer [†]		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
ETThm	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



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

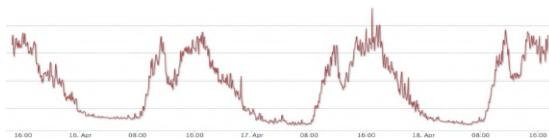




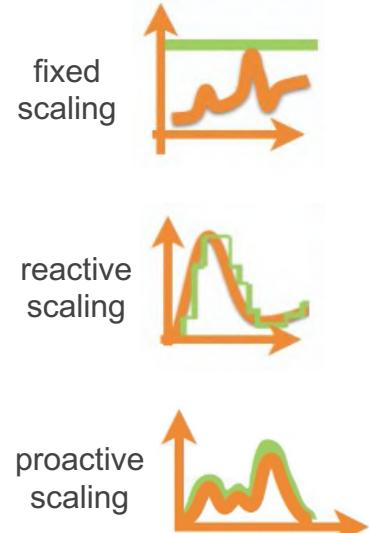
From Forecasting to *Decision-making*: Autoscaling

- 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



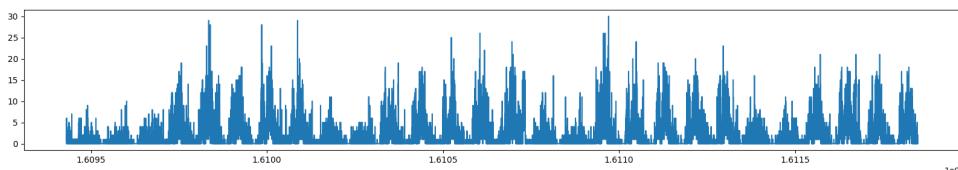
nonperiodic & random scenario
less potential for proactive autoscaling



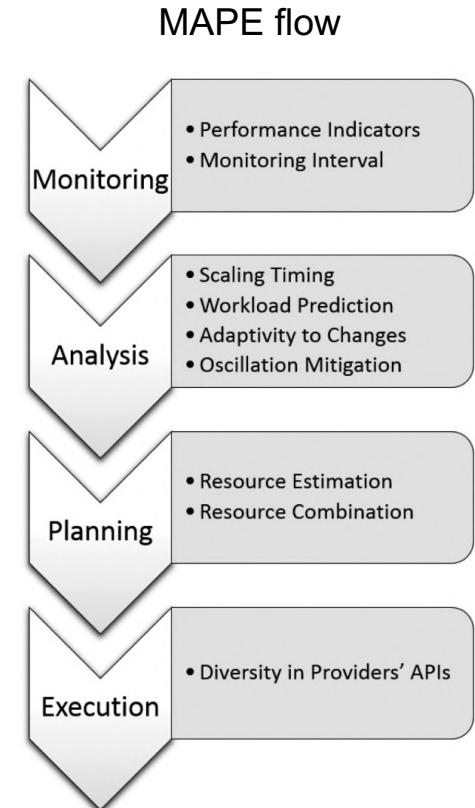


Autoscaling

- 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



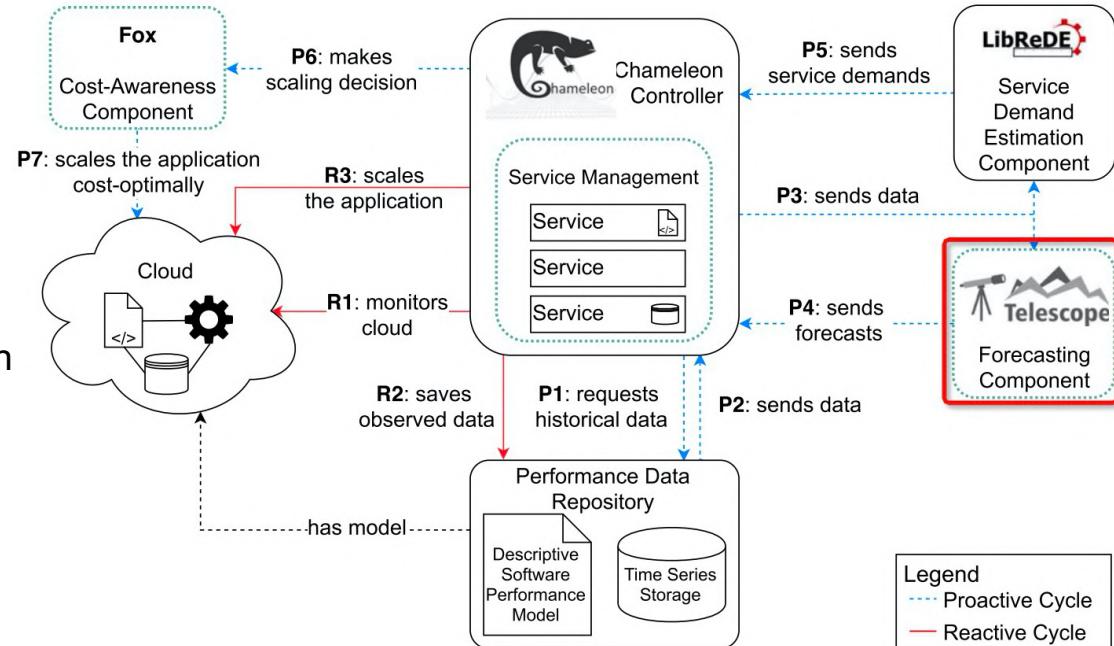
real-world QPS series from Alibaba container registry service



Autoscaling: Chamulteon

- System:

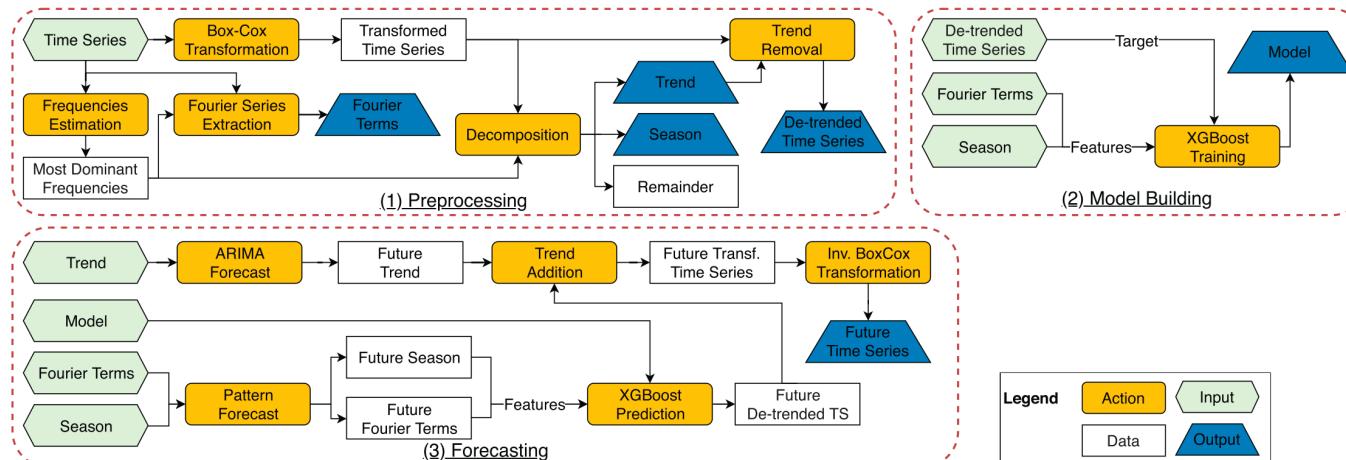
- Reactive
- Proactive
 - Monitor
 - Forecast
 - Demand estimation
 - Decision
 - Optimization





Autoscaling: Chamalteon

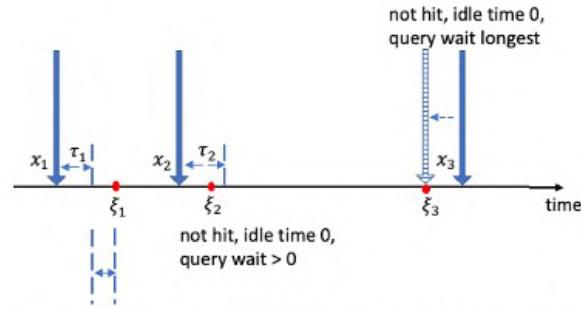
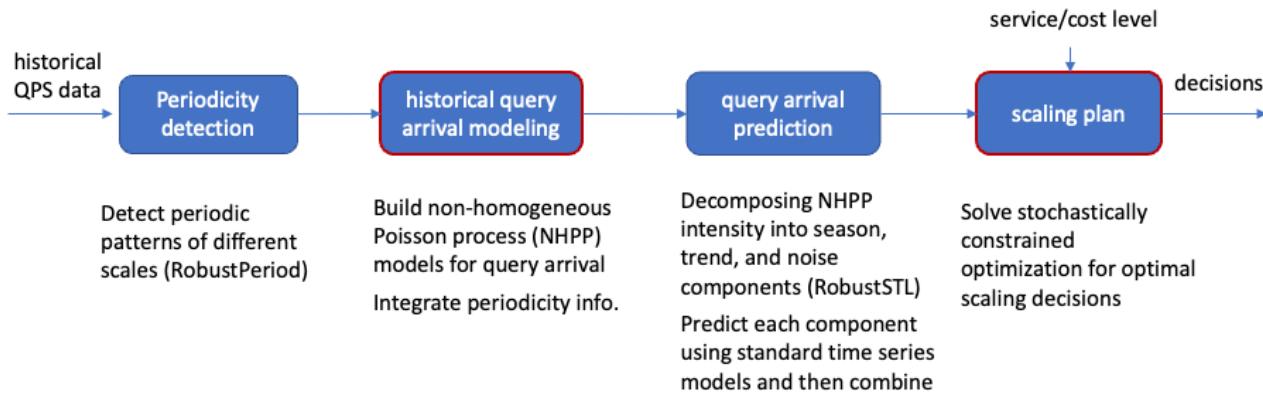
- Forecasting: Telescope
 - Designed for *seasonal* time series forecasting
 - Periodicity detection, seasonal-trend decomposition
 - Forecasting trend by ARIMA, seasonality by XGBoost





Autoscaling: RobustScaler

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



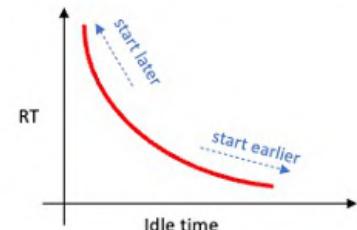
$x_1 < x_2 < \dots < x_l < \dots$ (Blue arrows): time of start of each instance
 $\tau_1, \tau_2, \dots, \tau_l, \dots$: startup delay of each instance
 $\xi_1 < \xi_2 < \dots < \xi_i < \dots$ (Red dots): time of arrival of each query

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





RobustScaler: Arrival Modeling as NHPP

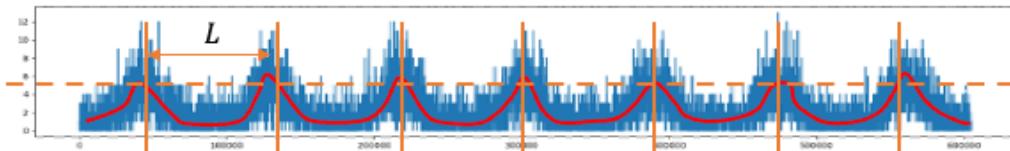
- Non-homogeneous Poisson process + periodicity regularization for complex periodic pattern

Input:

$Q = (Q_t, t = 1, \dots, T)$: query count within each time step Δt

Δt : length of each time step

L : known period length (from periodicity detection)



effect of periodicity penalty

Output:

$r = (r_t, t = 1, \dots, T)$: $\log(\text{intensity})$ for each time step, assumed piece-wise constant

Poisson likelihood for each interval t : $\frac{\exp(-\lambda_t \Delta t) (\lambda_t \Delta t)^{Q_t}}{Q_t!}$ with $\lambda_t = \exp(r_t)$

$$\min_r -\underbrace{\mathbf{Q}^T \mathbf{r} + \Delta t \cdot \mathbf{1}^T \exp(\mathbf{r})}_{\text{log likelihood}} + \underbrace{\beta_1 \|D^2 \mathbf{r}\|_1}_{\text{smoothness penalty}} + \underbrace{\frac{\beta_2}{2} \|D_L \mathbf{r}\|_2}_{\text{periodicity penalty}}$$

$$D^2 = \begin{bmatrix} 1 & -2 & 1 & 0 & \cdots & \cdots & \cdots & 0 \\ 0 & 1 & -2 & 1 & 0 & \cdots & \cdots & 0 \\ \vdots & \ddots & \ddots & \ddots & \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & \cdots & 0 & 1 & -2 & 1 \end{bmatrix} \in \mathbb{R}^{(T-2) \times T}$$

L entries apart

$$D_L = \begin{bmatrix} 1 & 0 & \cdots & 0 & -1 & \cdots & \cdots & 0 \\ 0 & 1 & 0 & \cdots & 0 & -1 & \cdots & 0 \\ \vdots & \ddots & \ddots & \ddots & \vdots & \vdots & \ddots & 0 \\ 0 & \cdots & \cdots & 1 & 0 & \cdots & 0 & -1 \end{bmatrix} \in \mathbb{R}^{(T-L) \times T}$$



RobustScaler: Constrained Optimization Model

- Optimize cost (i.e., idle time) while QoS (i.e., hitting probability (HP) or response time (RT)) is controlled, or vice versa.

- RobustScaler-HP: Minimize average **idle time** while controlling **HP** above a threshold $1 - \alpha$:

$$\begin{aligned} & \min_{x_1, \dots, x_K \geq 0} \sum_{i=1}^K E[(\xi_i - \tau_i - x_i)_+], \quad i = 1, \dots, K \\ & \text{s.t. } P(\xi_i > x_i + \tau_i) \geq 1 - \alpha \end{aligned} \quad \longrightarrow \quad x_i^* = \alpha \text{ quantile of } (\xi_i - \tau_i)$$

- RobustScaler-RT: Minimize average **idle time** while controlling average **RT** below a threshold d :

$$\begin{aligned} & \min_{x_1, \dots, x_K \geq 0} \sum_{i=1}^K E[(\xi_i - \tau_i - x_i)_+] \\ & \text{s.t. } \mu_s + E[(\tau_i - (\xi_i - x_i)_+)_+] \leq d, \quad i = 1, \dots, K \end{aligned} \quad \longrightarrow \quad x_i^* \text{ such that } E[(\tau_i - (\xi_i - x_i^*)_+)_+] = d - \mu_s$$

- RobustScaler-cost: Minimize average **RT** while controlling average **idle time** below a threshold B :

$$\begin{aligned} & \min_{x_1, \dots, x_K \geq 0} \sum_{i=1}^K E[(\tau_i - (\xi_i - x_i)_+)_+] \\ & \text{s.t. } E[(\xi_i - \tau_i - x_i)_+] \leq B, \quad i = 1, \dots, K \end{aligned} \quad \longrightarrow \quad \begin{cases} x_i^* = 0, & \text{if } E[(\xi_i - \tau_i)_+] \leq B \\ x_i^* \text{ s.t. } E[(\xi_i - \tau_i - x_i^*)_+] = B & \text{otherwise} \end{cases}$$

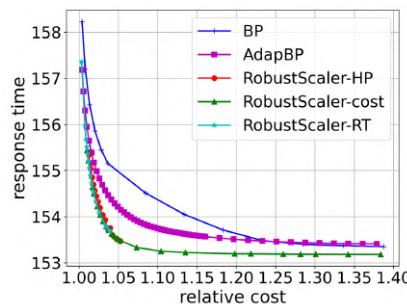
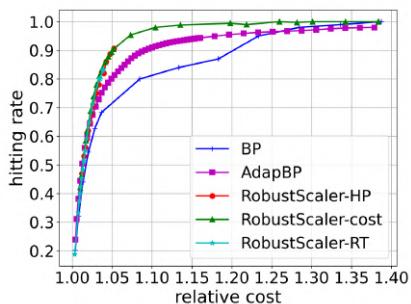


RobustScaler: QoS-cost Pareto Plots

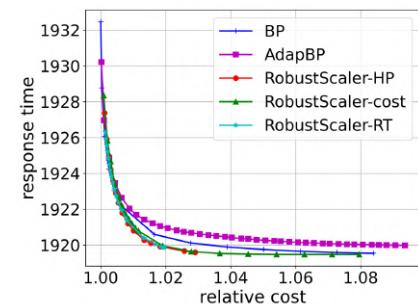
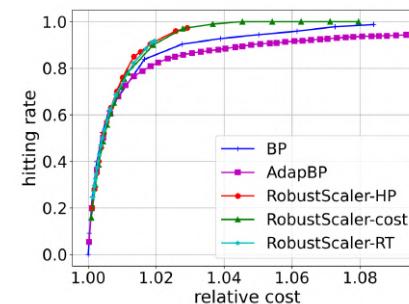
- Baseline autoscaling methods
- **Backup Pool (BP)**: A warm instance pool of fixed size B is maintained throughout, whenever one of the instances is consumed, fill with a new one immediately. ($B=0$ is purely reactive)
- **Adaptive Backup Pool (AdapBP)**: Adaptive version of BP with $B=\text{coeff}^* \text{QPS}$ for some fixed coefficient coeff .

	#queries	duration
Alibaba cluster 2018 trace (https://github.com/alibaba/clusterdata)	503,850	5 days
Google cluster 2019 trace (https://github.com/google/cluster-data)	20,254	24 hours

QoS-cost Pareto plots on Alibaba trace



QoS-cost Pareto plots on Google trace

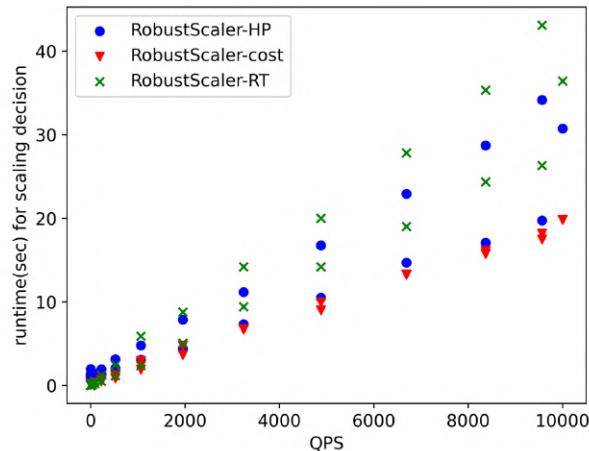
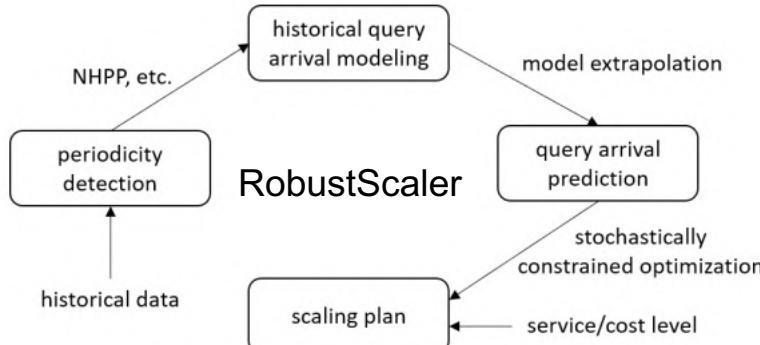


Under the same relative cost, higher hitting rate or lower response time is better.
RobustScaler is better than others (AdapBP, BP)



RobustScaler: Scalability

- Scalability by modules
 - Periodicity detection & query arrival prediction: RobustPeriod & RobustSTL known to be scalable
 - Historical query arrival modeling: 7 secs for QPS series of four days, 100 secs for QPS series of three weeks
 - Scaling plan: linearly growing runtime, remain efficient under high QPS





Time Series Anomaly Detection: Background

- Time-series anomalies

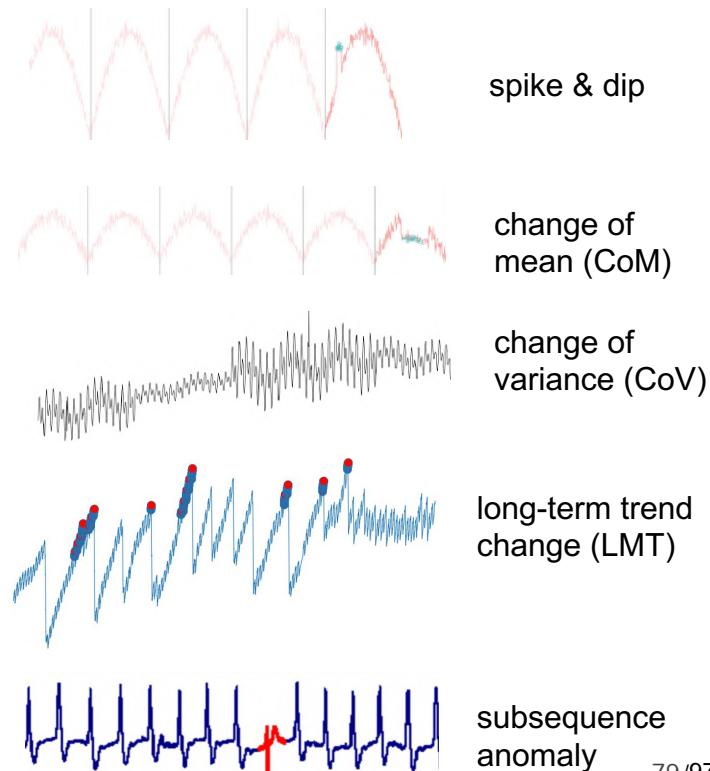
- Points/sequences that significantly deviate from the normal pattern
- **Differences** from static data anomaly detection:
 - Anomalies have temporal context
 - Noise is non-stationary (i.e., time-varying noise)
 - Lack of anomaly labels

- Types of anomalies

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

- Models

- *Traditional methods, deep models*



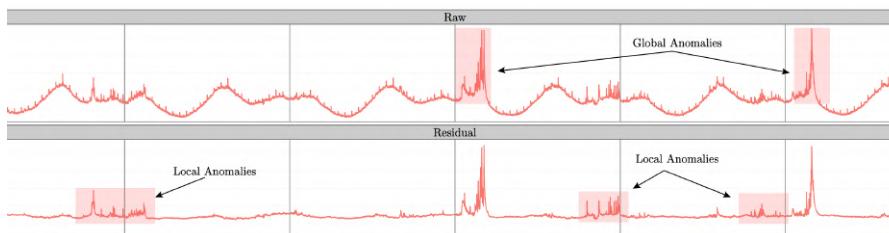
subsequence
anomaly



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:

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

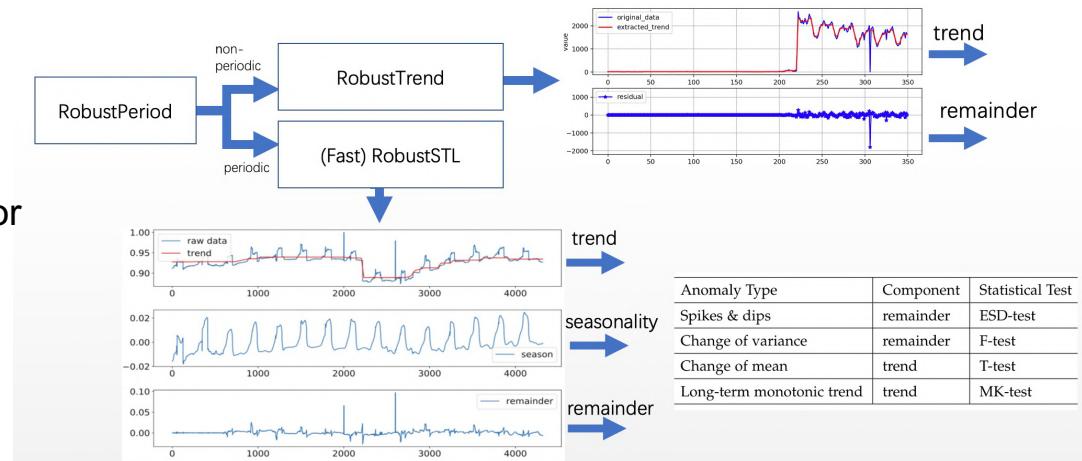
1. Extract seasonal component S_X using STL Variant
2. Compute median \tilde{X}
- /* Compute residual */
3. $R_X = X - S_X - \tilde{X}$
- /* Detect anomalies vector X_A using ESD */
4. $X_A = \text{ESD}(R, k)$

return X_A

Anomaly Detection: Decomposition based Model

- Robust Decomposition + Statistics

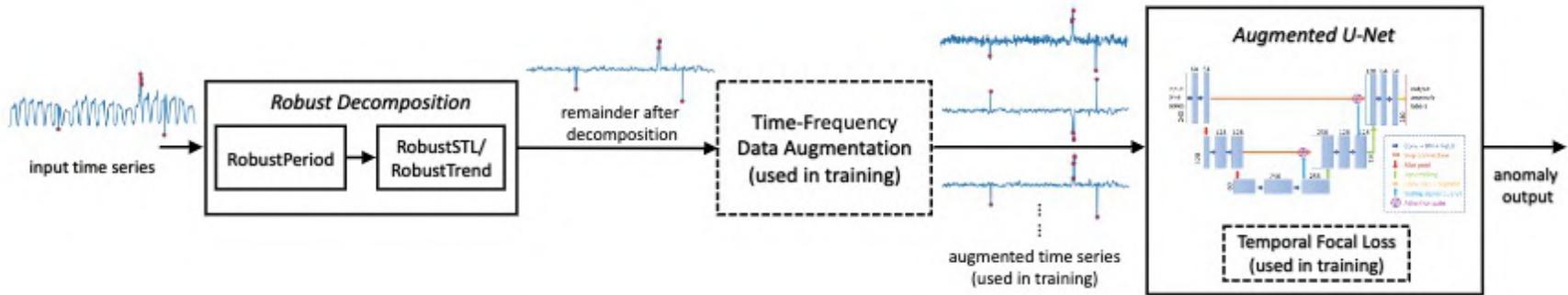
- Robust time series decompositions
reduce complexity and bring
explainability
- Robust statistical tests on each
components lead to high accuracy for
different types of anomalies





Anomaly Detection: RobustTAD

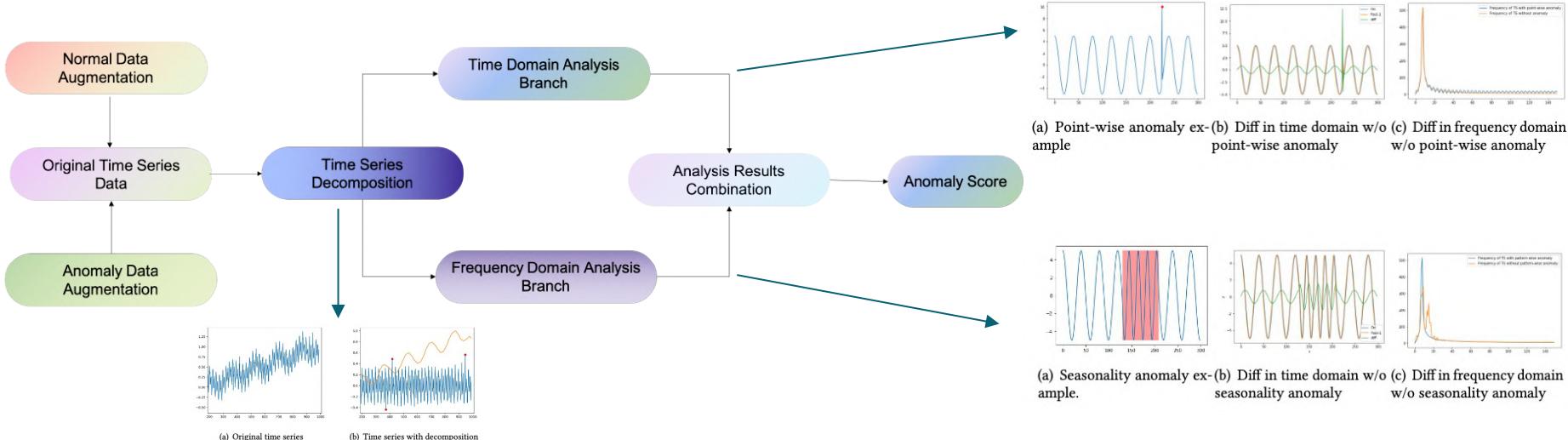
- Decomposition + Data Augmentation + U-Net (deep model)
 - Robust decomposition: handle complicated patterns, and simplify neural network
 - Data augmentation: mitigate the effects of limited labeled data
 - U-Net: capture multi-scale information of time series
 - Temporal focal loss: label-based weight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





Anomaly Detection: TFAD

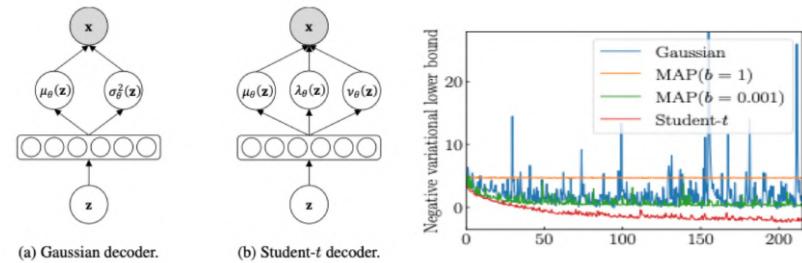
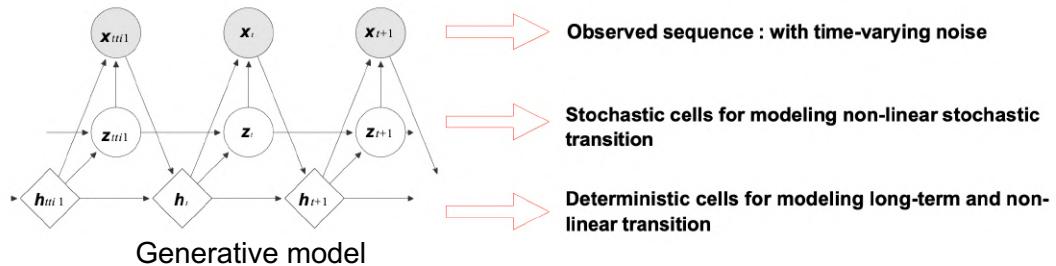
- Data Augmentation + Decomposition + TCN + *Time-Frequency Processing*
 - Anomalies types: seasonal anomaly, trend anomaly, global/context point anomaly
 - Point anomaly easier to detect in **Time**; seasonal anomaly easier to detect in **Frequency**



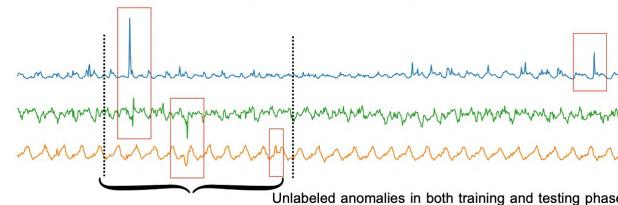


Anomaly Detection: RDSSM

- 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

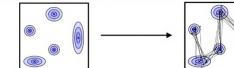


A. Slide window over noisy multi-variate time series contaminated with unlabeled anomalies



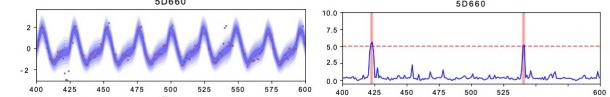
B. Get time series frames $t+1, \dots, t+T$

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

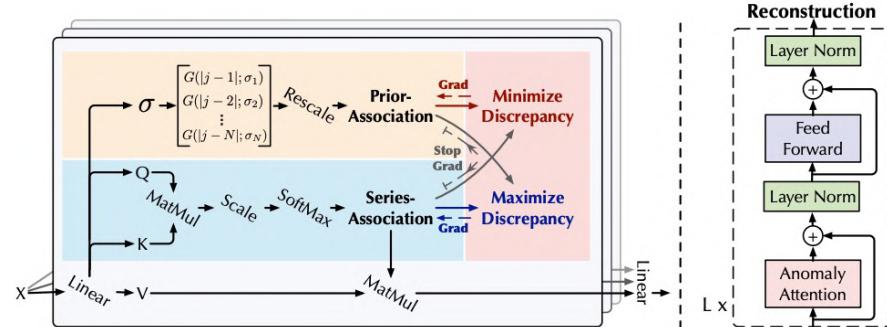
- Architecture: Anomaly-Attention

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

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

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

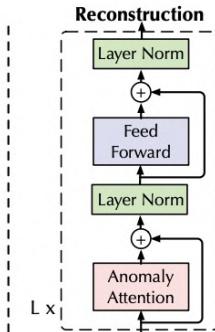
Adjacent-concentration inductive bias
Find the most effective associations



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

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

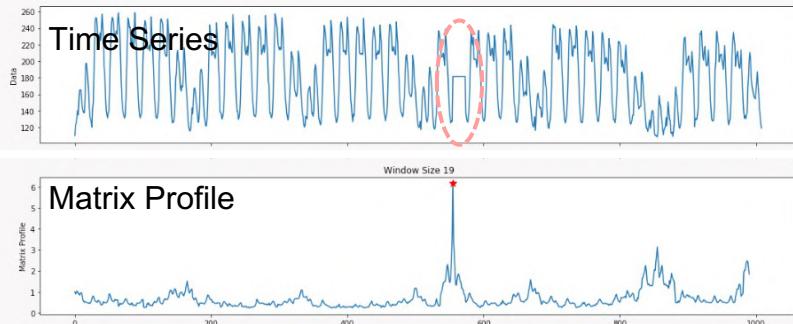
$$\mathcal{L}_{\text{Total}}(\hat{\mathcal{X}}, \mathcal{P}, \mathcal{S}, \lambda; \mathcal{X}) = \|\mathcal{X} - \hat{\mathcal{X}}\|_F^2 - \lambda \times \|\text{AssDis}(\mathcal{P}, \mathcal{S}; \mathcal{X})\|_1$$



Subsequence Anomaly Detection

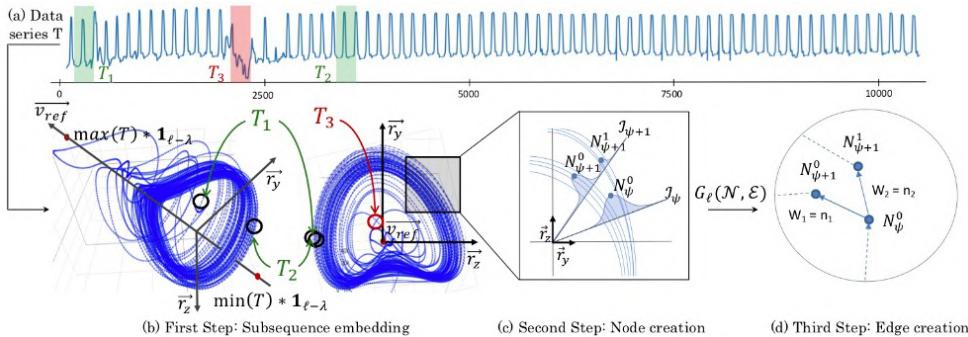
- Matrix Profile (MP)

- 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



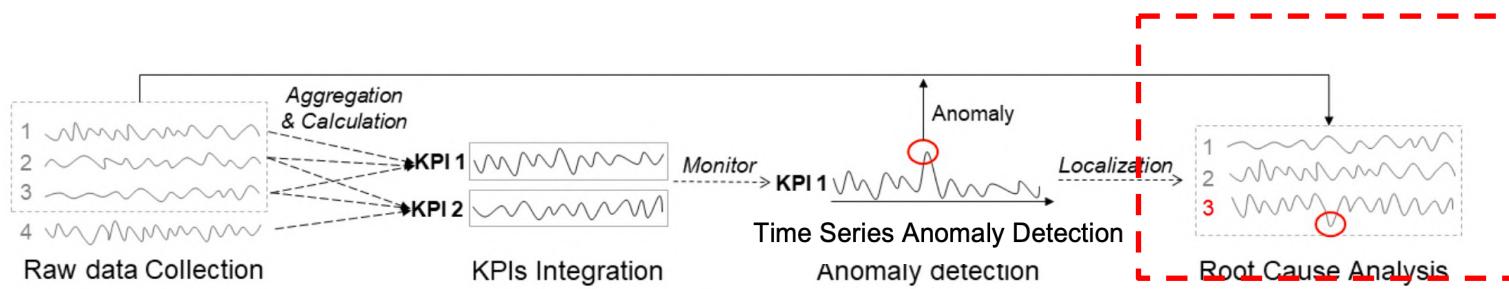
- Series2Graph

- 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.

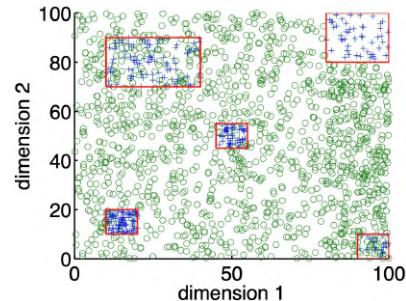
[4] Li, Zeyan, et al. "Generic and robust localization of multi-dimensional root causes." 2019 IEEE 30th International Symposium on Software Reliability Engineering (ISSRE). IEEE, 2019.

[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

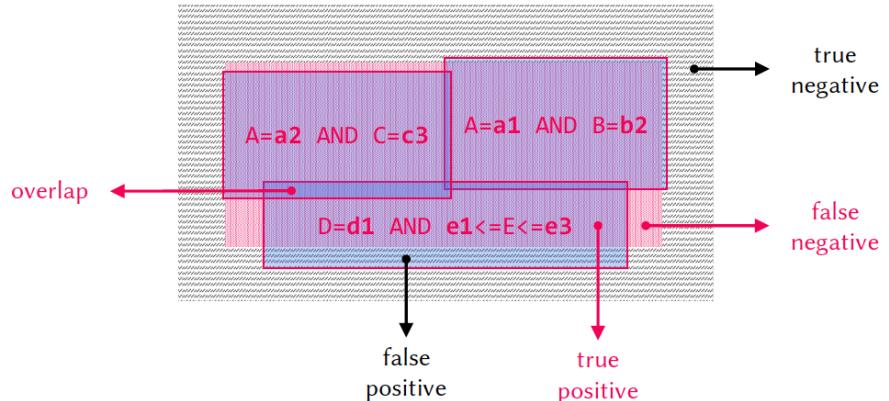
The rule set learning methods aim to find discriminative patterns that only covers the samples of interest class



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

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

Minorize-Maximization

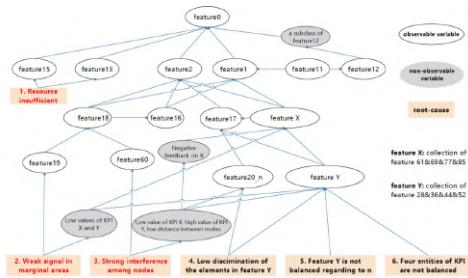
Linearization

• 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

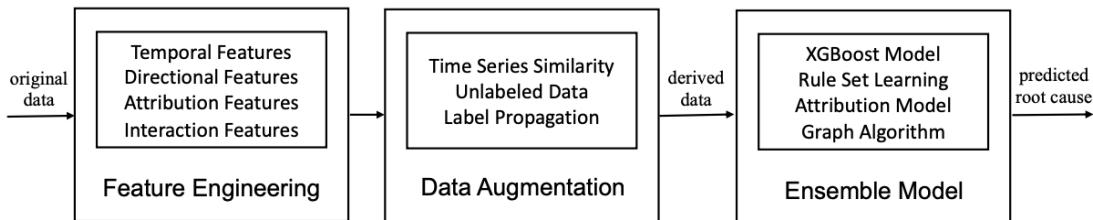


NetRCA



- 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



Rank	Team	Prize	Award
1	DMLAB	\$5000	Gold Prize: \$5000
2	Team Modulus		Paper link: https://arxiv.org/abs/2204.08003
3			Code link: https://github.com/Modulus-AI/NetRCA
4			Silver Prize: \$2000
5	Team DMLAB		Paper link: https://arxiv.org/abs/2204.08003
6			Code link: https://github.com/Modulus-AI/NetRCA
7			Bronze Prize: \$1000
8	Team 000000		Paper link: https://arxiv.org/abs/2204.08003
9			Code link: https://github.com/Modulus-AI/NetRCA
10	Accepted ICASSP SPQC Papers		
11	Team name: 000000	Title: ACCURATE INFERENCE OF UNKNOWN COMBINATIONS OF MULTIPLE ROOTCAUSES WITH CLASSIFIER ENSEMBLE	Link: https://arxiv.org/abs/2204.08003
12	Team name: DMLAB	Title: CAUSAL ALIGNMENT BASED ROOT CAUSE LOCALIZATION FOR WIRELESS NETWORK	Link: https://arxiv.org/abs/2204.08003
13	Team name: Modulus	Title: NetRCA: An Effective Network Fault Cause Localization Algorithm	Link: https://arxiv.org/abs/2204.08003



Tutorial Summary

- ❑ **Introduction:** Real-world Challenges and Needs for Robust Time Series Processing
- ❑ **Preliminaries:** Multiple Disciplines: Statistics, Signal Processing, Optimization, Deep Learning
- ❑ **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



Thanks!

Q&A

Tutorial Website: <https://github.com/DAMO-DI-ML/KDD2022-Tutorial-Time-Series>

Alibaba DAMO Academy - Decision Intelligence Lab: <https://damo.alibaba.com/labs/decision-intelligence>

Hiring: Time Series, XAI

*Research Interns, Postdocs, and
Research & Applied Scientists*

Seattle (US), Hangzhou (China)

Email CV to

qingsong.wen@alibaba-inc.com
liang.sun@alibaba-inc.com