



# Explaining Time Series via Contrastive and Locally Sparse Perturbations

Zichuan Liu<sup>1,2</sup> Yingying Zhang<sup>2</sup> Tianchun Wang<sup>3</sup> Zefan Wang<sup>2,4</sup> Dongsheng Luo<sup>5</sup>  
Mengnan Du<sup>6</sup> Min Wu<sup>7</sup> Yi Wang<sup>8</sup> Chunlin Chen<sup>1</sup> Lunting Fan<sup>2</sup> Qingsong Wen<sup>2</sup>

<sup>1</sup>Nanjing University, <sup>2</sup>Ailibaba Group,

<sup>3</sup>Pennsylvania State University, <sup>4</sup>Tsinghua University,

<sup>5</sup>Florida International University,

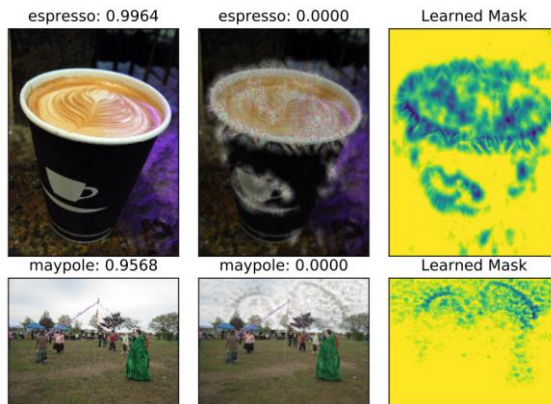
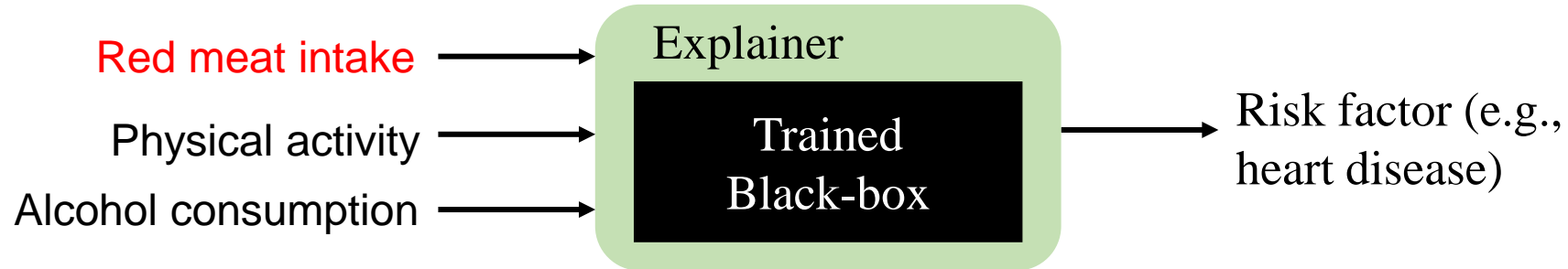
<sup>6</sup>New Jersey Institute of Technology,

<sup>7</sup>A\*STAR, <sup>8</sup>The University of Hong Kong



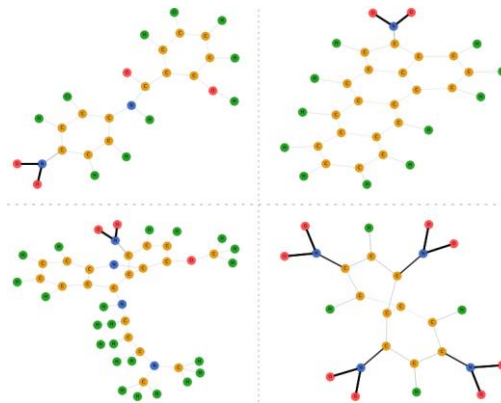
# Background

Black-box models with post-hoc explanation techniques: *Find salient features!*



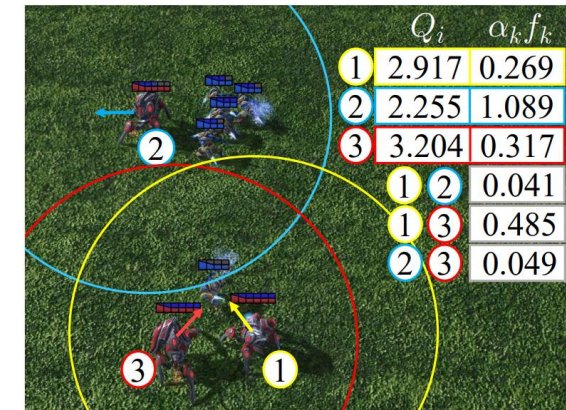
Visual Explanation

Source: [Fong et al.](#)



Graph Explanation

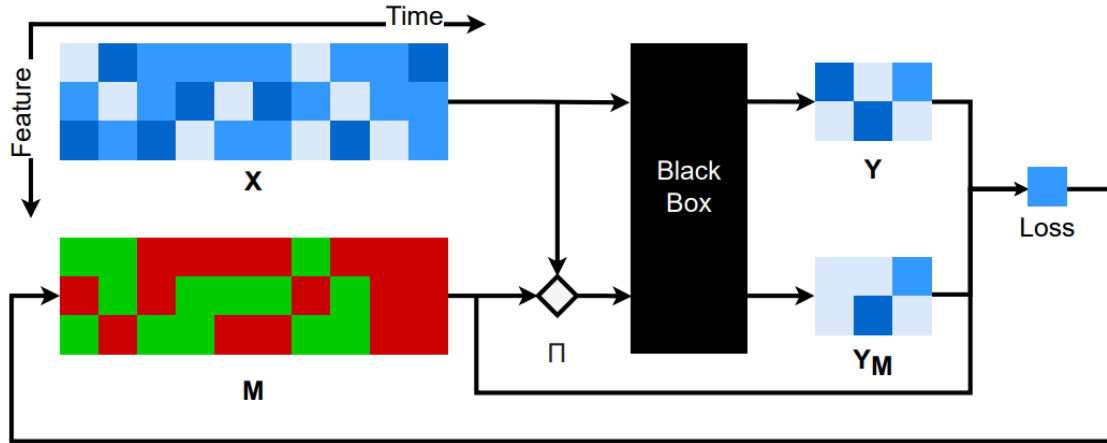
Source: [Miao et al.](#)



Game Explanation

Source: [Liu et al.](#)

# Challenges for Explaining Time Series



Dynamask, [Crabbe' et al.](#)

$$\Phi(x, m) = m \times x + (1 - m) \times u$$

$$\arg \min \underbrace{\mathcal{L}(f(x), f \circ \Phi(x, m))}_{\text{Predictive consistency}} + \underbrace{\mathcal{R}(m)}_{\text{regular}} + \underbrace{\mathcal{A}(m)}_{\text{smooth}}$$

## ➤ Fail to interpret visually

- Dense salient features (unlike the image and text)
- Noisy samples in time series

## ➤ Hard find temporal patterns

- The time series is smoothed

## ➤ Perturbations matter

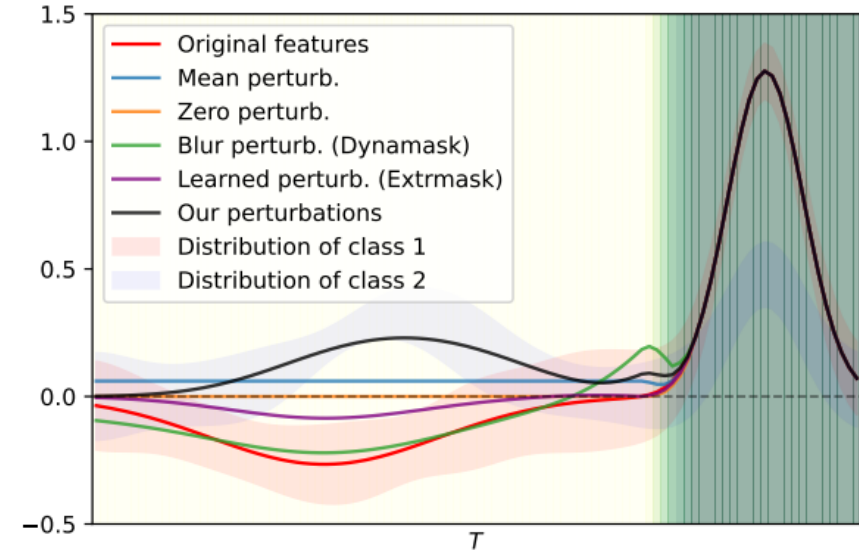
- Setting a more uninformative values is important
- Give only instance-based explanations

# Existing Perturbations are Inadequate

$$\Phi(x, m) = m \times x + (1 - m) \times u$$

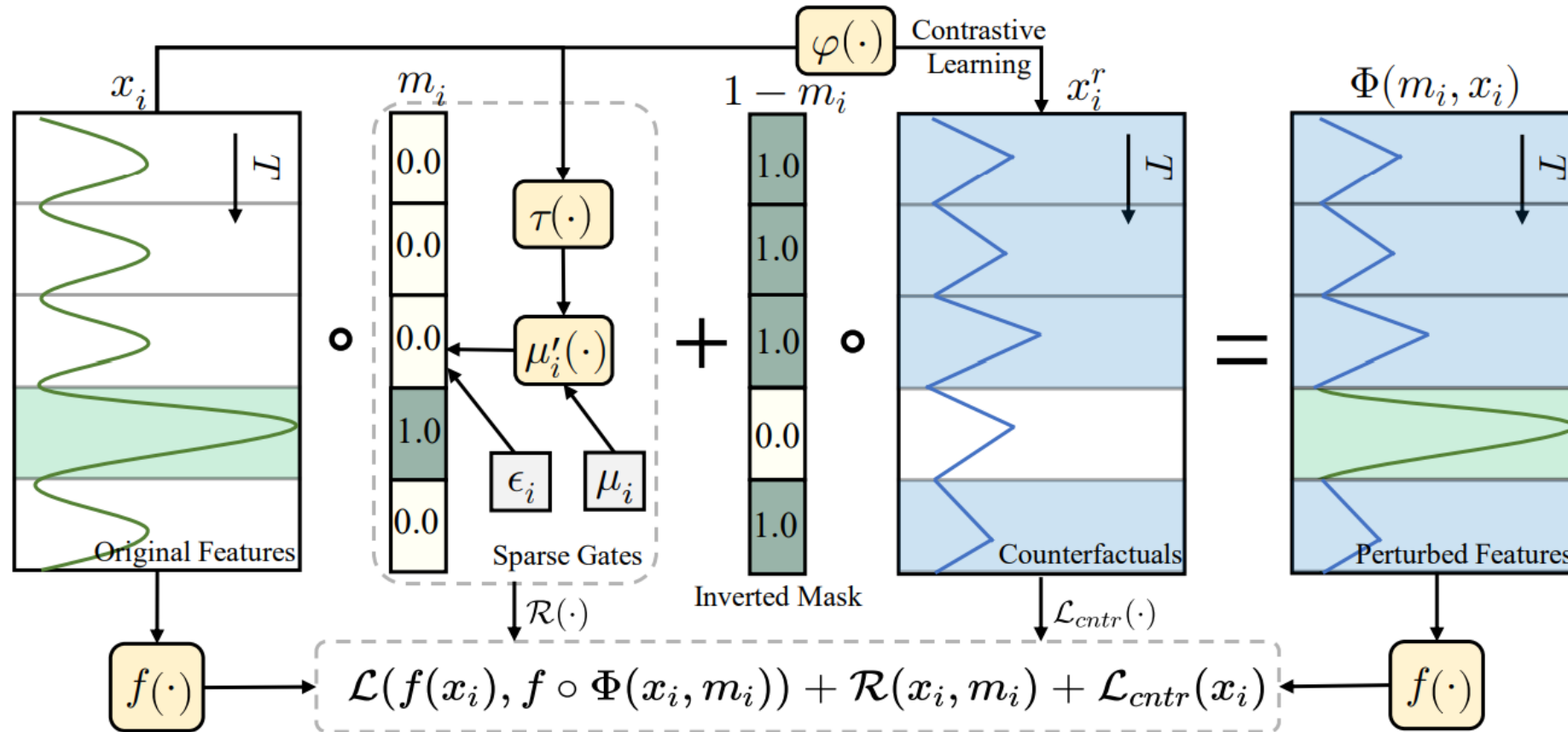
where  $u = \begin{cases} 0 \\ \frac{1}{w+1} \sum_{t-w}^t x_i \\ \text{Gaussian blur} \\ \text{NN}(x) \\ \dots \end{cases}$

- Those perturbations may *out of distribution* or *label leakage*
- Cannot relate temporal patterns *across samples*



Illustrating different styles of perturbation. Other perturbations could be either not uninformative or not in-domain, while ours is counterfactual that is toward the distribution of negative samples.

# ContraLSP Architecture



**Perturbation:**  $\Phi(x, m) = m \times x + (1 - m) \times \varphi_{cntr}(x)$

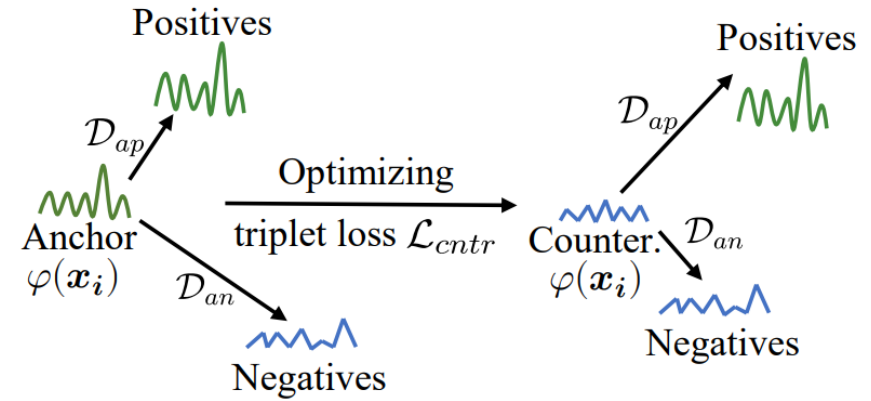
How to learn the *sparse mask*  $m$  and *uninformative*  $\varphi_{cntr}(x)$  ?

# Two Main Contributions

## ➤ Learning counterfactuals from contrastive loss

- Step1: Find positive and negative samples
- Step2: Optimizing via Manhattan distance

$$\mathcal{L}_{cntr}(\mathbf{x}_i) = \max(0, \mathcal{D}_{an} - \mathcal{D}_{ap} - b) + \|\mathbf{x}_i^r\|_1,$$



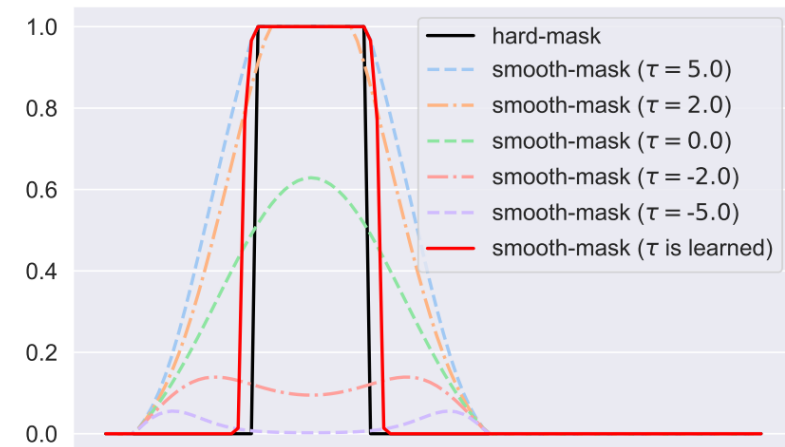
Learning counterfactuals

## ➤ Learning sparse gates with smooth constraint



If not smooth, predictor  $f$  may error!

$$\mathcal{R}(\mathbf{x}_i, \mathbf{m}_i) = \|\mathbf{m}_i\|_0 = \sum_{t=1}^T \sum_{d=1}^D \left( \frac{1}{2} + \frac{1}{2} \operatorname{erf} \left( \frac{\mu'_i[t, d]}{\sqrt{2}\delta} \right) \right),$$



Binary-skewed masks



# Synthetic Experiments (with label)

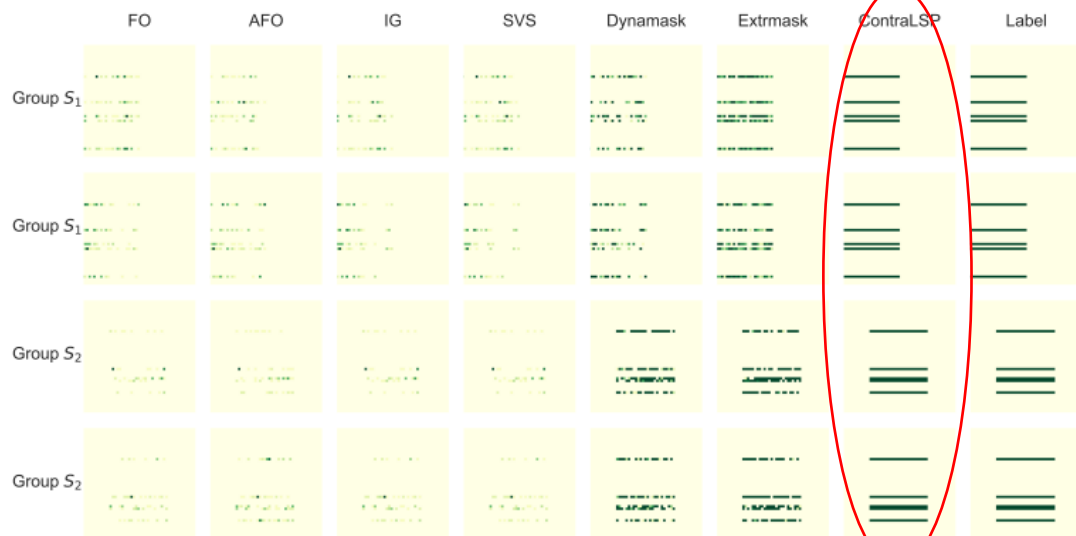
## 1. White-box Regression

Table 1: Performance on Rare-Time and Rare-Observation experiments w/o different groups.

METHOD	RARE-TIME				RARE-TIME (DIFFGROUPS)			
	AUP $\uparrow$	AUR $\uparrow$	$I_m/10^4 \uparrow$	$S_m/10^2 \downarrow$	AUP $\uparrow$	AUR $\uparrow$	$I_m/10^4 \uparrow$	$S_m/10^2 \downarrow$
FO	<b>1.00</b> $\pm 0.00$	0.13 $\pm 0.00$	0.46 $\pm 0.01$	47.20 $\pm 0.61$	<b>1.00</b> $\pm 0.00$	0.16 $\pm 0.00$	0.53 $\pm 0.01$	54.89 $\pm 0.70$
AFO	<b>1.00</b> $\pm 0.00$	0.15 $\pm 0.01$	0.51 $\pm 0.01$	55.60 $\pm 0.85$	<b>1.00</b> $\pm 0.00$	0.16 $\pm 0.00$	0.54 $\pm 0.01$	57.76 $\pm 0.72$
IG	<b>1.00</b> $\pm 0.00$	0.13 $\pm 0.00$	0.46 $\pm 0.01$	47.61 $\pm 0.62$	<b>1.00</b> $\pm 0.00$	0.15 $\pm 0.00$	0.53 $\pm 0.01$	54.62 $\pm 0.85$
SVS	<b>1.00</b> $\pm 0.00$	0.13 $\pm 0.00$	0.47 $\pm 0.01$	47.20 $\pm 0.61$	<b>1.00</b> $\pm 0.00$	0.15 $\pm 0.00$	0.52 $\pm 0.02$	54.28 $\pm 0.84$
DYNAMASK	<u>0.99</u> $\pm 0.01$	0.67 $\pm 0.02$	8.68 $\pm 0.11$	37.24 $\pm 0.48$	<u>0.99</u> $\pm 0.01$	0.51 $\pm 0.00$	5.75 $\pm 0.13$	47.33 $\pm 1.02$
EXTRMASK	<b>1.00</b> $\pm 0.00$	<u>0.88</u> $\pm 0.00$	<u>16.40</u> $\pm 0.13$	<u>13.10</u> $\pm 0.78$	<b>1.00</b> $\pm 0.00$	<u>0.83</u> $\pm 0.03$	<u>13.37</u> $\pm 0.78$	<u>27.44</u> $\pm 3.68$
CONTRALSP	<b>1.00</b> $\pm 0.00$	<b>0.97</b> $\pm 0.01$	<b>19.51</b> $\pm 0.30$	<b>4.65</b> $\pm 0.71$	<b>1.00</b> $\pm 0.00$	<b>0.94</b> $\pm 0.01$	<b>18.92</b> $\pm 0.37$	<b>4.40</b> $\pm 0.60$

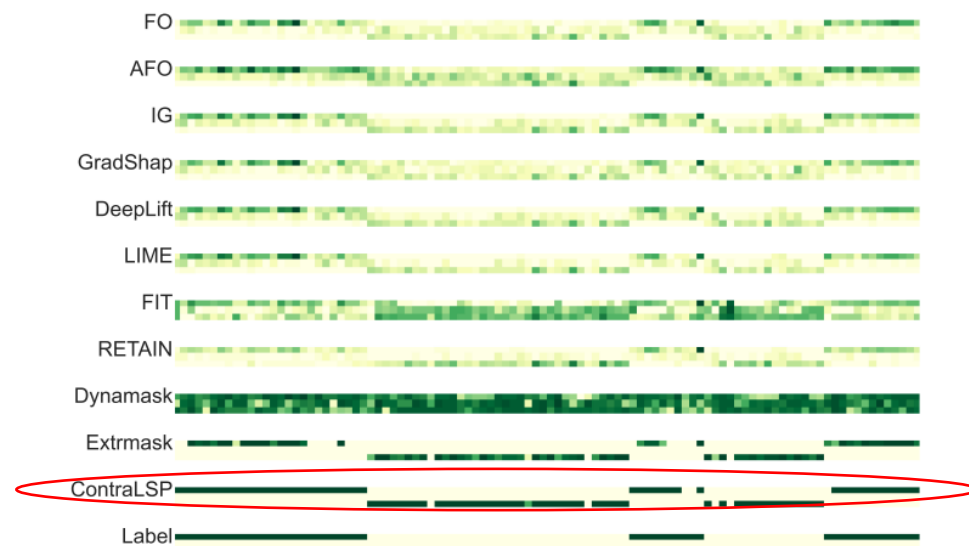
METHOD	RARE-OBSERVATION				RARE-OBSERVATION (DIFFGROUPS)			
	AUP $\uparrow$	AUR $\uparrow$	$I_m/10^4 \uparrow$	$S_m/10^2 \downarrow$	AUP $\uparrow$	AUR $\uparrow$	$I_m/10^4 \uparrow$	$S_m/10^2 \downarrow$
FO	<b>1.00</b> $\pm 0.00$	0.13 $\pm 0.00$	0.46 $\pm 0.00$	47.39 $\pm 0.16$	<b>1.00</b> $\pm 0.00$	0.14 $\pm 0.00$	0.50 $\pm 0.01$	52.13 $\pm 0.96$
AFO	<b>1.00</b> $\pm 0.00$	0.16 $\pm 0.00$	0.55 $\pm 0.01$	56.81 $\pm 0.39$	<b>1.00</b> $\pm 0.00$	0.16 $\pm 0.01$	0.54 $\pm 0.02$	56.92 $\pm 1.24$
IG	<b>1.00</b> $\pm 0.00$	0.13 $\pm 0.00$	0.46 $\pm 0.00$	47.82 $\pm 0.15$	<b>1.00</b> $\pm 0.00$	0.13 $\pm 0.00$	0.47 $\pm 0.00$	49.90 $\pm 0.88$
SVS	<b>1.00</b> $\pm 0.00$	0.13 $\pm 0.00$	0.46 $\pm 0.00$	47.39 $\pm 0.16$	<b>1.00</b> $\pm 0.00$	0.13 $\pm 0.00$	0.47 $\pm 0.01$	49.53 $\pm 0.84$
DYNAMASK	<u>0.97</u> $\pm 0.00$	0.65 $\pm 0.00$	8.32 $\pm 0.06$	22.87 $\pm 0.58$	<u>0.98</u> $\pm 0.00$	0.52 $\pm 0.01$	6.12 $\pm 0.10$	<u>30.88</u> $\pm 0.70$
EXTRMASK	<b>1.00</b> $\pm 0.00$	<u>0.76</u> $\pm 0.00$	<u>13.25</u> $\pm 0.07$	<u>9.55</u> $\pm 0.39$	<b>1.00</b> $\pm 0.00$	<u>0.70</u> $\pm 0.04$	<u>10.40</u> $\pm 0.54$	32.81 $\pm 0.88$
CONTRALSP	<b>1.00</b> $\pm 0.00$	<b>1.00</b> $\pm 0.00$	<b>20.68</b> $\pm 0.03$	<b>0.32</b> $\pm 0.16$	<b>1.00</b> $\pm 0.00$	<b>0.99</b> $\pm 0.00$	<b>20.51</b> $\pm 0.07$	<b>0.57</b> $\pm 0.20$



## 2. Black-box Classification

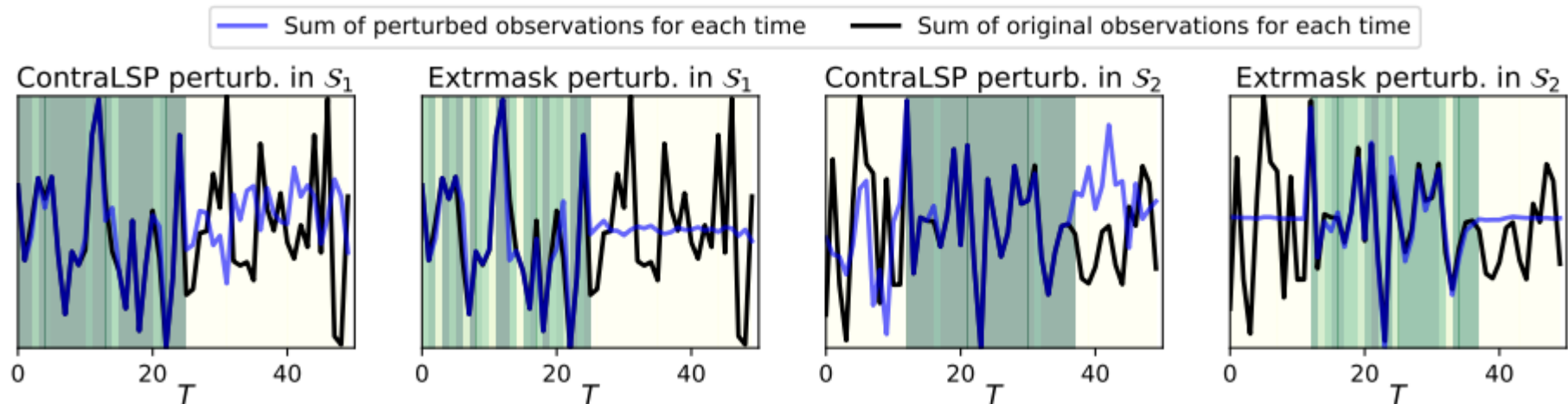
Table 2: Performance on Switch Feature and State data.

METHOD	SWITCH-FEATURE				STATE			
	AUP $\uparrow$	AUR $\uparrow$	$I_m/10^4 \uparrow$	$S_m/10^3 \downarrow$	AUP $\uparrow$	AUR $\uparrow$	$I_m/10^4 \uparrow$	$S_m/10^3 \downarrow$
FO	0.89 $\pm 0.03$	0.37 $\pm 0.02$	1.86 $\pm 0.14$	15.60 $\pm 0.28$	0.90 $\pm 0.05$	0.30 $\pm 0.01$	2.73 $\pm 0.15$	28.07 $\pm 0.54$
AFO	0.82 $\pm 0.06$	0.41 $\pm 0.02$	2.00 $\pm 0.14$	17.32 $\pm 0.29$	0.84 $\pm 0.08$	0.36 $\pm 0.03$	3.16 $\pm 0.27$	34.03 $\pm 1.10$
IG	0.91 $\pm 0.02$	0.44 $\pm 0.03$	2.21 $\pm 0.17$	16.87 $\pm 0.52$	<u>0.93</u> $\pm 0.02$	0.34 $\pm 0.03$	3.17 $\pm 0.28$	30.19 $\pm 1.22$
GRADSHAP	0.88 $\pm 0.02$	0.38 $\pm 0.02$	1.92 $\pm 0.13$	15.85 $\pm 0.40$	0.88 $\pm 0.06$	0.30 $\pm 0.02$	2.76 $\pm 0.20$	28.18 $\pm 0.96$
DEEPLIFT	0.91 $\pm 0.02$	0.44 $\pm 0.02$	2.23 $\pm 0.16$	16.86 $\pm 0.52$	<u>0.93</u> $\pm 0.02$	0.35 $\pm 0.03$	3.20 $\pm 0.27$	30.21 $\pm 1.19$
LIME	0.94 $\pm 0.02$	0.40 $\pm 0.02$	2.01 $\pm 0.13$	16.09 $\pm 0.58$	<b>0.95</b> $\pm 0.02$	0.32 $\pm 0.03$	2.94 $\pm 0.26$	28.55 $\pm 1.53$
FIT	0.48 $\pm 0.03$	0.43 $\pm 0.02$	1.99 $\pm 0.11$	17.16 $\pm 0.50$	0.45 $\pm 0.02$	0.59 $\pm 0.02$	7.92 $\pm 0.40$	33.59 $\pm 0.17$
RETAIN	0.93 $\pm 0.01$	0.33 $\pm 0.04$	1.54 $\pm 0.20$	15.08 $\pm 1.13$	0.52 $\pm 0.16$	0.21 $\pm 0.02$	1.56 $\pm 0.24$	25.01 $\pm 0.57$
DYNAMASK	0.35 $\pm 0.00$	<u>0.77</u> $\pm 0.02$	5.22 $\pm 0.26$	12.85 $\pm 0.53$	0.36 $\pm 0.01$	<u>0.79</u> $\pm 0.01$	10.59 $\pm 0.20$	25.11 $\pm 0.40$
EXTRMASK	<u>0.97</u> $\pm 0.01$	0.65 $\pm 0.05$	<u>8.45</u> $\pm 0.51$	<u>6.90</u> $\pm 1.44$	0.87 $\pm 0.01$	0.77 $\pm 0.01$	<u>29.71</u> $\pm 1.39$	<u>7.54</u> $\pm 0.46$
CONTRALSP	<b>0.98</b> $\pm 0.00$	<b>0.80</b> $\pm 0.03$	<b>24.23</b> $\pm 1.27$	<b>0.91</b> $\pm 0.26$	0.90 $\pm 0.03$	<b>0.81</b> $\pm 0.01$	<b>50.09</b> $\pm 0.78$	<b>0.50</b> $\pm 0.05$



# Synthetic Experiments (with label)

## ➤ Counterfactual information



## ➤ Distribution analysis of perturbations

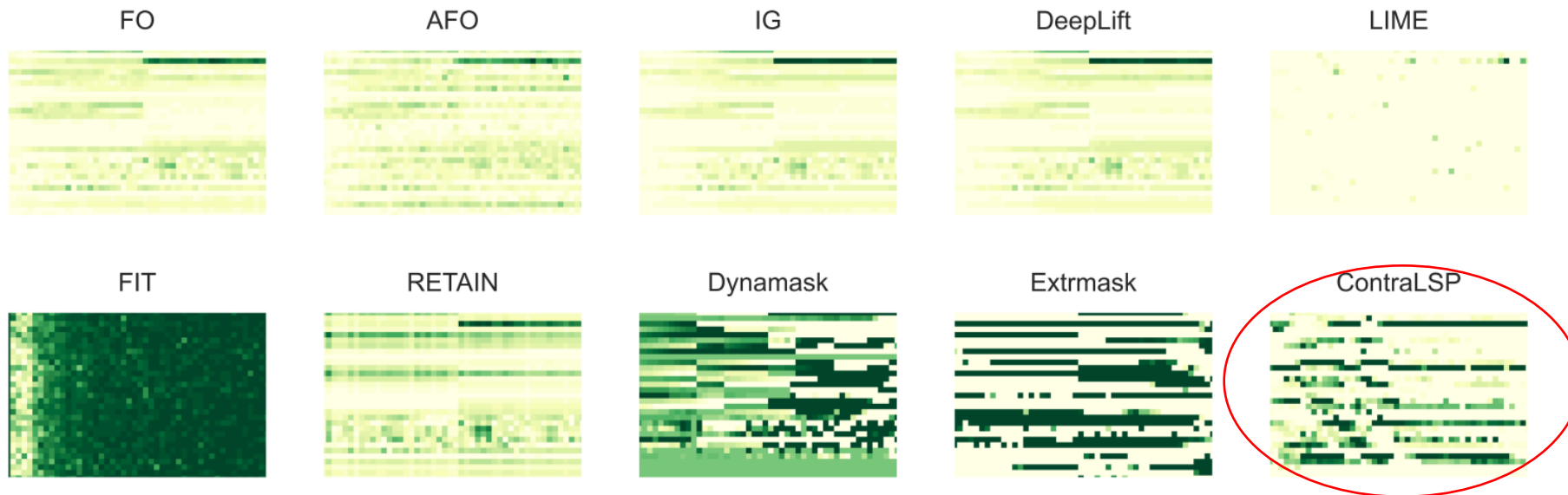
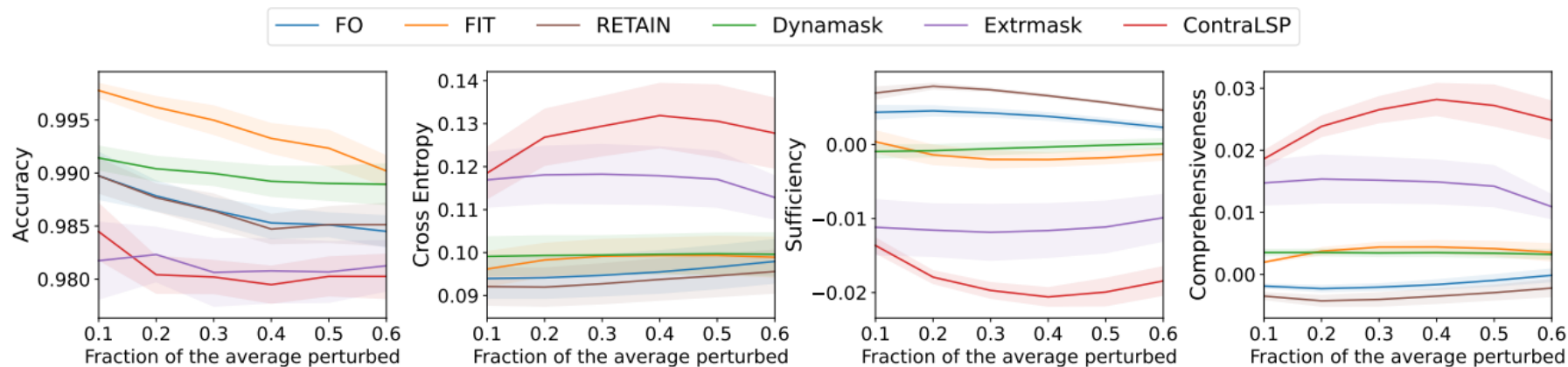
Table 12: Difference between the distribution of different perturbations and the original distribution.

PERTURBATION TYPE	RARE-TIME		RARE-OBSERVATION	
	KDE-SCORE $\uparrow$	KL-DIVERGENCE $\downarrow$	KDE-SCORE $\uparrow$	KL-DIVERGENCE $\downarrow$
ZERO PERTURBATION	-25.242	0.0523	-23.377	0.0421
MEAN PERTURBATION	-30.805	0.0731	-26.421	0.0589
EXTRMASK PERTURBATION	-22.532	0.0219	-19.102	0.0104
CONTRALSP PERTURBATION	-23.290	0.0393	-22.732	0.0386



# Real-world Experiments (without label)

## 3. MIMIC-III Mortality Data



# Closing Remarks

---

- We propose ContraLSP as a time series explainer, which incorporates counterfactual samples to build uninformative in-domain perturbation.
- We incorporate sample-specific sparse gates to generate more binary-skewed and smooth masks.
- The code is available at <https://github.com/zichuan-liu/ContraLSP>.

**Thanks for your listening!**