



Explaining Time Series via Contrastive and Locally Sparse Perturbations

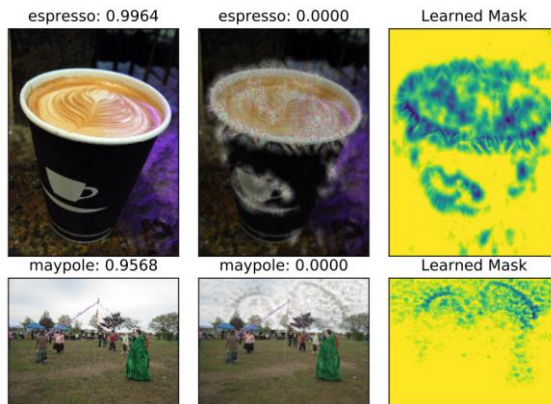
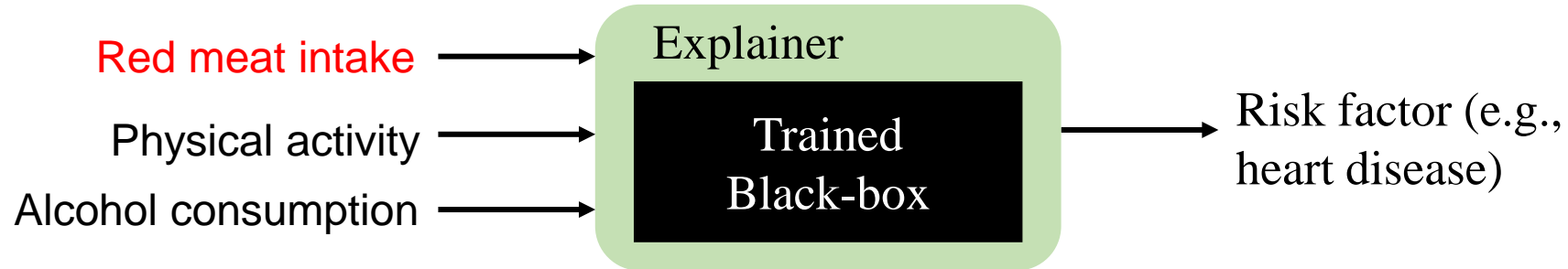
Zichuan Liu^{1,2} Yingying Zhang² Tianchun Wang³ Zefan Wang^{2,4} Dongsheng Luo⁵
Mengnan Du⁶ Min Wu⁷ Yi Wang⁸ Chunlin Chen¹ Lunting Fan² Qingsong Wen²

¹Nanjing University, ²Ailibaba Group,
³Pennsylvania State University, ⁴Tsinghua University,
⁵Florida International University,
⁶New Jersey Institute of Technology,
⁷A*STAR, ⁸The University of Hong Kong



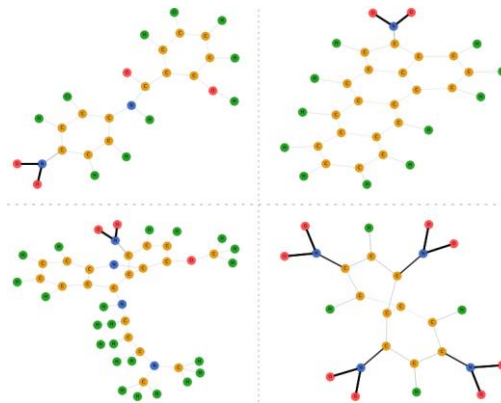
Background

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



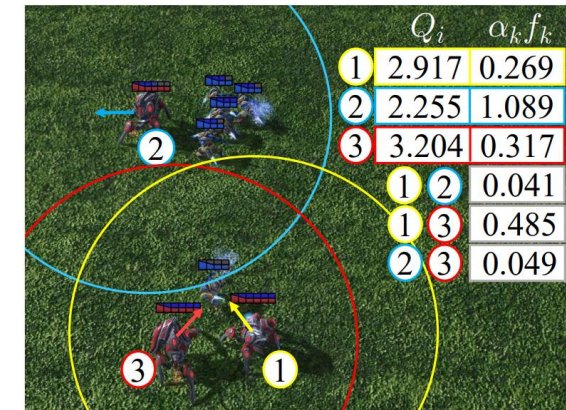
Visual Explanation

Source: [Fong et al.](#)



Graph Explanation

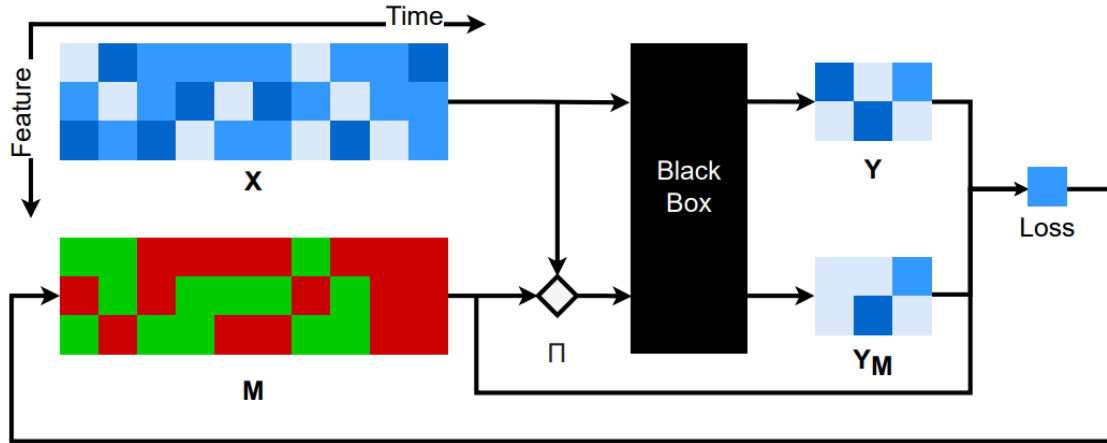
Source: [Miao et al.](#)



Game Explanation

Source: [Liu et al.](#)

Challenges for Explaining Time Series



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

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

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

➤ 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

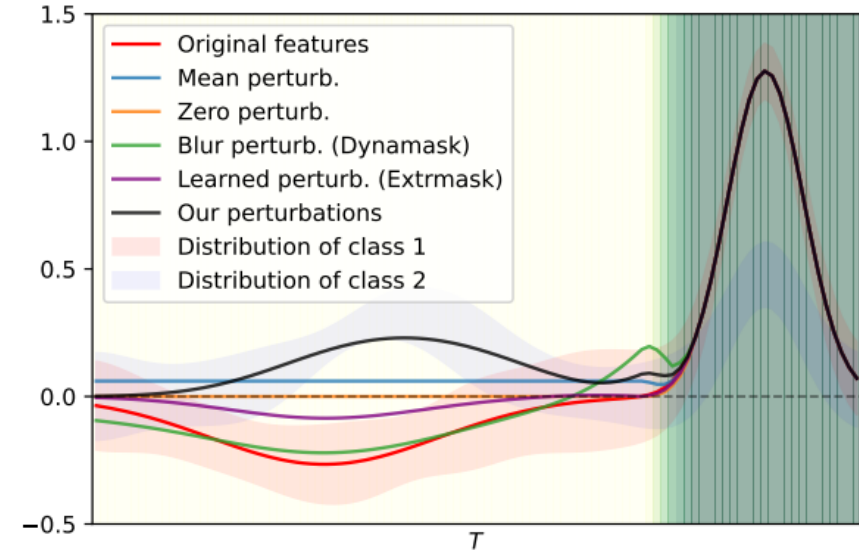
- Setting a more uninformative values is important
- Give only instance-based explanations

Existing Perturbations are Inadequate

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

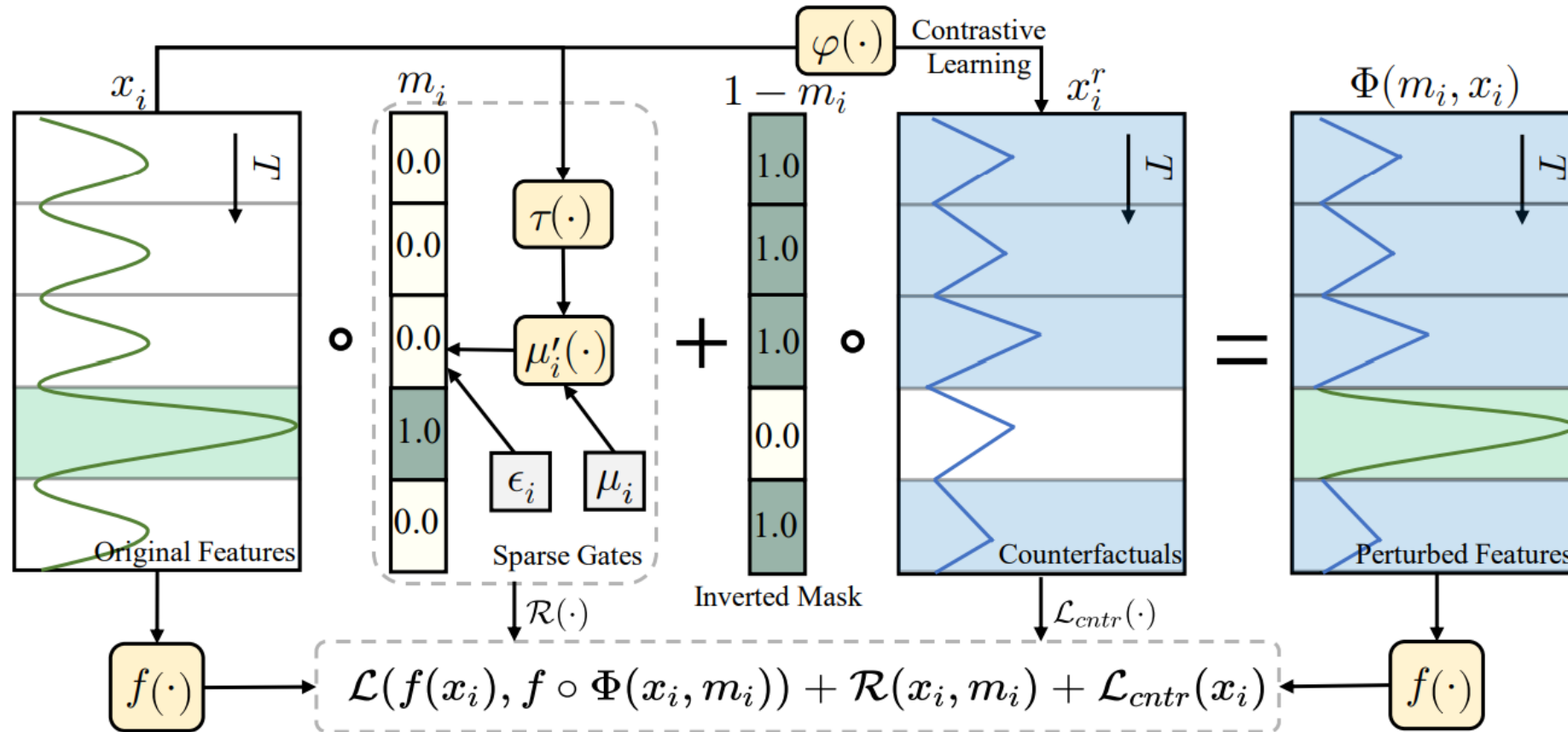
where $u = \begin{cases} 0 \\ \frac{1}{w+1} \sum_{t-w}^t x_i \\ \text{Gaussian blur} \\ \text{NN}(x) \\ \dots \end{cases}$

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

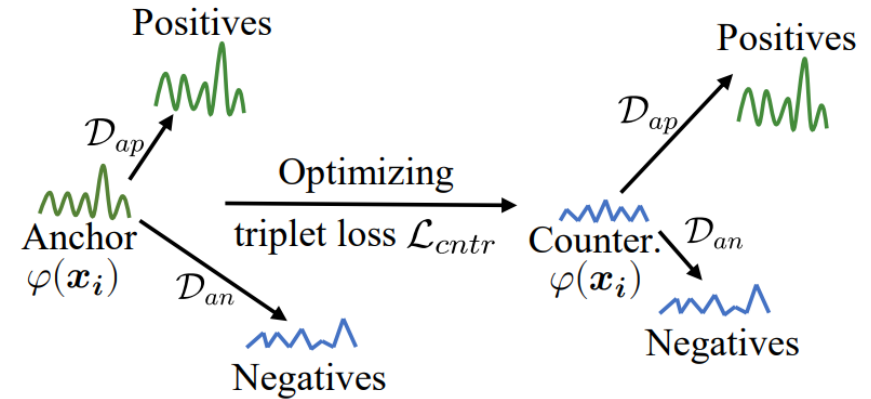
How to learn the *sparse mask* m and *uninformative* $\varphi_{cntr}(x)$?

Two Main Contributions

➤ Learning counterfactuals from contrastive loss

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

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



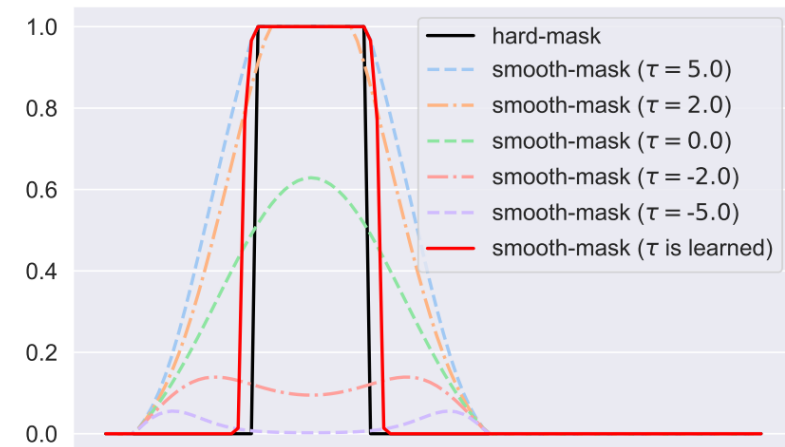
Learning counterfactuals

➤ 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

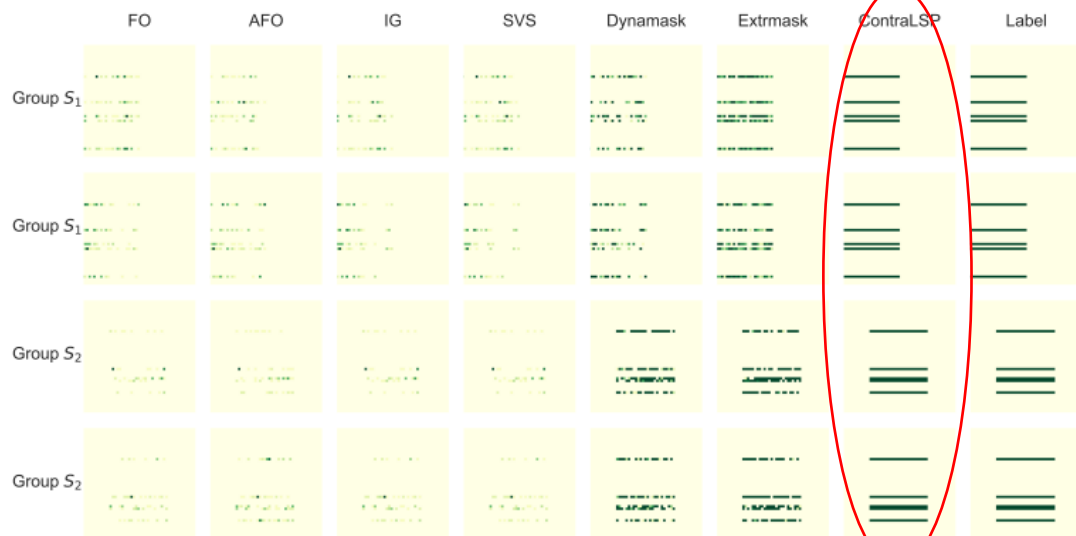
Synthetic Experiments (with label)

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	1.00 ± 0.00	0.13 ± 0.00	0.46 ± 0.01	47.20 ± 0.61	1.00 ± 0.00	0.16 ± 0.00	0.53 ± 0.01	54.89 ± 0.70
AFO	1.00 ± 0.00	0.15 ± 0.01	0.51 ± 0.01	55.60 ± 0.85	1.00 ± 0.00	0.16 ± 0.00	0.54 ± 0.01	57.76 ± 0.72
IG	1.00 ± 0.00	0.13 ± 0.00	0.46 ± 0.01	47.61 ± 0.62	1.00 ± 0.00	0.15 ± 0.00	0.53 ± 0.01	54.62 ± 0.85
SVS	1.00 ± 0.00	0.13 ± 0.00	0.47 ± 0.01	47.20 ± 0.61	1.00 ± 0.00	0.15 ± 0.00	0.52 ± 0.02	54.28 ± 0.84
DYNAMASK	<u>0.99</u> ± 0.01	0.67 ± 0.02	8.68 ± 0.11	37.24 ± 0.48	<u>0.99</u> ± 0.01	0.51 ± 0.00	5.75 ± 0.13	47.33 ± 1.02
EXTRMASK	1.00 ± 0.00	<u>0.88</u> ± 0.00	<u>16.40</u> ± 0.13	<u>13.10</u> ± 0.78	1.00 ± 0.00	<u>0.83</u> ± 0.03	<u>13.37</u> ± 0.78	<u>27.44</u> ± 3.68
CONTRALSP	1.00 ± 0.00	0.97 ± 0.01	19.51 ± 0.30	4.65 ± 0.71	1.00 ± 0.00	0.94 ± 0.01	18.92 ± 0.37	4.40 ± 0.60

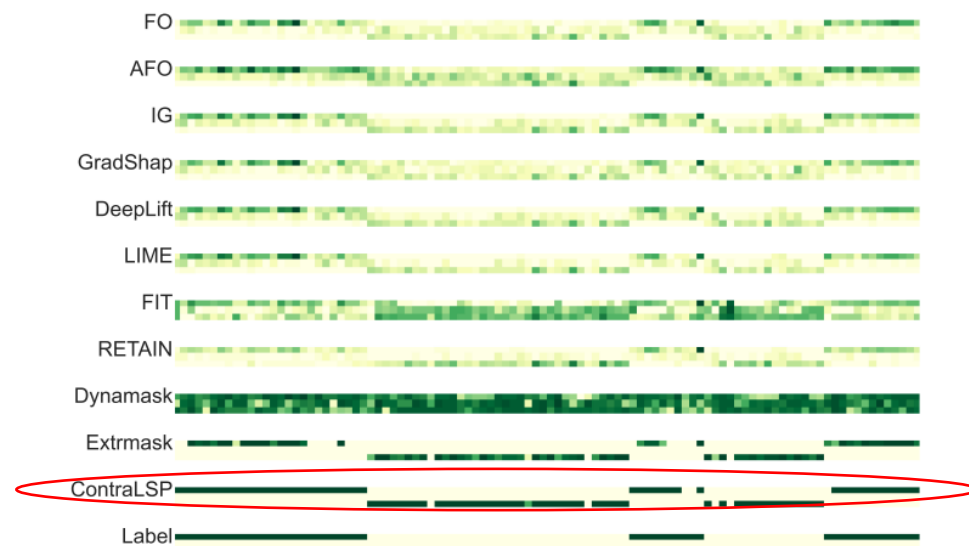
METHOD	RARE-OBSERVATION				RARE-OBSERVATION (DIFFGROUPS)			
	AUP \uparrow	AUR \uparrow	$I_m/10^4 \uparrow$	$S_m/10^2 \downarrow$	AUP \uparrow	AUR \uparrow	$I_m/10^4 \uparrow$	$S_m/10^2 \downarrow$
FO	1.00 ± 0.00	0.13 ± 0.00	0.46 ± 0.00	47.39 ± 0.16	1.00 ± 0.00	0.14 ± 0.00	0.50 ± 0.01	52.13 ± 0.96
AFO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



2. Black-box Classification

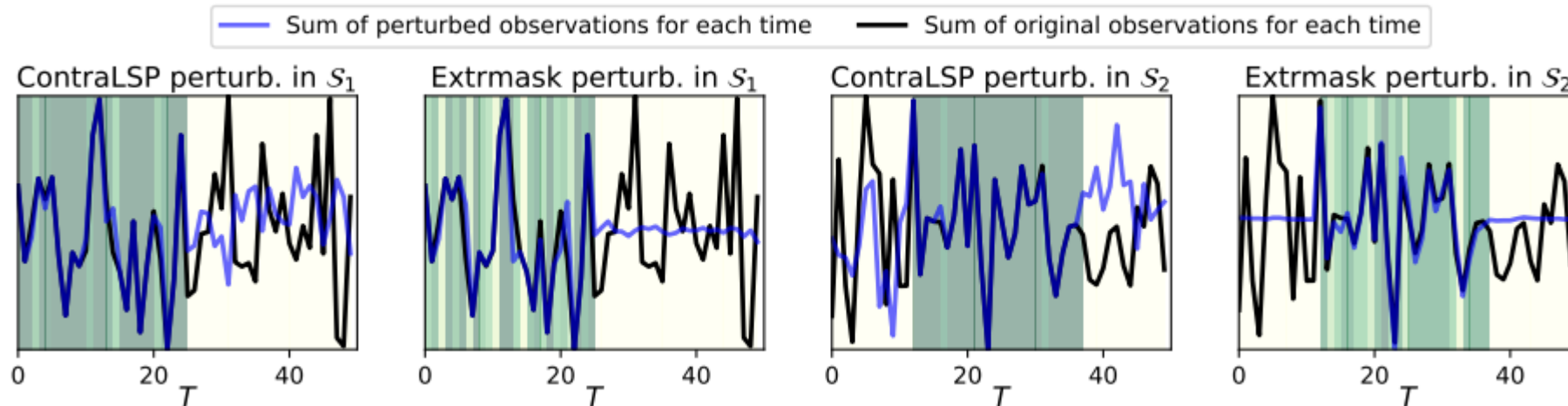
Table 2: Performance on Switch Feature and State data.

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



Synthetic Experiments (with label)

➤ Counterfactual information



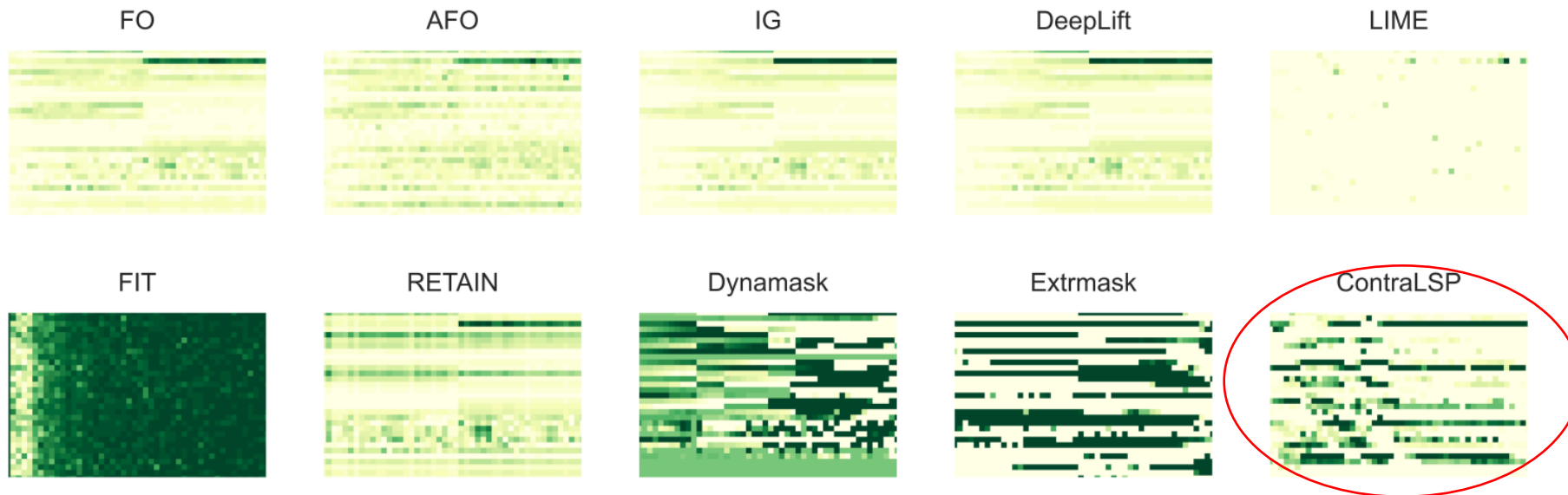
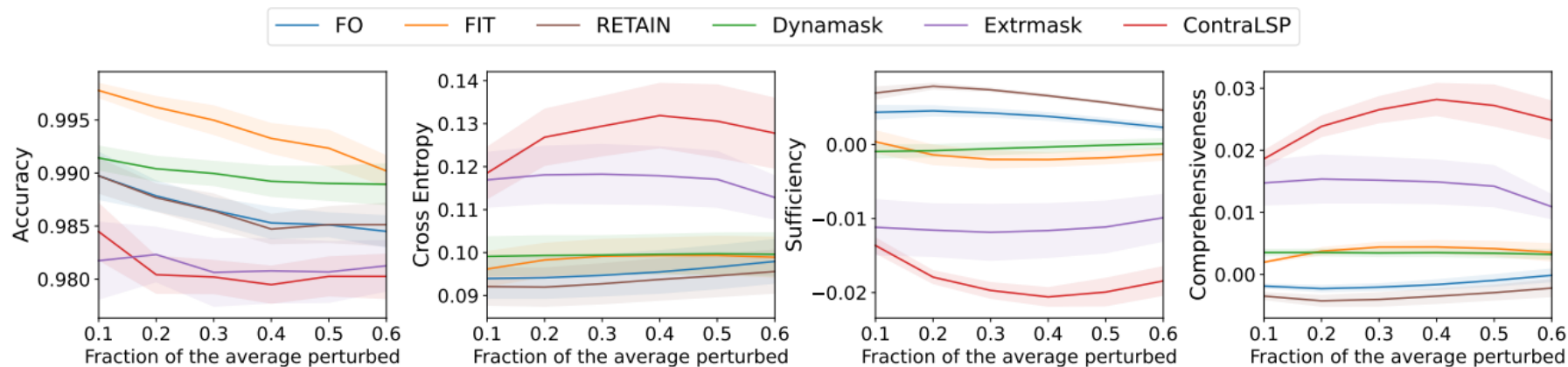
➤ Distribution analysis of perturbations

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!