
Delta-XAI: A Unified Framework for Explaining Prediction Changes in Online Time Series Monitoring



Changhun Kim*



Yechan Mun*



Hyeongwon Jang



Eunseo Lee



Sangchul Hahn[†]



Eunho Yang[†]

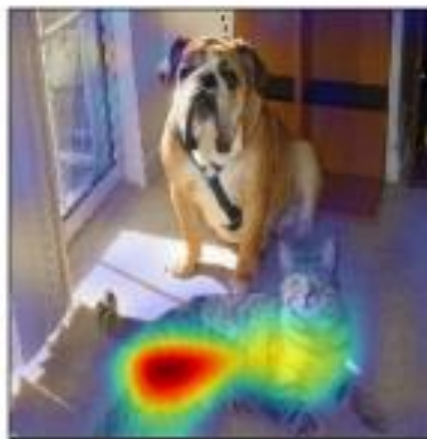
*Equal Contribution [†]Equal Advising

Motivation: Explainable Artificial Intelligence (XAI)

- Modern deep learning models achieve strong predictive performance, but their decision-making processes are often **difficult to interpret**.
- **Explainable artificial intelligence (XAI)** methods such as SHAP, LIME, and Grad-CAM have been proposed to provide human-understandable explanations for model predictions.
- Trustworthy AI systems are increasingly important in real-world applications, especially in **high-stakes domains** such as healthcare, and autonomous driving.



Original Image



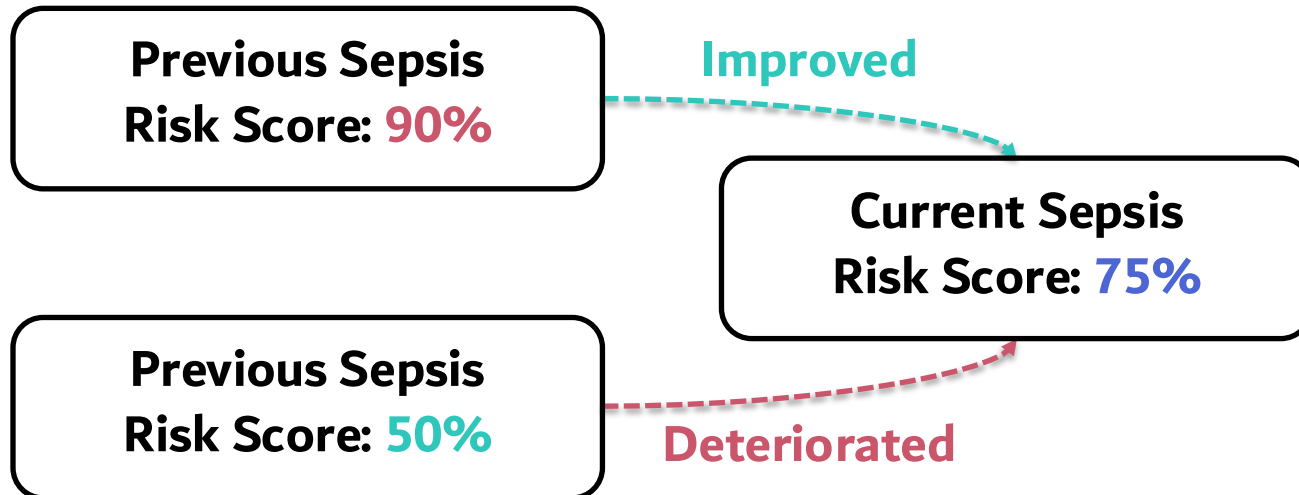
Grad-CAM 'Cat'



Grad-CAM 'Dog'

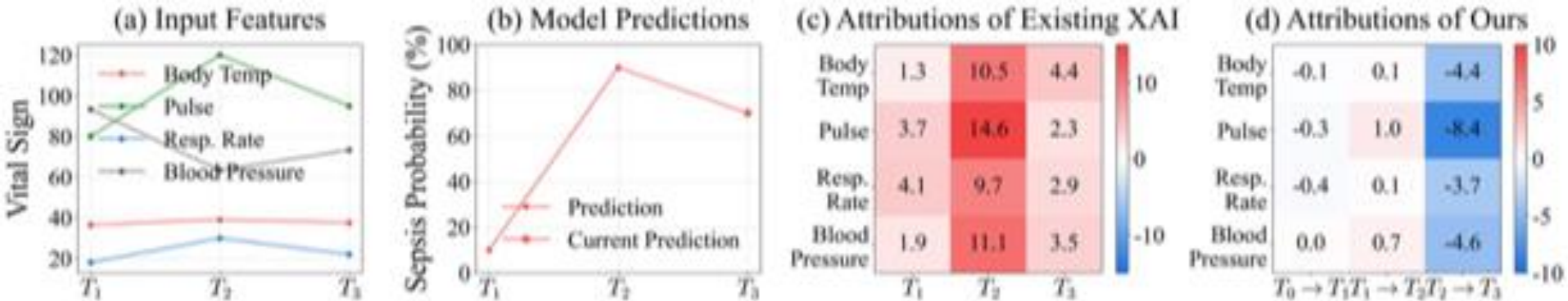
Motivation: XAI for Online Time Series Monitoring

- In online time series monitoring, **explaining prediction changes** between nearby time steps is more important than **explaining a single prediction**.
- For example, in sepsis monitoring, clinicians require explanations for both **improvement** (90% → 75%) and **deterioration** (50% → 75%), even if the **final prediction is the same** (75%).
- Existing XAI methods are mostly designed for **single-time predictions** and cannot adequately explain **temporal prediction changes**.



Motivation: XAI for Online Time Series Monitoring

- A patient's risk may rise from 10% at T_1 to 90% at T_2 , then decrease to 70% at T_3 . Clinicians want to explain the **recovery** from $T_2 \rightarrow T_3$.
- Existing XAI methods explain only the prediction at T_3 , often still highlighting features responsible for the earlier deterioration at T_2 .
- Our Delta-XAI framework instead explains prediction **changes** between consecutive time steps, such as $T_1 \rightarrow T_2$ and $T_2 \rightarrow T_3$.



Proposed Method: Adopting Static XAI to Prediction Changes

- We formalize **explanations for prediction change** in online time series monitoring by introducing a **prediction difference wrapper** that **unifies existing 14 single-time XAI methods** and enables them to directly explain temporal prediction changes.

$$g : \mathbb{R}^{(T_2 - T_1 + W) \times D} \rightarrow [0, 1]^C, \quad g(\mathbf{X}_{T_1 - W + 1 : T_2}) := f(\mathbf{X}_{T_2 - W + 1 : T_2}) - f(\mathbf{X}_{T_1 - W + 1 : T_1}),$$

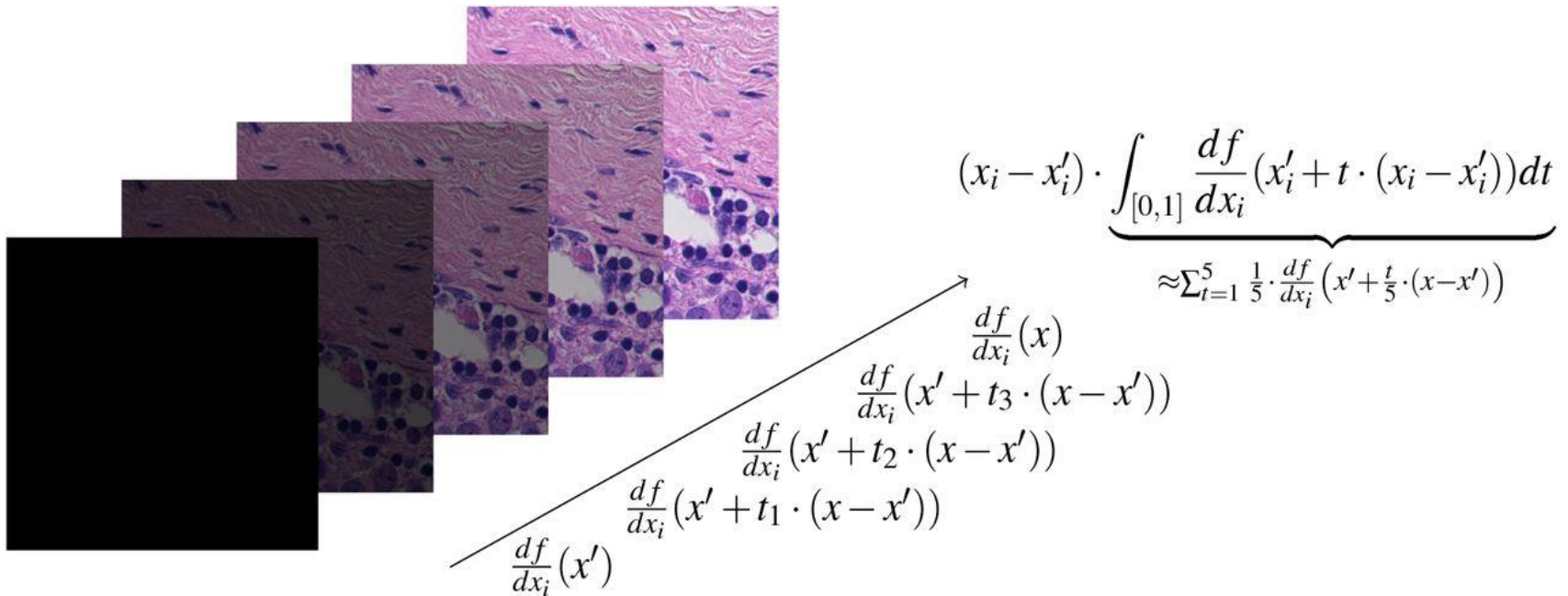
- For **linear and complete methods** such as SHAP and Integrated Gradients, explanations for prediction change reduce to differences of single-time attributions.

Theorem 1 (Attribution Decomposition Theorem for Online Completeness). *Given a linear and complete attribution method φ with a fixed baseline, the following decomposition holds:*

$$\begin{aligned}
 f(\mathbf{X}_{T_2 - W + 1 : T_2})_{\hat{c}} - f(\mathbf{X}_{T_1 - W + 1 : T_1})_{\hat{c}} &= \underbrace{\sum_{t=T_1+1}^{T_2} \sum_{d=1}^D \varphi(f, \mathbf{X}_{t,d} \mid T_2)}_{\text{Addition of newest features}} \\
 + \underbrace{\sum_{t=T_2-W+1}^{T_1} \sum_{d=1}^D [\varphi(f, \mathbf{X}_{t,d} \mid T_2) - \varphi(f, \mathbf{X}_{t,d} \mid T_1)]}_{\text{Delayed effect of intermediate features}} &- \underbrace{\sum_{t=T_1-W+1}^{T_2-W} \sum_{d=1}^D \varphi(f, \mathbf{X}_{t,d} \mid T_1)}_{\text{Removal of oldest features}}.
 \end{aligned} \tag{4}$$

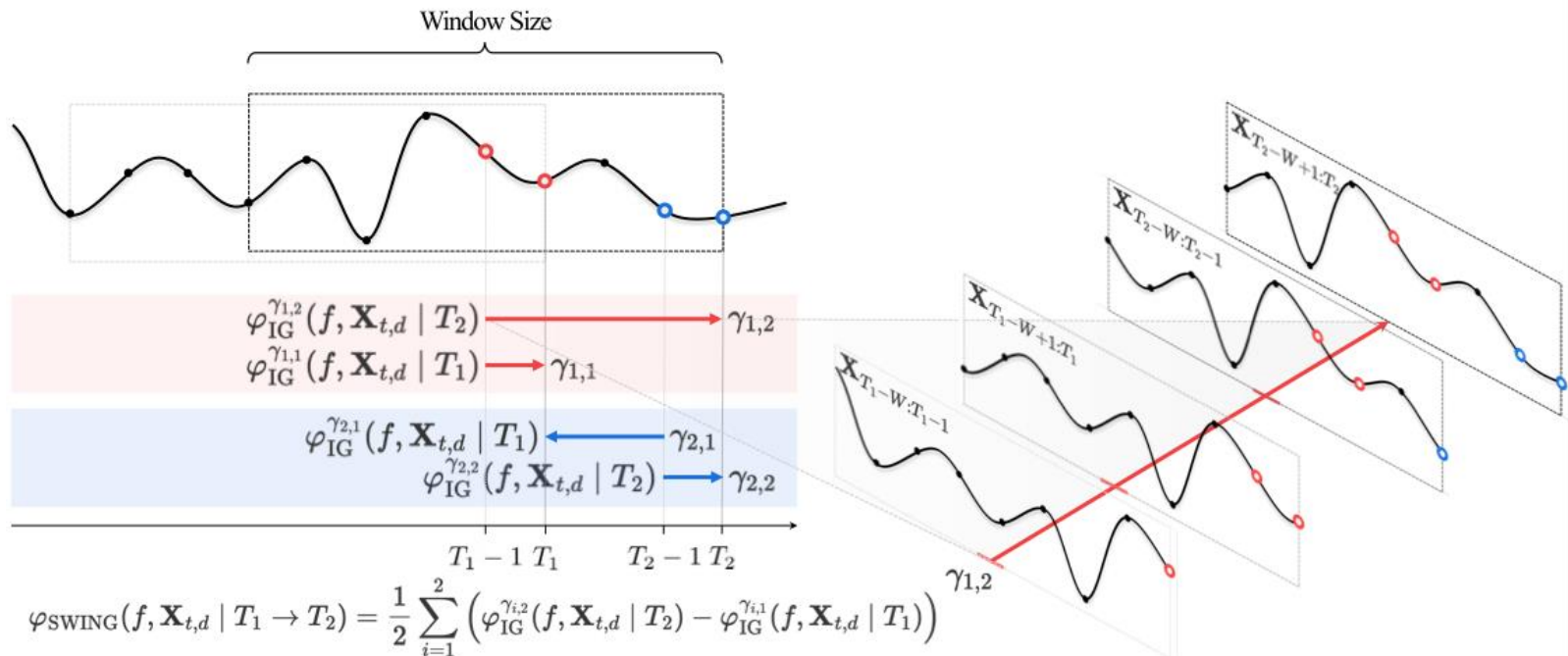
Proposed Method: Background of Integrated Gradients (IG)

- **Integrated Gradients (IG)** uses linear interpolation from baseline; intermediate points can be **Out-of-Distribution (OOD)** region.
- IG also fails to consider **temporal dependencies**; it merely scales the input.



Proposed Method: Shifted Window Integrated Gradients (SWING)

- We extend **Integrated Gradients (IG)** to explain prediction changes by integrating along temporally shifted windows rather than static input-baseline pairs.
- SWING uses **recent past observations as realistic baselines** and follows **piecewise-linear historical paths** to reduce **off-manifold artifacts**.
- By combining **retrospective baseline selection**, **dual-path integration**, and **historical integration**, SWING produces temporally consistent explanations while preserving key theoretical properties.



Proposed Evaluation Metrics

- **Area-based ranking quality:** AUPD/AUPP aggregate CPD/CPD over all prefixes, evaluating whether highly ranked features consistently affect predictions.

$$\text{AUPD}(g, \mathbf{X}, K) = \frac{1}{2K} \sum_{k=1}^K [\text{CPD}(g, \mathbf{X}, k) + \text{CPD}(g, \mathbf{X}, k - 1)]$$

- **Macro-level temporal dynamics:** MPD/MPP first aggregate attribution scores across nearby online windows to capture evolving explanations.

$$\varphi(\mathbf{X}_{t,d} | T) = \frac{1}{2W - 1} \sum_{T'=t-W+1}^{t+W-1} \varphi(\mathbf{X}_{t,d} | T')$$

- **Magnitude-aware attribution quality:** Corr. measures whether attribution magnitudes align with actual prediction changes after feature removal.

$$\text{Corr.} = \text{Corr} \left([\varphi^{(1)}, \dots, \varphi^{(K)}], [|\mathbf{g}_1^\uparrow - \mathbf{g}_0^\uparrow|, \dots, |\mathbf{g}_K^\uparrow - \mathbf{g}_{K-1}^\uparrow|] \right)$$

Experiments: Main Table

- SWING consistently outperforms existing XAI methods.

Table 2: Performance comparison of XAI methods on clinical prediction tasks: MIMIC-III decompensation benchmark using LSTM as backbone architecture. Evaluation is performed by removing the most or least salient 50 feature points per time step, using forward-fill substitution.

Algorithm	Removal of Most Salient 50 Points				Removal of Least Salient 50 Points				Corr. \uparrow
	CPD \uparrow	AUPD \uparrow	MPD \uparrow	AUMPD \uparrow	CPP \downarrow	AUPP \downarrow	MPP \downarrow	AUMPP \downarrow	
LIME (Ribeiro et al., 2016)	2.26 \pm 0.04	1.72 \pm 0.03	13.78 \pm 0.08	7.70 \pm 0.02	32.46 \pm 0.16	14.26 \pm 0.10	33.45 \pm 0.14	15.30 \pm 0.10	0.02 \pm 0.00
GradSHAP (Lundberg & Lee, 2017)	13.73 \pm 0.06	9.05 \pm 0.04	16.68 \pm 0.08	11.19 \pm 0.05	32.97 \pm 0.13	13.96 \pm 0.06	30.13 \pm 0.17	11.95 \pm 0.09	0.14 \pm 0.00
IG (Sundararajan et al., 2017)	13.42 \pm 0.06	9.10 \pm 0.05	16.14 \pm 0.07	11.31 \pm 0.04	33.55 \pm 0.12	13.97 \pm 0.04	29.46 \pm 0.17	10.85 \pm 0.05	0.17 \pm 0.00
DeepLIFT (Shrikumar et al., 2017)	13.58 \pm 0.06	9.42 \pm 0.04	16.03 \pm 0.08	11.25 \pm 0.06	35.96 \pm 0.16	14.61 \pm 0.07	31.53 \pm 0.15	11.41 \pm 0.06	0.19 \pm 0.00
FO (Suresh et al., 2017)	13.14 \pm 0.10	9.92 \pm 0.07	17.79 \pm 0.10	12.92 \pm 0.06	44.02 \pm 0.19	22.32 \pm 0.09	16.92 \pm 0.09	6.24 \pm 0.03	0.26 \pm 0.00
AFO (Tonekaboni et al., 2020)	13.24 \pm 0.07	9.30 \pm 0.05	17.16 \pm 0.07	11.95 \pm 0.03	36.13 \pm 0.22	16.64 \pm 0.09	24.13 \pm 0.16	9.32 \pm 0.07	0.28 \pm 0.00
FIT (Tonekaboni et al., 2020)	3.40 \pm 0.04	2.70 \pm 0.03	7.11 \pm 0.04	6.14 \pm 0.04	35.52 \pm 0.11	17.55 \pm 0.08	12.07\pm0.05	10.19 \pm 0.05	0.06 \pm 0.00
WinIT (Leung et al., 2023)	19.64 \pm 0.07	12.25 \pm 0.04	24.87\pm0.13	15.45 \pm 0.08	29.22 \pm 0.07	13.05 \pm 0.05	26.11 \pm 0.12	11.92 \pm 0.06	0.21 \pm 0.00
Dynamask (Crabbé & Van Der Schaar, 2021)	11.72 \pm 0.08	7.56 \pm 0.04	13.15 \pm 0.08	8.25 \pm 0.04	53.07 \pm 0.24	26.22 \pm 0.08	49.80 \pm 0.16	24.26 \pm 0.06	0.04 \pm 0.00
Extrmask (Enguehard, 2023)	16.66 \pm 0.11	10.47 \pm 0.06	17.51 \pm 0.12	10.63 \pm 0.05	29.91 \pm 0.17	15.13 \pm 0.09	29.64 \pm 0.17	14.84 \pm 0.12	0.08 \pm 0.00
ContraLSP (Liu et al., 2024b)	12.88 \pm 0.36	8.69 \pm 0.26	18.00 \pm 0.16	11.11 \pm 0.18	41.62 \pm 0.30	21.17 \pm 0.10	42.67 \pm 0.29	21.94 \pm 0.10	0.03 \pm 0.00
TimeX (Queen et al., 2024)	16.99 \pm 0.09	11.45 \pm 0.06	19.45 \pm 0.10	12.45 \pm 0.06	50.34 \pm 0.10	24.11 \pm 0.04	51.06 \pm 0.09	24.50 \pm 0.05	0.03 \pm 0.00
TimeX++ (Liu et al., 2024a)	11.12 \pm 0.05	7.00 \pm 0.04	13.14 \pm 0.02	7.76 \pm 0.02	34.21 \pm 0.17	13.72 \pm 0.07	32.34 \pm 0.11	13.08 \pm 0.05	0.03 \pm 0.00
TIMING (Jang et al., 2025)	14.99 \pm 0.07	9.71 \pm 0.05	16.50 \pm 0.08	11.53 \pm 0.04	31.22 \pm 0.16	13.36 \pm 0.05	27.19 \pm 0.19	10.24 \pm 0.07	0.19 \pm 0.00
SWING	23.87\pm0.16	16.23\pm0.10	22.27\pm0.19	15.52\pm0.12	17.76\pm0.04	5.85\pm0.04	18.20 \pm 0.06	6.06\pm0.05	0.40\pm0.00

Experiments: Ablation Study

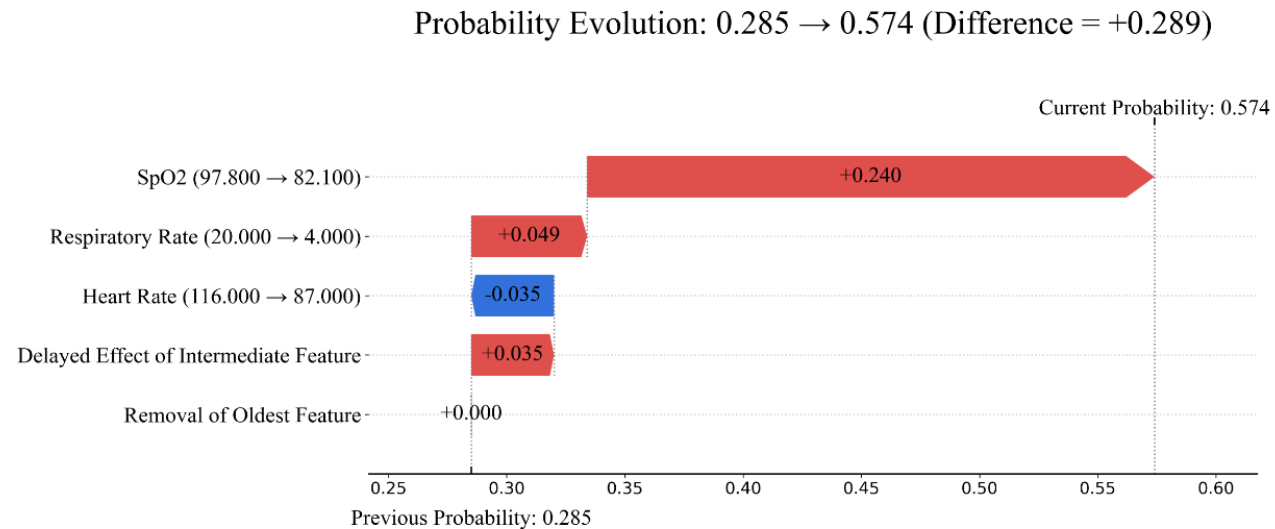
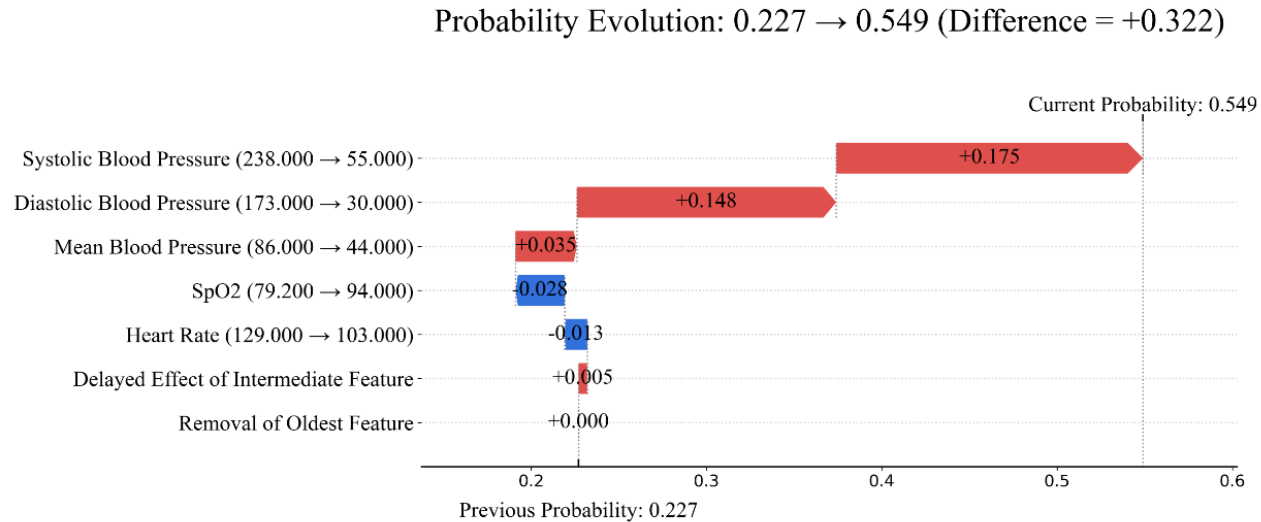
- Each component contributes to reliable prediction-change explanations.

Table 3: Ablation study of SWING examining retrospective baseline selection (RBS), dual-path integration (DPI), and piecewise-linear historical integration (PHI) on the MIMIC-III decompensation benchmark (Johnson et al., 2016), with LSTM (Hochreiter & Schmidhuber, 1997) backbone and interval $T_2 - T_1 = 24$. We vary the baseline distance d (default: 1) and remove the most or least salient 50 feature points per time step, using forward-fill substitution.

Algorithm	Removal of Most Salient 50 Points				Removal of Least Salient 50 Points				Corr. \uparrow
	CPD \uparrow	AUPD \uparrow	MPD \uparrow	AUMPD \uparrow	CPP \downarrow	AUPP \downarrow	MPP \downarrow	AUMPP \downarrow	
w/o RBS, PHI	45.82 \pm 0.20	26.85 \pm 0.13	40.85 \pm 0.22	26.31 \pm 0.14	104.05 \pm 0.34	48.43 \pm 0.19	77.94 \pm 0.30	29.56 \pm 0.11	0.20 \pm 0.00
w/o RBS	40.25 \pm 0.22	24.39 \pm 0.12	46.71 \pm 0.35	29.10 \pm 0.21	79.70 \pm 0.30	31.29 \pm 0.07	82.59 \pm 0.35	33.33 \pm 0.11	0.16 \pm 0.00
w/o DPI ($\gamma_{1,1}, \gamma_{2,2}$)	55.01\pm0.35	32.39\pm0.21	53.09\pm0.38	33.08\pm0.21	71.01 \pm 0.30	26.12 \pm 0.10	76.26 \pm 0.34	29.29 \pm 0.09	0.19 \pm 0.00
$d = 0$	33.80 \pm 0.23	20.11 \pm 0.13	48.98 \pm 0.44	29.78 \pm 0.24	82.28 \pm 0.29	33.99 \pm 0.10	83.69 \pm 0.35	33.31 \pm 0.12	0.11 \pm 0.00
$d = 3$	42.00 \pm 0.33	26.00 \pm 0.28	47.49 \pm 0.44	29.74 \pm 0.34	68.53 \pm 0.49	25.60 \pm 0.18	71.88 \pm 0.62	27.54 \pm 0.27	0.19 \pm 0.00
$d = 5$	41.98 \pm 0.31	25.65 \pm 0.25	47.18 \pm 0.45	29.44 \pm 0.33	72.31 \pm 0.66	27.88 \pm 0.22	74.62 \pm 0.71	29.11 \pm 0.26	0.18 \pm 0.00
$d = 10$	41.02 \pm 0.28	24.94 \pm 0.18	46.89 \pm 0.39	29.23 \pm 0.25	76.55 \pm 0.48	29.96 \pm 0.16	78.11 \pm 0.44	31.04 \pm 0.16	0.17 \pm 0.00
SWING	41.07 \pm 0.22	26.46 \pm 0.18	<u>50.58\pm0.28</u>	<u>32.29\pm0.20</u>	60.29\pm0.14	21.60\pm0.08	64.43\pm0.19	23.87\pm0.10	0.21\pm0.00

Experiments: Qualitative Analysis

- SWING is not only accurate, but also provides clinically coherent explanations.



Conclusion

- We formulate **explanations for prediction change** in online time series monitoring and **adapt 14 existing XAI methods** through a **unified prediction wrapper** with **tailored evaluation metrics**.
- We propose **SWING**, which extends IG with **historical integration paths** to capture temporal dynamics, mitigate OOD effects, and preserve key theoretical properties.
- Extensive experiments show that **SWING consistently outperforms state-of-the-art time series XAI methods** across diverse benchmarks, backbones, and metrics.



Paper



Code