



Multivariate Time Series Anomaly Detection and Interpretation using Hierarchical Inter-Metric and Temporal Embedding

Zhihan Li, Youjian Zhao, Jiaqi Han, Ya Su, Rui Jiao, Xidao Wen, and Dan Pei. 2021.
In Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '21)
2025년 2월 18일

이규원

Department of Computer Science and Engineering
Chung-Ang University

Index

- **Backgrounds**
- **Motivation**
- **Model Architecture**
- **Experimental Results**
- **Conclusion**

Backgrounds

- Multivariate Time Series

- Involves more than one variable over time
- High dimensional data

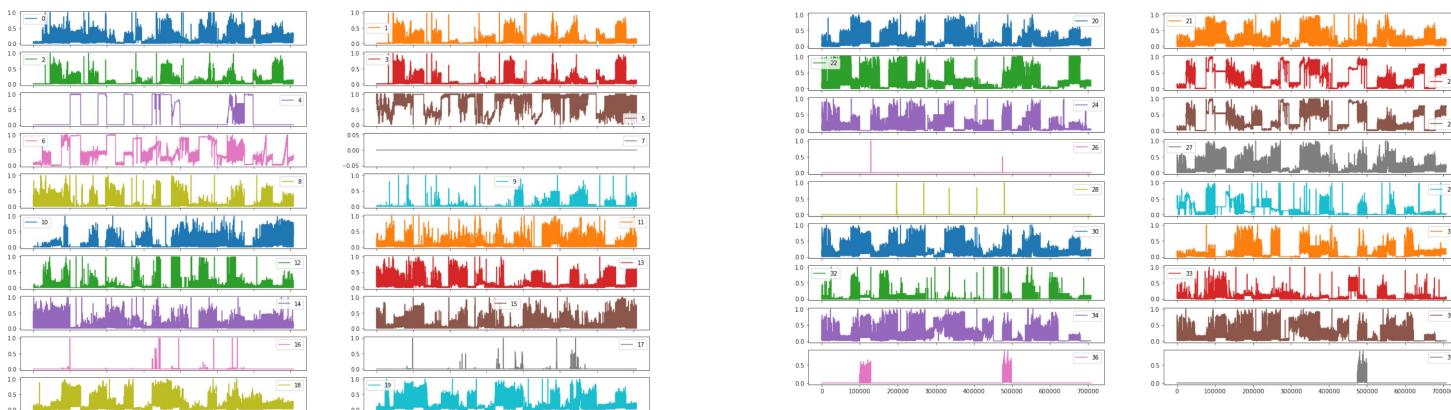


Figure 6: Waveform of SMD dataset.

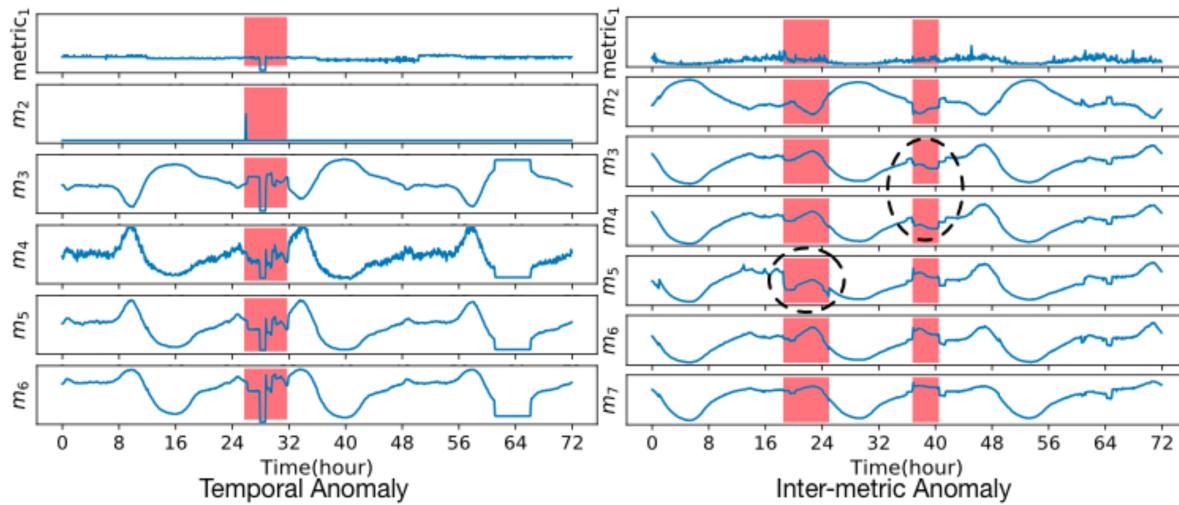
Backgrounds

- **Time Series Anomaly Detection**

- Identifying patterns or events in time-ordered data that deviate significantly from expected behavior or norms

- **Time Series Anomaly**

- Temporal: Significantly deviate from their historical patterns
- Inter-metric: The relationships among different metrics violate their expected correlations



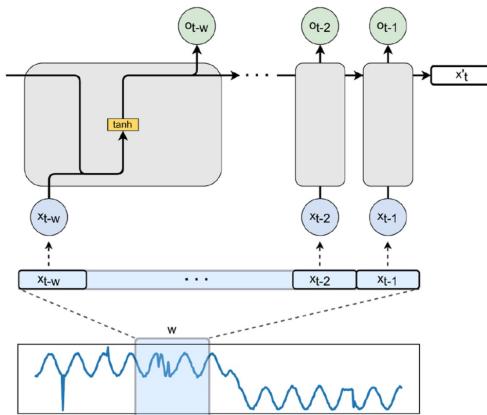
Backgrounds

- **Forecasting-based Model**

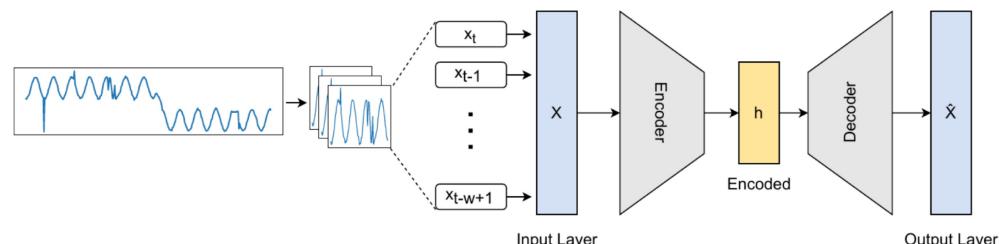
- To predict a future point or subsequence based on a point or a recent window

- **Reconstruction-based Model**

- To reconstruct normal data well while failing to reconstruct anomalous data



Forecasting



Reconstruction

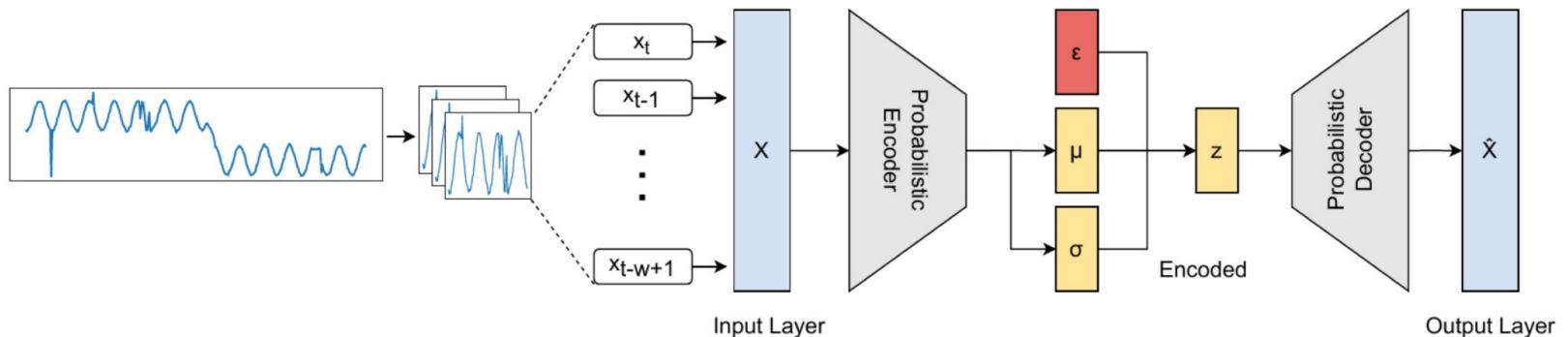
Backgrounds

- **Stochasticity**

- Uncertainty (factors we cannot provide as input) affects time-series data, causing variations

- **Variational Auto Encoder**

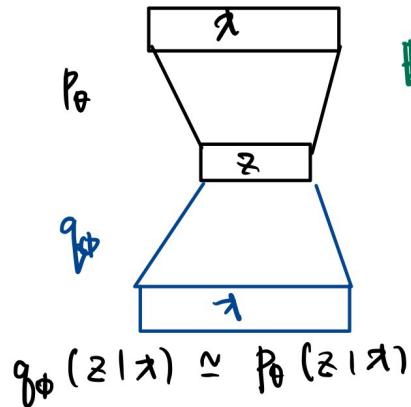
- Decoder: Finds the generative model $p_\theta(x|z)$ that makes data most probable from latent space z
- Encoder: Uses $q_\phi(z|x)$ to approximate $p_\theta(z|x)$ due to its intractability
- Estimate parameters of data distributions at each timestamp (e.g., μ, σ of a Gaussian)



(b) Variational Auto-Encoder

Backgrounds

- Variational Auto Encoder – ELBO



$$\mathbb{E}_{z \sim q_\phi(z|x)} \left[\log \frac{q_\phi(z|x^{(i)})}{p_\theta(z)} \right]$$

$$= \int_z q_\phi(z|x^{(i)}) \log \frac{q_\phi(z|x^{(i)})}{p_\theta(z)} dz$$

$KL(P||Q)$

$$= \sum_x P(x) \log \frac{P(x)}{Q(x)}$$

$\mathbb{E}_{z \sim q_\phi(z|x)} [\log(p_\theta(x^{(i)})]$ 은 최대화해보자.

$$p_\theta(x^{(i)}|z) = \frac{p_\theta(z|x^{(i)}) p_\theta(x^{(i)})}{p_\theta(z)}$$

$$p_\theta(x^{(i)}) = \frac{p_\theta(x^{(i)}|z) p_\theta(z)}{p_\theta(z|x^{(i)})}$$

$$= \mathbb{E}_z \left[\log \frac{p_\theta(x^{(i)}|z) p_\theta(z)}{p_\theta(z|x^{(i)})} \right]$$

$$= \mathbb{E}_z \left[\log \frac{p_\theta(x^{(i)}|z) p_\theta(z)}{p_\theta(z|x^{(i)})} \cdot \frac{q_\phi(z|x^{(i)})}{q_\phi(z|x^{(i)})} \right]$$

$$= \mathbb{E}_z \left[\log p_\theta(x^{(i)}|z) \right] - \mathbb{E}_z \left[\log \frac{q_\phi(z|x^{(i)})}{p_\theta(z)} \right] + \mathbb{E}_z \left[\log \frac{q_\phi(z|x^{(i)})}{p_\theta(z|x^{(i)})} \right]$$

(lower bound)

$$= \mathbb{E}_z \left[\log p_\theta(x^{(i)}|z) \right] - D_{KL}(q_\phi(z|x^{(i)}) || p_\theta(z)) + D_{KL}(q_\phi(z|x^{(i)}) || p_\theta(z|x^{(i)}))$$

Reconstruction loss.

*시뮬레이션 prior를 최적화해라.
→ ex? 오버파팅 방지*

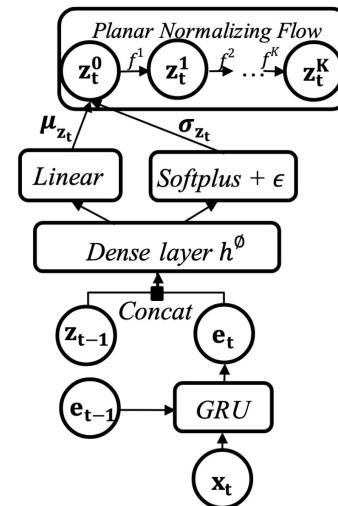
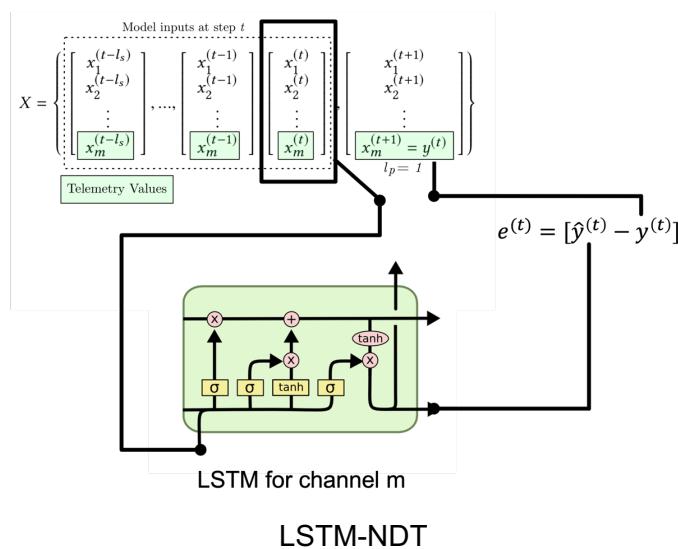
$D_{KL} \geq 0$

p_θ 는 Gaussian, Bernoulli decoder

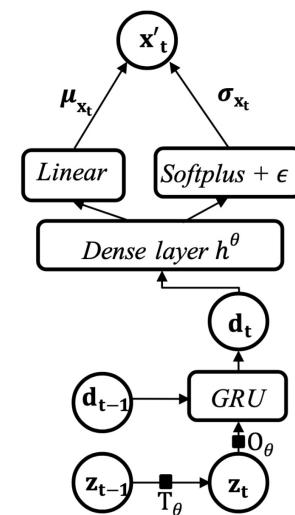
Motivation

- Previous Methods

- LSTM-NDT
- MSCRED
- LSTM-VAE, OmniAnomaly
- MADGAN



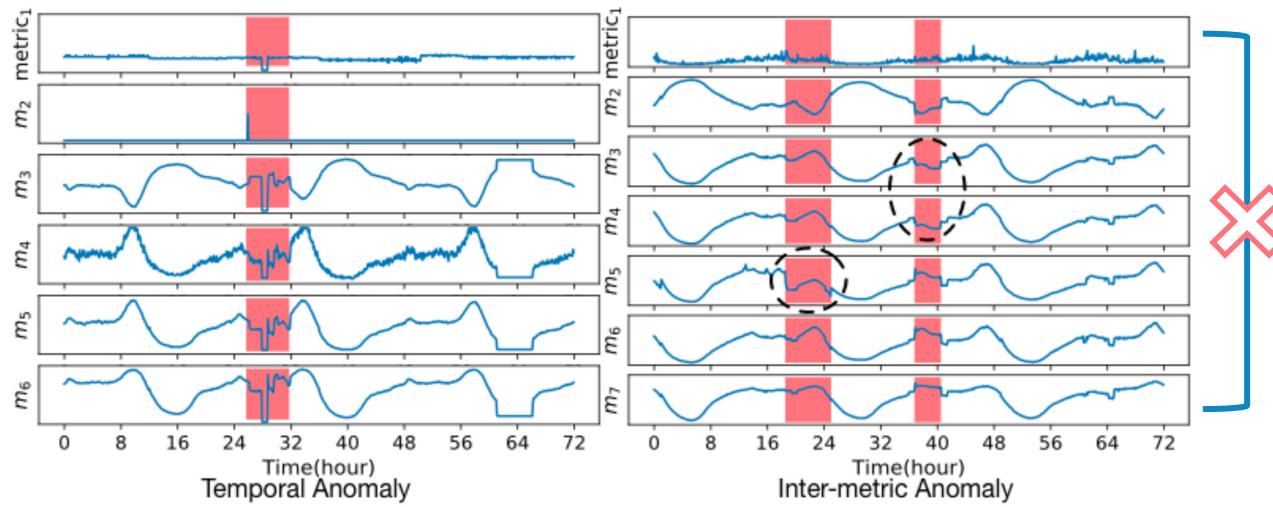
OmniAnomaly



Motivation

- **Limitation**

- Modeled either **temporal dependency** or **intermetric dependency**
- MSCRED Captures intermetric **correlation** rather than learning explicit intermetric embeddings



Motivation

- **Purpose**

- Explicitly learn the low-dimensional intermetric and temporal representations
- Robustness to noise and anomalies in training data
- Improving interpretability of anomaly detection

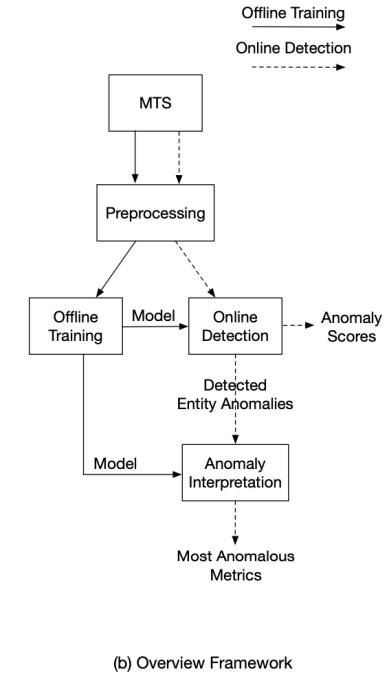
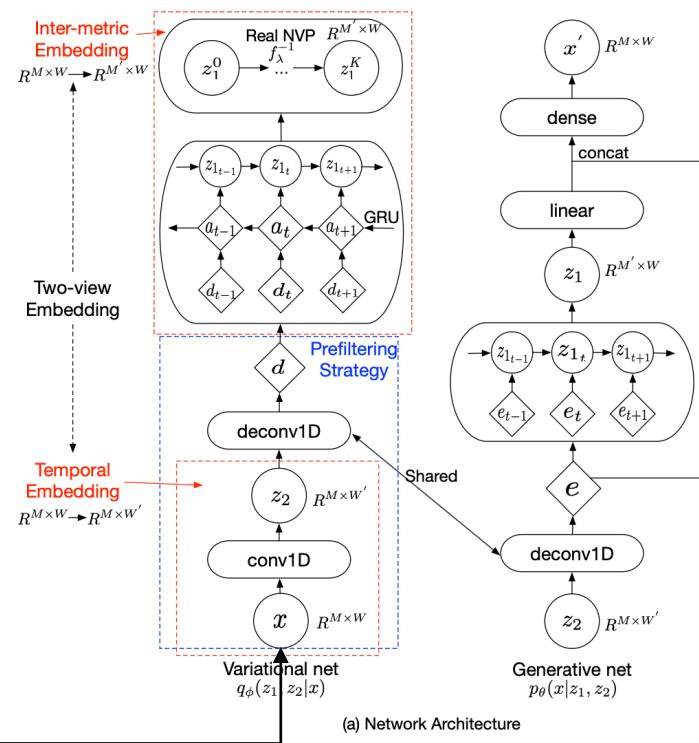
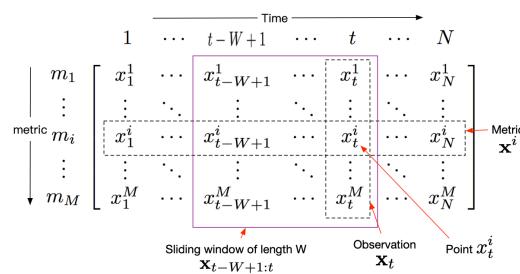
- **Idea**

- Model based on HVAE and Two-view Embedding
- Prefiltering Strategy
- MCMC based Interpretation

Model Architecture

- Idea

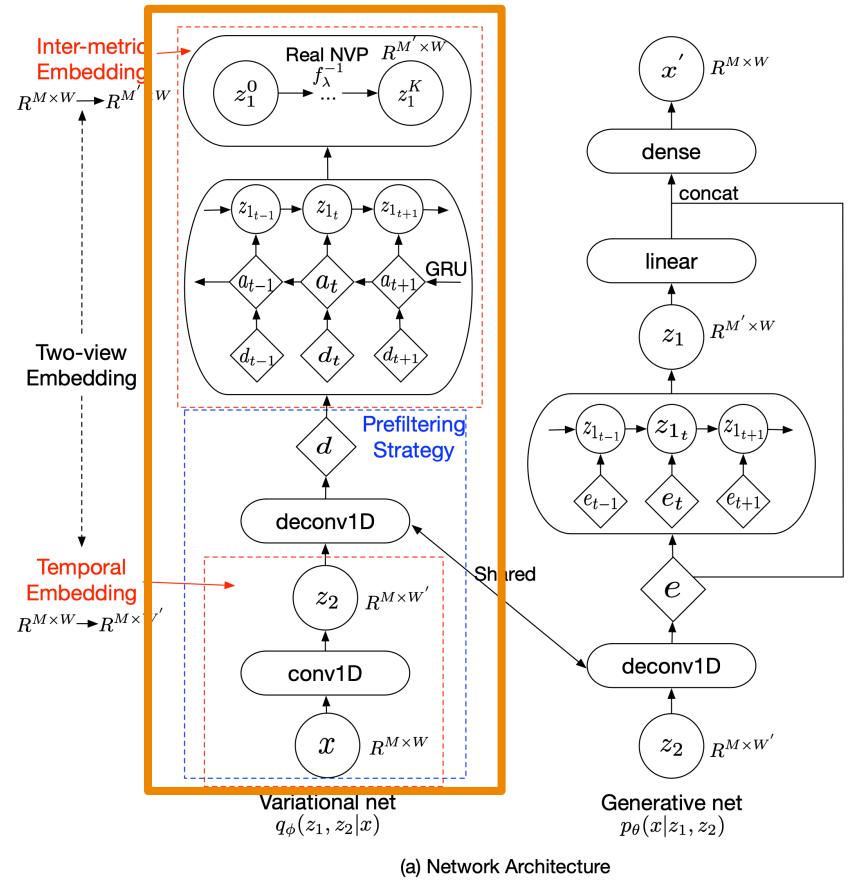
- Model based on HVAE and Two-view Embedding
- Prefiltering Strategy
- MCMC based Interpretation



Model Architecture

- HVAE

- Temporal → Intermetric
- Rather than learning independently

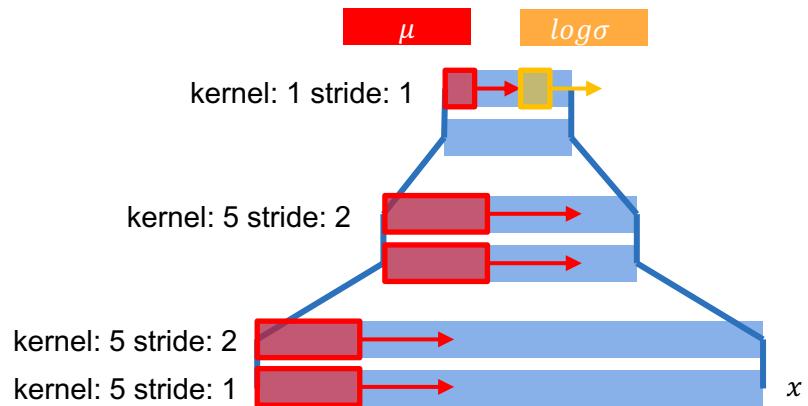
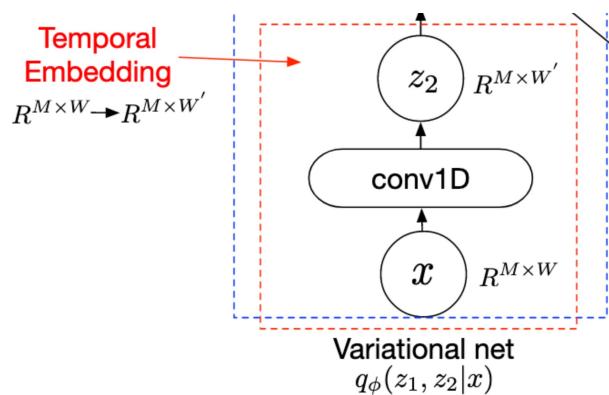


(a) Network Architecture

Model Architecture

- HVAE – Temporal Embedding

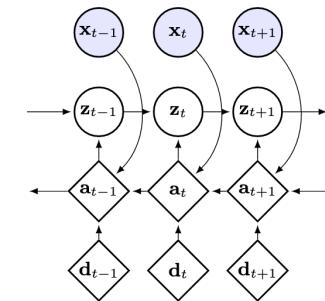
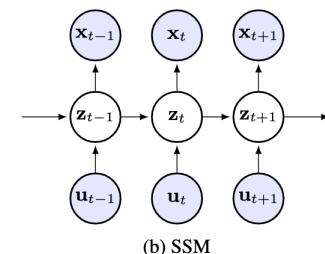
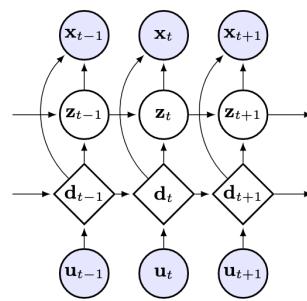
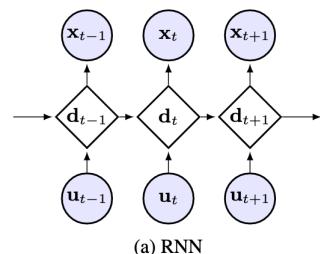
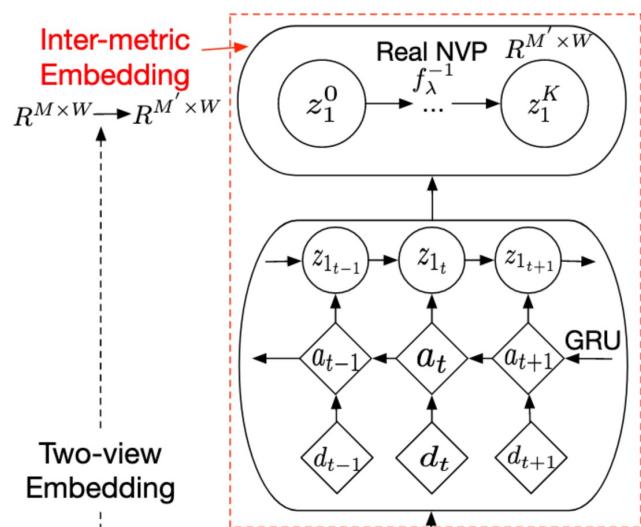
- Captures the temporal dependency using Conv1D layers along the time dimension



Model Architecture

- HVAE – Intermetric Embedding

- Represents the intermetric dependency using an SRNN-like architecture



Stochastic RNN

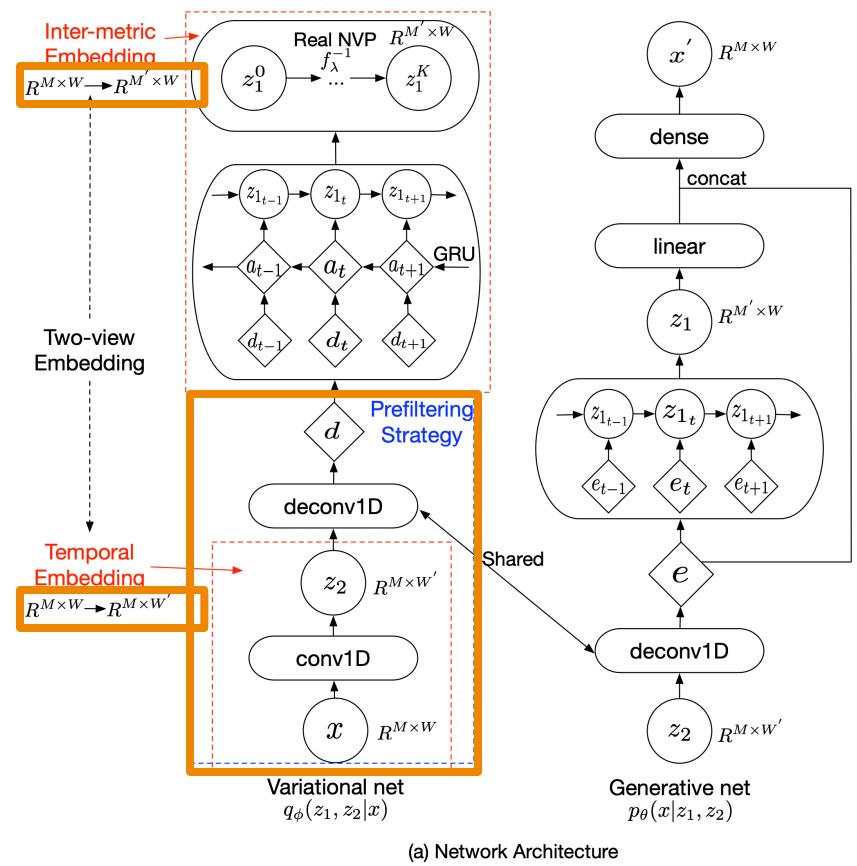
Model Architecture

- **Two-view Embedding**

- Derives intermetric embeddings using the reconstructed d $R^{M \times W}$
- To preserve time consistency

- **Prefiltering Strategy**

- Filtering noise and anomalies in real world training data
- z_1 derived from reconstructed d
- d is pretrained for initial reconstruction

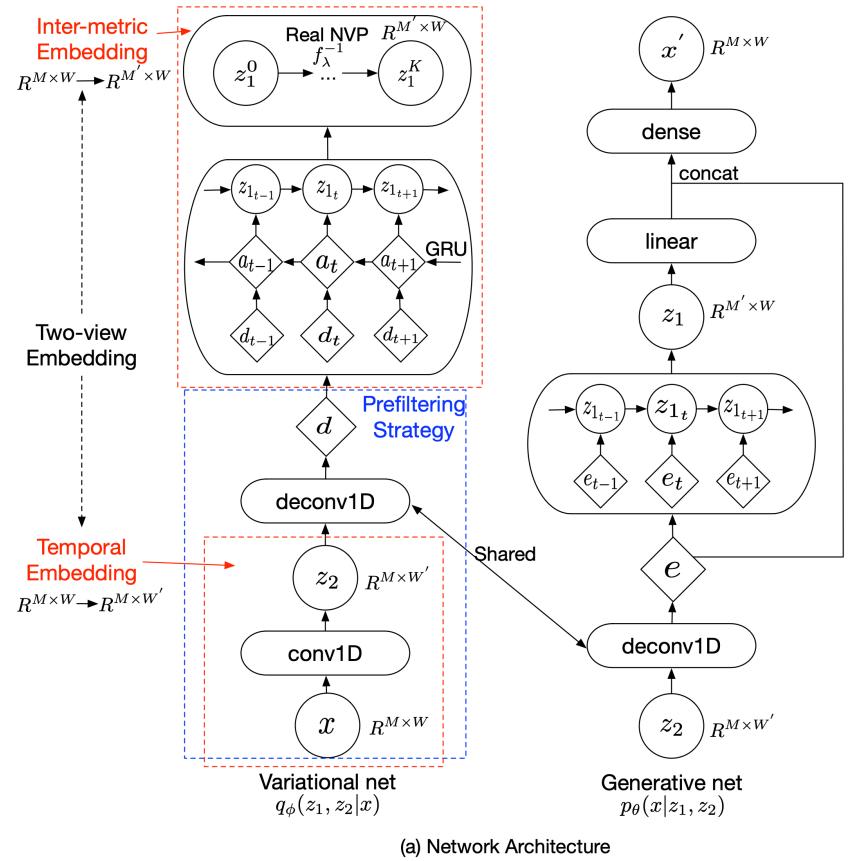


(a) Network Architecture

Model Architecture

- **Generative net**

- p_θ uses z_1, z_2 sampled from q_ϕ
- $p_\theta(z_1)$ without backward GRU
- $p_\theta(z_2)$ is Normal distribution
- $p_\theta(z_1), p_\theta(z_2)$ are used in ELBO calculation



$$\mathcal{L}(\mathbf{x}, \theta, \phi) = \mathbb{E}_{q_\phi(\mathbf{z}_1, \mathbf{z}_2, \mathbf{d}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z}_1, \mathbf{z}_2, \mathbf{e})] - D_{KL}(q_\phi(\mathbf{z}_1, \mathbf{z}_2, \mathbf{d}|\mathbf{x}) || p_\theta(\mathbf{z}_1, \mathbf{z}_2, \mathbf{e}))$$

Model Architecture

- Interpretation (MCMC imputation)

- Masking anomalous data points
(metric, time)
- Infer the normal values of severe anomalies
- \tilde{x} : Anomalous parts have been replaced with inferred normal values
- Perform $\mathbb{E}_{q_\phi(z_1, z_2 | \tilde{x})} [\log p_\theta(\textcolor{red}{x} | z_1, z_2)]$
- Detect both severe and subtle anomalies

Algorithm 1: *InterFusion Anomaly Interpretation*

Input: input sequence $\mathbf{x} \in R^{M \times W}$, original reconstruction probability \mathbf{r}^0 , normal baseline b , window length W , number of metrics M , small constant ratio $\beta_{init}, \beta_{inc}$

Output: revised anomaly score $AS \sim R^{M \times W}$ for interpretation

$n_p \leftarrow$ number of points (\mathbf{x}_m, t) where $\mathbf{r}_{m,t}^0 < \frac{b}{M*W}$;

$n_{init} \leftarrow \beta_{init} n_p$, $n_{inc} \leftarrow \beta_{inc} n_p$, $n \leftarrow n_{init}$, $r^a = \sum_{m,t} \mathbf{r}_{m,t}^0$;

while not $(r^a \geq b \text{ or } n > n_p)$ **do**

- $\mathbf{x}_m \leftarrow$ top n points in \mathbf{x} that have the lowest $\mathbf{r}_{m,t}^0$;
- $\mathbf{x}_o \leftarrow$ other points in \mathbf{x} but not in \mathbf{x}_m ;
- Denote $\mathbf{x}' = \mathbf{x} = (\mathbf{x}_o, \mathbf{x}_m)$;
- for** $s \leftarrow 1$ to S **do** // MCMC imputation for S times
 - sample $(\mathbf{z}_1, \mathbf{z}_2)$ from $q_\phi(\mathbf{z}_1, \mathbf{z}_2 | \mathbf{x}_o, \mathbf{x}_m)$;
 - reconstruct $(\mathbf{x}'_o, \mathbf{x}'_m)$ from $p_\theta(\mathbf{x}'_o, \mathbf{x}'_m | \mathbf{z}_1, \mathbf{z}_2)$;
 - update $\mathbf{x}' \leftarrow (\mathbf{x}'_o, \mathbf{x}'_m)$;
- end**
- /* Approximate the true reconstruction prob of the input window using revised \mathbf{x}' */
- $r^a = \frac{M*W}{M*W-n} \mathbb{E}_{q_\phi(\mathbf{z}_1, \mathbf{z}_2 | \mathbf{x}')} [\sum_{\mathbf{x}_i \in \mathbf{x}_o} \log p_\theta(\mathbf{x}_i | \mathbf{z}_1, \mathbf{z}_2)]$;
- add r^a to rlist, $n \leftarrow n + n_{inc}$;
- end**

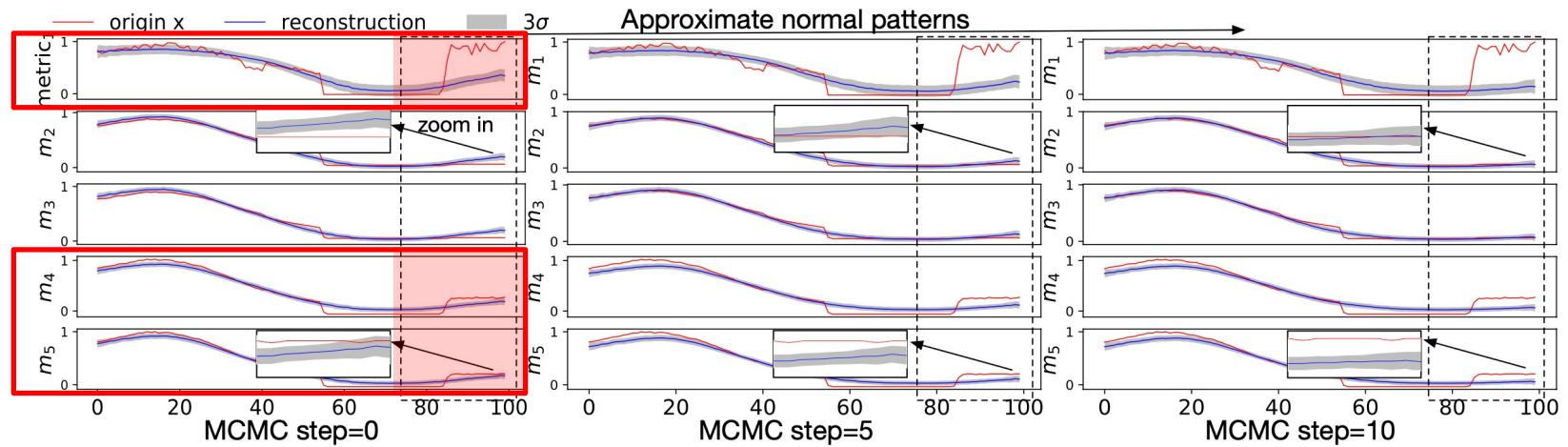
$\tilde{\mathbf{x}} \leftarrow \mathbf{x}'$ that achieves the highest r^a in rlist;

$\mathbf{r}^f = \mathbb{E}_{q_\phi(\mathbf{z}_1, \mathbf{z}_2 | \tilde{\mathbf{x}})} [\log p_\theta(\mathbf{x} | \mathbf{z}_1, \mathbf{z}_2)]$, $AS = -\mathbf{r}^f$;

Model Architecture

- Interpretation (MCMC imputation)

- Masking anomalous data points (metric, time)
- Infer the normal values of severe anomalies
- \tilde{x} represents the input where anomalous parts have been replaced with inferred normal values
- Perform $\mathbb{E}_{q_\phi(z_1, z_2 | \tilde{x})} [\log p_\theta(x | z_1, z_2)]$
- Detect both severe and subtle anomalies



Experimental Results

- **Evaluation Metric**

- Point Adjust Precision/Recall
- Interpretation Score(Proposed)

- **Thresholding Method**

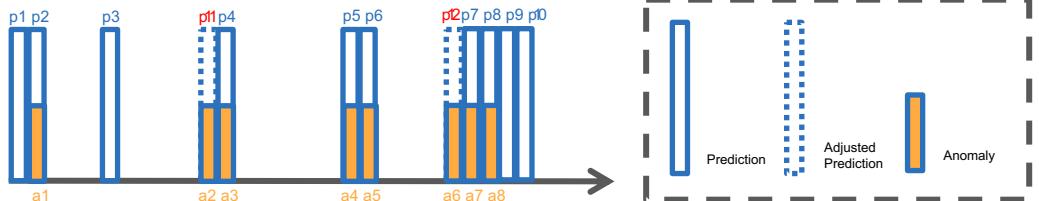
- Best F1 score

$$IPS = \sum_{a=1}^A w_a \frac{|G_{\Phi_a} \cap I_{\Phi_a}|}{|G_{\Phi_a}|}, \quad w_a = \frac{N_{\phi_a}}{\sum_{a=1}^A N_{\phi_a}}$$

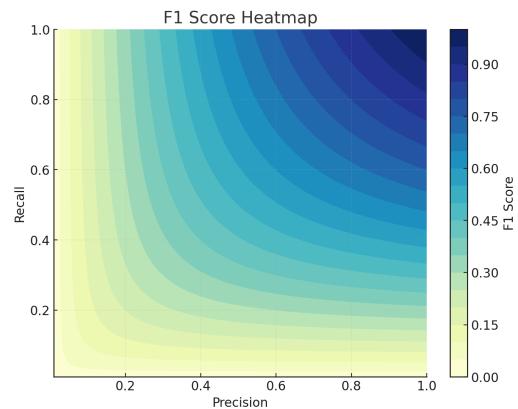
Interpretation Score

PA Precision : $\frac{\text{Correct Predictions}}{\text{Total Predictions}}$
 $(\text{Predictions} = \text{Predictions} + \text{Adjusted Predictions})$

PA Recall : $\frac{\text{Detected Anomalies}}{\text{Total Anomalies}}$



PA Precision = 0.75, PA Recall = 1
F1 Score = 0.857



ϕ_a : An anomaly segment

G_{ϕ_a} : Ground-truth anomalous metric set for segment ϕ_a

I_{ϕ_a} : The set of top k anomalous metrics in the segment ϕ_a

$k = |G_{\phi_a}|$, $w = \text{length weight}$

Experimental Results-F1 Score

Table 1: Average best-F1 for *InterFusion* and baselines.

Methods	SWaT	WADI	SMD	ASD	Avg.
LSTM-NDT	0.8133	0.5067	0.7687	0.4061	0.6237
MSCRED	0.8346	0.5469	0.8252	0.5948	0.7004
MAD-GAN	0.8431	0.7085	0.8966	0.6325	0.7702
OmniAnomaly	0.7344	0.7927	0.9628	0.8344	0.8311
DSANet	0.8924	0.8739	0.9630	0.8740	0.9008
USAD	0.8227	0.4275	0.9024	0.7987	0.7378
VAEpro	0.8369	0.8200	0.8693	0.8522	0.8446
<i>InterFusion</i>	0.9280	0.9103	0.9817	0.9531	0.9433

Experimental Results-Interpretation Score

Table 2: Interpretation IPS for *InterFusion* and baselines.

Methods	SMD	ASD	Avg.
LSTM-NDT	0.5751	0.8619	0.7185
MSCRED	0.6421	0.7652	0.7037
OmniAnomaly	0.8008	0.8029	0.8019
DSANet	0.6713	0.8123	0.7418
VAEpro	0.5681	0.8236	0.6959
VAEpro*	0.7433	0.8916	0.8175
InterFusion-nI	0.7752	0.8881	0.8317
<i>InterFusion</i>	0.8340	0.9107	0.8724

No MCMC imputation
Original reconstruction probability

Experimental Results-Ablation Study

Model Variant	Description
TimeVAE	Only Temporal
m-SRNN	Only Intermetric
IF-p	No HVAE
IF-s	No Two-view Embeddings
IF-x	No Prefiltering Strategy
PureAE	AutoEncoder
IF-AERNN	AutoEncoder + RNN

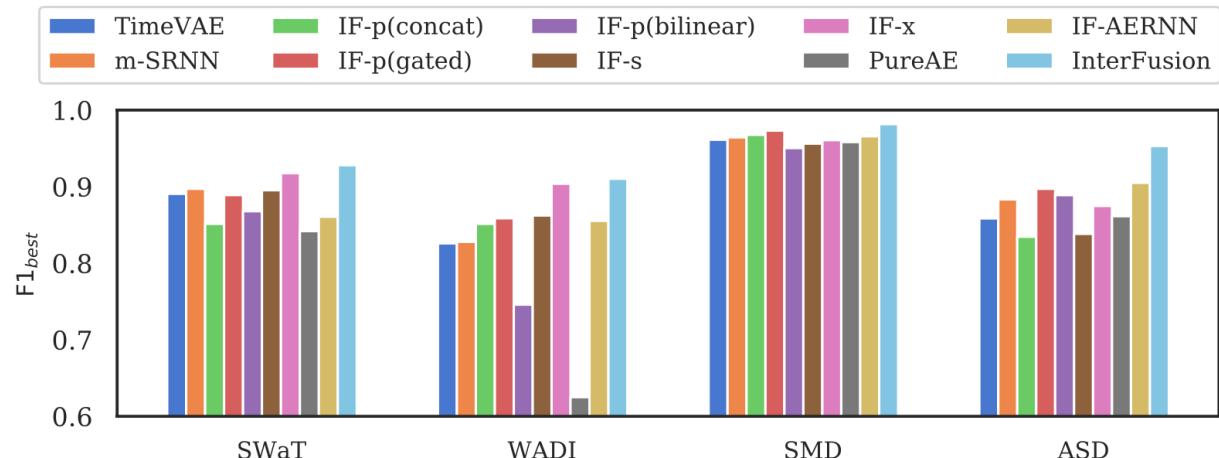


Figure 4: Average anomaly detection best-F1 for *InterFusion* and its variants. ‘IF’ denotes *InterFusion* for short.

Conclusion

- **Capturing both temporal and intermetric dependencies**
 - HVAE, Two-view Embedding
- **Filtering noise and anomalies in training data**
 - Prefiltering Strategy
- **Enhancing interpretability**
 - MCMC-based Interpretation



Thank you