

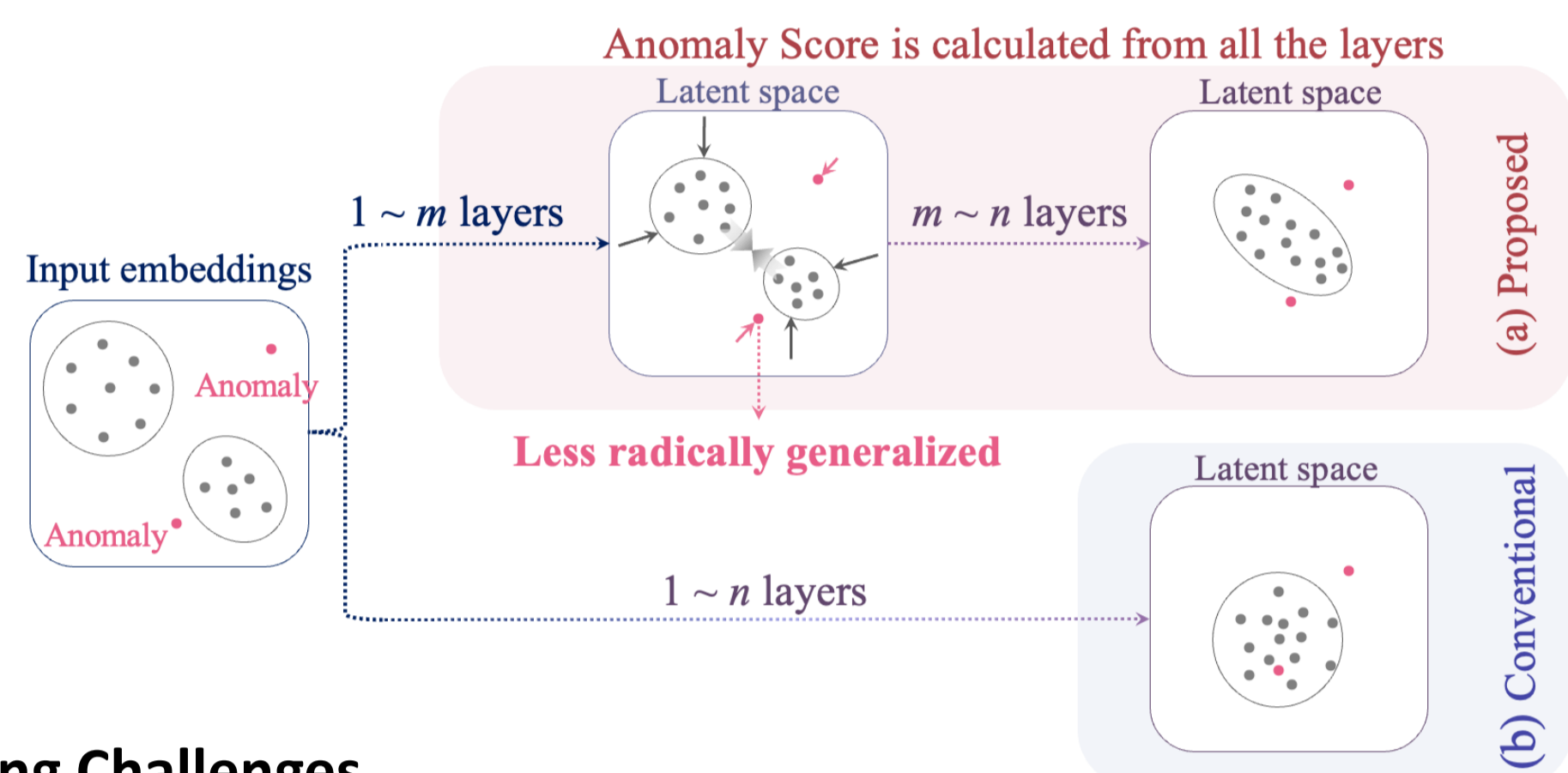
Introduction

Video Anomaly Detection

- Anomaly detection is identifying **atypical patterns** that diverge from the majority
- Recent advances in deep learning have sparked increasing interest in utilizing **video data for anomaly detection**
- Video anomaly detection methods [1,2] can identify **abnormal behaviors in video footage**, and reduce the workload of human operators

Main Challenges

- Challenge 1: Extracting Spatio-temporal Representations from Video Footage**
 - The identification of specific situations often hinges on a **complex mix of spatial and temporal data**
- Challenge 2: Navigating the Complexity of Data in Anomaly Detection Models**
 - Video recordings capture a broad spectrum of real-world dynamics and result in complex and often **non-uniform data distributions**
 - Traditional anomaly detection methods, which typically **generalize data into a singular distribution**, are prone to high false positive and false negative rates



Overcoming Challenges

- We leverage a **pre-trained 3D-CNN** (Convolutional Neural Network) [3] to extract spatio-temporal representations
- We propose a **novel Transformer [4]-based Autoencoder (AE)** to deal with complex data distributions
 - The encoder in our framework progressively generalizes the training data employing **self-attention mechanism**
 - The level of generalization for individual data points is indirectly assessed through **attention weight** at each self-attention layer

Remarks

- Discerning the context of **human behavior from a single frame is challenging**, so we segment continuous video data into uniform intervals along the temporal axis
- We define a **single data point as a segment** and an anomaly data point as a segment that contains at least one anomaly frame

Related Work

Self-attention Mechanism

- By employing **parallelizable self-attention mechanism**, Transformer [4] can consider **global dependencies** within the data and reflect the complexity and context of input
- Multi-head attention performs self-attention using multiple sets of weights to facilitate attention operations from **various perspectives**

$$MHA(X) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

$$\text{head}_i = \text{Self-Attention}_i(X)$$

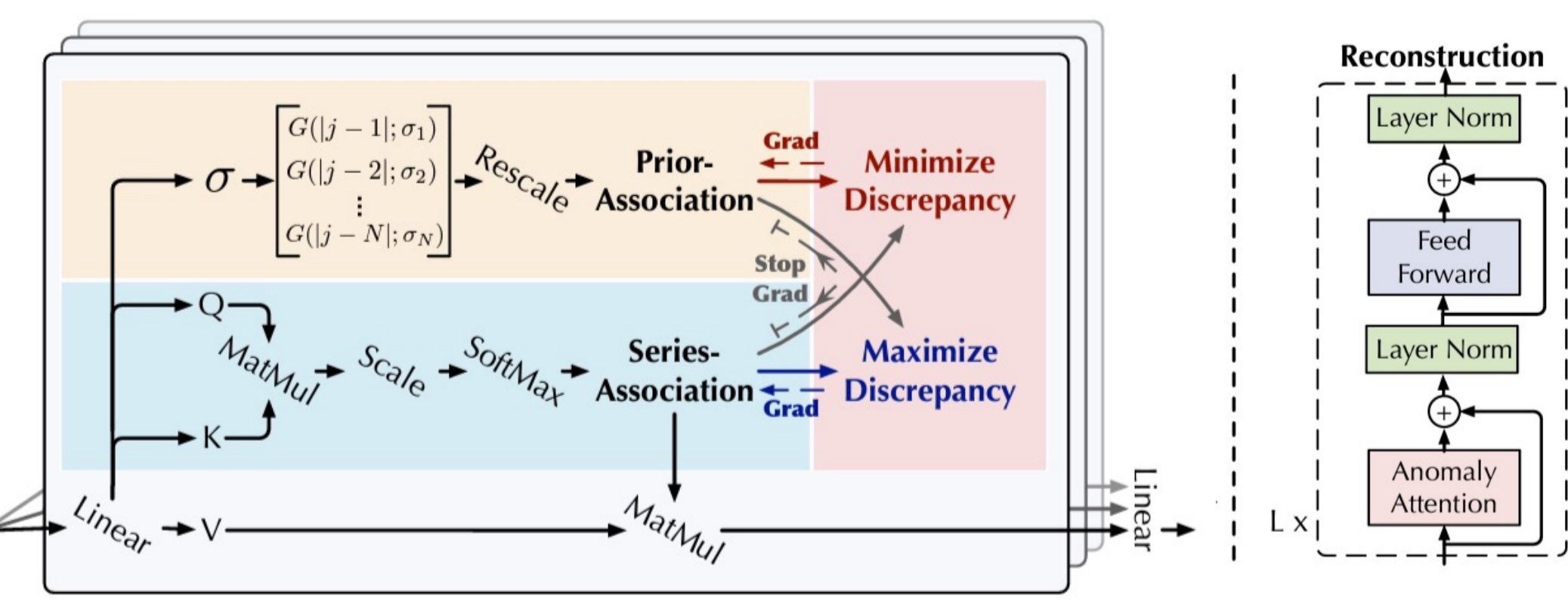
$$\text{Self-Attention}_i(X) = \text{softmax}(QK^T / \sqrt{d_k})V$$

$$Q = XW_i^Q, K = XW_i^K, V = XW_i^V$$

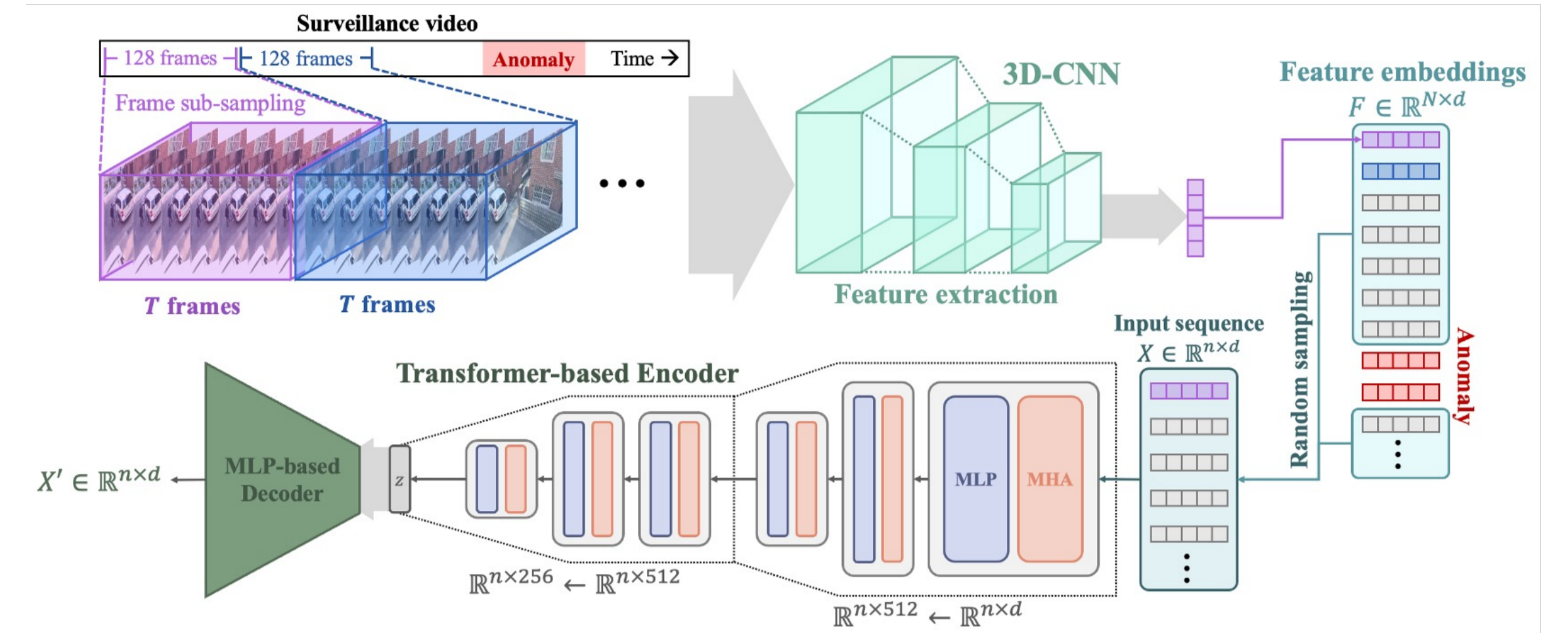
Anomaly Transformer

- Anomaly Transformer [5] is **anomaly detection approach in time-series data**
- This approach tends to **make input data points appear more similar** by applying the self-attention (weighted sum) operation across multiple windows
- This may result in **overly simplified generalization**

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



Proposed Method



Feature Extraction

- To effectively capture spatio-temporal information from the segments, we employ a **3D-CNN pre-trained on large-scale action recognition datasets**, such as Kinetics-400 [6] and Charades [7]

Transformer-based Autoencoder

- Key Features Distinguished from Conventional Transformer
 - The MLP (Multi-layer Perceptron) layers within our encoder are designed to **gradually reduce the dimensionality** of the output space at the 3rd and 6th layers
 - We employ **random sampling** from all input segments to generate input sequences
 - The decoder is designed as an **MLP structure** to simply match the output dimensions

Loss Function

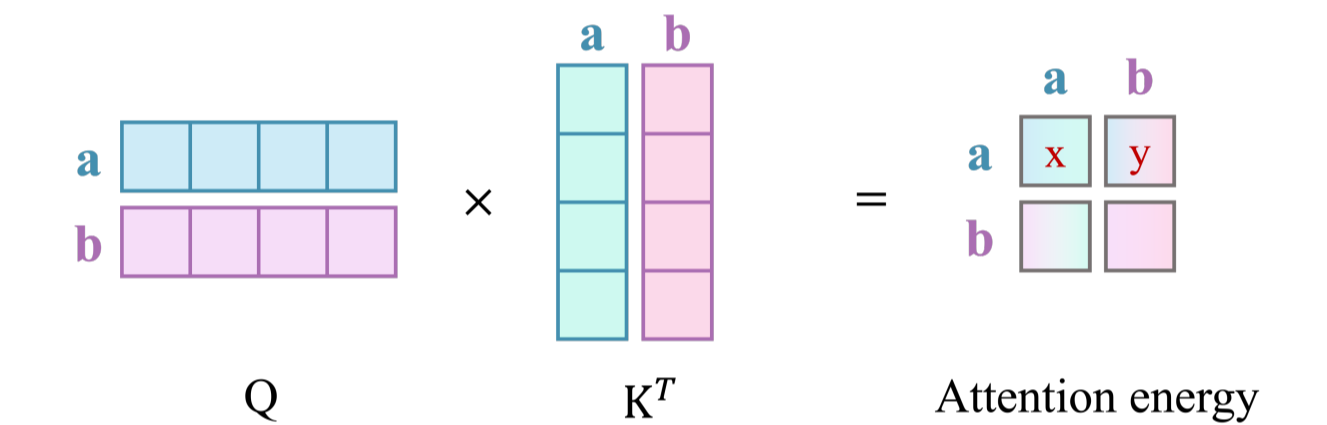
$$\text{Reconstruction loss} = \frac{1}{n} \sum_{i=1}^n \|x_i - x'_i\|^2$$

Detecting Anomalies

- Self-attention is designed to allocate **higher weights to input segments that exhibit greater similarities** within the sequence
- Due to their uniqueness and the scarcity of similar samples, **anomalies tend to draw high attention among themselves**
- By exploiting this, we propose a novel method for determining the anomaly score, which involves **aggregating the attention weights** across all encoder layers

$$X = \{x_0, x_1, \dots, x_{n-1}\}$$

$$\text{Score} = \frac{1}{S} \cdot \frac{1}{L} \cdot \frac{1}{H} \sum_{s=1}^S \sum_{l=1}^L \sum_{h=1}^H \text{Attention}(x_0)^{(s,l,h)}$$



Experiments and Results

Dataset

- Abnormal Behavior CCTV Video Dataset** [8] provided by the South Korean National Information Society Agency (NIA)
- Three abnormal behavior types: trespassing, fighting and vandalism

Video specification of each behavior type		
time of day	length (s)	frame rate (fps)
daytime	~ 3600	30

Data Preprocessing

- We segment the video data into intervals of 128 frames, with the number of data points (N) for each type of behavior
- The input frames are resized and sampled at regular intervals to match the input shape required by each pre-trained 3D-CNN model used in the experiments

Performance Comparison

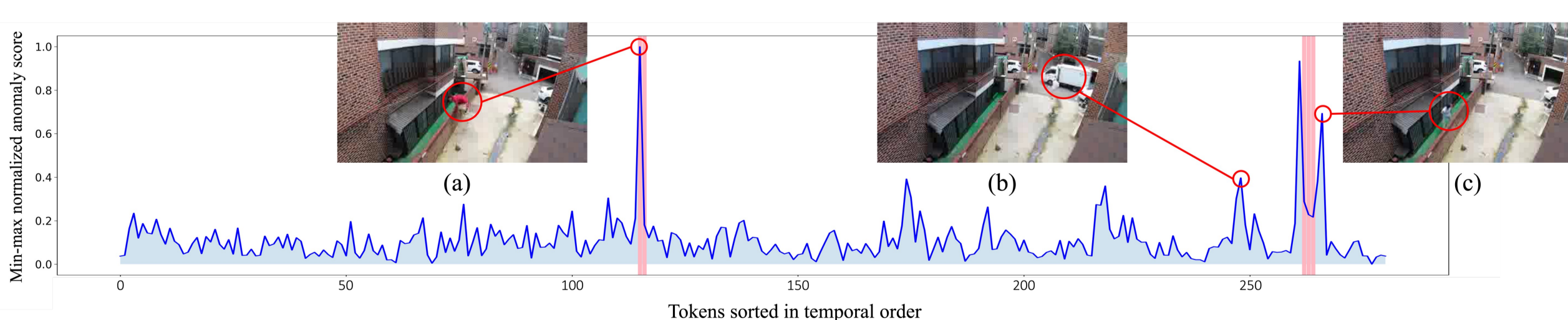
- Our method **outperforms the existing solutions**
- This indicates that our method not only precisely identifies the anomalies but also minimize false alarms and false negatives across different datasets

Visualization of Attention Weights

- 1-3rd layers**: suggests that our approach can extract meaningful anomaly scores **at each stage of the generalization process**
- 6th layer**: while the data begin to show less varied distributions compared to earlier layers, **higher attention weights continue to be assigned to anomalies** and their vicinity

Case Analysis

- (a): the precise moment when a person dressed in red jumps over the fence **was successfully pinpointed**
- (b): the detection was triggered by a white truck passing by, likely due to the scarcity of such events in the training data
- (c): despite being labeled as normal in the dataset, it can be considered as a potential anomaly



Dataset	Trespassing	Vandalism	Fighting
N	1,419	1,432	1,389
<i>anomaly ratio</i>	0.016	0.049	0.040
AE	0.907 / 0.343	0.839 / 0.224	0.895 / 0.330
VAE	0.907 / 0.300	0.834 / 0.219	0.892 / 0.321
AE-SVDD [17]	0.864 / 0.263	0.892 / 0.226	0.911 / 0.262
VAE-SVDD [22]	0.845 / 0.388	0.802 / 0.197	0.731 / 0.119
Anomaly Transformer [21]	0.737 / 0.326	0.724 / 0.144	0.790 / 0.162
VATMAN (ours)	0.913 / 0.278	0.899 / 0.274	0.944 / 0.455

Table 1. AUROC/AUPRC performances of compared methods (N : total number of data points)

	Input shape ($T \times \text{height} \times \text{width}$)	Dataset		
		trespassing	vandalism	fighting
P3D [16]	(16×160×160)	0.431 / 0.015	0.762 / 0.128	0.502 / 0.076
I3D [5]	(64×224×224)	0.842 / 0.324	0.878 / 0.305	0.826 / 0.245
S3D [20]	(64×224×224)	0.772 / 0.294	0.919 / 0.357	0.594 / 0.078
X3D [7]	(16×224×224)	0.880 / 0.367	0.897 / 0.236	0.915 / 0.383
TimeSformer [3]	(16×224×224)	0.690 / 0.172	0.755 / 0.173	0.684 / 0.142

Table 2. AUROC/AUPRC performances with difference features

