

Learning Prompt-Enhanced Context Features for Weakly-Supervised Video Anomaly Detection

Yujiang Pu¹, Xiaoyu Wu¹, Lulu Yang, and Shengjin Wang², *Senior Member, IEEE*

Abstract—Video anomaly detection under weak supervision presents significant challenges, particularly due to the lack of frame-level annotations during training. While prior research has utilized graph convolution networks and self-attention mechanisms alongside multiple instance learning (MIL)-based classification loss to model temporal relations and learn discriminative features, these methods often employ multi-branch architectures to capture local and global dependencies separately, resulting in increased parameters and computational costs. Moreover, the coarse-grained interclass separability provided by the binary constraint of MIL-based loss neglects the fine-grained discriminability within anomalous classes. In response, this paper introduces a weakly supervised anomaly detection framework that focuses on efficient context modeling and enhanced semantic discriminability. We present a Temporal Context Aggregation (TCA) module that captures comprehensive contextual information by reusing the similarity matrix and implementing adaptive fusion. Additionally, we propose a Prompt-Enhanced Learning (PEL) module that integrates semantic priors using knowledge-based prompts to boost the discriminative capacity of context features while ensuring separability between anomaly sub-classes. Extensive experiments validate the effectiveness of our method’s components, demonstrating competitive performance with reduced parameters and computational effort on three challenging benchmarks: UCF-Crime, XD-Violence, and ShanghaiTech datasets. Notably, our approach significantly improves the detection accuracy of certain anomaly sub-classes, underscoring its practical value and efficacy. Our code is available at: <https://github.com/yujiangpu20/PEL4VAD>.

Index Terms—Anomaly Detection, Context Aggregation, Prompt Learning, Weak Supervision.

I. INTRODUCTION

VIDEO Anomaly Detection (VAD) is the process of identifying unusual events or behaviors that diverge from normal patterns in video streams [1]–[3]. With the increasing ubiquity of surveillance cameras, relying solely on human supervision is inadequate for meeting the demands of practical applications. Consequently, there’s a burgeoning need for efficient and accurate automated VAD methods. These methods are crucial across various sectors, including security, industrial monitoring, healthcare, and social media analysis, where they help mitigate potential threats and enhance operational efficiency, thereby holding significant practical value.

Unlike classical action recognition tasks, gathering relevant videos for anomaly detection is challenging due to the rarity

and sensitive nature of anomalous events. The prevailing approach treats anomaly detection as a semi-supervised learning task [4], where the model learns the patterns of normal behavior during training and identifies deviations as anomalies [5]–[10]. These approaches don’t require detailed annotations but often mistakenly identifies unseen samples as anomalies due to the incomplete representation of normal patterns, resulting in high false alarm rates.

Recently, there’s been a notable shift toward developing more robust algorithms for anomaly detection, with a specific focus on weakly supervised methods [11]–[15]. These offer several advantages over semi-supervised approaches: 1) They incorporate both normal and abnormal videos in training, enabling more discriminative representation learning. 2) They require only video-level annotations, indicating the presence but not the exact timing of anomalies. 3) They facilitate the creation of large-scale datasets due to less stringent annotation requirements. Owing to these benefits, weakly supervised methods have demonstrated superior performance, significantly outshining their semi-supervised counterparts.

To effectively model the temporal dynamics of anomalous events, which vary significantly in duration, it’s crucial to understand the contextual information within video snippets. Existing methods, such as graph convolutional networks [16] and self-attention mechanisms [17], aim to capture these temporal relationships. However, each has its limitations: graph convolutions only consider local node aggregates, potentially missing the broader context of anomalies, while self-attention mechanisms, despite capturing global correlations, might introduce irrelevant noise. To refine this understanding, Wu *et al.* [18] and Tian *et al.* [13] devised networks that integrate local and global contexts through parallel structures. Although these approaches enhance context modeling by mitigating long-range noise interference, they come at the cost of increased computational complexity and parameter count, which can lead to overfitting and hinder practical deployment.

Furthermore, the inherent diversity and complexity of anomalies pose additional challenges. Anomaly datasets often showcase large intra-class variance and minimal inter-class variation, with examples ranging from fixed-camera surveillance footage to artistically varied film and TV content. Addressing this requires a robust detection pipeline capable of navigating these complexities. While recent methods aim to construct a compact normal manifold [12] and introduce additional binary constraints [19] for discriminative learning, they often neglect the distinct visual characteristics within anomaly classes and fail to capture the specific semantics of various anomalies. This oversight leads to poor interpretability

Corresponding author: Xiaoyu Wu.

Yujiang Pu, Xiaoyu Wu and Lulu Yang are with the School of Information and Communication Engineering, Communication University of China, Beijing 100024, China (e-mail: pyj2020@cuc.edu.cn; wuxiaoyu@cuc.edu.cn; yangll@cuc.edu.cn).

Shengjin Wang is with the Department of Electrical Engineering, Tsinghua University, Beijing 100084, China (e-mail: wsgsj@tsinghua.edu.cn).

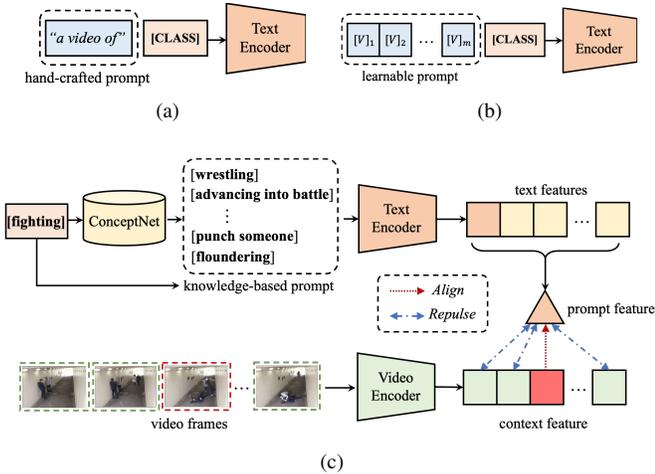


Fig. 1. Comparison of different methods for prompt construction. (a) Hand-crafted prompt. (b) Learnable prompt. (c) Knowledge-based prompt (Ours). The context features are enhanced by cross-modal alignment with the corresponding prompt features.

and an increased risk of misclassification.

In light of these limitations, our motivation is twofold: 1) to develop an approach that balances the need for detailed temporal modeling with computational efficiency, thereby addressing the high parameter count and low efficiency issues inherent in multi-branch structures. 2) to construct a method that recognizes and interprets the diversity and complexity of anomaly semantics more effectively than existing models.

In this paper, we firstly introduce a Temporal Context Aggregation (TCA) module designed to capture temporal relations across video snippets more efficiently than the current parallel architectures [13] [18] [20] [21]. This module distinguishes itself by reutilizing the similarity matrix, thereby reducing the computational load and parameter count. It also incorporates a learnable factor to adaptively fuse local and global contexts and introduces a Dynamic Position Encoding (DPE) with a parametric kernel function to accurately model the relative positions of video snippets. As evidenced by the results in Section IV, our TCA module demonstrates superior performance with a more streamlined and cost-effective design compared to existing methods.

Secondly, we propose a Prompt-Enhanced Learning (PEL) module, drawing inspiration from prompt learning [22]. Unlike existing methods [23]–[26] that rely on hand-crafted or learnable prompts, our approach leverages an external knowledge base, i.e., ConceptNet [27], to construct prompt templates, as demonstrated in Figure 1. These templates are contextually rich and possess a degree of interpretability, offering a nuanced understanding of the specific semantics of anomalies. By aligning anomalous contexts with corresponding prompt features and distancing non-anomalous contexts, our PEL module injects rich semantics into visual features, enhancing fine-grained discriminability and inter-class separability. This innovative approach fosters the development of a more discerning decision boundary in the embedding space.

Finally, considering that anomalies in videos generally extend across several frames, we utilize a plug-and-play score

smoothing (SS) strategy during the testing phase to address potential mis-classifications due to transient disturbances such as frame jitter or switching. This module applies average pooling to consecutive anomaly scores, thereby suppressing individual biases and mitigating the impact of ephemeral noise inconsistent with the typical duration of an anomaly.

The main contributions of this paper are succinctly outlined as follows:

- 1) We introduce a Temporal Context Aggregation (TCA) module that leverages a reused similarity matrix and adaptive fusion for efficient and effective video anomaly detection. This module demonstrates enhanced encoding performance while significantly reducing the parameter count and computational expense.
- 2) We present a Prompt-Enhanced Learning (PEL) module designed to augment fine-grained contextual understanding. This is achieved through the construction of knowledge-based prompts, precise context separation, and cross-modal alignment, thereby enriching the semantic discrimination of the model.
- 3) Our model exhibits competitive performance across three benchmark datasets: UCF-Crime, XD-Violence, and ShanghaiTech. Notably, it shows an approximate 10% improvement in detection accuracy for certain fine-grained anomalies, underscoring its practical efficacy.

The remainder of the paper is structured as follows: Section II reviews related work in video anomaly detection and prompt learning. Section III delineates the proposed method in detail. Section IV presents comparative results with state-of-the-art methods on benchmark datasets and validates the effectiveness of each component through ablation studies. Finally, Section V concludes the paper, reflecting on findings and future directions.

II. RELATED WORK

A. Video Anomaly Detection

Early studies [28]–[31] treated anomaly detection as a semi-supervised task, utilizing statistical models with hand-crafted features. These approaches, however, suffered from limited representational capacity, leading to suboptimal robustness and generalization. Recent advancements in deep learning have spurred the development of more sophisticated anomaly detection techniques. Some methods [32]–[35] aim to create a one-class classifier that delineates normality in a latent space, identifying deviations as anomalies. Yet, without prior knowledge of anomalies, these models may mislabel complex or unseen normal samples as anomalous. Other strategies [36]–[39] postulate that anomalies manifest as larger reconstruction or prediction errors, assuming normality can be accurately reconstructed or predicted. However, these methods are prone to overfitting and may reconstruct or predict anomalies well.

Recent research has pivoted towards weakly supervised anomaly detection, leveraging both normal and abnormal samples with video-level annotations. These methods generally categorize into two paradigms: one-stage methods utilizing Multiple Instance Learning (MIL) and two-stage self-training strategies. Sultani *et al.* [11] introduced a deep MIL model

optimizing a margin ranking loss to enhance the differentiation between positive and negative snippets. Tian *et al.* [13] developed Robust Temporal Feature Magnitude (RTFM) learning, detecting subtle anomalies through feature similarity and temporal continuity. Cho *et al.* [40] devised Class-Activate Feature Learning (CLAV) to extract distinctive features and employed relative distance learning to accentuate their differences, alongside the Context-Motion Interrelation Module (CoMo) for analyzing scene-related motion anomalies. Conversely, Lv *et al.* [41] proposed an unbiased MIL approach that considers both confident and ambiguous snippets, enhancing the detector’s ability to differentiate between normal and abnormal events without relying on contextual biases.

Two-stage self-training methods have emerged to generate high-confidence pseudo-labels for video snippets, recasting weakly supervised anomaly detection as a supervised task with noisy labels. Zhong *et al.* [20] employed a graph convolution network to refine these labels and iteratively enhance the anomaly classifier. Feng *et al.* [42] proposed utilizing multi-instance pseudo label generation to fine-tune feature encoders, thereby yielding task-specific discriminative features. Li *et al.* [43] introduced Multi-Sequence Learning (MSL), a method that incrementally optimizes reduced sample lengths to refine localization boundaries. Zhang *et al.* [15] introduced a multi-head module to generate varied pseudo labels and an iterative uncertainty-based training strategy that refines the process by prioritizing clips with lower uncertainty.

In parallel, several studies have leveraged multimodal information to enhance video anomaly detection. Wu *et al.* [18] proposed HL-Net, integrating appearance, motion, and audio for a comprehensive multimodal approach. Yu *et al.* [44] utilized a dual-stream network to improve instance clustering and employ self-distillation to minimize the semantic gap between unimodal and audio-visual features. Wei *et al.* [45] introduced the MSAF framework, employing multimodal MIL ranking loss to generate pseudo clip-level labels and a supervise-attention regression loss in feature learning to foster implicit alignment between different modalities. Pu *et al.* [46] enhanced contextual video representation through cross-modal interactions between audio and visual elements, utilizing temporal convolution to compute high-confidence scores for specific events like violence.

While numerous methods have explored temporal relation modeling, many depend on parallel branches that introduce additional parameters and computational demands. In contrast, our approach presents two significant advantages. Firstly, the TCA module capitalizes on the reuse of the similarity matrix to concurrently capture local-global dependencies. This strategy outperforms methods requiring extra branches [13], [18], [21], and [47], while minimizing both parameters and computational load. Secondly, diverging from techniques that construct normality prototypes from the data itself, like [12] and [19], our PEL module introduces anomaly priors into the model using external knowledge. This approach not only enriches the semantics of visual features but also notably advances the fine-grained anomaly detection capabilities of our model.

B. Prompt Learning in Video Understanding

Prompt learning, initially introduced in natural language processing, aims to utilize pre-trained models for few-shot or zero-shot scenarios by designing appropriate prompt templates. This technique has recently been extended to video understanding tasks. Wang *et al.* [24] demonstrated this by aligning video clips with prompt embeddings of category labels for action recognition. They incorporated category labels as prefixes, fill-in-the-blanks, and suffixes into prompt templates, thereby significantly outperforming single category prompts. Similarly, Ju *et al.* [26] crafted prompt templates by merging category labels with learnable vectors, examining the impact of label positioning within the templates. Despite these advancements, manual prompt design remains laborious, and models are highly sensitive to template content [48]. While parametric prompts eliminate the need for linguistic priors, they often result in abstract and cryptic strings.

Recently, Yao *et al.* [49] leveraged WordNet’s [50] definitions to formulate prompt templates and created additional category label descriptions using a concept dictionary. This method significantly advanced open-vocabulary object detection tasks. Inspired by their work, we introduce a novel prompt-enhanced learning strategy for video anomaly detection. Unlike previous methods, we extract anomaly-relevant, class-specific concepts from ConceptNet [27] to augment contextual information. Our objective is to improve visual features through external prompts. We achieve this by infusing semantic information into visual features via cross-modal alignment, rather than bidirectional matching, ensuring reliance solely on visual features during testing.

III. THE PROPOSED METHOD

In this section, we will introduce the proposed method for weakly-supervised video anomaly detection. Firstly, we present the overall framework, followed by a detailed elaboration of the core components.

A. Overall Framework

The overall framework of our method is shown in Figure 2. Specifically, an untrimmed video is first divided into non-overlapping snippets by a sliding window of 16 frames. Subsequently, we use a pre-trained I3D network [51] to extract the snippet features \mathbf{X} , which are fed into the TCA module to generate context features \mathbf{X}^c . After that, feature reduction is achieved through a two-layer multilayer perceptron (MLP), with the PEL module applied to the middle layer to learn discriminative features \mathbf{X}^e via knowledge-based prompts. Finally, the classifier predicts snippet-level anomaly scores \mathbf{S} . During the training phase, the MIL-based loss function converts the snippet-level scores into bag-level predictions to learn high activations for anomalous cases, where the subscripts a and n denote abnormal and normal videos, respectively.

B. Temporal Context Aggregation Module

As previously discussed, the Temporal Context Aggregation (TCA) module, depicted in Figure 2, diverges from the conventional multi-branch parallel frameworks used in existing

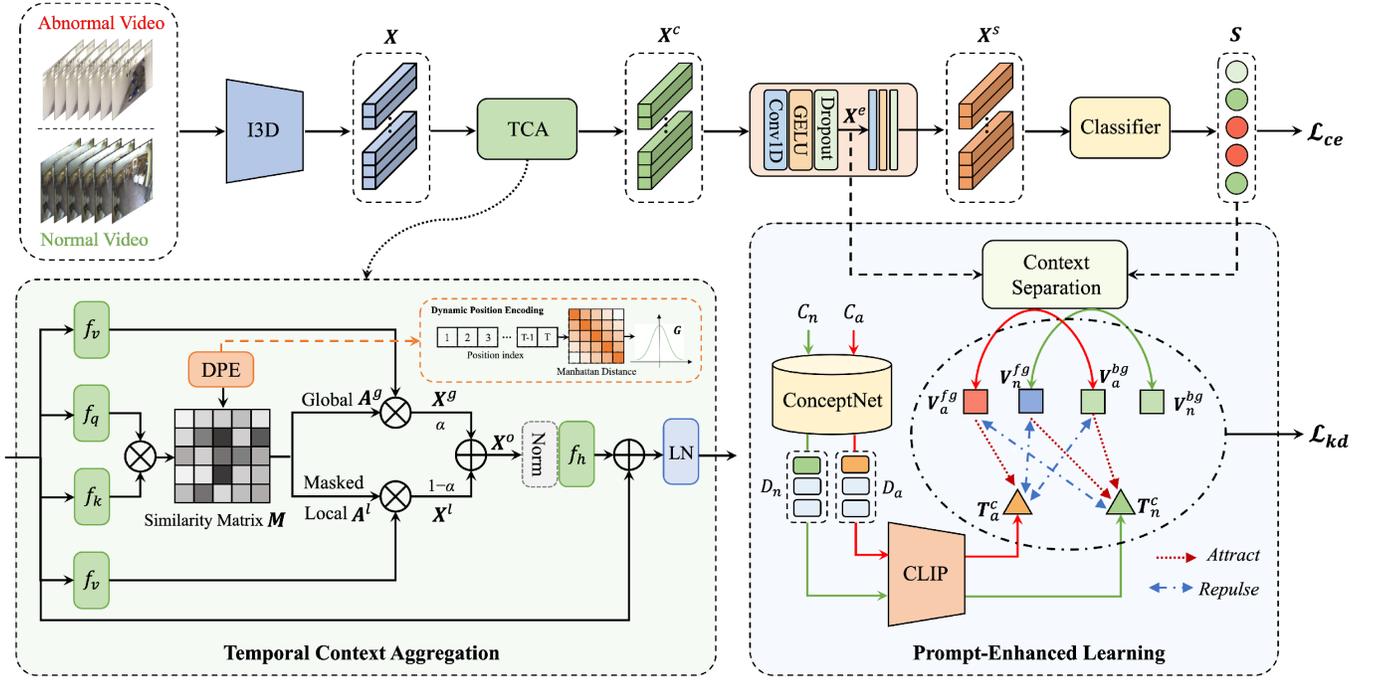


Fig. 2. Overview of the proposed framework. We first process untrimmed videos using a pre-trained I3D network to extract snippet features. The TCA module then simultaneously captures both local and global contexts. High-level representations are derived using a two-layer MLP, with the PEL module operating in the middle layer to learn fine-grained semantics for context features. Finally, a causal convolution layer functions as a classifier to predict snippet-level anomaly scores. During the training stage, the model is optimized using both cross entropy loss \mathcal{L}_{ce} and KL-divergence loss \mathcal{L}_{kd} .

methods. It concurrently models global-local dependencies by leveraging a similarity matrix and adaptive fusion. The snippet feature \mathbf{X} is first projected to the latent space by different linear layers, and the similarity matrix \mathbf{M} is obtained by computing the inner product as follows:

$$\mathbf{M} = f_q(\mathbf{X}) \cdot f_k(\mathbf{X})^\top, \quad (1)$$

$$\mathbf{A}^g = \text{softmax}\left(\frac{\mathbf{M}}{\sqrt{D_h}}\right), \quad (2)$$

$$\mathbf{X}^g = \mathbf{A}^g \cdot f_v(\mathbf{X}), \quad (3)$$

where $f_q(\cdot)$, $f_k(\cdot)$ and $f_v(\cdot)$ are three different linear layers, \top refers to the transpose operation, and D_h is the hidden dimension in the latent space. Softmax normalization is then applied to generate the global attention map \mathbf{A}^g . The projected snippet feature is re-weighted by the attention map to obtain the global context feature \mathbf{X}^g .

Although the above operation facilitates global context modeling, it inevitably introduces long-range noise. To address this issue, we implement local context calibration by reusing the similarity matrix with a masking window, which can be formulated as:

$$\tilde{\mathbf{M}}_{ij} = \begin{cases} \mathbf{M}_{ij}, & \text{if } j \in [\max(0, i - \lfloor \frac{w}{2} \rfloor), \min(i + \lfloor \frac{w}{2} \rfloor, T)] \\ -\infty, & \text{otherwise} \end{cases} \quad (4)$$

$$\mathbf{A}^l = \text{softmax}\left(\frac{\tilde{\mathbf{M}}}{\sqrt{D_h}}\right), \quad (5)$$

$$\mathbf{X}^l = \mathbf{A}^l \cdot f_v(\mathbf{X}), \quad (6)$$

where w is the window size of the mask and T is the maximum length of the input sequence. Eq.(4) ensures that the i^{th} snippet interacts only with its neighborhood of window size w . The lower bound of this window is the earliest moment that can be historically observed, and the upper bound is the maximum length of the sequence. Similarly, the projected snippet feature is re-weighted with the attention map \mathbf{A}^l to obtain local calibration features \mathbf{X}^l , which can effectively capture slight changes and achieve feature enhancement in the local neighborhood.

Subsequently, we resort to a learnable fashion rather than direct concatenation or average pooling to achieve local-global contextual adaptive fusion, allowing the model to dynamically balance the importance of global temporal patterns and local nuances. The fusion process is formulated as follows:

$$\mathbf{X}^o = \alpha \cdot \mathbf{X}^g + (1 - \alpha) \cdot \mathbf{X}^l, \quad (7)$$

$$\mathbf{X}^c = \text{LN}(\mathbf{X} + f_h(\text{Norm}(\mathbf{X}^o))), \quad (8)$$

where α and $1 - \alpha$ denote the global and local weight for context fusion, respectively, and $\text{Norm}(\cdot)$ denote a combination of power normalization [52] and L2 normalization. Then a linear layer $f_h(\cdot)$ followed by residual connection and layer normalization $\text{LN}(\cdot)$ is applied to obtain context feature \mathbf{X}^c .

In addition, considering the importance of position information, we introduce Dynamic Position Encoding (DPE) to model the relative distances of snippets, formulated as:

$$\mathbf{G} = \exp(-|\gamma(i - j)^2 + \beta|), \quad (9)$$

where i and j denote the absolute positions of two snippets, and γ and β are learnable weights and bias terms. In particular,

the DPE is embedded into the similarity matrix \mathbf{M} as a location prior, i.e., $\mathbf{M} \leftarrow \mathbf{M} + \mathbf{G}$, thus avoiding affecting the original feature distribution. Unlike fixed position encodings in [17], DPE adapts to varying video lengths due to its dynamic nature. Also, the Gaussian-like kernel of DPE inherently suppresses the influence of long-distance noise, which emphasizes closer snippet relationships over distant ones implies an innate sensitivity to non-linear patterns.

C. Multilayer Perceptron and Classifier

To obtain high-level semantic representations, a two-layer MLP is utilized for feature reduction. Each Conv1D layer is followed by a GELU activation and dropout operation. This process is denoted as follows:

$$\begin{aligned} \mathbf{X}^e &= \text{Dropout}(\text{GELU}(\text{Conv1D}(\mathbf{X}^c))), \\ \mathbf{X}^s &= \text{Dropout}(\text{GELU}(\text{Conv1D}(\mathbf{X}^e))). \end{aligned} \quad (10)$$

Finally, a causal convolution layer is employed to predict snippet-level anomaly scores, which considers both current and historical observations to obtain more reliable results. The classifier is formulated as:

$$\mathbf{S} = \sigma(f_t(\mathbf{X}^s)), \quad (11)$$

where $f_t(\cdot)$ is the causal convolution layer with kernel size Δt , $\sigma(\cdot)$ is the sigmoid function, and s_i is the anomaly score of the i^{th} snippet.

Following [12], we adopt the MIL-based loss as the basic objective function. Specifically, we determine the video-level prediction p_i by taking the mean value of the top- k anomaly scores. For positive bags, we set $k = \lfloor T/16 + 1 \rfloor$, and for negative bags, we set $k = 1$. Given a mini batch containing B samples with video-level ground-truth y_i , the binary cross-entropy is formulated as:

$$\mathcal{L}_{ce} = \sum_{n=1}^B -y_i \log(p_i). \quad (12)$$

D. Prompt-Enhanced Learning

The PEL module aims to enrich the visual representations by incorporating knowledge-based contextual information, thereby enhancing the model’s ability to identify anomalies in intricate scenarios. Specifically, the PEL module comprises three steps, namely prompt construction, context separation, and cross-modal alignment, as depicted in Figure 2.

1) *Prompt Construction*: Considering the versatility of prompts, we first select 12 common relations¹ from ConceptNet as the pre-retrieval semantics, followed by choosing the top five items with the highest occurrence frequency across all categories as the retrieval relations, denoted as $\{r_j\}_{j=1}^R$. The concept dictionary \mathcal{D} is then constructed by retrieving all edges of the relation r_j established with a given class c as the head or tail node, as illustrated in Figure 3. The non-category node in each edge is used as the key, and the relevance score

¹The relations include: *IsA*, *PartOf*, *HasA*, *UsedFor*, *CapableOf*, *Causes*, *HasSubevent*, *HasPrerequisite*, *HasProperty*, *DefinedAs*, *MannerOf* and *SimilarTo*. The detailed description of each relation can be found at: <https://github.com/commonsense/conceptnet5/wiki/Relations>.



Fig. 3. An example of concept dictionary given the anomaly class *fighting*. The arrows point from the head node to the tail node with relevance scores, and the colors indicate different relationships. The entire graph constitutes a concept dictionary, where the bold items are those retained after node filtering.

of the edge is taken as the value. To remove noisy entries that may interfere with the semantics of anomalies, we first eliminate concepts with relevance scores less than or equal to 0 (Step 1). Next, the remaining entries are filtered by using the average of the relevance scores as the threshold (Step 2). In general, the higher the relevance score, the stronger the semantic relationship between the nodes.

After obtaining the concept dictionary, we use a pre-trained CLIP model [23] to extract the corresponding prompt representations. For a given class c , a set of keys $\{k_i^c\}_{i=1}^N$ from the concept dictionary is first extracted as the context prompt, which is then separately fed into the text encoder to extract 512-dimensional feature vectors $\{\mathbf{K}_i^c\}_{i=1}^N$. Finally, the average of all feature vectors is regarded as the knowledge-based prompt feature, denoted as follows:

$$\mathbf{T}^c = \frac{1}{N} \sum_{i=1}^N \mathbf{K}_i^c, \quad (13)$$

where N denotes the number of keys. The prompt representation incorporates relevant concepts from multiple relations and facilitates the semantic enhancement of visual features.

2) *Context Separation*: Since snippet-level features don’t encapsulate complete contextual details, aligning them directly with prompt features could lead to inaccuracies in anomaly localization. Therefore, it’s imperative to differentiate between class-specific foreground and class-agnostic background in successive snippets. We employ scaled anomaly scores as activations to generate video-level foreground and background features. These are represented as follows:

$$\mathbf{A}^{fg} = \frac{\exp(\mu \mathbf{S}) - 1}{\sum_t (\exp(\mu \mathbf{S}_t) - 1)}, \quad \mathbf{V}^{fg} = \mathbf{A}^{fg} \mathbf{X}^e, \quad (14)$$

$$\mathbf{A}^{bg} = \frac{\exp(\mu \bar{\mathbf{S}}) - 1}{\sum_t (\exp(\mu \bar{\mathbf{S}}_t) - 1)}, \quad \mathbf{V}^{bg} = \mathbf{A}^{bg} \mathbf{X}^e, \quad (15)$$

where \mathbf{S}_t represents the anomaly score at the snippet-level, and μ is a predefined scaling factor that cooperates with the $\exp(\cdot)$ operation to enhance high-confidence activations. In contrast, $\bar{\mathbf{S}} = 1 - \mathbf{S}_t$ represents the normal confidence of the current snippet, and \mathbf{X}^e denotes the middle-layer output of the MLP, which shares the same dimension as \mathbf{T}^c .

3) *Cross-modal Alignment*: Finally, we present our approach to enhancing the fine-grained semantics of visual features by aligning them with the prompt features, as illustrated in Figure 2. For abnormal videos, foreground features \mathbf{V}_a^{fg} are matched with the corresponding abnormal prompt \mathbf{T}_a^c , increasing the chances of accurate abnormal behavior identification. In contrast, background features \mathbf{V}_a^{bg} are paired with the normal prompt \mathbf{T}_n^c , maintaining normal classification within an abnormal setting. For semantically inconsistent video-prompt pairs, repulsion is achieved by minimizing the cosine distance of all negative pairs, further forming discriminative decision boundaries. The process is formulated as:

$$\psi(\mathbf{V}, \mathbf{T}) = \frac{\mathbf{V} \cdot \mathbf{T}^\top}{\|\mathbf{V}\| \|\mathbf{T}\|}, \quad (16)$$

$$p_i^{v2t}(\mathbf{V}) = \frac{\exp(\psi(\mathbf{V}, \mathbf{T})/\tau)}{\sum_{k=1}^{C+1} \exp(\psi(\mathbf{V}, \mathbf{T}_k)/\tau)}, \quad (17)$$

where $\psi(\cdot)$ measures the cosine similarity between the visual representation $\mathbf{V} \in \{\mathbf{V}_a^{fg} \cup \mathbf{V}_a^{bg} \cup \mathbf{V}_n^{fg}\}$ and the textual representation $\mathbf{T} \in \{\mathbf{T}_a^c \cup \mathbf{T}_n^c\}$. Probability $p^{v2t}(\cdot)$ estimates how likely a visual feature matches a specific prompt across C anomaly classes and 1 normal class, where τ is a temperature coefficient. Finally, the cross-modal alignment loss is computed using the Kullback-Leibler divergence, which forces the network to learn to distinguish between the visual content of the video that represents abnormal behavior (foreground) and the content that is not relevant to the abnormal behavior (background). The loss function is formulated as follows:

$$\mathcal{L}_{kd} = \mathbb{E}_{p \sim p(v)} [\log p^{v2t}(v) - \log q^{v2t}(v)], \quad (18)$$

where $p^{v2t}(v)$ and $q^{v2t}(v)$ denote the similarity score and semantic consistency label of the video-prompt pair, respectively. If it is a positive pair, $q = 1$; otherwise, $q = 0$.

E. Training and testing procedures

In the training phase, the overall objective function of our model is denoted as:

$$\mathcal{L} = \mathcal{L}_{ce} + \lambda \mathcal{L}_{kd}, \quad (19)$$

where the coefficient λ is utilized to adjust the alignment loss. By optimizing this objective function, our model acquires the ability to generate discriminative representations of positive and negative snippets, while also effectively capturing the fine-grained semantics of anomalies. As a result, the generalizability of our model in complex situations is improved.

Considering the temporal consistency, we utilize a score smoothing (SS) strategy in the testing phase to mitigate the impact of transient noise that does not conform to the expected duration of an anomaly. This strategy employs distinct pooling operations to achieve optimal results. Given an anomaly score sequence $\{s_i\}_{i=1}^T$ and a pooling window of size κ , the smoothing process can be expressed as follows:

$$\tilde{s}_i = \frac{1}{\kappa} \sum_{i=\kappa_s}^{\kappa_e} s_i, \quad (20)$$

where κ_s and κ_e denote the start and end positions of pooling, respectively. For a moving window, we set $\kappa_s = \lfloor i/\kappa \rfloor \kappa$ and

$\kappa_e = (\lfloor i/\kappa \rfloor + 1) \kappa - 1$. For a sliding window, we define $\kappa_s = i$ and $\kappa_e = i + \kappa - 1$. In cases where the score sequence is shorter than the window size, the remainder is padded with zeroes. By smoothing the prediction scores, individual biases can be effectively suppressed while further reducing the occurrence of false alarms.

IV. EXPERIMENTS

In this section, we first introduce the datasets and implementation details of our framework. Subsequently, we make a comparison between our approach and state-of-the-art methods. Finally, we assess the contribution of each component via extensive ablation studies and qualitative analysis.

A. Datasets and Evaluation Metric

We perform experiments on three challenging anomaly benchmarks, i.e., UCF-Crime [11], XD-Violence [18], and ShanghaiTech [53] datasets. Details are given as follows:

1) *UCF-Crime*: It comprises 1900 surveillance videos with a total duration of 128 hours, covering 13 real-world anomalies, including abuse, robbery, explosion, and road accidents. For the weakly supervised setting, the dataset is divided into 1610 training and 290 test videos, with only video-level annotations available for training and frame-level annotations provided for testing.

2) *XD-Violence*: It is the latest and largest multimodal violence dataset containing 4754 untrimmed videos with a total duration of 217 hours. The dataset includes 3954 training and 800 test videos from various sources such as surveillance, movies, car cameras, and games. The dataset covers six types of violence, including abuse, car accidents, explosions, fighting, riots, and shooting. Most of the videos in this dataset contain artistic expressions such as camera movements and scene switching, which poses a challenge for video anomaly detection.

3) *ShanghaiTech*: It comprises 437 videos from 13 campus scenes. The original version was used for semi-supervised anomaly detection, where the training set consisted of only normal videos. Zhong *et al.* [20] reorganized the dataset for weakly supervised setting, with 238 videos in the training set and 199 videos in the test set.

4) *Evaluation Metric*: For the UCF-Crime and ShanghaiTech datasets, we use the area under the frame-level receiver operating characteristic curve (AUC) as the evaluation metric. For the XD-Violence dataset, the area under the precision-recall curve, also known as the average precision (AP), is utilized as the standard evaluation metric. In addition, to guarantee the reliability of anomaly prediction, a false positive rate with a threshold value of 0.5, or false alarm rate (FAR), is employed for evaluation. In general, the lower the FAR, the more robust the model is to normal samples.

B. Implementation Details

1) *Data Pre-processing*: Following existing work [12], we utilize the I3D network [51] pre-trained on Kinetics-400 for feature extraction. For a fair comparison, we utilize a 10-crop augmentation strategy, consisting of the center, four

corners, and their mirrored counterparts, for the UCF-Crime and ShanghaiTech datasets. In contrast, a 5-crop augmentation strategy, consisting of the center and four corners, is employed for the XD-Violence dataset as in [18].

2) *Hyperparameter settings*: The hidden dimension of the TCA module is set to 128, and the fusion weight α is initialized to 0.5. The two Conv1D layers of MLP have 512 and 300 nodes, respectively, both with a dropout rate of 0.1. For the UCF-Crime, XD-Violence, and ShanghaiTech datasets, we adopt local window sizes of 9, 9, and 5, respectively. The causal convolution kernel size Δt is set to 9, 3 and 3, respectively, while the temperature coefficient τ is initialized to 0.09, 0.05, and 0.2, respectively. The scaling factor μ for context separation is set to 10, and the loss coefficient λ is set to 1 for the UCF-Crime and XD-Violence datasets, while it is set to 9 for the ShanghaiTech dataset.

3) *Training and Test Details*: During the training phase, our model is trained using the Adam optimizer [54] with a batch size of 128 and 50 epochs in total. The initial learning rate for the three datasets is set to 5×10^{-4} , with a cosine decay strategy. For the balance of computational efficiency and detection performance, the snippet sampling threshold in the training phase is set to 200 as in [18] [12] [47] [55]. In the testing phase, we apply sliding pooling with window sizes of 7 and 3 to the UCF-Crime and ShanghaiTech datasets, respectively. Moving pooling with a window size of 9 is adopted for the XD-Violence dataset. More details can be found in our supplementary materials.

C. Comparison With State-of-the-Art Methods

We report the state-of-the-art results on the three benchmarks in Tables I-III. Notably, the weakly supervised-based methods exhibit superior performance compared to the semi-supervised learning methods across all datasets. The former, which solely relies on normal videos during the training phase, is susceptible to higher false alarms. This is primarily due to the weak generalization of these methods to unseen samples. On the other hand, the latter enables the model to learn discriminative patterns of both normal and abnormal samples, thus leading to better detection performance.

Our model achieves competitive results across three datasets among all weakly supervised methods. Although our model’s AUC value on the UCF-Crime dataset is lower than that of UR-DMU [21] by 0.21%, we outperform UR-DMU on the XD-Violence dataset with an absolute gain of 3.93% AP and lower false alarm rates on both datasets. UR-DMU employs an additional encoder layer to learn local features with temporal masks, whereas our TCA module achieves local-global context modeling by reusing the similarity matrix. Additionally, UR-DMU introduces a dual memory unit to store normal and abnormal patterns and enlarges the decision boundary through a magnitude distance loss. In contrast, our proposed PEL module enhances fine-grained detection performance by considering the intra-class discriminability of exceptions while ensuring inter-class separability. Notably, our model even outperforms multimodal-based approaches, namely CMA-LA [46] and MACIL-SD [44], with 2.19%

TABLE I
PERFORMANCE COMPARISON OF STATE-OF-THE-ART METHODS ON THE UCF-CRIME DATASET

Supervision	Method	Feature	AUC (%)	FAR (%)
Semi-supervised	Conv-AE [5]	-	50.60	27.2
	Lu <i>et al.</i> [30]	-	65.51	3.1
	GODS [35]	BoW+TCN	70.46	2.1
Weakly-supervised	MIL-Rank [11]	C3D RGB	75.41	1.9
	IBL [56]	C3D RGB	78.66	-
	Motion-Aware [57]	PWC Flow	79.00	-
	GCN [20]	TSN RGB	82.12	0.1
	MIST [42]	I3D RGB	82.30	<u>0.13</u>
	HL-Net [18]	I3D RGB	82.44	-
	MS-BSAD [47]	I3D RGB	83.53	-
	RTFM [13]	I3D RGB	84.30	-
	CRFD [12]	I3D RGB	84.89	0.72
	DDL [55]	I3D RGB	85.12	-
	MSL [43]	I3D RGB	85.30	-
	MLAD [58]	I3D RGB	85.47	7.47
	NL-MIL [19]	I3D RGB	85.63	-
	S3R [14]	I3D RGB	85.99	-
	Cho <i>et al.</i> [40]	I3D RGB	86.10	-
	CUPL [15]	I3D RGB	86.22	-
	UML [41]	X-CLIP RGB	86.75	-
UR-DMU [21]	I3D RGB	86.97	1.05	
	Ours	I3D RGB	<u>86.76</u>	0.43

TABLE II
PERFORMANCE COMPARISON OF STATE-OF-THE-ART METHODS ON THE XD-VIOLENCE DATASET

Supervision	Method	Feature	AP (%)	FAR (%)
Semi-supervised	SVM baseline	-	50.78	-
	OCSVM [59]	-	27.25	-
	Conv-AE [5]	-	30.77	-
Weakly-supervised	MIL-Rank [11]	C3D RGB	73.20	-
	HL-Net [18]	I3D RGB	75.44	-
	CA-VAD [60]	I3D RGB	76.90	-
	RTFM [13]	I3D RGB	77.81	-
	CRFD [12]	I3D RGB	75.90	-
	DDL [55]	I3D RGB	80.72	-
	MSL [43]	I3D RGB	78.28	-
	NL-MIL [19]	I3D RGB	78.51	-
	S3R [14]	I3D RGB	80.26	-
	UR-DMU [21]	I3D RGB	<u>81.66</u>	<u>0.65</u>
	Cho <i>et al.</i> [40]	I3D RGB	81.30	-
	ACF [61]	I3D+VGGish	80.13	-
	MSAF [45]	I3D+VGGish	80.51	-
	CUPL [15]	I3D+VGGish	81.43	-
	CMA-LA [46]	I3D+VGGish	83.54	-
MACIL-SD [44]	I3D+VGGish	83.40	-	
	Ours	I3D RGB	85.59	0.57

and 2.05% gains on the XD-violence dataset, respectively. Although the PEL module requires extra textual information in the training phase, our model only requires appearance features as input during inference. Furthermore, our model achieves a new state-of-the-art performance on the ShanghaiTech dataset, with a 0.54% improvement in AUC compared to Cho *et al.* [40]. For the first time, the false alarm rate is also reduced to 0, indicating that the model has sufficient confidence and robustness for normal samples.

TABLE III
PERFORMANCE COMPARISON OF STATE-OF-THE-ART METHODS ON THE SHANGHAI TECH DATASET

Supervision	Method	Feature	AUC (%)	FAR (%)
Semi-supervised	Mem-AE [6]	-	71.20	-
	HF ² -VAD [62]	-	76.20	-
	AMP-Net [39]	-	78.80	-
Weakly-supervised	MIL-Rank [11]	C3D RGB	86.30	0.15
	IBL [56]	C3D RGB	82.50	0.10
	GCN [20]	TSN RGB	84.44	-
	CLAWS [63]	C3D RGB	89.67	-
	AR-Net [64]	RGB+Flow	91.24	0.10
	MIST [42]	I3D RGB	94.83	0.05
	CRFD [12]	I3D RGB	97.48	-
	RTFM [13]	I3D RGB	97.21	-
	MSL [43]	VideoSwin	97.32	-
	NL-MIL [19]	I3D RGB	97.43	-
	S3R [14]	I3D RGB	97.48	-
	UML [41]	X-CLIP RGB	96.78	-
	Cho <i>et al.</i> [40]	I3D RGB	97.60	-
	Ours	I3D RGB	98.14	0.00

TABLE IV
THE CONTRIBUTION OF EACH COMPONENT. TCA: TEMPORAL CONTEXT AGGREGATION. PEL: PROMPT-ENHANCED LEARNING. SS: SCORE SMOOTHING.

Baseline	TCA	PEL	SS	UCF AUC (%)	XD AP (%)	SHTech AUC (%)
✓	✗	✗	✗	82.50	73.51	92.25
✓	✓	✗	✗	85.72	83.28	97.82
✓	✓	✓	✗	86.36	85.26	98.00
✓	✓	✓	✓	86.76	85.59	98.14

D. Ablation Studies

In this subsection, we perform extensive ablation studies to verify the contribution of each component of our model.

1) *Contribution of Each Component:* Table IV details the impact of each model component. The ‘Baseline’, comprising MLP and the Classifier, is notably enhanced by the TCA module, which yields substantial improvements of 9.77% and 5.57% on XD-Violence and ShanghaiTech datasets, respectively. This underscores the importance of context modeling. The PEL module further advances detection across all datasets, particularly notable with a 1.98% AP increase on XD-Violence. This reflects the varied content sources of XD-Violence compared to UCF-Crime and ShanghaiTech. The SS strategy also demonstrates consistent, albeit smaller, improvements across the benchmarks, confirming its effectiveness and versatility.

2) *Contribution of Components of the TCA module:* Table V shows that local context modeling significantly improves baseline detection. In contrast, global context modeling decreases performance on UCF-Crime and ShanghaiTech but is more effective on XD-Violence. We argue that global context may introduce noise in fixed surveillance footage, reducing feature discrimination, but is beneficial for diverse film and TV content. Combining local and global contexts outperforms the baseline by 9.54% and 5.37% on XD-Violence and ShanghaiTech, respectively. This indicates the complementary roles of these approaches. The improvement from DPE highlights

TABLE V
CONTRIBUTION OF INTERNAL COMPONENTS OF THE TCA MODULE. DPE: DYNAMIC POSITION ENCODING.

Local	Global	DPE	UCF AUC (%)	XD AP (%)	SHTech AUC (%)
✓	✗	✗	84.64	77.29	92.79
✗	✓	✗	82.30	80.80	85.65
✓	✓	✗	85.18	83.05	97.62
✓	✓	✓	85.72	83.28	97.82

TABLE VI
THE SUPERIORITY OF ADAPTIVE FUSION IN THE TCA MODULE.

Weight	Global	Local	UCF AUC (%)	XD AP (%)	SHTech AUC (%)
Fixed	0.1	0.9	84.82	77.57	97.74
	0.3	0.7	85.07	81.08	97.57
	0.5	0.5	85.17	82.45	97.71
	0.7	0.3	85.24	82.53	97.56
	0.9	0.1	84.85	80.57	97.60
Learnable	α	$1 - \alpha$	85.72	83.28	97.82

TABLE VII
COMPARISON OF DIFFERENT NORMALIZATION METHODS FOR ADAPTIVE FUSION

Normalization	UCF AUC (%)	XD AP (%)	SHTech AUC (%)
None	83.44	83.28	94.06
Power Norm	82.88	79.04	94.52
L2 Norm	85.09	73.92	97.48
Power + L2 Norm	85.72	71.13	97.82

the importance of relative distance in feature location.

Furthermore, Table VI demonstrates adaptive fusion’s superiority over fixed weights, showing better results across all datasets. This supports the dynamic adjustment of local-global context fusion to suit diverse videos. Table VII compares normalization methods. Power and L2 normalization combined yield optimal results on UCF-Crime and ShanghaiTech. However, on XD-Violence, either normalization method significantly reduces performance, likely affecting the dataset’s diverse video source distribution. Consequently, we omit normalization for XD-Violence in adaptive fusion to preserve content diversity.

Finally, we visualize the correlation of context features to demonstrate the encoding performance of the TCA module. As shown in Figure 4, the heatmap demonstrates the cosine similarity between the snippets. From this, we observe: 1) UCF-Crime videos show less variability compared to XD-Violence due to their fixed surveillance nature versus XD-Violence’s complex scene and camera movements; 2) Local features offer more discriminability than global features, which tend to introduce more long-range noise; and 3) The aggregation of local and global features suppresses redundancies and enhances snippet discriminability.

3) *Analysis of Time-Space Complexity:* This subsection evaluates the Temporal Context Aggregation (TCA) module’s

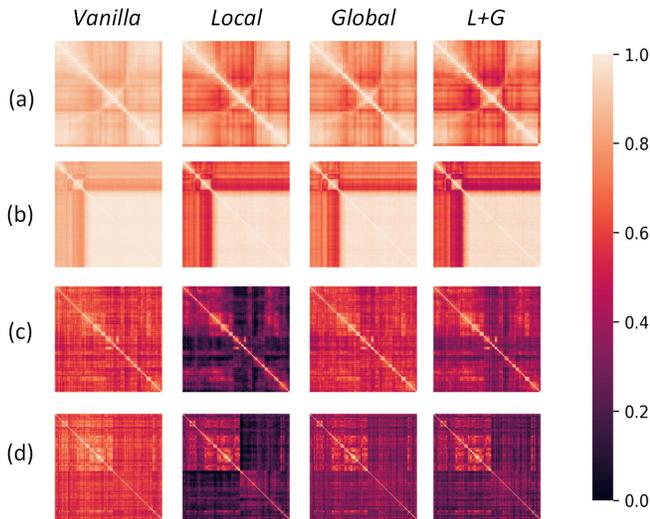


Fig. 4. Visualization of cosine similarity of context features. (a) and (b) are videos from the UCF-Crime dataset. (c) and (d) are videos from the XD-Violence dataset.

TABLE VIII
COMPARISON OF THE SPATIO-TEMPORAL COMPLEXITY OF DIFFERENT METHODS

Method	# Params	# FLOPS	UCF-Crime AUC(%)	XD-Violence AP(%)
HL-Net [†] [18]	0.66M	133M	83.12	79.58
MTCM [†] [47]	13.27M	2.653G	83.15	78.92
RTFM [†] [13]	24.72M	791M	84.30	77.81
LA-Net [55]	2.69M	538M	83.67	79.18
LGTE [†] [65]	1.21M	241M	84.34	79.07
TCA	1.21M	241M	85.72	83.28

[†] refers to re-implementation in our baseline.

performance and efficiency against recent temporal modeling techniques, shown in Table VIII. For a fair comparison, we used open-source codes to reproduce existing modules and integrated them into our baseline. HL-Net [18] uses dual graph convolutional networks for global and local relationship modeling. MTCM [47] employs multiple 1D convolution layers for multi-scale context, and RTFM [13] uses self-attention with pyramid dilated convolution for understanding temporal relations. These methods significantly increase parameters and computational demands. However, our TCA module, with similar space-time complexity, outperforms LGTE [65] by 1.38% and 4.21% on UCF-Crime and XD-Violence datasets, respectively. The TCA module’s success is due to its effective use of the similarity matrix, maintaining contextual continuity. This contrasts with LGTE’s channel grouping, which might lose some feature channel details. TCA’s adaptive aggregation and dynamic positional coding also improve snippet feature uniqueness. Notably, the feature extraction model I3D can achieve a processing speed of 263 FPS on a single Tesla A40 GPU when the input frame resolution is set to 224×224 . Subsequently, the anomaly detection model outputs anomaly scores at 350 FPS (2.9 ms per frame), striking a good balance between detection performance and efficiency.

TABLE IX
CONTRIBUTION OF NODE FILTERING IN PROMPT CONSTRUCTION

Method	UCF-Crime AUC (%)	XD-Violence AP (%)
None	85.34	84.35
step 1	85.87	84.61
step 1 + step 2 (fixed threshold)	86.31	85.10
step 1 + step 2 (dynamic threshold)	86.36	85.26

TABLE X
PERFORMANCE COMPARISON OF DIFFERENT PROMPT TEMPLATE

Prompt Template	UCF-Crime AUC (%)	XD-Violence AP (%)
{label}	85.95	84.69
{‘a video of’ + label}	85.55	84.51
{‘a long video of’ + label}	85.63	84.59
{label + WordNet Definition} [49]	85.92	84.14
{Learnable Prompt + label} [26]	85.48	84.48
{label + ConceptNet Relation}	86.36	85.26

TABLE XI
CONTRIBUTION OF CONTEXT SEPARATION IN THE PEL MODULE

Method	UCF-Crime AUC (%)	XD-Violence AP (%)
w/o context separation	85.71	79.79
w/ context separation	86.36	85.26

TABLE XII
CONTRIBUTION OF SCORE SMOOTHING STRATEGY

Method	UCF AUC (%)	XD AP (%)	SHTech AUC (%)
None	86.36	85.26	98
moving window	86.65	85.59	98.03
sliding window	86.76	85.56	98.14

4) *Contribution of Components of the PEL:* Table IX first evaluates the impact of node filtering in prompt construction. Step 1 involves removing entries with relevance scores at or below zero, while Step 2 applies a specific threshold to the remaining entries. Step 1 enhances detection by 0.53% and 0.26% on UCF-Crime and XD-Violence datasets, respectively. Step 2 with a dynamic threshold further improves performance on both datasets, indicating effective reduction of redundant entries. The dynamic threshold maintains relevant semantic relationships, while a fixed threshold might include low-relevance, potentially anomalous entries.

a) *Prompt Templates Comparison:* Different prompt templates are compared in Table X. Category labels as prompts increase detection by 0.23% and 1.41% on UCF-Crime and XD-Violence, respectively, surpassing hand-crafted templates and WordNet definitions. ConceptNet-based prompts further improve performance by 0.41% and 0.57% on these datasets. Using class-specific concepts based on semantic relations aids in identifying co-occurring anomalous items. Threshold filtering ensures the selection of highly relevant semantic concepts, enhancing fine-grained visual information capture.

b) *Context Separation*: As shown in Table XI, ‘w/o context separation’ denotes the mean pooling operation of snippet features and alignment with the corresponding prompt features, leading to notable performance drops. Mean pooling dilutes the anomalous region’s impact and introduces background noise, skewing cross-modal alignment. Therefore, context separation effectively minimizes background noise, ensuring better semantic consistency between visual foreground and prompt representations.

5) *PEL’s Impact on Fine-Grained Anomaly Detection*: As shown in Figure 5(a), our PEL method demonstrates superior anomaly detection capabilities on the UCF-Crime dataset, outperforming RTFM [13] and UMIL [41] in most categories with significant leads in *Abuse* (76.9%) and *Assault* (96.2%). Its robust average AUC of 72.2% underscores its effectiveness and potential applicability in security and surveillance. While PEL excels in complex anomaly scenarios, it lags slightly in *Explosion*, suggesting a need for refinement in detecting abrupt anomalies. In Figure 5(b), the performance of the PEL method in detecting various sub-classes of violence is contrasted with a baseline (w/o PEL). PEL consistently improves detection across all categories, with significant AP gains in *Fighting* (83.8% vs. 81.0%) and *Abuse* (70.5% vs. 61.5%). These improvements suggest that PEL effectively captures the nuances of complex, dynamic events. The overall average performance increase to 70.3% from 68.1% also underscores PEL’s robustness across diverse scenarios. However, in *Car Accident*, the impact of PEL is minimal, indicating potential areas for refinement.

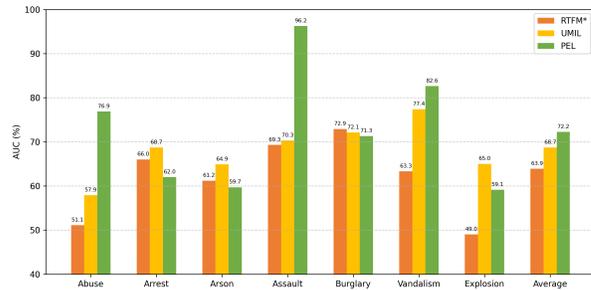
6) *Impact of Score Smoothing*: Table XII illustrates the effects of different pooling window sizes in the SS (Sliding and Sliding) strategy. Implementing both moving and sliding pooling techniques has shown to not only improve the model’s performance but also decrease the false alarm rate. Specifically, the sliding pooling technique yields better results for ShanghaiTech and UCF-Crime datasets. Conversely, the moving pooling technique is more effective for the XD-Violence dataset.

E. *Parameter Evaluation*

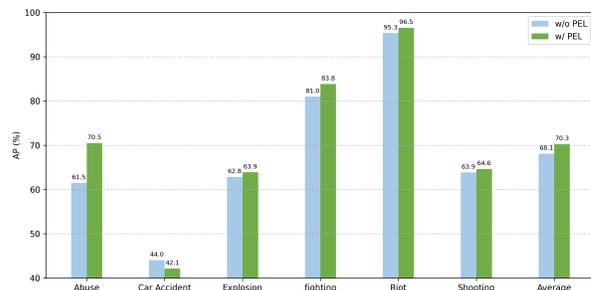
In this subsection, we explore the effect of different hyperparameter settings on model performance.

1) *Local Window Size in TCA module*: Figure 6(a) demonstrates the effect of different local window sizes in the TCA module. The performance of the model rises and then falls on both datasets as the window continues to increase, with the best results achieved on both datasets when the window w is set to 9. In general, smaller windows lack sufficient contextual information, while larger windows introduce long-range noise, leading to oscillations in false alarm rates.

2) *Kernel Size Δt in Classifier*: Figure 6(b) illustrates the impact of varying causal convolution kernel sizes on the UCF-Crime dataset. It is observed that the best results are achieved with a kernel size of 9. However, at this kernel size, the false alarm rate is not optimal, indicating a higher number of normal samples being misclassified. This suggests a need for further refinement in the model to balance detection accuracy and false alarm rates effectively.



(a) Class-wise AUC results of three methods on UCF-Crime.



(b) Class-wise AP results of before/after PEL on XD-Violence.

Fig. 5. Contribution of the PEL module to fine-grained anomaly detection.

TABLE XIII
PERFORMANCE COMPARISON OF DIFFERENT TEMPERATURE COEFFICIENTS τ

τ	0.01	0.03	0.05	0.07	0.09	0.1
UCF@AUC (%)	85.49	85.94	85.84	85.61	86.36	85.98
XD@AP (%)	84.11	83.56	85.26	84.16	84	84.63

TABLE XIV
PERFORMANCE COMPARISON OF DIFFERENT LOSS WEIGHT λ

λ	0.01	0.1	0.5	1	5	10
UCF@AUC (%)	85.11	85.19	85.52	86.36	84.51	84.42
XD@AP (%)	82.77	83.3	84.48	85.26	82.65	83.43

3) *Pooling Size κ for Score Smoothing*: We also investigate the effect of different pooling window sizes in the SS strategy, as presented in Figure 6 (c) and (d). The results show that our model is more sensitive to the size of the sliding window compared to the moving window. The optimal performance on the UCF-Crime dataset is achieved when the sliding window size is set to 7, resulting in a false alarm rate of 0.43%. For the XD-Violence dataset, a moving window size of 9 works best, with a further reduction in the FAR to 0.57%.

4) *Temperature Coefficient τ and Loss Weight λ* : Tables XIII and XIV report the model’s performance for different loss weights and temperature coefficients. Our model achieves optimal performance on both datasets when the weight is set to 1, which ensures a good balance of classification and alignment terms. However, the temperature coefficient varies for each dataset, with the best performance on UCF-Crime achieved at 0.09 and on XD-Violence at 0.05.

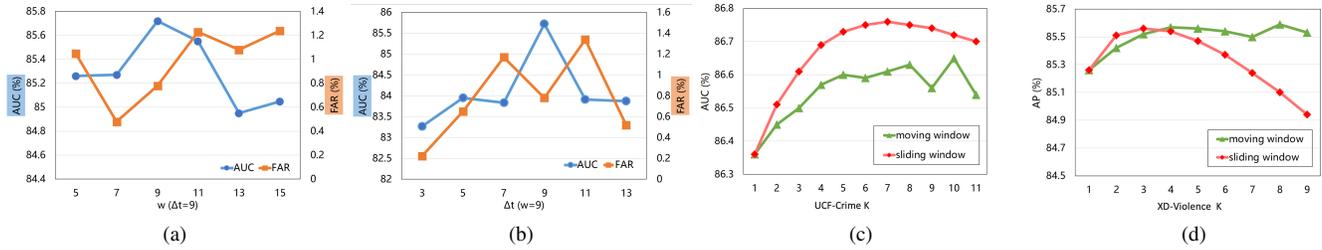


Fig. 6. Model performance under different hyperparameter settings. (a) and (b) show the AUC and FAR results of our model for different local windows w and causal convolution size Δt , respectively. (c) and (d) illustrate the effects of two different score smoothing methods on the model’s performance as the window size varies.

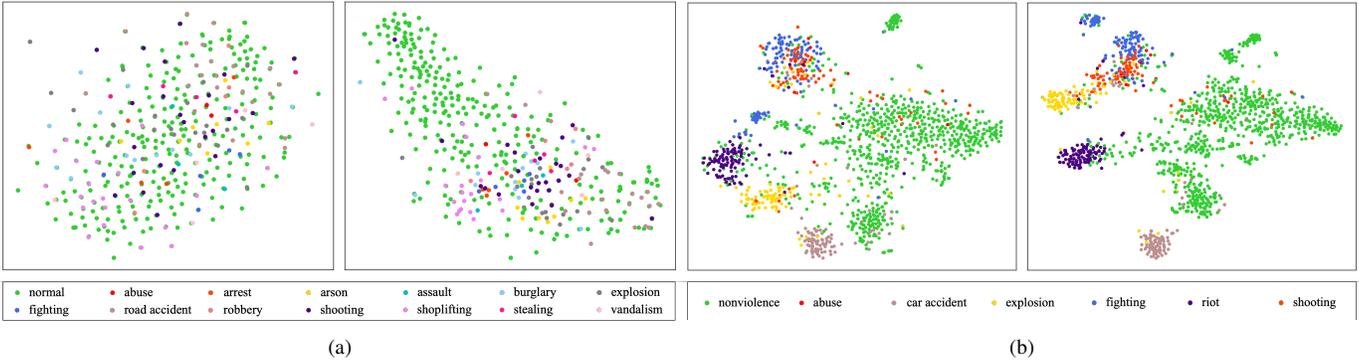


Fig. 7. Distribution of discriminative features before and after PEL using t-SNE. (a) the UCF-Crime dataset. (b) the XD-Violence dataset.

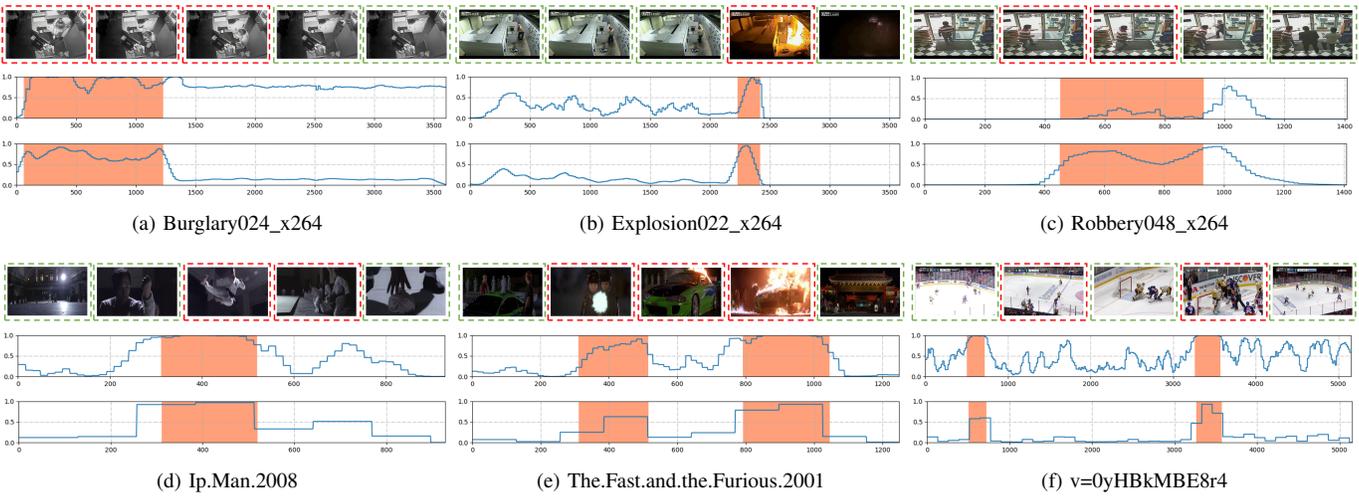


Fig. 8. Anomaly score visualization of the proposed method. Orange regions indicate ground-truths and blue curves indicate anomaly scores. (a)-(c) are videos from UCF-Crime and (d)-(f) are from XD-Violence, where the second and third rows indicate the results of ‘w/ TCA’ and ‘w/ PEL & SS’, respectively.

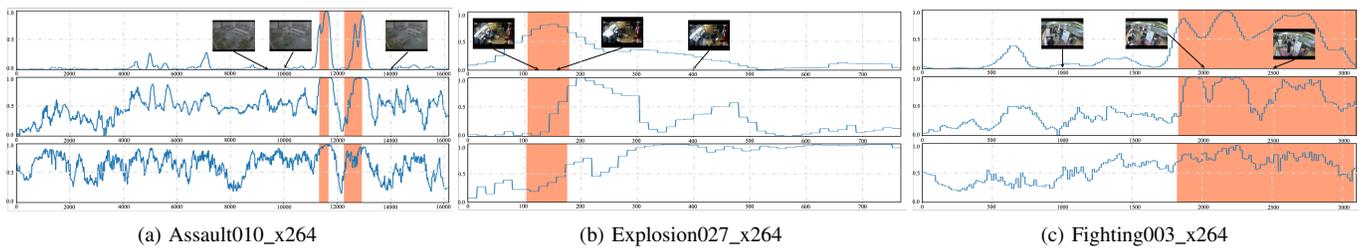


Fig. 9. Anomaly score visualization of different methods on the UCF-Crime dataset. The first row is from our proposed framework, and the second and third rows show the detection results of DDL [55] and CUPL [15], respectively.

F. Qualitative Analysis

We first use t-SNE visualization for the MLP’s middle layer features, as shown in Figure 7. Before introducing PEL, the UCF-Crime dataset samples appear disordered, while the XD-Violence dataset shows clearer clustering. This difference is attributed to the temporal coherence in surveillance videos and distinct shot boundaries in XD-Violence content. Cross-modal alignment loss makes abnormal foregrounds converge towards prompt features, forming tighter clusters. It also aligns hard-negative snippets and backgrounds with a normal center, increasing the distance between abnormal and normal snippets in the embedding space. This results in a more effective differentiation of fine-grained anomalous samples compared to traditional binary classification models.

Figure 8 illustrates the anomaly scores obtained using our methods. The TCA module’s detection results are shown in the second row, while the combined effect of PEL and SS is presented in the third row. PEL and SS, when used together, achieve more precise anomaly localization than TCA alone. This increased accuracy is attributed to the effective suppression of non-anomaly noise by PEL’s context separation and the mitigation of false alarms caused by frame changes and luminance variations. Moreover, the integration of rich semantic information from prompts significantly enhances the model’s versatility in detecting various anomalous events across different scenarios.

Finally, we compare the detection effectiveness of various methods in Figure 9. The PEL method exhibits superior sensitivity and specificity in identifying anomalies within the UCF-Crime dataset, particularly evident in scenarios like Assault, Explosion, and Fighting. PEL’s anomaly scoring is characterized by sharp peaks and clear distinctions of anomalous events, in stark contrast to the broader, less defined detection of DDL and the delayed or dispersed responses of CUPL. This demonstrates PEL’s exceptional temporal accuracy and event localization capabilities, establishing it as a highly effective tool for real-world surveillance applications.

V. CONCLUSION

In this paper, we focus on efficient temporal context modeling and semantic enhancement of visual features for weakly supervised video anomaly detection. We introduce a temporal context aggregation module that reuses the similarity matrix to capture local-global dependencies simultaneously. This method not only cuts down on parameters and computational load but also boosts detection capabilities, marking a significant shift from conventional paralleled systems and enhancing real-world applicability and efficiency. Recognizing the complex nature of anomalies, our model doesn’t just detect but understand them by integrating external knowledge for semantic interpretability. This aligns the model with human cognition, offering nuanced understanding beyond typical binary constraints. Our proposed prompt-enhanced learning module further refines this by using class-specific prompts derived from external knowledge, improving the distinction between abnormal sub-classes and maintaining clear inter-class separation. In the future, multimodal information such

as motion and audio remains to be explored, and open-set anomaly detection also deserves attention.

REFERENCES

- [1] A. Adam, E. Rivlin, I. Shimshoni, and D. Reinitz, “Robust real-time unusual event detection using multiple fixed-location monitors,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 30, no. 3, pp. 555–560, Mar. 2008.
- [2] B. Ramachandra, M. J. Jones, and R. R. Vatsavai, “A survey of single-scene video anomaly detection,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 44, no. 5, pp. 2293–2312, 2020.
- [3] Y. Liu, D. Yang, Y. Wang, J. Liu, and L. Song, “Generalized video anomaly event detection: Systematic taxonomy and comparison of deep models,” *arXiv preprint arXiv:2302.05087*, 2023.
- [4] G. Pang, C. Shen, L. Cao, and A. V. D. Hengel, “Deep learning for anomaly detection: A review,” *ACM Comput. Surv.*, vol. 54, no. 2, pp. 1–38, Mar. 2021.
- [5] M. Hasan, J. Choi, J. Neumann, A. K. Roy-Chowdhury, and L. S. Davis, “Learning temporal regularity in video sequences,” in *IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2016, pp. 733–742.
- [6] D. Gong, L. Liu, V. Le, B. Saha, M. R. Mansour, S. Venkatesh, and A. V. D. Hengel, “Memorizing normality to detect anomaly: Memory-augmented deep autoencoder for unsupervised anomaly detection,” in *Proc. IEEE/CVF Int. Conf. Comput. Vis. (ICCV)*, Oct. 2019, pp. 1705–1714.
- [7] H. Park, J. Noh, and B. Ham, “Learning memory-guided normality for anomaly detection,” in *IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2020, pp. 14 360–14 369.
- [8] W. Luo, W. Liu, D. Lian, J. Tang, L. Duan, X. Peng, and S. Gao, “Video anomaly detection with sparse coding inspired deep neural networks,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 43, no. 3, pp. 1070–1084, Mar. 2021.
- [9] C. Chen, Y. Xie, S. Lin, A. Yao, G. Jiang, W. Zhang, Y. Qu, R. Qiao, B. Ren, and L. Ma, “Comprehensive regularization in a bi-directional predictive network for video anomaly detection,” in *Proc. AAAI Conf. Artif. Intell.*, vol. 36, no. 1, Jun. 2022, pp. 230–238.
- [10] L. Wang, J. Tian, S. Zhou, H. Shi, and G. Hua, “Memory-augmented appearance-motion network for video anomaly detection,” *Pattern Recognit.*, vol. 138, p. 109335, Jun. 2023.
- [11] W. Sultani, C. Chen, and M. Shah, “Real-world anomaly detection in surveillance videos,” in *IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2018, pp. 6479–6488.
- [12] P. Wu and J. Liu, “Learning causal temporal relation and feature discrimination for anomaly detection,” *IEEE Trans. Image Process.*, vol. 30, pp. 3513–3527, 2021.
- [13] Y. Tian, G. Pang, Y. Chen, R. Singh, J. W. Verjans, and G. Carneiro, “Weakly-supervised video anomaly detection with robust temporal feature magnitude learning,” in *Proc. IEEE/CVF Int. Conf. Comput. Vis. (ICCV)*, Oct. 2021, pp. 4955–4966.
- [14] J.-C. Wu, H.-Y. Hsieh, D.-J. Chen, C.-S. Fuh, and T.-L. Liu, “Self-supervised sparse representation for video anomaly detection,” in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, 2022, pp. 729–745.
- [15] C. Zhang, G. Li, Y. Qi, S. Wang, L. Qing, Q. Huang, and M.-H. Yang, “Exploiting completeness and uncertainty of pseudo labels for weakly supervised video anomaly detection,” in *IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2023, pp. 16 271–16 280.
- [16] M. Defferrard, X. Bresson, and P. Vandergheynst, “Convolutional neural networks on graphs with fast localized spectral filtering,” in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 29, 2016, pp. 3837–3845.
- [17] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 30, 2017, pp. 5998–6008.
- [18] P. Wu, J. Liu, Y. Shi, Y. Sun, F. Shao, Z. Wu, and Z. Yang, “Not only look, but also listen: Learning multimodal violence detection under weak supervision,” in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, 2020, pp. 322–339.
- [19] S. Park, H. Kim, M. Kim, D. Kim, and K. Sohn, “Normality guided multiple instance learning for weakly supervised video anomaly detection,” in *Proc. IEEE Winter Conf. Appl. Comput. Vis. (WACV)*, Jan. 2023, pp. 2664–2673.
- [20] J.-X. Zhong, N. Li, W. Kong, S. Liu, T. H. Li, and G. Li, “Graph convolutional label noise cleaner: Train a plug-and-play action classifier for anomaly detection,” in *IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2019, pp. 1237–1246.

- [21] H. Zhou, J. Yu, and W. Yang, "Dual memory units with uncertainty regulation for weakly supervised video anomaly detection," *arXiv preprint arXiv:2302.05160*, 2023.
- [22] P. Liu, W. Yuan, J. Fu, Z. Jiang, H. Hayashi, and G. Neubig, "Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing," *ACM Comput. Surv.*, vol. 55, no. 9, pp. 1–35, Jan. 2023.
- [23] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, "Learning transferable visual models from natural language supervision," in *Proc. IEEE Int. Conf. Mach. Learn. (ICML)*, 2021, pp. 8748–8763.
- [24] M. Wang, J. Xing, and Y. Liu, "Actionclip: A new paradigm for video action recognition," *arXiv preprint arXiv:2109.08472*, 2021.
- [25] K. Zhou, J. Yang, C. C. Loy, and Z. Liu, "Learning to prompt for vision-language models," *Int. J. Comput. Vis.*, vol. 130, no. 9, pp. 2337–2348, Jul. 2022.
- [26] C. Ju, T. Han, K. Zheng, Y. Zhang, and W. Xie, "Prompting visual-language models for efficient video understanding," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, 2022, pp. 105–124.
- [27] R. Speer, J. Chin, and C. Havasi, "ConceptNet 5.5: An open multilingual graph of general knowledge," in *Proc. AAAI Conf. Artif. Intell.*, vol. 31, no. 1, Feb. 2017, pp. 4444–4451.
- [28] R. Mehran, A. Oyama, and M. Shah, "Abnormal crowd behavior detection using social force model," in *IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2009, pp. 935–942.
- [29] B. Zhao, L. Fei-Fei, and E. P. Xing, "Online detection of unusual events in videos via dynamic sparse coding," in *IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2011, pp. 3313–3320.
- [30] C. Lu, J. Shi, and J. Jia, "Abnormal event detection at 150 FPS in MATLAB," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, Dec. 2013, pp. 2720–2727.
- [31] W. Li, V. Mahadevan, and N. Vasconcelos, "Anomaly detection and localization in crowded scenes," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 36, no. 1, pp. 18–32, Jan. 2014.
- [32] D. Xu, Y. Yan, E. Ricci, and N. Sebe, "Detecting anomalous events in videos by learning deep representations of appearance and motion," *Comput. Vis. Image Underst.*, vol. 156, pp. 117–127, Mar. 2017.
- [33] M. Sabokrou, M. Khalooei, M. Fathy, and E. Adeli, "Adversarially learned one-class classifier for novelty detection," in *IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2018, pp. 3379–3388.
- [34] P. Wu, J. Liu, and F. Shen, "A deep one-class neural network for anomalous event detection in complex scenes," *IEEE Trans. Neural Netw. Learn. Syst.*, pp. 1–14, 2019.
- [35] J. Wang and A. Cheria, "GODS: Generalized one-class discriminative subspaces for anomaly detection," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, Oct. 2019, pp. 8200–8210.
- [36] W. Liu, W. Luo, D. Lian, and S. Gao, "Future frame prediction for anomaly detection - a new baseline," in *IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2018, pp. 6536–6545.
- [37] M. Z. Zaheer, J.-H. Lee, M. Astrid, and S.-I. Lee, "Old is gold: Redefining the adversarially learned one-class classifier training paradigm," in *IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2020, pp. 14 171–14 181.
- [38] J. R. Medel and A. Savakis, "Anomaly detection in video using predictive convolutional long short-term memory networks," *arXiv preprint arXiv:1612.00390*, 2016.
- [39] Y. Liu, J. Liu, K. Yang, B. Ju, S. Liu, Y. Wang, D. Yang, P. Sun, and L. Song, "Amp-net: Appearance-motion prototype network assisted automatic video anomaly detection system," *IEEE Trans. Ind. Inform.*, 2023.
- [40] M. Cho, M. Kim, S. Hwang, C. Park, K. Lee, and S. Lee, "Look around for anomalies: Weakly-supervised anomaly detection via context-motion relational learning," in *IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2023, pp. 12 137–12 146.
- [41] H. Lv, Z. Yue, Q. Sun, B. Luo, Z. Cui, and H. Zhang, "Unbiased multiple instance learning for weakly supervised video anomaly detection," in *IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2023, pp. 8022–8031.
- [42] J.-C. Feng, F.-T. Hong, and W.-S. Zheng, "MIST: Multiple instance self-training framework for video anomaly detection," in *IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2021, pp. 14 009–14 018.
- [43] S. Li, F. Liu, and L. Jiao, "Self-training multi-sequence learning with transformer for weakly supervised video anomaly detection," in *Proc. AAAI Conf. Artif. Intell.*, vol. 36, no. 2, 2022, pp. 1395–1403.
- [44] J. Yu, J. Liu, Y. Cheng, R. Feng, and Y. Zhang, "Modality-aware contrastive instance learning with self-distillation for weakly-supervised audio-visual violence detection," in *Proc. 30th ACM Int. Conf. Multimedia*, Oct. 2022, pp. 6278–6287.
- [45] D. Wei, Y. Liu, X. Zhu, J. Liu, and X. Zeng, "Msaf: Multimodal supervise-attention enhanced fusion for video anomaly detection," *IEEE Signal Process. Lett.*, vol. 29, pp. 2178–2182, 2022.
- [46] Y. Pu and X. Wu, "Audio-guided attention network for weakly supervised violence detection," in *Proc. IEEE Int. Conf. Consum. Electron. Comput. Eng. (ICCECE)*, Jan. 2022, pp. 219–223.
- [47] Y. Zhen, Y. Guo, J. Wei, X. Bao, and D. Huang, "Multi-scale background suppression anomaly detection in surveillance videos," in *Proc. IEEE Int. Conf. Image Process. (ICIP)*, Sep. 2021, pp. 1114–1118.
- [48] T. Gao, A. Fisch, and D. Chen, "Making pre-trained language models better few-shot learners," *arXiv preprint arXiv:2012.15723*, 2020.
- [49] L. Yao, J. Han, Y. Wen, X. Liang, D. Xu, W. Zhang, Z. Li, C. Xu, and H. Xu, "Detclip: Dictionary-enriched visual-concept paralleled pre-training for open-world detection," *arXiv preprint arXiv:2209.09407*, 2022.
- [50] G. A. Miller, "WordNet," *Commun. ACM*, vol. 38, no. 11, pp. 39–41, Nov. 1995.
- [51] J. Carreira and A. Zisserman, "Quo vadis, action recognition? a new model and the kinetics dataset," in *IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jul. 2017, pp. 4724–4733.
- [52] Z. Yu, J. Yu, J. Fan, and D. Tao, "Multi-modal factorized bilinear pooling with co-attention learning for visual question answering," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, Oct. 2017, pp. 1839–1848.
- [53] W. Luo, W. Liu, and S. Gao, "A revisit of sparse coding based anomaly detection in stacked RNN framework," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, Oct. 2017, pp. 341–349.
- [54] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [55] Y. Pu and X. Wu, "Locality-aware attention network with discriminative dynamics learning for weakly supervised anomaly detection," in *Proc. IEEE Int. Conf. Multimedia Expo (ICME)*, Jul. 2022, pp. 1–6.
- [56] J. Zhang, L. Qing, and J. Miao, "Temporal convolutional network with complementary inner bag loss for weakly supervised anomaly detection," in *Proc. IEEE Int. Conf. Image Process. (ICIP)*, Sep. 2019, pp. 4030–4034.
- [57] Y. Zhu and S. Newsam, "Motion-aware feature for improved video anomaly detection," in *Proc. Brit. Mach. Vis. Conf. (BMVC)*, 2019, p. 270.
- [58] C. Zhang, G. Li, Q. Xu, X. Zhang, L. Su, and Q. Huang, "Weakly supervised anomaly detection in videos considering the openness of events," *IEEE Trans. Intell. Transp. Syst.*, vol. 23, no. 11, pp. 21 687–21 699, Nov. 2022.
- [59] B. Schölkopf, R. C. Williamson, A. Smola, J. Shawe-Taylor, and J. Platt, "Support vector method for novelty detection," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 12, 1999, pp. 582–588.
- [60] S. Chang, Y. Li, S. Shen, J. Feng, and Z. Zhou, "Contrastive attention for video anomaly detection," *IEEE Trans. Multimedia*, vol. 24, pp. 4067–4076, 2022.
- [61] D.-L. Wei, C.-G. Liu, Y. Liu, J. Liu, X.-G. Zhu, and X.-H. Zeng, "Look, listen and pay more attention: Fusing multi-modal information for video violence detection," in *Proc. IEEE Int. Conf. Acoust. Speech Signal Process. (ICASSP)*. IEEE, 2022, pp. 1980–1984.
- [62] Z. Liu, Y. Nie, C. Long, Q. Zhang, and G. Li, "A hybrid video anomaly detection framework via memory-augmented flow reconstruction and flow-guided frame prediction," in *Proc. IEEE/CVF Int. Conf. Comput. Vis. (ICCV)*, Oct. 2021, pp. 13 568–13 577.
- [63] M. Z. Zaheer, A. Mahmood, M. Astrid, and S.-I. Lee, "Claws: Clustering assisted weakly supervised learning with normalcy suppression for anomalous event detection," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, 2020, pp. 358–376.
- [64] B. Wan, Y. Fang, X. Xia, and J. Mei, "Weakly supervised video anomaly detection via center-guided discriminative learning," in *Proc. IEEE Int. Conf. Multimedia Expo (ICME)*, Jul. 2020, pp. 1–6.
- [65] Z. Qing, H. Su, W. Gan, D. Wang, W. Wu, X. Wang, Y. Qiao, J. Yan, C. Gao, and N. Sang, "Temporal context aggregation network for temporal action proposal refinement," in *IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2021, pp. 485–494.