

Highlights

A Lightweight Video Anomaly Detection Model with Weak Supervision and Adaptive Instance Selection

Yang Wang, Jiaogen Zhou, Jihong Guan

- A new Lightweight video anomaly detection model is proposed .
- Weakly labeled data problem is mitigated by an adaptive sampling strategy .
- A lightweight multi-level temporal correlation attention module is designed.
- A lightweight hourglass-shaped fully connected layer is designed.
- Extensive experiments have shown that the proposed method is both lightweight and effective.

A Lightweight Video Anomaly Detection Model with Weak Supervision and Adaptive Instance Selection

Yang Wang^a, Jiaogen Zhou^b and Jihong Guan^{a,*} (Corresponding author)

^aTongji University, No. 4800 Cao'an Highway, Shanghai, 201804, China

^bJiangsu Province Engineering Research Center for Intelligent Monitoring and Management of Small Water Bodies, Huaiyin Normal University, Huaian 223300, China

ARTICLE INFO

Keywords:

Video Anomaly Detection
Weak Supervision
Adaptive Instance Selection
Lightweight Model

ABSTRACT

Video anomaly detection is to determine whether there are any abnormal events, behaviors or objects in a given video, which enables effective and intelligent public safety management. As video anomaly labeling is both time-consuming and expensive, most existing works employ unsupervised or weakly supervised learning methods. This paper focuses on weakly supervised video anomaly detection, in which the training videos are labeled whether or not they contain any anomalies, but lack information about the specific frames and quantities of anomalies. However, the uncertainty of weakly labeled data and the large model size prevent existing methods from wide deployment in real scenarios, especially the resource-limit situations such as edge-computing. In this paper, we develop a lightweight video anomaly detection model. On the one hand, we propose an adaptive instance selection strategy, which is based on the model's current status to select confident instances, thereby mitigating the uncertainty of weakly labeled data and subsequently promoting the model's performance. On the other hand, we design a lightweight multi-level temporal correlation attention module and an hourglass-shaped fully connected layer to construct the model, which can reduce the model parameters to only 0.56% of the existing methods (e.g. RTFM). Our extensive experiments on two public datasets UCF-Crime and ShanghaiTech show that our model can achieve comparable or even superior AUC score compared to the state-of-the-art methods, with a significantly reduced number of model parameters.

1. Introduction

surveillance serves as a critical tool for identifying unexpected or abnormal events in many scenarios such as traffic monitoring and public safety management. Traditionally, video surveillance is heavily dependent on manual operations and of low intelligence [1–3]. For example, many cameras are installed in public venues such as stations and parks, to monitor unexpected or abnormal events, which generate huge amounts of videos. To check these videos manually is time-consuming and laborious. The rapid development of computer vision and deep learning technologies has spurred more and more research on video abnormal event detection or video anomaly detection (VAD), which enables the applications of automatic scene monitoring and intelligent early warning.

Generally, there are three types of video anomaly detection methods, supervised [4], unsupervised [5–9] and weakly-supervised [3, 10–15]. Supervised video anomaly detection typically requires frame-level or even pixel-level labels, which incurs expensive training cost. Hence, there is a little related research in this direction. Unsupervised video anomaly detection uses unlabeled data to train models, with the lowest training cost, but exhibits poor performance. Weakly-supervised video anomaly detection (WVAD) uses weakly-labeled data to train the model, where the training videos are labeled whether or not they contain any anomalies, but there is no information about which frames the anomalies are lo-

cated. Thus, WVAD inherits the advantages of supervised and unsupervised methods. On the one hand, WVAD has better performance than unsupervised methods as some supervision is exploited. On the other hand, WVAD is much cheaper and more efficient in acquiring the training data than supervised methods, because the latter requires to label each frame whether it contains anomalies. Therefore, WVAD becomes a hot topic of video anomaly detection.

WVAD is typically based on Multiple Instance Learning (MIL) [3, 10]. Under the MIL framework, a video is viewed as a bag that consists of various clips, each of which is considered as an instance. For the training videos, the annotations are on the video level. That is, we know which videos have anomalies, but we do not know which clips (or instances) and frames have anomalies. In this context, WVAD methods face two major challenges. The first challenge is the uncertainty of the weakly labeled data: we do not know both the number and the locations of anomalous clips in each anomalous video, which limits the full exploitation of the training anomalous data, thus resulting in unsatisfactory performance. The second challenge is the huge model size. The models of existing methods have too many parameters, making them difficult to be applied in resource-pressing scenarios, such as edge-computing applications. Existing methods have been mainly trying to tackle these two challenges.

Concretely, existing MIL based methods [3, 10–12, 16] adopt a strategy to maintain balance between the number of anomalous videos and that of normal videos in the training set, and compute the loss by selecting the instance (or clip) of the highest anomaly score. This strategy aims to minimize the uncertainty of weakly labeled data. Typically, in real sce-

*Corresponding author

✉ tongji_wangyang@tongji.edu.cn (Y. Wang); zhoujg@hytc.edu.cn (J. Zhou); jhguan@tongji.edu.cn (J. Guan)

ORCID(s): 0000-0001-7511-2910 (Y. Wang)

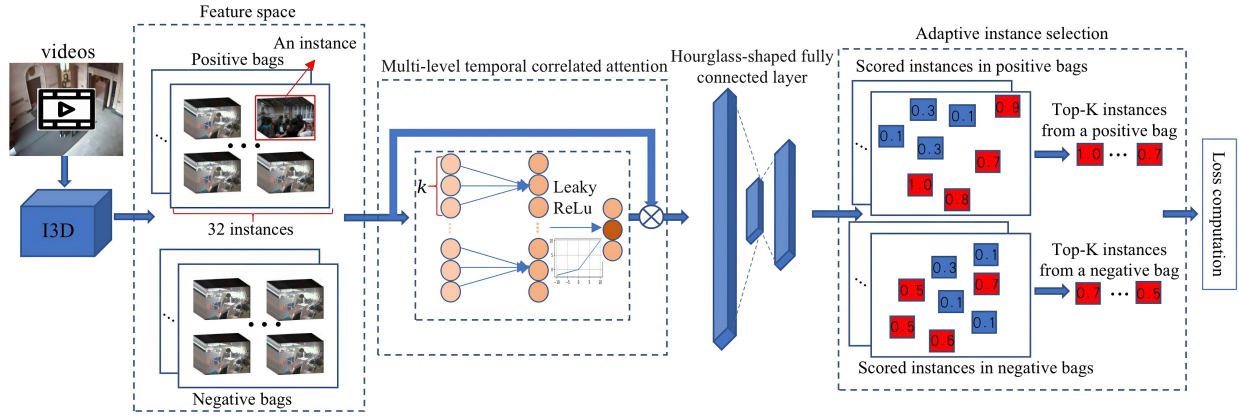


Figure 1: The framework of our method Light-WVAD. Our model is based on the multi-instance learning (MIL) framework. Each video is divided into 32 consecutive clips (or instances), which are grouped into a positive instance bag (for abnormal videos) or a negative instance bag (for normal videos). Video features are extracted by I3D. A *Multi-level Temporal correlation Attention* (MTA) module is designed to capture time-related information, which is then input to a *Hourglass-shaped Fully Connected layer* (HFC) to calculate the score of each instance. The top- K reliable instances are selected based on an *Adaptive Instance Selection* (AIS) strategy for subsequent loss calculation.

narios, the time duration of anomalous videos used for training usually exceeds 30 seconds, which suggests the presence of multiple anomalous clips. However, the selection of only the instance (clip) of the highest anomaly score might lead to under-utilization of all the available anomalous data. Statistical finding from RTFM [13] reveals that the mean number of anomalous clips within an anomalous video is approximately three. Therefore, choosing the three instances (clips) of the largest feature magnitude in both anomalous and normal videos for loss calculation can more effectively utilize the training data, thereby improving performance. However, our point of view is that different numbers of instances should be chosen for different datasets. Therefore, we propose an *Adaptive Instance Selection* (AIS) strategy, which selects the number of normal or anomalous instances adaptively based on the current status of the model in training, thereby it is able to accommodate different training datasets.

On the other hand, there are relatively fewer lightweight models for VAD in the literature [14–16], which usually compromise performance to embrace smaller models. In this paper, to devise a lightweight yet high-performance anomaly detection model, we develop a lightweight *Multi-level Temporal correlation Attention* (MTA) module. This module emphasizes the relationships between video instances (clips) of different time spans, thus making the model focus on the important instances in the videos. Furthermore, we design an *Hourglass-shaped Fully Connected layer* (HFC), which contains only half the parameters of a conventional fully connected layer (FC), yet does not degrade the model’s performance. Finally, Observing the limitation of the sparsity loss widely used in existing methods, we specifically introduce a more effective antagonistic loss.

In summary, in this paper we propose a novel, MIL-based lightweight video anomaly detection model, whose framework is shown in Fig. 1. First, we use I3D to decompose each video into 32 continuous clips, which are encoded

and organized into the positive or negative instance bag. The instances in the bags are input into the Multi-level Temporal correlation Attention (MTA) module, to make the model focus on the important features in the time dimension of the videos. Next, instance features are input into an Hourglass-shaped Fully Connected (HFC) layer to obtain the anomaly score for each instance. Finally, we dynamically determine the number K of positive and negative instances used for loss calculation based on the adaptive instance selection strategy. The main contributions of this study include:

1. We propose an adaptive sample selection strategy that alleviates the uncertainty problem of weakly-labeled data and improves model performance.
2. We design a lightweight multi-level temporal correlation attention module and an hourglass-shaped fully connected layer, which leads to a lightweight video anomaly detection method with only 0.56% of the parameters of the existing SOTA method RTFM [13].
3. We analyze the limitations of using sparsity loss in weakly-supervised video anomaly detection, and develop a more suitable antagonistic loss for this problem.
4. We conduct extensive experiments on two benchmark datasets UCF-Crime and ShanghaiTech, which show that our model achieves competitive or even superior performance compared to state-of-the-art methods of video anomaly detection.

The rest of this paper is organized as follows: Sec. 2 reviews the related works and highlights the differences of our work from the existing ones. Sec. 3 introduces the proposed method in detail. Sec. 4 presents the results of performance evaluation, including performance comparison with existing works and ablation studies. Finally, Sec. 5 concludes the paper and pinpoints the future works.

2. Related Work

Here we review the related work from two aspects: weakly supervised video anomaly detection and attention network.

2.1. Weakly Supervised Video Anomaly Detection

Traditional anomaly detection methods typically assume that only normal training data is available and use hand-crafted features for one-class classification to solve the problem [17–20]. With the development of deep learning techniques, some unsupervised learning methods utilize deep neural networks to learn features such as human posture and optical flow [5, 21–24], or utilize the difference in feature distribution between normal and abnormal samples [9] to perform abnormality detection. The essence of unsupervised video anomaly detection methods [5, 21–24] lies in the assumption that anomalies are rare events, and therefore the model is trained to learn mainly the feature distribution of normal samples. Based on this assumption, the model then determines whether a test sample is an anomaly by calculating the reconstruction errors or feature distributions of the test sample and the normal samples. However, due to the lack of prior knowledge about anomalies, these methods are prone to overfitting to training data and unable to distinguish normal and anomalous events.

Some works [3, 10–13, 25, 26] have shown that utilizing partially labeled abnormal samples can achieve better performance than unsupervised methods. However, the cost of obtaining a large number of frame-level labels is prohibitively high. Therefore, some video anomaly detection methods apply video-level labels for weakly supervised training. Sultani et al. [3] proposed a method that uses video-level labels and introduced the large-scale *weakly supervised video anomaly detection* (WVAD) dataset UCF-Crime. This makes WVAD one of the mainstream research directions in video anomaly detection [10, 13, 27, 28].

Existing Weakly supervised video anomaly detection methods are mainly based on multiple-instance learning (MIL). Given that the training data possesses only video-level labels, these approaches typically employ the instances of the highest anomaly prediction score from both positive and negative bags for loss computation during training, leading to under-utilization of the training data. To address this issue, Zhong et al. [29] transformed WVAD into a binary classification problem in the presence of label noise. They employed a *Graph Convolutional Neural Network* (GCN) [30] to eliminate label noise, thus enhancing data utilization and model performance. While this method improves model performance, the training computation cost associated with GCN and MIL is considerably high. Furthermore, it may make features unconstrained in the feature space, resulting in unstable performance. Furthermore, Tian et al. [13] amalgamated representation learning and anomaly score learning by devising *Robust Temporal Feature Magnitude* (RTFM) learning. They separately selected the three instances of the largest feature magnitude, in both abnormal and normal bags for loss calculation. Such an approach can effectively utilize the training data, thus achieving better performance. How-

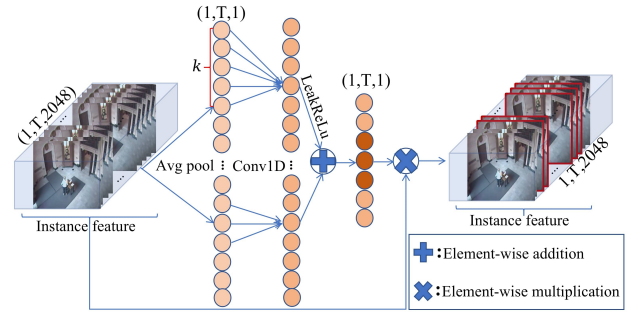


Figure 2: The structure of the multi-level temporal correlation attention (MTA) module. Here, each video is divided into T (32 in this paper) clips, each of which corresponds to an instance.

ever, this training data utilization strategy does not consider the diverse characteristics of different datasets.

In this paper, we propose an adaptive instance selection strategy, which adaptively selects the numbers of normal and abnormal instances for subsequent optimization, based on the current training status of the model, aiming to boost the utilization of weak label datasets.

2.2. Attention Network

Attention networks are initially used for machine translation [31] and later widely used for various computer vision tasks such as image classification [32], object detection [33], image segmentation [34], image captioning [35], and action recognition [36] etc., and have achieved excellent performance. Recently, attention networks have also been employed in weakly supervised tasks. Choe et al. [37] proposed an attention-oriented dropout layer, leveraging self-attention mechanisms to address the weakly supervised object localization problem. W-TALC [38] combines MIL and common activity similarity loss to train an attention module to solve the weakly supervised action localization problem in videos. Zhou et al. [6] proposed an attention that focuses on the foreground of an image to alleviate the foreground-background imbalance problem in anomaly detection. Li et al. [39] improved the ability of weakly supervised anomaly detection models to extract relationships between video frames by using SE-attention [40]. However, this method, not specifically designed for video data, fails to enable the model to concentrate on the temporal correlation between consecutive instances (clips), and has a relatively large number of parameters.

In this paper, We devise a Multi-level Temporal Correlation Attention module for video data with temporal relationships. This helps the model to focus attention on important instances (clips) in the video. Furthermore, its limited number of parameters makes it more suitable for video anomaly detection models in resource-constrained scenarios.

3. Methodology

This section presents the proposed *Lightweight Weakly Supervised Video Anomaly Detection* (Light-WVAD in short)

method. The framework of Light-WVAD is shown in Fig. 1. Light-WVAD mainly comprises a lightweight multi-level temporal correlation attention module, an hourglass-shaped fully connected layer, and an adaptive instance selection strategy to alleviate the uncertainty of weakly labeled data. In addition, we employ a more robust antagonistic loss to further optimize the model's performance. In what follows, we first give a formal definition of the problem, then introduce the major modules of our method in detail.

3.1. Problem Statement

In the context of weakly-supervised learning, anomaly detection can typically be regarded as a multiple-instance learning problem. Given a set of videos with video-level annotations (i.e., just labeling whether or not a video in the training dataset contains anomalous content), Weakly-supervised video anomaly detection (WVAD) aims to train a model with the annotated data, which is able to predict whether there are anomalies in any new videos.

3.2. Feature Extraction

In the data preprocessing stage, the widely used networks for feature extraction are I3D [41] and C3D [42]. Some existing research works [13, 27, 43] have shown that I3D can more effectively extract sample features. Therefore, in this paper we use I3D for data preprocessing and convert videos into feature vectors. Concretely, Given a video V_i , we divide it into 32 consecutive and non-overlapping clips, each of which is regarded as an instance. The clips are grouped into a positive or negative bag based on the video-level labels Y . Here, the positive bag ($Y = 1$) contains at least one anomalous instance (clip), while the negative bag ($Y = 0$) consists only of normal instances (clips).

3.3. Multi-level Temporal Correlation Attention

The structure of the Multi-level Temporal Correlation Attention (MTA) module is shown in Fig. 2. In this module, we first use a global average pooling layer to convert the T (T is 32 in this paper) instance features in a bag into a T -dimensional vector representing T channels. Then, we evaluate cross-channel interactions via convolution, and ultimately determine the weight of each channel. This makes the model focus on the important instances of each video in the time dimension.

To reduce parameters, we use a one-dimensional convolution $Conv1D$ of kernel size k ($k \geq 3$) to capture the cross-channel interactions of adjacent instances. In addition, in order to capture the interaction information between adjacent channels at different time spans, we jointly use convolution kernels of different sizes. The calculation process is as follows:

$$T_{attention} = (Conv1D[k], Conv1D[k-2], \dots, Conv1D[3]) * G(\chi) \quad (1)$$

In Equ. (1), $Conv1D[k]$ represents the convolution kernel function of size k ($k \leq T$), $*$ denotes the convolution operation, χ is the feature vector of a video and $G(\chi)$ is the fea-

ture vector after global pooling. In the one-dimensional convolution layer, the convolution operation between $Conv1D[k]$ and $G(\chi)$ is equivalent to sliding the convolution kernel, and multiplying it with the feature, then adding the convolution results. The calculation of $G(\chi)$ is as follows:

$$G(\chi) = \frac{1}{WH} \sum_{i=1, j=1}^{W, H} \chi_{ij} \quad (2)$$

Finally, by feeding $T_{attention}$ into the activation layer LeakRelu [44] and performing a concatenation operation with the original input $G(\chi)$, the attention module outputs $Y_{attention}$:

$$Y_{attention} = \lambda_1 Sum(L_Relu(T_{attention})) \cdot \chi \quad (3)$$

Here, $L_Relu()$ [44] is the activation function, λ_1 is a hyperparameter set to 0.1, and \cdot represents the corresponding multiplication operation. $Y_{attention}$ is the video feature with multi-level temporal correlations obtained by MTA, which can enhance the features of important instances in each video and weaken the features of less important instances.

3.4. Hourglass-shaped Fully Connected Layer

In the fully connected layer, each node connects with all the nodes of the preceding layer, resulting in a substantial number of parameters in the fully connected layer. To cut off the parameters in this layer, we propose a novel *hourglass-shaped fully connected* (HFC) layer. Fig. 3 illustrates the HFC structure on the right side, for comparison the traditional fully connected (FC) layer is illustrated on the left. HFC and FC have the same number of layers and layer dimensions, however, HFC has a *2048-64-128* structure, in contrast to the traditional FC's *2048-128-64*. Consequently, the number of parameters in HFC is approximately half of that in FC. With HFC, we get the anomaly scores of all instances in the positive and negative bags.

3.5. Adaptive Instance Selection

Here, we design an adaptive instance selection (AIS) strategy to choose important instances based on the current training status of the model and autonomously determines the number of instances used for loss computation. This strategy is primarily based on two important facts of MIL-based weakly supervised video anomaly detection. First, the negative bag contains only normal instances, i.e., each instance in the negative bag is normal. Second, as continuity exists amongst instances in the videos, the predicted anomaly scores of neighboring instances should also be continuous. These two factors enable us to infer the model's training status by examining its predictions of negative instances during the training process. With this, the confident positive (anomalous) instances can be determined from the set of current positive instances.

The workflow of AIS is shown in Fig. 4, which consists of three steps:

Step 1. With the anomaly scores of all instances from the HFC module, the first step of AIS is to calculate the confidence score ω that measures the mature degree of the model

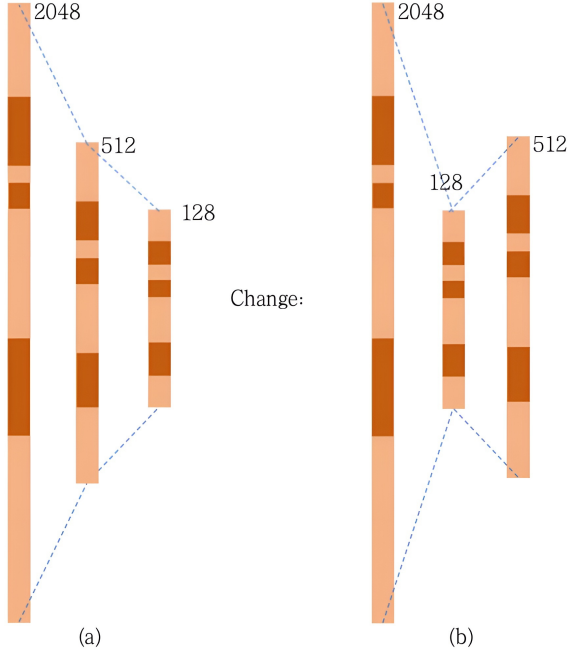


Figure 3: The structures of (a) the traditional fully connected layer (FC), and (b) our hourglass-shaped fully connected layer (HFC).

(i.e., how well the model is trained), based on the average anomaly scores of the negative instances and the average score difference between consecutive positive and negative instances. Concretely,

$$\omega = 1 - \frac{1}{T} \sum_{i=1}^T S_i^N - \frac{1}{2T-2} \sum_{i=1}^{T-1} (|S_{i+1}^N - S_i^N| + |S_{i+1}^P - S_i^P|) \quad (4)$$

Above, T is the number of instances in each positive/negative bag, S^P and S^N are the scores of positive and negative instances respectively. In Equ. (4), the 2nd item computes the mean of anomaly scores of negative instances, while the 3rd item calculates the mean of the differences in anomaly scores between consecutive positive and negative instances. As training goes, the 2nd term approaches 0 as negative instances should have 0 anomaly score, and the 3rd term also be close to 0 considering the continuity of anomaly scores of consecutive instances. Thus, the confidence score ω becomes close to 1.

Step 2. We select confident instances for loss computation based on the confidence score ω . The number of confident instances to be selected is determined as follows:

$$K = \omega * \sum_{i=1}^T f(S_i^P) \quad (5)$$

where the function $f(\cdot)$ is 1 if $S_i^P \geq 0.9$, otherwise is 0. That is, we first count the number of instances with score over 0.9, and then multiply it by the confidence score ω to obtain the final number of confident anomaly instances K .

Step 3. The feature magnitudes [13] of the top- K instances with the highest scores are selected separately from the positive and negative sets, and corresponding optimizations are performed accordingly. Thus, the loss function is as follows:

$$loss_{AIS} = \sum_{\chi \in U_{\chi_{top-K}}} (y \log(\text{mean}(s_{\theta}(\chi))) + (1-y) \log(1 - \text{mean}(s_{\theta}(\chi)))) \quad (6)$$

where s_{θ} is the feature extraction part of the model, mean is the mean function. U is the set of video instance (clip) features extracted by I3D, and $U_{\chi_{top-K}}$ is the set of features of the top- K instances selected from the positive and negative bags, χ is the feature of an instance. $y \in \{0, 1\}$, and $y = 0$ when χ is a normal video clip feature, otherwise $y = 1$.

3.6. Antagonistic Loss Function

Existing methods of weakly supervised anomaly detection predominantly employ smooth loss and sparsity loss for model optimization. Smooth loss assumes that there exists only slight discrepancy in features between consecutive instances, and thus the variation in anomaly scores is tiny, which conforms to the real scenario. On the other hand, the sparsity loss assumes that anomaly events in videos occur infrequently, therefore the mean anomaly score of instances in an anomalous video should approach zero. Although the assumption of sparsity is reasonable, it does not hold in actual model training. Owing to computational resource constraint, typically 32 consecutive instances are used to form a bag. In training, when anomalies are present in a bag, the proportion of abnormal instances cannot be overlooked. According to our preliminary analysis conducted on ShanghaiTech [29], the ratio of abnormal instances approximates 20%, and the mean score is around 0.2. Based on the above analysis, the use of smoothing loss is retained in this study. The implementation of smoothing loss is defined by the following formula:

$$loss_{smooth} = \frac{1}{T-1} \sum_{i=1}^{T-1} (\|S_{i+1} - S_i\|^2) \quad (7)$$

where S represents the anomaly scores predicted by the model for instances, and $\|\cdot\|^2$ is the L_2 norm. The L_2 norm assigns larger losses and gradients to anomalous instances.

Besides the smooth loss, we also devise an antagonistic loss, which is more suitable for weakly supervised MIL. The antagonistic loss is based on the antagonistic assumption of scores for positive and negative instances, capable of gauging the model's predictive performance on both positive and negative instances. As the negative bag exclusively contains normal instances, the scores of the most normal instances should approach zero. On the contrary, for the positive bag that includes some abnormal instances, the scores of the most anomalous instances should be close to 1. With this in mind, we have the following antagonistic loss:

$$loss_{antagonistic} = S_{top-1}^N + (1 - S_{top-1}^P) \quad (8)$$

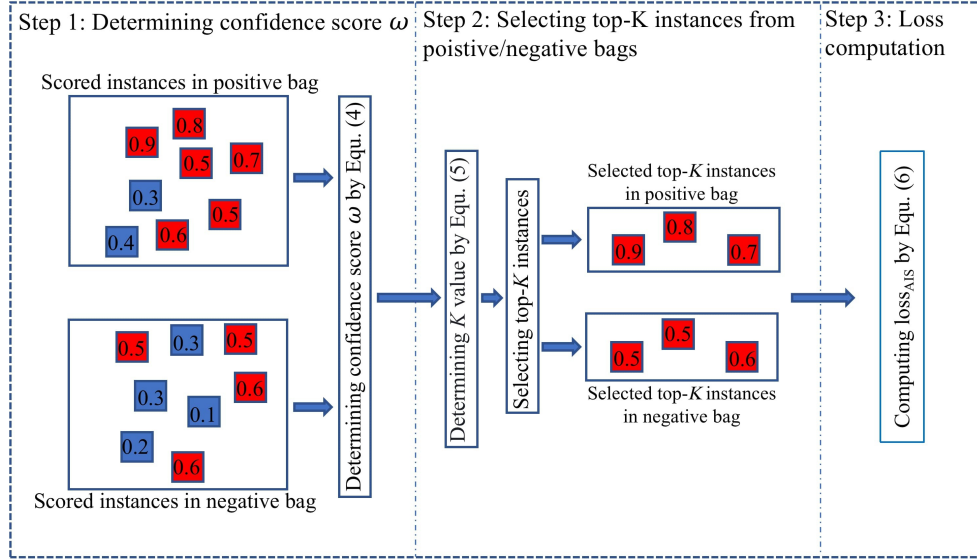


Figure 4: The workflow of our adaptive instance selection (AIS) strategy on a pair of positive and negative bags. It consists of three steps. Here, each red or blue square is an instance. Top- K instances are selected from both the positive and negative bags.

Here, S_{top-1}^P and S_{top-1}^N respectively represent the highest (top-1) abnormality score of instances in the positive and negative bags predicted by the model. $loss_{antagonistic}$ not only embodies the antagonistic constraint between positive and negative instances, but also provides optimization constraints for them respectively. Finally, the total loss function of our method Light-WVAD is as follows:

$$loss_{all} = loss_{AIS} + loss_{smooth} + loss_{antagonistic} \quad (9)$$

4. Performance Evaluation

4.1. Datasets and Evaluation Metrics

We evaluate our model on two commonly used video anomaly detection benchmark datasets, ShanghaiTech [45] and UCF-Crime [3]. The videos in both datasets were collected by fixed recording devices (surveillance cameras).

1. ShanghaiTech [45]: This dataset consists of 437 street surveillance videos captured from fixed angles, with 13 different background scenes. It contains 307 normal videos and 130 abnormal videos. Originally, it was used as a benchmark for unsupervised video anomaly detection. However, Zhong et al. [29] reorganized the dataset by selecting a subset of anomalous testing videos to create a weakly supervised training set, so that all 13 background scenes are covered in both the training and testing sets. We follow exactly the same procedure as in Zhong et al. [29] to convert ShanghaiTech for the weakly supervised setting.
2. UCF-Crime [3]: The UCF-Crime dataset is a real-world surveillance video dataset that consists of 950 anomalous videos belonging to 13 anomalous categories, and 950 normal videos. The training set provides only video-level labels, while the testing set pro-

vides both video-level labels and frame-wise annotations for evaluation. Each anomalous video in the testing set contains one or two anomalous events. Due to significant variations in the time duration of abnormal events across different videos, this dataset poses a serious challenge for video anomaly detection models.

For performance evaluation, we adopt the frame-based receiver operating characteristic (ROC) curve and the corresponding area under the curve (AUC) as performance metrics, which is consistent with that used in previous studies [3, 14, 15, 46, 47].

4.2. Compared Methods

We compare our method with a number of existing works, including thirteen general (non-lightweight) models and three lightweight methods. For some models, we consider two versions of using different feature extractors.

Among the thirteen non-lightweight models, the work of Sultani et al. [3] is the first to use the MIL framework and weakly supervision for video anomaly detection, and most of the subsequent works follow their settings. The method of Zhang et al. [10] defines a new inner-bag loss in the MIL framework to limit the function space and expand the differences in feature distributions for different types of instances. ARNet [12] uses a dynamic multi-instance learning loss to enlarge the interclass distance between anomalous and normal instances, and a center loss to narrow the intraclass distance of normal instances to boost the model's ability to distinguish anomalies. CLAWS Net+ [48] employs a clustering loss to mitigate labeling noise, which improves the learning of representations for both abnormal and normal videos. MIST [43] adopts a multi-instance self-training framework to effectively refine task-specific discriminative representations using only video-level annotations. Similarly, MSL [27]

proposes a self-training strategy and designs a transformer-based Multi-Sequence Learning network to further learn video-level anomaly probabilities and clip-level anomaly scores. RTFM [13] trains the feature magnitude learning function to effectively identify positive instances and improves the robustness of the model. BN-SVP [4] introduces a novel approach called Bayesian non-parametric submodular video partitioning to enhance the training of MIL models, and provides a reliable solution for robust anomaly detection in practical scenarios that involve outlier instances or multiple types of abnormal events. DAR [46] and the method of Wu et al. [11] design the model structure using a multi-branching and multi-stage approach, respectively, to improve the model's ability to understand and fuse different modal data in the videos. Mu et al.'s method [49] and GCN-Anomaly [29] innovatively employ graph convolutional networks for enhancing the model's ability to understand spatio-temporal correlation features in videos. The NTCN-ML [51] model extracts temporal representations of video data to construct a time-series pattern to optimize the multi-instance learning process. MGFN [52] propose a novel glance and focus network to effectively integrate spatial-temporal information for accurate anomaly detection. In addition, MGFN [52] propose the Feature Amplification Mechanism and a Magnitude Contrastive Loss to enhance the discriminativeness of feature magnitudes for detecting anomalies. NG-MIL [53] encodes diverse normal patterns from noise-free normal videos into prototypes to construct a similarity-based classifier. By combining predictions from classifiers, this approach can refine anomaly scores, reducing training instability from weak labels. CU-Net [54] introduces an enhanced framework with a two-stage self-supervised model. This model can generate and iteratively refine pseudo-labels by leveraging the completeness and uncertainty properties of weakly labeled data. HSN [55] propose a Human-Scene Network to learn discriminative representations by capturing both subtle and strong cues in a dissociative manner. Despite their remarkable performance in anomaly detection, these models do not consider factors such as model parameter size and runtime speed in their designs. In the three lightweight methods, the work of Chang et al. [14] proposes a lightweight MIL model incorporating a comparative attention module to improve model performance. The method of [15] introduces a self-attention mechanism to enhance the performance of the lightweight MIL model. BE-WVAD [16] proposes short-input inference modes, which can significantly reduce the required length of input videos and therefore greatly cut down memory and computational costs. Although these two methods achieve considerable advance, there remains substantial room to enhance the design of the model and boost the performance.

4.3. Implementation Details

Following the experimental setup in [3], we first divide each video into 32 ($T=32$) video clips. For the model parameters, we have the following settings:

1. The number of nodes in the *fully connected* (FC) layer is set to 2048, 64, and 128, respectively. The thresholds

of the Leaky ReLU activation function and dropout function are set to 0.5.

2. The input features are extracted from the "mix 5c" layer of the pre-trained I3D [41] network, and the multi-level temporal correlation attention module is used with a 1×1 Conv1D.
3. Our model is trained end-to-end using the Adam optimizer [50], with a weight decay coefficient of 0.0005, a batch size of 32, and a total of 200 epochs. The learning rate for both ShanghaiTech and UCF-Crime is set to 0.001. Each batch consists of 32 normal and abnormal instances respectively.

We implement our model using PyTorch. To ensure fairness in the comparison of model performance, we use the same benchmark settings as [3, 14, 15], and report the results of all baselines with the same backbone network as ours.

4.4. Performance Comparison

To demonstrate our method's effectiveness, we compare it with existing anomaly detection methods, including general WVAD models and lightweight WVAD methods.

Table 1 presents the frame-level AUC results on the ShanghaiTech dataset. We can see that our method obtains the highest performance of 95.9% among the lightweight weakly-supervised methods, and surpasses most of general (non-lightweight) methods, whilst maintaining a parameter size of merely 0.14M.

Table 2 gives the frame-level AUC performance on the UCF dataset. Our method not only achieves the highest AUC of 84.7% among the lightweight weakly-supervised methods, but also outperforms most non-lightweight weakly-supervised methods, standing at the third place among all general models.

Table 3 gives the frame-level AUC performance on the XD-Violence dataset. Our method not only achieves the highest AUC of 77.3% among the lightweight weakly-supervised methods, but also performs comparably to most of the non-lightweight weakly-supervised methods.

Besides comparing the AUC results of different methods on the two datasets, we also compare the number of parameters among all methods, the results are presented in Table 4. We can see that our method utilizes only 0.14M parameters, which is the lowest among the compared models. Particularly, our method has only less than 1% of the parameters of the RTMF method, and only half of the parameters of the currently smallest resource-intensive lightweight method [14].

Combined the comparison results from Table 1 to Table 4, it is obvious that our method is a lightweight yet good-performance method for weakly supervised video anomaly detection.

4.5. Ablation Studies

Effect of major modules in Light-WVAD. To evaluate the effectiveness of the major modules in our method, we conduct ablation experiment on the ShanghaiTech dataset. The results are presented in Table 5. Here, the baseline model

Table 1

Performance comparison of frame-level AUC on ShanghaiTech. For lightweight models, the best result is in red and the second best in blue; For general (Non-lightweight) models, the best result is in bold.

Model type	Method	Feature extractor	AUC (%)
General Models	GCN-Anomaly [29]	TSN	84.4
	AR-Net [12]	I3D	91.2
	CLAWS Net+ [48]	C3D	89.7
	MIST [43]	C3D	93.1
	MIST [43]	I3D	94.8
	RTFM [13]	C3D	91.5
	RTFM [13]	I3D	97.2
	MSL [27]	C3D	94.8
	MSL [27]	I3D	97.3
	BN-SVP [4]	C3D	96.0
	Mu et al. [49]	I3D	92.3
	NTCN-ML [51]	I3D+TCN	95.3
	DAR [46]	I3D	97.5
	NG-MIL [53]	I3D	97.4
	HSN [55]	I3D	96.2
Lightweight Models	Watanabe et al. [15]	I3D	95.7
	Chang et al. [14]	C3D	87.3
	Chang et al. [14]	I3D	92.3
	BE-WVAD [16]	I3D	95.0
	Light-WVAD (ours)	I3D	95.9

is a simple network consisting of fully connected layers, with an AUC of 93.8%. The results in Table 5 show that employing MTA, AIS, and the antagonistic loss (A-Loss in short) to the baseline individually can obviously improve the model's performance. Furthermore, our proposed lightweight HFC structure does not have negative impact on the model's performance. In summary, by combining MTA, HFC, AIS, and A-Loss into our method, a 2.1% performance improvement on the ShanghaiTech dataset is achieved.

Effect of parameter k in MTA. MTA is to capture multi-level temporal correlations of consecutive instances by integrating inter-instance feature relationships across multiple time intervals. As described in Section 3.3, MTA has a hyperparameter k , which represents the maximum number of consecutive k instances, from which MTA can extract temporal correlation information. The value of k affects the model's performance. We change the value of k from 3 to 15 with a stepsize of 2, and report the performance results in Table 6. We can see that when k is set to 5, the model achieves the best performance. By analyzing the performance change when setting different values of k for MTA, we can see that abnormal events typically have a relatively short duration. When k is too large, the proportion of normal instances is too high, which may dampen the influence of abnormal instances so that MTA fails to provide discriminative temporal correlation information. On the other hand, a too small k may result in insufficient coverage of anomalous instances, so that MTA cannot obtain effective temporal correlation information by enough inter-instance information.

Effect of the antagonistic loss function. To further verify the advantage of the proposed antagonistic loss, we compare the performance of three configurations: not using the

Table 2

Performance comparison of frame-level AUC on UCF. For lightweight models, the best result is in red and the second best in blue; For general (Non-lightweight) models, the best result is in bold.

Model type	Method	Feature extractor	AUC (%)
General Models	Sultani et al. [3]	TSN	75.4
	Zhang et al. [10]	TSN	78.7
	Wu et al. [11]	I3D	82.4
	GCN-Anomaly [29]	TSN	82.1
	CLAWS Net+ [48]	C3D	83.4
	MIST [43]	C3D	81.4
	MIST [43]	I3D	82.3
	RTFM [13]	C3D	83.3
	RTFM [13]	I3D	84.3
	MSL [27]	C3D	82.9
	MSL [27]	I3D	85.3
	BN-SVP [4]	C3D	83.4
	Mu et al. [49]	I3D	84.2
	NTCN-ML [51]	I3D+TCN	85.1
	DAR [46]	I3D	85.2
	CU-Net [54]	I3D	86.2
	MGFN [52]	I3D	87.0
	NG-MIL [53]	I3D	85.6
	HSN [55]	I3D	85.5
Lightweight Models	Watanabe et al. [15]	I3D	84.7
	Chang et al. [14]	C3D	83.4
	Chang et al. [14]	I3D	84.6
	BE-WVAD [16]	I3D	84.1
	Light-WVAD (ours)	I3D	84.7

Table 3

Performance comparison of frame-level AP on XD-Violence. For lightweight models, the best result is in red and the second best in blue; For general (Non-lightweight) models, the best result is in bold.

Model type	Method	Feature extractor	AP(%)
General Models	Sultani et al. [3]	C3D	73.2
	Wu et al. [11]	I3D	75.4
	RTFM [13]	C3D	75.9
	RTFM [13]	I3D	77.8
	MSL [27]	C3D	75.5
	MSL [27]	I3D	78.3
	DAR [46]	I3D	78.9
	MGFN [52]	I3D	79.2
	NG-MIL [53]	I3D	78.5
	CU-Net [54]	I3D	78.7
Lightweight Models	Chang et al. [14]	I3D+flow	71.5
	Chang et al. [14]	I3D	76.9
	BE-WVAD [16]	I3D	74.9
	Light-WVAD(ours)	I3D	77.3

sparsity loss, using the sparsity loss, and using our antagonistic loss (our method) on the ShanghaiTech dataset, the results are presented in Table 7. We can see that the antagonistic loss can obviously improve model performance, whereas the sparsity loss results in a decrease in model performance.

To delve deeper into why the sparsity loss causes model performance degradation, we train models employing the antagonistic loss and the sparsity loss respectively on the ShanghaiTech dataset, and illustrate their loss curves in Fig. 5,

Table 4

Model size comparison among weakly supervised methods. The best result is in red and the second best is in blue.

Method	#Parameters (M)
Sultani et al. [3]	2.11
Wu et al. [11]	0.76
RTFM [13]	24.72
Mu et al. [49]	13.20
BE-WVAD [16]	2.49
Watanabe et al. [15]	0.33
Chang et al. [14]	0.26
Light-WVAD (ours)	0.14

Table 5

Ablation study on ShanghaiTech.

Baseline	MTA	HFC	AIS	A-Loss	AUC (%) -SH
✓					93.8
✓	✓				94.8
✓		✓			93.9
✓	✓	✓			94.9
✓			✓		95.0
✓		✓	✓		95.1
✓	✓		✓		95.4
✓	✓	✓	✓		95.4
✓				✓	95.4
✓	✓	✓	✓	✓	95.9

Table 6

Ablation study on parameter k in MTA on ShanghaiTech.

Model	Hyperparameter k	AUC (%)
Baseline	-	93.8
Baseline+MTA	3	94.4
	5	94.8
	7	94.5
	9	93.9
	11	93.3
	13	93.3
	15	93.0

from which we can see a rapid drop of the antagonistic loss as the training goes, conforming to our expectation, whereas the sparsity loss curve exhibits an upward trend, which suggests a gradual increase in the average anomaly score of abnormal videos during the training. This indicates the unsuitability of the sparsity loss for model training, and further indicates that the proportion of abnormal instances in abnormal videos should not be ignored.

4.6. Visual Analysis

Visualization of test results. Here, we visualize the test results of our method and the baseline model on the ShanghaiTech and UCF-Crime datasets in Fig. 6 and Fig. 7 respectively, to further demonstrate the performance of our

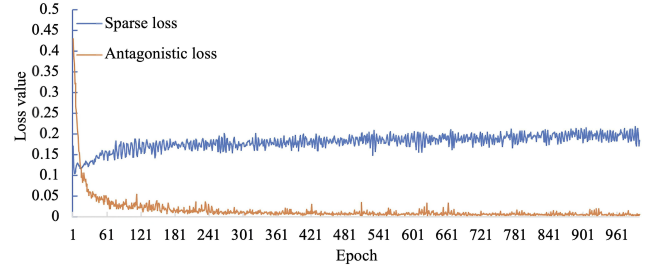


Figure 5: The loss curves in model training when using (a) the sparsity loss and (b) our antagonistic loss.

Table 7

Ablation study on loss functions in our method on ShanghaiTech.

Model	AUC (%) -SH
Ours ($loss_{AIS} + loss_{smooth}$)	95.4
Ours ($loss_{AIS} + loss_{smooth} + loss_{sparse}$)	94.8
Ours ($loss_{AIS} + loss_{smooth} + loss_{antagonistic}$)	95.9

method.

In Fig. 6(a), the anomalous event is a man falling to the ground, and in Fig. 6(b) the anomalous event is someone playing skateboard on pavement. However, in Fig. 6(c) and Fig. 6(d), there is no anomaly.

As illustrated in Fig. 6(a) and Fig. 6(b), our method can accurately identify the anomalous frames (e.g. No. 161, 205 and 241 in Fig. 6(a), and No. 127, 218 and 344 in Fig. 6(b)), i.e., assigning very high anomaly scores to these anomalous frames, while the baseline model identifies some normal frames as anomaly (e.g. Frame No. 337 in Fig. 6(a) and Frame No. 423 in Fig. 6(b)), and assigns low scores to some anomalies (e.g. Frame 127 in Fig. 6(b)). This indicates that our method has stronger capability of anomaly detection, and is more accurate in detecting the starting and ending of anomalies than the baseline model. Furthermore, Fig. 6(c) and Fig. 6(d) show that our method is more stable on normal video detection than the baseline model by consistently maintaining lower anomaly scores for normal frames.

In Fig. 7(a) and Fig. 7(b), the anomaly is an arrest action and a robbery event, respectively, while in Fig. 7(c) and Fig. 7(d), there is no anomaly.

The visualization results in Fig. 7(a) indicate that our method can accurately identify the anomalous event of police performing an arrest action. Nonetheless, we also observe incorrect responses at Frame No. 621 and No. 2272 in the normal screenshots, which are attributable to the drastic movements of individuals in the scene. Thus, we guess that our model faces challenge in delineating certain normal-anomaly boundaries, primarily due to the absence of detailed annotations. In Fig. 7(b), the anomaly corresponds to a gun-

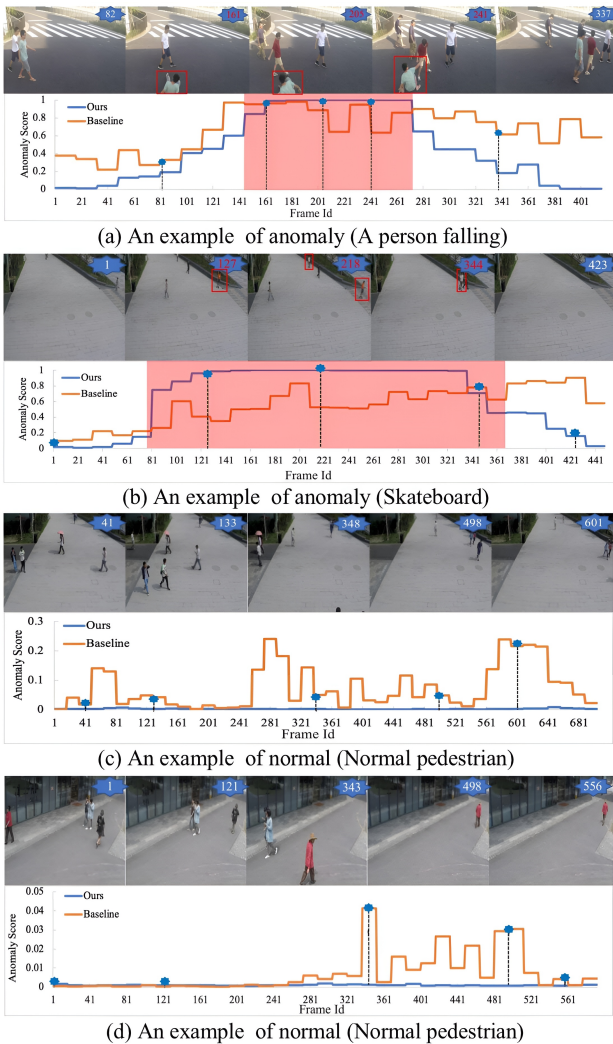


Figure 6: Visualization of test results of our method Light-WVAD and the baseline model on the ShanghaiTech dataset. The pink area denotes the time period that anomalies occur in the video, and the numerical labels on the line correspond to the labels in the video frames. In the anomalous video screenshots, the numbers are highlighted in red, and the anomalous objects are marked with red boxes in the frames.

man robbing a car, and a noticeable score fluctuation occurs at Frame No. 1351 in the video, due to the gunman being obscured by the car door, preventing the model from recognizing the anomaly. Nevertheless, outside the scope of the car door’s obstruction, our model can still successfully detect the anomaly. Ultimately, as per the visualization results in Fig. 7(c) and Fig. 7(d), our model still exhibits obvious stability on normal videos.

Detection performance of each anomalous class in the UCF-Crime dataset. Fig. 8 presents the AUC of our method on each anomalous class in the UCF-Crime dataset. Compared to RTFM and the baseline, our method yields superior or equivalent detection accuracy across 11 anomaly classes. Specifically, AUC is improved over 12% for the “Assault” and “Stealing” classes. Typically, without long-term analysis of object and human motion, detecting these two types of

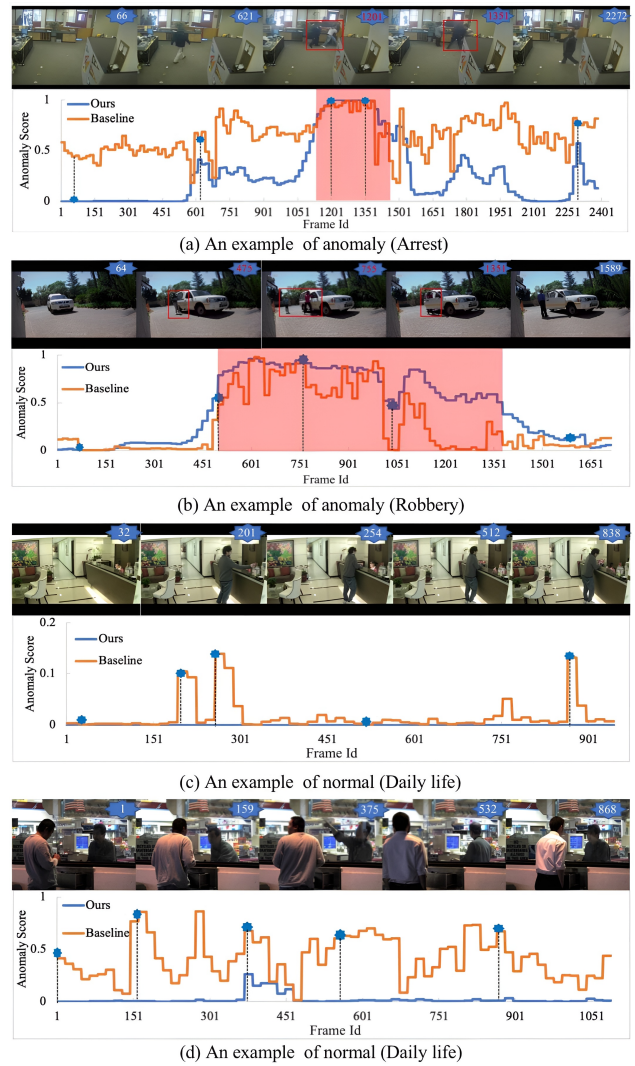


Figure 7: Visualization of test results of our method Light-WVAD and the baseline model on the UCF-Crime dataset. The pink area denotes the time period that anomalies occur in the video, and the numerical labels on the line correspond to the labels in the video frames. In the anomalous video screenshots, the numbers are highlighted in red, and the anomalous objects are marked with red boxes in the frames.

anomalies is very difficult. Nevertheless, our method achieves high detection accuracy, indicating that our MTA module is effective. Similar to the other methods, our model exhibits suboptimal performance on the “Explosion”, “Road Accidents”, “Vandalism”, and “Abuse” classes. For the sudden anomalies without enough warning signals, they remains a challenge for video anomaly detection.

5. Conclusion

This study develops a lightweight weakly-supervised video anomaly detection method (Light-WVAD) that can effectively addresses the uncertainty and high-parameter issues associated with the existing WVAD methods. Compared with existing lightweight models, Light-WVAD has the small-

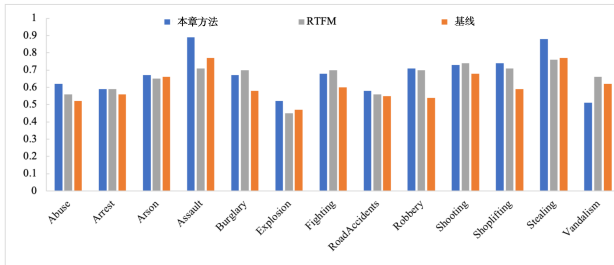


Figure 8: AUC results of three methods on different anomalous classes in the UCF-Crime dataset.

est number of parameters and the best performance. However, predicting sudden anomalies without warning signals is still a serious challenge for current video anomaly detection methods. And the predicted results on different anomalous classes in the UCF-Crime dataset also show that our model is not effective enough in detecting sudden anomaly classes. In the future, we plan to employ multi-modal techniques to enhance the model's understanding of video content, and thus develop more effective video anomaly detection models.

CRedit authorship contribution statement

Yang Wang: Writing—original draft, Conceptualization, Methodology, Formal analysis. **Jiaogen Zhou:** Review, editing, Conceptualization. **Jihong Guan:** Conceptualization, Resources, Review, editing.

Acknowledgement

This work was supported in part by National Natural Science Foundation of China (No. U1936205, No. 6217230 0), National Key R&D Program of China (No. 2021YFC330030 0), Open Research Projects of Zhejiang Lab (No. 2021KH0A B04), and the Fundamental Research Funds for the Central Universities (No. ZD-21-202101). Finally, we would like to thank Ang Li and Jun Yan for their corrections to the writing of this paper.

References

- [1] A.A. Sodemann, M.P. Ross, and B.J. Borghetti, "A review of anomaly detection in automated surveillance," *IEEE Trans. Syst. Man Cybern. Part C*, vol. 42, no. 6, pp. 1257–1272, 2012.
- [2] P. P. Oluwatoyin and K. Wang, "Video-based abnormal human behavior recognition - A review," *IEEE Trans. Syst. Man Cybern. Part C*, vol. 42, no. 6, pp. 865–878, 2012.
- [3] W. Sultani, C. Chen, and M. Shah, "Real-world anomaly detection in surveillance videos," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pp. 6479–6488.
- [4] H. Sapkota and Q. Yu, "Bayesian nonparametric submodular video partition for robust anomaly detection," in *Conference on Computer Vision and Pattern Recognition, CVPR, New Orleans, LA, USA, June 18-24, 2022*, pp. 3212–3221.
- [5] R. Colque, C. Caetano, M. de Andrade, W. Schwartz, "Histograms of Optical Flow Orientation and Magnitude and Entropy to Detect Anomalous Events in Videos," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 27, no. 3, pp. 673–682, 2017.
- [6] J. Zhou, L. Zhang, Z. Fang, J. Du, X. Peng, Y. Xiao, "Attention-Driven Loss for Anomaly Detection in Video Surveillance," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 30, no. 12, pp. 4639–4647, 2020.
- [7] G. Pang, C. Yan, C. Shen, A. van den Hengel, and X. Bai, "Self-trained deep ordinal regression for end-to-end video anomaly detection," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*, pp. 12 170–12 179.
- [8] Y. Zhang, X. Nie, R. He, M. Chen, Y. Yin, "Normality Learning in Multispace for Video Anomaly Detection," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 31, no. 9, pp. 3694–3706, 2021.
- [9] Y. Lu, C. Cao, Y. Zhang, Y. Zhang, "Learnable Locality-Sensitive Hashing for Video Anomaly Detection," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 33, no. 2, pp. 963–976, 2023.
- [10] J. Zhang, L. Qing, and J. Miao, "Temporal convolutional network with complementary inner bag loss for weakly supervised anomaly detection," in *IEEE International Conference on Image Processing, ICIP 2019, Taipei, Taiwan, September 22-25, 2019*, pp. 4030–4034.
- [11] P. Wu, J. Liu, Y. Shi, Y. Sun, F. Shao, Z. Wu, and Z. Yang, "Not only look, but also listen: Learning multimodal anomaly detection under weak supervision," in *Computer Vision - ECCV European Conference, Glasgow, UK, August 23-28, 2020*, pp. 322–339.
- [12] B. Wan, Y. Fang, X. Xia, and J. Mei, "Weakly supervised video anomaly detection via center-guided discriminative learning," in *IEEE International Conference on Multimedia and Expo, ICME 2020, London, UK, July 6-10, 2020*, pp. 1–6.
- [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 *IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*, pp. 4955–4966.
- [14] S. Chang, Y. Li, S. Shen, J. Feng, and S. Z. Zhou, "Contrastive attention for video anomaly detection," *IEEE Trans. Multim.*, vol. 24, pp. 4067–4076, 2022.
- [15] Y. Watanabe, M. Okabe, Y. Harada, and N. Kashima, "Real-world video anomaly detection by extracting salient features," in *IEEE International Conference on Image Processing (ICIP)*. IEEE, 2022, pp. 891–895.
- [16] Z. Yang, Y. Guo, J. Wang, D. Huang, X. Bao, Y. Wang, "Towards Video Anomaly Detection in the Real World: A Binarization Embedded Weakly-Supervised Network," *IEEE Trans. Circuits Syst. Video Technol.*, Early Access Article, 2023.
- [17] G. G. Medioni, I. Cohen, F. Brémond, S. Hongeng, and R. Nevatia, "Event detection and analysis from video streams," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 23, no. 8, pp. 873–889, 2001.
- [18] A. Basharat, A. Gritai, and M. Shah, "Learning object motion patterns for anomaly detection and improved object detection," in *IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2008), 24-26 June 2008, Anchorage, Alaska, USA*.
- [19] T. Zhang, H. Lu, and S. Z. Li, "Learning semantic scene models by object classification and trajectory clustering," in *IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2009), 20-25 June 2009, Miami, Florida, USA*, pp. 1940–1947.
- [20] J. Wang, Y. Song, T. Leung, C. Rosenberg, J. Wang, J. Philbin, B. Chen, and Y. Wu, "Learning fine-grained image similarity with deep ranking," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2014, Columbus, OH, USA, June 23-28, 2014*, pp. 1386–1393.
- [21] S. Sun, X. Gong, "Hierarchical Semantic Contrast for Scene-aware Video Anomaly Detection," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR, Vancouver, BC, Canada, June 17-24, 2023*, pp. 22846–22856.
- [22] S. Yu, Z. Zhao, H. Fang, A. Deng, H. Su, D. Wang, W. Gan, C. Lu, W. Wu, "Regularity Learning via Explicit Distribution Modeling for Skeletal Video Anomaly Detection," *IEEE Trans. Circuits Syst. Video Technol.*, Early Access Article, 2023.
- [23] X. Zeng, Y. Jiang, W. Ding, H. Li, Y. Hao, Z. Qiu, "A Hierarchical Spatio-Temporal Graph Convolutional Neural Network for Anomaly

- Detection in Videos,” *IEEE Trans. Circuits Syst. Video Technol.*, vol. 33, no. 1, pp. 200–212, 2023.
- [24] Z. Fang, J. Liang, J. T. Zhou, Y. Xiao, and F. Yang, “Anomaly detection with bidirectional consistency in videos,” *IEEE Trans. Neural Networks Learn. Syst.*, vol. 33, no. 3, pp. 1079–1092, 2022.
- [25] F. Hong, X. Huang, W. Li, and W. Zheng, “Mini-net: Multiple instance ranking network for video highlight detection,” in *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020*, pp. 345–360.
- [26] M. Z. Zaheer, A. Mahmood, H. Shin, and S. Lee, “A self-reasoning framework for anomaly detection using video-level labels,” *IEEE Signal Process. Lett.*, vol. 27, pp. 1705–1709, 2020.
- [27] S. Li, F. Liu, and L. Jiao, “Self-training multi-sequence learning with transformer for weakly supervised video anomaly detection,” in *AAAI Conference on Artificial Intelligence, AAAI 2022, February 22 - March 1, 2022*, pp. 1395–1403.
- [28] D. Zhang, C. Huang, C. Liu, and Y. Xu, “Weakly supervised video anomaly detection via transformer-enabled temporal relation learning,” *IEEE Signal Process. Lett.*, vol. 29, pp. 1197–1201, 2022.
- [29] J. 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 Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pp. 1237–1246.
- [30] R. Li, S. Wang, F. Zhu, and J. Huang, “Adaptive graph convolutional neural networks,” in *AAAI Conference on Artificial Intelligence, (AAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018*, pp. 3546–3553.
- [31] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” in *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pp. 5998–6008.
- [32] F. Wang, M. Jiang, C. Qian, S. Yang, C. Li, H. Zhang, X. Wang, and X. Tang, “Residual attention network for image classification,” in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pp. 6450–6458.
- [33] Y. Cao, K. Chen, C. C. Loy, and D. Lin, “Prime sample attention in object detection,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*, pp. 11 580–11 588.
- [34] L. Chen, Y. Yang, J. Wang, W. Xu, and A. L. Yuille, “Attention to scale: Scale-aware semantic image segmentation,” in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pp. 3640–3649.
- [35] K. Xu, J. Ba, R. Kiros, K. Cho, A. C. Courville, R. Salakhutdinov, R. S. Zemel, and Y. Bengio, “Show, attend and tell: Neural image caption generation with visual attention,” in *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, pp. 2048–2057.
- [36] H. Wu, X. Ma, and Y. Li, “Convolutional networks with channel and strips attention model for action recognition in videos,” *IEEE Trans. Multim.*, vol. 22, no. 9, pp. 2293–2306, 2020.
- [37] J. Choe and H. Shim, “Attention-based dropout layer for weakly supervised object localization,” in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pp. 2219–2228.
- [38] S. Paul, S. Roy, and A. K. Roy-Chowdhury, “W-TALC: weakly-supervised temporal activity localization and classification,” in *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018*, pp. 588–607.
- [39] Q. Li, R. Yang, F. Xiao, B. Bhanu, and F. Zhang, “Attention-based anomaly detection in multi-view surveillance videos,” *Knowl. Based Syst.*, vol. 252, p. 109348, 2022.
- [40] J. Hu, L. Shen, and G. Sun, “Squeeze-and-excitation networks,” in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pp. 7132–7141.
- [41] J. Carreira and A. Zisserman, “Quo vadis, action recognition? A new model and the kinetics dataset,” in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pp. 4724–4733.
- [42] D. Tran, L. D. Bourdev, R. Fergus, L. Torresani, and M. Paluri, “Learning spatiotemporal features with 3d convolutional networks,” in *IEEE International Conference on Computer Vision, ICCV 2015, Santiago, Chile, December 7-13, 2015*, pp. 4489–4497.
- [43] J.-C. Feng, F.-T. Hong, and W.-S. Zheng, “Mist: Multiple instance self-training framework for video anomaly detection,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 14009–14018.
- [44] X. Glorot, A. Bordes, and Y. Bengio, “Deep sparse rectifier neural networks,” in *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics, AISTATS 2011, Fort Lauderdale, USA, April 11-13, 2011*, pp. 315–323.
- [45] W. Liu, W. Luo, D. Lian, and S. Gao, “Future frame prediction for anomaly detection - A new baseline,” in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pp. 6536–6545.
- [46] T. Liu, C. Zhang, K. Lam, and J. Kong, “Decouple and resolve: Transformer-based models for online anomaly detection from weakly labeled videos,” *IEEE Trans. Inf. Forensics Secur.*, vol. 18, pp. 15–28, 2023.
- [47] Y. Wang, T. Liu, J. Zhou, and J. Guan, “Video anomaly detection based on spatio-temporal relationships among objects,” *Neurocomputing*, 2023.
- [48] M. Z. Zaheer, A. Mahmood, M. Astrid, and S.-I. Lee, “Clustering aided weakly supervised training to detect anomalous events in surveillance videos,” *arXiv preprint arXiv:2203.13704*, 2022.
- [49] H. Mu, R. Sun, M. Wang, and Z. Chen, “Spatio-temporal graph-based cnns for anomaly detection in weakly-labeled videos,” *Inf. Process. Manag.*, vol. 59, no. 4, p. 102983, 2022.
- [50] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *3rd International Conference on Learning Representations, ICLR, San Diego, CA, USA, May 7-9, 2015*.
- [51] W. Shao, R. Xiao, P. Rajapaksha, M. Wang, N. Crespi, “Video anomaly detection with NTCN-ML: A novel TCN for multi-instance learning,” *Pattern Recognit.*, vol. 143, p. 109765, 2023.
- [52] Y. Chen, Z. Liu, B. Zhang, W. Fok, X. Qi, and YC. Wu, “MGFN: Magnitude-Contrastive Glance-and-Focus Network for Weakly-Supervised Video Anomaly Detection,” in *Thirty-Seventh AAAI Conference on Artificial Intelligence, AAAI 2023, Washington, DC, USA, February 7-14, 2023*, pp. 387–395.
- [53] S. Park, H. Kim, M. Kim, D. Kim, and K. Sohn, “Normality Guided Multiple Instance Learning for Weakly Supervised Video Anomaly Detection,” in *IEEE/CVF Winter Conference on Applications of Computer Vision, WACV 2023, Waikoloa, HI, USA, January 2-7, 2023*, pp. 2664–2673.
- [54] C. Zhang, G. Li, Y. Qi, S. Wang, L. Qing, Q. Huang, and MH. Yang, “Exploiting Completeness and Uncertainty of Pseudo Labels for Weakly Supervised Video Anomaly Detection,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada, June 17-24, 2023*, pp. 16271–16280.
- [55] S. Majhi, R. Dai, Q. Kong, L. Garattoni, G. Francesca, and F. Brémond, “Human-Scene Network: A novel baseline with self-rectifying loss for weakly supervised video anomaly detection,” *Computer Vision and Image Understanding.*, vol. 241, pp. 103955–103965, 2024.



Yang Wang received the M.S. degree in Computer technology from the Nanchang University, Nanchang, China, in 2017. He is currently working toward the Ph.D. degree in Computer Science with the Tongji University, Shanghai, China. His research interests include anomaly detection, object detection and deep learning.



Jiaogen Zhou received the bachelor's degree from University of Chinese Academy of Sciences, China and the PhD degree from Wuhan University, China. His research interests include Image classification, object detection, animal pose estimation and machine learning.



Jihong Guan is now a professor of Department of Computer Science & Technology, Tongji University, Shanghai, China. She received his Bachelor's degree from Huazhong Normal University in 1991, her Master's degree from Wuhan Technical University of Surveying and Mapping (merged into Wuhan University since Aug. 2000) in 1991, and her Ph.D. from Wuhan University in 2002. Before joining Tongji University, she served in the Department of Computer, Wuhan Technical University of Surveying and Mapping from 1991 to 1997, as an assistant professor and an associate professor (since August 2000) respectively. She was an associate professor (Aug. 2000-Oct. 2003) and a professor (Since Nov. 2003) in the School of Computer, Wuhan University. Her research interests include spatial databases, artificial intelligence and bioinformatics. She has published more than 200 papers in domestic and international journals and conferences.