

# WEAKLY SUPERVISED VIDEO ANOMALY DETECTION BASED ON CROSS-BATCH CLUSTERING GUIDANCE

Congqi Cao, Xin Zhang, Shizhou Zhang, Peng Wang, and Yanning Zhang

ASGO National Engineering Laboratory, School of Computer Science  
Northwestern Polytechnical University

{congqi.cao, szzhang, peng.wang, ynzhang}@nwpu.edu.cn; zhangxin\_@mail.nwpu.edu.cn

## ABSTRACT

Weakly supervised video anomaly detection (WSVAD) is a challenging task since only video-level labels are available for training. In previous studies, the discriminative power of the learned features is not strong enough, and the data imbalance resulting from the mini-batch training strategy is ignored. To address these two issues, we propose a novel WSVAD method based on cross-batch clustering guidance. To enhance the discriminative power of features, we propose a batch clustering based loss to encourage a clustering branch to generate distinct normal and abnormal clusters based on a batch of data. Meanwhile, we design a cross-batch learning strategy by introducing clustering results from previous mini-batches to reduce the impact of data imbalance. In addition, we propose to generate more accurate segment-level anomaly scores based on batch clustering guidance further improving the performance of WSVAD. Extensive experiments on two public datasets demonstrate the effectiveness of our approach.

**Index Terms**— anomaly detection, weakly supervised learning, cross-epoch learning

## 1. INTRODUCTION

An efficient and accurate video anomaly detection algorithm can help maintain social security and stability. Therefore, video anomaly detection has high practical value and broad application prospects. With the development of weakly supervised learning, the weakly supervised video anomaly detection (WSVAD) method is an effective method for detecting anomalies, which uses weakly labeled training data containing both normal and anomalous videos to train the model. Recently, WSVAD has been formulated as a multiple instance learning (MIL) task. Sultani et al. [1] constructed a large-scale anomaly dataset and proposed a deep MIL ranking based method for WSVAD. Wan et al. [2] replaced the max anomaly score selection policy with a k-max value selection policy. Li et al. [3] selected the sequence with the highest sum of anomaly scores instead of selecting the instance with the highest anomaly score. Gong et al. [4] introduced temporal continuity of multiple neighboring instances at different

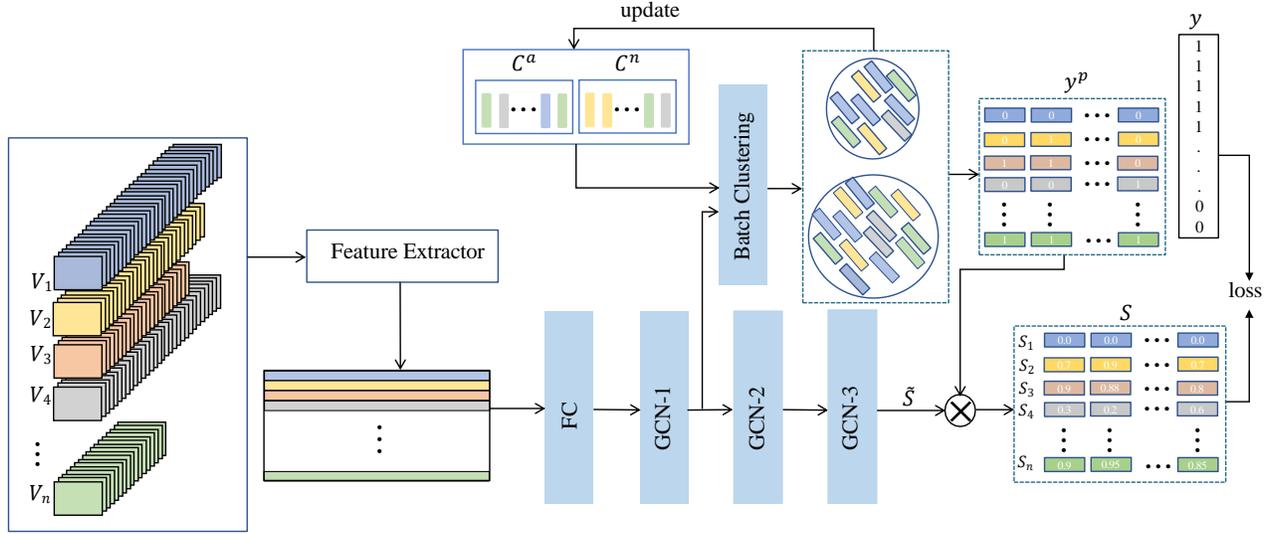
time scales.

However, most of the existing work [1][2][3] have only used MIL-based classification loss. Although MIL-based classification loss ensures the inter-class separability of the learned features to some extent, it is not sufficient for accurate anomaly detection. Therefore, we propose a batch clustering based loss to further increase the discriminative power of the features, and our framework is shown in Figure 1. The abnormal videos in a batch are clustered into two clusters, and then the loss encourages the network to maximize the distance between these two clusters. Meanwhile, for normal videos in a batch, the loss encourages the network to minimize the distance between these two clusters. Compared to clustering each video individually [5][6], batch clustering improves the robustness of the model and reduces the influence of the model affected by the noise.

Considering the data in WSVAD task is highly unbalanced, it may negatively affect model training when only using a small portion of data. Previous mini-batches can provide valuable knowledge enabling the model to better understand the underlying distribution of the data [7][8]. Therefore, we use a cross-batch learning strategy to provide guidance for the current batch clustering by introducing clustering results from previous batches. The introduced cross-batch learning strategy can make the clustering results more accurate, model the temporal-spatial distribution of the data better, and improve the adaptability of the model to unbalanced samples.

In addition, the knowledge that the model can learn in WSVAD directly from supervised learning is limited, since only video-level labels are available. The clustering branch can obtain potential information reflecting the similarity of video segments. Therefore, we propose anomaly score generation based on batch clustering guidance to generate more discriminative anomaly scores and further improve the model performance. We first generate pseudo labels for video segments based on the batch clustering results, and then uses the pseudo labels to guide the backbone network to rectify the estimated anomaly scores for video segments. In summary, our main contributions are as follows:

1) We propose a loss based on batch clustering to comple-



**Fig. 1.** Overview of our proposed method. The feature extractor extracts features from video segments. The extracted features are fed into a fully connected layer and three graph convolution layers to generate segment-level anomaly scores. Simultaneously, the batch clustering branch uses intermediate representations of a batch of videos learned from the GCN-1 layer to create clusters. And we design a cross-batch learning strategy to store the clustering results of previous batches and introduce them into the batch clustering of the current batch. Finally, the pseudo labels generated by the batch clustering branch guide the generation of anomaly scores.

ment and enhance the separability of the features guided by MIL-based classification loss.

2) We propose a cross-batch learning strategy to generate more accurate clustering results.

3) We propose to use the knowledge learned from batch clustering to guide the prediction of more discriminative anomaly scores.

## 2. PROPOSED METHOD

### 2.1. Backbone network

Our approach employs a backbone network based on graph convolutional neural network (GCN) to model video sequences. The backbone network consists of a feature extraction module and a graph convolution module [9]. The Inflated 3D (I3D) [10] pretrained on the Kinetics dataset is used as the feature extraction network to extract the appearance and motion information of the video segments. Before each video  $V_i$  is fed into the feature extraction module, the video is divided into non-overlapping segments containing 16 consecutive frames, and we denote the number of segments by  $T_i$ . The graph convolution module consists three graph convolution layers, where the first two layers are followed by a ReLU activation function and a dropout layer, and the last layer is followed by a Sigmoid activation function. For each video  $V_i$ , the input layer receives the temporal-spatial features extracted from the feature extraction module and the adjacency

matrix of a global graph constructed based on feature similarity and temporal proximity of the video segments. The output layer produces the anomaly score vector of the video  $\tilde{S}_i = \{\tilde{s}_{i,j}\}_{j=1}^{T_i}$ . The network is trained using video-level labels.  $y_i \in \{0, 1\}$  is the video-level label of video  $V_i$ , where  $y_i = 0$  indicates that video  $V_i$  is a normal video and  $y_i = 1$  indicates that  $V_i$  is an abnormal video.

### 2.2. Batch clustering based on K-means

Although MIL-based classification loss ensures the inter-class separability of learned features to a certain extent, it cannot ensure a more discriminative power of features since there is no explicit supervision. Several studies on unsupervised anomaly detection [11] have enlightened us, in which normal samples are compulsorily clustered in a compact space such that they can be kept away from the anomaly space. Therefore, we have reasonable grounds to believe that the normal activities should be compact in the feature space. Consequently, we use batch clustering to cluster normal video segments to enhance the intra-class compactness of normal features. A larger inter-class distance in abnormal video indicates that normal and abnormal are separated by a higher probability. Therefore, we perform batch clustering on abnormal video segments to enhance the inter-class dispersion of normal and abnormal features.

Here, we propose batch clustering to provide supervision to enhance the discriminative power of features as shown in

Figure 1. For all normal videos in a batch, the feature representations of each video are clustered into two clusters. Since all segments in normal videos are normal, we try to close the distance of the two clusters to ensure the intra-class compactness of the normal features. All the abnormal videos in a batch are also clustered into two clusters. Since there are both abnormal and normal video segments in the abnormal videos, we try to push the centers of the two clusters away from each other to achieve the inter-class dispersion of normal and abnormal features. Specifically, for the abnormal or normal videos in a batch, we use the K-means algorithm to cluster the normalized intermediate feature representations, which is the output of the first layer of GCN. The loss based on batch clustering is shown below:

$$L_{bc} = \begin{cases} \min(d, \mu), & \text{if } \{V_i\}_{i=1}^b \text{ are normal videos} \\ \frac{1}{d}, & \text{if } \{V_i\}_{i=1}^b \text{ are abnormal videos} \end{cases} \quad (1)$$

where  $d = c_1 - c_2$  is the distance between two cluster centers, and  $c_1, c_2$  are the two cluster centers.  $\mu$  is an upper bound that helps the model to be robust to different videos, and  $b$  is the batch size.

We train the model using the  $k$ -max loss function to expand the inter-class distance between abnormal and normal segments, denoted as follows:

$$L_{k\text{-max}} = -\frac{1}{k_i} \sum_{s_{i,j} \in p_i} [y_i \log(s_{i,j}) + (1 - y_i) \log(1 - s_{i,j})] \quad (2)$$

where  $p_i$  is the first  $k_i$  large elements in  $S_i$  of video  $V_i$ ,  $k_i = \lfloor \frac{T_i}{8} + 1 \rfloor$ ,  $y_i$  is the video-level label of video  $V_i$ . So, the total loss function is expressed as:

$$L = L_{k\text{-max}} + \lambda_1 L_{bc} \quad (3)$$

where  $\lambda_1$  is a trade-off hyperparameter to keep the balance between two losses.

### 2.3. Cross-batch learning strategy

The data in WSVAD task is highly imbalanced. However, only a mini-batch of samples can be accessed in each iteration. The performance of anomaly detection models can be improved when the batch size becomes larger on large-scale datasets. However, simply scaling up the batch is not an ideal solution because batch size is limited by GPU memory and computational cost. A simple solution to collect rich information is to introduce information from previous batches at each training iteration to enable the model to better understand the underlying distribution of the data. Thus, video segments from past batches can also serve as an important reference when performing  $K$ -means clustering on video segments of the current batch. Previously, we perform batch clustering based on  $K$ -means by directly selecting any two video segments from all video segments as the initial clustering centers, but the final clustering results of the  $K$ -means depend

on the selection of the initial clustering centers heavily. In order to select the most appropriate initial clustering centers, we introduce the clustering results of previous batches to provide guidance for the current batch clustering, which can help model training.

Specifically, at the  $t$ -th epoch,  $C^a$  and  $C^n$  are constructed in the iteration process to store the learned knowledge. That is, during the  $i$ -th iteration, we cluster all abnormal video segments in a batch, and add the two clustering centers into  $C^a$ . We will obtain  $C^a = \{c_{i,1}^a, c_{i,2}^a\}_{i=1}^m$  at the end of the  $t$ -th epoch, where  $m$  is the number of iterations in an epoch. We also do the above operation for all normal video segments in a batch to get  $C^n = \{c_{i,1}^n, c_{i,2}^n\}_{i=1}^m$ .

At the  $(t + 1)$ -th epoch, we use the stored information to guide the batch clustering. We cluster the data in  $C^a$  to obtain the clustering centers  $c_1^a, c_2^a$ , which are used as the initial clustering centers for the batch clustering of all abnormal video segments in each batch of the current epoch. Meanwhile, the data in  $C^n$  are clustered to obtain the clustering centers  $c_1^n, c_2^n$ , which are used as the initial clustering centers for all normal video segments in each batch of the current epoch. Dynamically updating the data in  $C^a$  and  $C^n$  among different epochs and introducing them into the batch clustering as a priori information can improve the effect of batch clustering and further improve the performance of the anomaly detection task.

### 2.4. Anomaly score generation based on batch clustering guidance

When the available labels are limited, the labeled samples often do not provide sufficient supervised information for the model, so the deep model is prone to overfitting. Since WSVAD only considers video-level labels, the knowledge that the model can learn will be limited. To address this problem, we use cluster labels obtained from batch clustering to guide backbone network to predict the anomaly scores of segments. The unsupervised information is effectively transferred to the weakly supervised learning process to improve the performance of the anomaly detection task.

For the labeled normal videos, each segment of these videos can simply be annotated as normal because there are no abnormal events in it. However, in the case of abnormal videos, there are also some normal events, so we use batch clustering results to generate pseudo label for each segment of the abnormal video. All segments of abnormal videos are clustered into two clusters assuming that one cluster would contain normal segments, while the other would contain abnormal. Therefore, we need to analyze which of the two clusters contains mostly normal segments and which contains mostly abnormal segments so that we can assign the appropriate pseudo labels for video segments. We obtain the pseudo labels of the segments by the similarity score between the anomaly scores and the clustering labels. Specifically, we compute the cosine similarity score  $S_1$  between the anomaly

score  $s_i \in [0, 1]$  of video  $V_i$  predicted by the backbone network and the label  $y_i^c \in \{0, 1\}$  generated by clustering, and the cosine similarity score  $S_2$  between  $s_i$  and the inverted clustering label  $\neg y_i^c$ . Finally, the pseudo label  $y_{i,j}^p$  of the  $j$ -th segment of video  $V_i$  is given by the following equation:

$$y_{i,j}^p = \begin{cases} y_{i,j}^c, & \text{if } s_1 > s_2 \\ \neg y_{i,j}^c, & \text{otherwise} \end{cases} \quad (4)$$

The pseudo labels generated by batch clustering for video segments and the anomaly scores generated by adaptive graph convolutional networks for video segments can complement each other in the anomaly detection task. In an anomalous video, if the pseudo label of a video segment is 1, the segment has a high probability of being abnormal. Therefore, we expand the anomaly score of the segments with a pseudo label 1 in the abnormal video. In particular, in an abnormal video, if the pseudo label of the video segment is 1, the abnormal score of the video segment becomes  $\min(\alpha \times s_{i,j}, 1)$ ; if the pseudo label of the video segment is 0, the abnormal score of the video segment remains unchanged. The anomaly score of each video segment is given by the following equation:

$$s_{i,j} = \begin{cases} \min(\alpha \times \tilde{s}_{i,j}, 1), & \text{if } y_i = 1 \text{ and } y_{i,j}^p = 1 \\ \tilde{s}_{i,j}, & \text{otherwise} \end{cases} \quad (5)$$

where  $\alpha$  is the expansion factor.

### 3. EXPERIMENTAL RESULTS

#### 3.1. Datasets and evaluation metrics

UCF-Crime [1] is a large-scale dataset which spans over 128 hours of surveillance videos. It covers 13 realistic anomalies. The entire dataset contains 1,900 long untrimmed videos, of which 1610 videos with video-level label are used for training and the rest for testing.

ShanghaiTech [12] is a medium-scale campus surveillance dataset containing 437 videos. Following Zhong et al. [13], we split the data into two subsets: the training set consisting of 175 normal videos and 63 anomalous videos, and the test set containing 155 normal videos and 44 anomalous videos.

Following previous work [1], we use the area under curve (AUC) of the receiver operating characteristic curve (ROC) at the frame level as the criterion for the model, and a higher AUC value indicates a better detection of the model.

#### 3.2. Experimental details

We extract the 2048-dimensional features from the ‘‘mix 5c’’ layer of I3D [10]. Following previous work [14], we extract  $T$  segments from the video uniformly to represent the whole video. By default, we set  $T$  to 150 for UCF-Crime and 100 for ShanghaiTech. The fully connected layer in the model has 512 nodes, and the graph convolutional network layer has

**Table 1.** Frame-level AUC performance comparisons.

Model	Feature	UCF-Crime	ShanghaiTech
Sultani et al. [1]	I3D RGB	76.92	86.30
Zhong et al. [13]	C3D RGB	81.08	76.44
AR_Net [2]	I3D RGB	78.96	85.38
SRF [5]	I3D RGB	79.54	84.17
Wu et al. [14]	I3D RGB	82.44	/
CLAWS [6]	C3D RGB	83.03	89.67
MIST [15]	I3D RGB	82.30	94.83
BN-SVP [16]	I3D RGB	83.39	96.00
MCR [4]	I3D RGB	81.0	90.10
MSLNet [3]	I3D RGB	85.30	96.08
Ours	I3D RGB	<b>85.87</b>	<b>96.45</b>

**Table 2.** Ablation study on UCF-Crime dataset ( $L_{bc}$ : batch clustering based loss, CBL: cross-batch learning strategy, BCG: anomaly score generation based on batch clustering guidance).

backbone	$L_{bc}$	CBL	BGG	AUC(%)
✓				84.67
✓	✓			85.21
✓	✓	✓		85.56
✓	✓	✓	✓	85.87

128, 32 and 1 nodes, respectively. Our model is trained with a mini-batch size of 64 using the Adam optimizer. For hyper-parameters,  $\alpha$  is set to 1.3.

#### 3.3. Experimental results

We compare our method with the current methods on two datasets, and show results in Table 1. For the UCF-Crime dataset, comparison results show that our method outperforms all comparison methods. Using the same I3D-RGB features, our method outperforms the previous GCN-based method by 4.79% over Zhong et al. [13] and 3.43% over Wu et al. [14]. Our method also surpasses previous methods using clustering by 6.33% over SRF [5] and 2.84% over CLAWS [6]. It also achieves better performance compared to previous weakly supervised methods on the ShanghaiTech dataset. It is 20.01% higher than the GCN-based method [13], 12.29% and 6.78% higher than the clustering-based method [5] and [6], respectively.

#### 3.4. Ablation experiments

We perform ablation study to investigate the contribution of each component of our proposed method. We start by evaluating only the backbone network and observe the performance gain while adding different modules. The result on UCF-Crime is shown in Table 2. The backbone network achieves 84.67% AUC while the addition of batch clustering based

**Table 3.** Performance comparison of batch clustering with different network layer outputs on UCF-Crime.

clustering layer	AUC(%)
FC	84.96
GCN-1	85.21
GCN-2	84.27

**Table 4.** Performance comparison of different cross-batch learning strategies on UCF-Crime.

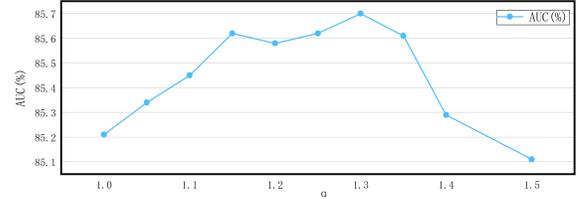
different strategies	AUC(%)
way1	85.56
way2	85.42
way3	85.51
way4	85.34

loss improves the performance to 85.21% and the addition of cross-batch learning strategy further improves the performance to 85.56% which validates their effectiveness. Addition of the anomaly score generation based on batch clustering guidance improves the performance to 85.87%, demonstrating that the pseudo labels generated by batch clustering can guide the backbone network.

We select different network layer outputs for batch clustering to investigate their effects on model performance. We choose the outputs of FC layer (FC), the first GCN layer (GCN-1), and the second GCN layer (GCN-2), respectively. As shown in Table 3, using GCN-1 achieves the best result with 0.25% higher than that using FC, probably because GCN-1 exploits the temporal relationship between video segments. The performance using GCN-2 is 0.94% lower than that using GCN-1 because GCN-1 contains more information due to its higher dimensionality compared to GCN-2.

We conduct experiments and show comparison results to evaluate the effects of different cross-batch learning strategies. The first method is the one mentioned in Section 2.3. The second method is to save the clustering centers obtained from the previous batch clustering as the initial clustering centers when clustering the current batch. The third method is to save the clustering centers obtained from all previous batch clustering and use the clustering centers obtained by clustering them again as the initial clustering centers for the current batch clustering. The fourth method is to save the clustering centers of all the batches from the previous epoch as clustering samples to participate in the clustering of each batch in the current epoch. The Table 4 shows that all our proposed cross-batch learning strategies improve the model performance. The first method has the best performance, and the fourth method performing less well, with only a small improvement.

We also conduct experiment to further investigate the ex-



**Fig. 2.** Performance comparison with different expansion factors  $\alpha$  on UCF-Crime.

pansion factor  $\alpha$ , and show the change of performance with different expansion factors  $\alpha$  on UCF-Crime in Figure 2. We can observe that the AUC first increases and then decreases as the value of  $\alpha$  increases. The model achieves the best performance when  $\alpha$  is 1.3.

### 3.5. Qualitative Result and Analysis

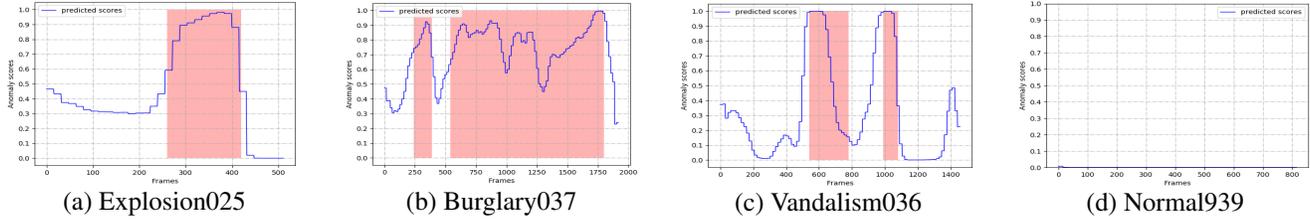
To further evaluate the performance of our method, we visualize the anomaly score curves, as shown in Figure 3. The figure shows the ground truth and the predicted anomaly scores of three abnormal videos and one normal video from UCF-Crime. It is obviously that our model exactly localizes the anomalous events, showing the effectiveness and robustness of our model. Our method successfully predicts both short-term anomalous events and long-term anomalous events. In addition, our method can also detect multiple anomalous events in a video.

## 4. CONCLUSION

In this work, we propose a novel WSVAD model based on cross-batch clustering guidance. The method enhances feature discrimination by binary batch clustering for normal and abnormal videos within a batch separately. In addition, a cross-batch learning strategy is incorporated to solve the data imbalance problem caused by the mini-batch training strategy, allowing the model to better capture the potential distribution of the data. Finally, the pseudo labels generated by batch clustering guide the backbone network to generate the anomaly scores, which further enhances the separability of normal and anomalous. The experimental results show that the proposed method achieves significant improvements on commonly-used WSVAD datasets.

## 5. REFERENCES

- [1] W. Sultani, C. Chen, and M. Shah, “Real-world anomaly detection in surveillance videos,” in *Proc. IEEE Conf. Comput. Vision Pattern Recognit.*, 2018, pp. 6479–6488.
- [2] Y. Wan, B., X. Xia, and J. Mei, “Weakly supervised video anomaly detection via center-guided discrimina-



**Fig. 3.** Visualization of the testing results on UCF-Crime. Red regions are ground truths of anomalous events.

- itive learning,” in *Proc. IEEE Int. Conf. Multimedia Expo.*, 2020, pp. 1–6.
- [3] S. Li, F. Liu, and L. Jiao, “Self-training multi-sequence learning with transformer for weakly supervised video anomaly detection,” *Proceedings of the AAAI, Virtual*, vol. 24, 2022.
- [4] Y. Gong, C. Wang, X. Dai, S. Yu, L. Xiang, and J. Wu, “Multi-scale continuity-aware refinement network for weakly supervised video anomaly detection,” in *2022 IEEE International Conference on Multimedia and Expo*, 2022, pp. 1–6.
- [5] M. Z. Zaheer, A. Mahmood, H. Shin, and S.-I. Lee, “A self-reasoning framework for anomaly detection using video-level labels,” *IEEE Signal Process. Lett.*, vol. 27, pp. 1705–1709, 2020.
- [6] 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.*, 2020, pp. 358–376.
- [7] X. Wang, H. Zhang, W. Huang, and M. R. Scott, “Cross-batch memory for embedding learning,” in *Proc. IEEE Conf. Comput. Vision Pattern Recognit.*, 2020, pp. 6388–6397.
- [8] C.-Y. Wu, C. Feichtenhofer, H. Fan, K. He, P. Krahenbuhl, and R. Girshick, “Long-term feature banks for detailed video understanding,” in *Proc. IEEE Conf. Comput. Vision Pattern Recognit.*, 2019, pp. 284–293.
- [9] C. Cao, X. Zhang, S. Zhang, P. Wang, and Y. Zhang, “Adaptive graph convolutional networks for weakly supervised anomaly detection in videos,” *arXiv preprint arXiv:2202.06503*, 2022.
- [10] J. Carreira and A. Zisserman, “Quo vadis, action recognition? a new model and the kinetics dataset,” in *Proc. IEEE Conf. Comput. Vision Pattern Recognit.*, 2017, pp. 6299–6308.
- [11] P. Perera and V. M. Patel, “Learning deep features for one-class classification,” *IEEE Transactions on Image Processing*, vol. 28, no. 11, pp. 5450–5463, 2019.
- [12] W. Luo, W. Liu, and S. Gao, “A revisit of sparse coding based anomaly detection in stacked rnn framework,” in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 341–349.
- [13] 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 *Proc. IEEE Conf. Comput. Vision Pattern Recognit.*, 2019, pp. 1237–1246.
- [14] 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.* Springer, 2020, pp. 322–339.
- [15] J.-C. Feng, F.-T. Hong, and W.-S. Zheng, “Mist: Multiple instance self-training framework for video anomaly detection,” in *Proc. IEEE Conf. Comput. Vision Pattern Recognit.*, 2021, pp. 14 009–14 018.
- [16] H. Sapkota and Q. Yu, “Bayesian nonparametric submodular video partition for robust anomaly detection,” in *Proc. IEEE Conf. Comput. Vision Pattern Recognit.*, 2022, pp. 3212–3221.