# Hierarchical Semantic Contrast for Scene-aware Video Anomaly Detection

Shengyang Sun        Xiaojin Gong*

College of Information Science & Electronic Engineering,
Zhejiang University, Hangzhou, Zhejiang, China

{sunshy,gongxj}@zju.edu.cn

## Abstract

*Increasing scene-awareness is a key challenge in video anomaly detection (VAD). In this work, we propose a hierarchical semantic contrast (HSC) method to learn a scene-aware VAD model from normal videos. We first incorporate foreground object and background scene features with high-level semantics by taking advantage of pre-trained video parsing models. Then, building upon the autoencoder-based reconstruction framework, we introduce both scene-level and object-level contrastive learning to enforce the encoded latent features to be compact within the same semantic classes while being separable across different classes. This hierarchical semantic contrast strategy helps to deal with the diversity of normal patterns and also increases their discrimination ability. Moreover, for the sake of tackling rare normal activities, we design a skeleton-based motion augmentation to increase samples and refine the model further. Extensive experiments on three public datasets and scene-dependent mixture datasets validate the effectiveness of our proposed method.*

## 1. Introduction

With the prevalence of surveillance cameras deployed in public places, video anomaly detection (VAD) has attracted considerable attention from both academia and industry. It aims to automatically detect abnormal events so that the workload of human monitors can be greatly reduced. By now, numerous VAD methods have been developed under different supervision settings, including weakly supervised [13, 50, 55, 58, 64, 76], purely unsupervised [69, 72], and ones learning from normal videos only [20, 24, 33, 44, 45]. However, it is extremely difficult or even impossible to collect sufficient and comprehensive abnormal data due to the rare occurrence of anomalies, whereas collecting abundant normal data is relatively easy. Therefore, the setting of learning from normal data is more
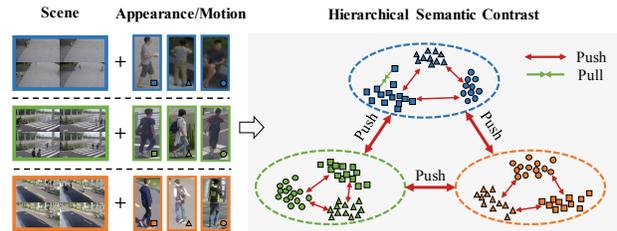


Figure 1. An illustration of hierarchical semantic contrast. The encoded scene-appearance/motion features are gathered together with respect to their semantic classes. Best viewed in color.

practical and plays the dominant role in past studies.

Although a majority of previous techniques learn their VAD models from normal data, this task has still not been well addressed due to the following reasons. First, some anomalies are scene-dependent [46, 51], implying that an appearance or motion may be anomalous in one scene but normal in other scenes. How to detect scene-dependent anomalies while preventing background bias (*i.e.* learning the background noise rather than the essence of anomaly [31]) is a challenging problem. Second, normal patterns are diverse. How to enable a deep VAD model to represent the diverse normality well but not generalize to anomalous data is also a challenge [18, 44]. Last but not least, samples collected from different normal patterns are imbalanced because some normal activities may appear very sparsely [46]. How to deal with rare but normal activities is challenging as well.

Previous VAD methods mainly perform learning at frame-level [20, 47, 75] or in an object-centric [17, 24, 78] way. The former is prone to suffer from the background bias [31] while most of the latter methods are background-agnostic. There are some attempts to address the above-mentioned challenges in one or another aspect. For instance, a spatio-temporal context graph [51] and a hierarchical scene normality-binding model [1] are constructed to discover scene-dependent anomalies. Memory-augmented autoencoders (AE) [18, 44] are designed to represent diverse normal patterns while lessening the powerful capacity of

---

*Corresponding author.

AEs. An over-sampling strategy [32] is adopted but to solve the imbalance between normal and abnormal data. Contrastively, in this work we address all of these challenges simultaneously and in distinct ways.

The primary objective of our work is to handle scene-dependent anomalies. An intuition behind scene-dependent anomalies is that, if a type of object or activity is never observed in one scene in normal videos, then it should be viewed as an anomaly. It implies that we can first determine the scene type and then check if an object or activity has occurred in normal patterns of this scene. Based on this observation, we propose a hierarchical semantic contrast method to learn a scene-aware VAD model. Taking advantage of pre-trained video parsing networks, we group the appearance and activity of objects and background scenes into semantic categories. Then, building upon the autoencoder-based reconstruction framework, we design both scene-level and object-level contrastive learning to enforce the encoded latent features to gather together with respect to their semantic categories, as shown in Fig. 1. When a test video is input, we retrieve weighted normal features for reconstruction and the clips of high errors are detected as anomalies.

The contributions of this work are as follows:

- We build a scene-aware reconstruction framework composed of scene-aware feature encoders and object-centric feature decoders for anomaly detection. The scene-aware encoders take background scenes into account while the object-centric decoders are to reduce the background noise.

- We propose hierarchical semantic contrastive learning to regularize the encoded features in the latent spaces, making normal features more compact within the same semantic classes and separable between different classes. Consequently, it helps to discriminate anomalies from normal patterns.

- We design a skeleton-based augmentation method to generate both normal and abnormal samples based on our scene-aware VAD framework. The augmented samples enable us to additionally train a binary classifier that helps to boost the performance further.

- Experiments on three public datasets demonstrate promising results on scene-independent VAD. Moreover, our method also shows a strong ability in detecting scene-dependent anomalies on self-built datasets.

## 2. Related Work

### 2.1. Video Anomaly Detection

Most previous VAD studies can be grouped into weakly supervised category [13, 50, 55, 58, 64, 76] that learns with video-level labels, or the one learning from normal videos only [18, 20, 44, 47, 75]. In this work, we focus on the latter category, which is mainly addressed by reconstruction- or distance-based techniques. The reconstruction-based techniques use autoencoder (AE) [20, 38, 75], memory-augmented AE [18, 34, 44], or generative models [42, 47] to reconstruct current frame [18, 20, 44, 75] or predict future frames [33, 42], by which the frames of high reconstruction errors are detected as anomalies. The distance-based techniques often adopt one-class SVMs [24, 25], Gaussian mixture models [49, 52], or other classifiers [17] to compute a decision boundary and those deviating from the normality are screened out as anomalies.

A majority of reconstruction- and distance-based techniques [18, 20, 25, 42, 44, 47, 75] learn their models at frame-level, which may suffer from the background bias [31] and lack of explainability. To this end, various object-centric methods have been developed, leveraging appearance and motion [16, 17, 24, 67, 78], or skeleton [28, 39, 41, 67] of objects to promote the performance. However, the VAD models learned by most of them are background-agnostic. Considering that some anomaly events are scene-dependent, a few scene-aware methods [1, 3, 51, 52] have been proposed recently. For instance, Sun *et al.* [51] construct a spatio-temporal context graph to represent both objects and the background, Sun *et al.* [52] and Bao *et al.* [1] learn memory-augmented AEs to encode scene and objects, and Cao *et al.* [3] design a network with context recovery and knowledge retrieval streams. Our work adopts the autoencoder-based reconstruction framework like [1]. But differently, we build scene-aware encoders and object-centric decoders for reconstruction and propose the hierarchical semantic contrast to regularize the encoded latent features.

### 2.2. Contrastive Learning

Contrastive learning has been successfully applied to various vision tasks, such as representation learning [7, 21, 27, 66], person re-identification [15, 61], and semantic segmentation [9, 77]. It performs learning via contrasting anchor instances with their positive and negative instances or prototypes, which are sampled either from a large batch [7] or from an external memory bank [21, 61]. Recently, contrastive learning has also been exploited in anomaly detection [23, 26, 36, 60, 63]. Most methods [23, 26, 36, 63] perform contrast between an instance and its augmented version, focusing on the instance level only. An exceptional work is HSCL [60], which takes into account sample-to-sample, sample-to-prototype, and normal-to-abnormal contrasts to implement semi-supervised anomaly detection. In our work, we design a hierarchical contrastive learning strategy that performs contrast at the scene-level and object-level, enforcing instances to gather together according to their semantic categories.
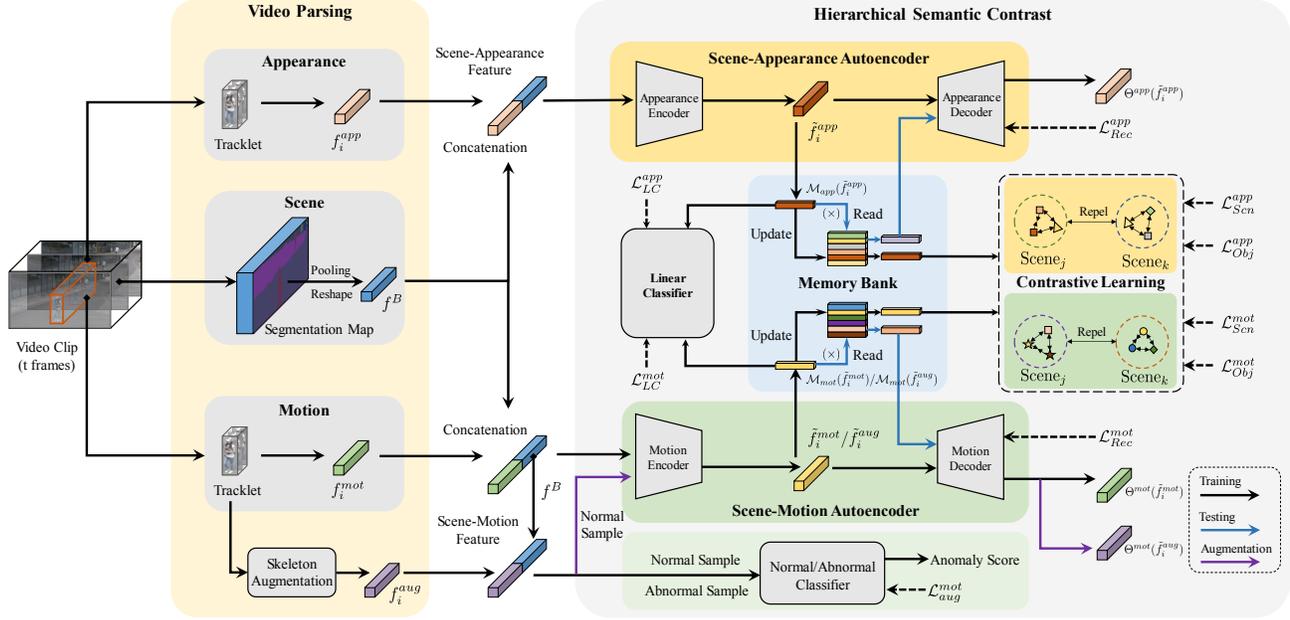
Figure 2. An overview of the proposed method. It consists of video parsing, scene-aware autoencoders, memory-based contrastive learning, and motion augmentation modules. Best viewed in color.

## 2.3. Data Augmentation

Data augmentation is extensively used in contrastive learning and other class-imbalanced learning tasks [30]. The works most related to ours are skeleton augmentation methods. For instance, Meng *et al.* [40] design a transformation network to generate new skeleton samples. Thoker *et al.* [57] design spatial and temporal skeleton augmentation based on shear transformation and joint jittering. Guo *et al.* [19] apply shear, spatial flip, rotate, and axis mask to generate extreme augmentations. These methods apply skeleton-based augmentation to generate positive samples for the action recognition task. In contrast, we design a skeleton-based augmentation to produce both normal and abnormal samples of rare activities, helping the learning of imbalanced anomaly detection.

## 3. The Proposed Method

Figure 2 presents an overview of the proposed method. When a video clip (*i.e.* a set of consecutive frames) is input, we first parse it to get high-level semantic features, including the appearance and motion of objects, together with the background scene. Then, the appearance or motion feature of each object is incorporated with the scene feature. The obtained scene-appearance and scene-motion features are fed into scene-aware encoders and object-centric decoders for feature encoding and reconstruction. All encoded latent features are stored in external memory banks, based on which we perform scene- and object-level semantic contrastive learning. The hierarchical contrastive learning enforces the diverse latent normal features to be compact within the same semantic classes and separable between different classes, which consequently increases the discrimination ability of normal patterns. During inference time, normal features stored in memory are retrieved and weighted to reconstruct the features of objects in a test clip, and those with high errors are detected as anomalies.

## 3.1. Video Parsing

Pre-trained video parsing networks are extensively used in many VAD methods [1, 16, 24, 39, 41, 51, 52, 67, 67, 78] to extract different visual cues. In this work, we take advantage of several pre-trained networks to extract high-level features while introducing semantic labels.

Given a video clip $\mathcal{C}$ composed of $T$ consecutive frames, we first adopt the pre-trained YOLOv3 [48] and Fair-MOT [74] to detect and track objects, which produce several object tracklets and their semantic class labels such as *pedestrian*, *bicycle*, *etc*. Then, we extract both appearance and motion features for each object tracklet and extract a scene feature for the remaining background as follows.

**Appearance feature extraction.** Appearance information plays an important role in detecting appearance anomalies. Therefore, for an object tracklet $\mathcal{O}_i$ in the clip, we employ ViT [8] to extract an appearance feature for each frame of the tracklet, and the features of all frames are averaged to generate one appearance feature $f_i^{app} \in \mathbb{R}^{1024}$.

**Motion feature extraction.** Motion information is of equal importance in VAD. Considering that human-related anomalies are dominant in non-traffic surveillance, we opt

3

to extract action information as a motion feature instead of using optical flow. More specifically, for an object tracklet $\mathcal{O}_i$, we use a pre-trained HRNet [54] to extract a skeleton feature for each frame. The features of all frames are further fed into PoseConv3D [10] to produce one motion feature $f_i^{mot} \in \mathbb{R}^{512}$, together with an action class label such as *walking*, *jumping*, *kicking*, *etc*.

**Scene feature extraction.** In pursuit of scene-awareness, we also extract a scene feature for the clip background. For each clip frame, we employ DeepLabV3+ [6] to generate a segmentation map while masking out the foreground object categories. Then, we perform max-pooling, reshape, averaging, and $l_2$ normalization on all segmentation maps to obtain one scene feature $f^B \in \mathbb{R}^{D_B}$, where $D_B$ depends on the size of the video frame. To discriminate different scenes at a fine-grained level, we utilize DBSCAN [11] for clustering and generating pseudo labels of scene classes.

### 3.2. Semantic Feature Reconstruction

In this work, we adopt the extensively used reconstruction framework for our anomaly detection. For each appearance or motion feature, we design an autoencoder composed of a scene-aware encoder and an object-centric decoder for feature reconstruction.

**Scene-aware feature encoder.** To correlate foreground objects with the background scene, we incorporate each appearance/motion feature with its corresponding scene feature. The obtained scene-appearance or scene-motion feature is fed into a scene-aware feature encoder. Formally, it is represented by

$$\tilde{f}_i^* = \Phi^*([f^B, f_i^*]), \quad (1)$$

where $\tilde{f}_i^* \in \mathbb{R}^{D_E}$ is the encoded latent feature of object $\mathcal{O}_i$ in clip $\mathcal{C}$, '$*$' denotes either $app$ or $mot$, and $D_E$ is the feature dimension. Moreover, $[\cdot, \cdot]$ denotes the concatenation and $\Phi^*(\cdot)$ is the feature encoder, which is implemented by a two-layer MLP followed with a $l_2$ normalization.

**Object-centric feature decoder.** The reconstruction-based framework assumes that anomalies cannot be represented well by normal patterns. To reduce the background bias [31] in reconstruction, we opt to reconstruct the feature of each foreground object instead of the incorporated scene-aware feature. That is, given a latent code $\tilde{f}_i^*$, we enforce the decoder to reconstruct a feature close to the appearance/motion feature $f_i^*$, which is

$$\mathcal{L}_{Rec}^* = \|f_i^* - \Theta^*(\tilde{f}_i^*)\|_2^2, \quad (2)$$

where $\|\cdot\|_2$ is the $l_2$ norm. $\Theta^*$ is the feature decoder implemented by a two-layer MLP as well.

### 3.3. Hierarchical Semantic Contrast

Due to the diversity of normal patterns as well as the large capacity of deep networks, the model learned from

normal data may also reconstruct anomalies well [18, 44]. To address this problem, we propose a hierarchical semantic contrast (HSC) strategy to regularize the encoded normal features in the latent space, by which diverse normal patterns can be represented more compactly and therefore be more discriminative to anomalies. HSC conducts contrastive learning at the scene- and object-level by taking advantage of the semantic labels introduced in video parsing.

**Scene-level contrastive learning.** The scene-level contrastive learning aims to attract the latent features within the same scene class and repel the features of different scenes. To this end, we adopt the InfoNCE loss [7, 66] to conduct learning, assisted by an external memory bank. The scene-level contrastive loss is defined by

$$\mathcal{L}_{Scn}^* = - \sum_{\tilde{f}_j^* \in \mathcal{X}_*(\tilde{f}_i^*)} log \frac{exp(sim(\tilde{f}_i^*, \tilde{f}_j^*)/\tau)}{\sum_{k=1}^{N} exp(sim(\tilde{f}_i^*, \tilde{f}_k^*)/\tau)}, \quad (3)$$

where $N$ is the number of all encoded latent features, $\mathcal{X}_*(\tilde{f}_i^*)$ indicates the set of features sharing the same pseudo scene label with $\tilde{f}_i^*$, $\tau$ is the temperature hyperparameter, and $sim(\cdot, \cdot)$ denotes the cosine similarity.

Besides, we also build a linear classification (LC) head to classify each latent feature into its pseudo scene class by using the cross-entropy loss:

$$\mathcal{L}_{LC}^* = -log < \mathcal{Y}(\tilde{f}_i^*), \Lambda^*(\tilde{f}_i^*) >, \quad (4)$$

where $< \cdot, \cdot >$ denotes dot product, $\Lambda^*(\cdot)$ is the linear classifier, and $\mathcal{Y}$ represents the pseudo scene label of $\tilde{f}_i^*$.

**Object-level contrastive learning.** Within each scene class, the object-level contrastive learning pulls the latent features of the same appearance/motion category together and pushes away those from different appearance/motion categories. Therefore, the object-level contrastive loss is defined by

$$\mathcal{L}_{Obj}^* = - \sum_{\tilde{f}_j^* \in \mathcal{N}_*(\tilde{f}_i^*)} log \frac{exp(sim(\tilde{f}_i^*, \tilde{f}_j^*)/\tau)}{\sum_{\tilde{f}_k^* \in \mathcal{X}_*(\tilde{f}_i^*)} exp(sim(\tilde{f}_i^*, \tilde{f}_k^*)/\tau)}, \quad (5)$$

where $\mathcal{N}_*(\tilde{f}_i^*)$ represents the set of latent features sharing the same appearance/motion class and same scene class with $\tilde{f}_i^*$. Note that in this loss only the features within the same scene class are considered and all others are ignored.

**Memory banks.** In contrast to the memory-augmented AEs [18, 44] that utilize memory for the learning of autoencoders, we use memory mainly for our contrastive learning. To this end, two memory banks are built for storing the latent scene-appearance and scene-motion features respectively. Each entry is updated by

$$\mathcal{M}_*(\tilde{f}_i^*) \leftarrow (1-m)\tilde{f}_i^* + m\mathcal{M}_*(\tilde{f}_i^*), \quad (6)$$

followed with a $l_2$ normalization, where $m \in [0, 1)$ is a momentum coefficient.
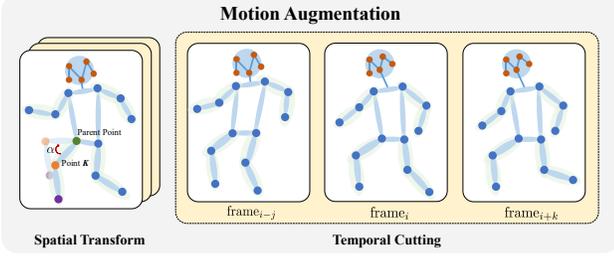
Figure 3. An illustration of our skeleton-based motion augmentation, which consists of spatial transformation and temporal cutting.

## 3.4. Motion Augmentation

The occurrence of rare but normal activities is a challenge in VAD [46]. This challenge stands out in scene-dependent anomaly detection when compared to the scene-agnostic case. The reason is that normal samples collected from different scenes are not counted together anymore. To address this problem, we design a skeleton-based augmentation to produce more samples, which includes spatial transformation and temporal cutting as shown in Fig. 3.

**Spatial transformation.** A skeleton feature extracted from one object frame contains a set of human anatomical keypoints including *shoulder*, *elbow*, *wrist*, *etc*. In this work, we design a rotation-based augmentation scheme. For each keypoint $K$ except those on *head*, we set a probability $P_{st}$ to decide if the keypoint is rotated or not. If the keypoint $K$ is chosen to rotate, it rotates around its parent node and the new coordinates $K_{rot}$ are obtained by

$$K_{rot} = (K - P(K)) \begin{bmatrix} cos(\alpha) & sin(\alpha) \\ -sin(\alpha) & cos(\alpha) \end{bmatrix} + P(K), \quad (7)$$

where $P(K)$ is the parent keypoint of $K$, and $\alpha$ is a rotation angle randomly selected within a pre-defined range. Moreover, when $K$ is rotated, its descendant keypoints are all rotated consequently.

**Temporal cutting.** An action is identified not only by the spatial distribution of keypoints but also by the temporal distribution. In this work, we simply adopt the cutting strategy for temporal augmentation. That is, given the frames of an object tracklet, we set a probability $P_{tc}$ for each frame to decide if it is left out or not.

**Spatio-temporal augmentation.** To increase the diversity of motion samples, we combine spatial transformation and temporal cutting together as our spatio-temporal augmentation. Given an object tracklet, we apply the spatio-temporal augmentation to produce a new set of skeleton features and then feed them into PoseConv3D [10] to obtain the motion feature of the augmented sample.

## 3.5. Training and Test

**Training.** The training loss of our full model contains a loss $\mathcal{L}^{app}$ for the appearance stream and a loss $\mathcal{L}^{mot}$ for the motion stream. That is, the total loss is defined by

$$\mathcal{L} = \mathcal{L}^{app} + \mathcal{L}^{mot}. \quad (8)$$

Here, the loss for each stream consists of two contrastive losses, together with a classification loss and a reconstruction loss. That is,

$$\mathcal{L}^* = \mathcal{L}^*_{Scn} + \mathcal{L}^*_{Obj} + \mathcal{L}^*_{LC} + \mathcal{L}^*_{Rec}, \quad (9)$$

in which $*$ denotes either $app$ or $mot$ as before.

At the first stage of training, we use the loss $\mathcal{L}$ to train our model on an original dataset without motion augmentation. Once the model is trained, we take augmented samples into consideration for refinement. Since the samples generated in motion augmentation are not guaranteed to be normal, we apply our trained model to discriminate normal and abnormal samples based on their reconstruction errors defined in Eq. (11). Then, we leverage both normal and abnormal samples to additionally train a binary classifier on the motion stream using a cross-entropy loss $\mathcal{L}^{mot}_{aug}$.

**Test.** During inference time, we apply video parsing to obtain high-level features for each test video clip. Then, each test feature $f_t^*$ is fed into the appearance/motion stream for encoding and reconstruction. Let us denote the encoded latent feature as $\tilde{f}_t^*$. Different from training that directly reconstructs the latent feature, we calculate the similarity between it to each entry stored in the memory $\mathcal{M}_*$ by

$$w_i = \frac{exp((\tilde{f}_t^*)^T \mathcal{M}_*(i))}{\sum_{i=1}^{N} exp((\tilde{f}_t^*)^T \mathcal{M}_*(i))}, \quad (10)$$

and get a weighted average of all stored normal features for reconstruction.

The reconstruction error of one stream is therefore defined by

$$\mathcal{S}^*(f_t^*) = \|f_t^* - \Theta^*(\sum_{i=1}^{N} w_i \mathcal{M}_*(i))\|_2^2. \quad (11)$$

The final anomaly score of an object is defined as the average reconstruction error of two streams, which is

$$\mathcal{S}(f_t^{app}, f_t^{mot}) = \frac{1}{2}(\mathcal{S}^{app}(f_t^{app}) + \mathcal{S}^{mot}(f_t^{mot})). \quad (12)$$

When motion augmentation is considered, the anomaly score of the motion stream is replaced by the anomaly probability output by the binary classifier. Moreover, the anomaly score of a clip is decided by the highest final anomaly score of objects in this clip. Finally, we apply a Gaussian filter for temporal smoothing over all video clips.

# 4. Experiments

## 4.1. Datasets and Evaluation Metrics

We evaluate the proposed method on three public datasets: UCSD Ped2 [29], Avenue [35], and ShanghaiTech [33]. UCSD Ped2 [29] is a single-scene dataset collected from pedestrian walkways, including anomalies such as *bikers*, *skaters*, small *carts* across a walkway. Avenue [35] is a single-scene dataset as well. It is captured in CUHK campus avenue, containing anomalies like *running*, *bicycling*, *etc*. It also contains some rare normal patterns [35]. ShanghaiTech [33] is a challenging multi-scene dataset containing 13 campus scenes with various light conditions and camera angles. The statistics of these datasets are summarized in Table 1.

However, these three datasets contain very few scene-dependent anomalies. And as far as we know, there is no public scene-dependent anomaly dataset available. In order to investigate the performance of our method on scene-dependent anomaly detection, we additionally create three mixture datasets based on ShanghaiTech. The mixture set $[01, 02]$ consists of videos taken from scenes 01 and 02. We select a part of test videos of scene 01 containing the *cyclist* events into the mixture training set and delete them from the test set. It implies that *cyclist* is normal in scene 01, but it is still abnormal in scene 02. Likewise, we get a mixture set $[04, 08]$ and a set $[10, 12]$, in which some events are normal in one scene but abnormal in the other scene. More details are provided in our supplementary material.

For performance evaluation, we adopt the area under the curve (AUC) of the frame-level receiver operating characteristics (ROC) as the evaluation metric following the common practice [2, 4, 12, 17, 18, 38, 45, 62]. It concatenates all frames and then computes the score, also known as micro-averaged AUC [17].

## 4.2. Implementation Details

We implement the proposed method in Pytorch. The hyper-parameters involved in our model are set as follows. The dimension of encoded latent features is $D_E = 1280$. The temperature factor in contrastive learning is $\tau = 0.5$ and the momentum coefficient in memory updating is $m =$

Table 1. The statistics of three public datasets and self-built scene-dependent datasets.

| Dataset | Training Frame | Test Frame | Scene | Resolution |
|---|---|---|---|---|
| UCSD Ped2 [29] | 2,550 | 2,010 | 1 | 360×240 |
| CUHK Avenue [35] | 15,328 | 15,324 | 1 | 640×360 |
| ShanghaiTech [33] | 274,515 | 42,883 | 13 | 856×480 |
| Mixture $[01, 02]$ | 14,080 | 5,488 | 2 | 856×480 |
| Mixture $[04, 08]$ | 37,600 | 5,104 | 2 | 856×480 |
| Mixture $[10, 12]$ | 33,856 | 3,584 | 2 | 856×480 |

Table 2. The AUC(%) performance of our model variants.

| MemCL | SA-AE | SM-AE | MA | Avenue | ShanghaiTech |
|---|---|---|---|---|---|
| | ✓ | | | 90.6 | 78.4 |
| | | ✓ | | 81.3 | 77.6 |
| | | ✓ | ✓ | 82.6 | 77.8 |
| | ✓ | ✓ | | 91.1 | 80.7 |
| | ✓ | ✓ | ✓ | 91.5 | 81.2 |
| ✓ | ✓ | | | 92.1 | 79.3 |
| ✓ | | ✓ | | 82.9 | 78.1 |
| ✓ | | ✓ | ✓ | 84.9 | 78.3 |
| ✓ | ✓ | ✓ | | 92.4 | 83.0 |
| ✓ | ✓ | ✓ | ✓ | **93.7** | **83.4** |

0.9. The probabilities used for motion augmentation are set as $P_{st} = P_{tc} = 0.5$. In addition, our model is trained using the AdaGrad optimizer with a learning rate of 0.01 and a batch size of 128 for both UCSD Ped2 and Avenue and 512 for ShanghaiTech. Some other details are provided in our supplementary material.

## 4.3. Ablation Studies

Although the proposed method aims at scene-dependent VAD, it works for scene-independent anomalies as well. Therefore, we conduct ablation studies mostly on Avenue and ShanghaiTech and partially on the mixture sets.

**Effectiveness of the proposed components.** We first validate the effectiveness of our proposed components. We decompose the full model into scene-appearance autoencoder (SA-AE), scene-motion autoencoder (SM-AE), and memory-based contrastive learning (MemCL), together with scene-motion augmentation (MA) components. The performance of the model variants holding different components is reported in Table 2. From the results, we observe that SA-AE outperforms SM-AE or SM-AE+MA when only a single stream is learned and the combination of both streams performs better. Besides, memory-based contrastive learning enables the models to outperform their counterparts by a considerable margin. Motion augmentation also improves the performance on both datasets, especially on the Avenue dataset that contains rare normal activities.

**Effectiveness of scene-aware AEs and HSC.** We here go deeper into the above-mentioned components for investigation. More specifically, we check the effectiveness of the scene-aware feature encoder (SA-E) and object-centric feature decoder (OC-D) in our autoencoders, together with the contrastive losses used in hierarchical semantic contrast (HSC). We conduct a series of experiments on the model without using motion augmentation. The results are presented in Table 3. It shows that, when contrastive learning is not applied, the scene-aware feature encoder slightly degenerates the performance on scene-independent Avenue and ShanghaiTech but improves the performance on the scene-dependent mixture sets. Moreover, the object-centric de-

**(a) Scene-appearance (w/o HSC)**  **(b) Scene-appearance (w/ HSC)**  **(c) Scene-motion (w/o HSC)**  **(d) Scene-motion (w/ HSC)**
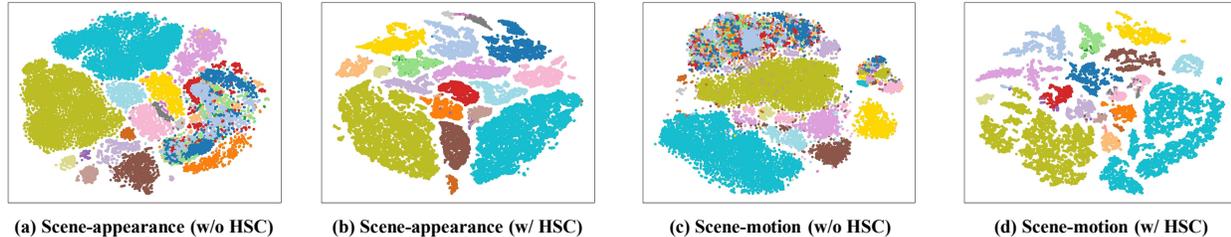
Figure 4. t-SNE [59] visualization of the scene-appearance/motion features encoded by our models without or with hierarchical semantic contrast. The points with the same color belong to the identical scene. Best viewed in color.
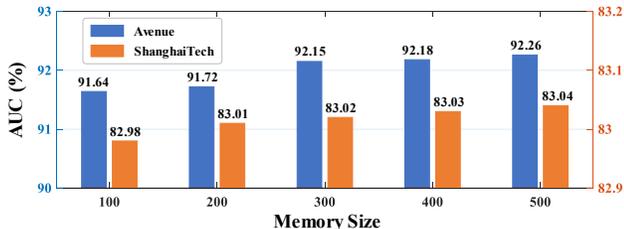


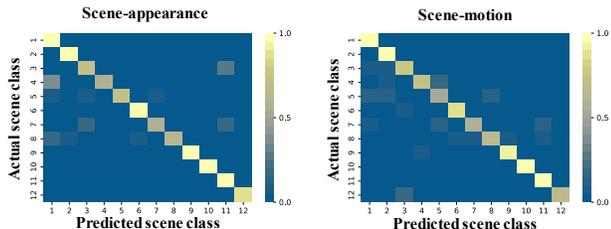Figure 5. The AUC(%) performance varies with respect to the



Figure 6. The confusion matrices of encoded scene-appearance and scene-motion features. Best viewed in color.

coder improves the performance of all datasets since the background noise in reconstruction is avoided. In HSC, the individual contrastive learning at either scene- or object-level can consistently boost the performance, indicating the necessity of regularizing encoded features in the latent space. And the best performance is achieved when the losses work together.

**Impact of the memory size at test time.** The memories in our work are used for hierarchical semantic contrast during training and feature reconstruction at test time. In order to make our model more compact and efficient for inference, we may reduce the memory size by reserving a small portion of normal patterns. In this experiment, we randomly select a number of entries and discard the remaining at test time. Fig. 5 illustrates the performance varies with the memory size. It shows that the performance is maintained well even if only 500 entries are reserved, and the performance only degenerates a little bit when only 100 entries are kept.

## 4.4. Visualization

To investigate how well the hierarchical semantic contrast strategy works, we further analyze the scene classifi-

Table 3. The AUC(%) performance of more detailed variants on CUHK Avenue (Avenue), ShanghaiTech (SHT), and the scene-dependent mixture datasets (*i.e.* [01,02] and [04,08]). When SA-E is not checked, only appearance/motion features are input to the encoders. When OC-D is not checked, the decoders reconstruct both scene and appearance/motion features.

| SA-E | OC-D | $\mathcal{L}_{Scn}$ | $\mathcal{L}_{Obj}$ | $\mathcal{L}_{LC}$ | Avenue | SHT | [01,02] | [04,08] |
|------|------|---------------------|---------------------|--------------------|--------|------|---------|---------|
|      | ✓    |                     |                     |                    | 91.0   | 80.8 | 78.6    | 76.4    |
| ✓    |      |                     |                     |                    | 90.9   | 80.2 | 80.5    | 77.9    |
| ✓    | ✓    |                     |                     |                    | 91.1   | 80.7 | 81.0    | 78.2    |
| ✓    | ✓    | ✓                   |                     |                    | 91.3   | 81.6 | 82.1    | 79.0    |
| ✓    | ✓    |                     | ✓                   |                    | 91.8   | 82.0 | 81.6    | 78.7    |
| ✓    | ✓    | ✓                   | ✓                   |                    | 91.9   | 82.2 | 82.5    | 79.4    |
| ✓    | ✓    | ✓                   |                     | ✓                  | 91.6   | 81.8 | 82.3    | 79.3    |
| ✓    | ✓    |                     | ✓                   | ✓                  | 92.2   | 82.4 | 81.8    | 78.9    |
| ✓    | ✓    | ✓                   | ✓                   | ✓                  | **92.4** | **83.0** | **82.8** | **80.0** |

cation results and the distribution of encoded latent features for data on ShanghaiTech.

**The confusion matrix of scene classification.** We first investigate whether the encoded scene-aware features correctly fall into the actual scene clusters they belong to. To this end, we check the confusion matrix of scene classification for all test samples on ShanghaiTech, which contains 12 scenes. Fig. 6 (a) and (b) visualize the confusion matrices of encoded scene-appearance and scene-motion features, respectively. We observe that most encoded scene-aware features are correctly grouped.

**The distribution of encoded latent features.** We further investigate the distribution of encoded scene-aware normal features stored in the memory banks. Fig. 4 visualizes the distribution of them in the latent space, obtained by the models without or with HSC. We observe that the features distribute more compactly within classes and more separately between classes, consequently helping to discriminate anomalies from these normal patterns.

## 4.5. Comparison to State-of-the-Art

Finally, we compare our method with state-of-the-art. The comparison is first made on three public datasets which barely contain scene-dependent anomalies. To validate the effectiveness of our method on scene-dependent anomaly detection, we additionally make a comparison on the mixture datasets created upon ShanghaiTech.

Table 4. Comparison results on UCSD Ped2 (Ped2), CUHK Avenue (Avenue), and ShanghaiTech (SHT). Besides the frame-level micro-averaged AUC(%) performance, we also list the inputs of the methods, in which 'F' denotes the frame-level input and 'O' is object-centric. The subscript 'A' is appearance, 'F' is optical flow, 'S' is skeleton, and 'M' is other motion information. Besides, in our HSC model, $MA^{-,+}$ denotes using motion augmentation to generate both normal and abnormal samples.

| Method | Reference | Input | Ped2 | Avenue | SHT |
|---|---|---|---|---|---|
| AMC [43] | ICCV19 | F | 96.2 | 86.9 | - |
| Mem-AE [18] | ICCV19 | F | 94.1 | 83.3 | 71.2 |
| DeepOC [65] | TNNLS19 | F | 96.9 | 86.6 | - |
| r-GAN [37] | ECCV20 | F | 96.2 | 85.8 | 77.9 |
| CDAE [4] | ECCV20 | F | 96.5 | 86.0 | 73.3 |
| MNAD [44] | CVPR20 | F | 97.0 | 88.5 | 72.8 |
| IPR [56] | PRL20 | F | 96.3 | 85.1 | 73.0 |
| LDF [45] | WACV20 | F | 94.0 | 87.2 | - |
| CAC [63] | MM20 | F | - | 87.0 | 79.3 |
| CT-D2GAN [14] | MM21 | F | 97.2 | 85.9 | 77.7 |
| AMMCN [2] | AAAI21 | F | 96.6 | 86.6 | 73.7 |
| MPN [38] | CVPR21 | F | 96.9 | 89.5 | 73.8 |
| AEP [71] | TNNLS21 | F | 97.9 | 90.2 | - |
| SIGnet [12] | TNNLS22 | F | 96.2 | 86.8 | - |
| IAAN [73] | TCSVT22 | F | 92.9 | 80.5 | 80.3 |
| ROADMAP [62] | TNNLS22 | F | 96.4 | 88.3 | 76.6 |
| GEPC [39] | CVPR20 | $O_S$ | - | - | 76.1 |
| STGformer [22] | MM22 | $O_S$ | - | 88.8 | 82.9 |
| HSNBM [1] | MM22 | $F+O_A$ | 95.2 | 91.6 | 76.5 |
| STC-Graph [51] | MM20 | $O_A+O_M$ | - | 89.6 | 74.7 |
| $SSMTL^1$ [16] | CVPR21 | $O_A+O_M$ | 97.5 | 91.5 | 82.4 |
| VEC [70] | MM20 | $O_A+O_F$ | 97.3 | 90.2 | 74.8 |
| $HF^2$-VAD [34] | ICCV21 | $O_A+O_F$ | **99.3** | 91.1 | 76.2 |
| BAF [17] | TPAMI22 | $O_A+O_F$ | 98.7 | 92.3 | 82.7 |
| CAFE [68] | MM22 | $O_A+O_F$ | 98.4 | 92.6 | 77.0 |
| DERN [53] | MM22 | $O_A+O_F$ | 97.1 | 92.7 | 79.3 |
| BDPN [5] | AAAI22 | $O_A+O_F$ | 98.3 | 90.3 | 78.1 |
| HSC | This work | $O_A+O_S$ | $98.1^\dagger$ | 92.4 | 83.0 |
| HSC w/ $MA^{-,+}$ | | | | **93.7** | **83.4** |

Table 5. Comparison results on the scene-dependent mixture datasets built upon ShanghaiTech. $MA^{-,+}$ denotes using motion augmentation to generate both normal and abnormal samples, $MA^-$ denotes augmenting normal samples only.

| Method | Reference | Input | [01,02] | [04,08] | [10,12] |
|---|---|---|---|---|---|
| Mem-AE [18] | ICCV19 | F | 77.7 | 60.2 | 50.2 |
| MNAD [44] | CVPR20 | F | 77.8 | 68.6 | 50.0 |
| MPN [38] | CVPR21 | F | 78.4 | 61.5 | 45.3 |
| $HF^2$-VAD [34] | ICCV21 | $O_A+O_F$ | 74.8 | 75.2 | 66.8 |
| HSC | This work | $O_A+O_S$ | 82.8 | 80.0 | 87.3 |
| HSC w/ $MA^-$ | | | 85.7 | 81.8 | 90.1 |
| HSC w/ $MA^{-,+}$ | | | **86.9** | **82.6** | **91.0** |

for testing, which still achieves a performance higher than a great number of methods.

**Comparison on scene-dependent VAD.** Finally, we investigate the performance of our method on scene-dependent anomaly detection based on the mixture datasets introduced above. The results are presented in Table 5. We also test other SOTA methods [18, 34, 38, 44] using their released codes for comparison. Unfortunately, we are not able to compare with the scene-aware methods [1,5,51] since their codes are not available. The results show that the performance of the other methods, especially those with frame-level inputs, degenerates dramatically. In contrast, all model variants of our proposed method consistently demonstrate promising performance.

## 4.6. Limitations

Since our proposed method takes skeleton-based motion features as one of the inputs, the full model is restricted to human-related datasets and requires the datasets are not very low in resolution. Otherwise, only the scene-appearance stream works, which inevitably degenerates the performance. A possible extension is replacing the skeleton-based motion features with optical flow and conducting contrastive learning based on the clustering results of optical flow features. Besides, other components of this framework can be replaced by other advanced modules, e.g. substituting another advanced background parsing model for the simple segmentation map.

## 5. Conclusion

In this work, we have presented a hierarchical semantic contrast method to address scene-dependent video anomaly detection. The design of our hierarchical semantic contrastive learning, together with scene-aware autoencoders and motion augmentation, enables the proposed model to achieve promising results on both scene-independent and scene-dependent VAD. Experiments on three public datasets and self-created datasets have validated the effectiveness of our method.

**Comparison on scene-independent VAD.** We compare our method with recent VAD methods that learn from normal data as well. The comparison results are presented in Table 4, in which the inputs of all methods are also provided for reference. Generally speaking, benefiting from the extracted high-level features, most of the object-centric methods perform better than the methods with frame-level inputs, although some of the latter methods also use motion information like optical flow. In addition, the proposed method outperforms all other methods on both Avenue and ShanghaiTech, validating the superiority of our design. Note that we are not able to test our full model on the UCSD Ped2 dataset due to its low resolution, in which no high-quality skeleton keypoints can be detected. Therefore, we only use the scene-appearance stream of our model

---

[1] Here the micro-average AUC is reported from the officially released website https://github.com/lilygeorgescu/AED-SSMTL.

# References

[1] Qianyue Bao, Fang Liu, Yang Liu, Licheng Jiao, Xu Liu, and Lingling Li. Hierarchical scene normality-binding modeling for anomaly detection in surveillance videos. In *ACM MM*, pages 6103–6112, 2022. 1, 2, 3, 8

[2] Ruichu Cai, Hao Zhang, Wen Liu, Shenghua Gao, and Zhifeng Hao. Appearance-motion memory consistency network for video anomaly detection. In *AAAI*, volume 35, pages 938–946, 2021. 6, 8

[3] Congqi Cao, Yue Lu, and Yanning Zhang. Context recovery and knowledge retrieval: A novel two-stream framework for video anomaly detection. *arXiv preprint arXiv:2209.02899*, 2022. 2

[4] Yunpeng Chang, Zhigang Tu, Wei Xie, and Junsong Yuan. Clustering driven deep autoencoder for video anomaly detection. In *ECCV*, pages 1–16, 2020. 6, 8

[5] Chengwei Chen, Yuan Xie, Shaohui Lin, Angela Yao, Guannan Jiang, Wei Zhang, Yanyun Qu, Ruizhi Qiao, Bo Ren, and Lizhuang Ma. Comprehensive regularization in a bi-directional predictive network for video anomaly detection. In *AAAI*, pages 1–9, 2022. 8

[6] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *ECCV*, pages 801–818, 2018. 4

[7] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, pages 1597–1607. PMLR, 2020. 2, 4

[8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 3

[9] Ye Du, Zehua Fu, Qingjie Liu, and Yunhong Wang. Weakly supervised semantic segmentation by pixel-to-prototype contrast. In *CVPR*, pages 4320–4329, 2022. 2

[10] Haodong Duan, Yue Zhao, Kai Chen, Dahua Lin, and Bo Dai. Revisiting skeleton-based action recognition. In *CVPR*, pages 2969–2978, 2022. 4, 5

[11] Martin Ester, Hans-Peter Kriegel, Jörg Sander, Xiaowei Xu, et al. A density-based algorithm for discovering clusters in large spatial databases with noise. In *kdd*, volume 96, pages 226–231, 1996. 4

[12] Zhiwen Fang, Jiafei Liang, Joey Tianyi Zhou, Yang Xiao, and Feng Yang. Anomaly detection with bidirectional consistency in videos. *IEEE Transactions on Neural Networks and Learning Systems*, 33(3):1079–1092, 2022. 6, 8

[13] Jia-Chang Feng, Fa-Ting Hong, and Wei-Shi Zheng. Mist: Multiple instance self-training framework for video anomaly detection. In *CVPR*, pages 14009–14018, 2021. 1, 2

[14] Xinyang Feng, Dongjin Song, Yuncong Chen, Zhengzhang Chen, Jingchao Ni, and Haifeng Chen. Convolutional transformer based dual discriminator generative adversarial networks for video anomaly detection. In *ACM MM*, page 5546–5554, 2021. 8

[15] Yixiao Ge, Dapeng Chen, Feng Zhu, Rui Zhao, and Hongsheng Li. Self-paced contrastive learning with hybrid memory for domain adaptive object re-id. In *NeurIPS*, volume 33, pages 11309–11321, 2020. 2

[16] Mariana-Iuliana Georgescu, Antonio Barbalau, Radu Tudor Ionescu, Fahad Shahbaz Khan, Marius Popescu, and Mubarak Shah. Anomaly detection in video via self-supervised and multi-task learning. In *CVPR*, pages 12742–12752, 2021. 2, 3, 8

[17] Mariana Iuliana Georgescu, Radu Tudor Ionescu, Fahad Shahbaz Khan, Marius Popescu, and Mubarak Shah. A background-agnostic framework with adversarial training for abnormal event detection in video. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(9):4505–4523, 2022. 1, 2, 6, 8

[18] Dong Gong, Lingqiao Liu, Vuong Le, Budhaditya Saha, Moussa Reda Mansour, Svetha Venkatesh, and Anton van den Hengel. Memorizing normality to detect anomaly: Memory-augmented deep autoencoder for unsupervised anomaly detection. In *ICCV*, pages 1705–1714, 2019. 1, 2, 4, 6, 8

[19] Tianyu Guo, Hong Liu, Zhan Chen, Mengyuan Liu, Tao Wang, and Runwei Ding. Contrastive learning from extremely augmented skeleton sequences for self-supervised action recognition. In *AAAI*, volume 36, pages 762–770, 2022. 3

[20] Mahmudul Hasan, Jonghyun Choi, Jan Neumann, Amit K Roy-Chowdhury, and Larry S Davis. Learning temporal regularity in video sequences. In *CVPR*, pages 733–742, 2016. 1, 2

[21] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, pages 9729–9738, 2020. 2

[22] Chao Huang, Yabo Liu, Zheng Zhang, Chengliang Liu, Jie Wen, Yong Xu, and Yaowei Wang. Hierarchical graph embedded pose regularity learning via spatio-temporal transformer for abnormal behavior detection. In *ACM MM*, pages 307–315, 2022. 8

[23] Chao Huang, Zhihao Wu, JieWen, Yong Xu, Qiuping Jiang, and Yaowei Wang. Abnormal event detection using deep contrastive learning for intelligent video surveillance system. *IEEE Transactions on Industrial Informatics*, 18(8):5171–5179, 2022. 2

[24] Radu Tudor Ionescu, Fahad Shahbaz Khan, Mariana-Iuliana Georgescu, and Ling Shao. Object-centric auto-encoders and dummy anomalies for abnormal event detection in video. In *CVPR*, pages 7842–7851, 2019. 1, 2, 3

[25] Radu Tudor Ionescu, Sorina Smeureanu, Marius Popescu, and Bogdan Alexe. Detecting abnormal events in video using narrowed normality clusters. In *WACV*, pages 1951–1960. IEEE, 2019. 2

[26] Okan Kopuklu, Jiapeng Zheng, Hang Xu, and Gerhard Rigoll. Driver anomaly detection: A dataset and contrastive learning approach. In *WACV*, pages 91–100, 2021. 2

[27] Junnan Li, Pan Zhou, Caiming Xiong, and Steven Hoi. Prototypical contrastive learning of unsupervised representations. In *ICLR*, 2020. 2

[28] Nanjun Li, Faliang Chang, and Chunsheng Liu. A self-trained spatial graph convolutional network for unsupervised human-related anomalous event detection in complex scenes. *IEEE Transactions on Cognitive and Developmental Systems*, 2022. 2

[29] Weixin Li, Vijay Mahadevan, and Nuno Vasconcelos. Anomaly detection and localization in crowded scenes. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(1):18–32, 2013. 6

[30] Swee Kiat Lim, Yi Loo, Ngoc-Trung Tran, Ngai-Man Cheung, Gemma Roig, and Yuval Elovici. Doping: Generative data augmentation for unsupervised anomaly detection with gan. In *ICDM*, pages 1122–1127. IEEE, 2018. 3

[31] Kun Liu and Huadong Ma. Exploring background-bias for anomaly detection in surveillance videos. In *ACM MM*, pages 1490–1499, 2019. 1, 2, 4

[32] Wen Liu, Weixin Luo, Zhengxin Li, Peilin Zhao, and Shenghua Gao. Margin learning embedded prediction for video anomaly detection with a few anomalies. In *IJCAI*, pages 3023–3030, 2019. 2

[33] Wen Liu, Weixin Luo, Dongze Lian, and Shenghua Gao. Future frame prediction for anomaly detection–a new baseline. In *CVPR*, pages 6536–6545, 2018. 1, 2, 6

[34] Zhian Liu, Yongwei Nie, Chengjiang Long, Qing Zhang, and Guiqing Li. A hybrid video anomaly detection framework via memory-augmented flow reconstruction and flow-guided frame prediction. In *ICCV*, pages 13588–13597, 2021. 2, 8

[35] Cewu Lu, Jianping Shi, and Jiaya Jia. Abnormal event detection at 150 fps in matlab. In *ICCV*, pages 2720–2727, 2013. 6

[36] Yue Lu, Congqi Cao, Yifan Zhang, and Yanning Zhang. Learnable locality-sensitive hashing for video anomaly detection. *IEEE Transactions on Circuits and Systems for Video Technology*, 2022. 2

[37] Yiwei Lu, Frank Yu, Mahesh Kumar Krishna Reddy, and Yang Wang. Few-shot scene-adaptive anomaly detection. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *ECCV*, pages 125–141, Cham, 2020. Springer International Publishing. 8

[38] Hui Lv, Chen Chen, Zhen Cui, Chunyan Xu, Yong Li, and Jian Yang. Learning normal dynamics in videos with meta prototype network. In *CVPR*, pages 15425–15434, 2021. 2, 6, 8

[39] Amir Markovitz, Gilad Sharir, Itamar Friedman, Lihi Zelnik-Manor, and Shai Avidan. Graph embedded pose clustering for anomaly detection. In *CVPR*, pages 10539–10547, 2020. 2, 3, 8

[40] Fanyang Meng, Hong Liu, Yongsheng Liang, Juanhui Tu, and Mengyuan Liu. Sample fusion network: An end-to-end data augmentation network for skeleton-based human action recognition. *IEEE Transactions on Image Processing*, 28(11):5281–5295, 2019. 3

[41] Romero Morais, Vuong Le, Truyen Tran, Budhaditya Saha, Moussa Mansour, and Svetha Venkatesh. Learning regularity in skeleton trajectories for anomaly detection in videos. In *CVPR*, pages 11996–12004, 2019. 2, 3

[42] Khac-Tuan Nguyen, Dat-Thanh Dinh, Minh N. Do, and Minh-Triet Tran. Anomaly detection in traffic surveillance videos with gan-based future frame prediction. In *ICMR*, pages 457–463, 2020. 2

[43] Trong-Nguyen Nguyen and Jean Meunier. Anomaly detection in video sequence with appearance-motion correspondence. In *ICCV*, pages 1273–1283, 2019. 8

[44] Hyunjong Park, Jongyoun Noh, and Bumsub Ham. Learning memory-guided normality for anomaly detection. In *CVPR*, pages 14372–14381, 2020. 1, 2, 4, 8

[45] Bharathkumar Ramachandra, Michael Jones, and Ranga Vatsavai. Learning a distance function with a siamese network to localize anomalies in videos. In *WACV*, pages 2598–2607, 2020. 1, 6, 8

[46] Bharathkumar Ramachandra, Michael J. Jones, and Ranga Raju Vatsavai. A survey of single-scene video anomaly detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(5):2293–2312, 2022. 1, 5

[47] Mahdyar Ravanbakhsh, Enver Sangineto, Moin Nabi, and Nicu Sebe. Training adversarial discriminators for cross-channel abnormal event detection in crowds. In *WACV*, pages 1896–1904, 2019. 1, 2

[48] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*, 2018. 3

[49] Mohammad Sabokrou, Mohsen Fayyaz, Mahmood Fathy, and Reinhard Klette. Deep-cascade: Cascading 3d deep neural networks for fast anomaly detection and localization in crowded scenes. *IEEE Transactions on Image Processing*, 26(4):1992–2004, 2017. 2

[50] Waqas Sultani, Chen Chen, and Mubarak Shah. Real-world anomaly detection in surveillance videos. In *CVPR*, pages 6479–6488, 2018. 1, 2

[51] Che Sun, Yunde Jia, Yao Hu, and Yuwei Wu. Scene-aware context reasoning for unsupervised abnormal event detection in videos. In *ACM MM*, pages 184–192, 2020. 1, 2, 3, 8

[52] Che Sun, Yunde Jia, and Yuwei Wu. Evidential reasoning for video anomaly detection. In *ACM MM*, pages 2106–2114, 2022. 2, 3

[53] Che Sun, Yunde Jia, and Yuwei Wu. Evidential reasoning for video anomaly detection. In *ACM MM*, pages 2106–2114, 2022. 8

[54] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep high-resolution representation learning for human pose estimation. In *CVPR*, pages 5693–5703, 2019. 4

[55] Shengyang Sun and Xiaojin Gong. Long-short temporal co-teaching for weakly supervised video anomaly detection. In *ICME*, pages 1–6. IEEE, 2023. 1, 2

[56] Yao Tang, Lin Zhao, Shanshan Zhang, Chen Gong, Guangyu Li, and Jian Yang. Integrating prediction and reconstruction for anomaly detection. *Pattern Recognition Letters*, 129:123–130, 2020. 8

[57] Fida Mohammad Thoker, Hazel Doughty, and Cees GM Snoek. Skeleton-contrastive 3d action representation learning. In *ACM MM*, pages 1655–1663, 2021. 3

[58] Yu Tian, Guansong Pang, Yuanhong Chen, Rajvinder Singh, Johan W Verjans, and Gustavo Carneiro. Weakly-supervised video anomaly detection with robust temporal feature magnitude learning. In *ICCV*, pages 4975–4986, 2021. 1, 2

[59] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008. 7

[60] Gaoang Wang, Yibing Zhan, Xinchao Wang, Mingli Song, and Klara Nahrstedt. Hierarchical semi-supervised contrastive learning for contamination-resistant anomaly detection. In *ECCV*, pages 110–128. Springer, 2022. 2

[61] Menglin Wang, Jiachen Li, Baisheng Lai, Xiaojin Gong, and Xian-Sheng Hua. Offline-online associated camera-aware proxies for unsupervised person re-identification. *IEEE Transactions on Image Processing*, 31:6548–6561, 2022. 2

[62] Xuanzhao Wang, Zhengping Che, Bo Jiang, Ning Xiao, Ke Yang, Jian Tang, Jieping Ye, Jingyu Wang, and Qi Qi. Robust unsupervised video anomaly detection by multipath frame prediction. *IEEE Transactions on Neural Networks and Learning Systems*, 33(6):2301–2312, 2022. 6, 8

[63] Ziming Wang, Yuexian Zou, and Zeming Zhang. Cluster attention contrast for video anomaly detection. In *ACM MM*, pages 2463–2471, 2020. 2, 8

[64] Jie Wu, Wei Zhang, Guanbin Li, Wenhao Wu, Xiao Tan, Yingying Li, Errui Ding, and Liang Lin. Weakly-supervised spatio-temporal anomaly detection in surveillance video. In Zhi-Hua Zhou, editor, *IJCAI*, pages 1172–1178. International Joint Conferences on Artificial Intelligence Organization, 8 2021. Main Track. 1, 2

[65] Peng Wu, Jing Liu, and Fang Shen. A deep one-class neural network for anomalous event detection in complex scenes. *IEEE transactions on neural networks and learning systems*, 31(7):2609–2622, 2019. 8

[66] Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. Unsupervised feature learning via non-parametric instance discrimination. In *CVPR*, pages 3733–3742, 2018. 2, 4

[67] Yuxing Yang, Zeyu Fu, and Syed Mohsen Naqvi. A two-stream information fusion approach to abnormal event detection in video. In *ICASSP*, pages 5787–5791, 2022. 2, 3

[68] Guang Yu, Siqi Wang, Zhiping Cai, Xinwang Liu, and Chengkun Wu. Effective video abnormal event detection by learning a consistency-aware high-level feature extractor. In *ACM MM*, pages 6337–6346, 2022. 8

[69] Guang Yu, Siqi Wang, Zhiping Cai, Xinwang Liu, Chuanfu Xu, and Chengkun Wu. Deep anomaly discovery from unlabeled videos via normality advantage and self-paced refinement. In *CVPR*, pages 13987–13998, 2022. 1

[70] Guang Yu, Siqi Wang, Zhiping Cai, En Zhu, Chuanfu Xu, Jianping Yin, and Marius Kloft. Cloze test helps: Effective video anomaly detection via learning to complete video events. In *ACM MM*, pages 583–591, 2020. 8

[71] Jongmin Yu, Younkwan Lee, Kin Choong Yow, Moongu Jeon, and Witold Pedrycz. Abnormal event detection and localization via adversarial event prediction. *IEEE Transactions on Neural Networks and Learning Systems*, 2021. 8

[72] M Zaigham Zaheer, Arif Mahmood, M Haris Khan, Mattia Segu, Fisher Yu, and Seung-Ik Lee. Generative cooperative learning for unsupervised video anomaly detection. In *CVPR*, pages 14744–14754, 2022. 1

[73] Sijia Zhang, Maoguo Gong, Yu Xie, AK Qin, Hao Li, Yuan Gao, and Yew-Soon Ong. Influence-aware attention networks for anomaly detection in surveillance videos. *IEEE Transactions on Circuits and Systems for Video Technology*, 2022. 8

[74] Yifu Zhang, Chunyu Wang, Xinggang Wang, Wenjun Zeng, and Wenyu Liu. Fairmot: On the fairness of detection and re-identification in multiple object tracking. *International Journal of Computer Vision*, 129:3069–3087, 2021. 3

[75] Yiru Zhao, Bing Deng, Chen Shen, Yao Liu, Hongtao Lu, and Xian-Sheng Hua. Spatio-temporal autoencoder for video anomaly detection. In *ACM MM*, pages 1933–1941, 2017. 1, 2

[76] Jia-Xing Zhong, Nannan Li, Weijie Kong, Shan Liu, Thomas H Li, and Ge Li. Graph convolutional label noise cleaner: Train a plug-and-play action classifier for anomaly detection. In *CVPR*, pages 1237–1246, 2019. 1, 2

[77] Tianfei Zhou, Meijie Zhang, Fang Zhao, and Jianwu Li. Regional semantic contrast and aggregation for weakly supervised semantic segmentation. In *CVPR*, pages 4299–4309, 2022. 2

[78] Wenhao Zhou, Yingxuan Li, and Chunhui Zhao. Object-guided and motion-refined attention network for video anomaly detection. In *ICME*, pages 1–6. IEEE, 2022. 1, 2, 3