

Review

Not peer-reviewed version

A Review of Crowd Abnormal Behavior Recognition Technology Based on Computer Vision

[Rongyong Zhao](#), [Feng Hua](#)^{*}, Bingyu Wei, [Cuiling Li](#), [Yulong Ma](#), Eric S. W. Wong, Fengnian Liu

Posted Date: 24 September 2024

doi: 10.20944/preprints202409.1879.v1

Keywords: Computer Vision; Abnormal Behavior Recognition; Crowd Abnormal Behavior; Deep Learning; Neural Network; Self-Attention Mechanism



Preprints.org is a free multidiscipline platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This is an open access article distributed under the Creative Commons Attribution License which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Review

A Review of Crowd Abnormal Behavior Recognition Technology Based on Computer Vision

Rongyong Zhao ¹, Feng Hua ^{1,*}, Bingyu Wei ¹, Cuiling Li ¹, Yulong Ma ¹, Eric S. W. Wong ² and Fengnian Liu ³

¹ School of Electronic and Information Engineering, Tongji University, Shanghai, 201804, China; zhaorongyong@tongji.edu.cn; licuiling@tongji.edu.cn; evanma@tongji.edu.cn; sw.wong@connect.polyu.hk; Fengnian.Liu@karon-valve.com;

² Hongkong Institute of Water and Sanitation Safety, Hong Kong, 999077, China;

³ Shanghai KARON ECO-VALVE Manufacturing Co., Ltd. Shanghai, 201804, China.

* Correspondence: 2341070@tongji.edu.cn

Abstract: Crowd abnormal behavior recognition is one of the research hotspots in computer vision. Its goal is to use computer vision technology and abnormal behavior detection models to accurately perceive, predict, and intervene in potential abnormal behaviors of the crowd, and monitor the status of the crowd system in public places in real-time, to effectively prevent and deal with public security risks and ensure public life safety and social order. To this end, focusing on the crowd abnormal behavior recognition technology in the computer vision system, a systematic review study of its theory and cutting-edge technology is conducted. First, the crowd level and abnormal behaviors in public places are defined, and the challenges faced by crowd abnormal behavior recognition are expounded. Then, from the dimensions based on traditional methods and based on deep learning, the mainstream technologies of abnormal behavior recognition are discussed, and the design ideas, advantages, and limitations of various methods are analyzed. Next, the mainstream software tools are introduced to provide a comprehensive reference for the technical framework; Secondly, typical abnormal behavior datasets at home and abroad are sorted out, and the characteristics of these datasets are compared in detail from multiple perspectives such as scale, characteristics and uses, and the performance indicators of different algorithms on the datasets are compared and analyzed; Finally, the full text is summarized and the future development direction of crowd abnormal behavior recognition technology is prospected.

Keywords: computer vision; abnormal behavior recognition; crowd abnormal behavior; deep learning; neural network; self-attention mechanism

1. Introduction

With the acceleration of urbanization in China and the improvement of safety requirements in public places, the importance of abnormal behavior recognition in dense crowds has become increasingly prominent. At the same time, various large-scale public activities and large-scale crowd gatherings show a continuous growth trend, including various exhibitions, entertainment performances, and sports events. In these scenarios, the flow of people surges, the personnel is highly dense and the composition is diverse and complex. Once an abnormal event or emergency occurs, it is very easy to cause panic within the crowd, which in turn leads to situations such as rapid movement of the crowd, crowded collisions, and even chaotic pushing and shoving. In this case, the crowd will fall into an unstable state and even serious stampede accidents may occur, causing a tragedy of a large number of casualties [1].

In the early computer vision tasks, artificial feature descriptors played an important role, especially in the recognition of abnormal behaviors in crowds [2]. Through artificial feature descriptors, key attributes with discrimination and invariance are extracted from image data to

represent the characteristics of pedestrians and the movement state of the crowd. However, when dealing with images with large type differences or processing complex visual scenes (such as complex backgrounds, target occlusion, and drastic changes in illumination), traditional manual feature methods can usually only capture low-level local texture and shape information in the image, and cannot maintain stability and effectiveness, resulting in the performance of most abnormal detection methods being restricted by the above factors. Compared with traditional recognition methods, advanced computer vision technologies such as deep learning can automatically extract image features and better model time series data, capture the temporal and spatial characteristics of the dynamic changes of the crowd, and automatically distinguish abnormal and normal behaviors of the crowd in video frames.

In recent years, automatic recognition of abnormal behaviors in crowds based on computer vision has become a research hotspot. Literature [3] uses the FCM clustering algorithm to group the key points of the trajectory and construct a feature histogram of the cluster group motion pattern, visualizing the motion pattern features in the form of high-dimensional coordinates. Literature [4] proposes a crowd abnormal behavior detection algorithm SFCNN-ABD based on the streamlined flow convolutional neural network. Literature [5] detects possible group abnormal events, such as panic, pushing, gathering, and other dangerous behaviors by real-time analysis of the characteristics of pedestrian movement, density changes, and interaction behaviors in the surveillance video. Literature [6] details the video-based human abnormal behavior recognition and detection technology, covering many aspects such as the classification of abnormal behaviors, feature extraction methods, and discrimination methods. Literature [7] reviews the non-invasive human fall detection system based on deep learning and elaborates in detail on the performance indicators of models such as CNN, Auto-Encoder, LSTM, and MLP. Literature [8] combs the research progress in the field of abnormal detection of surveillance videos, including many aspects such as abnormal classification, detection methods, feature representation, model construction, and evaluation criteria. Literature [3–8] represents the progress of different stages in the field of abnormal behavior recognition of crowds.

It is worth noting that there are three deficiencies in the existing literature research: 1) The analysis of abnormal behaviors in the literature is not comprehensive, and each literature only focuses on the research results of a specific development stage; 2) The pedestrian abnormal behavior detection technology has not been systematically included, and the key factors such as the characteristics and limitations of various methods have not been compared in detail; 3) For the visual problems that are common in complex scenes, such as dense crowds and severe occlusion between individuals, few existing literatures have proposed targeted and efficient visual occlusion resolution strategies.

Therefore, given the above deficiencies, this paper conducts a comprehensive and systematic analysis of the automatic recognition technology of abnormal behaviors of crowds in the field of computer vision, so that researchers in this field can better grasp the current situation and development direction of abnormal behavior detection of crowds. The main contributions of this paper are summarized as follows:

(1) Comprehensively summarize the traditional methods of abnormal behavior detection, and deeply explore how to use mathematical models to capture and analyze key characteristics such as the speed, direction, abnormal movement of individuals and the spatial layout, flow of people, flow speed, and direction of the group.

(2) Starting from the design concept and focus of the algorithm core, the abnormal behavior detection methods based on deep learning are divided into five categories. In addition, for the scale change and occlusion problems in the crowd, this paper provides a novel dynamic pedestrian centroid model and a visual occlusion resolution strategy.

(3) Provide four types of representative software tools and provide an experimental platform and practical tool reference for researchers in related fields.

(4) Summarize four international public datasets of typical abnormal behaviors of crowds, and compare these datasets in detail from multiple perspectives such as scale, annotation, and main uses;

and prospect the future development direction of abnormal behavior recognition from five aspects such as multi-modal fusion, multi-source multi-dimensional data fusion, and metaverse evolution model.

The structure of this paper is as follows: Section 1 gives the definitions of crowd levels and abnormal behaviors; Section 2 introduces the main challenges faced by the abnormal behavior recognition task in the dense crowd scene; Section 3 comprehensively summarizes the methods of abnormal behavior recognition in recent years from the two dimensions of traditional methods and deep learning; and further introduces the current mainstream software tools; Section 4 introduces the datasets widely used in the field of abnormal behavior detection at home and abroad and the performance indicators of each algorithm on these datasets. Section 5 summarizes this paper and presents the future development trend of this research field. The article framework is shown in Figure 1.

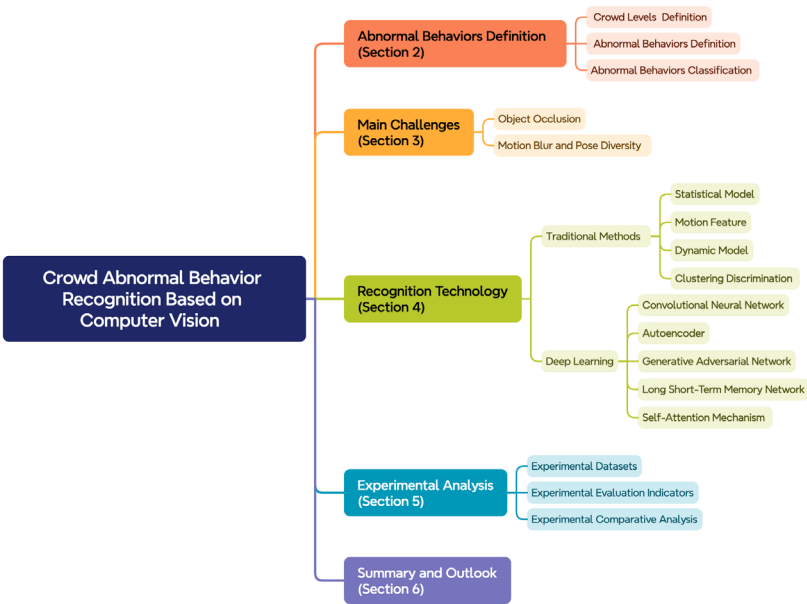


Figure 1. Paper structure diagram.

2. Abnormal Behaviors of Pedestrians in Public Places

2.1. Definition of Crowd Levels

High-density crowds usually refer to the dense concentration of pedestrians in a specific space or area, such that the distance between individuals is small, which may lead to restricted movement, blocked sight, or poor air circulation. In this case, the density of the crowd has reached a level that may affect public safety, increasing the risks of stampede, congestion and difficult evacuation. Specifically, the quantification of high-density crowds can be defined by “crowd density” [9], that is, the number of people in a unit area, often expressed as the number of people per square meter (p/m^2). According to the crowd density value, the stability state of the crowd can be divided into five different levels [10], as shown in Table 1: Very Low (VL), Low (L), Medium (M), High (H) and Very High (VH). Among them, VL and L are stable states, M is a critical stable state, and H and VH are unstable states. When the crowd is in a critical stable state, the safety management department should pay attention to the movement of the crowd. Once the crowd reaches an unstable state, all departments should take emergency safety management measures, such as restricting the flow of people, increasing protective fences and dispatching more on-site safety management personnel.

Table 1. Crowd states.

Crowd States	Density Range (p/m^2)	Crowd Level
Free Flow	< 0.5	Very Low (VL)
Restricted Flow	0.50 - 0.80	Low (L)

Dense Flow	0.81 - 1.26	Medium (M)
Very Dense Flow	1.27 - 2.00	High (H)
Congestion	> 2.00	Very High (VH)

2.2. *Definition and Classification of Abnormal Behaviors*

2.2.1. Definition of Abnormal Behaviors

Abnormal behaviors refer to the behaviors that individuals or groups exhibit that are significantly different from the normal social behavior patterns in crowded scenes. These behaviors not only include direct physical conflicts and dangerous actions, but also involve any atypical behaviors that may lead to public disorder, the spread of panic emotions, and the damage to the safety of individuals or collectives. In a high-density environment, the combined effects of the density, mobility, emotional state of the crowd and environmental factors make the identification and management of abnormal behaviors extremely complex and urgent. The ubiquitous characteristics of abnormal behaviors [11] include: ①Suddenness, abnormal behaviors of pedestrians often occur suddenly without warning; ②Harmfulness, causing harm to the pedestrians themselves and triggering abnormal reactions of the surrounding crowd; ③Abnormal emotions [12], when the crowd faces a certain emergency, panic emotions arise; ④Spatial abnormality: Irregular flow and aggregation, which may cause safety accidents such as stampede.

2.2.2. Classification of Abnormal Behaviors

The abnormal behaviors of the crowd can be divided into two major categories: violent and non-violent [13]. Violent abnormal behaviors usually lead to more serious casualties or property losses. Among them, common violent abnormal behaviors include group fighting, stampede, terrorist attacks, riots, panic and escape, etc. While non-violent abnormal events do not directly cause violent injuries, they may still bring a series of impacts, such as parades and onlookers.

3. **Challenges Brought by Crowd Density**

In high-density crowd scenes, crowd behavior recognition faces a series of challenges, which stem from the high complexity brought by the dense distribution of the crowd. Specifically, the frequent occlusion between individuals and the motion blur induced by rapid movement together constitute various obstacles in the recognition and analysis process. These factors are intertwined, greatly increasing the difficulty of accurately detecting, continuously tracking and deeply understanding crowd behaviors, and posing a severe test for realizing crowd density distribution, pedestrian trajectories, and abnormal behavior recognition [14].

3.1. *Target Occlusion Problem*

In a crowded crowd, the occlusion situation between individuals is complex and changeable, including not only partial occlusion, but also complete occlusion, and even the alternating occurrence of continuous occlusion and instantaneous occlusion [15], especially in public places such as railway stations, subways and popular squares. Occlusion causes the key features of the target to be hidden or changed, making it difficult for feature-based detection and recognition algorithms to extract complete and reliable visual information. For behavior recognition models, occlusion may cause the established target feature model to fail, and it is difficult to correctly associate the targets before and after occlusion, thereby causing target loss or incorrect association. In crowded scenes, the mutual occlusion between targets is often accompanied by dense interactions between targets, such as collisions and passing by, and these interactions further increase the complexity of the occlusion problem. The algorithm needs to consider the interaction between targets and potential group behavior patterns [16] while dealing with occlusion.

3.2. *Motion Blur and Postural Diversity*

The rapid change of individual movement in the crowd not only causes the target edge captured in the image to be blurred, affecting the accuracy of feature extraction based on clear contours and details but also may cause misrecognition. The blur effect can confuse the understanding of target types and behaviors in visual algorithms [17]. Especially for those individuals who accelerate instantaneously or suddenly turn, their performance in the image is often trailing and residual images, which further reduces the stability and accuracy of the behavior recognition model. In addition, the joints of the human body are highly flexible, and the combination of different bending degrees and directions of each joint will show completely different postures. In the context of movement, this flexibility makes the human posture extremely susceptible to changes due to external forces and other external conditions [18]. A series of actions, whether it is head tilting, arm waving, switching between standing and walking, or adjusting the walking speed, will affect the judgment of pedestrian abnormal behaviors by the behavior recognition model.

4. Research Status of Crowd Abnormal Behavior Recognition

In recent years, traditional learning methods have accumulated rich research results in the field of crowd abnormal behavior recognition. However, with the increasing complexity of crowd scenes, the recognition performance of traditional methods has certain limitations and it is often difficult to capture the subtle differences and dynamic changes of abnormal behaviors. Against this background, the development of deep learning technology and the method of integrating neural networks for efficient pedestrian abnormal behavior recognition have gradually become research hotspots. Depending on whether neural network elements are included in the model construction, the existing abnormal behavior recognition technologies can be divided into two main classifications [19]: traditional methods-based and deep learning-based methods. The recognition technologies based on traditional methods can be classified into 4 categories: methods based on statistical models [20], based on motion features [21], based on dynamic models [22], and based on clustering discrimination [23]. And the deep learning-based methods roughly include 5 categories: methods based on convolutional neural networks (CNN) [24], based on autoencoders (AE) [25], based on generative adversarial networks (GAN) [26], based on long short-term memory networks (LSTM) [27] and based on the self-attention mechanism (Self-Attention) [28].

4.1. Traditional Methods

This subsection classifies traditional abnormal behavior recognition methods into the following categories according to the detection principle: abnormal behavior recognition methods based on statistical models, based on motion features, based on dynamic models, and based on clustering discrimination. Most traditional abnormal behavior recognition methods first construct a normal behavior model through visual feature extraction and motion pattern analysis and then perform abnormal behavior recognition by calculating the statistics of features and setting thresholds.

4.1.1. Methods Based on Statistical Models

The core of the method based on statistical models lies in the in-depth analysis of the intrinsic distribution law of data by applying statistical theories and methods to accurately identify behaviors that deviate significantly from the normal behavior pattern. It builds statistical probability models to describe the statistical characteristics of normal behavior features and then uses these models to detect anomalies in the data.

The Gaussian Mixture Model (GMM) [29] is based on the Gaussian probability density function. By calculating parameters such as the mean and variance of pixel intensity, a model can be established, which can effectively distinguish the foreground and background, and then identify the crowd behavior pattern. In high-density crowd scenes, GMM can be used to capture and describe different dynamic characteristics of the crowd, such as density, movement speed, trajectory, etc., thereby identifying abnormal behaviors that do not conform to the normal behavior pattern. The parameter estimation of GMM usually adopts the Expectation Maximization (EM) algorithm, which

is an iterative algorithm including the E step (calculating the expected value) and the M step (maximizing the likelihood function), and continuously optimizes to approximate the true distribution of the data. Afig et al. [30] discussed and summarized four enhanced methods based on GMM, including basic GMM, the combination of GMM and Markov Random Field (MRF) (GMM-MRF) [31], Gaussian-Poisson Mixture Model (GPMM) [32] and the combination of GMM and Support Vector Machine (SVM) (GMM-SVM) [33], pointing out that when dealing with complex crowd scenes, a combination of multiple methods can be used to improve the ability to recognize abnormal behaviors.

4.1.2. Methods Based on Motion Features

The methods based on motion features focus on analyzing the motion features of individuals or groups in the video, such as speed, direction, acceleration, and trajectory, to identify abnormal behaviors.

Among them, the motion behavior of the crowd can be effectively captured by the optical flow that changes continuously over time. Optical flow [34] is a vector field that describes the motion information of objects in an image sequence, and it reflects the position change of each pixel point in the scene between adjacent frames. In high-density crowd scenes, the optical flow field can reveal the movement direction and speed distribution of individuals or groups, thereby identifying abnormal movements contrary to the normal behavior pattern, such as reverse movement and fast running.

Traditional motion description techniques are mostly based on the velocity information of optical flow, but the acceleration information often contains more abundant motion details. Especially when describing complex motion patterns, it can provide the information missed by the velocity descriptor and help to better understand the motion pattern. Wang et al. [35] studied an acceleration feature descriptor to improve the accuracy of abnormal behavior detection in videos. The process is shown in **Figure 2** below. First, the optical flow field information of each frame is extracted from the video. The inter-frame acceleration is calculated through the optical flow information of two consecutive frames. The acceleration histogram feature is constructed to form the acceleration descriptor (HAVA). Using the video training set containing only normal behaviors, the normal motion pattern is learned through the energy-optimized Restricted Boltzmann Machine (RBM) model [36]. In addition, Jiang et al. [37] proposed a motion descriptor of stripe flow acceleration (SFA), and the process is shown in **Figure 3** below. Stripe flow is a motion representation method that can effectively capture long-term spatiotemporal changes in crowded scenes. Compared with traditional optical flow, it performs better in terms of accuracy and robustness. By introducing stripe flow acceleration to explicitly model motion information, the spatiotemporal changes in crowded scenes can be accurately represented.

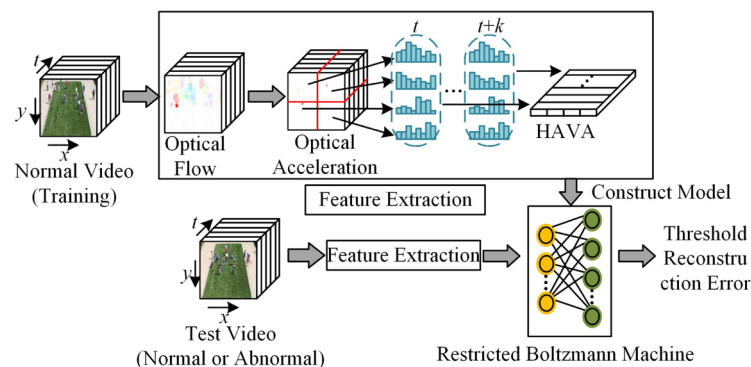


Figure 2. The framework of the acceleration optical flow method.

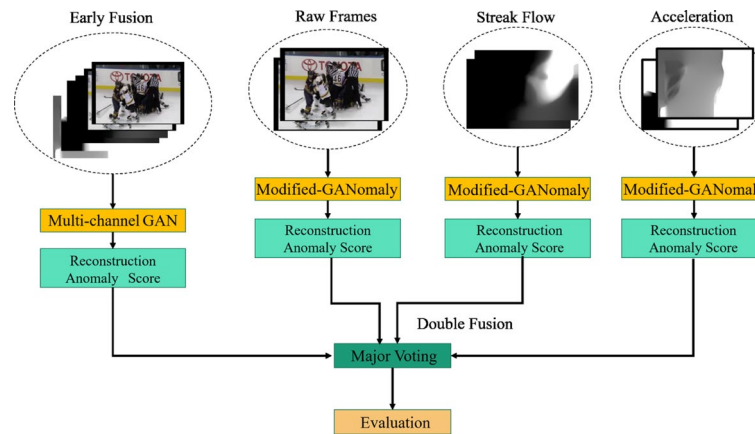


Figure 3. The framework of the SFA method.

4.1.3. Methods Based on Dynamic Models

Methods based on dynamic models, such as cellular automata models [38], particle system models [39], etc., predict the overall movement trend of the crowd by simulating the interaction between individuals and environmental constraints. These models can be used to establish a baseline behavior model for the movement of normal crowds. When the observed behavior significantly deviates from the model prediction, it can be regarded as abnormal. By adjusting model parameters, such as attractive force and repulsive force, it is possible to adapt to the dynamic characteristics of the crowd in different scenarios and improve the accuracy of abnormal behavior recognition.

Chang et al. [40] proposed a hybrid model of cellular automata and agents, dividing the space into multiple discrete units (cells), each cell representing an individual or a small part of the crowd in the crowd, and dynamically updating its state according to preset rules to simulate the dynamic behavior of the crowd in normal or emergencies, such as movement, aggregation, evacuation, etc. Combined with the Agent model to further refine individual behaviors, considering individual differences such as gender, age, physical conditions, and psychological factors such as panic and conformity on evacuation behaviors. This model can identify which individual behaviors deviate from the norm, that is, abnormal behaviors.

4.1.4. Methods Based on Clustering Discrimination

Abnormal behavior recognition technology based on clustering is an unsupervised learning method. The core idea is to group behavior data into several clusters, most of which represent normal behavior patterns, and behaviors deviating from these clusters are regarded as abnormal.

DBSCAN (Density-Based Spatial Clustering of Applications with Noise) is a density-based clustering algorithm that can identify clusters of any shape and can effectively mark noise points. For high-density crowd scenes, select appropriate ϵ (neighborhood radius) and MinPts (minimum number of neighborhood points), where ϵ determines the definition of density and MinPts determines the minimum number of points required to form a dense area. Then, DBSCAN clustering analysis is used to identify dense crowd areas under normal behavior patterns and detect abnormal points that cannot be assigned to any cluster. Chebi et al. [41] proposed a method combining DBSCAN and neural networks for the dynamic detection of abnormal crowd behaviors in video analysis. This method uses the density-based feature of the DBSCAN algorithm to identify the dense areas of the crowd in video surveillance and further analyzes the behavior patterns of these areas through neural networks to identify abnormal behaviors that are different from normal behavior patterns.

In summary, Table 2 summarizes the design ideas, advantages, and disadvantages of the above various abnormal behavior recognition technologies based on traditional methods. Among them, the method based on statistical models is suitable for scenarios that can capture time-series behavior data, especially for the analysis of dynamic behaviors with hidden states, such as crowd flow monitoring.

However, a certain number of training samples are required to evaluate the parameters of the background model. The method based on motion features is suitable for scenarios that require detailed analysis of the movement details of individuals or groups, but this method is susceptible to noise interference, and extracting high-precision motion features will increase the computational cost. The method based on dynamic models is suitable for simulating and predicting collective behaviors of the crowd, such as evacuation scene analysis, but the model complexity is high, so the universality is poor. The method based on clustering discrimination is suitable for unlabeled data sets, does not need to define abnormal standards in advance, and has higher universality, but it is more dependent on the selection of parameters.

Table 2. The classification and characteristics of traditional behavior anomaly recognition methods.

Method Categories	Design Ideas	Advantages	Limitations	References
Statistical Model	Build the background model through statistical methods	Strong pattern recognition ability and capable of handling the complex dynamics of time series data	Complex parameter estimation and sensitivity to initial values	[29–33]
Motion Feature	Identify abnormalities by analyzing the motion characteristics (such as speed, acceleration, etc.) of individuals or groups	Has high sensitivity and specificity for the recognition of motion patterns	Vulnerable to environmental factors interference (such as illumination changes, occlusion)	[34–37]
Dynamic Model	Build a microscopic model of group behavior by simulating the interaction between individuals and environmental constraints	Suitable for simulation and prediction, and model parameters can be adjusted to adapt to different scenarios	High model complexity and unable to fully simulate	[38–40]
Clustering Discrimination	Group the behavior data into different clusters and identify abnormalities based on the differences between clusters	Unsupervised learning is applicable in the case of lacking labeled data	Poor processing ability for time series data	[41]

4.2. Deep Learning-Based Methods

In recent years, deep learning has become the core driving force in the field of computer vision. Compared with traditional methods, deep learning mostly uses the powerful representational learning ability of deep neural networks to automatically extract high-level abstract features from raw data, which can effectively distinguish normal behaviors from abnormal behaviors. According to the network structure and flow processing method, deep learning abnormal behavior recognition methods are classified into the following categories: abnormal behavior recognition methods based on convolutional neural networks, based on autoencoders, based on generative adversarial networks, based on long short-term memory networks and based on self-attention mechanisms.

4.2.1. Methods Based on Convolutional Neural Network

The Convolutional Neural Network (CNN), due to its strong ability to process pixel data and efficient learning feature representation, has become an indispensable part of the field of computer vision, such as image classification, semantic segmentation, and target feature analysis. As a deep,

feed-forward neural network system, the operation process of CNN can be divided into two major steps: feature extraction and classification decision-making, and the structure is shown in **Figure 4**. In the field of crowd abnormal behavior recognition, the current CNN-based methods mainly focus on two core directions: spatial feature and spatiotemporal feature learning.

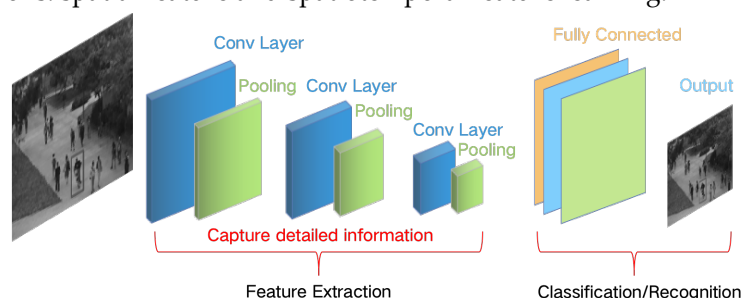


Figure 4. The structure of CNN.

1) Spatial feature learning focuses on extracting the appearance information of pedestrians from video frames to identify atypical behavior patterns in appearance, that is, appearance anomalies. Based on different video forms and annotations, crowd abnormal behavior recognition methods can be divided into three learning methods: supervised, semi-supervised, and unsupervised.

Supervised abnormal behavior recognition is a strategy based on annotated datasets. It extracts pedestrian features through CNN and builds a classification model to identify abnormal patterns in the original data based on normal and abnormal behavior labels [42]. Since explicit category labels are provided during training, CNN can learn precise feature representations for specific tasks, thereby achieving a high accuracy rate for labeled abnormal behavior recognition. However, this method highly depends on the annotated datasets of abnormal behaviors, and the cumbersome manual annotation limits the development of supervised algorithms.

Unsupervised/weakly supervised learning methods can use CNN to extract and analyze pedestrian feature representations without explicit abnormal behavior category labels or with only a small number of labels. This method can detect unknown or unlabeled abnormal behavior types, adapt to a wider range of behavior patterns, and improve the generalization ability of the model in practical applications. Unsupervised/weakly supervised abnormal behavior recognition is implemented based on two stages: feature extraction and abnormal recognition. In the feature extraction stage, pedestrian appearance features in the video sequence are extracted through pre-trained CNNs, commonly including VGG [43], ResNet50 [44], AlexNet [45], GoogLeNet [46], Inception [47], etc. After extracting the appearance features of pedestrians, abnormal classification is performed through abnormal behavior recognition algorithms, commonly including a one-class classifier [48], Gaussian classifier [49], Support Vector Machine (SVM) [50], etc. Singh et al. [51] proposed a method of Aggregation of Ensembles (AOE), which fused and fine-tuned the three networks of AlexNet, GoogLeNet and VGG, especially adjusted the number of output nodes of the fully connected layer to match the binary classification task, and adopted a hierarchical fine-tuning strategy, only updating the learning rate of specific layers and keeping other layers unchanged, thereby avoiding training the network from scratch and improving the efficiency of abnormal behavior recognition. At the same time, each sub-ensemble is composed of a specific classifier and the combination of three CNN models, providing multiple classification outputs for each video frame. To solve the problem of low efficiency of video frame block convolution, Sabokrou et al. [52] proposed a method based on a Fully Convolutional Neural Network (FCN), extracting features of all regions from the entire video frame, and performing convolution and pooling operations in parallel, which can achieve a processing speed of about 370 frames per second, meeting the requirements of real-time processing.

2) Spatiotemporal feature learning methods integrate dynamic information in the time series based on space. This is usually achieved by using optical flow data or upgrading to 3D convolution technology, aiming to capture the complex relationship between motion and time between video frames and identify abnormalities in the movement trajectories and speeds of pedestrians. Under this

framework, according to the different time-perception feature fusion strategies, it is mainly divided into two implementation approaches: two-stream CNN architecture and a three-dimensional convolutional neural network (3D CNN).

The two-stream convolutional neural network decomposes the processing of video data into spatial stream and temporal stream. The spatial stream focuses on the static content of the video frame, that is, the appearance feature in each frame image. The temporal stream focuses on analyzing the temporal changes and dynamic behaviors in the video sequence and usually uses optical flow images as input. Hu et al. [53] proposed a weakly supervised learning ABDL framework. Firstly, the Faster R-CNN network is used to identify objects (such as pedestrians) in the scene, then the large-scale optical flow histogram (HLSOF) is used to describe the behavior features of the objects, and finally, the behaviors are classified through the multi-instance support vector machine (MISVM) to distinguish normal or abnormal. Wang et al. [54] designed a two-stream convolutional neural network model (TS-CNN), as shown in **Figure 5**. Combining the above optical flow features and trajectory information, this network can process spatiotemporal information simultaneously and improve the accuracy of abnormal behavior detection. At the same time, the Kanade-Lucas-Tomasi (KLT) tracking algorithm is used to obtain the single-frame image of the crowd movement trajectory. This strategy enables the algorithm to better handle the occlusion problem in the crowd and further enhances the recognition of abnormal behaviors through the continuity feature of the behavior.

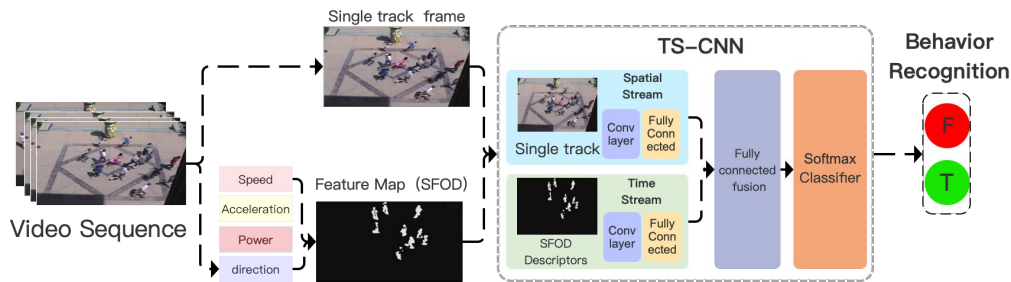


Figure 5. The structure of TS-CNN.

3D CNN adds a time dimension based on the original 2D convolution. By performing convolution on the time axis, 3D CNN can capture the change patterns of behaviors over time, which is extremely crucial for identifying abnormal behaviors (such as falling, running, gathering, etc.) that need to consider the behavior sequence and time context. Hu et al. [55] proposed a method based on a Deep Spatiotemporal Convolutional Neural Network (DSTCNN), extending two-dimensional convolution to three-dimensional space, and comprehensively considering the spatial features of static images and the temporal features between front and rear frames. Firstly, the video screen is divided into multiple sub-regions, and spatiotemporal data samples are extracted from these sub-regions and input into DSTCNN for training and classification, achieving accurate detection and location of abnormal behaviors. At the same time, to solve the problems of network dispersion and insufficient recognition ability, Gong et al. [56] proposed the LDA-Net framework, as shown in **Figure 6**. It consists of a human body detection module and an abnormal detection module. Among them, the YOLO algorithm is introduced in the human body detection module to make the abnormal detection module more concentrated; in the abnormal detection module, 3D CNN is used to simultaneously capture the motion information of labeled normal and abnormal action sequences from the spatial and temporal dimensions.

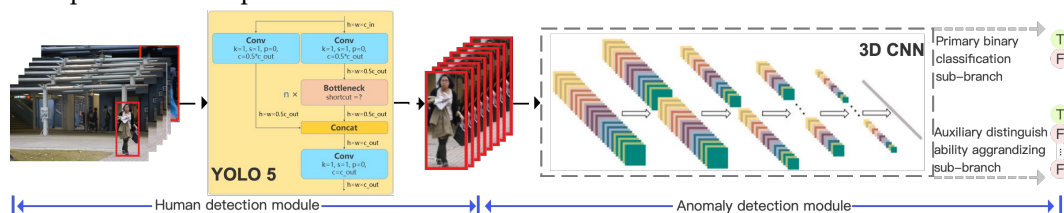


Figure 6. The structure of LDA-Net.

4.2.2. Methods Based on Autoencoders

Autoencoders (AE) [57] is an unsupervised learning method. The basic architecture consists of two parts, the encoder and the decoder. It reconstructs the original input by learning the effective compressed representation of the data. The structure is shown in **Figure 7**. In anomaly detection, AEs are trained to reconstruct normal behavior data. Thus, for those inputs significantly different from the training data (i.e., abnormal behaviors), their reconstruction errors will increase. During this process, the Mean Square Error (MSE) [58] is commonly used as the loss function. MSE is a statistical indicator that measures the difference between the predicted value and the true value, and the formula is expressed as:

$$MSE = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2 \quad (1)$$

Among them, \hat{y}_i is the predicted value of the model for the i th sample, y_i is the true value of this sample, and $\hat{y}_i - y_i$ is the residual (the difference between the predicted value and the true value). The smaller the value of MSE, the smaller the average gap between the predicted value and the true value of the model, and the higher the fitting degree of the model. In recent years, autoencoders have been widely used for crowd abnormal behavior recognition. These methods are mainly divided into two categories [59]: methods based on similarity measurement and methods based on hidden feature representation learning.

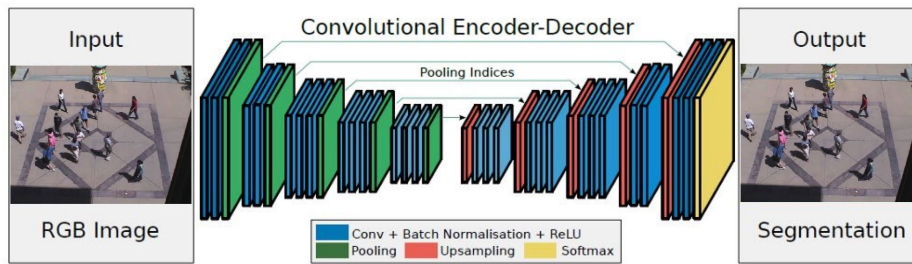
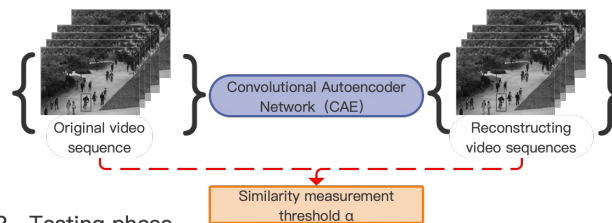


Figure 7. The structure of AutoEncoder.

1) Methods based on similarity measurement are techniques used to evaluate and quantify the degree of similarity between two data objects, samples, sets, or vectors. In the field of crowd abnormal behavior recognition, abnormal behavior recognition is carried out by comparing the similarity between the original input image and the reconstructed image by the Convolutional Autoencoder (CAE), and the process is shown in **Figure 8**. The similarity between the original image and the reconstructed image is compared through similarity measurement techniques (such as Euclidean Distance [60], Pearson Correlation Coefficient [61], etc.). If the similarity is low, it indicates that there is abnormal behavior in the original image.

1. Training phase



2. Testing phase

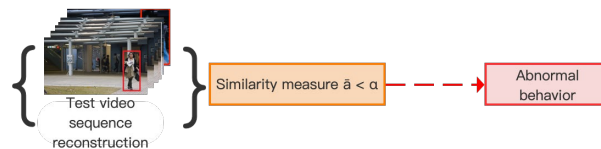


Figure 8. Anomaly behavior recognition process based on similarity measurement.

The original video sequence contains the appearance features of pedestrians, and the use of optical flow analysis can describe the dynamic trajectories of foreground targets. Information fusion can be carried out through the Convolutional Autoencoder (CAE) to extract the appearance attributes and dynamic features of pedestrians. Nguyen et al. [62] designed a CNN combining CAE and U-Net, which share the same encoder. CAE is responsible for learning the normal appearance structure, while U-Net attempts to associate these structures with the corresponding motion templates. This design aims to jointly improve the detection ability of abnormal frames through two complementary streams (one dealing with appearance and the other dealing with motion), and the model supports end-to-end training.

In addition to the fusion of optical flow analysis, some recent studies have utilized the processing ability of the Convolutional LSTM network (ConvLSTM) [63] for time series and combined it with the Convolutional Autoencoder (CAE), significantly improving the extraction ability of spatiotemporal features in videos. Xiao et al. [64] designed the Probabilistic Memory Autoencoder Network (PMAE), integrating 3D causal convolution and time dimension shared fully connected layers in the autoencoder network to extract spatiotemporal features in video frames, while ensuring the temporal sequence of information and avoiding the leakage of future information. Based on the spatiotemporal encoder, Nawaratne et al. [65] proposed the method of Incremental Spatiotemporal Learner (ISTL) to online learn the spatiotemporal normal behavior patterns in video surveillance streams, thereby achieving abnormal detection and location. At the same time, an active learning strategy of fuzzy aggregation is introduced to dynamically adapt to unknown or emerging normal behavior patterns in surveillance videos, allowing the system to continuously refine the understanding of normal and abnormal behaviors based on environmental changes and the acquisition of new information.

2) Methods based on hidden feature representation learning focus on using the encoder to map high-dimensional input data to a low-dimensional vector space, focusing on extracting the core information of normal behavior patterns. Subsequently, the decoder attempts to reconstruct the original data. During the optimization process, it ensures that the model captures the key structure of the behavior data. The process is shown in **Figure 9**. Common autoencoders can be divided into: Stacked Denoising Autoencoders (SDAE) [66], and Variational Autoencoders (VAE) [67].

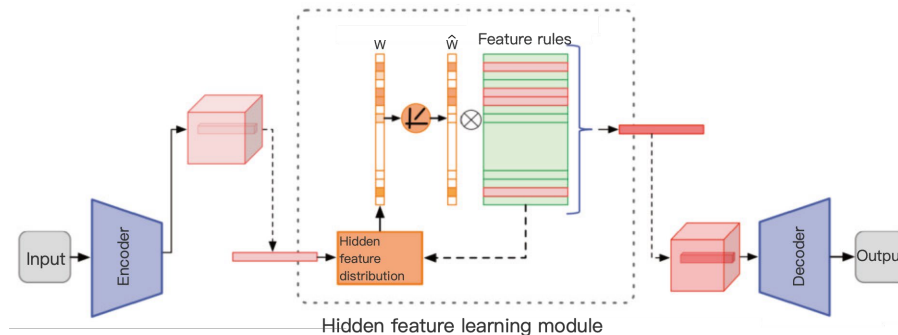


Figure 9. Abnormal behavior recognition process based on hidden feature representation learning.

SDAE extracts higher-level feature representations by combining multiple Denoising Autoencoders (DAE) [68]. The goal of a single DAE is to reconstruct the original data on noisy data. The main principle: ① The encoder maps the noisy input data to a low-dimensional hidden representation; ② The decoder reconstructs the original noise-free data based on this and minimizes the reconstruction error to learn effective information. During the stacking process, the output of the first DAE (i.e., the denoised representation) serves as the input of the second DAE, and so on. Each layer attempts to further remove noise or extract more abstract features from the output of the previous layer. Wang et al. [69] used two SDAEs to learn the appearance features and motion features of behaviors respectively. For the appearance feature, the input is the image block within the spatiotemporal volume around the dense trajectory; the motion feature is based on the optical flow

block. Each SDAE contains multiple encoding layers, and the number of nodes is halved layer by layer until the bottleneck layer is reached. The output of this bottleneck layer is regarded as the learned deep feature. The structure is shown in **Figure 10**.

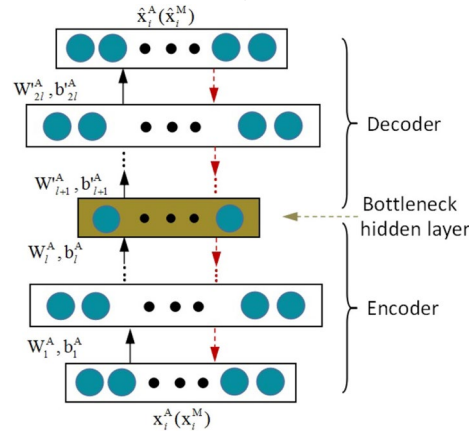


Figure 10. The structure of SDAE.

VAE is a generative model based on a probabilistic framework. Its encoder is similar to the traditional encoder, mapping the input x to a pair of vectors μ and σ , representing the mean and standard deviation of the posterior distribution $q(z|x)$, respectively. The decoder receives the latent variable z and attempts to reconstruct the original input x , that is, estimating $p(x|z)$. The output of the decoder can be regarded as the probability distribution of x given z . Wang et al. [70] designed a new method called S2-VAE. This method combines the stacked fully connected variational autoencoder (SF-VAE) and the skip convolutional variational autoencoder (SC-VAE). Among them, SF-VAE is a shallow generative network designed to simulate Gaussian mixture models to adapt to the actual data distribution. Filter obvious normal samples and improve the speed of abnormal detection. SC-VAE is a deep generative network that integrates low/middle/high-level features and reduces the loss in the information transmission process by fusing the features between the encoder and decoder layers through skip connections.

4.2.3. Methods Based on Generative Adversarial Networks

Generative adversarial networks (GAN) learn the data distribution through the adversarial process of the generator and the discriminator. In anomaly detection, a GAN can be trained to effectively generate only normal behavior data, and then use the reconstruction error of the generator or the output of the discriminator to determine whether the input data is normal. Abnormal behaviors are not within the distribution range of the training data and have poor generation or reconstruction effects, and thus are identified. The GAN model structure is shown in Figure 12. The abnormal behavior recognition strategies based on GAN can be roughly classified into two categories: direct reconstruction or prediction error detection method, and enhanced reconstruction method combined with autoencoder.

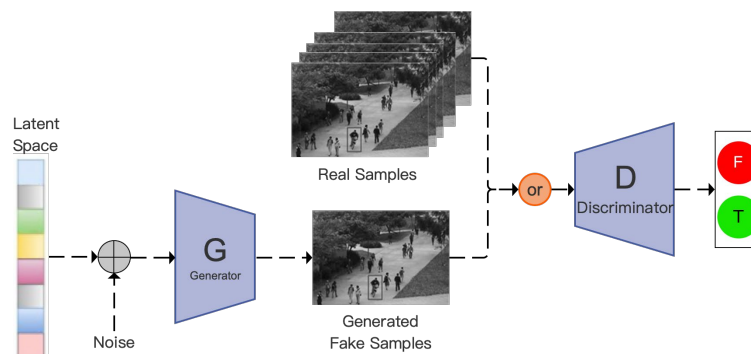


Figure 11. The structure of GAN.

The core idea of the reconstruction-based method is to train a model to learn the distribution representation of normal video data. In the testing stage, it is determined whether the test sample is abnormal based on its reconstruction error. Specifically, first, a deep learning model, that is, a neural network g , is constructed. Its goal is to learn a mapping relationship so that for any given video segment or video frame x , the output $g(x)$ processed by g is as close as possible to the original input x . Secondly, an error function f is designed, which can quantify the difference between x and $g(x)$, that is, the reconstruction error [71] $\varepsilon = f(x, g(x))$. Then, during the training process, efforts are made to find an optimal set of neural network parameters so that for all samples x in the entire dataset, the corresponding sum (or average, weighted sum, etc.) of the reconstruction errors ε reaches the minimum. Song et al. [72] proposed an Ada-Net network architecture that integrates an attention-based autoencoder with a GAN model. This architecture can adaptively learn normal behavior patterns from video data and enhance the reconstruction ability of the autoencoder through adversarial learning, making the reconstructed video frames indistinguishable from the original frames. Different from the traditional reconstruction error metric based on Euclidean distance, the researchers introduced adversarial loss for the frame discriminator to make the reconstructed frame highly similar to the original frame, thereby improving the reconstruction accuracy of the autoencoder. In addition, Chen et al. [73] proposed an anomaly detection model called NM-GAN, which consists of a reconstruction network R and a discrimination network D , forming a GAN-like architecture. NM-GAN builds an end-to-end framework, mainly consisting of three modules: (1) An image-to-image encoder-decoder reconstruction network with appropriate generalization ability; (2) A CNN-based discrimination network for identifying the spatial distribution pattern of the reconstruction error map; (3) An estimation model for quantitatively scoring anomalies. The entire model is trained in an unsupervised manner.

The prediction-based method is based on the fact that a continuous sequence of normal videos has a certain contextual dependence and regularity. Abnormal behaviors are identified by comparing the difference between the observed test frame and its predicted frame. Specifically, Given the consecutive t frames of video x_1, x_2, \dots, x_t , the goal of the prediction model is to generate the next frame \hat{x}_{t+1} and strive to make \hat{x}_{t+1} as close as possible to the actual next frame x_{t+1} . In the testing stage, it is determined whether the current video frame is abnormal by comparing the difference (prediction error) between the \hat{x}_{t+1} predicted by the model and the actual x_{t+1} . Specifically, let h represent the prediction model, then $\hat{x}_{t+1} = h(x_1, x_2, \dots, x_t)$. Prediction frameworks include unidirectional prediction and bidirectional prediction: Unidirectional prediction is usually based on predicting the current frame from the previous few frames. For example, Tang et al. [74] proposed the first method that combines future frame prediction and reconstruction, which is implemented through an end-to-end network. The network consists of two continuously connected U-Net modules. The first module performs future frame prediction based on the input image sequence, and the second module reconstructs the future frame based on the predicted intermediate frame, making normal and abnormal behaviors easier to distinguish in the feature space. Lee et al. [75] used the Bidirectional Multi-scale Aggregation Network (BMAN) for abnormal behavior recognition. By using bidirectional multi-scale aggregation and attention-based feature encoding, normal patterns including object scale changes and complex motions are learned. Based on the learned normal patterns, abnormal events are detected by simultaneously analyzing the appearance characteristics and motion characteristics of the scene through an appearance-motion joint detector.

The enhanced reconstruction method combining autoencoder embeds the adversarial training logic of GAN in the training framework of AE. The autoencoder combined with GAN can not only minimize the reconstruction error but also be linked with the discriminator in GAN to output a data distribution closer to the real one. The process is shown in **Figure 12**. A discriminator is integrated into the architecture of the autoencoder to compare and distinguish the reconstructed image of the autoencoder network and the original input image. The quality of the reconstructed image of the autoencoder network is enhanced through adversarial training. Schlegl et al. [76] proposed a method

of fast anomaly generative adversarial network (f-AnoGAN), training a Wasserstein GAN through normal samples, and then training an encoder to map the image to the latent space to achieve fast inference and anomaly detection. During the anomaly detection process, the input image is reconstructed through the encoder and generator, and the combined score of the image reconstruction residual and the discriminator feature residual is used as a reliable indicator of anomaly.

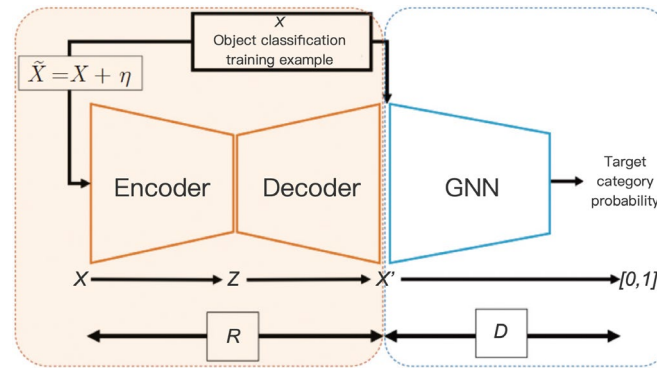


Figure 12. Enhanced Reconstruction Process Combined with Autoencoders.

4.2.4. Methods Based on Long Short-Term Memory Network

Compared with feedforward neural networks, Recurrent Neural Network (RNN) shows significant advantages in processing sequence data, and is particularly good at understanding and processing time series information. Long Short-Term Memory (LSTM) network, as an improvement and extension of RNN, is mainly designed to solve the problems of gradient vanishing and gradient explosion encountered in the training process of long sequence data. LSTM captures and analyzes the changing patterns of behaviors over time, and automatically identifies abnormal sequences that do not conform to it by learning the patterns of normal behavior sequences. LSTM introduces a gate mechanism for controlling the flow and loss of features, including three gating mechanisms: input gate, forget gate, and output gate. The structure is shown in **Figure 13**.

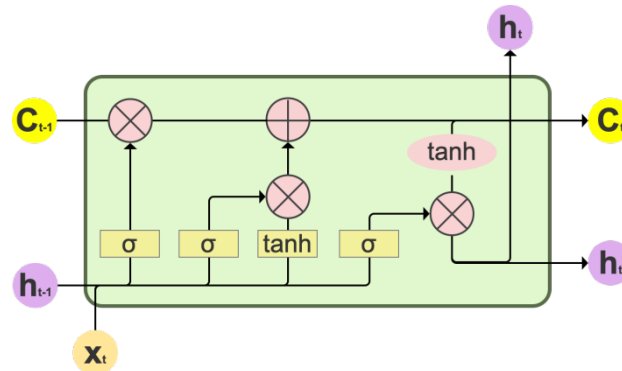


Figure 13. The structure of LSTM.

As shown in the figure above, x_t represents the input at the current moment, h_{t-1} and h_t represent the output of the previous unit and the current unit of LSTM, C_{t-1} , and C_t represent the state of the previous unit and the current unit, σ represents the sigmoid activation function, and \tanh represents the hyperbolic tangent activation function.

1) The forget gate determines the information to be retained. Based on h_{t-1} and x_t , the forget gate outputs a number between 0 and 1 for the state C_{t-1} . 0 means elimination and 1 means retention. The mathematical expression is:

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f) \quad (2)$$

2) The input gate determines the newly incoming information. Composed of the sigmoid function and the tanh function, the result of multiplying the output values is used to update the state. The mathematical expression is as follows:

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i) \quad (3)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C) \quad (4)$$

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t \quad (5)$$

3) The output gate determines the output information. The output state is determined through the sigmoid layer, and the output result of the state through the tanh layer is multiplied by the output result of the sigmoid layer. The mathematical expression is as follows:

$$O_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o) \quad (6)$$

$$H_t = O_t * \tanh(C_t) \quad (7)$$

In the field of crowd abnormal behavior recognition, researchers usually combine LSTM with other technologies to improve the recognition rate and processing speed of crowd abnormal behaviors. Meng et al. [77] integrated the focal loss function into the LSTM algorithm, optimized the model's attention to abnormal samples, and reduced the model loss. This improvement not only enhanced the sensitivity of the LSTM algorithm to abnormal behaviors but also effectively improved the detection accuracy. At the same time, the author also improved the ViBe image foreground extraction method by initializing the background model and randomly selecting the neighborhood sample set to accurately extract the foreground target, thereby simplifying the scene complexity. Wu et al. [78] combined the FCN with strong spatial feature extraction ability and the LSTM with excellent time series modeling ability. The structure is shown in **Figure 14**. FCN is responsible for extracting high-level semantic features from video frames, while LSTM is used to capture the changing patterns of these features over time. By simultaneously learning the spatial and temporal information in the video, it effectively distinguishes normal and abnormal behaviors in the video. Sabih et al. [79] proposed an improved end-to-end supervised learning method that combines Long Short-Term Memory Network (LSTM) and Convolutional Neural Network (CNN), especially in processing optical flow features. CNN is used to extract the frame-level optical flow features calculated based on the Lucas-Kanade algorithm from the video, while the bidirectional LSTM captures the time series information of these features to jointly perform crowd abnormal detection. This method helps to understand the normal behavior patterns in the video and identify abnormal instances.

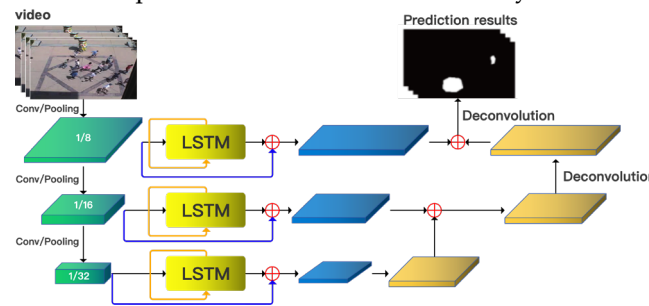


Figure 14. The structure of FCN-LSTM.

4.2.5. Methods Based on Self-Attention Mechanism

In the field of deep learning, the attention mechanism (Self-Attention SA) [80] can be interpreted as a transformation process involving the mapping between a query vector and a series of key-value pairs, thereby generating an output vector. This output vector is obtained by performing a weighted sum of each Value vector, where the weight attached to each Value vector is determined based on the quantitative evaluation of the correlation between its corresponding Key vector and the main query vector. The calculation of the attention mechanism can be divided into three stages, as shown in **Figure 15**.

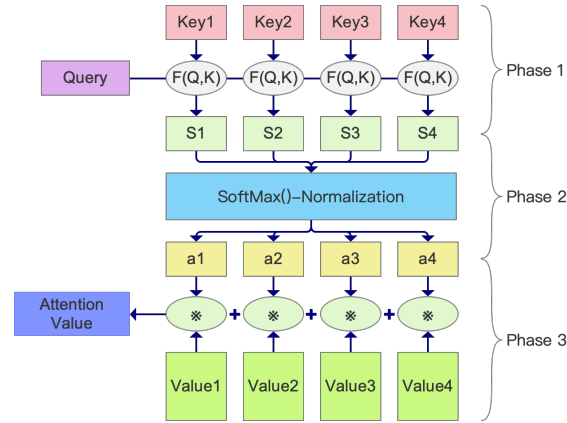


Figure 15. The process of attention mechanism computation.

Stage 1 is to calculate the similarity between Query and Key. Common methods include the vector dot product, Cosine similarity, and MLP network. The expressions are as follows:

Vector dot product

$$Sim(Query, key_i) = Query \cdot Key_i \quad (8)$$

1) Cosine similarity

$$Sim(Query, key_i) = \frac{Query \cdot Key_i}{||Query|| \cdot ||Key_i||} \quad (9)$$

2) MLP network

$$Sim(Query, key_i) = MLP(Query \cdot Key_i) \quad (10)$$

Stage 2 introduces SoftMax for normalization to sort out the probability distribution of all elements. The expression is as follows:

$$a_i = SoftMax(Sim_i) = \frac{e^{Sim_i}}{\sum_{j=1}^N e^{Sim_j}} \quad (11)$$

Stage 3 performs weighted summation to obtain the Attention value. The expression is as follows:

$$Attention(Q, K, V) = SoftMax(\frac{QK^T}{\sqrt{d_k}})V \quad (12)$$

Self-attention is a special form of the attention mechanism. Its uniqueness lies in that when calculating the attention weight, the query, key, and value all come from different parts of the same input sequence. It was first widely used in the “Transformer” model, which greatly promoted the development of the field of natural language processing, and is currently also widely used in crowd sequence modeling and abnormal behavior recognition tasks. Self-attention can not only capture the long-distance dependency within the sequence but also remain efficient in parallel computing, thereby overcoming some limitations of the Recurrent Neural Network (RNN) model.

Ye et al. [81] designed a Self-Attention Feature Aggregation (SAFA) module. The structure is shown in **Figure 16**. This module regenerates the feature map by aggregating the embedding representations with similar information according to the attention map. At the same time, the prior bias on normal data is minimized through the self-suppression strategy, and it is used in combination with the pixel-level frame prediction error to jointly detect abnormal frames, enhancing the ability to recognize abnormalities.

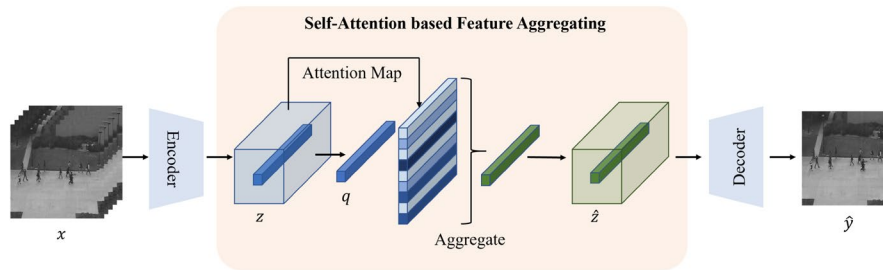


Figure 16. The structure of SAFA.

Zhang et al. [82] proposed an autoencoder model that integrates the global self-attention mechanism to capture the interaction information between the overall features of the image, facilitating the model to understand the context in which the behavior occurs and improving the accuracy of abnormal behavior determination. At the same time, a memory module is embedded in the bottleneck layer of the autoencoder to constrain the model's over-generalization of normal behaviors.

Zhang et al. [83] fused the attention mechanism and the bidirectional Long Short-Term Memory Autoencoder Network (SABiAE). The structure is shown in **Figure 17**. The encoder with the self-attention mechanism captures global appearance features, and the self-attention bidirectional LSTM network is used to reduce the loss of target features, thereby extracting the information between video frames globally. The model reconstructs the frame through the decoding process and determines abnormal behaviors based on the reconstruction error.

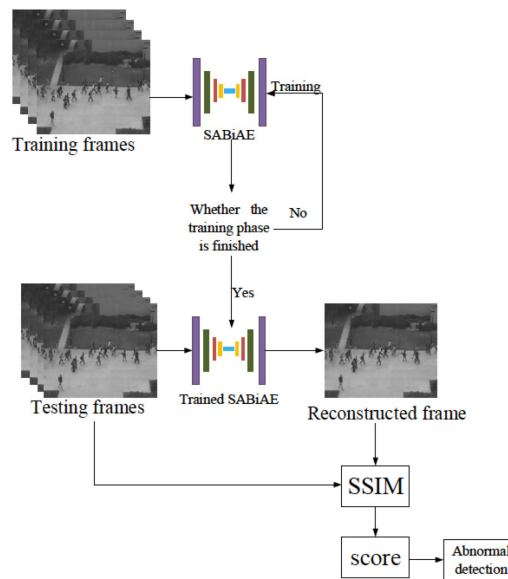


Figure 17. The structure of SABiAE.

Zhang et al. [84] introduced the self-attention mechanism in the Generative Adversarial Network (GAN), using the autoencoder containing the dense residual network and the self-attention mechanism as the generator. At the same time, a new discriminator is proposed, which integrates the self-attention module based on the relative discriminator, thereby avoiding the problem of gradient vanishing during the training process and ensuring that the frames generated by the generator are closer to the real frames.

Singh et al. [85] proposed an Attention-guided Generator and Dual Discriminator Adversarial Network (A2D-GAN), the structure is shown in Figure 19, for real-time video anomaly detection. A2D-GAN utilizes the encoder-decoder architecture, in which the encoder adds multi-level self-

attention to focus on key regions and capture context information, while the decoder uses channel attention to emphasize important features and identify abnormal patterns specific to each frame.

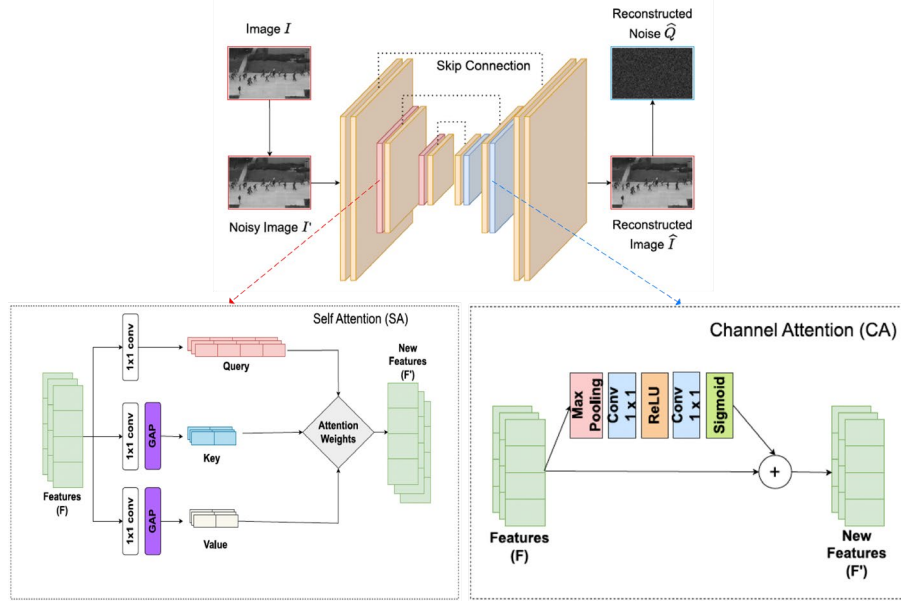


Figure 18. The structure of A2D-GAN.

In addition, to solve problems such as scale changes and occlusions in the crowd, fusion models such as multi-scale feature attention fusion networks and visual occlusion resolution have emerged for crowd abnormal behavior detection. For example, Sharma et al. [86] introduced a scale-aware attention module in CNN, using the cascade of multiple self-attention branches to enhance the recognition ability for scale changes. This architecture integrates motion information, motion influence maps, and features based on the energy level distribution to achieve frame-level abnormal behavior analysis. Zhao et al. [87] proposed a Dynamic Pedestrian Centroid Model (DCM). The key points of the human skeleton are extracted from the image, the pedestrian joint sub-segments are constructed accordingly, and the dynamic characteristics of pedestrians in the video sequence are described mathematically. Experiments show that DCM can detect the U-turn behavior 277ms earlier on average, and detect the fall behavior 562ms earlier on average to find abnormalities. At the same time, to solve the problem of partial occlusion of pedestrians, this paper also designed an anti-occlusion algorithm, as shown in the schematic in **Figure 19**. Firstly, cluster the key points of the pedestrian skeleton, divide the 21 key points into the head, upper body, and lower body, calculate the coordinates of the occluded part, estimate the effective position combined with the mechanical principle, and apply it in combination with DCM. This method has been verified by the fall behavior detection experiment, and the AUC on the 50 Volu. The U-turn dataset reaches 82.13%, and the AUC on the 50 Volu. Fall-down dataset reaches 86.70%.

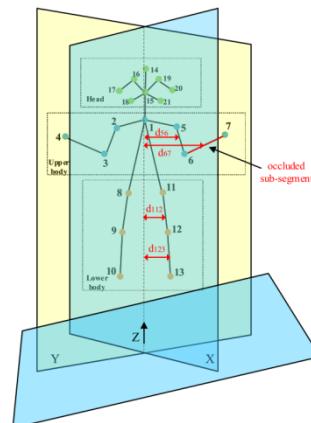


Figure 19. Diagram of the algorithm for occlusion removal based on skeleton points.

In summary, Table 3 summarizes the design ideas, advantages, and disadvantages of the above various deep learning-based abnormal behavior recognition technologies. Overall, these deep learning-based methods benefit from the powerful learning ability of neural networks and can be applied to the recognition of abnormal behaviors of pedestrians in various crowd scenarios. At the same time, a strategy of using multiple methods in combination can be adopted to further improve the generalization ability and recognition accuracy of the model.

Table 3. The classification and characteristics of deep learning behavior anomaly recognition methods.

Method	Design Ideas	Advantages	Limitations	References
CNN	Through local receptive fields and pooling operations, capture the spatial hierarchical features of the original data	Good at extracting local features and combining them into more complex patterns	Lack of time series processing capabilities	[Error! Reference source not found.- Error! Reference source not found.]
AutoEncoder	An unsupervised learning method that captures the main features of the data by encoding and decoding the original data	Reduce the data dimension and extract key features	Easily lead to overfitting of training data	[Error! Reference source not found.- Error! Reference source not found.]
GNN	It is composed of a generator and a discriminator, and the two play a game.	Reconstruct and generate new samples close to real data	Prone to model collapse and non-convergence of training	[Error! Reference source not found.- Error! Reference source not found.]
LSTM	Introduces forget gates, input gates, and output gates to solve the problems of gradient vanishing and gradient explosion	Good at processing sequence data of any length	The structure is complex, and training and reasoning take a long time	[Error! Reference source not found.- Error! Reference source not found.]
Self-Attention	Automatically assigns different attention weights according to different parts of the input sequence	Capture global interdependencies	Easy to overfit on small datasets	[Error! Reference source not found.- Error! Reference source not found.]

4.3. Mainstream Software Tools

In recent years, with the deep integration and breakthrough progress of big data and computer vision technology, the analysis of crowd abnormal behaviors has become an important research topic in the field of public safety. Algorithms and software tools in this field have also emerged continuously. Each representative software tool aims to reveal and identify changes and deviations

from the normal patterns of crowd behaviors in complex scenarios from different dimensions. Many algorithm manufacturers focusing on this field have also launched a series of advanced software tools specifically for crowd abnormal behavior recognition. The following is an overview of the mainstream tools.

The PP-Human crowd abnormal behavior recognition tool is an industrial-level open-source real-time pedestrian analysis tool that integrates core capabilities such as object detection, tracking, and key point detection. It can adapt to different light conditions, complex backgrounds, and cross-lens scenarios, as shown in Figure 21.

The SenseTime crowd abnormal behavior recognition tool analyzes the overall characteristics and individual behavior characteristics of the crowd in the surveillance images, masters the activity rules of the crowd, and realizes the automation or even intelligence of video surveillance.

The Hikvision crowd abnormal behavior recognition tool provides powerful analysis data for crowd control, prevents the occurrence of crowding and stampede incidents, and provides effective guarantees for social security.

The Keda crowd abnormal behavior recognition tool combines the general intelligent recognition ability of the large model and the precise recognition ability of the small model to achieve group scene analysis such as holding knives, fighting, running, gathering, etc.

5. Comparative Analysis of Different Algorithm Experiments

5.1. Experimental Datasets

In the field of computer vision, due to the diversity and complexity of real abnormal behaviors in the research of crowd abnormal behavior analysis, there is a shortage of high-quality datasets. To overcome this challenge and improve the performance of algorithms, researchers have constructed a series of datasets with different characteristics, focusing on parameters such as duration, size, resolution, etc., and covering various monitoring environments and scenarios, providing an important reference basis for the research of crowd anomaly recognition. Commonly used datasets mainly include UCSD, UMN, Shanghai-Tech, CUHK-Avenue, and other related datasets. Figure 25 shows some abnormal behavior samples in these four abnormal behavior datasets.

5.1.1. UCSD

The UCSD dataset [88] is a dataset for crowd behavior analysis and abnormal detection. The dataset contains two subsets, Ped1 and Ped2. The Ped1 subset includes 34 video clips, recording the pedestrian activities from the campus crosswalk scene; the Ped2 subset provides 12 video clips of the same specification, presenting similar but different pedestrian crossing area activities.

5.1.2. UMN

The UMN dataset [89] is mainly used for the evaluation and research of pedestrian detection and tracking algorithms. This dataset consists of 19 videos, including a series of pedestrian video sequences in complex backgrounds, providing various indoor and outdoor environments such as campuses, shopping malls, streets, and other pedestrian activity video clips in different backgrounds.

5.1.3. ShanghaiTech

The ShanghaiTech dataset [90] is a dataset for video abnormal detection and crowd counting. This dataset has 13 complex scenes, recording the pedestrian flow on the campus. It also contains 130 abnormal events (such as fighting, falling, running, etc.) and more than 270,000 training frames, with a duration ranging from about 30 seconds to 90 seconds.

5.1.4. CUHK-Avenue

The CUHK database [91] is a database about crowd behavior scenes. It includes crowd videos collected from many different environments with different densities and perspective ratios, such as

streets, shopping centers, airports, and parks. It consists of traffic datasets and pedestrian datasets. There are a total of 474 video clips in 215 scenes in the dataset.

Table 4 summarizes and compares the above commonly used datasets of crowd abnormal behaviors from the dimensions of scene description, scale, abnormal behavior, video resolution, objects, and limitations of the dataset.

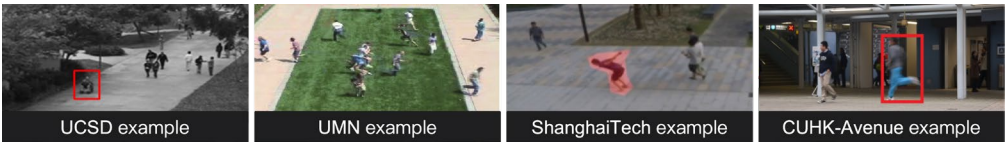


Figure 20. Examples of anomalies in the datasets.

Table 4. Comparison of crowd anomaly behavior datasets.

Name	Scene Description	Scale	Resolution	Abnormal Behavior	Object	Limitation
UCSD [Error! Reference source not found.]	The crowd movement on the sidewalk from the perspective of the surveillance camera	Ped1 14000frame 34 Training segments 36 Test segments	238*158 360*240	Fast moving, Reverse driving Riding a bicycle, Driving a car, Sitting in a wheelchair, skateboarding, etc.	individual	Low resolution; few types of abnormal behaviors; relatively simple background
		Ped2 4560 frame 12 Test segments				
UMN [Error! Reference source not found.]	Video clips of pedestrian activities in different backgrounds such as campuses, shopping malls, and streets	8010 frame 11 Video segment	320*240	The crowd suddenly scattered, ran, and gathered.	group	Simple background; Limited abnormal types
Shanghai-Tech [Error! Reference source not found.]	13 campus area scenes with complex lighting conditions and camera angles	317398 frame 130 Video segment	856*480	Crowd gathering, fighting, running, cycling.	group	Abnormal events are repetitive; The annotations have errors
CUHK-Avenue [Error! Reference source not found.]	Video surveillance clips of outdoor public places	30625 frame 16 Training segments 21 Test segments	640*360	Pedestrians fighting, throwing objects, running	individual Vehicle	The shooting angle is single; The resolution is low

5.2. Experimental Evaluation Indicators

The evaluation of crowd abnormal behavior recognition integrates multiple indicators to ensure that the recognition is not only accurate but also precisely located. During the evaluation, frame-level and pixel-level standards are usually considered. Frame-level detection focuses on determining

whether there is any abnormality within the video frame. Even if the specific location of the abnormality is not precise, as long as there is abnormal behavior in the frame, the entire frame is marked as abnormal; the pixel-level standard further requires precise positioning of the spatial location of the abnormal behavior, usually requiring that the predicted abnormal pixels cover at least 40% of the real abnormal area to evaluate the accuracy of abnormal positioning.

The performance of crowd abnormal behavior recognition is usually evaluated by indicators such as confusion matrix, ROC curve and AUC, and Equal Error Rate (EER).

Confusion matrix

Comprehensively considering the four basic indicators of True Positive (TP), True Negative (TN), False Positive (FP), and False Negative (FN), it is used for basic classification performance evaluation.

ROC curve and AUC

The ROC curve depicts the classification performance of the algorithm through the changes of True Positive Rate (TPR) and False Positive Rate (FPR). AUC, that is, the area under the ROC curve, reflects the average performance of the classifier under all thresholds. The closer the AUC value is to 1, the better the performance of the classifier.

Equal Error Rate (EER)

It is defined as the error classification rate when the True Positive Rate is equal to the False Positive Rate. The smaller the EER value, the better the performance of the abnormal detection method, because it balances missed alarms and false alarms and is an important indicator to measure the performance of a binary classification system.

5.3. Experimental Comparative Analysis

Table 5 summarizes the performance of abnormal behavior recognition methods based on deep learning on UCSD, UMN, ShanghaiTech, and Avenue datasets in recent years. This paper compares and shows through the evaluation indicators of AUC and EER. These evaluation data are all quoted from research papers in recent years.

Table 5. Comparison of abnormal behavior recognition experiments in crowd.

Classification	Method	Frame level AUC / EER / %										Year
		UCSDPed1		UCSDPed2		UMN		Shanghai Tech		Avenue		
		EER	AUC	EER	AUC	EER	AUC	EER	AUC	EER	AUC	
CNN	FCN [Error! Reference source not found.]	--	--	11.00	--	--	--	--	--	--	--	2019
	ABDL [Error! Reference source not found.]	22.00	--	16.00	--	5.80	98.90	--	--	21.00	84.50	2020
	TS-CNN [Error! Reference source not found.]	--	--	--	--	--	99.60	--	--	--	--	2020
	DSTCNN [Error! Reference source not found.]	--	99.74	--	99.94	--	--	--	--	----	--	2020
	LDA-Net [Error! Reference source not found.]	--	--	5.63	97.87	--	--	--	--	--	--	2020
	CAE-UNet [Error! Reference source not found.]	--	--	--	96.20	--	--	--	--	--	86.90	2019
	PMAE [Error! Reference source not found.]	--	--	--	95.90	--	--	--	72.90	--	--	2023
AE	ISTL [Error! Reference source not found.]	29.80	75.20	8.90	91.10	--	--	--	--	29.20	76.8	2019

GAN	S2-VAE [Error! Reference source not found.]	14.30	94.25	--	--	--	99.81	--	--	--	87.6	2019
	Ada-Net [Error! Reference source not found.]	11.90	90.50	11.50	90.70	--	--	--	--	17.60	89.20	2019
	NM-GAN [Error! Reference source not found.]	15.00	90.70	6.00	96.30	--	--	17.00	85.30	15.30	88.60	2021
	D-UNet [Error! Reference source not found.]		84.70		96.30	--	--	--	73.00		85.10	2019
	BMAN [Error! Reference source not found.]	--	--	--	96.60	--	99.60		76.20		90.00	2019
	FocalLoss-LSTM [Error! Reference source not found.]	--	--	--	--	--	99.83	--	--	--	--	2021
LSTM	FCN-LSTM [Error! Reference source not found.]	--	--	--	98.20	--	93.70	--	--	--	--	2021
	CNN-LSTM [Error! Reference source not found.]	--	94.83		96.50	--	--	--	--	--	--	2022
	SAFA [Error! Reference source not found.]	--	--	--	96.80	--	--	--	--	--	87.30	2023
SA	SA-AE [Error! Reference source not found.]	--	--	--	95.69	--	--	--	--	--	84.10	2023
	SABiAE [Error! Reference source not found.]	--	--	9.80	95.60	--	--	--	--	20.90	84.70	2022
	SA-GAN [Error! Reference source not found.]	--	--	--	--	--	--	--	75.70	--	89.20	2021
	A2D-GAN [Error! Reference source not found.]	9.70	94.10	5.10	97.40	--	--	25.20	74.20	9.00	91.00	2024
	SA-CNN [Error! Reference source not found.]	--	--	--	--	--	99.29	--	--	--	--	2023

Through the comparison of the experimental results, the following conclusions can be drawn:

(1) Although most methods have achieved high recognition rates on some datasets, on datasets such as UCSDPed1, the EER and AUC values of many methods indicate that there are certain false alarms or missed alarms, meaning that there is still room for further optimization of abnormal behavior recognition technology, especially when dealing with specific types of data or complex scenarios.

(2) Two-dimensional convolutional networks that only rely on spatial features often perform worse in abnormal behavior recognition tasks than those models that integrate optical flow features, 3D convolution, and LSTM. These models incorporate temporal information and can more comprehensively capture the dynamic changes in video sequences, not only considering spatial features but also in-depth analysis of the continuity and change trends of behavior patterns in the time dimension.

(3) The hidden feature representation learning method based on the autoencoder network is usually superior to the similarity measurement method in terms of automatic feature learning, adaptability, and generalization ability. It is suitable for processing large-scale, high-dimensional, and complex abnormal behavior recognition tasks, especially in the unsupervised learning paradigm, it can effectively learn the low-dimensional representation of high-dimensional data.

(4) The introduction of the self-attention mechanism significantly improves the accuracy and efficiency of abnormal behavior recognition. This mechanism enhances the model’s ability to capture long-distance dependencies in the sequence. By allowing the model to dynamically adjust the degree of attention to different parts of the input, key features can be extracted more effectively, promoting the abnormal behavior recognition technology to a new height.

6. Summary and Research Prospects

This paper conducts comprehensive and in-depth research and analysis of crowd abnormal behavior recognition technology from four dimensions: basic definitions, traditional methods, deep learning, and application scenarios, as shown in **Figure 21**.

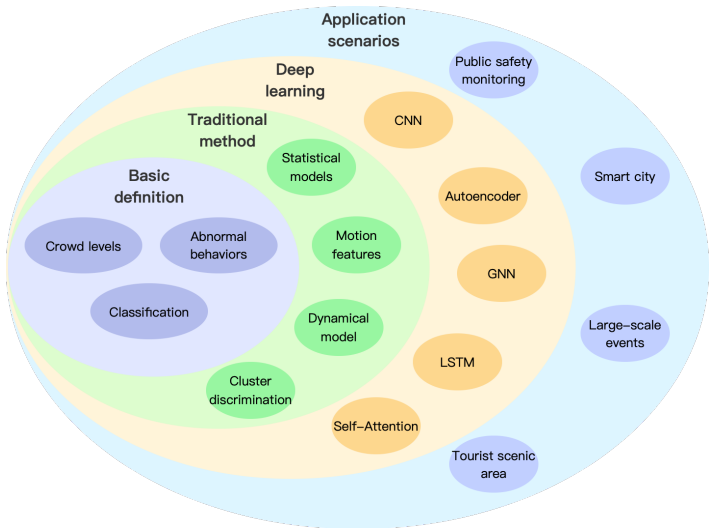


Figure 21. Crowd abnormal behavior recognition hierarchical structure diagram.

(1) Abnormal definition

This paper defines the crowd level and abnormal behavior in public places. By quantifying the “crowd density”, the stability state of the crowd is divided into five levels: very low (VL), low (L), medium (M), high (H), and very high (VH). Among them, the critical stable state (M) and unstable state (H/VH) require special attention and corresponding measures from the safety management department. And the crowd’s abnormal behavior can be further classified into two categories: violent and non-violent.

(2) Traditional methods

This paper deeply analyzes the abnormal behavior recognition technology based on traditional methods. Statistical models identify abnormalities by constructing a probabilistic framework; motion feature methods focus on dynamic indicators such as optical flow and speed for abnormal analysis; dynamic models predict the overall behavior trend by simulating the interaction between individuals; clustering discrimination technology automatically groups behavior data through unsupervised learning and marks abnormal patterns. However, in the face of the challenges of high-dimensional, nonlinear, and complex dynamic scenarios, the limitations of traditional methods are increasingly prominent.

(3) Deep learning

With its powerful feature learning ability, deep learning has gradually become a research hotspot in the field of abnormal behavior recognition. This paper elaborates in detail on 5 methods based on CNN, GAN, LSTM, AE, and SA. These methods have shown excellent performance in the field of abnormal behavior recognition, especially for complex scenarios and diverse abnormal behavior patterns. By using the representational learning ability of deep neural networks, high-level features are automatically extracted from the original data to effectively distinguish normal and abnormal behaviors.

(4) Application scenarios

The research field of this paper plays an important role in public safety monitoring, large-scale event security, and tourist attraction management, providing strong technical support for ensuring public safety and improving the efficiency of urban management. This technology accurately identifies various abnormal behaviors such as violent conflicts and stampede risks in dense crowd scenes through real-time monitoring of video streams, realizes early warning and rapid response to high-risk events, and helps relevant departments make scientific decisions and rationally dispatch resources to effectively prevent and handle emergencies.

Despite the progress made, this field still faces problems such as robustness in complex environments, generalization ability, real-time performance, and dependence on labeled data. Future research will focus on:

(1) Enhance the fusion ability of deep learning and multimodal

Future research will further strengthen the fusion ability of deep learning methods in processing multi-source information (such as video, audio, environmental data, etc.), and construct more sophisticated crowd behavior models, thereby improving the accuracy and reliability of abnormal behavior detection in different scenarios.

(2) Enhance context understanding and situational reasoning

The algorithm will develop in a more intelligent direction, not only being able to recognize the human body movement itself, but also being able to deeply analyze the intrinsic connection between individual behavior and its environment, social norms, and interpersonal relationships, and improve the intelligent level of abnormal detection by introducing situational awareness and causal reasoning mechanisms.

(3) Enhance the robustness and adaptive learning mechanism of the abnormal detection model

To enhance the robustness and generalization ability of the abnormal detection model, future research can develop more robust models, and use technologies such as adversarial learning and meta-learning to enhance the anti-interference and transfer learning ability of the model. At the same time, design an adaptive learning mechanism so that the model can dynamically adjust parameters and strategies according to changes in the environment and tasks.

(4) Expand multi-source and multi-dimensional data fusion recognition

Future research will integrate and deeply analyze multivariate data such as view data, crowd movement data, and UWB location information, and use advanced fusion algorithms to process time series features, spatial attributes, and even deeper psychophysiological indicators (such as brain waves, heart rate variability, etc.) captured by advanced technologies such as brain-computer interfaces. Through this cross-modal and cross-level multi-source and multi-dimensional data fusion technology, a more comprehensive and three-dimensional crowd behavior model will be constructed to further improve the comprehensiveness and accuracy of the abnormal behavior detection system.

(5) Construct a self-consistent metaverse virtual-real fusion evolution model

Generative AI and adaptive networks are adopted to study the metaverse self-consistent scene model to ensure that the virtual environment can be dynamically adjusted according to the behavior and feedback of pedestrians to maintain the consistency and rationality of the virtual environment. At the same time, the "perception-prediction-intervention-interconstruction" virtual-real fusion presentation mechanism is studied by adopting a time-delay system adaptive feedforward control network, which can ensure that the system can perceive the crowd behavior in real-time, predict potential violent conflicts or stampede risks in dense crowds, and take measures quickly to avoid the occurrence of such events.

Acknowledgements: This research is partially supported by the National Natural Science Foundation of China (No. 72374154). The authors deeply appreciate the supports.

References

1. Haghani M, Lovreglio R. Data-based tools can prevent crowd crushes[J]. *Science*, 2022, 378(6624): 1060-1061.

2. Luo L, Xie S, Yin H, et al. Detecting and Quantifying Crowd-level Abnormal Behaviors in Crowd Events[J]. IEEE Transactions on Information Forensics and Security, 2024.
3. Fadhel M A, Duham A M, Saihood A, et al. Comprehensive systematic review of information fusion methods in smart cities and urban environments[J]. Information Fusion, 2024: 102317.
4. Zhang Bingbing, Ge Shuyu, Wang Qilong, Li Peihua. Research on behavior recognition method based on multi-order information fusion. Acta Automatica Sinica, 2021, 47(3): 609-619 doi: 10.16383/j.aas.c180265
5. Jiang Jun, Zhang Zhuojun, Gao Mingliang, et al. A crowd abnormal behavior detection algorithm based on pulse line flow convolutional neural network[J]. Journal of Engineering Science and Technology, 2020, 52(6):215-222.
6. Xiao Jinsheng, Shen Mengyao, Jiang Mingjun, Lei Junfeng, Bao Zhenyu. Abnormal behavior detection in surveillance videos by fusing bag attention mechanism. Acta Automatica Sinica, 2022, 48(12): 2951-2959 doi: 10.16383/j.aas.c190805
7. Alam E, Sufian A, Dutta P, et al. Vision-based human fall detection systems using deep learning: A review[J]. Computers in biology and medicine, 2022, 146: 105626.
8. Yang Fan. Xiao Bin, Yu Zhiwen. Review of anomaly detection and modeling in surveillance videos[J]. Journal of Computer Research and Development, 2021(012):058.
9. Li Y C, Jia R S, Hu Y X, et al. A Weakly-Supervised Crowd Density Estimation Method Based on Two-Stage Linear Feature Calibration[J]. IEEE/CAA Journal of Automatica Sinica, 2024, 11(4): 965-981.
10. Zhao R, Dong D, Wang Y, et al. Image-based crowd stability analysis using improved multi-column convolutional neural network[J]. IEEE Transactions on Intelligent Transportation Systems, 2021, 23(6): 5480-5489.
11. Zhao Rongyong, Wei Bingyu, Zhu Wenjie, et al. Review of pedestrian abnormal behavior recognition methods in public places[J]. China Safety Science Journal, 2024, 34(2):83-93. DOI:10.16265/j.cnki.issn1003-3033.2024.02.1125.
12. Wang Duorui, Du Yang, Dong Lanfang, Hu Weiming, Li Bing. Small sample learning algorithm based on feature transformation and metric network. Acta Automatica Sinica, 2024, 50(7): 1305-1314 doi:
13. Ng S H, Platow M J. The violent turn in non-violent collective action: What happens?[J]. Asian Journal of Social Psychology, 2024.
14. Chen Chong, Bai Shuo, Huang Lida, et al. Research on passenger flow monitoring and early warning in crowded places based on video analysis [J]. Journal of Safety Science and Technology, 2020, 16 (04): 143-148.
15. Wang X, Li Y. Edge Detection and Simulation Analysis of Multimedia Images Based on Intelligent Monitoring Robot[J]. Informatica, 2024, 48(5).
16. Chen X, Yan B, Zhu J, et al. High-performance transformer tracking[J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2022, 45(7): 8507-8523.
17. Farooq M U, Saad M N B M, Malik A S, et al. Motion estimation of high-density crowd using fluid dynamics[J]. The Imaging Science Journal, 2020, 68(3): 141-155.
18. Zhou Zhiqiang, Sun Riming, Guo Chenglong, et al. Hierarchical motion estimation method for spatially unstable targets under Gaussian mixture model[J/OL]. Acta Optica Sinica:1-19[2024-05-26].<http://kns.cnki.net/kcms/detail/31.1252.O4.20240517.1540.029.html>.
19. Zhou Y, Liu C, Ding Y, et al. Crowd descriptors and interpretable gathering understanding[J]. IEEE Transactions on Multimedia, 2024.
20. Alhothali A, Balabid A, Alharthi R, et al. Anomalous event detection and localization in dense crowd scenes[J]. Multimedia Tools and Applications, 2023, 82(10): 15673-15694.
21. Nayan N, Sahu S S, Kumar S. Detecting anomalous crowd behavior using correlation analysis of optical flow[J]. Signal, Image and Video Processing, 2019, 13: 1233-1241.
22. Aziz Z, Bhatti N, Mahmood H, et al. Video anomaly detection and localization based on appearance and motion models[J]. Multimedia Tools and Applications, 2021, 80(17): 25875-25895.
23. Matkovic F, Ivasic-Kos M, Ribaric S. A new approach to dominant motion pattern recognition at the macroscopic crowd level[J]. Engineering applications of artificial intelligence, 2022, 116: 105387.
24. Ganga B, Lata B T, Venugopal K R. Object detection and crowd analysis using deep learning techniques: Comprehensive review and future directions[J]. Neurocomputing, 2024: 127932.
25. Madan N, Ristea N C, Ionescu R T, et al. Self-supervised masked convolutional transformer block for anomaly detection[J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2023.
26. Chai L, Liu Y, Liu W, et al. CrowdGAN: Identity-free interactive crowd video generation and beyond[J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2020, 44(6): 2856-2871.
27. Song X, Chen K, Li X, et al. Pedestrian trajectory prediction based on deep convolutional LSTM network[J]. IEEE Transactions on Intelligent Transportation Systems, 2020, 22(6): 3285-3302.
28. Huang Y, Abeliuk A, Morstatter F, et al. Anchor attention for hybrid crowd forecasts aggregation[J]. arXiv preprint arXiv:2003.12447, 2020.

29. He F, Xiang Y, Zhao X, et al. Informative scene decomposition for crowd analysis, comparison and simulation guidance[J]. *ACM Transactions on Graphics (TOG)*, 2020, 39(4): 50: 1-50: 13.
30. Afiq A A, Zakariya M A, Saad M N, et al. A review on classifying abnormal behavior in crowd scene[J]. *Journal of Visual Communication and Image Representation*, 2019, 58: 285-303.
31. Goel S, Koundal D, Nijhawan R. Learning Models in Crowd Analysis: A Review[J]. *Archives of Computational Methods in Engineering*, 2024: 1-19.
32. Rajasekaran G, Raja Sekar J. Abnormal Crowd Behavior Detection Using Optimized Pyramidal Lucas-Kanade Technique[J]. *Intelligent Automation & Soft Computing*, 2023, 35(2).
33. Gündüz M Ş, Işık G. A new YOLO-based method for real-time crowd detection from video and performance analysis of YOLO models[J]. *Journal of Real-Time Image Processing*, 2023, 20(1): 5.
34. Miao Y, Yang J, Alzahrani B, et al. Abnormal behavior learning based on edge computing toward a crowd monitoring system[J]. *IEEE Network*, 2022, 36(3): 90-96.
35. Wang Kunlun, Liu Wencan, He Xiaohai, et al. A motion feature descriptor for abnormal behavior detection[J]. *Computer Science*, 2020, 47(4):119-124. DOI:10.11896/jsjcx.190300392.
36. Qin Yue, Shi Yuexiang. Human behavior recognition based on two-stream network fusion and spatio-temporal convolution [J]. *Computing Technology and Automation*, 2021, 40 (02): 140-147. DOI:10.16339/j.cnki.jsjzyzh.202102027.
37. Jiang J, Wang X Y, Gao M, et al. Abnormal behavior detection using streak flow acceleration[J]. *Applied Intelligence*, 2022: 1-18.
38. Clarke K C. Cellular automata and agent-based models[M]//Handbook of regional science. Berlin, Heidelberg: Springer Berlin Heidelberg, 2021: 1751-1766.
39. Li Changhua, Bi Chenggong, Li Zhijie. Crowd evacuation model based on improved PSO algorithm [J]. *Journal of System Simulation*, 2020, 32 (06): 1000-1008. DOI:10.16182/j.issn1004731x.joss.18-0782.
40. Chang D, Cui L, Huang Z. A cellular-automaton agent-hybrid model for emergency evacuation of people in public places[J]. *IEEE Access*, 2020, 8: 79541-79551.
41. Chebi H, Dalila A, Mohamed K. Dynamic detection of abnormalities in video analysis of crowd behavior with DBSCAN and neural networks[J]. *Advances in Science, Technology and Engineering Systems Journal*, 2020, 1(5): 56-63.
42. Tay N C, Connie T, Ong T S, et al. A robust abnormal behavior detection method using convolutional neural network[C]//Computational Science and Technology: 5th ICCST 2018, Kota Kinabalu, Malaysia, 29-30 August 2018. Springer Singapore, 2019: 37-47.
43. Koonce B, Koonce B. Vgg network[J]. *Convolutional Neural Networks with Swift for Tensorflow: Image Recognition and Dataset Categorization*, 2021: 35-50.
44. Koonce B, Koonce B. ResNet 50[J]. *Convolutional neural networks with swift for tensorflow: image recognition and dataset categorization*, 2021: 63-72.
45. Behera N K S, Sa P K, Muhammad K, et al. Large-Scale Person Re-Identification for Crowd Monitoring in Emergency[J]. *IEEE Transactions on Automation Science and Engineering*, 2023.
46. Wang S, Pu Z, Li Q, et al. Estimating crowd density with edge intelligence based on lightweight convolutional neural networks[J]. *Expert Systems with Applications*, 2022, 206: 117823.
47. Dai F, Huang P, Mo Q, et al. ST-InNet: Deep spatio-temporal inception networks for traffic flow prediction in smart cities[J]. *IEEE Transactions on Intelligent Transportation Systems*, 2022, 23(10): 19782-19794.
48. Liu W, Luo W, Li Z, et al. Margin Learning Embedded Prediction for Video Anomaly Detection with A Few Anomalies[C]//IJCAI. 2019, 3: 023-3.
49. Bhuiyan M R, Abdullah J, Hashim N, et al. Hajj pilgrimage abnormal crowd movement monitoring using optical flow and FCNN[J]. *Journal of Big Data*, 2023, 10(1): 86.
50. Garg S, Sharma S, Dhariwal S, et al. Human crowd behaviour analysis based on video segmentation and classification using expectation-maximization with deep learning architectures[J]. *Multimedia Tools and Applications*, 2024: 1-23.
51. Singh K, Rajora S, Vishwakarma D K, et al. Crowd anomaly detection using aggregation of ensembles of fine-tuned convnets[J]. *Neurocomputing*, 2020, 371: 188-198.
52. Sabokrou M, Fayyaz M, Fathy M, et al. Deep-anomaly: Fully convolutional neural network for fast anomaly detection in crowded scenes[J]. *Computer Vision and Image Understanding*, 2019, 172: 88-97.
53. Hu X, Dai J, Huang Y, et al. A weakly supervised framework for abnormal behavior detection and localization in crowded scenes[J]. *Neurocomputing*, 2020, 383: 270-281.
54. Wang Hongyan, Zhou Mengxing. Crowd abnormal behavior detection based on optical flow and trajectory [J]. *Journal of Jilin University (Engineering Edition)*, 2020, 50 (06): 2229-2237. DOI:10.13229/j.cnki.jdxbgxb.20190665.
55. Hu Xuemin, Chen Qin, Yang Li, et al. Crowd abnormal behavior detection and location based on deep spatio-temporal convolutional neural network[J]. *Application Research of Computers*, 2020, 37(3):891-895. DOI:10.19734/j.issn.1001-3695.2018.09.0671.

56. Gong M, Zeng H, Xie Y, et al. Local distinguishability aggrandizing network for human anomaly detection[J]. *Neural Networks*, 2020, 122: 364-373.
57. Pinaya W H L, Vieira S, Garcia-Dias R, et al. Autoencoders[M]//Machine learning. Academic Press, 2020: 193-208.
58. Tyagi B, Nigam S, Singh R. A review of deep learning techniques for crowd behavior analysis[J]. *Archives of Computational Methods in Engineering*, 2022, 29(7): 5427-5455.
59. Xu Tao, Tian Chongyang. Review of abnormal behavior detection of crowd based on deep learning[J]. *Computer Science*, 2021, 48(9):125-134. DOI:10.11896/jsjx.201100015.
60. Zhang X, Ma D, Yu H, et al. Scene perception guided crowd anomaly detection[J]. *Neurocomputing*, 2020, 414: 291-302.
61. Zhou S, Shi R, Wang L. Extracting macroscopic quantities in crowd behaviour with deep learning[J]. *Physica Scripta*, 2024, 99(6): 065213.
62. Nguyen T N, Meunier J. Anomaly detection in video sequence with appearance-motion correspondence[C]//Proceedings of the IEEE/CVF international conference on computer vision. 2019: 1273-1283.
63. Wei H, Li K, Li H, et al. Detecting video anomaly with a stacked convolutional LSTM framework[C]//International Conference on Computer Vision Systems. Cham: Springer International Publishing, 2019: 330-342.
64. Xiao Jinsheng, Guo Haowen, Xie Honggang, et al. Probabilistic memory autoencoder network for abnormal behavior detection in surveillance videos[J]. *Journal of Software*, 2023, 34(9): 4362-4377.
65. Nawaratne R, Alahakoon D, De Silva D, et al. Spatiotemporal anomaly detection using deep learning for real-time video surveillance[J]. *IEEE Transactions on Industrial Informatics*, 2019, 16(1): 393-402.
66. Roka S, Diwakar M. Deep stacked denoising autoencoder for unsupervised anomaly detection in video surveillance[J]. *Journal of Electronic Imaging*, 2023, 32(3): 033015-033015.
67. Li J, Huang Q, Du Y, et al. Variational abnormal behavior detection with motion consistency[J]. *IEEE Transactions on Image Processing*, 2021, 31: 275-286.
68. Ashfahani A, Pratama M, Lughofer E, et al. DEV DAN: Deep evolving denoising autoencoder[J]. *Neurocomputing*, 2020, 390: 297-314.
69. Wang J, Xia L. Abnormal behavior detection in videos using deep learning[J]. *Cluster Computing*, 2019, 22(Suppl 4): 9229-9239.
70. Wang T, Qiao M, Lin Z, et al. Generative neural networks for anomaly detection in crowded scenes[J]. *IEEE Transactions on Information Forensics and Security*, 2019, 14(5): 1390-1399.
71. Lü Chengkan, Shen Fei, Zhang Zhengtao, et al. Review of the Research Status of Image Anomaly Detection[J]. *Acta Automatica Sinica*, 2022, 48(6):1402-1428. DOI:10.16383/j.aas.c200956.
72. Song H, Sun C, Wu X, et al. Learning normal patterns via adversarial attention-based autoencoder for abnormal event detection in videos[J]. *IEEE Transactions on Multimedia*, 2019, 22(8): 2138-2148.
73. Chen D, Yue L, Chang X, et al. NM-GAN: Noise-modulated generative adversarial network for video anomaly detection[J]. *Pattern Recognition*, 2021, 116: 107969.
74. Tang Y, Zhao L, Zhang S, et al. Integrating prediction and reconstruction for anomaly detection[J]. *Pattern Recognition Letters*, 2020, 129: 123-130.
75. Lee S, Kim H G, Ro Y M. BMAN: Bidirectional multi-scale aggregation networks for abnormal event detection[J]. *IEEE Transactions on Image Processing*, 2019, 29: 2395-2408.
76. Schlegl T, Seeböck P, Waldstein S M, et al. f-AnoGAN: Fast unsupervised anomaly detection with generative adversarial networks[J]. *Medical image analysis*, 2019, 54: 30-44.
77. Meng B, Wang L, Li D W. Detection Method for Crowd Abnormal Behavior Based on Long Short-Term Memory Network[C]//Advances in Intelligent Information Hiding and Multimedia Signal Processing: Proceeding of the 16th International Conference on IIHMSP in conjunction with the 13th international conference on FITAT, November 5-7, 2020, Ho Chi Minh City, Vietnam, Volume 1. Springer Singapore, 2021: 305-313.
78. Wu Guangli, Guo Zhenzhou, Li Leiting, et al. Video Anomaly Event Detection Fusing FCN and LSTM[J]. *Journal of Shanghai Jiao Tong University*, 2021, 55(5): 607.
79. Sabih M, Vishwakarma D K. Crowd anomaly detection with LSTMs using optical features and domain knowledge for improved inferring[J]. *The Visual Computer*, 2022, 38(5): 1719-1730.
80. Vaswani A, Shazeer N, Parmar N, et al. Attention is all you need[J]. *Advances in Neural Information Processing Systems*, 2017, 30.
81. Ye Z, Li Y, Cui Z, et al. Unsupervised Video Anomaly Detection with Self-Attention Based Feature Aggregating[C]//2023 IEEE 26th International Conference on Intelligent Transportation Systems (ITSC). IEEE, 2023: 3551-3556.
82. Zhang Hongmin, Fang Xiaobing, Zhuang Xu. Video Human Abnormal Behavior Detection Model of Autoencoder Fusing Attention Mechanism[J]. *Laser Journal*, 2023, 44 (02): 69-75. DOI:10.14016/j.cnki.jgzz.2023.02.069.

83. Zhang J, Qi X, Ji G. Self Attention Based Bi-Directional Long Short-Term Memory Auto Encoder for Video Anomaly Detection[C]//2021 Ninth International Conference on Advanced Cloud and Big Data (CBD). IEEE, 2022: 107-112.
84. Zhang W, Wang G, Huang M, et al. Generative Adversarial Networks for Abnormal Event Detection in Videos Based on Self-Attention Mechanism[J]. IEEE Access, 2021, 9: 124847-124860.
85. Singh R, Sethi A, Saini K, et al. Attention-Guided Generator with Dual Discriminator GAN for Real-Time Video Anomaly Detection[J]. Engineering Applications of Artificial Intelligence, 2024, 131: 107830.
86. Sharma V K, Mir R N, Singh C. Scale-Aware CNN for Crowd Density Estimation and Crowd Behavior Analysis[J]. Computers and Electrical Engineering, 2023, 106: 108569.
87. Zhao R, Wang Y, Jia P, et al. Abnormal Behavior Detection Based on Dynamic Pedestrian Centroid Model: Case Study on U-Turn and Fall-Down[J]. IEEE Transactions on Intelligent Transportation Systems 2023.
88. UCSD Anomaly Detection Dataset, <http://www.svcl.ucsd.edu/projects/anomaly/dataset.htm>.
89. UMN Crowd Dataset, <http://mha.cs.umn.edu/projevents.shtml#crowd>.
90. Cao C, Lu Y, Zhang Y. Context Recovery and Knowledge Retrieval: A Novel Two-Stream Framework for Video Anomaly Detection[J]. IEEE Transactions on Image Processing, 2024.
91. Wang L, Wang X, Li M, et al. Anomaly Detection Method Based on Temporal Spatial Information Enhancement[J]. Measurement Science and Technology, 2023, 35(3): 035410.

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.