

# Video Violence Recognition and Localization Using a Semi-Supervised Hard Attention Model

Hamid Mohammadi<sup>a</sup> (hamid.mohammadi@aut.ac.ir), Ehsan Nazerfard<sup>a</sup>  
(nazerfard@aut.ac.ir)

<sup>a</sup> Department of Computer Engineering and Information Technology, Amirkabir  
University of Technology, Valiasr Sq., 350 Hafez Ave. Tehran, Iran

**Corresponding Author:**

Ehsan Nazerfard

Department of Computer Engineering and Information Technology, Amirkabir  
University of Technology, Valiasr Sq., 350 Hafez Ave. Tehran, Iran

Tel: (+98) 6454-2707

Email: nazerfard@aut.ac.ir

# Video Violence Recognition and Localization Using a Semi-Supervised Hard Attention Model

Hamid Mohammadi<sup>a</sup>, Ehsan Nazerfard<sup>a,\*</sup>

<sup>a</sup>*Department of Computer Engineering and Information Technology, Amirkabir University of Technology, Valiasr Sq., 350 Hafez Ave. Tehran, Iran*

---

## Abstract

The significant growth of surveillance camera networks necessitates scalable AI solutions to efficiently analyze the large amount of video data produced by these networks. As a typical analysis performed on surveillance footage, video violence detection has recently received considerable attention. The majority of research has focused on improving existing methods using supervised methods, with little, if any, attention to the semi-supervised learning approaches. In this study, a reinforcement learning model is introduced that can outperform existing models through a semi-supervised approach. The main novelty of the proposed method lies in the introduction of a semi-supervised hard attention mechanism. Using hard attention, the essential regions of videos are identified and separated from the non-informative parts of the data. A model's accuracy is improved by removing redundant data and focusing on useful visual information in a higher resolution. Implementing hard attention mechanisms using semi-supervised reinforcement learning algorithms eliminates the need for attention annotations in video violence datasets, thus making them readily applicable. The proposed model utilizes a pre-trained I3D backbone to accelerate and stabilize the training process. The proposed model achieved state-of-the-art accuracy of 90.4% and 98.7% on RWF and Hockey datasets, respectively.

*Keywords:* deep reinforcement learning, violence detection, hard attention,

---

\*Corresponding author.

*Email addresses:* hamid.mohammadi@aut.ac.ir (Hamid Mohammadi), nazerfard@aut.ac.ir (Ehsan Nazerfard )

## 1. Introduction

A comprehensive and dedicated violence monitoring system is becoming increasingly necessary as social unrest, social violence, and homicide cases increase. A security system’s most significant weakness has always been the reliability of its human agents (Gill & Spriggs, 2005; Bugeja et al., 2018). The number of surveillance cameras installed worldwide is rapidly exceeding the data-load capacity of legacy human-operated surveillance and monitoring systems(Gill & Spriggs, 2005; Bugeja et al., 2018). Violence monitoring systems based on artificial intelligence have significant potential as a reliable and scalable alternative to human-based surveillance systems. As opposed to human-based surveillance systems, artificial intelligence and machine learning surveillance systems offer predictable downtime, reliability, and effortless scaling(Nguyen et al., 2020; Shidik et al., 2019; Sung & Park, 2021).

The quality of deep learning models depends on the quality of their data. It is difficult to obtain surveillance footage due to its security and privacy concerns. Almost all privately collected datasets in this field are covered by non-disclosure agreements, which prevent the public from accessing them, including artificial intelligence researchers. In light of this issue, it is valuable for the community to collect and publish a surveillance-based human violence dataset. The significant advantage of the RWF dataset (Cheng et al., 2021) over other video violence datasets (Hockey fights (Nievas et al., 2011) and Movie violence (Gong et al., 2008)) is collection of a relatively large set of real-world surveillance violence videos.

In video classification, temporal knowledge is essential (Yao et al., 2016; Murthy & Goecke, 2014). A single video frame provides little insight into human violence (Feichtenhofer et al., 2016). Having temporospatial understanding requires recurrent or 3D architectures in deep neural networks (Feichtenhofer et al., 2016; Yao et al., 2016; Algamdi et al., 2019; Du et al., 2017). In action and

violence recognition research, recurrent networks, 3D convolutional layers, and multi-stream architectures are common themes(Zong et al., 2021; Wang et al., 2021; Shi et al., 2020). However, recurrent neural networks, such as LSTMs or GRUs, cannot grasp the complete dependency between consecutive frames (Zeyer et al., 2019; Ezen-Can, 2020; Wang et al., 2019). While transformers have state-of-the-art accuracy in video classification and action recognition, they lack the computational agility and performance required for cost-effective large-scale video surveillance platforms (Arnab et al., 2021; Girdhar et al., 2019). The 3D convolutional layers can be used to capture the temporal information in video using 3D convolutional filters with parameters sharing capabilities (Feichtenhofer et al., 2016; Yao et al., 2016; Algamdi et al., 2019; Du et al., 2017). As a result of their effectiveness and efficiency, these layers are ideal for this study.

Video violence recognition accuracy is enhanced by the use of auxiliary features (Tu et al., 2018b). Aside from raw RGB video frames, additional features such as RGB-differences (Wang et al., 2017), optical flow (Sevilla-Lara et al., 2018), pose estimation (Luvizon et al., 2018; Pham et al., 2020), and deep-learning-based features (Li et al., 2016; Khan et al., 2018; Xiao et al., 2019) are utilized to capture different aspects of the input videos. When the network itself cannot comprehend temporally spatial patterns, motion estimation features, such as RGB-difference and optical flow characteristics, assist in capturing the temporal interdependence of video frames (Wang et al., 2018a). Additionally, pose estimation, and deep-learning-based features are intended to convey contextual information about each frame. Pose estimation features represent the position and movement of the human body (Luvizon et al., 2018; Pham et al., 2020). Deep features (extracted from a computer vision backbone network trained on tasks, such as image classification or object detection) represent the environment, background, and objects in a latent vector space suitable for action understanding and classification (Li et al., 2016; Khan et al., 2018; Xiao et al., 2019).

As valuable as these features are, their extraction incurs an additional computational cost detrimental to real-world applications. The extraction of ac-

curate optical flow requires specialized hardware or significant computational resources (equal to or exceeding the costs of the primary neural network in some cases) (Sun et al., 2018). As a result, it is advantageous to identify more cost-effective methods of enhancing video violence identification models without adding additional features.

The proposed method for improving video violence recognition relies on the premise that hard attention enhances model accuracy. Surveillance videos contain considerable information redundancy as humans are the only subjects in a pipeline for violence recognition. Therefore, the background information can be eliminated without compromising vital details(Sharma et al., 2015). By using hard attention, the video violence recognition model eliminates redundant information from the input frames. The reduction of redundancy and increased focus on the useful information in a video increases the model’s accuracy by reducing its search space and avoiding overfitting.

Reinforcement learning is an effective method to implement hard attention (Rao et al., 2017; Driessens, 2019; Shen et al., 2018; Mott et al., 2019). In reinforcement learning methods, partial signals (rewards) are used to optimize a global criterion (Sutton & Barto, 2018). A reinforcement learning implementation of the hard attention method can utilize the video level annotation information to learn the functionality of a region proposal model. Reinforcement learning methods can select a region of interest by expressing the region proposal information in action space. Furthermore, other partial information acquired from the video could provide additional learning signals for the reinforcement learning method (Sutton & Barto, 2018). In order to improve the precision of the reinforcement learning implementation of hard attention, information such as the region of motion, objects present in the environment, skeleton models, and similar data can be included in the reward shaping process.

The novelty of this study lies in the use of semi-supervised reinforcement learning techniques to create a hard attention mechanism to improve the purity of visual data in order to achieve state-of-the-art accuracy. Unlike the mainstream state-of-the-art approaches, this method takes a reductive approach to

model improvement by focusing on removing less helpful information. It is possible to apply the semi-supervised hard attention approach to the existing video violence datasets without making any modifications. The region of interest is learned based on the annotations at the video level.

Accordingly, the remainder of the paper is organized as follows: Section 2 reviews the literature on video violence detection and reinforcement learning. Section 3 introduces the proposed model, followed by the evaluation results in Section 4. In Section 5, the pros and cons of the proposed method are thoroughly discussed. Lastly, Section 6 discusses conclusions and future research.

## 2. Related Work

The background literature can be divided into two different perspectives. In the first perspective, research is being conducted to find supervised models to classify videos based on the presence of specific actions. The second perspective involves reinforcement learning approaches to improve computer vision tasks (such as classification, detection, or tracking). The research presented here combines the perspectives mentioned above.

### 2.1. Video action classification

In addition to the established image classification methods, video as three-dimensional data poses additional challenges. The addition of the third dimension requires the use of specialized features and representation learning techniques (Hara et al., 2017; Wang et al., 2018b; Nazir et al., 2018). The temporal dimension can be captured using techniques such as recurrent layers (Liu et al., 2016; Ullah et al., 2017; Majd & Safabakhsh, 2020; Liu et al., 2017), 3D convolutional layers (Ji et al., 2012; Yang et al., 2019; Zhou et al., 2018; Hara et al., 2017), and, more recently, transformers (Plizzari et al., 2021; Li et al., 2021; Mazzia et al., 2021; Girdhar et al., 2019). Recurrent layers have a reputation for inconsistent training and poor temporal learning (Vaswani et al., 2017); however, 3D convolutional layers and transformers are highly effective in

this field. Due to this, most cutting-edge techniques use 3D convolutional networks and transformers to map temporal information to latent features. Such networks provide the backbone for the extraction of features in a video classification model. The backbone is followed by a simple classifier, i.e., a fully connected neural network, to form an end-to-end video classification model.

A temporal feature may be as simple as the difference between consecutive frames or more complex such as optical flow. Using richer forms of temporal data (e.g., optical flow) will result in more accurate models (Sevilla-Lara et al., 2018; Sun et al., 2018). Conversely, the trade-off between accuracy and performance leads to the use of superficial temporal features (such as RGB-difference) in some applications (Zhang et al., 2016; Hu et al., 2018; Crasto et al., 2019; Wang et al., 2018c).

The accuracy of action recognition can be improved by including additional explicit information in addition to RGB frames and motion features. Visual cues (Tu et al., 2018a; Wang et al., 2016a) and skeleton estimations (Yan et al., 2018; Plizzari et al., 2021; Shi et al., 2019) are examples of such data. In spite of the increased accuracy, each additional data type adds two overheads to the performance of the activity recognition system. Firstly, each feature requires additional computation during extraction. As a result of the substantial processing power required for extracting motion vectors from RGB frames, researchers have removed, replaced, or estimated this feature (Sun et al., 2018; Wang & Schmid, 2013). Secondly, each additional data stream makes the neural network larger and more expensive to compute (Simonyan & Zisserman, 2014).

The integration of additional data types beyond the RGB video frames has created a unique architecture for action recognition models. Two-stream and multi-stream neural networks utilize neural networks with two or more backbones that can accept multiple types of data. Following the extraction of features from each backbone, the features are fused and classified. This architecture is common in video action recognition approaches (Simonyan & Zisserman, 2014; Tu et al., 2018b; Shi et al., 2020; Wang et al., 2021).

## 2.2. RL-based attention

Videos contain a great deal of redundant information. In most cases, categorizing a video based on human action does not require understanding what is taking place in the background. A neural network’s accuracy can be improved by removing excessive information and purifying the data (Song et al., 2022). Accordingly, the accuracy of action recognition can be improved through the use of soft or hard attention in network architectures.

Soft attention can be thought of as a weighted version of the original data. Weights for each data region are automatically learned during the neural network training process (Sharma et al., 2015; Liu et al., 2017; Song et al., 2017; Li et al., 2020). A general form of soft attention is represented in Equation 1.

$$c_i = \sum_{j=1}^T \underbrace{\text{softmax}(s_{i-1}.h_j)}_{\text{attention coefficient}} \cdot \underbrace{h_j}_{\text{input index } j} \quad (1)$$

Total number of inputs
hidden state  
attention output
step i-1  
 $T$

A sequence of input vectors is represented by  $h$ . At each step of the soft attention function ( $c_i$ ), the output is the weighted sum of the input matrices. The weight of each input matrix ( $h_j$ ) represents the attention coefficient calculated at each step for each input matrix. At the current index, the coefficient depends on the input matrix at the current index and the hidden state of the attention layer at the previous index ( $s_i$ ).

Hard attention is a binary form of soft attention in which zero-weighted values are completely omitted from the input. As a result, the output size of a hard attention function is not continuous and may differ from its input size. It is possible to implement hard attention using supervised, self-supervised, and semi-supervised techniques. The implementation of hard attention using supervised methods is costly since the localization annotations must be manually entered into the dataset. It is also possible to implement hard attention through self-supervised methods (Manchin et al., 2019). For example, as most actions are considered a type of motion, removing motionless parts of a video is a form of

hard attention (Crasto et al., 2019). Motion attention guarantees the inclusion of most activities in the attention area (except for activities with no motion, such as sleeping, sitting, pointing, and so on); however, it creates a lot of redundant data. Using motion attention methods, it is difficult to distinguish a particular action occurring in a busy street from other movements in the video.

Spatial transformer networks (STNs) are end-to-end solutions to implement hard attention in neural networks (Jaderberg et al., 2015). STNs are neural layers that compute matrix spatial transformations (for example, scaling, rotating, and cropping). These networks implicitly learn the transformation required for each input according to the global accuracy of the larger neural network. With a few minute design changes, convolutional neural networks (CNNs) can be transformed into hard attention CNNs utilizing STNs (Li et al., 2018; Malinowski et al., 2018). However, since their output is simply a modified version of the input image, these networks are constrained by the resolution of the input image. In order to gain the maximum benefit from the high-resolution input, it is necessary to utilize alternative methods.

Classification, localization, and tracking performance are good reward sources for reinforcement learning algorithms. As a result, the main objective is to assign the reinforcement learning model the task of improving the accuracy of the underlying model. Reinforcement learning models improve accuracy by focusing on important information in videos, given their attention ability. This approach eliminates the need for extra information (e.g., annotations) in the learning data. The use of neural networks in reinforcement learning models (or deep reinforcement learning) allows advanced computer vision algorithms to be applied to reinforcement learning. As a result, state-of-the-art computer vision neural architectures are paired with reinforcement learning techniques to improve performance on current tasks. A combination of these approaches is investigated in computer vision tasks, including object localization (Wang et al., 2018d; Jie et al., 2016; Caicedo & Lazebnik, 2015; Mathe et al., 2016), and visual tracking (Luo et al., 2019; Ren et al., 2018; Yun et al., 2018; Zhong et al., 2019; Cui et al., 2021). For a more innovative example, Rao et al. (2017) used

reinforcement learning in order to create temporal attention in the face recognition task. Identifying faces that are most likely to be correctly recognized in a timeline is an ideal task for deep reinforcement learning.

### 3. Method

A reinforcement learning model and a set of strategies for describing the observations and actions taken by an agent to detect violence in videos are presented in this study. Through the addition of hard attention and semi-supervised learning capabilities, the deep reinforcement learning agent improves the established deep video classification models. In order to train the semi-supervised model, the dataset annotations are converted into reward signals for the agent. The deep reinforcement learning model is improved by reward shaping and train stabilization.

#### 3.1. *Semi-Supervised Hard-Attention*

SSHA, short for **Semi-Supervised Hard Attention**, is based on two assumptions: (i) Given a high-quality dataset, a neural network with a larger size performs better than a similar neural network with smaller size. (ii) Removing redundant information from neural network input results in either a smaller model (fewer parameters) with roughly the same accuracy or a model of the same size (same number of parameters) the same size but with higher accuracy.

For the purpose of reducing the computational cost, input images are often shrunk to a smaller size. The information contained in an image is lost as it becomes smaller because details are removed, and key characteristics are represented in a smaller vector space (the number of pixels is reduced). Consequently, it is common in the computer vision community to propose a range of models with different input sizes (thus varying the number of parameters) in order to address the trade-off between accuracy and computational cost. The redundancy of information is also a common problem in computer vision applications. It is critical to note that redundant data increases the dimensionality

of input data. This fills the input data with information that could have been used to demonstrate valuable features. In action recognition tasks, data redundancy is evident since the majority of information is redundant in each video (i.e., environment, background objects, almost everything except the subject of the action).

The proposed method combines the above assumptions with minimal drawbacks. Hence, obtaining high accuracy and performance with a smaller model and fewer redundant inputs. With hard attention, redundant data can be removed from a network’s input, improving accuracy without the need for a larger network. A greater amount of computational power is available to process the valuable information once redundant information is removed.

Compared to using the auxiliary features, the hard-attention methodology is more generalizable and cost-effective. In addition to the computational overhead of auxiliary features, some features, such as skeleton estimation, are selected based on the video violence detection application (importance of human body dynamics). Due to their application-specific nature, such features cannot be used in a broader range of computer vision applications. The low computational overhead and application-neutral assumptions of the hard attention mechanism make it suitable for a wider range of applications in computer vision.

The concept of hard attention in computer vision can be understood as the process of cropping out redundant information from each frame. As a result, the hard attention task can be formulated as a method of determining the coordinates of a crop function. Equation 2 is an interpretation of such an approach.

$$c = f_{crop}(arg \max_j \underbrace{f_{score}(h_j)}_{\substack{\text{attention score function} \\ \text{DRL network}}}, h) \quad (2)$$

The diagram illustrates the components of Equation 2. It shows the following relationships:

- attention output** points to the variable  $c$ .
- crop function** points to the function  $f_{crop}$ .
- attention score function DRL network** points to the function  $f_{score}(h_j)$ .
- input at region j** points to the variable  $h_j$ .

Equation 2 is centered on deep reinforcement learning (DRL) attention scor-

ing model ( $f_{score}$ ). The DRL model is a function that scores regions based on their importance in the detection of violence ( $h_j$ ). The region with the maximum attention score is then cropped using the crop function ( $f_{crop}$ ) and considered to be the output of hard attention. The crop function is a light image processing function that returns a region of an input image as its output.

The use of reinforcement learning to create hard attention has the following main benefits: (i) Removing the need for localization annotations in training data and using the currently available datasets. (ii) Removing the need for a separate attention network through the design of a multi-task network based on reinforcement learning. Annotations regarding action localization are not standard in action recognition and violence detection datasets. Through the conversion of video violence detection tasks to reinforcement learning tasks, the network is able to use a partial signal (classification annotation) within a reward system to simultaneously learn recognition and localization.

### 3.2. Design

The hard attention is implemented as a multi-stage process. The reinforcement learning environment utilizes predefined regions called prior boxes, as shown in Figure 1. Through the use of prior boxes, the model can focus on different regions of a video within one or more region selection stages. In each stage, the chosen area replaces the current frame; as a result, the rest of the inference is carried out using the chosen region. In this manner, the model may continue to tighten the attention region by selecting the most appropriate area at each subsequent iteration. The purpose of the training is to teach the model how to select regions that contain individuals who exhibit violent behavior. Choosing a class (violent or non-violent) for the video concludes the violence recognition task. This process is depicted in Figure 2.

The implemented hard attention mechanism has only one region of interest at a given time. Although this may be a limitation, in the case of video violence detection, multiple regions of interest are not required. To correctly classify a video, it is sufficient to detect one of many acts of violence occurring within a

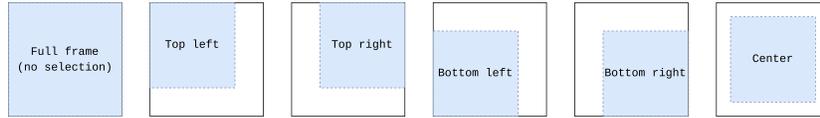


Figure 1: Prior boxes defined on the input frame.

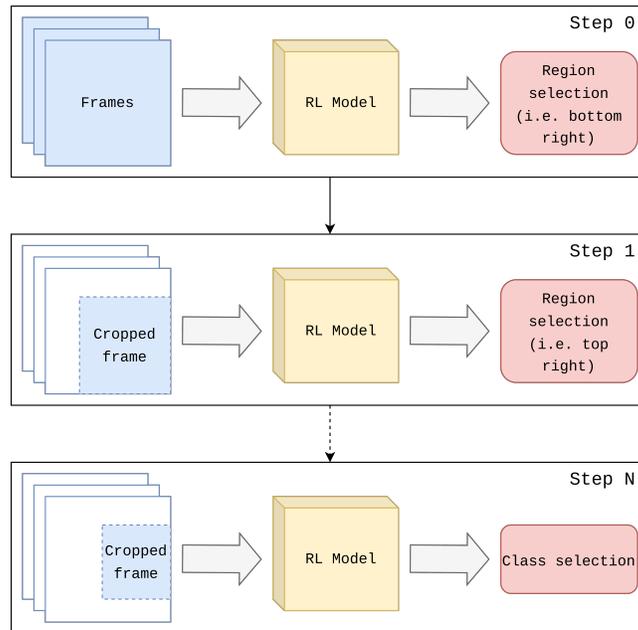


Figure 2: Model interaction with an input video.

single region. By transforming SSHA into a multi-agent reinforcement learning problem, SSHA can be extended to a larger number of attention regions. Attention agents work in collaboration to cover multiple areas of interest in order to maximize global precision.

Selection of the region is accomplished by choosing from a set of prior boxes. Using static regions rather than a free-moving attention bounding box in the training phase reduces the search space for the deep reinforcement learning model. As opposed to the limited number of predefined prior bounding boxes, a free-moving attention bounding box can possess all the possible bounding boxes on the input image. Due to the limited number of sample videos in the video violence recognition datasets, reducing the reinforcement learning search space enhances the convergence and generalization of the learned network weights.

A reward system complements the definition of a search space and an environment. The rewards are designed to encourage correct video classification and discourage incorrect classification. The rewards are defined to be +1 for a correct video class selection and -1 for an incorrect class selection. Furthermore, a diminishing +0.5 reward is associated with attention action in order to encourage the reinforcement learning model to experiment with region selection. As rewards below 0 act as punishments, the auxiliary reward should be greater than 0 to encourage the model. Additionally, since the reward for a correct class selection is +1, the auxiliary reward must not overwhelm the primary reward by being more than +1. According to empirical analysis, values such as 0.3, 0.4, 0.5, and 0.7 result in reasonably similar results, with 0.5 being the best.

RGB and optical flow frames are used as raw inputs to the SSHA model. RGB frames are  $224 \times 224 \times 3$ , and 79 RGB frames are sampled and fed to the model at each step. Optical flow input has similar dimensions, except the frames represent 2D motion vectors using two channels (instead of three RGB channels). The optical flow frames are calculated using the TV-L1 algorithm (Carreira & Zisserman, 2017). TV-L1 is the algorithm of choice for the I3D backbone (Carreira & Zisserman, 2017) for its highly accurate motion vectors. Equation 3

displays TV-L1 visual movement calculations.

$$\vec{v} = \min_{\vec{u}} \left\{ \underbrace{\int_{\Omega} |\nabla x| + |\nabla y| dt}_{\text{regularization term}} + \lambda \underbrace{\int_{\Omega} |\rho(x, y)| dt}_{\text{optical flow constraint}} \right\} \quad (3)$$

By finding the smallest displacement vector ( $\vec{u}$ ) in each region of an image, the motion vector ( $\vec{v}$ ) can be calculated. In the unconstrained form, motion vectors are calculated as the displacement of a pixel within a short period of time ( $\int_{\Omega} |\nabla x| + |\nabla y| dt$ ). Where  $\nabla x$  and  $\nabla y$  are the displacement amounts of the pixel. The TVL1 optical-flow equation is constrained as it considers the weighted ( $\lambda$ ) value difference of the tracked pixel. This difference is calculated as the derivative of the pixel value over a given period of time ( $|\rho(x, y)| dt$ ).

The use of a pretrained model in this research improves the final accuracy, stabilizes deep reinforcement learning training, and accelerates model convergence. I3D model (Carreira & Zisserman, 2017) is chosen as the backbone network. Aside from being one of the best action recognition models on the leaderboards, this model also has excellent source code and pre-trained weights that make it suitable for the application. Kinetics dataset (Smaira et al., 2020) is used to train the model’s pre-trained weights, including RGB and Optical-flow streams.

To simultaneously use RGB and optical flow features, a two-stream architecture (Feichtenhofer et al., 2016) with a multiplication fusion layer is implemented. As shown in Figure 5, the two-stream fusion network uses RGB and Optical-flow I3D backbones. The rest of the network has the same architecture as the single-stream configuration. The 3D feature map generated by the backbones is used to learn Q values. Feature maps are reduced in size using 3D convolution layers, which reduce the number of parameters and prevent over-fitting. The reduced feature is fed to the fully-connected layer with linear activations. Figures 3 to 5 demonstrate an overall view of various SSHA model architectures.

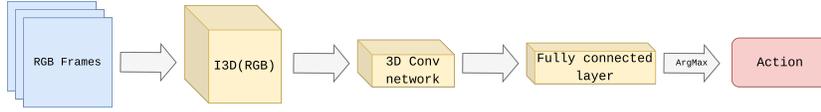


Figure 3: SSHA model architecture (RGB only).

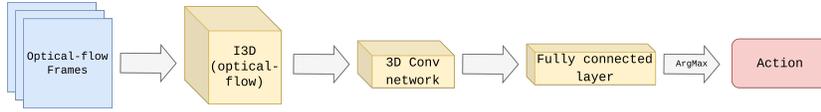


Figure 4: SSHA model architecture (Optical-flow only).

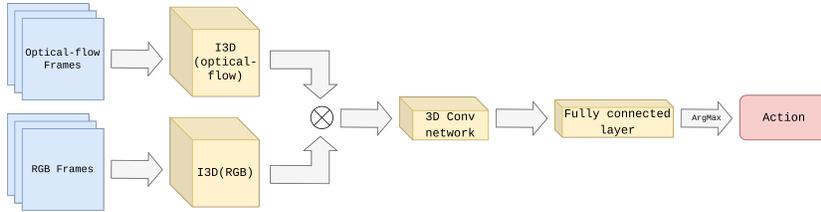


Figure 5: SSHA model architecture (Two-stream fusion).

### 3.3. Training

The SSHA model is trained using Q learning. A value iteration method (such as Q learning) uncovers the underlying value function of an environment. Comparatively, policy iteration approaches immediately identify the most effective actions for each state. The advantage of knowing the hidden value structure of the environment is the higher data efficiency (Hamadouche et al., 2021). In this study, value iteration methods are preferred because the dataset size is limited (e.g., 1600 training videos in the RWF dataset). The Q-learning equation is presented in Equation 4.

$$\underbrace{Q(s, a)}_{\text{New Q value}} = \underbrace{R(s, a)}_{\text{Reward}} + \underbrace{\gamma}_{\text{Discount factor}} \underbrace{\max_{a'} Q'(s', a')}_{\text{Max expected Q value obtained using the target network}} \quad (4)$$

Q value update given the current state and action ( $Q(s, a)$ ) is calculated using Equation 4. The new Q and observed reward ( $R(s, a)$ ) are functions of

current state ( $s$ ) and the applied action ( $a$ ). While the expected Q value ( $Q'$ ) is the function of expected state ( $s'$ ) and future action ( $a'$ ).

The network outputs include two distinct groups of nodes that select attention regions or video classes. The node with the highest Q-value indicates the network choice at each stage. When a region selection action is selected, the intended region is cropped and used as input for the next inference step. In contrast to region selection actions, class selection actions indicate the classification decision and complete the violence detection process.

SSHA does not use recurrent or memory layers, so it cannot remember how many steps it took. In order to compensate for the lack of memory, the number of steps is fed into the fully connected layer. A one-hot vector encodes the number of steps. Making the network aware of the number of actions taken previously prevents the infinite loop scenario of selecting regions indefinitely.

In general, vanilla deep reinforcement learning models are unstable during training. To stabilize deep models, dueling training (Wang et al., 2016b) is used along with regularization, small learning rates, and batch-size tuning. Moreover, an adaptive sampling method maintains the reward sparsity in the */epsilon - greedy* exploration algorithm (Sutton & Barto, 2018). The reward sparsity is defined as the low probability of receiving positive rewards, while the model explores the environment in the training phase. The adaptive sampling method keeps the reward sparsity constant throughout the training phase by calculating the current reward sparsity and tuning the sampling probability from positive and negative rewards accordingly. Adaptive sampling is especially beneficial in the early stages of training as the model is entirely random, and environment complexity causes the samples to be primarily negative.

The training is further stabilized and accelerated using a reward injection technique. Based on the known classification label for each video, we can determine the expected reward for correct and incorrect classifications. Since relevant information is provided without additional observations, the known reward accelerates the training. Since the output activation function is linear (pass-through) and the randomly initialized model outputs can have a large

range. This strategy stabilizes the training phase by reducing the search space size at the start of the training procedure. Large Q values disrupt the main network’s training stability because of the large learning gradients. The reward trimming method is also used to maintain training stability. Using this method, Q values that exceed or fall below a predefined range would be clipped to the maximum or minimum possible value. Applications and architectures define the predefined range. Considering the region selection reward, the expected Q value is between  $-1$  and  $1 + (N - 1) * .5$ , where  $N$  is the maximum number of steps the environment allows. In this study,  $N$  is equal to 5 as with more zooming, the quality of the input video falls below the SSHA model’s input size. SSHA’s training procedure is thoroughly demonstrated in Algorithm 1 to 3.

---

**Algorithm 1** SSHA training: train loop.

---

```

1: for num_episodes do
2:   exploration()
3:   network_update()
4: end for

```

---

## 4. Experiments

Using classification metrics reported in the respective studies, the SSHA model is compared to the previous state-of-the-art models for video violence detection. Further, information regarding the class-level performance and model actions is provided for a more detailed assessment of the SSHA model. Based on the evaluation results, the advantages and limitations of the SSHA model are discussed in section 5.

### 4.1. Experiments setup

The primary dataset used in this study is the RWF dataset (Cheng et al., 2021). The RWF dataset is one of the most comprehensive datasets available for the detection of video violence. A total of 2000 videos were included in the RWF dataset, divided into violent and non-violent categories of equal size. Video

---

**Algorithm 2** SSHA training:  $\epsilon$  – greedy exploration.

---

```
1:  $state_{i-1} = state_i$ 
2:  $action = -1$ 
3: while  $action == -1$  do
4:   if  $random\_number() < .5$  then
5:     if  $random\_number() < \epsilon$  then
6:        $selected\_action = select\_random\_action()$ 
7:     else
8:        $selected\_action = argmax(main\_network(state_i))$ 
9:     end if
10:     $state_i, reward, done = environment.act(action)$ 
11:  else
12:     $state_{i-1}, action, state_i, reward, done =$ 
13:       $replay\_buffer.random\_sample(action)$ 
14:  end if
15:   $prob = random\_number()$ 
16:  if ( $reward > 0$ 
17:    and  $prob < positive\_reward\_selection\_prob$ )
18:    or ( $reward < 0$  and  $prob >$ 
19:       $positive\_reward\_selection\_prob$ ) then
20:     $action = selected\_action$ 
21:     $reward\_history.append(reward)$ 
22:     $replay\_buffer.append($ 
23:       $state_{i-1}, action, state_i, reward, done)$ 
24:  end if
25: end while
```

---

---

**Algorithm 3** SSHA training: network update.

---

```
1: positive_reward_prob =  
    calculate_positive_reward_prob(reward_history)  
2: if positive_reward_prob > target_positive_reward_prob then  
3:   decrease_positive_reward_selection_prob()  
4: else  
5:   increase_positive_reward_selection_prob()  
6: end if  
7:  
8: if done then  
9:   statei = environment.reset()  
10: end if  
11:  
12: Qtarget = Q_learning_equation(  
    statei, reward, target_network)  
13: Qtarget = inject_known_Q_values(Qtarget)  
14: optimize_main_network(Qtarget)  
15:  $\epsilon = \epsilon - \frac{1}{\text{num\_episodes}}$ 
```

---

content in RWF is more abundant in quantity and diversity than in previously published datasets. Because this dataset is scraped from the YouTube <sup>1</sup> video streaming service, it contains videos with varying resolutions. The videos are divided into five-second clips, each with a frame rate of 30 frames per second. It is easier to compare competing models with the help of a predefined list of train and test videos. Classification accuracy is the reported metric in the baseline RWF paper is (Cheng et al., 2021).

Hockey and movie fight datasets are used in addition to the RWF for SSHA model evaluation. Despite containing 1000 and 200 videos, they do not meet the criteria for a practical dataset for building a video violence detection model in the real world. In these datasets, violence is only depicted in hockey games and Hollywood films. A generalizable violence detection model cannot be trained with the repeating situation of hockey players fighting on an icy hockey field and the cinematic quality and perspective of the cinematic film. Therefore, despite the high accuracy of models on the mentioned datasets, they are not useful for real-world and general applications. Dataset characteristics are presented in Table 1.

Dataset	# videos	Video length (seconds)	Video size (pixels)
RWF (Cheng et al., 2021)	2000	5	varied
Hockey fights (Nievas et al., 2011)	1000	2	360 x 288
Movie fights (Gong et al., 2008)	200	2	720 x 480

Table 1: Video violence datasets characteristics.

#### 4.2. Results

The SSHA model is trained and tested using the predefined training and testing videos included in the RWF dataset. Furthermore, the hockey and movie datasets are divided 80/20 between training (80%) and testing (20%). SSHA’s accuracy and class-level performance are presented in Table 2 and 3,

---

<sup>1</sup><https://youtube.com>

respectively. The class-level results provide insight into the SSHA model’s performance in detecting violent and non-violent videos in isolation. In addition, Table 4 provides some characteristics of the SSHA model, such as the model size and the average number of actions per video.

Model/Dataset	RWF (Cheng et al., 2021)	Hockey fights (Nievas et al., 2011)	Movie fights (Gong et al., 2008)
ConvLSTM (Sudhakaran & Lanz, 2017)	77.0 %	97.1 %	100.0 %
I3D (RGB only) (Carreira & Zisserman, 2017)	85.7 %	98.5 %	100.0 %
I3D (Optical-flow only) (Carreira & Zisserman, 2017)	75.5 %	84.0 %	100.0 %
I3D (Two-stream) (Carreira & Zisserman, 2017)	81.5 %	97.5 %	100.0 %
Cheng et al. (RBG only) (Cheng et al., 2021)	84.5 %	-	-
Cheng et al. (Optical-flow only) (Cheng et al., 2021)	75.5 %	-	-
Cheng et al. (C3D) (Cheng et al., 2021)	85.7 %	-	-
Cheng et al. (P3D) (Cheng et al., 2021)	87.2 %	98.0 %	100.0 %
SSHA model (RBG only no localization)	85.3 %	98.0 %	99.0 %
SSHA model (RBG only)	<b>90.4 %</b>	<b>98.7 %</b>	99.0 %
SSHA model (Optical-flow only)	76.0 %	86.2 %	98.5 %
SSHA model (Two-stream)	86.4 %	97.0 %	99.0 %

Table 2: Models accuracy banchmark.

Class	Precision	Recall	F1-score
Violent	0.88 %	0.92 %	0.9 %
Non violent	0.92 %	0.9 %	0.91 %
Total	0.9 %	0.91 %	0.9 %

Table 3: SSHA Class-level evluation results on RWF dataset.

Property	Value
# parameters (total)	13.2 million
# parameters (trained)	.9 million
# input frames	79
Input size	224 x 224 pixels
Avg. # actions per video (RWF)	1.8

Table 4: SSHA Class-level evluation results.

## 5. Discussion

As shown in Table 2 the SSHA model has achieved state-of-the-art accuracy on RWF and Hockey fights datasets, and fair accuracy on the Movie fights dataset using the RBG-only architecture. The accuracy of models has reached

saturation for the Hockey and Movie fights datasets, but state-of-the-art models continue to produce meaningful results on RWF. As a result, the SSHA model’s advantages are more evident in the RWF dataset. Furthermore, this study does not suffer from imbalance learning because the evaluation datasets are perfectly balanced. As shown in Table 3, the class-level results of the SSHA model on the RWF dataset present a proper balance between the violence and non-violence classes.

To assess the effects of hard attention on the accuracy of the SSHA model, the SSHA model is trained and evaluated with and without hard attention capability (RGB only and RGB only no localization in Table 2). The attention mechanism is detached from the SSHA model by removing the region selection actions from the model and converting it to a single-stage video classification model. Nonetheless, this model is still trained using reinforcement learning loss. The superior accuracy of the SSHA model with region selection capabilities demonstrates the effectiveness of the hard attention method.

Learned from the training phase, 1.8 is the optimized average number of actions per video learned by the SSHA model (including the classification action). When only one step is taken to classify a video, the model has classified the video without applying region selection. This scenario is reasonable when the region of interest is large; thus, no further attention is required for accurate classification. Many region selection actions result in a narrow viewpoint on the input video. A perspective with such a narrow field of view often lacks the necessary visual details. Consequently, the average number of actions per video indicates the model’s tendency toward one region selection action before classification.

The notable conclusion of experimenting with a two-stream fusion neural architecture is the drawback of having a larger model trained on a limited dataset. Even though an I3D network with a two-stream fusion architecture achieves higher accuracy on the Kinetics dataset (Carreira & Zisserman, 2017), results on the RWF dataset do not follow the same principle. According to Table 2, the I3D two-stream fusion has inferior accuracy on the RWF dataset compared

to the RGB-only type in previous and current research. The bad performance results from expanding the feature space and, subsequently, the neural network search space while having a fixed number of training data samples. However, the positive aspect of this outcome is the performance-wise preferability of an RGB-only architecture. The I3D backbone has 13 million parameters (Carreira & Zisserman, 2017). A two-stream architecture utilizes two I3D backbones for feature extraction from RGB and Optical-flow frames, doubling the number of parameters relative to the single-stream architecture. Thus, the computational overhead of extracting Optical-flow frames from RGB frames and a double-sized neural network adversely affects the two-stream architecture’s performance.

## 6. Conclusion

This paper presents a semi-supervised hard attention mechanism (SSHA) based on reinforcement learning. SSHA achieves state-of-the-art accuracy in video violence detection by analyzing the most critical region of the video in greater depth. It utilizes video violence datasets that are readily available and eliminates the need for specialized datasets or annotations. The multi-stage implementation of SSHA enables the proposed model to utilize high-definition surveillance footage by selecting attention regions according to the user’s preferences. The RGB-only version of the SSHA model achieved state-of-the-art 90.4 percent accuracy on the RWF dataset and 98.5 and 99.5 percent accuracy on the Hockey and Movies datasets, respectively.

Even though the proposed SSHA model significantly improved the accuracy of the existing state-of-the-art models, future research could apply the hard attention mechanism to action recognition to improve the accuracy of SSHA methods. Additionally, applying the proposed hard attention mechanism to multi-attention scenarios using collaborative agents would be an interesting future direction for this research. Contributing to the RWF dataset regarding the number and quality of available videos and annotations will provide the greatest benefit to state-of-the-art automated video violence detection in the short

term.

## References

- Algamdi, A. M., Sanchez, V., & Li, C.-T. (2019). Learning temporal information from spatial information using capsnets for human action recognition. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (pp. 3867–3871). IEEE.
- Arnab, A., Dehghani, M., Heigold, G., Sun, C., Lučić, M., & Schmid, C. (2021). Vivit: A video vision transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (pp. 6836–6846).
- Bugeja, J., Jönsson, D., & Jacobsson, A. (2018). An investigation of vulnerabilities in smart connected cameras. In *2018 IEEE international conference on pervasive computing and communications workshops (PerCom workshops)* (pp. 537–542). IEEE.
- Caicedo, J. C., & Lazebnik, S. (2015). Active object localization with deep reinforcement learning. In *Proceedings of the IEEE international conference on computer vision* (pp. 2488–2496).
- Carreira, J., & Zisserman, A. (2017). Quo vadis, action recognition? a new model and the kinetics dataset. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 6299–6308).
- Cheng, M., Cai, K., & Li, M. (2021). Rwf-2000: An open large scale video database for violence detection. In *2020 25th International Conference on Pattern Recognition (ICPR)* (pp. 4183–4190). IEEE.
- Crasto, N., Weinzaepfel, P., Alahari, K., & Schmid, C. (2019). Mars: Motion-augmented rgb stream for action recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (pp. 7882–7891).

- Cui, Y., Hou, B., Wu, Q., Ren, B., Wang, S., & Jiao, L. (2021). Remote sensing object tracking with deep reinforcement learning under occlusion. *IEEE Transactions on Geoscience and Remote Sensing*, .
- Driessens, J. (2019). *Focused imagination: Hard attention for reinforcement learning with imagination*. Ph.D. thesis Master’s thesis, University of Amsterdam.
- Du, W., Wang, Y., & Qiao, Y. (2017). Recurrent spatial-temporal attention network for action recognition in videos. *IEEE Transactions on Image Processing*, *27*, 1347–1360.
- Ezen-Can, A. (2020). A comparison of lstm and bert for small corpus. *arXiv preprint arXiv:2009.05451*, .
- Feichtenhofer, C., Pinz, A., & Zisserman, A. (2016). Convolutional two-stream network fusion for video action recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 1933–1941).
- Gill, M., & Spriggs, A. (2005). *Assessing the impact of CCTV* volume 292. Home Office Research, Development and Statistics Directorate London.
- Girdhar, R., Carreira, J., Doersch, C., & Zisserman, A. (2019). Video action transformer network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (pp. 244–253).
- Gong, Y., Wang, W., Jiang, S., Huang, Q., & Gao, W. (2008). Detecting violent scenes in movies by auditory and visual cues. In *Pacific-Rim Conference on Multimedia* (pp. 317–326). Springer.
- Hamadouche, M., Dezan, C., Espes, D., & Branco, K. (2021). Comparison of value iteration, policy iteration and q-learning for solving decision-making problems. In *2021 International Conference on Unmanned Aircraft Systems (ICUAS)* (pp. 101–110). IEEE.

- Hara, K., Kataoka, H., & Satoh, Y. (2017). Learning spatio-temporal features with 3d residual networks for action recognition. In *Proceedings of the IEEE International Conference on Computer Vision Workshops* (pp. 3154–3160).
- Hu, J.-F., Zheng, W.-S., Pan, J., Lai, J., & Zhang, J. (2018). Deep bilinear learning for rgb-d action recognition. In *Proceedings of the European conference on computer vision (ECCV)* (pp. 335–351).
- Jaderberg, M., Simonyan, K., Zisserman, A. et al. (2015). Spatial transformer networks. *Advances in neural information processing systems*, 28.
- Ji, S., Xu, W., Yang, M., & Yu, K. (2012). 3d convolutional neural networks for human action recognition. *IEEE transactions on pattern analysis and machine intelligence*, 35, 221–231.
- Jie, Z., Liang, X., Feng, J., Jin, X., Lu, W., & Yan, S. (2016). Tree-structured reinforcement learning for sequential object localization. In *Advances in Neural Information Processing Systems* (pp. 127–135).
- Khan, F. S., Van De Weijer, J., Anwer, R. M., Bagdanov, A. D., Felsberg, M., & Laaksonen, J. (2018). Scale coding bag of deep features for human attribute and action recognition. *Machine Vision and Applications*, 29, 55–71.
- Li, J., Liu, X., Zhang, W., Zhang, M., Song, J., & Sebe, N. (2020). Spatio-temporal attention networks for action recognition and detection. *IEEE Transactions on Multimedia*, 22, 2990–3001.
- Li, W., Zhu, X., & Gong, S. (2018). Harmonious attention network for person re-identification. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 2285–2294).
- Li, X., Hou, Y., Wang, P., Gao, Z., Xu, M., & Li, W. (2021). Trear: Transformer-based rgb-d egocentric action recognition. *IEEE Transactions on Cognitive and Developmental Systems*, .

- Li, Y., Li, W., Mahadevan, V., & Vasconcelos, N. (2016). Vlad3: Encoding dynamics of deep features for action recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 1951–1960).
- Liu, J., Shahroudy, A., Xu, D., & Wang, G. (2016). Spatio-temporal lstm with trust gates for 3d human action recognition. In *European conference on computer vision* (pp. 816–833). Springer.
- Liu, J., Wang, G., Hu, P., Duan, L.-Y., & Kot, A. C. (2017). Global context-aware attention lstm networks for 3d action recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 1647–1656).
- Luo, W., Sun, P., Zhong, F., Liu, W., Zhang, T., & Wang, Y. (2019). End-to-end active object tracking and its real-world deployment via reinforcement learning. *IEEE transactions on pattern analysis and machine intelligence*, *42*, 1317–1332.
- Luvizon, D. C., Picard, D., & Tabia, H. (2018). 2d/3d pose estimation and action recognition using multitask deep learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 5137–5146).
- Majd, M., & Safabakhsh, R. (2020). Correlational convolutional lstm for human action recognition. *Neurocomputing*, *396*, 224–229.
- Malinowski, M., Doersch, C., Santoro, A., & Battaglia, P. (2018). Learning visual question answering by bootstrapping hard attention. In *Proceedings of the European Conference on Computer Vision (ECCV)* (pp. 3–20).
- Manchin, A., Abbasnejad, E., & Hengel, A. v. d. (2019). Reinforcement learning with attention that works: A self-supervised approach. In *International conference on neural information processing* (pp. 223–230). Springer.

- Mathe, S., Pirinen, A., & Sminchisescu, C. (2016). Reinforcement learning for visual object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 2894–2902).
- Mazzia, V., Angarano, S., Salvetti, F., Angelini, F., & Chiaberge, M. (2021). Action transformer: A self-attention model for short-time human action recognition. *arXiv preprint arXiv:2107.00606*, .
- Mott, A., Zoran, D., Chrzanowski, M., Wierstra, D., & Jimenez Rezende, D. (2019). Towards interpretable reinforcement learning using attention augmented agents. *Advances in Neural Information Processing Systems*, 32.
- Murthy, O. R., & Goecke, R. (2014). The influence of temporal information on human action recognition with large number of classes. In *2014 International Conference on Digital Image Computing: Techniques and Applications (DICTA)* (pp. 1–8). IEEE.
- Nazir, S., Yousaf, M. H., & Velastin, S. A. (2018). Evaluating a bag-of-visual features approach using spatio-temporal features for action recognition. *Computers & Electrical Engineering*, 72, 660–669.
- Nguyen, M. T., Truong, L. H., Tran, T. T., & Chien, C.-F. (2020). Artificial intelligence based data processing algorithm for video surveillance to empower industry 3.5. *Computers & Industrial Engineering*, 148, 106671.
- Nievas, E. B., Suarez, O. D., García, G. B., & Sukthankar, R. (2011). Violence detection in video using computer vision techniques. In *International conference on Computer analysis of images and patterns* (pp. 332–339). Springer.
- Pham, H. H., Salmane, H., Khoudour, L., Crouzil, A., Velastin, S. A., & Zegers, P. (2020). A unified deep framework for joint 3d pose estimation and action recognition from a single rgb camera. *Sensors*, 20, 1825.

- Plizzari, C., Cannici, M., & Matteucci, M. (2021). Spatial temporal transformer network for skeleton-based action recognition. In *International Conference on Pattern Recognition* (pp. 694–701). Springer.
- Rao, Y., Lu, J., & Zhou, J. (2017). Attention-aware deep reinforcement learning for video face recognition. In *Proceedings of the IEEE international conference on computer vision* (pp. 3931–3940).
- Ren, L., Lu, J., Wang, Z., Tian, Q., & Zhou, J. (2018). Collaborative deep reinforcement learning for multi-object tracking. In *Proceedings of the European Conference on Computer Vision (ECCV)* (pp. 586–602).
- Sevilla-Lara, L., Liao, Y., Güney, F., Jampani, V., Geiger, A., & Black, M. J. (2018). On the integration of optical flow and action recognition. In *German Conference on Pattern Recognition* (pp. 281–297). Springer.
- Sharma, S., Kiros, R., & Salakhutdinov, R. (2015). Action recognition using visual attention. *arXiv preprint arXiv:1511.04119*, .
- Shen, T., Zhou, T., Long, G., Jiang, J., Wang, S., & Zhang, C. (2018). Reinforced self-attention network: a hybrid of hard and soft attention for sequence modeling. *arXiv preprint arXiv:1801.10296*, .
- Shi, L., Zhang, Y., Cheng, J., & Lu, H. (2019). Skeleton-based action recognition with directed graph neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (pp. 7912–7921).
- Shi, L., Zhang, Y., Cheng, J., & Lu, H. (2020). Skeleton-based action recognition with multi-stream adaptive graph convolutional networks. *IEEE Transactions on Image Processing*, 29, 9532–9545.
- Shidik, G. F., Noersasongko, E., Nugraha, A., Andono, P. N., Jumanto, J., & Kusuma, E. J. (2019). A systematic review of intelligence video surveillance: trends, techniques, frameworks, and datasets. *IEEE Access*, 7, 170457–170473.

- Simonyan, K., & Zisserman, A. (2014). Two-stream convolutional networks for action recognition in videos. *Advances in neural information processing systems*, 27.
- Smaira, L., Carreira, J., Noland, E., Clancy, E., Wu, A., & Zisserman, A. (2020). A short note on the kinetics-700-2020 human action dataset. *arXiv preprint arXiv:2010.10864*, .
- Song, H., Kim, M., Park, D., Shin, Y., & Lee, J.-G. (2022). Learning from noisy labels with deep neural networks: A survey. *IEEE Transactions on Neural Networks and Learning Systems*, .
- Song, S., Lan, C., Xing, J., Zeng, W., & Liu, J. (2017). An end-to-end spatio-temporal attention model for human action recognition from skeleton data. In *Proceedings of the AAAI conference on artificial intelligence*. volume 31.
- Sudhakaran, S., & Lanz, O. (2017). Learning to detect violent videos using convolutional long short-term memory. In *2017 14th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS)* (pp. 1–6). IEEE.
- Sun, S., Kuang, Z., Sheng, L., Ouyang, W., & Zhang, W. (2018). Optical flow guided feature: A fast and robust motion representation for video action recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 1390–1399).
- Sung, C.-S., & Park, J. Y. (2021). Design of an intelligent video surveillance system for crime prevention: applying deep learning technology. *Multimedia Tools and Applications*, 80, 34297–34309.
- Sutton, R. S., & Barto, A. G. (2018). *Reinforcement learning: An introduction*. MIT press.
- Tu, Z., Xie, W., Dauwels, J., Li, B., & Yuan, J. (2018a). Semantic cues enhanced multimodality multistream cnn for action recognition. *IEEE Transactions on Circuits and Systems for Video Technology*, 29, 1423–1437.

- Tu, Z., Xie, W., Qin, Q., Poppe, R., Veltkamp, R. C., Li, B., & Yuan, J. (2018b). Multi-stream cnn: Learning representations based on human-related regions for action recognition. *Pattern Recognition*, *79*, 32–43.
- Ullah, A., Ahmad, J., Muhammad, K., Sajjad, M., & Baik, S. W. (2017). Action recognition in video sequences using deep bi-directional lstm with cnn features. *IEEE access*, *6*, 1155–1166.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., & Polosukhin, I. (2017). Attention is all you need. In *Advances in neural information processing systems* (pp. 5998–6008).
- Wang, C., Li, M., & Smola, A. J. (2019). Language models with transformers. *arXiv preprint arXiv:1904.09408*, .
- Wang, H., & Schmid, C. (2013). Action recognition with improved trajectories. In *Proceedings of the IEEE international conference on computer vision* (pp. 3551–3558).
- Wang, H., Yu, B., Li, J., Zhang, L., & Chen, D. (2021). Multi-stream interaction networks for human action recognition. *IEEE Transactions on Circuits and Systems for Video Technology*, .
- Wang, L., Li, W., Li, W., & Van Gool, L. (2018a). Appearance-and-relation networks for video classification. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 1430–1439).
- Wang, L., Qiao, Y., & Tang, X. (2016a). Mofap: A multi-level representation for action recognition. *International Journal of Computer Vision*, *119*, 254–271.
- Wang, L., Xu, Y., Cheng, J., Xia, H., Yin, J., & Wu, J. (2018b). Human action recognition by learning spatio-temporal features with deep neural networks. *IEEE access*, *6*, 17913–17922.

- Wang, P., Li, W., Wan, J., Ogunbona, P., & Liu, X. (2018c). Cooperative training of deep aggregation networks for rgb-d action recognition. In *Proceedings of the AAAI Conference on Artificial Intelligence*. volume 32.
- Wang, P., Wang, S., Gao, Z., Hou, Y., & Li, W. (2017). Structured images for rgb-d action recognition. In *Proceedings of the IEEE International Conference on Computer Vision Workshops* (pp. 1005–1014).
- Wang, Y., Zhang, L., Wang, L., & Wang, Z. (2018d). Multitask learning for object localization with deep reinforcement learning. *IEEE Transactions on Cognitive and Developmental Systems*, 11, 573–580.
- Wang, Z., Schaul, T., Hessel, M., Hasselt, H., Lanctot, M., & Freitas, N. (2016b). Dueling network architectures for deep reinforcement learning. In *International conference on machine learning* (pp. 1995–2003). PMLR.
- Xiao, R., Hou, Y., Guo, Z., Li, C., Wang, P., & Li, W. (2019). Self-attention guided deep features for action recognition. In *2019 IEEE International Conference on Multimedia and Expo (ICME)* (pp. 1060–1065). IEEE.
- Yan, S., Xiong, Y., & Lin, D. (2018). Spatial temporal graph convolutional networks for skeleton-based action recognition. In *Thirty-second AAAI conference on artificial intelligence*.
- Yang, H., Yuan, C., Li, B., Du, Y., Xing, J., Hu, W., & Maybank, S. J. (2019). Asymmetric 3d convolutional neural networks for action recognition. *Pattern recognition*, 85, 1–12.
- Yao, L., Liu, Y., & Huang, S. (2016). Spatio-temporal information for human action recognition. *EURASIP Journal on Image and Video Processing*, 2016, 1–9.
- Yun, S., Choi, J., Yoo, Y., Yun, K., & Choi, J. Y. (2018). Action-driven visual object tracking with deep reinforcement learning. *IEEE transactions on neural networks and learning systems*, 29, 2239–2252.

- Zeyer, A., Bahar, P., Irie, K., Schlüter, R., & Ney, H. (2019). A comparison of transformer and lstm encoder decoder models for asr. In *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)* (pp. 8–15). IEEE.
- Zhang, J., Li, W., Ogunbona, P. O., Wang, P., & Tang, C. (2016). Rgb-d-based action recognition datasets: A survey. *Pattern Recognition*, *60*, 86–105.
- Zhong, Z., Yang, Z., Feng, W., Wu, W., Hu, Y., & Liu, C.-L. (2019). Decision controller for object tracking with deep reinforcement learning. *IEEE Access*, *7*, 28069–28079.
- Zhou, Y., Sun, X., Zha, Z.-J., & Zeng, W. (2018). Mict: Mixed 3d/2d convolutional tube for human action recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 449–458).
- Zong, M., Wang, R., Chen, X., Chen, Z., & Gong, Y. (2021). Motion saliency based multi-stream multiplier resnets for action recognition. *Image and Vision Computing*, *107*, 104108.