# Sparse Black-box Video Attack with Reinforcement Learning

**Xingxing Wei · Huanqian Yan · Bo Li**

**Abstract** Adversarial attacks on video recognition models have been explored recently. However, most existing works treat each video frame equally and ignore their temporal interactions. To overcome this drawback, a few methods try to select some key frames and then perform attacks based on them. Unfortunately, their selection strategy is independent of the attacking step, therefore the resulting performance is limited. Instead, we argue the frame selection phase is closely relevant with the attacking phase. The key frames should be adjusted according to the attacking results. For that, we formulate the black-box video attacks into a Reinforcement Learning (RL) framework. Specifically, the environment in RL is set as the recognition model, and the agent in RL plays the role of frame selecting. By continuously querying the recognition models and receiving the attacking feedback, the agent gradually adjusts its frame selection strategy and adversarial perturbations become smaller and smaller. We conduct a series of experiments with two mainstream video recognition models: *C3D* and *LRCN* on the public *UCF-101* and *HMDB-51* datasets. The results demonstrate that the proposed method can significantly reduce the adversarial perturbations with efficient query times.

**Keywords** Adversarial Examples · Black-box Video Attack · Reinforcement Learning · Sparse Attack

Xingxing Wei
Institute of Artificial Intelligence, Hangzhou Innovation Institute, Beihang University, Beijing, China
E-mail: xxwei@buaa.edu.cn

Huanqian Yan · Bo Li
Beijing Key Laboratory of Digital Media (DML) and State Key Laboratory of Virtual Reality Technology and Systems, School of Computer Science and Engineering, Beihang University, Beijing, China
E-mail: yanhq@buaa.edu.cn, boli@buaa.edu.cn

# 1 Introduction

Deep Neural Networks (DNNs) have achieved great success in a wide range of tasks [22,38,6,33,10,25]. Despite this fact, it is proved that deep neural networks are vulnerable to adversarial examples [11,36,9,18,42]. Recent works have shown that adding a carefully crafted, small human-imperceptible perturbation to a clean sample can make the deep neural models crash in image classification [1,11,37], object detection [44,43,31], semantic segmentation [43] and other tasks. Nowadays, more and more DNN models are deployed in various sectors with high-security requirements, which has caused the study of adversarial examples to gain increased attention. The existence of adversarial examples brings huge security risks to the deployment of deep learning systems, such as automatic driving [26], robotics [42], face recognition [2,9,48] and other aspects.

Due to many real-time video classification systems which are constructed based on the DNN models, it is crucial to investigate the adversarial examples for video models. On the one hand, video attacks can help researchers understand the working mechanism of time-series deep models. On the other hand, adversarial samples facilitate various deep neural network algorithms to assess the robustness by providing more varied video training data [39,9,49]. In this paper, we focus on video attacks, specifically attacking the video classification model under the black-box condition. Compared with the white-box setting, the black-box video attack is more realistic. Because the white-box attack needs to obtain the structures and parameters of the deep learning model, and it is usually difficult in the real applications. In contrast, black-box attacks don't need these information. A widely used way is to access the output of the target model when the input is given. In this way,
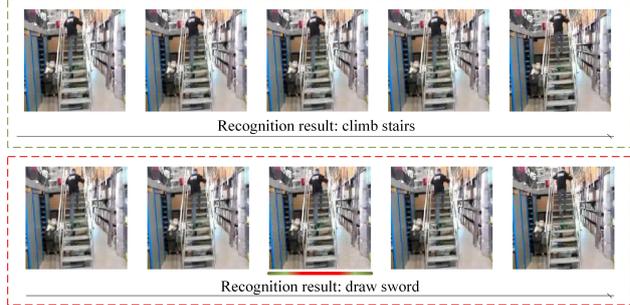
the gradients can be estimated to generate adversarial examples.

According to [17], the current video attacking methods can be roughly divided into two classes. The first class is called *dense attack* which pollutes each frame in a video [23, 19], and the second class is to select some key frames, and then generates perturbations on these selected frames [39, 40], called as *sparse attack*. Compared with the dense attack, the sparse attack is more reasonable because there are temporal interactions between adjacent frames in the video. Utilizing this relationship can help both reduce the adversarial perturbations and improve the efficiency of the generation process. For the former advantage, because the selected frames are the most important ones in a video, only adding small perturbations on these frames can fool the recognition model, leading to the reduction of adversarial perturbations on the whole video. For the second advantage, the selected key frames are usually sparse, compared with generating perturbations on the whole frames, the operation of dealing with a few frames is more efficient.

To better select key frames in the sparse attack, a heuristic black-box attack on video recognition models is proposed [40]. They first propose a heuristic algorithm to evaluate the importance of each frame, and then select key frames by sorting the importance scores, finally, the black-box attacks are performed on the selected frames. However, there are no interactions between the attacking process and the selecting key frames in their method. We argue that key frame selection should not only depend on the video itself, but also on the feedback from the recognition models. The frame selection and video attacking are characterized by mutual guidance and cooperation. The results of this way can produce more accurate key frames and smaller perturbations for adversarial videos.

To this end, we present a Sparse black-box Video Attack (SVA) method with Reinforcement Learning (RL) in this paper. Specifically, the environment in RL is set as the recognition model, and the agent in RL plays the role of frame selecting. By continuously querying the recognition models and receiving the feedback of predicted probabilities (rewards), the agent adjusts its frame selection strategy and performs attacks (actions). Step by step, the optimal key frames are selected and the smallest adversarial perturbations are achieved.
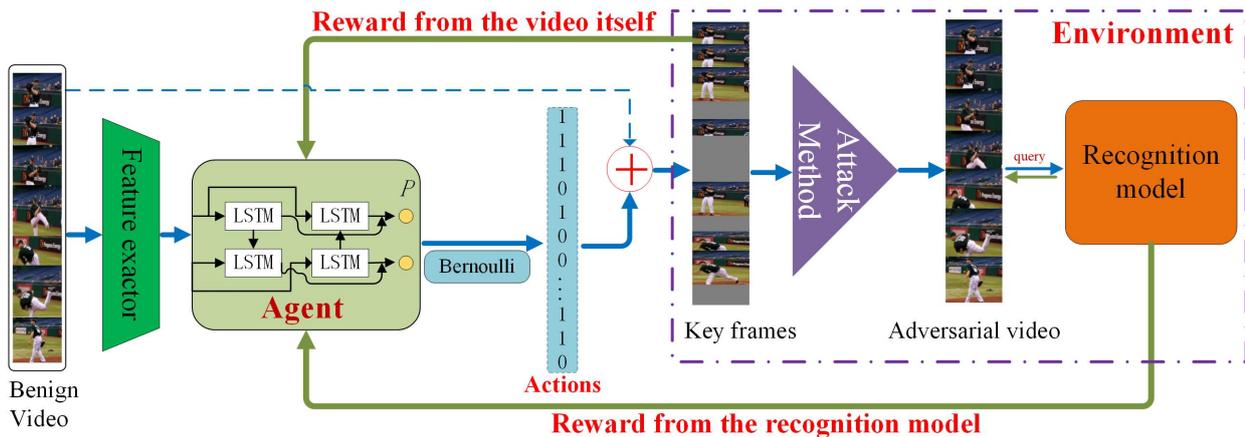
Technically, for the agent, we use an LSTM-based network [15] to measure the importance of each frame. After obtaining the importance degree $p_t$ of each frame $t$ in the video, Bernoulli sampling is used to determine whether the frames are critical or not. To perform the black-box attacks, the Natural Evolution Strategy



**Fig. 1** An example of sparse black-box video attack with Reinforcement Learning (RL). The clean video (top) can be recognized correctly. The adversarial video (bottom) produced by our proposed method is misclassified. Note that only one frame (green-red line annotation) is adaptively selected by the RL, and very small perturbations are added to the key frame.

(NES) [47, 16] is utilized to estimate the gradient from the recognition models, and then the adversarial videos are generated based on these gradients. For the reward, we design two kinds of functions, the first one comes from the video itself. For example, the frames with big action changes will have the high confidence to be the key frames. The other one comes from the feedback of attacking recognition models. The insight is that if the frames with tiny perturbations will lead to a big drop of predicted probability, these frames will have the high confidence to be the key frames. Figure 1 shows an adversarial video generated by our proposed method. Figure 2 overviews the proposed method. Our major contributions can be summarized as follows:

- We are the first one to use reinforcement learning to attack video recognition models in the black-box setting. In contrast, previous works perform adversarial attacks on the reinforcement learning model itself, and thus our approach differs fundamentally from those works.
- A novel algorithm is designed for selecting key frames from a video when attacking video recognition models, which is based on two factors including the visual features of the video itself and the feedback given by the recognition model. Video attacking and key frame selecting are cooperated and guided by each other.
- Extensive experiments on two widely used video recognition models (LRCN and C3D) and two benchmark video datasets (UCF-101 and HMDB-51) show that the proposed method can significantly reduce the adversarial perturbations while only needing a few query times compared with the state-of-the-art video attacking methods.

**Fig. 2** Overview of the black-box adversarial video attacking method. We formulate the key frame selection and attacking step into the reinforcement learning framework. Please see the texts for details.

The rest of this paper is organized as follows. In Section 2, we briefly review the related work. The proposed algorithm is described in Section 3. The experimental results and analysis are presented in Section 4. Finally, we summarize the conclusion in Section 5.

## 2 Related Work

**Adversarial Attack on Video Models:** Adversarial attacks on images have been extensively studied [12,5, 11]. Compared with images, the dimension of videos is very high. It is not easy for general attack algorithms to directly attack such high-dimensional data. A high dimension leads to large search space, and the algorithm needs more query times to find the optimal perturbations to accomplish the successful attack. Therefore, it is more difficult to design an efficient black-box video attack algorithm.

In the past years, many adversarial attacks for videos have been proposed. An $l_{2,1}$-norm regularization based optimization algorithm is the first method that is proposed to compute the sparse adversarial perturbations for video recognition [39]. $l_2$ is used to make the frame have small perturbations and $l_1$ is used to make adversarial frames fewer. The 3D universal perturbation [23] is generated by Generative Adversarial Networks offline and then used with unseen input for the real-time video recognition model. Unlike such white-box attack algorithms which need some knowledge about the video recognition models, Jiang *et al.* utilize tentative perturbations and partition-based rectifications to obtain good adversarial gradient estimates and high attack success rate in the black-box setting [19]. But attacking all frames of the video would cause more perturbations and poor robustness of the adversarial

video. The adversarial video can easily become ineffective when one of the frames is randomly replaced with a clean video frame. Another black-box method is proposed by Wei *et al.* [40], they heuristically search a subset of frames and the adversarial perturbations are only generated on those selected frames, but the attacking processes and key frames selection are separated from each other, the perturbations of adversarial videos are still unsatisfactory.

Unlike the algorithms mentioned above, our method generates adversarial perturbations on the key frames which are selected by an agent trained using visual features of the video and the feedback of attacking. Our method could generate smaller adversarial perturbations than state-of-the-art black-box video attack methods.

**Deep Reinforcement Learning:** Deep reinforcement learning is originally designed for learning and mimicking human decision-making processes, which aims to enable the agent to make appropriate behaviors according to the current environment through continuous interaction with the environment [24,28]. It doesn't require any supervisory information unlike supervised machine learning methods, but rather receives a reward signal to evaluate the performance of the action. Reinforcement learning has received a lot of attention since the AlphaGo [32,33] beats humans. Computer vision tasks have also benefited from deep reinforcement learning in recent years. For example, Zhou *et al.* have applied deep reinforcement learning to train a summarization network for video summary [45]. Dong *et al.* use reinforcement learning for action recognition [8]. The process of discarding some irrelevant frames is a kind of hard attention mechanism in their method. Besides,

it has been applied in some other fields like tracking, identification [50], and person search [24].

However, there is no example that deep reinforcement learning is applied in generating adversarial examples. Reinforcement learning algorithms have similar implementation mechanisms with adversarial attack algorithms, especially black-box attack algorithms. It is the first time that we attempt to apply reinforcement learning to the video black-box adversarial attacks. An agent is designed to select key frames while attacking a video using a novel reward function. The key frames selection and adversarial attacks are mutual guidance and cooperation in the whole attacking process.

## 3 Methodology

The adversary takes the video classifier $F(\cdot)$ as a black-box oracle and can only get its output of the top-1 class and its probability. Specifically, given a clean video $x$ and its ground-truth label $\bar{y}$, $F(\cdot)$ takes $x$ as an input and outputs the top-1 class label $F(x) = y$ and its probability $P(y|x)$. If the prediction is correct, then $y = \bar{y}$. The adversarial attack aims to find an adversarial example $x_{adv}$ which can make $F(x_{adv}) \neq \bar{y}$ in the un-targeted attack or $F(x_{adv}) = y_{adv}$ in the targeted attack with the targeted adversarial class $y_{adv}$, while keeping the adversarial example $x_{adv}$ satisfying the condition: $\| x_{adv} - x \|_\rho \leq \epsilon_{adv}$, where $\epsilon_{adv}$ is the bound of the perturbation $\epsilon$, the $\rho$ in $L_\rho$-norm can be set 0,2,$\infty$.

### 3.1 Video Attacking

The attack algorithm in our method is built based on Fast Gradient Sign Method (FGSM) [11], which is originally designed for image models. It is defined as:

$$x_{adv} = x + \alpha \cdot sign(\widehat{g}),  \qquad (1)$$

where $\alpha$ is the step size. $sign(\cdot)$ is sign function. $\widehat{g}$ is the gradient, in white-box attacking setting, $\widehat{g}$ can be computed using $\nabla_x l_{adv}(x)$. Here $l_{adv}(x)$ is abbreviated for adversarial loss function, which is described with $l_{adv}(x) = -l(F(x), \bar{y})$ in un-targeted attack and $l_{adv}(x) = l(F(x), y_{adv})$ in targeted attack. $l(\cdot)$ is a cross-entropy loss. Due to black-box settings, we cannot get the gradient from the recognition model directly, NES [47] algorithm is used as gradient estimator in the proposed method. For NES algorithm, it first generates $n/2$ values $\delta_i \backsim N(0, I), i \in \{1, 2...n/2\}$, where $n$ is the number of samples. Then, it sets $\delta_j = -\delta_{n-j+1}, j \in$ $\{(n/2 + 1), ..n\}$. Finally, the gradient $\widehat{g}$ is estimated as:

$$\widehat{g} \approx \frac{1}{\Delta n} \sum_{i=1}^{n} \delta_i P(y|x + \Delta \cdot \delta_i),  \qquad (2)$$

where $\Delta$ is the search variance.

We extend FGSM with NES from image models to video models as our attacking algorithm. As mentioned above, we use the agent to select the key frames and attack these key frames to achieve the attack of the entire video. Note that the core contribution in our method is to introduce RL to select key frames, rather than the attacking method module. Therefore, we choose a simple and widely used FGSM+NES method. Other methods like Opt-attack [3,4,46] can also be available.

In addition, for the targeted attack, we first initialize $x_{adv}$ by adding the random noises on the selected frames to make $F(x_{adv}) = y_{adv}$. In this time, the perturbation $\epsilon = ||x_{adv} - x||_\rho$ is large. Next, we will gradually decrease the perturbations' magnitude $\epsilon$ until the given bound $\epsilon_{adv}$ is achieved while keeping the prediction label not change, i.e., the prediction label is still $y_{adv}$. In this way, we obtain the minimal adversarial perturbation for the current selected frames. For the un-targeted attack, we don't need to initialize $x_{adv}$ like the targeted attack, instead, some key frames are randomly selected as victims and the perturbations are added. The key frames will be dynamically adjusted by the agent in the next section, and the optimal perturbation will be solved.

### 3.2 Key Frame Selection

Videos have successive frames in the temporal domain, thus, we consider searching key frames that contribute the most to the success of an adversarial attack. In our approach, key frames selection is considered as a one-step Markov decision process. Figure 2 provides a sketch map of this process. The agent learns to select the frames by maximizing the total expected reward by interacting with an environment that provides the rewards and updating its actions.

The input of the agent is a sequence of visual features of the video frames $\{v_t\}_{t=1}^T$ with the length $T$. The agent is a bidirectional Long-Short Term Memory network (BiLSTM) [15] topped with a fully connected (FC) layer. The BiLSTM produces corresponding hidden states $\{h_t\}_{t=1}^T$. We use the ResNet18 [14] to extract visual features and set the dimension of hidden state in the LSTM cell to 128 throughout this paper. Each $h_t$ contains both information from the forward hidden state $h_t^f$ and the backward hidden state $h_t^b$, which is

a good representation of the time domain information of its surrounding frames. The *FC* layer that ends with the sigmoid function $\sigma$ predicts a probability $p_t$ for each frame, and then the key frames $a_t$ are sampled via a Bernouli function:

$$p_t = \sigma(W \times h_t), \qquad (3)$$

$$a_t = Bernoulli(p_t), \qquad (4)$$

where $a_t \in \{0, 1\}$ indicates whether the $t^{th}$ frame is selected or not. $W$ is the weights of the full connection layer. $h_t = \phi(v_t)$, and $\phi(\cdot)$ is the BiLSTM network.

The reward reflects the quality of different actions. It contains two components in our method: the reward from the inherent attributes of the video itself and the reward from the feedback of the recognition model. The former reward includes diversity reward $R_{div}$ and representative reward $R_{rep}$ [45]. Let the indices of the selected frames be $K = \{t | a_t = 1, t = 1, ..., T\}$, the reward $R_{rep}$ and $R_{div}$ can be defined as:

$$R_{rep} = exp(-\frac{1}{T} \sum_{t=1}^{T} min_{t' \in K} \|v_t - v_{t'}\|_2), \qquad (5)$$

$$R_{div} = \frac{1}{|K|(|K| - 1)} \sum_{t \in K} \sum_{t' \in K, t' \neq t} d(v_t, v_{t'}), \qquad (6)$$

where $d(\cdot, \cdot)$ is the dissimilarity function calculated by

$$d(v_t, v_{t'}) = 1 - \frac{v_t^T v_{t'}}{\|v_t\|_2 \|v_{t'}\|_2}. \qquad (7)$$

The reward from the feedback of the video recognition model is defined as:

$$R_{attack} = \begin{cases} 0.999 \times exp(\frac{-\mathbb{P}}{0.05}) & 30000 > Q > 15000 \\ exp(\frac{-\mathbb{P}}{0.05}) & Q \leqslant 15000 \\ -1 & Q > 30000, \end{cases} \qquad (8)$$

where $Q$ is the number of queries, $\mathbb{P}$ is the mean perturbation of the adversarial video (MAP), 0.05 is a normalization factor, 0.999 is the penalty factor used for reducing the number of queries. The rewards $R_{div}$, $R_{rep}$ and $R_{attack}$ complement each other and work jointly to guide the learning of the agent:

$$R = R_{div} + \gamma_1 R_{rep} + \gamma_2 R_{attack}. \qquad (9)$$

The hyperparameters $\gamma_1$ and $\gamma_2$ are set according to the parameter tuning.

As shown by the reward function, the reward $R_{rep}$ and the reward $R_{div}$ only rely on the internal properties of the attack video. By the constraints of these two rewards, the agent can accurately identify the video frames that are representative and diverse. For the reward $R_{attack}$, it penalizes the actions versus perturbation and query number meanwhile. For query numbers, we set different rewards in three phases. For queries less than 15000 times, we think the query number is acceptable, and the reward is mainly related with the perturbation $exp(-\mathbb{P}/0.05)$, i.e., if the mean absolute perturbation is small, $exp(-\mathbb{P}/0.05)$ will be large. For queries more than 30000, we think the attack is unavailable, and we align a small reward -1, which is much less than $exp(-\mathbb{P}/0.05)$ for any $\mathbb{P}$. For queries more than 15000 and less than 30000, we add a penalty on $exp(-\mathbb{P}/0.05)$ to make its reward less than that in the first case. In this way, the rewards can encourage the agent to make an accurate decision to achieve the small perturbation and query number at the same time.

Since each frame corresponds to two actions, there are $2^T$ possible executions of a video, which is basically not feasible for deep Q learning. Thus, we employ the policy gradient method to make the agent learn a policy function $\pi_\theta$ with parameters $\theta$ by maximizing the expected reward $J(\theta) = E_{\tau \sim \pi_\theta}[R(\tau)]$. Following the REINFORCE algorithm proposed by Williams [41], we approximate the gradient by running the agent for $N$ episodes on the same video and then taking the average gradient:

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{n=1}^{N} \sum_{t=1}^{T} R_n \nabla_\theta log \pi_\theta(a_t | h_t), \qquad (10)$$

where $R_n$ is the reward computed at the $n^{th}$ episode. The number of episodes $N$ is set to 5 in our experiments. Although the above method can estimate the gradient well, it may contain high variance which will make the network hard to converge. A common countermeasure is to subtract the reward by a constant baseline $b$, so the gradient becomes

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{n=1}^{N} \sum_{t=1}^{T} (R_n - b) \nabla_\theta log \pi_\theta(a_t | h_t), \qquad (11)$$

where $b$ is a constant baseline that is used to alleviate the high variance. For computational efficiency, it is set as the moving average of rewards experienced so far. We optimize the policy function's parameter $\theta$ via Adam [20]. We update $\theta$ as: $\theta = \theta + ls \cdot \nabla_\theta J(\theta)$, where $ls$ is learning rate and set to $10^{-5}$ in our experiments.

---

**Algorithm 1:** Our SVA targeted attack

    **Input**       : Target class $y_{adv}$ and clean video $x$.
    **Output**     : Adversarial video $x_{adv}$.
    **Parameters:** Maximum $epochs$, episodes $N$
**1**  $\{v_t\}_{t=1}^T \leftarrow ResNet18(x)$; //extract visual feature
**2**  **for** $i = 1$ to $epochs$ **do**
**3**     $M \leftarrow 0$;
**4**     $p_t \leftarrow Agent(\{v_t\}_{t=1}^T)$ using Eq. (3);
**5**     **for** $j = 1$ to $N$ **do**
**6**         $a_t \leftarrow Bernoulli(p_t)$ using Eq.(4);
**7**         $M(t) = a_t, t = 1, ..., T$;
**8**         $Q, \mathbb{P}, x_{adv} \leftarrow$ Using **Algorithm** 2;
**9**         Compute $R_j \leftarrow$ according to Eq. (9);
**10**   **end**
**11**     Compute $\nabla_\theta J(\theta) \leftarrow$ according to Eq. (11);
**12**     Update the $Agent$: $\theta \leftarrow \theta + ls \cdot \nabla_\theta J(\theta)$;
**13** **end**
**14** **return** $x_{adv}$

---

### 3.3 Overall Framework

Here, the whole process of our method in the targeted setting is described in Algorithm 1, which is a continuous-learning algorithm. The epsilon decay $\triangle_\epsilon$ is used to control the reduction size of the perturbation bound. $\triangle_\epsilon$ and $FGSM$ step size $\alpha$ are dynamically adjusted as described in subsection 3.4. The binary vector $M$ has the same size as the frame number of the input video and some values in $M$ will be set to 1 after the agent is trained. The agent in the whole process of attacking is updated with the rewards' values, which is dynamic and interactive. This process does not need any human intervention.

### 3.4 Implementation Details.

To follow the query limited black-box settings and make our experiments more convenient, the maximum query number is set to $3 \times 10^4$ in the un-targeted mode and $6 \times 10^4$ in the targeted mode for all black-box attack algorithms in our experiments. For $NES$, we set the population size as 48, which works well on different datasets in terms of the success rate or the number of queries. For search variance $\Delta$ in $NES$, because the targeted attack needs to keep the target class as the top-1 position, while the un-targeted attack is to remove the current class from the top-1 position, we set it to $10^{-6}$ for the targeted attack setting and $10^{-3}$ for the un-targeted attack setting. The adjustment of step size $\alpha$ adopts the strategy in [27]. For the targeted attack, we adjust the step size $\alpha$ and epsilon decay $\triangle_\epsilon$ dynamically. If the proportion of the adversarial examples cannot be maintained above the threshold 50%, the step size $\alpha$ is
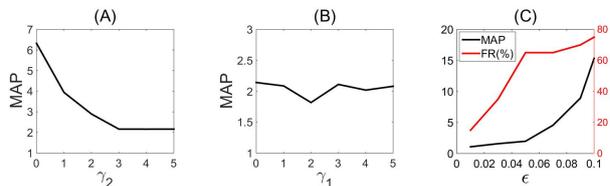
---

**Algorithm 2:** Video attack algorithm

    **Input**       : The classifier $F(\cdot)$, target class $y_{adv}$,
                   clean video $x$, key frame vector $M$.
    **Output**     : Mean pixel perturbation $\mathbb{P}$, query
                   numbers $Q$, adversarial video $x_{adv}$.
    **Parameters:** $FGSM$ step size $\alpha$, epsilon decay $\triangle_\epsilon$,
                       perturbation bound $\epsilon_{adv}$.
**1**  $x_{adv} \leftarrow$ initializing video with the label $y_{adv}$.
**2**  $\epsilon = ||x_{adv} - x||_\rho$; $Q \leftarrow 0$;
**3**  **while** $\epsilon > \epsilon_{adv}$ **do**
**4**     $\widehat{g} \leftarrow$ according to Eq. (2);
**5**     $Q \leftarrow Q + n$; // $n$ is the NES sample numbers.
**6**     $\widehat{g} \leftarrow \widehat{g} \times M$;
**7**     $\widehat{\epsilon} \leftarrow \epsilon - \triangle_\epsilon$; //reduce $\epsilon$ value
**8**     $x'_{adv} \leftarrow x_{adv} + \alpha \cdot sign(\widehat{g})$ using Eq.(1);
**9**     $\widehat{x}_{adv} \leftarrow Clip(x'_{adv}, x - \widehat{\epsilon}, x + \widehat{\epsilon})$;
**10**   **if** $y_{adv} = F(\widehat{x}_{adv})$ **then**
**11**      $Q \leftarrow Q + 1$; $x_{adv} \leftarrow \widehat{x}_{adv}$; $\epsilon \leftarrow \widehat{\epsilon}$;
**12**   **else**
**13**      $\widehat{x}_{adv} \leftarrow Clip(x'_{adv}, x - \epsilon, x + \epsilon)$.
**14**      **if** $y_{adv} = F(\widehat{x}_{adv})$ **then**
**15**        $Q \leftarrow Q + 1$; $x_{adv} \leftarrow \widehat{x}_{adv}$.
**16**      **end**
**17**   **end**
**18**   Adjust $\triangle_\epsilon$ according to the change of $\epsilon$.
**19** **end**
**20** $\mathbb{P} = \frac{||x_{adv} - x||}{|pixel_x|}$ // mean pixel perturbation
**21** **return** $\mathbb{P}, Q, x_{adv}$

---



**Fig. 3** Hyperparameters tuning on randomly selected 20 videos and the C3D model with our un-targeted SVA. (A) $\gamma_2$ tuning. (B) $\gamma_1$ tuning. (C) $\epsilon$ tuning.

halved. If we can't reduce the perturbation size $\epsilon$ after 20 times in a row, we cut the epsilon decay $\triangle_\epsilon$ in half.

There are three hyperparameters in our method, and we obtain their best values via parameter tuning experiments. The $\mathbb{P}$ in Eq.(8) is automatically updated in the algorithm, but we need to set its maximum value $\epsilon_{adv}$. In Eq.(9), we need tune $\gamma_1$ and $\gamma_2$. When $\gamma_1$ is tuning, we fix $\gamma_2$, and vice versa. The tuning results are given in Figure 3. From the figure, we see $\gamma_1 = 2$ and $\gamma_2 = 3$ are the reasonable choices. As for $\epsilon_{adv}$, we find that FR (the bigger, the better) and MAP (the lower, the better) reach a balance when $\epsilon = 0.05$. Therefore, we set the maximum adversarial perturbations magnitude to $\epsilon_{adv} = 0.05$ per frame, and this setting is also applied to other video attack algorithms in the experiment.

# 4 Experiments

In this section, two state-of-the-art video attack models are used to be compared with our proposed method on two video recognition models with two public datasets. We focus on the overall perturbations and the length of the frames selected in our experiments. Furthermore, a variant of our method is also designed as a comparison. A comprehensive evaluation of our method will be presented on those recognition models and datasets.

## 4.1 Datasets and Recognition Models

The datasets and recognition models used in the experiments are described in detail here.

**Datasets.** Two widely used datasets, UCF-101 [34] and HMDB-51 [21], are used in our experiments. UCF-101 is an action recognition dataset containing 13,320 videos with 101 action categories. HMDB-51 is a dataset for human motion recognition, which contains 51 action categories with a total of 7000 videos. Both datasets divide 70% of the video into training sets and 30% of the test sets. For convenience and fairness of comparison, 16-frame snippets evenly sampled from each video are used as input to the recognition models during the evaluation. The same approach was fist adopted in [13] and also used in [39]. In our experiments, we randomly sample 100 videos from the UCF-101 test dataset and 50 videos from the HMDB-51 test dataset for experimental comparison. It is worth noting that all the selected videos can be accurately classified by both recognition models.

**Recognition Models.** In our experiments, Long-term Recurrent Convolutional Networks (LRCN) [7] and C3D [13] are used as recognition models. The LRCN model uses a Recursive Neural Network to encode the Spatio-temporal features generated by CNNs. In our implementation, Inception V3 [35] is used to extract the features of video frames, and LSTM is used for video classification; The C3D model uses 3D convolution to learn Spatio-temporal features from video with Spatio-temporal filter for video classification. These models are all mainstream methods for video classification. In order to make the two video classification models better adapt to the input video, we train them with 16-frame snippets for 30 epochs via Adam in each dataset, and the learning rate is initialized as $10^{-5}$ and decreases to its 1/10 after 20 epochs. After training, both video classification models meet the requirements of our experiment. Table 1 summarizes the test accuracy of the recognition models with 16-frame snippets on the whole UCF101 and HMDB51 datasets.

**Table 1** The accuracy of the used video recognition models.

| Models | Datasets | |
|---|---|---|
| | UCF-101 | HMDB-51 |
| C3D | 85.88% | 59.57% |
| LRCN | 75.44% | 34.58% |

## 4.2 Evaluation

We select Fooling rate (FR), Query number (Q), Mean absolute perturbation (MAP), and Sparsity (S) as evaluation metrics to evaluate the performance of our method on various sides. Next, we will introduce and describe these metrics in detail.

**Fooling rate (FR).** The Fooling rate is defined as the percentage of attacks that successfully generate adversarial examples. In our experiments, it is influenced by two main factors: whether the adversarial video can successfully fool the video classifier and whether it is imperceptible. More concretely, with the upper limit of query numbers ($3 \times 10^4$ in the un-targeted mode and $6 \times 10^4$ in the targeted mode), if the adversarial video produced by the attack algorithm can make the classifier misclassify itself and the pixel mean perturbations of the whole video is less than the upper limit of perturbations (As described in subsection 3.4, the $\epsilon$ is set to 0.05 in our experiments), the adversarial video is considered successful. Fooling rate (FR) is a common metric for evaluating adversarial attacks, and a higher value means better performance.

**Query number (Q).** In the query-limited setting, To generate adversarial examples, fewer query times mean that the attack algorithm has a higher attack efficiency. Therefore, we measure the efficiency of the attack algorithm in the current dataset by using the average query number (Q) of all adversarial videos generated with the attack algorithm successfully.

**Mean absolute perturbation (MAP).** This metric is used to quantify the perceptibility of change in the video. Here, we follow some previous work [30, 29, 43] in getting an index (MAP) for an adversarial perturbation given by

$$MAP_i = \frac{\parallel x_{i,adv} - x_{i,orig} \parallel}{\mid pixel_{x_i} \mid}, \tag{12}$$

where $x_{i,adv}$ is the $i^{th}$ adversarial video, $x_{i,orig}$ is the $i^{th}$ original video, and $\mid pixel_{x_i} \mid$ is the total pixel number of this video. we normalize the $l_1$ norm of the video difference by the total number of pixels. In our experiments, we use the average MAP of all successful adversarial videos on one dataset to evaluate the attack algorithms, and a smaller value means better performance. Additionally, we have resized the value of MAP to 0-255 for the sake of clarity.

**Table 2** The results of SVAL on C3D with UCF-101 under different sparsity (S).

| | | S(%) | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | Metrics | 10 | 20 | 30 | 40 | 50 | 60 | 70 |
| Un-targeted | MAP | 5.5395 | 5.3805 | 5.3550 | - | **3.2895** | - | - |
| Attack | FR(%) | 100.0 | 100.00 | 100.00 | 80.0 | **100.00** | 80.0 | 60.0 |
| Targeted | MAP | 8.7538 | **6.6218** | - | - | - | - | - |
| Attack | FR(%) | 100.0 | **100.0** | 60.0 | 60.0 | 40.0 | 20.0 | 0.0 |

**Sparsity (S).** To show the sparsity effect of the attack algorithms, the sparsity (S), is used to represent the proportion of frames with no perturbations versus all frames in a specific video. It is defined as:

$$S_i = 1 - \frac{m_i}{T_i}, \tag{13}$$

where $m_i$ is the length of the key frames of the video $x_i$. A large sparsity value means that only a few frames are added the adversarial perturbations. In our experiments, we use the average $S$ of all successful adversarial videos on one dataset to show the sparse effect of the current attack algorithm.

### 4.3 Comparing Algorithms

We compare our Sparse Video Attack (SVA) method with Opt-attack [3] and Heuristic-attack [40]. For Opt-attack, it is originally proposed to attack image classification models under the black-box setting. The reason we select it as one competitor is that it can achieve smaller distortion compared with some other black-box attack algorithms. We directly extend Opt-attack to attack video models. For Heuristic-attack, it is also a time-domain sparse attack method like ours. We use the same parameter settings in their original papers and official implementation code, respectively.

Besides, one variant of our method, named SVAL, is joined to comparisons. For the SVAL algorithm, we replace the $R_{attack}$ reward in Eq.(9) with the reward $L_{percentage}$ that is used to limit the length of the key frames:

$$L_{percentage} = \| \frac{1}{T} \sum_{t=1}^{T} p_t + S - 1 \|, \tag{14}$$

where $S$ is the sparsity metric, a bigger $S$ value means the fewer frames will be selected. The definition of $p_t$ can be found in Eq.(3). The SVAL algorithm only needs some videos to train the agent and make it intelligent. So in different experiments, we first randomly select 50% videos to train agent, and then use the agent in our black-box attack method. We still optimize the policy function's parameter $\theta$ via Adam and update $\theta$ as: $\theta = \theta + ls \nabla_\theta (J_{\overline{R_{attack}}} + L_{percentage})$, where $J_{\overline{R_{attack}}}$ means

there is no reward $R_{attack}$. The epochs of training is set to 20, and the learning rate is initialized as $10^{-5}$ and decreases to its 1/10 after 15 epochs.

### 4.4 Performance Comparisons with SOTA methods

Because the parameter $S$ in SVAL needs to be set, a grid search method is used to select appropriate parameters with different experiments. We here only show the sparsity tuning for the C3D model with 10 randomly sampled videos from the UCF-101 dataset. The results are recorded in Table 2. The symbol "-" in Table 2 denotes that the agent cannot achieve the 100% fooling rate under the current sparsity. In this situation, the MAP cannot be computed across all the adversarial videos, and cannot compare with other sparsity fairly. Therefore, we use "-" to indicate it.

It is found that the fooling rate decreases with the rising of sparsity. In the un-targeted attack setting, when FR is 100%, the smallest MAP is 3.2895. Therefore, we set $S = 0.5$ in the following experiment. In the targeted attack setting, $S = 0.2$ is a good choice, so the setting is used in the following experiments. In the other sparsity settings, we use the same way to select the corresponding best results. As shown in Table 2, it can be found that attacking a part of the video frames is a feasible and effective way, which would significantly reduce the perturbations of the adversarial video.

The comparison results in the un-targeted setting are listed in Table 3 in different tasks. In each task, the best performance is emphasized with the bold number. As shown, SVAL and SVA have great advantages over other methods on the whole. On the MAP side, SVA ranks first in the 3/4 comparisons. The biggest gap between other algorithms and SVA occurs in the C3D model with the UCF-101 dataset, the MAP of SVA is only 2.4450, but others all exceed 3. Notice that there is no one case that our methods are not as effective as other methods. On the sparsity side, SVAL and SVA are all ahead of the others, the sparsity generated by them all exceeds 50%.

Usually, targeted attacks need more query numbers and perturbations than un-targeted attacks. For Opt-attack and Heuristic-attack methods, they don't successfully generate any adversarial video even after 60,000

**Table 3** The video attack results of four attack algorithms in the un-targeted mode. There are four metrics for measuring the effectiveness of the algorithms. For MAP, the smaller the value, the better. For S, the bigger the value, the better. For Q, the smaller the value, the better. For FR, the bigger the value, the better.

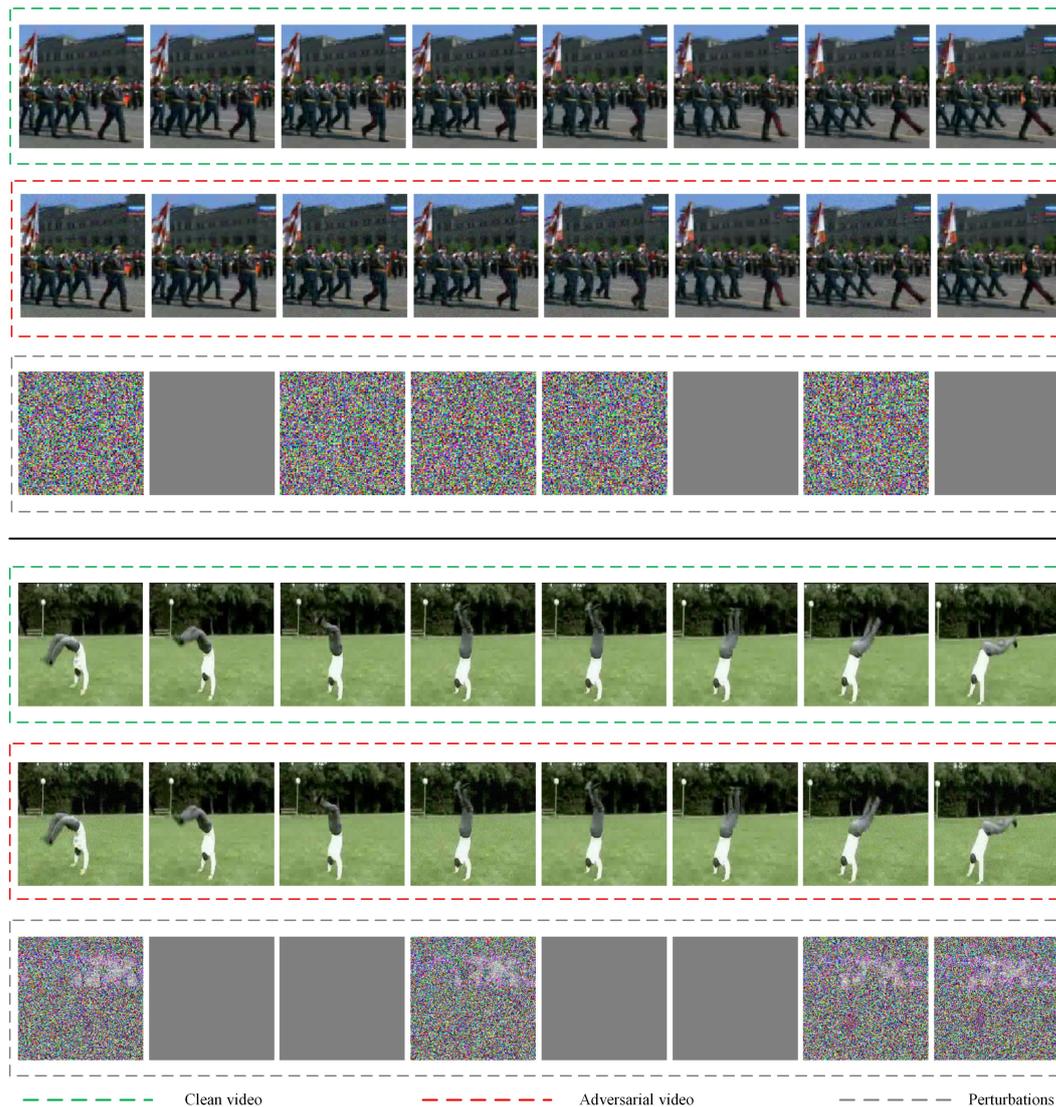| Dataset | Target Model | Attack Model | Metrics & Un-targeted Attack | | | |
|---|---|---|---|---|---|---|
| | | | MAP | S(%) | Q | FR(%) |
| UCF-101 | C3D | Opt-attack | 4.2540 | 0.00 | 15076.23 | 74.0 |
| | | Heuristic-attack | 3.2980 | 22.08 | 13609.91 | 79.0 |
| | | SVAL(ours) | 3.1765 | 50.00 | **8367.78** | 83.0 |
| | | SVA(ours) | **2.4450** | **63.14** | 9402.28 | **86.0** |
| | LRCN | Opt-attack | 2.8320 | 0.00 | 9032.68 | 57.0 |
| | | Heuristic-attack | 2.6940 | 17.19 | 9460.38 | 49.0 |
| | | SVAL(ours) | 2.4976 | 60.00 | **4131.57** | **68.0** |
| | | SVA(ours) | **2.396** | **62.14** | 6132.38 | 63.0 |
| HMDB-51 | C3D | Opt-attack | 2.8930 | 0.00 | 13274.14 | 76.0 |
| | | Heuristic-attack | 2.4960 | 25.68 | 11870.69 | 78.0 |
| | | SVAL(ours) | 2.4482 | **60.00** | 10727.93 | 94.0 |
| | | SVA(ours) | **2.3940** | 51.37 | 24948.67 | **98.0** |
| | LRCN | Opt-attack | 2.7586 | 0.00 | 18207.11 | 62.0 |
| | | Heuristic-attack | 2.6110 | 27.32 | 15663.41 | 66.0 |
| | | SVAL(ours) | **1.9479** | **70.00** | **10891.67** | **68.0** |
| | | SVA(ours) | 3.1570 | 62.50 | 18868.09 | 64.0 |

**Table 4** The video attack results of four attack algorithms in the targeted mode. There are four metrics for measuring the effectiveness of the algorithms. For MAP, the smaller the value, the better. For S, the bigger the value, the better. For Q, the smaller the value, the better. For FR, the bigger the value, the better.

| Dataset | Target Model | Attack Model | Metrics & Targeted Attack | | | |
|---|---|---|---|---|---|---|
| | | | MAP | S(%) | Q | FR(%) |
| UCF-101 | C3D | Opt-attack | - | - | > 60000 | - |
| | | Heuristic-attack | - | - | > 60000 | - |
| | | SVAL(ours) | 6.7672 | 20.00 | 43797.0 | **38.0** |
| | | SVA(ours) | **3.6450** | **57.24** | **36497.5** | 32.0 |
| | LRCN | Opt-attack | - | - | > 60000 | - |
| | | Heuristic-attack | - | - | > 60000 | - |
| | | SVAL(ours) | 5.8834 | 20.00 | **49065.3** | 39.0 |
| | | SVA(ours) | **3.270** | **56.64** | 57850.4 | **41.0** |
| HMDB-51 | C3D | Opt-attack | - | - | > 60000 | - |
| | | Heuristic-attack | - | - | > 60000 | - |
| | | SVAL(ours) | 6.9279 | 30.00 | 47190.3 | **40.0** |
| | | SVA(ours) | **3.8960** | **62.15** | 42900.3 | 38.0 |
| | LRCN | Opt-attack | - | - | > 60000 | - |
| | | Heuristic-attack | - | - | > 60000 | - |
| | | SVAL(ours) | 6.2861 | 20.00 | **43880.5** | 32.0 |
| | | SVA(ours) | **3.5170** | **66.77** | 47681.9 | **36.0** |

query times, which exceeds the experimental upper bound of query number pre-defined in section 3.4. Therefore we use "-" to represent their performance. The comparison results are recorded in Table 4. Obviously, the proposed methods are superior to other methods. The FR of our methods is at least 30% instead of 0% in the other competitive methods.

Two examples of the adversarial videos produced with our SVA un-targeted method are shown in Figure 4. For the first example (above the black line), the ground-truth label is "MilitaryParade", by adding the generated adversarial perturbations, the model tends to predict a wrong label "BandMarching". There are only 5 frames that have the perturbations in the whole 16 frames video. For the second example (below the black line), the ground-truth label is "flic flac", by adding the generated adversarial perturbations, the model tends to predict a wrong label "kick ball". The adversarial video has only 4 frames with perturbations, there are no perturbations on the other frames. Besides, two examples of the adversarial videos produced with our SVA targeted method are shown in Figure 5. The clean video, adversarial video, and the corresponding perturbations are shown in the green box, red box, and grey box respectively. For these two adversarial videos, the target classes are all "Archery". The adversarial video above the black line is from UCF-101 with the recognition model C3D, its original class is "RopeClimbing". There are 10 video frames that are polluted. The adversarial video below the black line is from UCF-101 with the

**Fig. 4** Two adversarial videos produced with SVA un-targeted attack. The clean video, adversarial video, and the corresponding perturbations are shown in the green box, red box, and grey box respectively. The perturbations have been rescaled into the range of 0-255. The adversarial video above the black line is from the UCF-101 dataset with the recognition model C3D, there are 8 frames distributed in video frames 4 to 8, 11, 13, and 14. There are only 5 frames that have the perturbations in the whole 16 frames video. The adversarial video below the black line is from the HMDB-51 dataset with the recognition model LRCN. The example is the front half of the adversarial video which has only 4 frames with perturbations, there are no perturbations on the other frames.
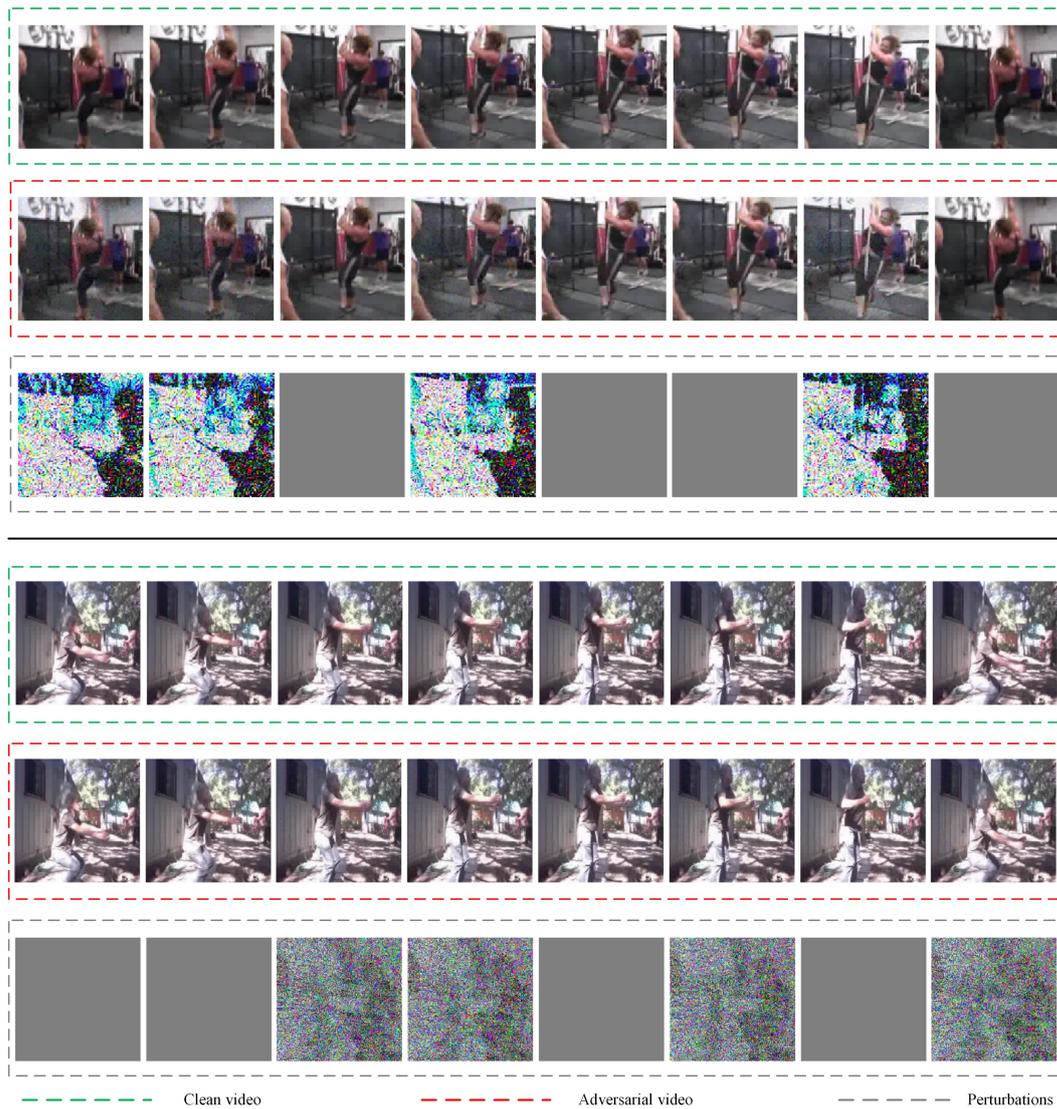
recognition model LRCN. Its original class is "Body-WeightSquats". The perturbations only exist in 8 video frames with a total of 16 frames. From these samples, we can conclude that the agent can select a small number of key and representative frames from the whole input video and the perturbations added on these key frames are human-imperceptible.

By observing and analysing all the results, we can draw the following conclusions: (1) Our SVA achieves the best performance in the majority of test tasks. (2) Attacking on key frames is an effective way to reduce perturbations of adversarial video. (3) Mutual guidance and cooperation between key frames selection and attacking is helpful to select the key frames for generating perturbations.

## 4.5 Ablation Study

We conduct a serial of ablation study experiments to analyze the proposed SVA in this subsection. Four experiments are conducted: the effectiveness of different reward functions, the influence of query numbers on ex-

**Fig. 5** Two examples of the adversarial videos produced with our targeted attack. The clean video, adversarial video, and the corresponding perturbations are shown in the green box, red box, and grey box respectively. The perturbations have been rescaled into the range of 0-255. We display the video frames corresponding to the odd number of their coordinates in two adversarial videos, so only 8 video frames are shown for each video. For these two adversarial videos, the target classes of them are all 'Archery'. The adversarial video above the black line is from UCF-101 with the recognition model C3D, its original class is 'RopeClimbing'. The MAP of this adversarial video is 5.4315, 62.5% of the total 16 video frames are polluted. The adversarial video below the black line is from UCF-101 with the recognition model LRCN. Its original class is 'BodyWeightSquats'. The MAP of this adversarial video is 1.479, the perturbations exist in 8 video frames with a total of 16 frames.

perimental results, the effectiveness of keyframe selection algorithms, and the influence of different gradient estimation algorithms on the final attack performance. We discuss the components in the following.
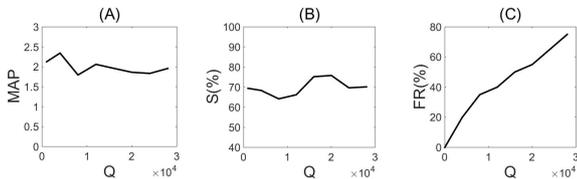
### 4.5.1 Different combination of reward functions

There are three rewards used to guide the agent in our experiments. In this subsection, the proposed method in the un-targeted setting on 20 randomly selected videos with the C3D model is used for the ablation study. The

ablation study for each reward function are given in Table 5, where "No RL" means that the results using FGSM+NES, and no RL module is used, "SVA$_{R_{attack}}$" means RL module is added, but only $R_{attack}$ is integrated, "SVA$_{R_{attack+rep}}$" means $R_{attack}$ and $R_{rep}$ are both integrated, and "SVA" means the full SVA model with $R_{attack}$, $R_{rep}$ and $R_{div}$. It can be shown that the attack reward $R_{attack}$ significantly improves the attacking performance. And the intrinsic rewards $R_{rep}$ and $R_{div}$ have relatively small contributions. The MAP and S in Table 5 are computed when FR meets the same

**Table 5** The ablation study of the proposed method SVA in an un-targeted setting.

| Metrics | Modules | | | |
|---------|---------|---------|--------------------|-----|
|         | No RL   | $\text{SVA}_{R_{attack}}$ | $\text{SVA}_{R_{attack+rep}}$ | SVA |
| MAP     | 6.5037  | 2.3723  | 2.0321             | 1.8624 |
| S(%)    | 0.00    | 62.35   | 68.75              | 74.65 |



**Fig. 6** The MAP, S, and FR indices change under the different query number (Q) on randomly selected 20 videos and the C3D model with our un-targeted SVA.

accuracy. It can be directly concluded that key frame selection is not only relative to the video itself, but also the feedback from the recognition model. The interactions between the attacking process and the keyframe selection can facilitate successful attacks.
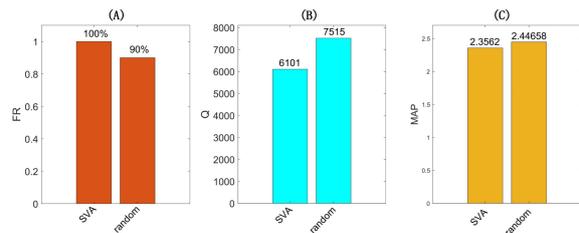
### 4.5.2 Sensitivity to query numbers

We investigate the changes of MAP, S, and FR versus different query numbers (Q). We also run a series of experiments on 20 randomly selected videos and the C3D model using our un-targeted SVA method. The results are shown in Figure 6. It can be found that FR is very relevant with the query number (Q) while MAP and S are relatively smooth versus the query times. It also indirectly proves that the results in our experiment have certain statistical significance and are relatively robust.

### 4.5.3 Effect of keyframe selection

To illustrate the effectiveness of our keyframe selection algorithm, we make a series of comparisons between our RL-based method and the random selection method. To conduct a fair comparison, we keep the other components consistent with the SVA and only replace the RL's agent with a random selection step. The number of selected frames is also the same between these two methods. Considering the randomness, each video is attacked five times by the random-based frame selection method, and the best attack result is recorded. The bar chart results of the comparisons are shown in Figure 7. Obviously, even the best attack result of the random algorithm is not the same as our SVA algorithm.
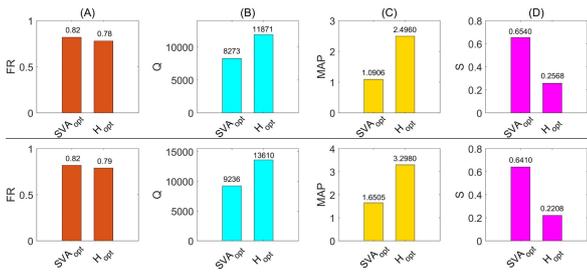
In addition, we also extend our keyframe selection module to the Heuristic-attack [40] algorithm to replace



**Fig. 7** The results of the proposed method SVA and random attack (random) on the C3D recognition model. The only difference between these two methods is the algorithm used to select keyframes. In the experiment, the random algorithm selected the same number of frames as the SVA algorithm, and we recorded the best results in five random attacks on 20 videos from UCF101 datasets. (A) the fooling rate of two algorithms (the bigger the value, the better). (B) the query number of two algorithms (the smaller the value, the better). (C) the MAP values of two algorithms (the smaller the value, the better). Obviously, even the best attack result of the random algorithm is not the same as the SVA algorithm.
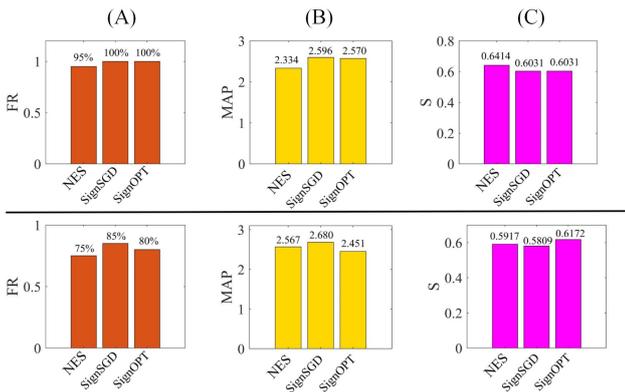
its frame selection part. The new algorithm ($\text{SVA}_{opt}$) and the Heuristic-attack algorithm ($\text{H}_{opt}$) have the same parameter settings, they are only different in keyframe selection. Here, we explore the performance of these two algorithms in terms of un-targeted attacks. We use the C3D recognition model, which has a wide application than other video recognition models. The results of these two methods are shown in Figure 8. Unsurprisingly, $\text{SVA}_{opt}$ has an overwhelming advantage over $\text{H}_{opt}$. For both datasets, HMDB-51 and UCF-101, $\text{SVA}_{opt}$ leads $\text{H}_{opt}$ in all evaluation metrics, which has less perturbation, fewer query numbers, higher sparsity, and higher fooling rate. We can conclude that the proposed keyframe selection algorithm is flexible and efficient.

### 4.5.4 Sensitivity to gradient estimation

Here, we make a series of comparisons between different gradient estimation methods, because the gradient estimation method is a key step in our SVA. We test additional gradient methods like SignSGD [46] and Sign-OPT [3] besides the NES method used in our method [47]. The attacking results of the proposed SVA with three different estimators on the C3D recognition model are shown in Figure 9. Here, we randomly select 20 videos from HMDB-51 and 20 videos from UCF-101, respectively. Each attack algorithm has the same upper query number limitation. From the observation, SVA with different gradient estimators will produce different attack results, but the differences are small. The big difference occurs in the fooling rate. On the UCF101 dataset, SVA with SignSGD and SVA with Sign-OPT obtain 10% and 5% better fooling rates than SVA with NES, respectively. For other evaluation metrics, those
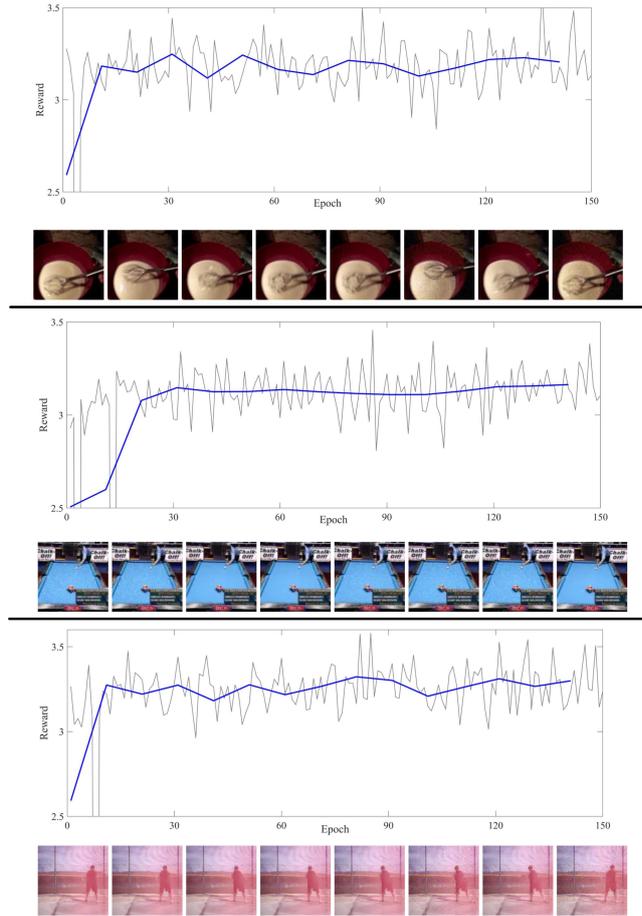
**Fig. 8** The un-targeted attack results of the Heuristic-attack ($H_{opt}$) and the variant algorithm $SVA_{opt}$ on the C3D recognition model. The proposed keyframe selection method is flexible and can integrate with Heuristic-attack on video attacking (aka $SVA_{opt}$). Here, we use datasets HMDB-51 and UCF-101, the adversarial results above the black line are from HMDB-51, and the adversarial results below the black line are from UCF-101. (A) The fooling rate (FR) of the two methods (the bigger the value, the better). (B) The query number of two methods (the smaller the value, the better). (C) The MAP values of two algorithms (the smaller the value, the better). (D) The sparsity of two algorithms (the bigger the value, the better).



**Fig. 9** The un-targeted attack results of the proposed SVA on the C3D recognition model with different gradient estimation methods. There are three different gradient estimation algorithms SignSGD, Sign-OPT, and NES. The only different setting between those video attacking methods is the gradient estimation method. Here, we randomly select 20 videos from HMDB-51 and 20 videos from UCF-101, respectively. The adversarial results above the black line are from HMDB-51, and the adversarial results below the black line are from UCF-101. (A) The fooling rate (FR) of the two methods (the bigger the value, the better). (B) The MAP values of two algorithms (the smaller the value, the better). (C) The sparsity of two algorithms (the bigger the value, the better). The max available query times are set to be consistent. It can be concluded from the observation that different gradient estimation methods do not produce significantly different results.
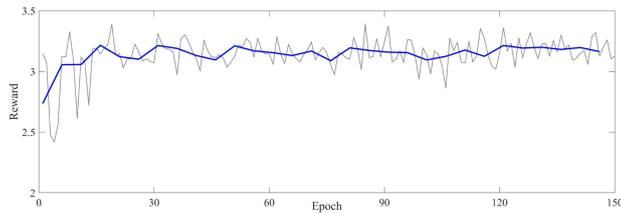
three methods have no obvious performance advantages. We can conclude that the proposed SVA is not sensitive to the gradient estimation method.



**Fig. 10** The change of reward value in the different epoch. Here, the C3D recognition model and three randomly selected videos from the UCF101 dataset are used to conduct the experiments. For clear observation, we record each reward value (gray lines) and the corresponding average reward value computed using every 10 epochs (blue lines). The corresponding video frames are also shown below the sub-figures.

## 4.6 Convergence of the proposed SVA

In this subsection, the convergence of the proposed SVA is discussed. We use the change of reward value (i.e., the value computed using Eq.(9)) in different epoch to test whether SVA can obtain a convergence during the learning, which is a widely used way to see the convergence in RL. We conduct experiments under the un-targeted attack setting. All experiments are based on the C3D model and UCF101 dataset. Due to the limitation of space, we show the change of reward value in the attack process of three separate videos. The reward results are recorded in Figure 10. To give a better description, we also show the corresponding video frames (Due to space constraints, only odd-numbered index frames are displayed) below the reward curve. As you can see from the figure, the reward value in each epoch

**Fig. 11** The change of average reward value computed on 20 randomly selected videos.

in SVA is choppy (gray lines), but the average reward value per 10 epochs is growing and eventually stabilizing (blue lines). A stable and high reward means that the frame selection strategy is optimal. At the beginning of the attack, attack failure is easy to occur, but this phenomenon will be better in the later stage. The convergence can be achieved after about 20 epochs for these three videos. Actually, this phenomenon is also suitable to other videos. We also give the change of average reward computed on 20 randomly selected video attacks. The result is given in Figure 11, where we can see the reward value becomes smooth and stable with the increasing epoch. It shows the good convergence of the proposed SVA.

## 5 Conclusion

In this paper, a sparse black-box adversarial video attack algorithm with reinforcement learning was proposed for video recognition models. Due to a large amount of temporal redundancy information of video data, we explored the sparsity of adversarial perturbations in the video frames through generating adversarial perturbations only on some key video frames. Considering that keyframe selection was not only relevant to the video itself but also the feedback from the recognition model, an agent based on attacking interaction and video intrinsic properties was designed for identifying key frames while attacking. As the perturbations were generated only for the selected frames, the proposed method could reduce the perturbations of adversarial examples significantly. The proposed algorithm was applicable to multiple target models and video datasets. Moreover, the experimental results demonstrated that the proposed algorithm achieved efficient query times to the recognition models. The most pertinent area of future work is to further investigate the black-box attack using fewer queries, such as modifying the update mechanism or designing new rewards. And we hope the query number can be reduced to meet the requirement in the real world.

## References

1. Naveed Akhtar and Ajmal Mian. Threat of adversarial attacks on deep learning in computer vision: A survey. *IEEE Access*, 6:14410–14430, 2018.
2. Avishek Joey Bose and Parham Aarabi. Adversarial attacks on face detectors using neural net based constrained optimization. pages 1–6, 2018.
3. Minhao Cheng, Thong Le, Pin-Yu Chen, Huan Zhang, Jinfeng Yi, and Cho-Jui Hsieh. Query-efficient hard-label black-box attack: An optimization-based approach. In *international conference on learning representations*, 2019.
4. Minhao Cheng, Simranjit Singh, Pin-Yu Chen, Sijia Liu, and Cho-Jui Hsieh. Sign-opt: A query-efficient hard-label adversarial attack. In *international conference on learning representations*, 2020.
5. Nilaksh Das, Madhuri Shanbhogue, Shangtse Chen, Fred Hohman, Siwei Li, Li Chen, Michael E Kounavis, and Duen Horng Chau. Shield: Fast, practical defense and vaccination for deep learning using jpeg compression. *knowledge discovery and data mining*, pages 196–204, 2018.
6. Lixi Deng, Jingjing Chen, Qianru Sun, Xiangnan He, Sheng Tang, Zhaoyan Ming, Yongdong Zhang, and Tat-Seng Chua. Mixed-dish recognition with contextual relation networks. *Proceedings of the 27th ACM International Conference on Multimedia*, 2019.
7. Jeff Donahue, Lisa Anne Hendricks, Marcus Rohrbach, Subhashini Venugopalan, Sergio Guadarrama, Kate Saenko, and Trevor Darrell. Long-term recurrent convolutional networks for visual recognition and description. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(4):677–691, 2017.
8. Wenkai Dong, Zhaoxiang Zhang, and Tieniu Tan. Attention-aware sampling via deep reinforcement learning for action recognition. In *national conference on artificial intelligence*, volume 33, pages 8247–8254, 2019.
9. Yinpeng Dong, Hang Su, Baoyuan Wu, Zhifeng Li, Wei Liu, Tong Zhang, and Jun Zhu. Efficient decision-based black-box adversarial attacks on face recognition. In *computer vision and pattern recognition*, pages 7714–7722, 2019.
10. Andre Esteva, Brett Kuprel, Roberto A. Novoa, Justin M. Ko, Susan M. Swetter, Helen M. Blau, and Sebastian Thrun. Dermatologist-level classification of skin cancer with deep neural networks. *Nature*, 542(7639):115–118, 2017.
11. Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *international conference on learning representations*, 2015.
12. Chuan Guo, Mayank Rana, Moustapha Cisse, and Laurens Van Der Maaten. Countering adversarial images using input transformations. *international conference on learning representations*, 2017.

13. Kensho Hara, Hirokatsu Kataoka, and Yutaka Satoh. Can spatiotemporal 3d cnns retrace the history of 2d cnns and imagenet? In *computer vision and pattern recognition*, pages 6546–6555, 2018.
14. Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *computer vision and pattern recognition*, pages 770–778, 2016.
15. Sepp Hochreiter and Jurgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997.
16. Andrew Ilyas, Logan Engstrom, Anish Athalye, Jessy Lin, Anish Athalye, Logan Engstrom, Andrew Ilyas, and Kevin Kwok. Black-box adversarial attacks with limited queries and information. In *international conference on machine learning*, 2018.
17. Xiaojun Jia, Xingxing Wei, and Xiaochun Cao. Identifying and resisting adversarial videos using temporal consistency. *arXiv preprint arXiv:1909.04837*, 2019.
18. Xiaojun Jia, Xingxing Wei, Xiaochun Cao, and Hassan Foroosh. Comdefend: An efficient image compression model to defend adversarial examples. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6084–6092, 2019.
19. Linxi Jiang, Xingjun Ma, Shaoxiang Chen, James Bailey, and Yu-Gang Jiang. Black-box adversarial attacks on video recognition models. In *acm multimedia*, pages 864–872, 2019.
20. Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *international conference on learning representations*, 2015.
21. Hildegard Kuehne, Hueihan Jhuang, Estíbaliz Garrote, Tomaso Poggio, and Thomas Serre. Hmdb: a large video database for human motion recognition. In *2011 International Conference on Computer Vision*, pages 2556–2563. IEEE, 2011.
22. Yann Lecun, Yoshua Bengio, and Geoffrey E Hinton. Deep learning. *Nature*, 521(7553):436–444, 2015.
23. Shasha Li, Ajaya Neupane, Sujoy Paul, Chengyu Song, Srikanth V Krishnamurthy, Amit K Roy-Chowdhury, and Ananthram Swami. Stealthy adversarial perturbations against real-time video classification systems. In *network and distributed system security symposium*, 2019.
24. Yuxi Li. Deep reinforcement learning: An overview. *arXiv: Learning*, 2017.
25. Geert J. S. Litjens, Thijs Kooi, Babak Ehteshami Bejnordi, Arnaud Arindra Adiyoso Setio, Francesco Ciompi, Mohsen Ghafoorian, Jeroen A. W. M. Van Der Laak, Bram Van Ginneken, and Clara I. Sánchez. A survey on deep learning in medical image analysis. *Medical Image Analysis*, 42:60–88, 2017.
26. Jiajun Lu, Hussein Sibai, Evan Fabry, and David Forsyth. No need to worry about adversarial examples in object detection in autonomous vehicles. *arXiv: Computer Vision and Pattern Recognition*, 2017.
27. Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *international conference on learning representations*, 2017.
28. Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin A. Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
29. Seyedmohsen Moosavidezfooli, Alhussein Fawzi, and Pascal Frossard. Deepfool: A simple and accurate method to fool deep neural networks. *computer vision and pattern recognition*, pages 2574–2582, 2016.
30. Omid Mohamad Nezami, Akshay Chaturvedi, Mark Dras, and Utpal Garain. Pick-object-attack: Type-specific adversarial attack for object detection. 2020.
31. Aaditya Prakash, Nick Moran, Solomon Garber, Antonella DiLillo, and James Storer. Deflecting adversarial attacks with pixel deflection. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8571–8580, 2018.
32. David Silver, Aja Huang, Christopher J. Maddison, Arthur Guez, Laurent Sifre, George van den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap, Madeleine Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. Mastering the game of go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016.
33. David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, Yutian Chen, Timothy P. Lillicrap, Fan Hui, Laurent Sifre, George van den Driessche, Thore Graepel, and Demis Hassabis. Mastering the game of go without human knowledge. *Nature*, 550(7676):354–359, 2017.
34. Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human actions classes from videos in the wild. *arXiv: Computer Vision and Pattern Recognition*, 2012.
35. Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. *computer vision and pattern recognition*, pages 2818–2826, 2016.
36. Florian Tramèr, Alexey Kurakin, Nicolas Papernot, Ian J. Goodfellow, Dan Boneh, and Patrick D. McDaniel. Ensemble adversarial training: Attacks and defenses. In *International Conference on Learning Representations*, 2018.
37. Xingxing Wei, Siyuan Liang, Ning Chen, and Xiaochun Cao. Transferable adversarial attacks for image and video object detection. *international joint conference on artificial intelligence*, pages 954–960, 2019.
38. Xingxing Wei, Jun Zhu, Sitong Feng, and Hang Su. Video-to-video translation with global temporal consistency. *Proceedings of the 26th ACM international conference on Multimedia*, 2018.
39. Xingxing Wei, Jun Zhu, Sha Yuan, and Hang Su. Sparse adversarial perturbations for videos. In *national conference on artificial intelligence*, volume 33, pages 8973–8980, 2019.
40. Zhipeng Wei, Jingjing Chen, Xingxing Wei, and Jiang Yugang. Heuristic black-box adversarial attacks on video recognition models. In *national conference on artificial intelligence*, 2020.
41. R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8(3-4):229–256, 1992.
42. Cihang Xie, Jianyu Wang, Zhishuai Zhang, Zhou Ren, and Alan L Yuille. Mitigating adversarial effects through randomization. *international conference on learning representations*.
43. Cihang Xie, Jianyu Wang, Zhishuai Zhang, Yuyin Zhou, Lingxi Xie, and Alan L Yuille. Adversarial examples for semantic segmentation and object detection. *international conference on computer vision*, pages 1378–1387, 2017.
44. Haichao Zhang and Jianyu Wang. Towards adversarially robust object detection. *international conference on computer vision*, pages 421–430, 2019.
45. Kaiyang Zhou, Yu Qiao, and Tao Xiang. Deep reinforcement learning for unsupervised video summarization with diversity-representativeness reward. In *national conference on artificial intelligence*, 2018.

46. Sijia Liu, Pin Yu Chen, Xiangyi Chen, and Mingyi Hong. Signsgd via zeroth-order oracle. In *7th International Conference on Learning Representations, ICLR 2019*, 2019.
47. Daan Wierstra, Tom Schaul, Tobias Glasmachers, Yi Sun, Jan Peters, and Jürgen Schmidhuber. Natural evolution strategies. *The Journal of Machine Learning Research*, 15(1):949–980, 2014.
48. Gaurav Goswami, Akshay Agarwal, Nalini Ratha, Richa Singh, and Mayank Vatsa. Detecting and mitigating adversarial perturbations for robust face recognition. In *International Journal of Computer Vision* 127.6 (2019): 719-742.
49. Francesco Croce, Jonas Rauber, and Matthias Hein. Scaling up the randomized gradient-free adversarial attack reveals overestimation of robustness using established attacks. In *International Journal of Computer Vision*, 2020, 128(4): 1028-1046.
50. Shangzhi Teng, Shiliang Zhang, Qingming Huang, and Nicu Sebe  Viewpoint and scale consistency reinforcement for UAV vehicle re-identification. *International Journal of Computer Vision* 129.3 (2021): 719-735.