

# Adversarial Attacks on SIFT Based Video Summarization techniques

Vireshwar Kumar  
Computer Science  
Indian Institute of Technology, Delhi  
viresh@cse.iitd.ac.in

Prateek Singh  
Maths and Computing  
Indian Institute of Technology, Delhi  
mt6180787@iitd.ac.in

Shreyansh Choudhary  
Maths and Computing  
Indian Institute of Technology, Delhi  
mt6180794@iitd.ac.in

**Abstract**—Recently our understanding of various concepts of computer vision gave rise to a wide variety of applications, be it face detection, movement detection and what not. But with the alarming threat of various adversarial attacks on such models, not only in computer vision but AI in general, a lot of research has been done in order to find out their weaknesses or try to make them more robust to previously known attacks. Below we have tried to find out how robust are state-of-the-art video summarizers vs various SIFT based adversarial attacks and find out how summaries are affected.

**Index Terms**—SIFT, Video Summarization, Adversarial Machine Learning

## I. INTRODUCTION

Video Summarization has been a topic of interest in computer vision ever since the success of various text summarizers. The idea of a video summarizer is similar, to summarise long videos into short and meaningful chunks. Given that a very large and high-quality video takes up a lot of space and a lot of bandwidth to transfer from one place to another, a viable solution is to shorten them sufficiently. Its application is very useful, especially in summarising hours of CCTV footage, where previously countless hours of footage was spent for review which is, however, inefficient and also, one can miss a lot of key details as one easily misses important information in a blink of an eye. In such cases, a shortened version of the same footage is very helpful as a good summarizer can ignore irrelevant information and only keep the relevant chunks. However, even the state-of-the-art applications are vulnerable to various adversarial attacks. In this paper, we would look into the robustness of one such state-of-the-art video summarization method. We tried various perturbations to a video before summarising it, and noted the results. Some perturbations gave bad summarised results, while one didn't affect at all.

## II. RELATED WORK

The most difficult challenge of video summarization is determining and separating the important content from the unimportant content. Most summarizers detect some important points(keypoints) from each frame, analysing their movements and then marking them important if they move significantly. As in Video Summarization technique proposed by Shruti Jadon and Mahmood Jasim [1] in their work. They use SIFT [7] and Image Histograms [8] to detect local features and then

analyse their movements in consecutive frames. If the change in pixel values is considerable, then that frame is considered important. For all the important frames, the summarizer then skims [9] some frames in their neighbourhood so that the final video is continuous.

Another approach, which is helpful for surveillance systems [3], separates foreground using background subtraction [4], dense optical flow [5], and gravitational clustering[6], associates detected foreground objects in consecutive frames through sliding temporal windows and then selects the keyframes.

## III. PROPOSED ATTACKS

We first tried to adjust the parameters of the video summarizer to work the best for CCTV footages. For this we asked people to find a summary for the videos, we have tested on and considered them as ground truth summary. We then tried to fine tune the parameters to get the best possible results from the summarizer. In this report, we will be using a white box attack to target the state-of-the-art video summarizer. We have adopted attacks on SIFT algorithm suggested in the paper by Zohaib et. al. The attacks fool the summarizer by perturbing the videos such that the Matcher provides incorrect matching of the keypoints. The following are the attacks mentioned in the work by Zohaib et. al.[2] on SIFT which worked fairly well.

### A. Attack-1: Blur

The first perturbation that we tried out was to blur the images at the corresponding detected key points. The idea was that if the SIFT detector would find it difficult to find the respective key points, then the corresponding matching would not be good. And as a consequence it would give us unsatisfactory results. The big problem with this attack was, although it did manage to affect the summary, the perturbations were visible to the naked eye. And hence this is not a great adversarial attack.

### B. Attack-2: Block-2-Block

The second attack was as Zohaib et. al.[2] named, Block-2-Block perturbation. This attack seemed much more successful as it disrupts the summary very well. In this, around every SIFT feature detected a patch (usually 3x3) is taken and blocks



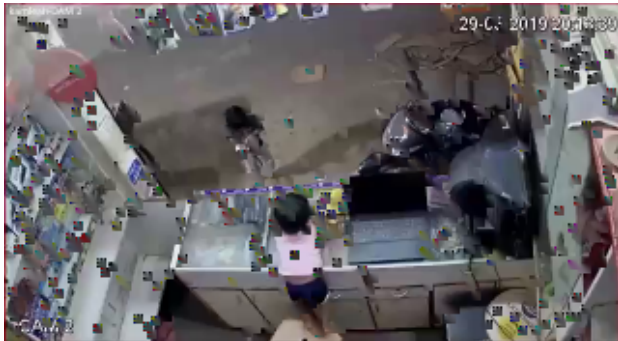
(a) Original



(b) Attack-1



(c) Attack-2



(d) Attack-3

Fig. 1: Perturbations Caused by the Attacks

are then perturbed according to a threshold. This threshold is a parameter that can be fine tuned to ones liking. Unlike the other two, this perturbation seemed impossible to detect from the naked eye.

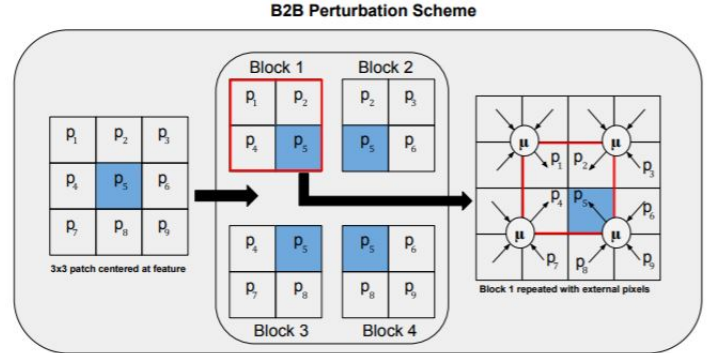


Fig. 2: Attack-2

### C. Attack-3: Pixel-2-Pixel

This was the last attack we tried, and this was not very successful. In this method, we would take the average of the left and upper pixel of the detected keypoint and replace it with the current pixel. This attack was not very effective, even as an adversary. The modified pixels were too obvious to the naked eye.

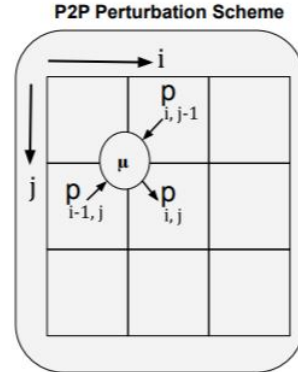


Fig. 3: Attack-3

## IV. RESULTS

The videos we used for testing were downloaded from Insecam[11] and trimmed to manageable sizes. We test the perturbed videos in 2 different ways. Firstly, it is important that the perturbations should not affect the video that can be easily perceived by human eye and secondly, the summarizer should fail to detect the critical frames from the required/ ground truth summary. For ground truth summary, we have asked people to provide the same for 10 small CCTV footages. Below we mention our results.

### A. Degree Of Perturbation

Here we have chosen SSIM(structural similarity) [10] as a metric for the similarity between the perturbed video frames and original frames. To present the results here, we have averaged them for all the videos.

Average SSIM scores	
Attack Type	SSIM Score
Attack-1	0.92
Attack-2	0.966
Attack-3	0.859

Here we can see that Attack-2(Block 2 Block) performs the best and the perturbations are not visible to human eye.

### B. Summarizer Performances

For this work, we have used SIFT based summarizer proposed by Shruti Jadon and Mahmood Jasim[2]. We mark the frames in the ground truth summary obtained from different people as critical and check how many of those are included by the summary of the original video and attack summary.

Summarizer performance on differently perturbed videos				
Attack Type	Original Video	Attack-1	Attack-2	Attack-3
1	0.8067	0.6992	0.7593	0.7222
2	0.99	0.972	0.962	0.921
3	0.941	0.90	0.856	0.901
4	0.8584	0.8159	0.8278	0.8398
5	0.9130	0.704	0.6393	0.5821
6	0.8144	0.7357	0.7663	0.7542
7	0.8983	0.8569	0.7974	0.8554

### FURTHER WORK

There is a lot of scope for developing security for these SIFT based attacks. SIFT is a widely used and robust technique, though it takes too much time per frame of the video. Hence, much work should be concentrated on the working time of both the summarizer and attacks. Some other attacks on this summarizer is using pepper noise on alternating frames to fool the motion threshold.

### REFERENCES

- [1] Shruti Jadon and Mahmood Jasim, "Unsupervised video summarization framework using keyframe extraction and video skimming".
- [2] Adversarial Examples for HandCrafted Features
- [3] Video Summarization of Surveillance Cameras
- [4] P. KaewTraKulPong and R. Bowden. An improved adaptive background mixture model for real-time tracking with shadow detection. In Video-based surveillance systems, pages 135–144. Springer, 2002
- [5] G. Farneback. Two-frame motion estimation based on polynomial expansion. In Image Analysis, pages 363–370. Springer, 2003
- [6] D. Dasguptal. A new gravitational clustering algorithm. In Proceedings of the Third SIAM International Conference on Data Mining, volume 112, page 83. SIAM, 2003
- [7] J.-M. Morel and G. Yu, "Is the scale invariant feature transform really scale invariant ?" 2010
- [8] Y. Wang, Q. Chen, and B. Zhang, "Image enhancement based on equal area dualistic sub-image histogram equalization

method," IEEE Transactions on Consumer Electronics, vol. 45, no. 1, pp. 68–75

- [9] Y. Gong and X. Liu, "Video summarization using singular value decomposition," in Proceedings IEEE Conference on Computer Vision and Pattern Recognition. CVPR 2000 (Cat. No.PR00662), vol. 2, 2000, pp. 174–180 vol.2

- [10] Structural Similarity

- [11] InSecam, Dataset for CCTV videos

- [12] Attentive and Adversarial Learning for Video Summarization

- [13] Video Summarization via Semantic Attended Networks

- [14] Quick Overview Of Adversarial Machine Learning

- [15] Evaluating and Improving Adversarial Robustness of Machine Learning-Based Network Intrusion Detectors

- [16] Sparse Adversarial Perturbations for Videos

- [17] Sparse Black-box Video Attack with Reinforcement Learning

- [18] A user attention model for video summarization

- [19] Video Summarization with Long Short-Term Memory

- [20] Discovering important people and objects for egocentric video summarization