# Video Image Detector: A Tool for Finding Similarity in Video Contents

Muhammad Imran Saeed[1,*], Intesab Hussain Sadhayo[2], Jawaid Shabbir[3], Nazar Hussain Phulpoto[4]

[1]Department of Computer Science, Nazeer Hussain University Karachi, Pakistan .
[2]Department of Telecom. Engineering, QUEST ,Nawabshah, Pakistan.
[3]Department of Computer Engineering, Sir Syed University of Engineering & Technology, Karachi, Pakistan.
[4]Department of Public Administration, Shah Abdul Latif University, Khairpur, Pakistan.
[*]Corresponding author: imran.saeed@nhu.edu.pk

## Abstract

Nowadays, when a sheer volume of multimedia data is being generated on daily basis, video piracy has become a genuine issue. In this paper, we propose a technique for matching video frames in two (or more) video files. Most of the work in this domain has been done on object detection, text detection, and spatio-temporal methods, however, the detection of copyright contents in videos has not been well-addressed. In this paper, we propose a technique to detect the copyright video frames in two or more videos. The given videos can be an advertisement or an especially worked-out video file by a journalist which is legally owned by the person who made it. Such a video files/clips can be matched with certain video streams or files to check if they contain the whole or a part of the given video file. The given video clip is composed of individual frames which could be matched on frame-to-frame basis with other (live) video streams to find the similarity extent between the successive images/frames. The method/technique to be proposed in this project will be mainly helpful for tracking or identifying the copyright digital video contents (e.g., songs, ads, news, etc) being played/transmitted illegally by a digital channel.

**Keywords**—Video similarity, content matching, feature matching

---◆---

## 1  Introduction

COPYRIGHT videos can be utilized maliciously by an association with no permission to the video's proprietor. The copyright material is more earnestly to recognize when it is duplicated, i.e., replicating a couple of casings from a video. The majority of past works in this context focus on object detection, text detection, spatio temporal methods etc. and so forth, however, to the best of our knowledge, there is no work on video-to-video content matching for piracy detection. Some of the video contents transmitted by the TV channels violate the copyright rules and do not properly acknowledge the videos' owners. In this regard, we propose a technique for automatic detection of the copyright contents in a video. Our proposed technique attempts to find the solution of the following questions /challenges.

- Finding a specific sequence of images in a video file.
- Deciding whether direct matching of two video files is possible.

- Determining the extent of color variation in a pixel to be accepted as identical.
- Determining the threshold/criteria based on the number of matched frames/images in a video to decide whether the contents are similar?

Our proposed technique has the ability to work with following types of videos.

1) Videos with same resolution for pixel-to-pixel comparisons
2) Videos with different frame rates, data rates and bit rates.
3) Videos with different resolutions with different data and bit rates.

In this technique, the movement calculation is utilized to satisfy the above parameters. The technique works well on small video contents, however, for larger videos, frames comparison requires huge storage.

Section II characterizes the related work in this domain. Section III gives an overview of video similarity detector. Section IV describes the motion detection algorithm. Section V depicts the diagram of the entire

application. Section VI characterizes the operational outline in which it finds the comparability substance between recordings. Section VII explains the tracking process for matching images. Section VIII describes the experimental results. Section IX concludes the paper.

## 2  Related Work

In [1], the authors propose a technique based on multi-frame end-to-end learning of image features and cross-frame motion. In some other techniques [2] [3], the authors present a programmed video subtitling model that joins spatio-temporal correlation and picture arrangement by neural network structures based on long short-term memory. The resulting system is demonstrated to produce state-of-the-art results in the standard YouTube captioning benchmark while also offering the advantage of localizing the visual concepts (subjects, verbs, objects), with no grounding supervision, over space and time. In [4], the authors address the issue of content based action recovery in video. Given a sentence portraying an action, the undertaking is to recover coordinating clasps from an untrimmed video. To capture the inherent structures present in both text and video, the authors present a multilevel model that coordinates vision and language. First, the authors inject text features early on when generating clip proposals to help eliminate unlikely clips and thus speeding up processing and boosting performance. Second, to learn a fine-grained similarity metric for retrieval, the authors use visual features to modulate the processing of query sentences at the word level in a recurrent neural network. In [5], the authors propose a content-based copy detection technique. This approach is based on the contents of media files. In this technique, the main focus is based on resolution, compression and digitization effects during detection of content based videos. In [6], the authors proposed a new motion signature. A different application of ordinal signature and experimental comparison of these methods to the color signature is proposed. This technique also matches content-based signatures to detect copies of videos as opposed to watermarking, which relies on inserting a distinct pattern into the video stream. In the end, the statistical features from this technique indicate that it has an impressive performance.

In [7], authors propose another copy detection technique which detects the key frames by using color histograms. This technique basically relies on color; however, dissimilarities in color are expected to be reasonable complications in this approach. In [8], the authors present a comparative study of background subtraction strategies. Methodologies extending from straightforward foundation subtraction with global thresholding to increasingly complex measurable techniques are implemented and tested on various videos with a ground truth. The objective of this investigation is to provide a strong systematic ground to highlight the qualities and shortcomings of the most widely used movement discovery strategies. The techniques are contrasted based on their vigor with various kinds of video, their memory requirement, and the computational exertion they require. In [9], the authors propose a new algorithm for motion detection. In this proposed scheme, the moving object is detected by using a stationary camera within a scene. Another successive result is to compare the frames on both sides with the calculation of $n$ consecutive frames. It finds out the percentage area in which the motion exists. In [10], the authors discuss the visual surveillance integration system that achieves better performance with respect to visual tracking in motion detection. However, the information and motion of tracking algorithm is combined into an appearance model and is used as a particle filter framework for tracking the object in subsequent frames. In [11] [12], the authors propose a simple recursive nonlinear operator, used along with a spatial temporal regularization algorithm. By using a static camera, these motions are performed by approximating the fixed part of the videos. The extensive range of motion is detected in a complex scene with different time constants.

Most of the work in this domain is done on background subtraction calculation which is focused on distinguishing an item in one frame to another frame. During the video playback, it is hard to identify moving objects from one video to another video. In [13], the authors present an algorithm for identifying moving objects from a static scene based on frame difference. Firstly, the first frame is captured through the static camera and the successive frames are captured at customary interims. Secondly, the absolute difference is calculated between the consecutive frames and the difference image is stored in the system. In [14], the authors present a survey on the most recent strategies for moving object detection in video sequences captured by a moving camera.

## 3  Video Similarity Detector

Video identification is a technique that discovers those edges that are coordinated with one video into another video. This technique contrasts with the current motion picture outlines and the prior casings or with

something that will be called as foundation. If an object in the frame is moving slickly, lesser variation is obtained by the smaller predefined threshold. This calculation additionally matches both the video outlines one by one amid the running video streams. Amid correlation of recordings, the entire moving edge is identified autonomously by its movement speed.

## 4  Motion Detection Algorithm

Motion detection [9] is the primary procedure in the abstraction of data concerning moving entities maintaining in efficient regions such as tracking, cataloging, acknowledging, etc. The background of video frames are calculated by taking means of $n$ successive frames and matching them with the existing frames using the sub blocks of the matching-based scheme.

BS techniques take the notion to experiential video sequence. Image $I$ is made up of a static background $B$ in front of which stirring entities are observed. By the observation of every moving object that is made up of color distribution different from the one in $B$, BS methods can be applied by the following formula,

$$X_{t(s)} = \begin{cases} 1, & \text{if } d(I_{(s,t)}, B_s) > \tau \\ 0, & \text{otherwise} \end{cases} \tag{1}$$

where $\tau$ is a threshold, $X_{t(s)}$ is the motion label field at time $t$, $d$ is the distance between $I_{(s,t)}$ and pixels, and $B_s$ is the background model at pixels.

The reasonable way to model the background $B$ is to conclude by a single gray scale color image void of moving objects [10]. By the instruction of handling with brightness changes and background adjustments, it can be iteratively updated as follows,

$$B(s, t+1) = (1-\alpha)B(s,t) + \alpha I(s,t) \tag{2}$$

where $\alpha$ is a constant whose value ranges between 0 and 1. In the case of mean method, background is the mean or average of the earlier frames and mathematically it is written by the following formulae.

$$B(x,y,t) = \frac{1}{n} \sum_{(i=0)}^{(n-1)} (x,y,t-1) \tag{3}$$

$$|(x,y,t) - B(x,y,t)| > Th \tag{4}$$

The background model is subtracted from the $n$ previous or existing frames. The threshold value used in this technique checks if the value of pixel is greater than the other pixel, in which case it is treated as a foreground pixel. Otherwise, if a pixel's value is smaller than the threshold value, it becomes a background pixel.

## 5  Overview of the Application

A given video clip which is required to be monitored in another video can be a notice or a particularly worked-out video document by a columnist which is legitimately claimed by the individual who made it. Such a video records/clasps can be coordinated with certain video streams to check in the event that they contain the entire or a piece of the given video. The recordings might be songs and ads, etc.

We initialize both the videos and start the tracking process. If the frames are found, they are stored into the resultant directory. This comparison will run till it completes the last $n^{th}$ frame to be matched with the whole video. If the initial frame is not matched, the recognition procedure analyzes the whole moving frame independent of its motion speed. Frames and background are intended to be calculated by taking mean of $n$ successive frames and comparing them with the existing frame. The subsequent frames use the threshhold values. These qualities demonstrate the variation of colors which is increased up to 10%, whereas, if the frames are not compared during detection process, then the comparison process itself continues and checks the next frame to compare with the previous or current frame. Figure 1 gives the logical view of this application.

## 6  Finding Similarity in Video Contents

To determine the similarity between two or more videos, some small or large videos are collected through different media programs. First video is contrasted and the second video in the given time span can be chosen from the client's decision. Figure 2 depicts the overall framework of our application.

Stage 1: Choose two videos that have diverse time frames and begin correlation.

Stage 2: Two directories are created automatically before the start of the comparison process. First directory is made when the video document is stacked into the video stream. This directory is named as "Resulting Matched Frame Folder".

Stage 3: The second directory is created and named as "Advertisement Image". The frames are separated and stored into this directory. We wait until the application completes its procedure.

Stage 4: In Figure 2, two cases are conceivable during the video matching phase. The first frame of the advertisement video clip is matched with the first frame of the video file. If the frame is matched, it is stored into
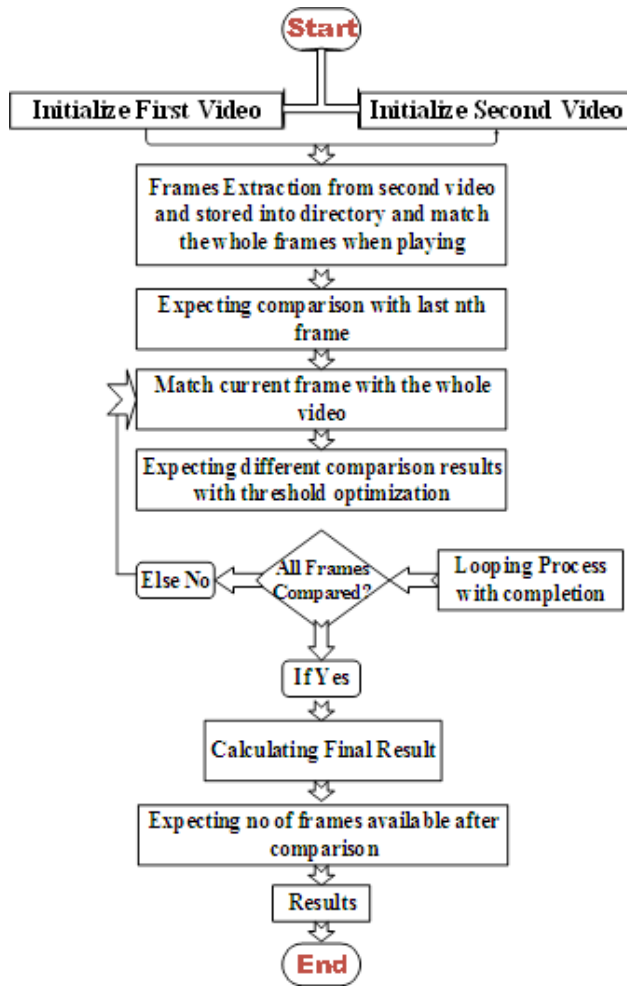
Fig. 1: Logical view of this application (operational model)

the directory "Resulting Matched Frame Folder".

Step 5: After consummation, the coordinated casing will likewise check amid or after the correlation.

Step 6: After completion, the matched frame will also check during or after the comparison. This application uses the motion detection technique in which the examination relies upon the shading variety and blend of different calculations in which the comparison relies upon the variety of changing nature between videos and blend of different calculations.

## 7    Tracking Process for Matching Images

The accompanying procedure shows the tracking frames between two recorded videos. Figure 3 demonstrates the tracking process while the threshold values
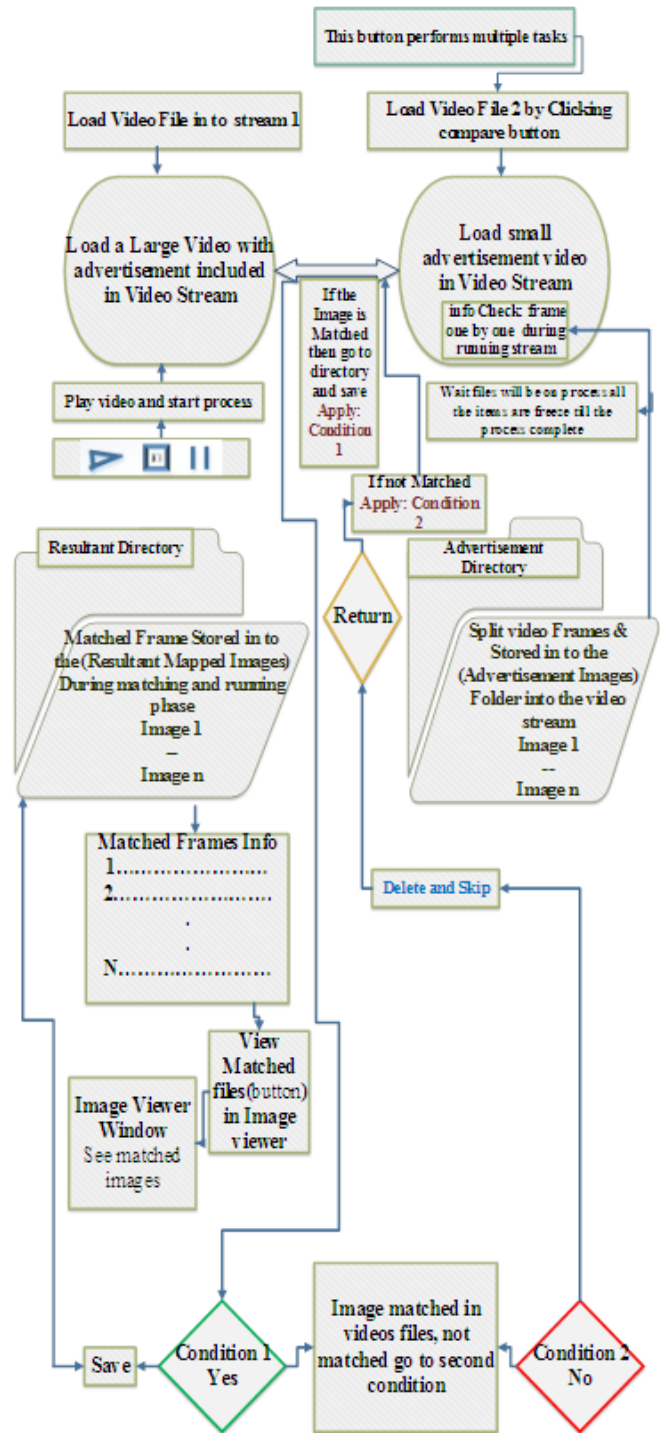


Fig. 2: Flow Chart for finding similarity in video contents

are compared with all the values of the running videos. In figure 3, the center box shows the threshold values and by using these parameters, following cases are tested. If the condition is more prominent than the threshold values, it can not distinguish the frames of the video, the values of comparable videos are less than or equal to thresh hold values or the application detects the related frames in both the videos, if and only if the related video frames are available in both the recorded videos.

This following procedure additionally utilizes the limit esteems, yet it likewise utilizes the framework's current date and time that recognizes the related frame. Without the utilization of framework date and time, the following procedure does not track the frames correctly. In this procedure, the frames are isolated into milliseconds, for instance, one edge is extracted into thirty frames for each second and th

The comparison frames are detected in 0 to 5 milliseconds. It detects the frames one by one, next frames are detected between 5 to 10 milliseconds and so on. By utilizing the present framework's date and time with the video length, it distinguishes the edges in milliseconds.

## 8 Experimental Results

The proposed procedure is feasible to detect the copyright frames by utilizing diverse clasps/video. Examination is done on various kinds of situations.

- Videos with same resolution
- Videos with different frame rates, data rates and bit rates.
- Videos with different resolutions with different data and bit rates.

### 8.1 Videos With Same Resolution

We pick two recorded videos that have distinctive time periods and begin comparison. Two directories are made on the beginning of the application. First directory is made amid the video document is stacked into the video stream. This Directory is named as "Resulting Matched Frame Folder". Second directory is made and named as "Advertisement Image". Figure 4 demonstrates the correlation procedure where two cases are conceivable amid the video coordinating stage. The primary frame of the commercial video cut is matched with the first frame of the recorded video. If the frame is matched, it is stored into the directory. Then again on the off chance that the frame isn't matched, the value returns and afterward goes to the next frame to compare with the entire video.
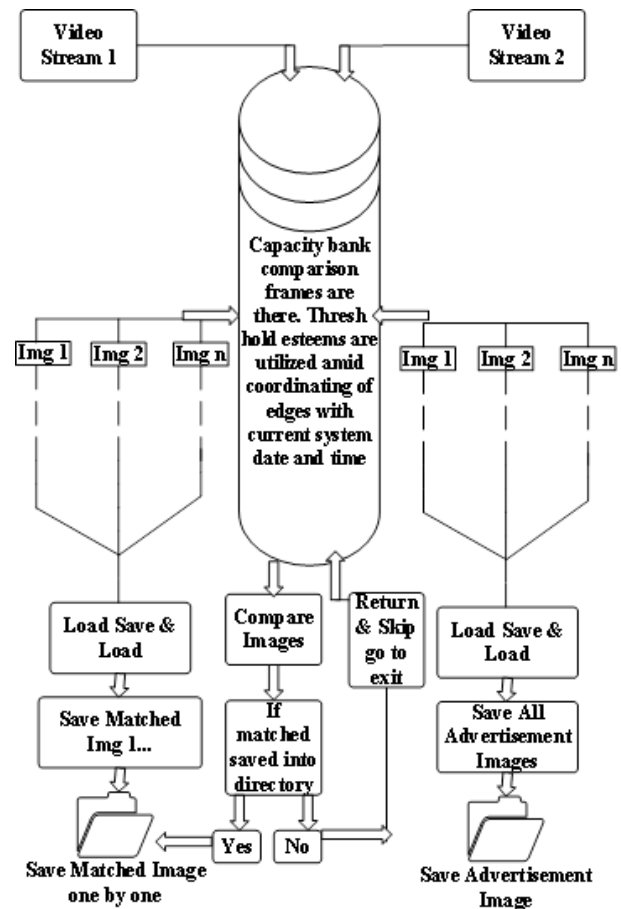


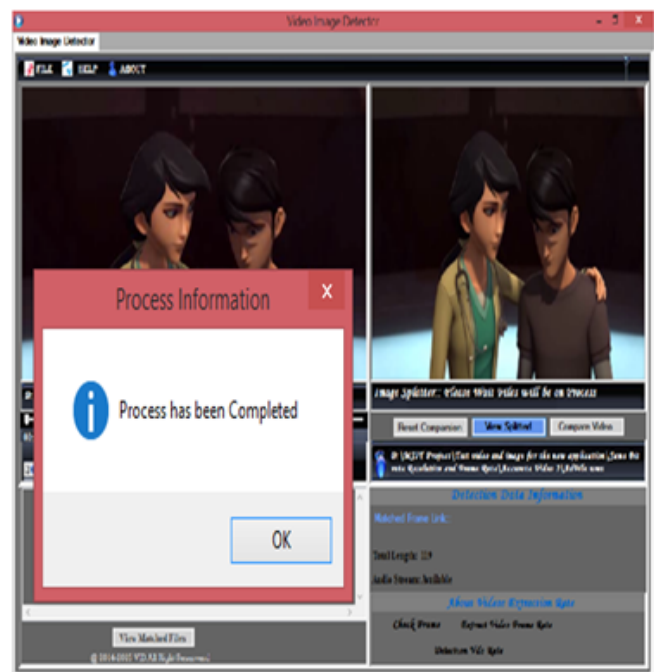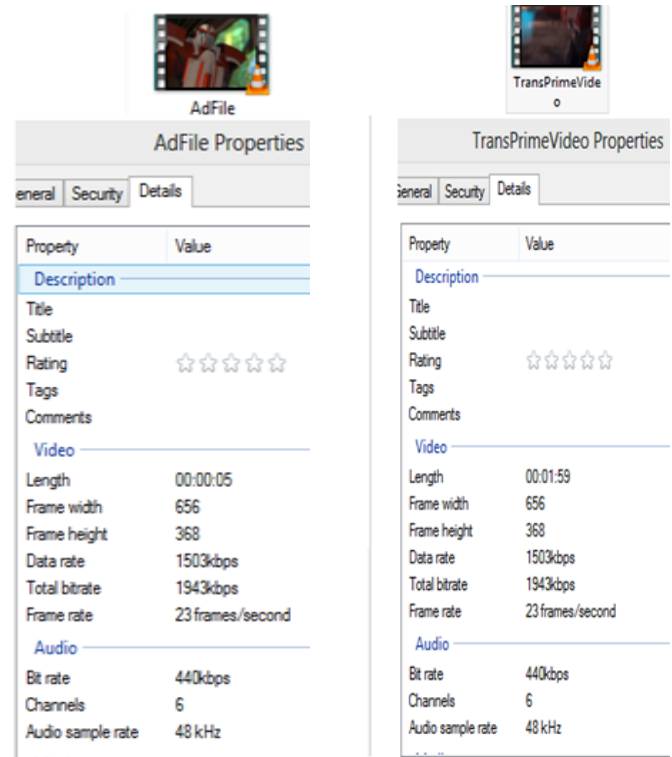Fig. 3: Tracking Process by using threshold



Fig. 4: Extracting video clips

Fig. 5: Comparison process



(a) Advertisement video



(b) Full video with advertisement clip

Fig. 6: Videos with similar resolution

Table 1 shows the original video data. Note that the resolution, data, frame rate and total bit rate of both the videos have the same data but their duration times are different from one another. Figure 5 demonstrates the correlation procedure in which the frames of a video files are coordinated with each other. The data is stored with frame numbers and is made available in the application. It can be seen in Figure 6(a) and 6(b) that the advertisement video clip are compared and the length, frame width, height and other related characteristics of both the videos are matched. During the comparison of both videos, the frames are detected and the parameters in this figure show that the contents of the video are copied. Figure 7 focuses on both the directories of the matched video frame. After comparison, the outcome among Adframe_0 and the matched frame TransPrimeVideo_0 has same

properties. Just a single frame of both the directories is pronounced to exhibit the properties after matching. Table 2 shows the results of the comparison. Out of 165 frames, the number of matched frames are 14. The resolution of both the videos is 720x480. The frame rates of both the videos is 23 frames/second, and their bit rate are 1943 kbps.
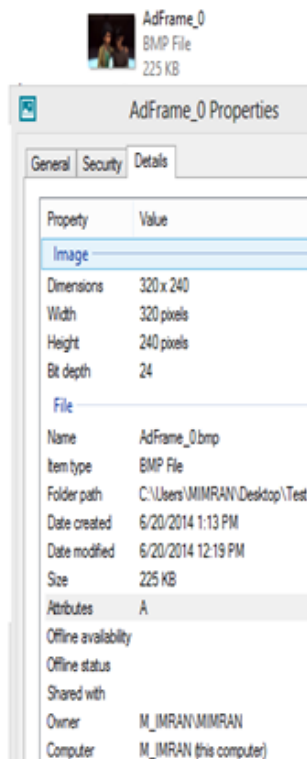
## 8.2 Video With Different Frame, Data & Bit Rates

We also test different resolution, frame rates and bit rates. The threshold values must be higher or like the perfect qualities. In the case when the threshold is higher than the other values during video comparison, the resultant is not identified and the procedure will proceed until the next frame is matched. Figure 8 shows video comparison of different length, date rate, and bit rate. Figure 9 delineates that these edges are distinguished effectively. Figure 10 shows that these frames are coordinated by utilizing diverse data rate and bit rate. Hence, our proposed technique distinguishes the video frames with different date rate and bit rate. Figure 10 also shows the matched (copyright) frames of the videos. Table 3 shows that out of 249 frames, the numbers of matched frames are 9. The

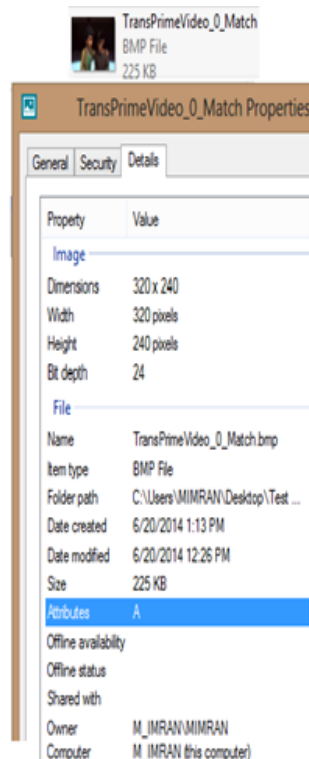| Original Video Data | | | |
|---|---|---|---|
| Video 1 (Advert. video clip) | | Video 2 ( Normal video with Advert) | |
| Original Video | | Original Video | |
| Video Resolution | 720 * 480 | Video Resolution | 720 * 480 |
| Video Data rates | 1503 kbps | Video Data rates | 1503 kbps |
| Video Frame rates | 23 frames /sec | Video Frame rates | 23 frames /sec |
| Total bit rates | 1943 Kbps | Total bit rates | 1943 Kbps |

TABLE 1: Video resolution, data, frame and total bit rate

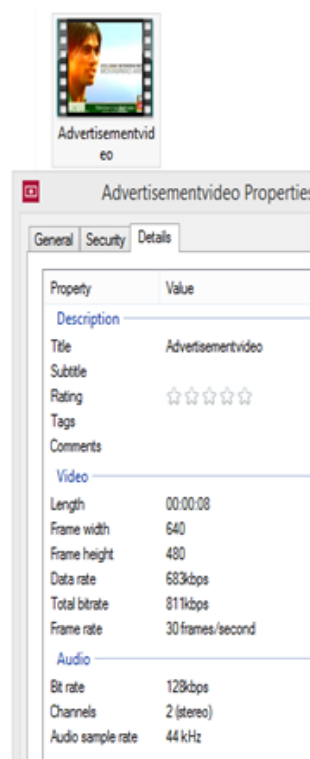| Matched Image Data b/w two videos | | | |
|---|---|---|---|
| Advertisement video clip | | Normal video with Advert | |
| Experimental Results | | | |
| No of original matching frames | 165 | No of frames matched | 14 |
| Video Resolution | 720 * 480 | Video Resolution | 720 * 480 |
| Video Data rates | 1503 kbps | Video Data rates | 1503 kbps |
| Video Frame rates | 23 frames /sec | Video Frame rates | 23 frames /sec |
| Total bit rates | 1943 Kbps | Total bit rates | 1943 Kbps |

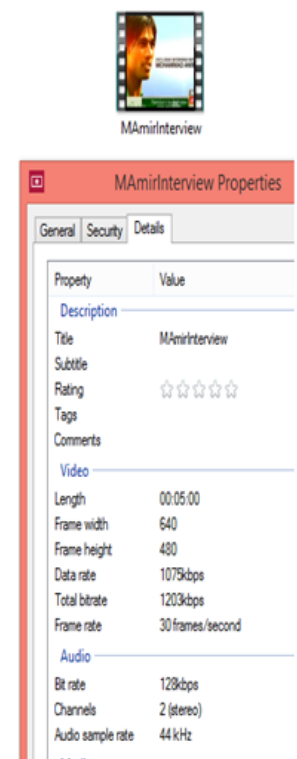TABLE 2: Total 165 frames 80 frames are similar, No of matched frames are 4.



(a) Advertisement video

(b) Full video with advertisement clip

Fig. 7: Advertisement & Matched frame properties



(a) Advertisement video

(b) Full video with advertisement clip

Fig. 8: Video Frame detection rate

Fig. 9: Frames are detected & displayed in the directory



Fig. 10: Resultant Frames

resolution of both the videos is 640x480. The frame rate of both the videos is 30 frames/second. The bit rates of the two videos are 811 kbps and 1203 kbps, respectively.

## 8.3 Video With Different Resolutions, Data Rate & Bit Rates
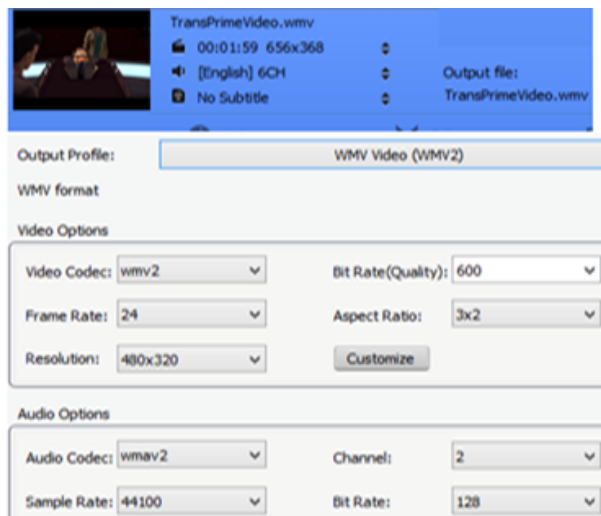
In the third scenario, the experiment is done on different resolutions, data and bit rates. The outcome appeared in Figure 11(a) and 11(b) demonstrates that the two videos are different and there bit rates are entirely different with one another. The comparison between the two videos begins from different parameters mentioned above. In this figure, the conceivable results may be 8-10 frames.
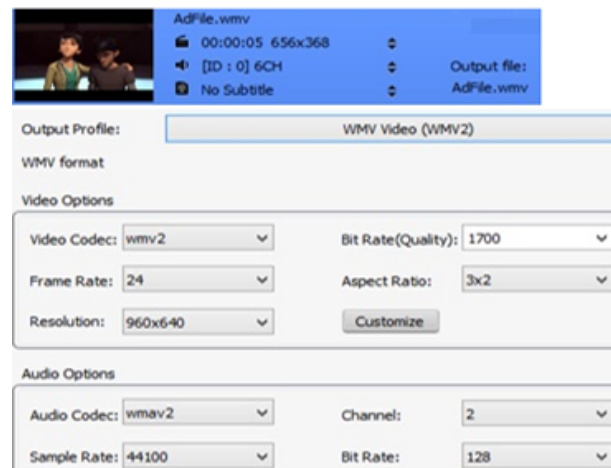
## References

[1] Zhu, Xizhou, Jifeng Dai, Lu Yuan, and Yichen Wei. "Towards high performance video object detection.", In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 7210–7218. 2018.

[2] Zanfir, Mihai, Elisabeta Marinoiu, and Cristian Sminchisescu. "Spatio-temporal attention models for grounded video captioning.", In asian conference on computer vision, pp. 104–119. Springer, Cham, 2016.

[3] Guo, Dashan, Wei Li, and Xiangzhong Fang. "Capturing temporal structures for video captioning by spatio-temporal contexts and channel attention mechanism.", Neural Processing Letters 46, no. 1 (2017): 313-328.

[4] Xu, Huijuan & He, Kun & Sigal, Leonid & Sclaroff, Stan & Saenko, Kate."Text-to-Clip Video Retrieval with Early Fusion and Re-Captioning.", ArXiv, 2018.

[5] Sand, Peter, and Seth Teller. "Video matching.", In ACM Transactions on Graphics (TOG), vol. 23, no. 3, pp. 592-599. 2004.

[6] Hampapur, Arun, and Ruud M. Bolle. "Comparison Of Distance Measures For Video Copy Detection.", In ICME, pp. 737-740. 2001.

[7] Hampapur, Arun, Kiho Hyun, and Ruud M. Bolle. "Comparison of sequence matching techniques for video copy detection.", In Storage and Retrieval for Media Databases, vol. 4676, pp. 194-202. International Society for Optics and Photonics, 2001.

[8] Benezeth, Yannick, Pierre-Marc Jodoin, Bruno Emile, HÃ¯lÃ¨ne Laurent, and Christophe Rosenberger. "Comparative study of background subtraction algorithms.", Journal of Electronic Imaging, vol. 19, no. 3, 2010.

[9] Kohli, Kamna, and Jatinder Pal Singh. "Motion Detection Algorithm.", International Journal of Computer Science & Applications (TIJCSA), vol. 1, no. 12, pp.1-5, 2013.

[10] Abdelkader, Mohamed F., Rama Chellappa, Qinfen Zheng, and Alex L. Chan. "IEEE International Conference on, Integrated motion detection and tracking for visual surveillance.", pp. 28-28, 2006.

[11] Manzanera, Antoine, and Julien C. Richefeu. "A new motion detection Algorith based on ÎčâĂŞÎ̂Ť background estimation.", Pattern Recognition Letters, vol. 28, no. 3, pp.320- 328, 2007.

[12] Deori, Barga, and Dalton Meitei Thounaojam. "A survey on moving objects Tracking in video.", International Journal on Information Theory (IJIT), vol. 3, no. 3, pp.1-16, 2014.

[13] Singla, Nishu. "Motion detection based on frame difference method."International Journal of Information & Computation Technology, vol. 4, no. 15, pp. 1559–1565, 2014.

[14] Yazdi, Mehran, and Thierry Bouwmans. "New trends on moving object detection in video images captured by a moving camera: A survey.", Computer Science Review, vol. 28, pp. 157–177, 2018.

| Original video data | | | |
|---|---|---|---|
| Video 1 (Advertisement video clip) | | Video 2 ( Normal video with Advert) | |
| Original Video | | Original Video | |
| Video Resolution | 640 * 480 | Video Resolution | 640 * 480 |
| Video Data rates | 683 kbps | Video Data rates | 1079 kbps |
| Video Frame rates | 30 frames /sec | Video Frame rates | 23 frames /sec |
| Total bit rates | 811 Kbps | Total bit rates | 1203 Kbps |

TABLE 3: Total 249 frames, 73 frames are similar, No of matched frms 9



(a) Different Resolution, Frame rate, bit rate of a Full Video



(b) Full video with advertisement clip

Fig. 11: Video with different resolutions, data rate & bit rates