# Multi-Modal Video Summarization

## P.Geetha[1]*, Vasumathi Narayanan[2]

[1]*Anna University, Chennai, India*
[2]*St.Joseph's College of Engineering, India*
\* Corresponding author's Email: geethap@annauniv.edu

**Abstract:**   Most previous work on summarization targeted either only on video or on its description of single video document. To efficiently manage huge video documents, this paper addresses a new prototype for summarizing video documents that automatically processes text and key-frames from the given video documents, indexes and provides mechanism for search and retrieval. The system separately summarizes the text and video related information. Both summarizations will enhance the performance as well the quality of the video document retrieval processes. The text summarizer provides a summarized view of a large video document to improve the learning speed and to give a glimpse of the document for better understanding. The video summarization consists of the key frames that effectively express the content and context of the video. This is achieved by quantifying visual attention of each frame in a video. This framework provides meaningful two-way summarization for video documents.

**Keywords:**   WAN; SAN; Text Summarization, Visual Attention Index; Video Adaptive Summary; Color -Motion Extraction.

## 1. Introduction

Most web videos and video articles that consist of related information are frequently used by users. Video summarization is an emerging task among the researchers. It becomes the main source for efficient browsing, access and manipulation of large video document collections. The previous work is mainly on the summarization of single video based on various features such as motion [1, 2], audio [3] or multi-modality [4]. A survey on this work is given clearly [5, 6].

Video document summarization is a process of summarizing a video using attention models [7, 9] and its descriptions by producing a summary of salient key frames that could convey the overall content of the video and its textual information is retrieved based on ranking and presented to the viewer [10]. Because of this two-way summary-zation, a viewer could be able to understand the content of the video without watching the video fully.

A video summary conveys the actual content of the video using a few most representative frames. Lots of research is going on in this field and many approaches are proposed to extract potential key frames from the video and to summarize them.

Video Document summarization involves the three basic steps.
(1) Extracting key frames from the video
(2) Summarizing video by selecting representative key frames from the extracted potential key frames.
(3) Text summarization include text pre-processing, association network building, textual unit weighting, sentence extraction and summary generation.

Key frame extraction refers to finding a set of salient images taken from a full length video that represent the visual content in the video efficiently. The key frames extracted should be content based and the process of extraction should be automatic. Previously key frames were extracted randomly from the video to summarize it. But randomly extracted frames will not represent the content of the video and information in the video efficiently. There are various approaches towards extraction of key frames.

[11] uses an algorithm to find a temporarily maximum occurrence frame (TMOF) in each shot. The key frames are extracted from the peak of the distance curve in [11]. [12] addresses this problem of lacking temporal consistency in the color based approaches of video summarization by combining both the color and motion based approaches. By doing so, more meaningful summarization can be achieved. [12] shows that summarization based on local estimation does not produce an image map that meaningfully summarizes the video content. The result is a summarizing map where the dominant object appears clearly and rest of the objects appears blurred. [13] proposes another motion based key frame extraction technique. It quantifies motion involved in each video segment and accordingly decides the number of key frames required to represent the video in the summary. [14] presents an idea to capture visual content and the dynamism in the video as well. It uses color based key frame extraction technique but takes account of activities in the video for determining the number of key frames to be represented in the summary. [15] extracts key frames based on a three-step algorithm. The three major steps are preprocessing, temporal filtering and post processing. [16] proposes a framework for summarizing story oriented videos based on narrative based abstraction. [17] proposed a method to obtain summary by considering degree of progress in the video.

Text summarization is the process of creating automatically a concise version of a text document carrying the main information of the original document. [18] uses a new term weighting approach based on Words Association Network (WAN) / Sentence Association Network (SAN) and used for text summarization.

[19] introduces a stochastic graph-based method for computing relative importance of textual unit for natural Language processing. [20] describes a new approach for estimating term weights in a document and shows how the new weighting scheme can be used to improve the accuracy of a text classifier.

The main objective of this proposed approach is combining the WAN/SAN [18] and spatio-temporal analysis [17] to summarize the video document based on salient text and key-frame extraction. Our proposed work is experimented over more number of video documents and revealed based on user evaluation that proposed work produces high quality summaries.

The rest of this paper is organized as follows. Section 2 gives system architecture of the proposed work. Section 3 & 4 gives a description of text and video summarization in detail. In section 5 experiments that have been conducted are presented and

discussed. Finally, conclusions and future research directions are presented in Section 6.

## 2. System architecture

Figure 1 describes the overall system architecture of the proposed work. It describes the overall functionality and processing of the application. It has two ways of video document summarization. One way is key frame based video summarization and the other way is text based summarization on video documents. This two-way summarization is then stored in a database for further retrieval. In key frame based video summarization phase, existing systems mostly use spatial low level features like color histogram of video that reflects only spatial changes between frames. The proposed system basically extracts key frames and summarizes them, based on a visual attention model. The system helps achieving a meaningful video summary by bridging the semantic gap between the low level descriptors used by computer systems and high level concepts perceived by human based on neurobiological concepts of human perception of vision.

From psychological studies, it is known that, human beings pay more attention to moving objects than static objects. However, at times interesting static objects at the back ground which is visually appealing might capture the viewer's attention. So to find the frames that are visually attractive and meaningful becomes essential to take both dynamic and static attention into account. Our system models both static and dynamic attention of human. Both the static and dynamic attention are quantified using a descriptor called 'Visual Attention Index' $p_i$ ($I_i$, $H_i$, $x_i$, $y_i$, $dx_i$, $dy_i$), i Є N. N refers to the number of blocks in each frame. Motion intensity and orientation in a video are computed to model dynamic attention of human on a video. Orientation histogram is developed. VAI of dynamic attention is calculated based on coherence of orientation that is calculated using Gaussian kernel density estimation [17]. Static attention of each frame is further calculated based on the Red-Green opponency and blue- Yellow opponency in each frame.

In text summarization phase, the first part of system architecture upload the document, it goes to preprocessing step which includes segmentation of the each sentences, tokenization of each word there after comparing with the stop words to remove the stop words exist in the document. After preprocessing, the sentence identifies the vertices and the edges (pair of words). After finding the vertices and the edges, co-occurrence probability is calculated. Co-occurrence probability is necessary for WAN (Words Association Network) and SAN (Sentence

Association Network). It gives the weight of word and weight of sentences. The sentence weights are used to select the most relevant sentences from the document which represent the whole document. Finally text and video summarization is stored in a database for easy retrieval.

## 3. Key frame based summarization

The proposed system basically extracts key frames and summarizes them, based on a visual attention model. The system helps achieve a meaningful video summary by bridging the semantic gap between the low level descriptors used by computer systems and high level concepts perceived by human illustrated by the following modules.

### 3.1 Static attention module

Figure 2 shows the operation flow of the static attention detection module that designed and implemented in this paper. The video which is to be summarized is splitted into individual frames and the individual frames using are given as input into the static Attention detection model. Each frame is further splitted into 64 blocks that are further operated upon by the following algorithm to get the Visual Attention Index of each frame.

**LMS Color model and human color perception**

LMS color model is used to extract the key frames from videos. These three classes of cones of human eyes are the short-wavelength sensitive (S-cones), middle-wavelength sensitive (M-cones) and long- wavelength sensitive (L-cones), and all have differ- ent but overlapping spectral sensitivities. LMS value ranges from -1 to 1. So, black is (-1, -1, -1), white is (1, 1, 1) and other colors lie in the range from -1 to 1. The sum of L+M is a measure of luminance. According to perception of color theory, each color can be coded by three principal color receptors rather than thousands of color receptors coding for individual colors. The algorithm to calculate Static Visual Attention Index is as follows:

1. Divide each frame into smaller macro blocks in such a way that each frame is subdivided into 64 blocks.

2. Take one block at a time. The color space from RGB to LMS is converted.

3. RGB to LMS conversion is achieved in two steps

    a) First RGB to XYZ conversion is carried out.

b) Followed by XYZ to LMS conversion.

4. RGB to XYZ conversion is done using the transformation matrix

$$[XYZ] = [RGB] \begin{bmatrix} 0.5767 & 0.2974 & 0.0270 \\ 0.1855 & 0.6273 & 0.0707 \\ 0.1882 & 0.0752 & 0.9911 \end{bmatrix} \quad (1)$$

5. Similarly, XYZ to LMS is also achieved through the transformation matrix

$$[LMS] = [XYZ] \begin{bmatrix} 0.7328 & -0.7036 & 0.0030 \\ 0.4296 & 1.6975 & 0.0136 \\ -0.1624 & 0.0061 & 0.9834 \end{bmatrix} \quad (2)$$

6. Calculate Red-Green, Blue-Yellow Opponency using the following formula,

$$\text{Red-Green Opponency} = (L - M)/(L + M) \quad (3)$$

$$\text{Blue-yellow Opponency} = \frac{(S - 0.5*(L+M))}{(S + 0.5*(L+M))} \quad (4)$$

7. Calculate color contrast (H) and intensity feature (I) of each block which are combined by both red-green and blue-yellow opponency.

8. Calculate contrast-based distance using the following formula,

$$d(p_i, q) = \left(0.5|p_i(I) - q(I)|\right) + \left(0.5|p_i(H) - q(H)|\right) \quad (5)$$

Where q represents the 64 blocks of neighborhood, $p_i$ is the descriptor of each block in a frame and is defined as $p_i(I_i, H_i, x_i, y_i, dx_i, dy_i)$, iЄN blocks of frame.

9. Calculate Static Visual attention index as follows:

$$A_S = \frac{1}{N}\left(\sum_{i=1}^{N} W_i C_i\right) \quad (6)$$

Where $W_i$ is Gaussian fal-off weight, $C_i$ is the center-surround differences of block and is obtained by the summation of contrast-based distance of a block d(pi, q).

10. Output the Static Visual Attention Index (VAI) for each frame.
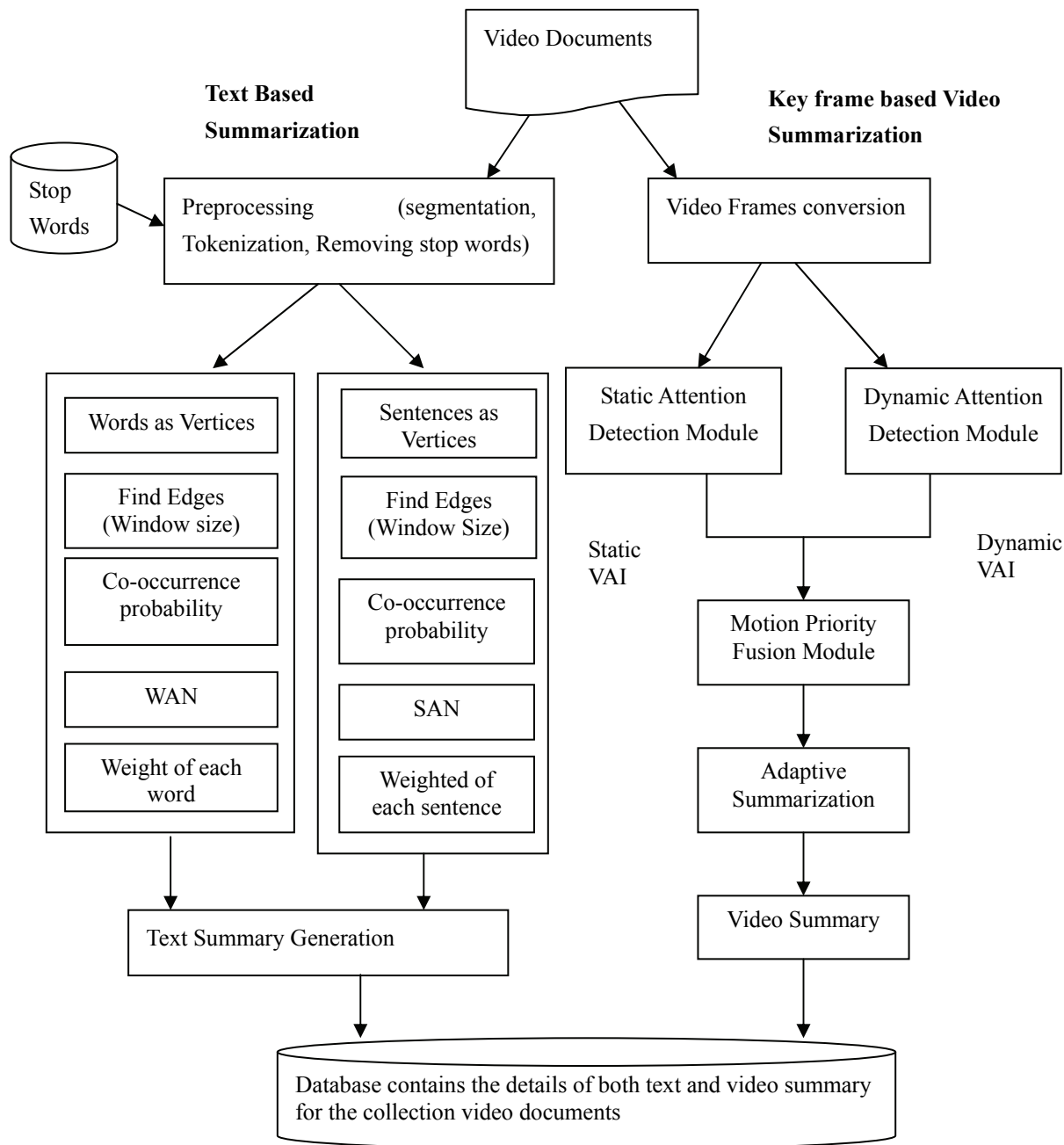
Figure 1    Simplified architecture of the system.
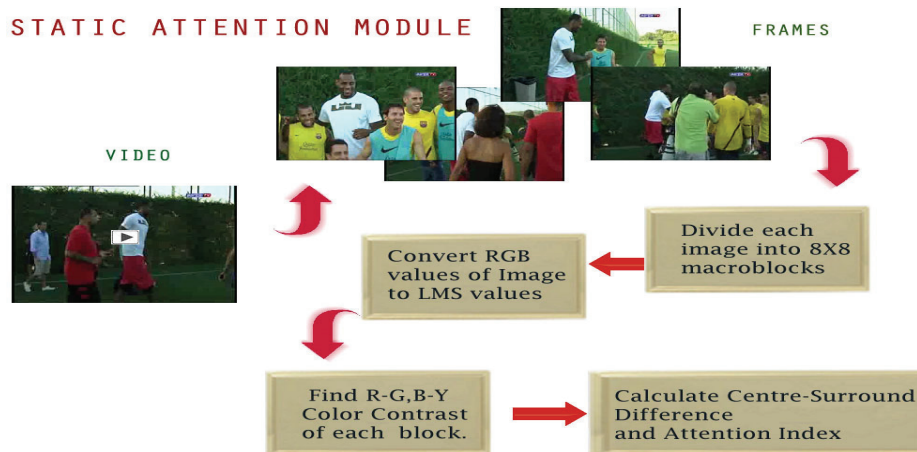


Figure 2    Operation flow of Static Attention detection module

Using the above algorithm, Static Visual Attention Index of each frame is calculated. The process involves conversion of each frame from RGB Color space to LMS color model. Using the LMS values Center-surround difference of each block is calculated and then Visual Attention Index of entire frame is obtained

## 3.2 Dynamic attention module

Figure 3 shows the operation flow of the dynamic attention detection module designed and imple- mented in this paper. In the Dynamic attention detection model, the motion associated with each frame of the video is calculated. The frame is divided into 64 blocks and the motion associated with each block is calculated by calculating the Motion intensity and the Motion orientation.
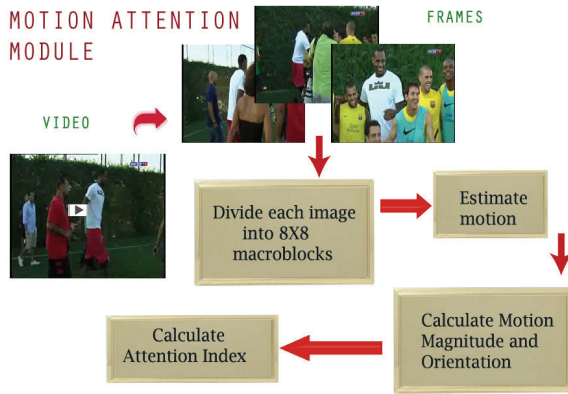


Figure 3    Operation flow of Motion Attention

The motion estimation is done on the basis that the patterns corresponding to objects and background in a frame of video sequence move within the frame. The idea behind block matching is to divide the current frame into a matrix of macro blocks that are then compared with corresponding blocks and its adjacent neighbors in the previous frame to create a vector that quantifies the movement of a macro block from one location to another in the previous frame.

This movement calculated for all the macro blocks comprising a frame constitutes the motion estimated in the current frame. The search area for a good macro block match is constrained up to p pixels on all four sides of the corresponding macro block in previous frame. This 'p' is called as the search parameter. The matching of one macro block with another is based on the output of a cost function. The macro block that results in the least cost is the one that matches the closest to current block.

There are various cost functions, Mean Absolute Difference (MAD) and Mean Squared Error (MSE), used to calculate Motion Attention Index for each frame based on the following algorithm:

1.    Divide the frame into 64 sub-blocks

2.    Compare each block of a frame with the corresponding block in the next frame and calculate the motion vectors dx and dy by block matching technique.

3.    The motion intensity of each block is given by the formula

$$\text{Motion intensity} \quad \gamma_i = \sqrt{dx_i^2 + dy_i^2} \qquad (7)$$

Where $x_i$, $y_i$ are the location information of block I and $dx_i$, $dy_i$ are the motion vectors.

4.    The orientation of each block is calculated

$$\text{Motion orientation} \; \Theta_i = \arctan(\; dy_i / dx_i\;) \qquad (8)$$

5.    Then, orientation histogram is built and the motion attention index is calculated from the following formula:

$$A_T(i) = 1 - \frac{v(b(i))}{\sum_{i=1}^{H} v(j)} \qquad (9)$$

Where b(i) is the bin index of block i, the histogram value at bin index j is v(j) and H is the maximum bin index.

6.    Motion Attention of the block is given by

$$A_T(i) = \gamma_i * A_T(i) \qquad (10)$$

$A_T(i)$ gives the Motion attention Index of frame i. Motion with high intensity attracts more attention. Frames with larger attention index are given much importance.

## 3.3 Motion priority fusion module

The motion priority fusion is done the static VAI and dynamic VAI of frame. The algorithm applied in Motion priority fusion Module is as follows:
1. Input the Static and Dynamic Visual Attention Indexes of each frame
2. The dynamic attention and static attention weights for each frame is computed by using the following formulae

$$W_T = A_T'.\exp(1 - A_T') \qquad (11)$$

$$W_S = 1 - W_T \qquad (12)$$

$$\text{where} \quad A'_T = Max\,(A_T) - Mean\,(A_T) \quad (13)$$

3. Calculate final Visual Attention Index

$$VAI = W_S.A_S + W_T.A_T \qquad (14)$$

4. Output the Total Visual Attention Index

### 3.4 Adaptive summarization module

The video is summarized with the important frames depending on the importance of the frame determined with the Visual Attention Index of each frame using Adaptive Summarization algorithm.

The adaptive summarization algorithm when applied over all the frames of the video gives the set of frames which have high Visual Attention Index and do not contain redundant information, which could be used in the Video Summary.

1. Input the video and get all the frames.
2. Calculate Color histogram for each frame.
3. Cluster the frames using K-Means clustering algorithm based on Euclidean distance.
4. Select Representative key frames (frames having high VAI) from each cluster.
5. Compare potential keyframes by calculating focus point shift between them.

$$D_{ij} = \left| \sum_{n=1}^{N} A_n^i - A_n^j \right| \qquad (15)$$

6. Calculate variation of attention index in a shot using

$$\overline{d}_s = \frac{\sum_{i=2}^{M_s} D_{i\,i-1}}{M_s} \qquad (16)$$

where $M_s$ represents the number of frames in shots.

7. Filter out adjacent keyframes being selected as keyframes using the inequality.

$$D_{key} > D_{ave} + \delta.D_{div} \,, (\delta=1.5 \text{ here}) \qquad (17)$$

Where $D_{ave}$ is the average VAI difference of the keyframe, $D_{div}$ is the standard deviation.

Thus by using the above algorithm, the most representative key frames are extracted from the pool of video frames.

## 4. Text based video summarization

Each video document has its description as in the form of document. The system consists of six modules. Te operation and the result of each module are discussed below.

*Text Preprocessing*
Text preprocessing module is used to preprocess the test before SAN and WAN can be applied. This module contains the three sub modules.

1. Sentence segmentation
2. Tokenization
3. Stop-words

*Sentence segmentation*
This module separates each individual sentence i.e. the sentence is separated after each full stop.

*Tokenization*
The functionality of this module is to separate the words from the given text document which are passed to another module.

*Stop-words*
In computing, stop-words are words which are filtered out prior to, or after, processing of natural language data (text). It is controlled by human input and not automated. There is not one definite list of stop words which all tools use, if even used. But for that we are using the more than 300 stop-words for the text summarization.

*Find vertices*
All the unique words in the documented are extracted and are called vertices. To generate the summary, it is necessary to identify the unique words form the text file.

*Find undirected edges*
All the possible pairs of the vertices are formed and those pairs which lie in the same window size are called undirected edges. Here window size refers to within a sentence.

*Co-occurrence probability*
Using the count of vertices i.e. Count (i) and the count of the pairs i.e. Count (i, j) is calculated.

*WAN (Word Association Network)*
System proposes a new term weighting approach based on WAN and is used for text summarization. We first treat all remaining individual words after preprocessing as vertices. An undirected edge is associated with any two vertices if the correspond-

ing two words co-occur within a sentence (sentence is the window size).

*SAN (Sentence Association Network)*

System [18] also investigates SAN. We first define co-occurrence probabilities between sentences and compute their co-occurrence information. We observe that a sentence is more salient if it has more co-occurrence information. Therefore, the weight of a sentence is evaluated by its sum of co-occurrence information.

This sentence association network is build on the basis of co-occurrence information of two sentences which share words.

*Pseudo Code*

Here the pseudo code for calculating the text summary is mentioned.

*To calculate the summary:*

Step 1.   Input the sentence and sentence weight
          ps_wn_i sentence length s_len.
Step 2.   Set index = order_sentence_index
          (ps_wan,slen).
Step 3.   Initialize string str=" ",i=0;
Step 4.   str=str+sentence[index[i]].
Step 5.   Increase i by one.
Step 6.   IF i=s_len THEN go to step 7.
          ELSE     Go to step 4.
          ENDIF
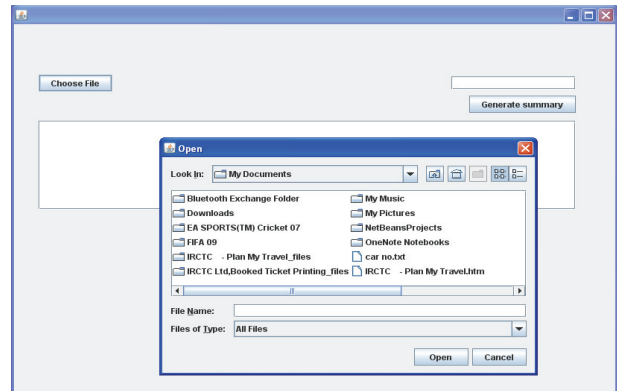Step 7.   Return str.

# 5. Experimental results

The key frame based video summarization methodology described in section III. The implementation is done in MATLAB R2009b provides the summary of the video as a sequence of salient frames of the video as shown in Figure 4.
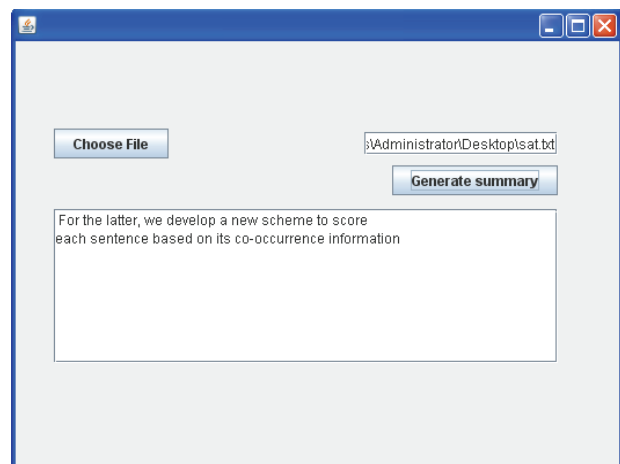


Figure 4    Video summary examples.

The experiments result is shown below for the text based video summarization. For generating the

summary of selected file, the generating summary button is button. After pressing the generating summary button, the summary of the contents of selected file is shown in Figure 5.



**(a)**



**(b)**

Figure 5    Text-based video summary steps
(a) Selection of video document (b) Summarized.

## 5.1  Performance evaluation

This work is tested on more than 200 short videos downloaded from different websites. The results showcased here with 3 short videos belong to different genres involving varying degrees of motion in them. The results are represented in Table 1. The summary thus produced is presented before *20 users*, who are later made to watch the video and their feedback on the effectiveness of the summary that is obtained. The feedback is presented in Table 2.

All the three video summaries are shown to the users and then they are made to watch the video. Their rating for the video summary is as excellent, good and satisfactory as obtained. The evaluation is based on above mentioned 3 main criteria. Content coverage refers to the extent to which the summary effectively conveys the content of the video. Presentation indicates how effective the presentation is.

Table1. Video summary details

| Video Id | No. of Sentences | | No. of words | | Total No. of frames | No. of Key frames extracted |
|---|---|---|---|---|---|---|
| | (Before Summarization) | (After Summarization) | (Before Summarization) | (After Summarization) | | |
| Video 1 | 25 | 3 | 220 | 28 | 447 | 10 |
| Video 2 | 64 | 8 | 539 | 62 | 120 | 7 |
| Video 3 | 110 | 15 | 640 | 76 | 135 | 7 |

Table2. User based evaluation results

| Video Id | Content Coverage | Present ation | Total effectiveness |
|---|---|---|---|
| Video 1 | Good | Good | 8 |
| Video 2 | Satisfactory | Good | 7 |
| Video 3 | Excellent | Good | 9 |

It refers to the entertainment quotient of the summary. Total effectiveness that is *rated out of 10* is based on *total satisfaction* provided by the summary.

# 6. Conclusion and future work

With the exponential growth of the information on the Internet, a level of abstraction of information from the results of IR becomes necessary. That is, the large number of documents returned by IR system need to be summarized. Currently this is the primary application of summarization.

This work is developed in order to summarize the video documents by using both key frame and textual Information present in it. Thus the work undertaken helps us to understand the given video by means of extracted key frames. The key frames extracted by the techniques that take both color and motion information into account are more relevant than the ones that use either color or motion information only. The effectiveness of the key frames extracted by this approach is more when compared to other extraction methods using either only color or only motion based techniques. From the user responses, it is found that our system produces a very effective summary. However in the system the number of key frames in the summary is independent of the duration and dynamicity involved in the video. This work is further used to summarize the given video document based on text. This included text preprocessing, word association network and sentence association network. Through extensive experiments over more number of video document data sets, we have showed that our

approaches can produce high quality summaries.

In the future, we plan to apply our schemes to multi-document summarization, text classification, text clustering, keyword extraction and information retrieval, in order to extend their utility. Besides, more priori-knowledge on text semantics or sentence semantic structures will be integrated into our system. And the effectiveness can be further improved by determining the sufficient amount of key frames required to summarize the video based on the dynamicity of the video. Also the system can be enhanced by allocating more key frames representing shots that involve more action. This requires efficient shot boundary detection method. These are the enhancements that could be added to the system in future.

## References

[1] T. Liu, H. Zhang, J, and F. Qi, "A Novel Video Keyframe Extraction Algorithm based on Perceived Motion Energy Model", *IEEE Trans. on Circuits and Systems for Video Technology,* Vol. 13(10), pp.1006 -1013, 2003.

[2] Ioannis Patras, Emile A. Hendriks, Reginal L. Lagendijk, "Probablistic Confidence Measures for Block Matching Motion Estimation", *IEEE Transaction on Circuits and Systems for Video Technology*, Vol. 17, No. 8, pp. 989-995, 2007.

[3] Y. F. Ma, L. Lu, H. J .Zhang, and M. Li, "A User Attention Model for Video Summarization", *ACM Multimedia*, pp. 533-542, 2002.

[4] Ani Nenkova, "A Compositional Context Sensitive Multi document Summarizer: Exploring the Factors that Influence Summarization", *SIGIR*, 2006.

[5] C. M. Taskiran, Z. Pizlo, A. Amir, D. Ponceleon, and E. Delp, "Automated Video Program Summarization using Speech Transcripts", *IEEE Trans. on Multimedia*, Vol. 8(4), pp. 775-791,2006,

[6] Stergos Afantenos, Vangelis Karkaletsis, Anagiotis Stamatopoulos, "Summarization from medical documents: a survey", *Journal Artificial Intelligence in*

*Medicine, Elsevier Science Publishers Ltd.* Vol. 33(2), pp.157-177, February, 2005.

[7] M. Guironnet et al., ''Spatial-Temporal Attention Model for Video Content Analysis'', *Proc IEEE Int'l Conf. Image Processing (ICIP),* pp. 1156-1159, 2005.

[8] Y.F. Ma and X.S. Hua, ''A Generic Framework of User Attention Model and Its Application in Video Summarization'', *IEEE Trans. Multimedia*, Vol. 7 (5), pp. 907-919, 2005.

[9] L. Liu and G. Fan, ''Combined Key-Frame Extraction and Object-Based Video Segmentation'', *IEEE Trans. Circuits and Systems for Video Tech.*, Vol. 15(7), pp. 869-884, 2005.

[10] R. Mihalcea and P. Tarau, "Text rank: Bringing order into texts", *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP 2004)*, Barcelona, Spain, 2004.

[11] Zhonghua Sun, Kebin Jia, Hexin Chen, "Video Key Frame Extraction Based on Spatial-temporal Color Distribution", *International Conference on Intelligent Information hiding and multimedia signal processing,* pp.196-199, 2008.

[12] Nuno Vasconcelos, Andrew Lippman, "A Spatio-temporal Motion Model for Video Summarization", *IEEE computer society conference on computer vision and pattern recognition*, pp.361-366, 1998.

[13] Ajay Divakaran., Regunathan Radhakrishnan. and Kadir A. Peker. "Motion Activity-Based Extraction of Key-Frames From Video Shots", *International Conference on Image processing*, Vol. 1, pp.932-935, 2002.

[14] Sathish kumar l.Varma and Sanjay N. Talbar. "Video summarization using dynamic Threshold", *International conference on Emerging trends in Engineering and Technology*, pp. 120-123, 2010.

[15] Tiecheng Liu and Ravi Katpelly., "Content-Adaptive Video Summarization Combining Queueing and Clustering", *International Conference on Image Processing*, pp.145-148, 2006.

[16] Unghee Jungjunehwa Song And Yoonjoon Lee, "Narrative Based Abstraction Framework For Story Oriented Video", *ACM Transactions on Multimedia Computing Communications and Applications*, Vol. 3(2), pp. 1-28, 2007.

[17] Jiang Pen And Qin Xiao Lin, "Key Frame Based Video Summary Using Visual Attention Clues", *IEEE Transactions on Multimedia*, 2010.

[18] Yuhui Tao, Shuigeng Zhou, Wai Lam, Jihong Guan, "Towards More Effective Text Summarization Based on Textual Association Networks", *Fourth International Conference on Semantics, Knowledge and Grid,* pp. 235 – 240, 2008.

[19] G. Erkan,, and D. Radev, "Lexrank: Graph-based lexical centrality assalience in text summarization", *Journal of Artificial Intelligence Research*, Vol.22, pp.457– 479, 2004.

[20] S. Hassan, R. Mihalcea, and C. Banea, "Random-walk term weighting for improved text classification", *International Conference on Semantic Computing, Irvine, CA*, pp. 1-8, 2007.