# An optimized approach for multi-user facial expression-based video summarization.

**Ashokkumar S[1], Suresh A[2], Priya S[2*], Dhanasekaran R[3]**

[1]Anna University, Chennai, India

[2]Department of Electrical and Electronics Engineering, S.A. Engineering College, Chennai, India

[3]Syed Ammal Engineering College, Ramanathapuram, India

## Abstract

**Video Summarization (VS), is one of the recent attractive research studies in huge amount of video contents news, sports videos, TV reality shows, Movie reviews etc., Based on user's facial expression, VS is the challenging task in digital video technology. An automatic recognition of facial expressions is the necessary task in VS used in behavioural science and in clinical practice. The reliable recognition of facial expression by a machine is the critical compared to by a human. The non-uniqueness and the deformations in the human face require the toleration in human action recognition system. This paper proposes an effective toleration of facial variations themselves and addresses the problems in facial expression and video summarization and focuses the Multi-user Facial Expressions for Video Summarization (MFE-VS). The utilization of quantitative measures in proposed work estimates the user's viewing interest based on the human emotional perspective. Initially, Gaussian filter eliminates the unnecessary things in an input image. Then, Viola-Jones face detector detects the face and eye regions. The utilization of Histogram of Gradients (HoG) and the Local Land scape Portrait (LLP) for extraction of shape and texture features. Once the labelling of features is over, the Probabilistic Neural Network (PNN) classifier predicts the exact expression of the user in the frames. Based on the obtained expression, the proposed work extracts the interesting frames and summarizes them that constitute final stage. The variations of facial expressions of viewers constitute attention states and the classification of positive and neutral states constitutes emotional states. The comparative analysis of proposed method with the traditional methods on the parameters of success rate, precision, recall, and interest scores assures the suitability of MFE-VS in facial expression based video summarization.**

## Introduction

Digital video technology, a fast growing field in multimedia processing in recent years with the enormous amount of digital contents such as reality shows, special TV programs, sports, news etc. The time lag of the user to watch the video requires the abstraction instead of whole. The video abstraction or summarization has the capability to convey the information related to the important incidental occurrences in video. In general, the collection of shots rather than scenes termed as single layer video abstraction model. The summarization process constitutes the sequential procedures of shot boundary detection, computation of summarization parameters (motion, color frames, audio levels), selection of important shots and writing video contents. The comfort and the story telling quality decide the summarization quality and personalization. The reduction of flickering caused by the abrupt transitions and

the utilization of narrative user preferences increased the level of personalization and quality. There are variety of methods are available to increase the summarization quality described in detail in next sessions.

Research studies discuss the components required for video indexing and retrieval, how they are arranged in hierarchical manner, inter-relationship between the components. The review has been conducted with the focus on different processes namely, key frame extraction, segmentation, video classification, video mining conveys the future developments in video summarization. User interest is one of the key attribute in the conduction of video summarization process. In view of attention and emotion, quantitative measures are available to estimate the user interest and such a systematic approach is called as Interest Meter (IM). The application of fuzzy fusion scheme combines the characteristics of both

attention and emotion and creates the abstract based on interest scores. The computer vision applications include important tasks namely, human face detection and moving objects detection. An effective use of dynamic programming in computer vision applications estimates scene skimming length and offer optimal skimming. But, the accurate generation of key-tags is the difficult task in skimming approaches. The evolution of event driven approaches improves the operational speed of key-shot identification and an accurate key-tags generation. The utilization of color based methods does not consider motion information. But, the extraction of high action content depends upon motion analysis. The summary generation depends upon the key frame extraction process that conveys the overall message in motion videos. The visual saliency detection mechanism, adaptive fast forwarding and content truncation are the major processes for key frame extraction and generation of comfort video summary respectively.

The fast face detection and recognition is an important stage in video-surveillance applications. The first object detection framework in research studies is that Viola-Jones (VJ) face detection method. The increase in number of pixels makes the conventional Viola-Jones method was unsuitable. Hence, the VJ integrated with the openCL to improve the performance. Training and detection are the two stages in VJ face detection method. The adaptation engine on the basis of VJ method includes user and environmental requirements as constraints. Based on the consideration of constraints and intermediate information an adaptive decision maker estimates the best coding parameters. Due to the implementation capability of VJ in openCV, the unknown size object can be easily identified. Three techniques are used in VJ method such as extraction of features in an integral image, learning and construction of strong classification by using the integration of Adaboost with weight value, an efficient combination of features by using the cascade classifier and isolation of face and non-face region. The retrieval, index, and content analysis require the text and caption detection.

The evolution of corner based approaches detects the text and caption of the videos. The presence of orderly and dense based corner points in characters in corner points approaches extracted the important occurrences in the video. The making of comics from video is an attractive research area. Even though various tools and methods are available, the conversion of video into comics is the time consuming process. The evolution of new method called automatic turning of movie into comics follows two principles namely, optimization of information preservation, output generation with rules and styles of comics. In motion analysis, background moving objects caused the distracting motions. Hence, a new mechanism is required to isolate the action motions from background motions and track the trajectory key points. The separate extraction of spatial and temporal information from original and Motion History Image (MHI). The MHI based information extraction efficiently avoids the unreliable tracking. The utilization of Histogram of Oriented Gradients (HOG) characterizes the temporal features and the spatial

features in an intensity model. On the other hand, the Gaussian Mixture Model (GMM) recognizes the action by the combination of both the spatial and temporal features. The built up of semantic gap between feature space and user perception requires the visual attention based key frame extraction. The computation of static and dynamic visual attention and their non-linear combination extracts the key frames.

The classification of extracted features is an important process in video summarization. On the basis of learning from the samples, the use of Neural Network (NN) classifies the extracted patterns. The time consumption for training is more in and the use of incremental adaptation leads to high false minima. The use of back propagation neural networks with statistical derived function instead of sigmoidal function referred the Probabilistic Neural Network (PNN). The PNN contains input, hidden, summation, output layers. The neuron in the first layer is the predictor variable. The hidden neuron in the second layer use the sigmoidal function in RBF kernel function estimates the distance between the neuron center point and test case. The research works conveys that all the traditional video summarization methods employs the single user interest since all the users has not same interest and high dimensionality.

This paper discusses the video summarization based on facial expression. The proposed detection and recognition scheme has the capability to tolerate the variations. The main objective of proposed work is to make the system consciousness about the user's reactions to measure the user interest and conduct the video summarization. Initially, the proposed work identifies the facial expression of the user and extracts the interesting patterns from the video. Depends upon the facial emotion perspectives such as attention and emotion states, the proposed work estimates the user viewing interest. The technical contributions of proposed work are listed as follows:

- A novel video summarization method based on user facial expressions with the optimized steps by using Local Landscape Portrait (LLP) texture pattern.
- The summarization of video is done directly based on human psychological states
- The effective analysis of viewing behaviours from facial expressions and the utilization of fuzzy scheme offer better video summarization compared to traditional methods.

The paper organized as follows: Section II discusses the related works on the video summarization, influence of Viola-Jones method on summarization, Local Landscape Portrait (LLP) texture extraction and the PNN- based emotional and attentions states classification. Section III presents the proposed MFE-VS implementation. Section IV presents the comparative analysis of performance parameters between proposed and existing methodologies. Finally, section V describes the conclusion to show the effectiveness of proposed MFE-VS in video summarization.

## Related Work

This section describes the various related works on video summarization methods and associated merits and demerits to provide the way of proposed work. Due to the various characteristics such as richer content, huge amount of data and little priority structure and user time lag, indexing and retrieval of information from video were the difficult process. Hu et al. presented an overview of strategies used in visual content based video indexing and retrieval process [1]. They discussed the components used for index and retrieval process, the inter relationship between the components and the future directions for video indexing and retrieval. The evolution of security systems considered the video surveillance scheme as an important process. Jasmine et al. discusses the modern approaches [2] in Human Action Recognition (HAR) in real-time videos. Sangeetha et al. recognized the unauthorized person in different angles in long-take video shots [3]. They presented the types of video summarization such as static and dynamic. In static summarization, the number of key frames were extracted whereas in dynamic skimming of original video into shorter version with preservation of original information with the timing properties. They formulated the skimming problem as a two stage optimization problem with the audio-channel analysis. The prediction of influence of emotion decided the people's attitude. Peng et al. proposed an Interest Meter (IM) to make the system consciousness about the user interest [4]. They included the spontaneous reactions during the viewing process in quantitative measures to measure the user interest. They compared the interesting scores for various movies with the existing method of multi-level video summarization [5], Random selected summaries and NOVOICE. Bio-metric authentication plays the major alternative role to passwords remembering. Kumar et al. utilized the video based retina recognition model in which more information was available rather than single image [6]. The growth of several videos in internet exploring of videos is the time consuming process and leads to performance degradation. Wang et al. discussed the web-video summarization with the event driven approach [7]. The tags associated with the shots were localized initially and then by a matching the query with the tags, relevant shots predicted. The presence of high action content in sports video needed a motion activity. Mendi et al. presented a motion analysis based automatic video summarization [8]. Using optical flow algorithms, motion parameters were estimated with different key frame selection criteria.

Research works addressed the crucial automatic video summarization due to multiple video libraries and rapid browsing. Dang et al. presented the entropy based measure of Heterogeneity Image Patch (HIP) index [9] and the HIP curve formed. Based on HIP curve, the key frame extraction and the dynamic video skimming problems were solved. HIP based video summarization offered the low complexity. The generation of comfortable, optimal and compact video summaries are the challenging tasks due to the resource allocation problem. Chen et al. considered the content associated to the scene activity to provide the optimal trade-off between playback speed and comfort [10]. The face detection is an initial stage in video surveillance applications. The optimization in the face detection steps improved the performance of video summarizing applications. Wang et al. discussed the parallel execution of Viola-Jones with the openCL [11] to improve the performance of the system. The parallel execution effectively made the time cost hidden and offered the high throughput. The accommodation of large number of users on diverse applications was not handled by conventional coding schemes. Zhang et al. included the multiple level video signal components and semantic features in the intermediate video description called 'inter-media' [12]. A novel-video adaptation based on inter-media created the bit stream with the low complexity. The challenges in VS based on facial expression are the variations of pose, illumination, and expression. Rowden et al. investigated the multi-frame sore level fusion to improve the accuracy of facial recognition [13].

The effective video navigation is the problematic area, since the video database was increased. Shao et al. proposed the content based video retrieval approach for searching of large video databases included the human actions and spatio-temporal localization [14]. The presence of complex backgrounds in the images degraded the performance of video summarization based on facial recognition. El et al. presented the study of novel and simple approach based on combination of techniques and algorithms on Viola-Jones face detection [15]. For a better content analysis, the detection of text and caption is an important stage. Zhao et al. proposed the corner based approaches for detection of both text and caption [16]. They utilized the various discriminative features for the description of text region. The algorithm was developed to detect the moving captions in videos depends upon extracted text features. The major advantage of proposed method is language independence. Thanikachalam et al. analysed the various new methods for HAR with the single space time approach that speeded up the Human Machine Interactions (HMI). The authors analyse the usage of Ontology on image [17]. The extension of HAR approach to 3D spatiotemporal volumes by correlation height filter [18] and defined Accumulated Frame difference (AFD) to find the gradient direction in motion appearance model respectively [19]. The modelling and detection of complex events in unconstrained real videos required the semantic vectors. Merler et al. used the semantic model vectors [20] that were created by discriminative classifiers. Also, they combined the semantic model vectors into visual descriptors for an end-end video detection system. The creation of comics from the movie is the time consuming process even though the various tools are available. Nagarajan et al. provided the semantic meaning to the images visual words creation with low level features [21]. Wang et al. presented the scheme for automatic comic generation from movie clips [22] based on two principles optimization of preserved information, and styles based output generation. The making of comics included various

components such as script-face mapping, descriptive picture extraction and cartoonization.

The recognition of human action in moving background is the difficult process. Tian et al. analysed crowded videos and detected the abnormal actions by using Motion History Image (MHI)-based Hierarchical Filter Motion (HFM) [23]. The estimation of recent motion detected the interest points and the use of local motion filter smoothed the gradients of MHI. Thanikachalam et al. utilized the MHI based on correlation filter model to find out the direction of motion [24]. The effective video summarization model is the visual attention model. But, the computational cost is high. An extra-effort was required to capture the required information in user-centric models. Gkonela et al. discussed implicit user interactions analysis [25] with the playing options in web video player. Interactions based method offered more accurate results in the detection of events with short duration. Ejaz et al. reduced the computational cost by using temporal gradient models [26,27]. The non-linear weighing fusion method detected the static and dynamic attention states in video summarization. The incorporation of object semantics improved the video summarization, indexing and retrieval compared to existing methodologies of DT [28], STIMO [29], and VSUMM [30]. Yong et al. presented the framework to model the semantic contexts for key frame extraction [31]. The utilization of one-class classifier in semantic contexts modelling located the significant novelties. The utilization of semantic key frames in the formation of co-occurrence matrix achieved the better key-frame extraction.

Research studies pointed out one of the challenging tasks namely, pedestrian detection approach with the view and posture problem. Ye et al. discussed the Error Correcting Output Code (ECOC)-based pedestrian detection approach by using [32] classification. The detailed description about the concepts and techniques of multi-view and multi-posture pedestrian patterns. Less detection speed and the detection rate are the major limitations observed in ECOC. The detection of abnormal activities is the major problem in ECOC classifier. Lin et al. discussed the patch based methods for modelling of normal activities. By this way, the isolation of abnormal activities from normal provided and avoided the blob extraction errors [33]. They compared the performance of blob optimization method with the existing methods of particle, TLD, CT, KSP, KCF, blob associate [34-39]. A new emerging field of wearable technologies action cameras and smart glasses increased the viewer interest in the view of first person. Betancourt et al. presented the various approaches in which different image features and quantitative methods were combined [40]. The image feature combination achieved the specific objectives of activity recognition, object detection and user-machine interaction. The ability of predicting the on-going user activity in earlier stage referred early recognition. The prediction of intended activities against abnormal events was the difficult task. Ryoo et al. presented an algorithm to perform activity recognition [41] at an early stage. They introduced the effective summarization of pre-activity observations referred 'onset' by using histogram of time series

gradients. They offered better human activities recognition from the first person videos. Video classification is one of the important applications in structure analysis. Sigari et al. proposed ensemble classifier [42] for sports video classification with the features namely color, dominant gray level, cut and motion rate. The extracted features were classified by using various simple classifiers (Probabilistic Neural Network (PNN), Decision Tree (DT)) and weighted majority vote is used to combine the outputs of each classifier to show the correct classification rate.

The PNN classifier is used in various real time image recognition applications. Saftoiu et al. discussed the post processing software analysis [43] which computed the individual frames and retrieved the hue histogram from numeric matrix on the basis of neural network analysis. A fast and accurate medical diagnosis based on medical decision process utilized the artificial intelligence. Researchers focused on the maximization of user experience and minimization of service cost in Internet Protocol Television (IPTV). The accurate prediction of popularity in the videos was the major problem. Li et al. improved the popularity prediction accuracy by using neural network model [44]. The utilization of historical logs effectively validated the popularity prediction of IPTV programs. The results obtained from the neural network model improved the accuracy. The existence of discrete nature of patterns in PNN and PNN-homogeneity methods leads to low classification speed and high memory requirements. Savchenko modified the conventional PNN scheme that required the histogram of training input samples [45]. The detailed description of text authorship attribution with n-gram features, face recognition provided the better accuracy with less memory consumption. The recognition of semantic contents from both image and video was the major problems in content based video analysis techniques. Karmokar et al. the composite feature vectors of extracted semantic features from low level features on the basis of color, text, and shape [46]. They used Manhattan distance and neural network classifier to validate the efficiency of recognition. The classification problem of assigning an observation to different disjoint groups played an important role in decision making process. Khashei et al. discussed the novel hybridization in VS by using the combination of Artificial Neural Network (ANN) with multiple linear regression models [47]. The ANN utilization effectively improved the classification accuracy compared to traditional classification models. Various research studies in this section conveyed that the video summarization based on multiple user interest is the difficult task. Hence, this paper overcomes the limitation by using recent face detection, feature extraction and the classification provides the better video summarization.

## Multi-User Facial Expression-Based Video Summarization

The video summarization is the difficult process, since the face of the each user in the video is not the unified object. The variety in deformations causes the difficulty in facial

expression based abstraction. This paper proposes the detection and recognition scheme with the capability of tolerating the variations in the human faces. The proposed Multi-user Facial Expression-based Video Summarization (MFE-VS) contains four significant processes to make the computer consciousness about the user viewing behaviour. They are listed as follows:

- Pre-processing
- Viola-Jones face/eye detection
- Feature extraction
- PNN-classification

Initially, the video to be summarized is given as the input to the system. Then, the proposed frameworks decompose the whole video into several numbers of frames. By using Gaussian filtering techniques, the unnecessary and noise oriented components are filtered. Then, the Viola-Jones (VJ) method detects the individual user faces and eye in each frame. The application of VJ to the proposed work is to summarize the video based on the Multi-user behaviour and facial expressions. Then, three methods such as Histogram of Gradient (HoG), Local Landscape Portrait (LLP), and color feature extraction are employed to extract shape, texture, color (mean, standard deviation, kurtosis, and skewness) features from each face and eye region Flow process of MFE-VS is shown in Figure 1.

Based on the features related to user behaviour, Probabilistic Neural Network (PNN) classifies two states namely, attention and emotional. The variation in facial expressions decides attention states and the recognition of positive or negative viewing behaviour decides the emotional states. The single user view based video summarization provided the non-satisfactory level since; each viewer has the different opinion during watching. If we include many users, then the summarization model was not suitable. Hence, Multi-user viewing based VS proposed and the corresponding parameters are also analysed and compared with the traditional approaches. The detailed implementation of all the processes discussed in this section. The prediction of exact expression of the users and based on the user expressions, the indexes of interesting frame are collected and the video is reconstructed which is called summarized output.

## A. Pre-processing

The frames in video considered as the two-dimensional image. The existence of unwanted components in the image affects the quality of recognition. Hence, in image processing applications, pre-processing is the first and foremost stage. The employment of filtering techniques removes the unwanted or noise components in spatial domain. The interoperability of information exist in an image requires the quality improvement. On the basis of convolution theorem, the frequency domain technique decomposes the image into frequency domain components.
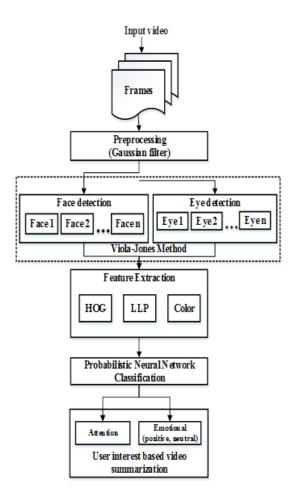


*Figure 1. Flow diagram of MFE-VS.*

The output of $g$ $(x, y)$ convolution depends upon position variation $h$ $(x, y)$ of an input image $f$ $(x, y)$ described by Equation 1

$$g \ (x, y) = f \ (x, y) \times h \ (x, y) \rightarrow (1)$$

If the $G$ $(u, v)$, $F$ $(u, v)$ and $H$ $(u, v)$ are the transformed output of variables, then the convolution in Equation 1 is described by Equation 2

$$G \ (u, v) = F \ (u, v), H \ (u, v) \rightarrow (2)$$

The working of spatial filters depends upon the neighbourhood pixels referred as sub-image. The noise smoothing for a given input image. The weighted sum of intensities of the pixels constitutes the intensity of the pixels in noise smoothened image. The operation of noise smoothing is governed by Equation 3

$$f(x, y) = \sum_{i=-k}^{k} \sum_{j=-1}^{i} f(x + i, y + j) h(i, j) \rightarrow (3)$$

If the variation of weight values on the kernel function represents the Gaussian coefficients, then the corresponding filter is formed with Gaussian distribution curve. The peak value of Gaussian curve specifies the center of the kernel. The increase in distance from the center pixel decreases the weight values correspondingly. The heavy biased weights towards the center cell and the immediate neighbours in the Gaussian filter

preserve the edges. The Sigma parameter refers the standard deviation (in numerals) decides the shape of the distribution curve. The large value of sigma assigns more weight to distance cells for greater smoothing and less edge preservation. The function to govern the noise elimination and image smoothening is described as follows:

$$h(x, y) = \frac{1}{2\pi\sigma^2} e^{-\frac{x^2 + y^2}{2\sigma^2}} \rightarrow (4)$$

The Equation 4 describes Gaussian filter is separable. The Gaussian filtering technique effectively removes the noise samples compared other spatial filters (low and average). The zero-mean noise in an image leads to an effective Gaussian filter.

The convolution operation performed in Gaussian filter performed in two ways of either directly or by using FFT. In FFT computation, the difference between the borders (left, right and top, bottom) will raise the artifacts in an image. Hence, direct computation to be preferred.

## B. Viola Jones Face/Eye detection

The wide utilization method for real time object detection is Viola-Jones face detector. The detection algorithm contains two stages training and testing. The training by machine learning frameworks consumes more execution time due to more number of iterations involved.

The illumination, pose and the non-uniqueness are the characteristics of a human face that leads to detection work as a complex one. But, the Viola-Jones face detection is the fast robust real time face detection to accurately identify the face and the parts in it. The operational stages in Viola-Jones face detection are as follows:

1. Haar feature selection
2. Integral image creation
3. Adaboost training
4. Classifier cascading

**Haar feature selection:** The properties common to all the human faces are the eye region is darker than the upper cheeks and the nose is brighter than the eyes which raises the facial features of location, size and value. The selection of Haar feature used the rectangular features as shown in Figure 2 with the following value

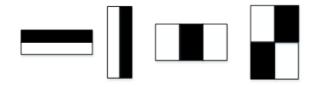$Diff_H = \sum$ *black area pixels* $-\sum$ *white area pixels* $\rightarrow (5)$
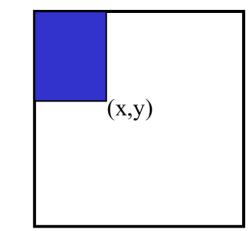


***Figure 2.*** *Haar features.*

The computed feature specifies the location of sub-window. The number of rectangle and corresponding computation is described in Table 1.
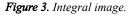
***Table 1.*** *HAAR features.*

| Rectangles | Computation |
| --- | --- |
| Two-sided rectangle | Difference between sum of pixels with in two regions |
| Three-sided rectangle | Sum of pixels within two outside rectangular region and subtract with the pixels in center rectangle |
| Four-sided rectangle | Difference between the diagonal pairs of rectangles. |

**Integral image:** The calculation speed in Haar features is slow. To speed up the process, an integral image is formed. The value is computed by using the sum of pixels in above and left of position as shown in Figure 3 refers an integral image. The computation of integral image by Equation 6

$ii\ (x,\ y) = \sum_{x}^{i} \leq x;\ _{y}^{i} \leq y\ value\ (x_i,\ y_i) \rightarrow (6)$



***Figure 3.*** *Integral image.*

**Adaboost learning algorithm:** The computation of whole features is an expensive process and the less number of computed features form an effective classifier. The selection of single rectangle feature that suits the best positive and negative samples in the image database. Several machine learning classifier are used for selection. But, the Adaboost classifier algorithm provides the best results since the training error is zero for exponential rounds. The utilization of weak classifier decides the threshold function in such a way that number of misclassified examples is minimum. The direction of inequality is expressed by classifier, feature and threshold as follows:

$$h_j(x) = \begin{cases} -1 & f_j(x) < \theta_j \\ 1 & Otherwise \end{cases} \rightarrow (7)$$

The strong classifier is obtained by combining all the weak classifiers by using Equation 8

$H\ (x) = sign\ (\sum N_j = 1 W_j h_j\ (x)) \rightarrow (8)$

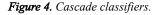The algorithm for computation of strong classifier is described as follows:

---

**Adaboost learning algorithm**

Input: Initialize the number of positive and negative training images ($x_j$, $y_j$) as N

Output: Strong classifier $h(x)$

Steps

| | |
|---|---|
| 1 | Assign weight $W^i_1$ to image I ($W^i_1$=1/N) |
| 2 | For each feature (j=1,2,...M) |
| 3 | Normalize weight $W_M$= $W^i_1$/$\sum^M_{j=1}w$ |
| 4 | Train the classifier $h_j$ with respect to $W_M$ |
| 5 | Estimate the error $E_j$=$\sum_i W_j\mid h_j(x_i)-y_j\mid$ |
| 6 | Choose the classifier with the lowest error |
| 7 | Assign the weight $W_j$ to the classifier $h_j$ |
| 8 | Update the weights |
| 9 | End*for* |
| 10 | Estimate the strong classifier $h(x)$ |

---

The detection of rectangle features by adaboost is meaningful and provides the easy interpretation. The detection performance needs to be further improvement. Hence, the classifiers are cascaded in the next stage to improve the detection performance with less computation time. The construction of boosted classifiers assures an effective removal of negative sub-windows and detects all the positive instances.

**Cascade classifier:** The employment of simple classifier rejected the majority of sub-window prior to complex classifier achieves the low false positive rates. The detection rate is improved to high by the sequential triggering of next classifier with the present classifier. The positive response from each classifier triggered the next classifier. If any negative response is obtained from any stage, the corresponding sub-window is rejected as shown in Figure 4.



*Figure 4. Cascade classifiers.*

The utilization of Adaboost training classifiers constructs each stage and adjusts the threshold value for minimization of false negatives. The initial classifiers in the cascade form do the elimination process that removes the large negative samples. The consecutive layers in each stage perform the additional computation task for the elimination of additional negatives. The sequential stages are applicable for the radial reduction of sub-windows. The detection and false positive rates are high due to the low threshold value. The definition of optimum

cascade framework is depends upon the principle of provision of trade-off requires among the number of classifier stages, number of features and threshold in each stage. The cropped face from the input video is shown in Figures 5a-5c.
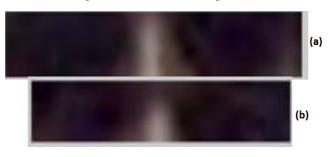


*Figure 5. Viola-Jones face detection for (a) Both users, (b) User 1, and (c) User 2.*

## Eye detection

The region of eye is darker than the other parts in human face. Based on this property, Viola-Jones algorithm searches the small patches in an input image which are darker than the other neighbourhood parts. The model uses the property that the center part of eye region is greater than the other parts in discarding of regions corresponds to eye brow. The utilization of Iris geometrical information in determination of whole region improved the alignment of two eyes in a same line. The employment of pupil detection improves the performance of iris recognition system. The unique feature of an individual governs the automatic identification in bio-metric system.

Once the face detection by Viola-Jones is over, the upper half of the face is cropped from an image. Then, the face is split up into two parts as left and right. Then, apply the Viola-Jones algorithm for left and right eye detection based on pupil based iris recognition scheme. Among various eye detection methods in research studies, we prefer the Viola-Jones method, since the detection process is easy and accurate one. But, the pupil of the eye is not necessarily being a circle in iris recognition. A small error in eye detection may cause the required information to be loss. The rotation of head or eye also creates the problems in pupil detection that leads to introduction of failures in iris recognition scheme. The eye detection output of both user 1 and 2 from the input video is shown in Figures 6a and 6b.



*Figure 6. Eyes of (a) User 1, (b) User 2.*

## Feature extraction

The human interpretation of informative, non-redundant derived values from the initial set of measured data requires the feature extraction. In general, the data input to the image

processing algorithm is large, which are transformed into the reduced set of features called feature extraction. Three types of feature extractions are employed as follows:

1. Histogram of Oriented Gradient (HoG) features
2. Local Landscape Portrait (LLP)
3. Color based features

**HOG features:** Histogram of oriented gradient features plays a vital role in classification and object recognition framework. The inclusion of edge direction and intensity gradients described the appearance and shape of an object. Initially, the image is decomposed into number of cells and computes the histogram of each pixel within the cell. The contrast-normalization with the measure of large region in an image is performed and by using this value all the cells are normalized. The extraction process includes sequential stages such as gradient computation, orientation binning, descriptor blocks creation and block normalization. The first stage requires the filtering governed by the following two kernels

$$K_x = [-1\ 0\ 1] \rightarrow (9)$$

$$K_y = \begin{matrix} 1 \\ 0 \\ -1 \end{matrix} \rightarrow (10)$$

The $X$ and $Y$ coordinates of a given image is computed from the convolution of image with the gradients

$$l_x = l \times K_X \rightarrow (11)$$

$$l_Y = l \times K_Y \rightarrow (12)$$

The magnitude and the orientation of image gradients computed by following equations

$$|M_G| = \sqrt{(l_X)^2 + (l_Y)^2} \rightarrow (13)$$

$$|O_G| = \tan^{-1}\frac{l_Y}{l_X} \rightarrow (14)$$

Based on the gradient values, each pixel in a cell caste a weighted vote and spread the HOG channel in two oriental directions such as 0˚ to 180˚ or 0˚ to 360˚ constitutes orientation binning. The normalization of gradient strengths in which grouping of cells into larger spatial connected blocks created the two descriptor blocks namely, rectangular or circular HOG blocks. Let us consider the non-normalized vector containing all the histograms in a given block. Then the normalization of given block with the additional constant value *(e)* described as follows:

$$L_2 norm = \frac{v}{\sqrt{\||v\||_2^2 + e^2}} \rightarrow (15)$$

The cosine similarity measurement between the vectors $V_1$ and $V_2$ corresponding to pixel in each block and other block describes the necessary HOG features described by

$$Sim = \cos(\theta) = \frac{v_1 \cdot v_2}{\||v_1\|| \cdot \||v_2\||} \rightarrow (16)$$

The orientation features obtained from gradient perspective are given to the classification module to measure the user interest.

**Local landscape portrait (LLP) features:** Here, the landscape and portrait features are extracted from the image by using proposed Local Landscape Portrait (LLP) algorithm based on the orientations. In this work, the local landscape and portrait patterns are extracted for each image.

---

**LLP algorithm**

*Input: Image*

*Output: Texture pattern image*

*Divide Image in to 3 × 3 window*

*Calculate average of diagonal pixels*

*for i=2: m-1*

*for j=2: n-1*

*temp=im (i-1: i+1, j-1: j+1, cc);*

*ts=mean ([[temp (2) temp (4) temp (6) temp (8)]);*

*temp1 (1)=im (i-1, j-1)>ts;*

*temp1 (2)=im (i-1, j)>ts;*

*temp1 (3)=im (i-1, j+1)>ts;*

*temp1 (4)=im (i, j+1)>ts;*

*temp1 (5)=im (i+1, j+1)>ts;*

*temp1 (6)=im (i+1, j)>ts;*

*temp1 (7)=im (i+1, j-1)>ts;*

*temp1 (8)=im (i, j-1)>ts;*

*temp1 (9)=im (i, j)>ts;*

*Endfor*

*Endfor*

*Replace all pixels by the new pattern*

*Get the texture pattern image*

*End*

---

The following steps illustrate the procedure for feature vector creation,

Initially, the window represents the image is divided into 3 × 3 cells and calculate the average for both pixels in landscape and portrait portions.

- Depends upon the status of average pixel value either greater than the neighbor pixel value or equal, the cell is filled with either 1 or 0 respectively.
- The cell representation provides the 8-bit number.
- The number of ones and zeros in the cell generates pattern in such a way that the overall count of ones placed in left part and count of zeros in left part.
- On the basis of combination of cells with smaller and greater values with respect to center, the histogram is computed.

- Then, do the normalization and concatenation in the computed histogram to extract the feature vector for each window.
- The utilization of LLP algorithm accurately predict the texture features with the less bin combinations and also supports the multi-user facial behaviour analysis, which is the novelty of this paper.
- Finally, $1 \times 256$ values from the image are obtained. The graphical representation of working process of LLP algorithm is depicted as follows

For instance, the diagrammatical representation is shown in below:

Input: $3 \times 3$ window.

| 1 | 4 | 3 |
|---|---|---|
| 2 | 5 | 7 |
| 3 | 9 | 5 |

**Step 1:** Find the average values for the horizontal, vertical and center pixels 2, 4, 7, and 9. 2+4+7+9 = 22/4 = 5.5. Here, the round of value is considered as 5.

| 1 | 4 | 3 |
|---|---|---|
| 2 | 5 | 7 |
| 3 | 9 | 5 |

**Step 2:** If the cell value is greater than the average value, then put the value as 1; otherwise, put the value as 0.

| 0 | 0 | 0 |
|---|---|---|
| 0 | 1 | 1 |
| 0 | 1 | 1 |

**Step 3:** The bit combinations obtained from previous step is 0101 which provides ($2^4$=16) combinations. The first 2 bits represents landscape and the next two bits represent portrait orientations of an image. The Local Binary Pattern (LB) requires 256 combinations. The reduction of combinations by LLP results in the reduction of steps involved with maximum accuracy. The HOG based features describe the shape of an image and the LLP features describe the texture of an image. To extract color based features, various color models and descriptors are available. The next section illustrates the detailed description of color feature extraction. The LLP output of both the users by using LLP method is shown in Figures 7a and 7b.

**Color based features:** We consider RGB model to extract the color features by using several color descriptors. This section discusses the color descriptors one by one. The color feature extraction includes following key components color space, quantification and similarity measurement. Consider the image of size. Mean, variance and standard deviation are the color moments for an image. They are defined with the individual and overall pixel values ($p_{ij}$, $P_{ij}$) as follows

$$\mu = \sum_{i=1}^{n} \sum_{j=1}^{m} \left(\frac{P_{ij}}{mn}\right) \rightarrow (17)$$

$$var = \sigma^2 = \frac{1}{nm} \sum_{i=1}^{n} \sum_{j=1}^{m} \left(P_{ij} - mean\right)^2 \rightarrow (18)$$

*Std dev=$\sigma \rightarrow (19)$*

The distribution of several pixels for an image referred as histogram. The elements of histogram depend upon the bits in pixels. Here, the 8-bit representation is considered. Hence, the mean and variance with color histogram for the maximum values of ($2^8$) 256 are described as follows

$$\mu_H = \sum_{k=0}^{255} k \times H(k) / \sum_{k=0}^{255} h(k) \rightarrow (20)$$

$$var_h = \sqrt{\frac{\sum_{k=0}^{255} H(k) * (k - mean)^2}{\sum_{k=0}^{255} H(k)}} \rightarrow (21)$$

Then, the skewness and Kurtosis is estimated by following equations

$$skewness = E\left[\left(\frac{x-\mu}{\sigma}\right)^3\right] \rightarrow (22)$$

$$Kurtosis = E\left[\left(\frac{x-\mu}{\sigma}\right)^4\right] \rightarrow (23)$$

After all the features are extracted, the classification is the next process for perfect recognition of an image.
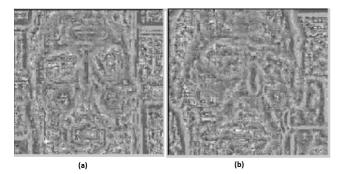


*Figure 7. LLP patterns for (a) User 1, and (b) User 2.*

### *PNN classification*

The predominant classifier that maps all the input patterns into various classifications referred as Probabilistic Neural Network (PNN) classifier. The multi-layered feed forward network performs the classification as shown in Figure 8. The results obtained from the Viola-Jones method and feature extraction methods for each user provided as the input to PNN network.

The weight function corresponding to the user interest measures assigned to each neuron in the pattern layer. The integration of all the individual scores done in the summarization layer. Finally, the emotional and attention states are provided by the output layer. The Gaussian function to be convoluted with the grayscale iris image input to find the deviation from the center of the eye region is described as follows:
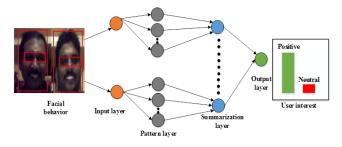


**Figure 8.** *PNN layers.*

$$g(x, y) = \frac{1}{2\pi^{p/2}\sigma^p} e^{-\frac{(x - X_0)^2 + (y - y_0)^2}{\sigma^2}} \rightarrow (24)$$

The normal distribution parameter is chosen with respect to iris size for the length of vector. The function in Equation 24 corresponds to single population. For the number of populations, the Equation 24 modified as

$$g(x, y) = \frac{1}{2\pi^{p/2}\sigma^p} \sum_{i=1}^{N} e^{-\frac{(x - X_{i0})^2 + (y - y_{i0})^2}{\sigma^2}} \rightarrow (25)$$

The larger head motion results higher user interest. Hence, the head motion score is evaluated by the displacement measure *(m (t))* between two consecutive frames with the constant factor (*w*=200) described as follows:

$$H_m(t) = e^{-(m(t))^2}/w \rightarrow (26)$$

The saccade and blink detection scores are very useful in classification of user interest. The distance between the center of eyeball and the corner influence in the estimation of degree of eye movement. The comparison between the eye movement between previous and present frame is larger than the threshold value means saccade is detected. To estimate the score, we define the status of eye whether it is open or closed for the distances of upper boundary to the iris center *(d_i)* and to eyelid point *(d_e)* as follows

$$b(t) = \begin{cases} open & d_i \geq d_e \\ close & else \end{cases} \rightarrow (27)$$

The transformation of blink *(b (t))* and saccade status *(s (t))* into quantitative scores *(B_S (t), S_S (t))* depends upon whether blink/saccade are detected in a window *(W)* over the time *t*. The relationship is described by the Equation 28

$$B_S(t) = \begin{cases} 1 & \sum_{t \in W} b(t) \leq 1 \\ 0 & otherwise \end{cases} \rightarrow (28)$$

$$S_S(t) = \begin{cases} 1 & \sum_{t \in W} s(t) = 0 \\ 0 & otherwise \end{cases} \rightarrow (29)$$

The combination of head motion, blink and saccade scores constitutes the attentive and emotional states. The part in the frame is viewed consider as an attentive for more than two scores are high. Otherwise, it can be regarded as inattentive states. By using the deviation in the positions of eyes and head motion for the consecutive frames by a PNN we classify the human expressions as positive and neutral. The facial expression recognition transformed into the quantitative measure called emotional score of each frame $E_f$ *(t)* over the time *t*. The integration of all the scores by a PNN denoted as $S_R$. Depends upon the integrated score or interest score, the most appropriate clip is selected in $i^{th}$ video shot. The sum of interest score for $j^{th}$ frame for each video shot of the length (*l*) is depicted as follows:

$$\sum_{k=j}^{j=1} S_R(k) \rightarrow (30)$$

The frames with the maximum score (max ($z_j$) are arranged to form a video that is the required summarized form based on user interest. Figures 9 and 10 show the summarized frames and corresponding user behaviour.



**Figure 9.** *Summarized video.*

## Performance Analysis

This section discusses the analysis of performance parameters of success rate, precision, recall, F-measure, and interest score for proposed Multi-user Facial Expression based Video Summarization (MFE-VS) and compare with the existing methods. The optimal way of texture computation using Local

Landscape Portrait (LLP) provides the less number of computations. Besides, a fast and accurate classifier called Probabilistic Neural Network (PNN) used to show the Multi-user facial expressions i.e. positive or neutral with less computation time.

### A. Success rate

The accuracy of proposed work is defined as the ratio of number of correct classified emotion states to the total number of states based on detection rate of eye and head motions in the video summarization. The more success rate provides the betters video summarization based on user behaviour due to the optimal texture computation using LLP algorithm in the bit oriented models. Figure 6 depicts the variation of success rate for proposed MFE-VS, and traditional methods namely, particle, TLD, CT, KSP, KCF, blob associate and blob optimization methods [33-39].
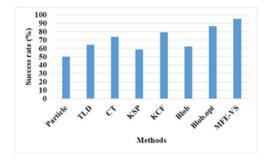


*Figure 10. User positive emotion.*



*Figure 11. Success rate analysis.*

Figure 11 shows the effectiveness of proposed MFE-VS over the traditional methods. The proposed MFE-VS model provided the 8.42% more success rate than the blob optimization methods that ensures the suitability of MFE-VS based on user behaviour.

### B. Precision, recall and F-measure

In this section, we investigate the variation of precision, recall, and F-measure for proposed and existing models. Figures 12-14 depicts the variation of measures for various frames. In general, the calculation of precision and recall depends upon the positive and negative rates (TP, TN, FP, and FN) as follows:

$Precision=TP/(TP+FP) \rightarrow (31)$

$Recall=TP/(TP+FN) \rightarrow (32)$

$\text{F}-\text{measure} = 2\frac{Precision \times Recall}{Precision + Recall} \rightarrow (33)$

Figures 12-14 describes that, the precision of proposed MFE-VS is 7% more compared to existing Feature Aggregation (FA) method which provides better values than DT, STIMO, and VSUMM [28-30]. Also, the percentage difference of recall and F-measure between MFE-VS and FA is 6% and 6% respectively that assures an effective classification.
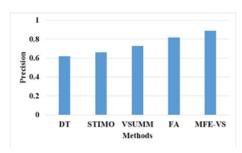

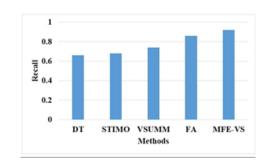
*Figure 12. Precision analysis.*
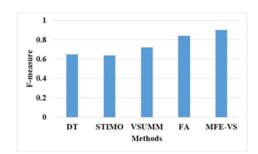


*Figure 13. Recall analysis.*



*Figure 14. F-measure analysis.*

The better classification results influenced in the video summarization. Thus, the MFE-VS are more suitable in the platforms where the number of computation steps is minimum.

### C. Interest score

The interest score of proposed MFE-VS is defined as the integration of attention, emotional states scores based on detection of eye and head motions in the video summarization. The more interest score influences the better summarization compared to existing methodologies due to the optimal texture computation using LLP algorithm in the bit oriented models.

Figure 15 describes the variation of interest scores for proposed MFE-VS with random selected summaries, NOVOICE and multi-level video summarization provided by You et al. [5].
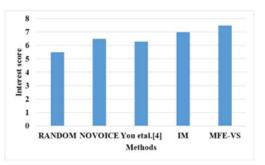


***Figure 15.*** *Interest score analysis.*

Figure 15 describes that the proposed MFE-VS provides 0.5 more interest score than the interest meter method which provides maximum score compared to other existing lead to an effective summarization of video based on multi-user facial expressions.

## Conclusion

This paper addressed the necessity of Video Summarization (VS) in multimedia with the huge amount of video contents. VS based on facial expression are the critical task in digital video technology. The making of better interpersonal relations and human machine interface with automatic recognition of facial expressions effectively used in behavioural science. Due to the non-uniqueness and the deformations in the human face, this paper proposed an effective toleration of facial variations themselves and addressed the problems in facial expression and video summarization. The proposed work provided the consciousness about the viewer reactions for the measurement of viewer's interest by using quantitative measures in based on the human emotional perspective. In this paper, the face and eye regions are detected by using Viola-Jones face detector. The utilization of Histogram of Gradients (HoG) and the Local Landscape Portrait pattern (LLP) extracted shape and texture features. Then, Probabilistic Neural Network (PNN) classifier predicted the exact user expression in the frames. Based on the obtained expression, the proposed work extracted the interesting frames and summarized referred as final stage. Two states are observed in the human expression such as attention and emotional. The first one denoted by the variations of facial expressions and the second one defined by the classification of positive and neutral states. The comparative analysis of proposed method with the traditional methods on the parameters of accuracy, sensitivity, specificity and recognition rate assured the capability to effectively summarize the video in multi-user environment.

## References

1. Weiming H, Nianhua X, Li L, Xianglin Z, Maybank S. A Survey on visual content-based video indexing and retrieval. IEEE Trans Sys Man Cybern Appl Rev 2011; 41: 797-819.

2. Jasmine RR, Thyagharajan K. Study on recent approaches for human action recognition in real time. Int J Eng Res Technol 2015; 4: 660-664.

3. Sangeetha S, Deepa S. A survey on video summarization using face recognition methods. Int J Adv Res Comp Sci Manag Stud 2013; 52.

4. Wei-Ting P, Wei-Ta C, Chia-Han C, Chien-Nan C, Wei-Jia H, Wen-Yan C. Editing by viewing: automatic home video summarization by viewing behaviour analysis. IEEE Trans Multimed 2011; 13: 539-550.

5. You J, Liu G, Sun L, Li H. A multiple visual models based perceptive analysis framework for multilevel video summarization. IEEE Trans Circ Sys Vid Technol 2007; 17: 273-285.

6. Ashokkumar S, Thyagharajan K. Retina biometric recognition in moving video stream using visible spectrum approach. Green Comp Commun Conserv Energy 2013; 180-187.

7. Meng W, Hong R, Guangda L, Zheng-Jun Z, Shuicheng Y, Tat-Seng C. Event driven web video summarization by tag localization and key-shot identification. IEEE Trans Multimed 2012; 14: 975-985.

8. Mendi E, Clemente HB, Bayrak C. Sports video summarization based on motion analysis. Comp Electr Eng 2013; 39: 790-796.

9. Dang CT, Radha H. Heterogeneity image patch index and its application to consumer video summarization. IEEE Trans Imag Proc 2014; 23: 2704-2718.

10. Fan C, De Vleeschouwer C, Cavallaro A. Resource allocation for personalized video summarization. IEEE Trans Multimed 2014; 16: 455-469.

11. Wang W, Zhang Y, Yan S, Zhang Y, Jia H. Parallelization and performance optimization on face detection algorithm with OpenCL: A case study. Tsinghua Sci Technol 2012; 17: 287-295.

12. Zhang D, Li B, Li H. Inter-media-based video adaptation system: Design and implementation. Tsinghua Sci Technol 2012; 17: 113-127.

13. Best-Rowden L, Klare B, Klontz J, Jain AK. Video-to-video face matching: Establishing a baseline for unconstrained face recognition. Sixth International Conference on iometrics IEEE Appl Sys 2013; 1-8.

14. Ling S, Jones S, Xuelong L. Efficient search and localization of human actions in video databases. IEEE Trans Circ Sys Video Technol 2014; 24: 504-512.

15. ElMaghraby A, Abdalla M, Enany O, El Nahas MY. Detect and analyze face parts information using Viola-Jones and geometric approaches. Int J Comp Appl 2014; 101: 23-28.

16. Zhao X, Lin KH, Fu Y, Hu Y, Liu Y. Text from corners: a novel approach to detect text and caption in videos. IEEE Trans Image Process 2011; 20: 790-799.

17. Nagarajan G, Thyagharajan KK. Rule-based semantic content extraction in image using fuzzy ontology. Int Rev Comput Softw 2014; 9: 266-277.

18. Thanikachalam KKTV. Human action recognition using temporal partitioning of activities and maximum average correlation height filter. Res J Appl Sci Eng Technol 2015; 11: 105-111.

19. Nagarajan G, Minu RI. Multimodal fuzzy ontology creation and knowledge information retrieval. Proc Int Conf Soft Comp Sys Springer India 2016.

20. Merler M, Huang B, Xie L, Gang H, Natsev A. Semantic model vectors for complex video event recognition. IEEE Trans Multimed 2012; 14: 88-101.

21. Nagarajan G, Minu RI. Fuzzy ontology based multi-modal semantic information retrieval. Procedia Comp Sci 2015; 48: 101-106.

22. Meng W, Hong R, Xiao-Tong Y, Shuicheng Y, Chua TS. Movie2Comics: towards a lively video content presentation. IEEE Trans Multimed 2012; 14: 858-870.

23. YingLi T, Liangliang C, Zicheng L, Zhengyou Z. Hierarchical filtered motion for action recognition in crowded videos. IEEE Trans Sys Man Cybern Appl Rev 2012; 42: 313-323.

24. Thanikachalam T. Human action recognition using motion history image and correlation filter. Int J Appl Eng Res 2015; 10: 361-363.

25. Gkonela C, Chorianopoulos K. Video skip: event detection in social web videos with an implicit user heuristic. Multimed Tools Appl 2014; 69: 383-396.

26. Ejaz N, Mehmood I, Baik SW. Feature aggregation based visual attention model for video summarization, Comp Electr Eng 2014; 40: 993-1005.

27. Ejaz N, Mehmood I, Wook Baik S. Efficient visual attention based framework for extracting key frames from videos. Sig Proc Imag Commun 2013; 28: 34-44.

28. Mundur P, Rao Y, Yesha Y. Keyframe-based video summarization using Delaunay clustering. Int J Dig Libr 2006; 6: 219-232.

29. Furini M, Geraci F, Montangero M, Pellegrini M. STIMO: STIll and moving video storyboard for the web scenario. Multimed Tools Appl 2010; 46: 47-69.

30. de Avila SEF, Lopes APB, da Luz A, de Albuquerque Araujo A. VSUMM: A mechanism designed to produce static video summaries and a novel evaluation method. Pat Recogn Lett 2011; 32: 56-68.

31. Yong SP, Deng J, Purvis M. Wildlife video key-frame extraction based on novelty detection in semantic context. Multimed Tools Appl 2013; 62: 359-376.

32. Qixiang Y, Jixiang L, Jianbin J. Pedestrian detection in video images via error correcting output code classification of manifold subclasses. IEEE Trans Intel Transp Sys 2012; 13: 193-202.

33. Lin W, Zhang Y, Lu J, Zhou B, Wang J, Zhou Y. Summarizing surveillance videos with local-patch-learning-based abnormality detection, blob sequence optimization, and type-based synopsis. Neurocomput 2015; 155: 84-98.

34. Hess R, Fern A. Discriminatively trained particle filters for complex multi-object tracking. IEEE Conf Comp Vis Pat Recogn 2009; 240-247.

35. Kalal Z, Mikolajczyk K, Matas J. Tracking-learning-detection. IEEE Trans Pat Anal Mach Intel 2012; 34: 1409-1422.

36. Zhang K, Zhang L, Yang MH. Real-time compressive tracking. Comp Vis 2012; 864-877.

37. Berclaz J, Fleuret F, Turetken E, Fua P. Multiple object tracking using k-shortest paths optimization. IEEE Trans Pat Anal Mach Intel 2011; 33: 1806-1819.

38. Henriques JF, Caseiro R, Martins P, Batista J. High-speed tracking with kernelized correlation filters. IEEE Trans Pattern Anal Mach Intell 2015; 37: 583-596.

39. Su X, Lin W, Zheng X, Han, Chu H, Zhang X. A new local-main-gradient-orientation hog and contour differences based algorithm for object classification. IEEE Int Symp Circ Sys 2013; 2892-2895.

40. Betancourt A, Morerio P, Regazzoni CS, Rauterberg M. The evolution of first person vision methods: a survey. IEEE Trans Circ Sys Video Technol 2015; 25: 744-760.

41. Ryoo M, Fuchs TJ, Xia L, Aggarwal J, Matthies L. Early recognition of human activities from first-person videos using onset representations. ArXiv 2014; 5309.

42. Sigari MH, Sureshjani SA, Soltanian-Zadeh H. Sport video classification using an ensemble classifier. 7th Iranian Mach Vis Imag Proc 2011; 1-4.

43. Saftoiu A, Vilmann P, Gorunescu F, Janssen J, Hocke M, Larsen M. Efficacy of an artificial neural network–based approach to endoscopic ultrasound elastography in diagnosis of focal pancreatic masses. Clin Gastroenterol HepatoL 2012; 10: 84-90.

44. Li J, Hong S, Xia S, Luo S. Neural network based popularity prediction for IPTV system. J Netw 2012; 7: 2051-2056.

45. Savchenko AV. Probabilistic neural network with homogeneity testing in recognition of discrete patterns set. Neur Netw 2013; 46: 227-241.

46. Karmokar PR, Parekh R. Recognition of semantic content in image and video. Int J Comp Appl 2013; 73: 31-35.

47. Khashei M, Zeinal Hamadani A, Bijari M. A novel hybrid classification model of artificial neural networks and multiple linear regression models. Exp Sys Appl 2012; 39: 2606-2620.

## *Correspondence to

Priya S

Department of Electrical and Electronics Engineering

S.A. Engineering College

India