

Improving Video Summarization based on User Preferences

Rajkumar Kannan

Department of Computer Science
Bishop Heber College (Autonomous)
Tiruchirappalli, India
dr.rajkumarkannan@gmail.com

Gheorghita Ghinea

Information Systems and Computing
Brunel University
London, United Kingdom
george.ghinea@brunel.ac.uk

Sridhar Swaminathan¹,
Suresh Kannaiyan²

Department of Computer Science
Bishop Heber College (Autonomous)
Tiruchirappalli, India
¹sridarah@gmail.com,
²sureshk.naga@gmail.com

Abstract—Although in the past, several automatic video summarization systems had been proposed to generate video summary, a generic summary based only on low-level features will not satisfy every user. As users' needs or preferences for the summary vastly differ for the same video, a unique personalized and customized video summarization system becomes an urgent need nowadays. To address this urgent need, this paper proposes a novel system for generating unique semantically meaningful video summaries for the same video, that are tailored to the preferences or interests of the users. The proposed system stitches video summary based on summary time span and top-ranked shots that are semantically relevant to the user's preferences. The experimental results on the performance of the proposed video summarization system are encouraging.

Keywords— Video summarization; video semantics; personalization; user preferences

I. INTRODUCTION

The past decade had witnessed an explosive growth of multimedia data such as videos and still images, both in the internet and in home computers. So, browsing through lengthy and voluminous video collection becomes tedious to the user, if the videos have little or no relevant content. Also, searching for interesting segments or shots within the videos is time consuming. Video summarization aims at producing compact version of a full-length video while preserving the significant content of the original video. For generating a video summary, most of the *automatic video summarization* methodologies detect significant segments of a video based on certain criteria which are mostly low-level audio-visual features. However users always try to summarize videos based on the semantic content of a video, rather than low-level features alone. Generic video summarization will not be sufficient when the users' needs and interests change over a time. Users are seldom satisfied by common video summary produced by a video summarization system. This is because, the produced video summary may not contain content of particular semantic concept or genre liked by the user. So the criterion used to summarize a video should be the user's preferences and interests over the video semantics. Personalized video summarization is a useful technique for producing customized video summaries to the users based on their needs.

Therefore our hypothesis is that a generic video summary does not satisfy every user and multiple video summaries for

the same video should be generated depending on the preferences and needs of individual users. Hence, this paper proposes a system for personalized video summarization that produces customized video summaries by adapting to the user's interest. The proposed system uses high-level feature extraction to reduce the manual video annotation. Different summarization techniques are employed to the test videos and evaluated with real users.

II. RELATED WORK

A wide number of contributions have been found in the area of video summarization. In [1] important people and objects in video obtained from a wearable camera is detected using region based regression for generating a storyboard summary. The authors in [2] proposed an SVM based structure learning for transferring knowledge between video and textual features for summarization. Two important criterions of a summary - *coverage* and *diversity* are considered for summarizing unconstrained videos, where summarization is treated as combinatorial optimization problem [3]. A multi video summarization framework is proposed for videos retrieved by event based queries from social media websites [4]. It utilizes the near duplication among visual content and user ratings of retrieved videos for calculating the *informative scores* to rank shots. Most of the previous methodologies proposed for summarization are based only on low-level features. Also, methods which used high-level features of a video ignored users' interests and preferences. In contrast to previous works, this paper presents a methodology for personalized semantic video summarization.

III. SYSTEM OVERVIEW

The architecture of the proposed video summarization system is shown in fig. 1. The system consists of three modules: *preprocessing*, *user interface* and *video summarization*. Here, the database contains collection of videos and their metadata. This proposed personalized video summarization system can generate summary for videos of any genre from any domain.

In preprocessing module, a video is segmented using shot boundary detection method proposed in [5], and a single key frame is extracted from each shot. Since a keyframe can represent a shot, near duplicate shots in a video are identified

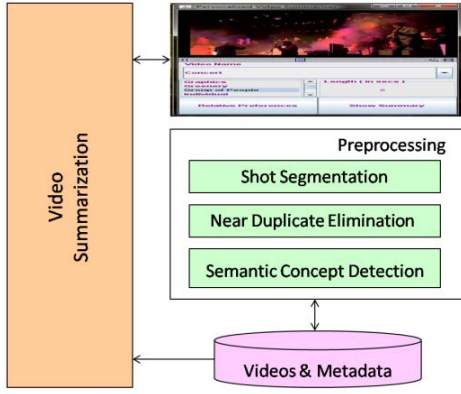


Fig. 1. Architecture of personalized video summarization.

using these keyframes. In a set of visually similar shots, shot that appears first in the video is kept and other near duplicate shots are removed, so that the summary will not contain more than one shot with similar visual content. High-level features are extracted from each shot and they constitute a set of preferences for the users. A set of 25 semantic concept detectors were developed for concepts such as *beach*, *flower scene*, *indoors*, etc., using method proposed in [6]. *Relevance scores* to each video shot for a set of semantic concepts are assigned. This score ranging from -1 to +1, shows the relevance between a shot and a particular semantic concept, where -1 implies *highly irrelevant* and +1, *highly relevant*.

Fig. 2 depicts the user interface of the video summarization system. User Interface module allows the user to construct their profile with multiple preferences. They are set of high-level semantic concepts used in preprocessing. Set of semantic concepts selected by the user as preferences will be considered as *user profile*. For the input video, users can select semantic concepts from the set of 25 semantic concepts and are also allowed to give length or duration of summary via the user interface. The users can choose preferences using either *list* or

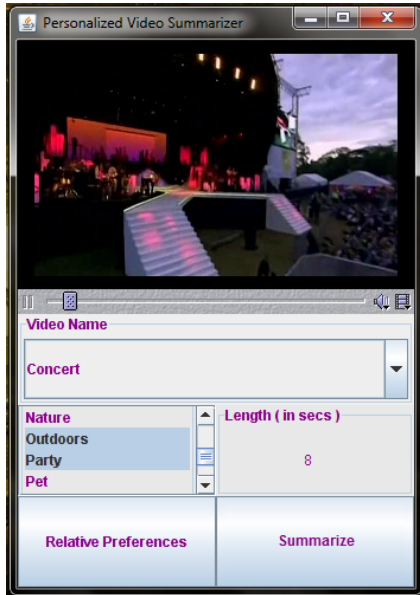


Fig. 2. User interface of the video summarization system.

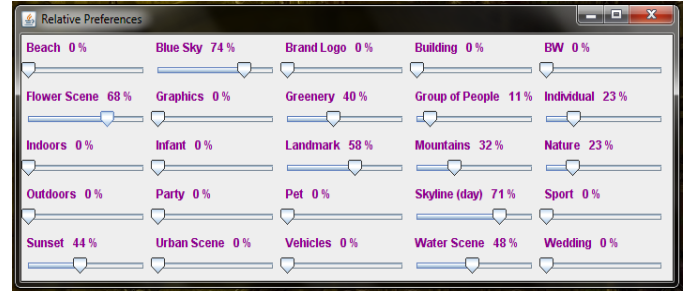


Fig. 3. Relative user preference panel of the video summarization system.

slider. When using list, all chosen preferences will assume a numeric value 1 and rest will assume a numeric value 0. Since list allows only the binary preferences, sliders can be used to select real values from 0 to 1 (fig. 3).

In video summarization module, *Dot Product* and *Cosine Similarity* metrics are employed to determine likeness between user's profile and actual shot content. Summary is generated based on ranking the video shots using similarity measures and shown to the user.

IV. SUMMARIZATION METHODOLOGY

Let video $V = \{u_i, 1 \leq i \leq n\}$ consists of n shots, where each shot u_i has duration d_i seconds. Let $C = \{c_j, 1 \leq j \leq m\}$ denote a set of m semantic concepts that were used in preprocessing. For a shot u_i , let $R_i = [r_{i1}, r_{i2}, r_{i3}, \dots, r_{im}]$ denote *model vector*, containing relevance scores of m semantic concept detectors. Let $P = [p_1, p_2, p_3, \dots, p_m]$ denote a vector, which consist of a set of high-level semantic concepts that are selected as preferences by the user. Each preference p_j takes value between 0 and 1. Let T denote summary time (in seconds) given by the user. Let S_i be *similarity score* computed for shot u_i .

A. Shot Ranking

Shots are ranked based on the similarity between relevance scores and user profile using either Dot Product or Cosine Similarity. For a shot u_i , the Dot Product between model vector R_i and preference vector P gives the similarity score S_i . That is,

$$S_i = \vec{R}_i \cdot \vec{P} \quad (1)$$

For a shot u_i , Cosine Similarity can be computed between model vector R_i and the preference vector P . That is,

$$S_i = \frac{\vec{R}_i \cdot \vec{P}}{\|\vec{R}_i\| \times \|\vec{P}\|} \quad (2)$$

B. Shot Selection

The objective of shot selection is to select shots which maximize the cumulative similarity score for the summary while not exceeding the time constraint T . This can be considered as an instance of 0-1 knapsack problem which is defined as,

$$\max \sum_{i=1}^n S_i \cdot x_i$$

$$\text{subject to } \sum_{i=1}^n d_i \cdot x_i \leq T \quad (3)$$

S_i is similarity score computed for shot i and x_i is a binary decision variable that takes value 1 if u_i is selected for summary, otherwise 0. Selected shots are decreasingly ordered based on their similarity scores. Summary is skimmed by concatenating the selected shots, and showed to the user with their corresponding audio.

V. EXPERIMENTS AND RESULTS

The proposed summarization system is evaluated both quantitatively and subjectively. Experiments were conducted on 10 song videos of total duration of 52 mins which were collected from various web sources. Total of 1240 video shots were manually labeled for validation.

A. Experimental Setup

Since the relevance scores are between -1 and +1, and the system allows multiple preferences, a higher negative relevance score for a preferred semantic concept will reduce the similarity score in the Dot Product, even though a shot has many higher positive scores for other semantic concepts (*false negatives*). This will also increase the chances for the shots with lower negative relevance scores of semantic concepts to enter in the summary (*false positives*). So, four types of summarization techniques are experimented. They are,

- Dot Product without negative relevance score (DP+)
- Dot Product with negative relevance score (DP-)
- Cosine Similarity without negative relevance score (CS+)
- Cosine Similarity with negative relevance score (CS-)

Techniques 2 and 4 do not use negative relevance scores and assumed zero.

B. Quantitative Evaluation

For each video, four different preferences (single preference or multiple preferences) were considered. The results are evaluated using *Ranked Precision*. Precisions at intervals n are averaged for all the preferences.

Fig. 4 shows comparison of different similarity measures and their average of precisions when using single semantic concept as preference. Dot Product similarity without negative score (DP+) performs better than others. As n increases, DP+ maintains high precision by limiting the number of shots displayed to the user.

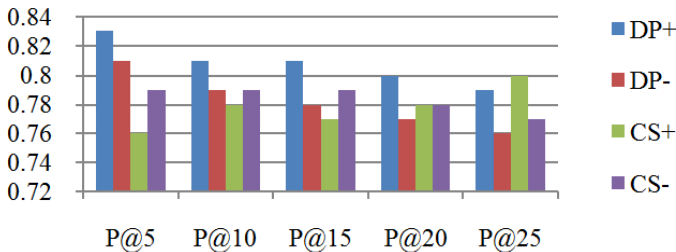


Fig. 4. Average of precisions for different single preferences.

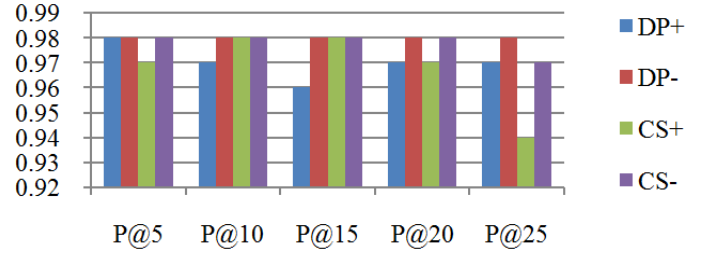


Fig. 5. Average of precisions for different multiple preferences.

The average of precisions when using multiple semantic concepts as preferences is shown in fig. 5. It shows that Dot Product with negative relevance (DP-) score performs better than the other summarization techniques.

Different multiple preferences for different videos were given to the system by using similarity measure as DP-. Averages of precisions at different recalls are calculated for the system (fig. 6). Precision falls abruptly when recall reaches 0.3.

To measure ranking efficiency, top ranked shots are manually graded using a scale of 0-3 to measure the average of *Normalized Discounted Cumulative Gain (nDCG)* at various positions. Fig. 7 shows that as the result set increases, nDCG decreases gradually.

Fig. 8 shows the keyframes of the top 25 shots that are retrieved for a user preference “flower scene” from a test video. It also shows some of the false positives in the resultant summary. This misprediction happens because, color values and distribution of the actor’s costume in that particular keyframe somehow resembles flower.

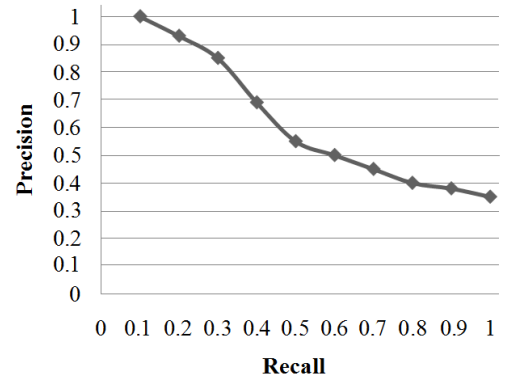


Fig. 6. Precision-Recall curve for different multiple preferences.

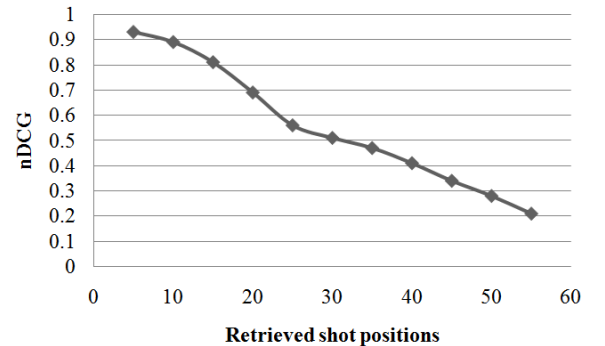


Fig. 7. Normalized Discounted Cumulative Gain for different multiple preferences.



Fig. 8. Keyframes of the top 25 shots for the preference ‘flower scene’.

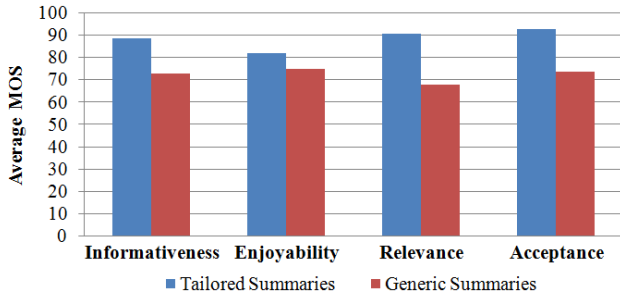


Fig. 9. Average mean opinion scores for tailored and generic summaries.

C. Subjective Evaluation

Performance of the system was evaluated subjectively using questionnaire by 20 test subjects. Among them, 12 were males and the rest were females. Generic summaries consisting of significant shots for each video was created. Subjects were asked to use the system. Summarization methodology was set as DP-. Subjects were not informed about the methodology used for summarization. When using the system, subjects were also shown generic summaries for the videos. Questions were asked about the *informativeness*, *enjoyability*, *relevance* and *acceptance* of personalized and generic summaries. Users were asked to rate each performance criterion on a scale of 1 to 100.

Fig. 9 shows the comparison of *Mean Opinion Score* (MOS) given by the subjects to evaluate likability of the summaries generated by the corresponding summarization type. Quality of summaries was assessed by two *Quality of Perception* measures, *QoP-IA* (the user’s ability to assimilate information) and *QoP-S* (the user’s satisfaction) [7]. *QoP-IA* is measured by averaging the scores for relevance and information acquired from summaries. *QoP-S* is calculated by averaging the scores of enjoyability and acceptance. For tailored summaries *QoP-IA* is 90.4 and *QoP-S* is 87.1. For generic summaries *QoP-IA* is 70.2 and *QoP-S* is 74.9. So, it can be seen that tailored summaries are both informative and satisfactory than the generic summaries.

The attractiveness and usability of the system was analyzed with MOS that measures the user satisfaction and comfortability of the user interface. *Computer System Usability*

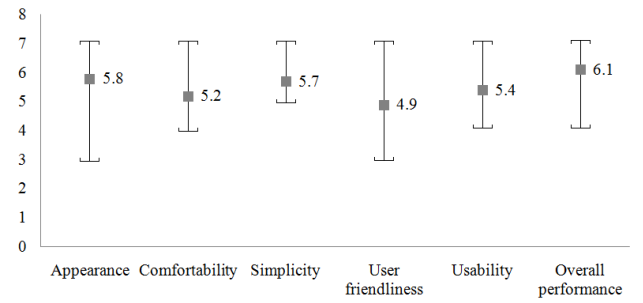


Fig. 10. Subjects’ opinion on summarization system usability.

Questionnaire (CSUQ) [8] was used to measure *Degree of Satisfaction* (on a scale of 1 to 7) for the system and user interface. Fig. 10 show the box plot indicating the average of MOS given by the subjects for system usability under each usability criteria such as *Appearance*, *Comfortability*, *Simplicity*, *User friendliness*, *Usability* and *Overall performance*.

VI. CONCLUSION

This paper presented a novel preference aware video summarization system that produces semantically meaningful personalized video summaries by adapting to the user’s interest. Experimental results on personalized video summarization demonstrate the effectiveness of the proposed system. Though the system summarizes video based on the high-level features, users may also be interested in choosing low-level features as preferences such as color and motion. In future, the system would also consider low-level features as users’ preferences.

REFERENCES

- [1] Y. J. Lee, J. Ghosh and K. Grauman, “Discovering important people and objects for egocentric video summarization,” In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 1346-1353, June 2012.
- [2] L. Li, K. Zhou, G. R. Xue, H. Zha, and Y. Yu, “Video summarization via transferrable structured learning,” In Proceedings of the 20th international conference on World wide web, pp. 287-296, March 2011.
- [3] N. Shroff, P. Turaga, and R. Chellappa, “Video précis: Highlighting diverse aspects of videos,” IEEE Transaction on Multimedia, vol. 12, no. 8, pp. 853-868, 2010.
- [4] R. Hong, J. Tang, H.K Tan, S. Yan, C. W. Ngo, and T. S. Chua, “Event driven summarization for web videos,” In WSM '09: Proceedings of the First SIGMM workshop on Social Media, pp. 43-48, 2009.
- [5] J. Mas, and G. Fernandez, “Video shot boundary detection based on color histogram,” Notebook Papers TRECVID2003, Gaithersburg, Maryland, NIST, 2003.
- [6] Y. G. Jiang, C. W. Ngo, and J. Yang, “Towards optimal bag-of-features for object categorization and semantic video retrieval,” In Proceedings of the 6th ACM international conference on Image and Video Retrieval, pp. 494-501, July 2007.
- [7] S. R. Gulliver and G. Ghinea, “Defining the users perception of distributed multimedia quality”, ACM Transactions on Multimedia Computing, Communications, and Application, vol. 2, no. 4, pp 241-257, Nov. 2006.
- [8] J. R. Lewis, “IBM computer usability satisfaction questionnaires: psychometric evaluation and instructions for use”, International Journal of Human-Computer Interaction, vol. 7, no. 1, pp. 57-78. 1995.