

International Journal of Innovative Research in Science, Engineering and Technology

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Visit: www.ijirset.com

Vol. 8, Issue 2, February 2019

A Survey on Automatic Summarization Using Multi-Modal Summarization System for Asynchronous Collections

Shilpadevi Vasant Bhagwat, Sheetal .S. Thokal

M. E Students, Department of Computer Engineering, JSPMs Imperial College of Engineering, Wagholi, Pune, India

Professor, Department of Computer Engineering, JSPMs Imperial College of Engineering, Wagholi, Pune, India

ABSTRACT: In recent years, much work has been performed to summarize meeting recordings, sport videos, movies, pictorial storylines and social multimedia. Automatic text summarization is an essential natural language processing (NLP) application that goals to summarize a given textual content into a shorter model. The fast growth in multimedia information transmission over the Internet demands multi-modal summarization (MMS) from asynchronous combination of text, image, audio and video. This paper represents an MMS framework that utilizes the techniques of NLP, speech processing and OCR watermarking technique to examine the elaborative information contained in multi-modal statistics and to enhance the aspects of multimedia news summarization. The basic concept is to bridge the semantic gaps among multi-modal content. For audio scheme, convert the audio signals into textual format. For visual scheme, extracts text from images using OCR technique. After, the generated summary for important visual information through content-picture matching or multi-modal topic modeling. Finally, all the multi-modal factors are considered to generate a textual summary by maximizing the importance, non-redundancy, credibility and scope through the allocated accumulation of submodular features. The contribution work is to identify theme of visual scheme. The experimental result shows that Multi-Modal Summarization framework outperforms other competitive techniques.

KEYWORDS: Summarization, Multimedia, Multi-Modal, Cross-Modal, Natural Language Processing, Computer Vision, OCR Technique

I. INTRODUCTION

TEXT summarization plays a vital role in our daily life and has been studied for several decades. From information retrieval to text mining, we are frequently exposed to text summarization. Multimedia data (including text, image, audio and video) have increased dramatically recently, which makes it difficult for users to obtain important information efficiently. Multi-modal summarization (MMS) can provide users with textual summaries that can help acquire the gist of multimedia data in a short time, without reading documents or watching videos from beginning to end. When summarizing meeting recordings, sport videos and movies, such videos consist of synchronized voice, visual and captions. For the summarization of pictorial storylines, the input is a set of images with text descriptions. None of these applications focus on summarizing multimedia data that contain asynchronous information about general topics. Intuitively, readers can grasp the gist of the event more easily by scanning the image or the video than by only reading news document, and thus we believe that the multi-modal data will also reduce the difficulty for machine to understand a news event. While most summarization systems focus on only natural language processing (NLP), the opportunity to jointly optimize the quality of the summary with the aid of automatic speech recognition (ASR) and computer vision (CV) processing systems is widely ignored. On the other hand, given a news event (i.e., news topic), multimedia data are generally asynchronous in real life, which means there is no given explicit description for images and no subtitles for videos. Thus, multi-modal summarization (MMS) faces a major challenge in understanding the

International Journal of Innovative Research in Science, Engineering and Technology

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Visit: www.ijirset.com

Vol. 8, Issue 2, February 2019

semantics of visual information. In this work, we present an MMS system that can provide users with textual summaries to help to acquire the gist of asynchronous multimedia data in a short time without reading documents or watching videos from beginning to end. The purpose of this work is to unite the NLP, ASR and CV techniques to explore a new framework for mining the rich information contained in multi-modal data to improve the quality of multimedia news summarization.

Motivation

- The Multi-modal Topic Modeling is that textual descriptions of images often provide important information about semantic aspects (topics), and image features are often correlated with semantic topics.
- Automatically generate a fixed-length textual summary to represent the principle content of the multi-modal data.

II. RELATED WORK

P. Li, J. Ma, and S. Gao, studies the problem of learning to summarize images by text and visualize text utilizing images, which is called Mutual-Summarization [2]. After divide the web image-text data space into three subspaces, namely pure image space (PIS), pure text space (PTS) and image-text joint space (ITJS). Advantages are: In image summarization procedure, map images from PIS to ITJS via image classification model and describe these images utilizing several high level semantic sentences. In text visualization procedure, map text from PTS to ITJS via text categorization model and then give a visual display utilizing images with high confidences in ITJS. Disadvantages are: Need to improve the Mutual-Summarization performance.

J. Bian, Y. Yang, and T.-S. Chua, proposes a multimedia microblog summarization [3] framework to automatically generate visualized summaries for trending topics. Specifically, a novel generative probabilistic model, termed multimodal-LDA (MMLDA) is proposed to discover subtopics from microblogs by exploring the correlations among different media types. Advantages are: Well organized the messy microblogs into structured subtopics. It generates high quality textual summary at subtopic level. It selects images relevant to subtopic that can best represent the textual contents. Disadvantages are: Only focus on summarizing synchronous multi-modal content.

P. Goyal, L. Behera, and T. M. McGinnity, proposes the novel idea of using the context sensitive document indexing to improve the sentence extraction-based document summarization task. In [4] paper, proposes a context sensitive document indexing model based on the Bernoulli model of randomness. Advantages are: The new context-based word indexing gives better performance than the baseline models. Disadvantages are: Need to calculate the lexical association over a large corpus.

G. Evangelopoulos, A. Zlatintsi, A. Potamianos, P. Maragos, K. Rapantzikos, G. Skoumas, and Y. Avrithis, proposes multimodal saliency representations of audiovisual streams [5], in which signal (audio and visual) and semantic (linguistic/textual) cues are integrated hierarchically. Each modality is independently analyzed in individual saliency representations: spectro-temporal for the audio channel, spatio-temporal for the visual channel, and syntactic for the transcribed subtitle text. Advantages are: The produced summaries, based on low-level features and content-independent fusion and selection, are of subjectively high aesthetic and informative quality. Minimum and inverse variance schemes performed best in terms of informativeness, while adaptive shot- and scene-based inverse variance in terms of enjoyability. Disadvantages are: Need to work on non-linear feature correlation algorithms.

T. Hasan, H. Bo'ril, A. Sangwan, and J. H. Hansen, proposes a generic video highlights generation scheme [6] based on information theoretic measure of user excitability was presented. Excitability is computed based on the likelihood of the segmental features residing in certain regions of their joint probability density function space which are considered both exciting and rare. The proposed measure is used to rank order the partitioned segments to compress the overall video sequence and produce a contiguous set of highlights. Advantages are: This scheme utilizes audio excitement and low-level video features. Effectively combine the multi-modal features in video segments and found to be highly correlated with a perceptual assessment of excitability. Disadvantages are: Only work on summarizing synchronous multi-modal content.

Z. Li, J. Liu, J. Tang, and H. Lu, proposes a novel Robust Structured Subspace Learning (RSSL) algorithm [7] by integrating image understanding and feature learning into a joint learning framework. The learned subspace is

International Journal of Innovative Research in Science, Engineering and Technology

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Visit: www.ijirset.com

Vol. 8, Issue 2, February 2019

adopted as an intermediate space to reduce the semantic gap between the low-level visual features and the high-level semantics. Advantages are: The proposed RSSL enables to effectively learn a robust structured subspace from data. The proposed framework can reduce the noise-induced uncertainty.

J. Bian, Y. Yang, H. Zhang, and T.-S. Chua, proposes a multimedia social event summarization framework [8] to automatically generate visualized summaries from the microblog stream of multiple media types. Specifically, the proposed framework comprises three stages: 1) A noise removal approach is first devised to eliminate potentially noisy images. 2) A novel cross-media probabilistic model, termed Cross-Media-LDA (CMLDA), is proposed to jointly discover subevents from microblogs of multiple media types. 3) Finally, based on the cross-media knowledge of all the discovered subevents. Advantages are: Eliminates the potentially noisy images from raw microblog image collection. It generates multimedia summary for social events utilizing the cross-media distribution knowledge of all the discovered subevents. Disadvantages are: Need to extend the cross-media framework for automatically detecting social events and retrieving related candidate microblogs. Need to personalized microblog summarization based on user profile.

W. Y. Wang, Y. Mehdad, D. R. Radev, and A. Stent, proposes a novel matrix factorization approach [9] for extractive summarization, leveraging the success of collaborative filtering. First to consider representation learning of a joint embedding for text and images in timeline summarization. Advantages are: It is easy for developers to deploy the system in real-world applications. Scalable approach for learning low-dimensional embedding's of news stories and images. Disadvantages are: Only work on summarizing synchronous multi-modal content.

Z. Li, J. Tang, X. Wang, J. Liu, and H. Lu, develop a novel approach of multimedia news summarization [10] for searching results on the Internet, which uncovers the underlying topics among query-related news information and threads the news events within each topic to generate a query-related brief overview. hLDA is adopted to discover the hierarchical topic structure from the query-related news articles, and an approach based on the weighted aggregation and max pooling to identify the typical news article for each topic is proposed. A time-bias MST method is developed to thread the subtopics within one topic to give a news summary on each topic in terms of temporal and spatial development. Advantages are: Proposed system can present vivid and comprehensive information conveniently. Readers can quickly understand the information that they require via the multimedia summarization in this system. Disadvantages are: Need to apply on news video.

H. Li, J. Zhu, C. Ma, J. Zhang, and C. Zong, proposes an extractive multi-modal summarization method [11] that can automatically generate a textual summary given a set of documents, images, audios and videos related to a specific topic. The key idea is to bridge the semantic gaps between multi-modal content. Advantages are: It avoids redundant information. It provides good readability. Disadvantages are: This works only limited dataset.

III. OPEN ISSUES

The existing applications related to MMS include meeting record summarization, sport video summarization, movie summarization, pictorial storyline summarization, timeline summarization and social multimedia summarization etc. Meeting recordings, sport videos and movies consist of synchronized voice, visual and captions and pictorial storylines consist of a set of images with textual descriptions. Most of the existing applications focus on synchronous multimodal content, in which images are paired with textual descriptions and videos, are paired with subtitles.

Disadvantages are:

1. The existing applications focus on only synchronous multimodal content.
2. MMS discovers the semantic gap between different modalities.
3. MMS shows less quality of generated summaries.

IV. SYSTEM OVERVIEW

This paper proposes an approach to generate a textual summary from a set of asynchronous documents, images, audios and videos on the same news event, as shown in Fig. 1. Because multimedia data are heterogeneous and contain more complex information than that contained in pure text. The MMS framework shows in Fig. 1. For the audio information contained in videos, obtain speech transcriptions through ASR and design a method to selectively use these

International Journal of Innovative Research in Science, Engineering and Technology

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Visit: www.ijirset.com

Vol. 8, Issue 2, February 2019

transcriptions. For visual information, including the key frames extracted from videos and the images that appear in documents learn the joint representations of text and images with a neural network; then identify the text that is relevant to the image based on text-image matching or multi-modal topic modeling. In this way, audio and visual information can be integrated into a textual summary by joint optimization. Contribution work is, to design an MMS method that can automatically generate a textual summary from a set of asynchronous documents, images, audios and videos related to a specific event. Consider four criteria like, salience, non-redundancy, readability and coverage for visual information that are jointly optimized by the budgeted maximization of submodular functions to select the representative sentences. Bridge the semantic gap between the textual and visual data.

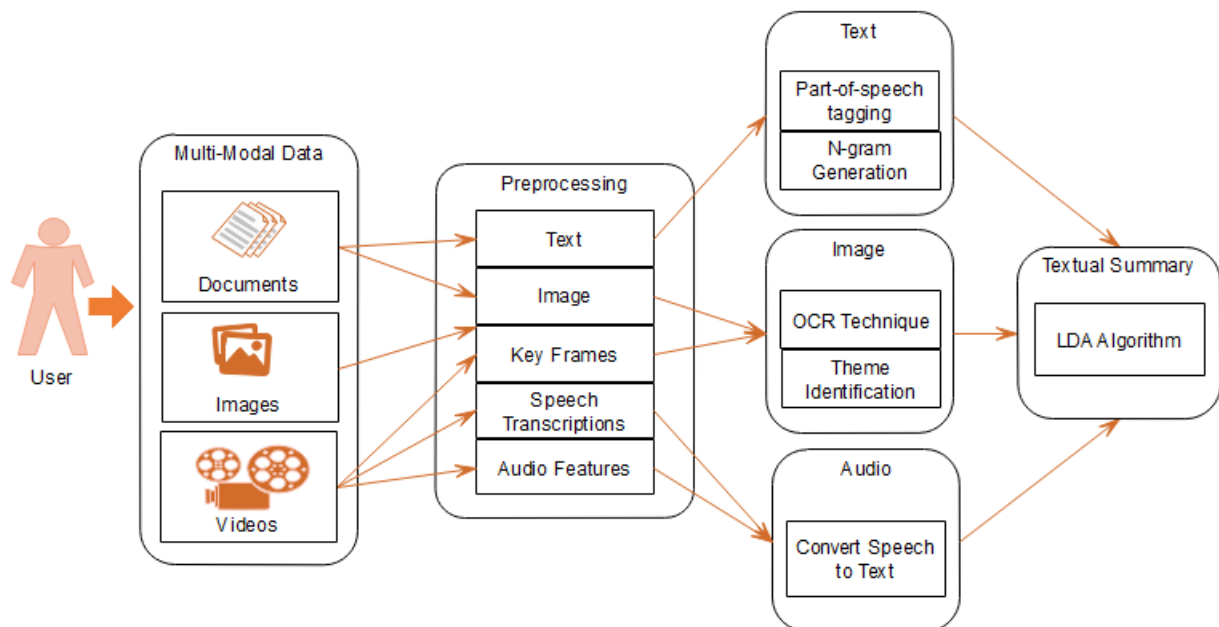


Fig. 1 System Architecture

Advantages are:

1. Provides users with textual summaries to help to acquire the gist of asynchronous multimedia data in a short time without reading documents or watching videos from beginning to end.
2. Automatically generate a fixed-length textual summary to represent the principle content of the multi-modal data.

V. MATHEMATICAL MODEL

5.1 Mathematical Formulae

The input is a collection of multi-modal data $M = \{D_1, \dots, D_{|D|}, V_1, \dots, V_{|V|}\}$ related to a news event T , where each news document $D_i = \{T_i, I_i\}$ consists of text T_i and image I_i (there may be no image for some documents). V_i denotes a news video and $|\cdot|$ denotes the cardinality of a set. The objective of our work is to automatically generate a fixed-length textual summary to represent the principle content of the multi-modal data M .

An extractive summarization method in which all these aspects can be jointly optimized through the budgeted maximization of submodular functions defined as follows:

International Journal of Innovative Research in Science, Engineering and Technology

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Visit: www.ijirset.com

Vol. 8, Issue 2, February 2019

$$\max_{S \subseteq T} \{F(S) : \sum_{s \in S} l_s \leq L\} \quad (1)$$

where T is a set of sentences, S is a summary, l_s is the length (number of words) of sentence s, L is budget, i.e., length constraint for the summary, and submodular function F(S) is the summary score.

- Text Saliency of Summarization

The text saliency of summary S by a diversity-aware objective:

$$F_s(S) = \sum_{i=1}^K \sqrt{\sum_{s_j \in P_i \cap S} Sa(s_j)} \quad (2)$$

where $P_i, i = 1, \dots, K$ is a disjoint partition of the set of all the sentences and speech transcriptions V into separate clusters, and the saliency scores $s(t_j)$ are normalized to [0,1] by dividing by the maximum value among all the sentences.

- Matching-based Image Coverage of Summarization

We model the summary S coverage for the image set I as follows:

$$F_m(S) = \sum_{p_i \in I} Im(p_i) \cdot m(p_i, C_j) \quad (3)$$

where the $Im(p_i)$ is the weight for image p_i . For keyframe p_i , $Im(p_i)$ is the average saliency score of the speech transcriptions within the shot to which p_i belongs. For document image p_i , $Im(p_i)$ is the average saliency score of the sentences in the document in which p_i is embedded. C_j is a sentence or a sentence segment obtained based on semantic framing, chunking or word tokenizing.

5.2 Text Summary Generator Algorithm

Step 1 - It takes a text as input.

Step 2 - Splits it into one or more paragraph(s).

Step 3 - Splits each paragraph into one or more sentence(s).

Step 4 - Splits each sentence into one or more words.

Step 5 - Gives each sentence weight-age (a floating point value) by comparing its words to a pre-defined dictionary called "stopWords.txt".

If some word of a sentence matches to any word with the pre-defined Dictionary, then the word is considered as Low weighted.

Step 6 - An ordered list of weighted sentences is then prepared (Relatively High weighted sentences comes first and low weighted sentences comes At last position).

Step 7 - Now, we have the ordered list of weighted sentences, it continues to Store each sentence (from ordered weighted sentences) in the output Variable (i.e. a list) until it reaches the reduction ratio (It uses a formula to determine max number of sentences to put in the output List).

Step 8 - The output list is then returned summary.

VI. CONCLUSION

This paper addresses an asynchronous MMS task, namely, how to use related text, audio and video information to generate a textual summary. We formulate the MMS task as an optimization problem with a budgeted maximization of submodular functions. We address readability by selectively using the transcription of audio through guidance

International Journal of Innovative Research in Science, Engineering and Technology

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Visit: www.ijirset.com

Vol. 8, Issue 2, February 2019

strategies. The experimental results shows that giving input any combination like, .txt, .docx, .jpg, .mp3, .mp4, .avi etc. file and output is automatic generated textual summary showing within time.

REFERENCES

- [1] H. Li, J. Zhu, C. Ma, J. Zhang and C. Zong, "Read, Watch, Listen and Summarize: Multi-modal Summarization for Asynchronous Text, Image, Audio and Video," in *IEEE Transactions on Knowledge and Data Engineering*.
- [2] P. Li, J. Ma, and S. Gao, "Learning to summarize web image and text mutually," in Proceedings of the 2nd ACM International Conference on Multimedia Retrieval. ACM, 2012, p. 28.
- [3] J. Bian, Y. Yang, and T.-S. Chua, "Multimedia summarization for trending topics in microblogs," in CIKM. ACM, 2013, pp. 1807–1812.
- [4] P. Goyal, L. Behera, and T. M. McGinnity, "A context-based word indexing model for document summarization," *IEEE Transactions on Knowledge & Data Engineering*, vol. 25, no. 8, pp. 1693–1705, 2013.
- [5] G. Evangelopoulos, A. Zlatintsi, A. Potamianos, P. Maragos, K. Rapantzikos, G. Skoumas, and Y. Avrithis, "Multimodal saliency and fusion for movie summarization based on aural, visual, and textual attention," *IEEE Transactions on Multimedia*, vol. 15, no. 7, pp. 1553–1568, 2013.
- [6] T. Hasan, H. Bořil, A. Sangwan, and J. H. Hansen, "Multi-modal highlight generation for sports videos using an information theoretic excitability measure," *EURASIP Journal on Advances in Signal Processing*, vol. 2013, no. 1, p. 173, 2013.
- [7] Z. Li, J. Liu, J. Tang, and H. Lu, "Robust structured subspace learning for data representation," *IEEE transactions on pattern analysis and machine intelligence*, vol. 37, no. 10, pp. 2085–2098, 2015.
- [8] J. Bian, Y. Yang, H. Zhang, and T.-S. Chua, "Multimedia summarization for social events in microblog stream," *IEEE Transactions on Multimedia*, vol. 17, no. 2, pp. 216–228, 2015.
- [9] W. Y. Wang, Y. Mehdad, D. R. Radev, and A. Stent, "A low-rank approximation approach to learning joint embeddings of news stories and images for timeline summarization," in NAACL-HLT, 2016, pp. 58–68.
- [10] Z. Li, J. Tang, X. Wang, J. Liu, and H. Lu, "Multimedia news summarization in search," *ACM Transactions on Intelligent Systems and Technology*, vol. 7, no. 3, p. 33, 2016.
- [11] H. Li, J. Zhu, C. Ma, J. Zhang, and C. Zong, "Multi-modal summarization for asynchronous collection of text, image, audio and video." in EMNLP, 2017, pp. 1092–1102.
- [12] H. Li, J. Zhang, Y. Zhou, and C. Zong, Guiderank: A guided ranking graph model for multilingual multi-document summarization, in International Conference on Computer Processing of Oriental Languages. Springer, 2016, pp. 608620.
- [13] K. Simonyan and A. Zisserman, Very deep convolutional networks for large-scale image recognition, arXiv preprint arXiv: 1409.1556, 2014.
- [14] M. Kshirsagar, S. Thomson, N. Schneider, J. Carbonell, A. N. Smith, and C. Dyer, Frame-semantic role labeling with heterogeneous annotations, in ACL, 2015, pp. 218224.
- [15] I. Vulic, W. De Smet, and M.-F. Moens, Identifying word translations from comparable corpora using latent topic models, in ACL, 2011, pp. 479484.
- [16] Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey et al., Google's neural machine translation system: Bridging the gap between human and machine translation, arXiv preprint arXiv:1609.08144, 2016.
- [17] G. Erkan and D. R. Radev, Lexrank: Graph-based lexical centrality as salience in text summarization, *Journal of Qiqihar Junior Teachers College*, vol. 22, p. 2004, 2011.