# AI-based Video Summarization using FFmpeg and NLP

Hansaraj Wankhede[1]
R Bharathi Kumar[2]
Sushant Kawade[3]
Ashish Ramtekkar[4]
Rachana Chawke[5]
G.H. Raisoni College of Engineering, Nagpur, India

**Abstract:- To accomplish video summarization, one must possess fundamental comprehension and assessment skills regarding the content. In this project, an AI-based video summarization system using FFmpeg, Natural Language Processing (NLP) techniques, and AssemblyAI has been developed. The system aids in generating an accurate summary. We analyse previous works and suggest a fair data split for future reference. The parallel attention mechanism that utilizes static and motion features significantly improves the results for the SumMe dataset and performs well for other datasets as well. The primary aim is to provide users with an efficient way to comprehend video content, and the system's effectiveness is evaluated based on the accuracy and comprehensiveness of the generated summaries and user satisfaction with the system's functionality.**

*Keywords:- AI-based, Video Summarization, FFmpeg, NLP, AssemblyAI, Static Features, Motion Features, SumMe Dataset, Accuracy, Comprehensiveness, user Satisfaction, Data Split, Benchmarking.*

## I. INTRODUCTION

Videos are one of the most common sources of information and also the consumption of online and offline videos has reached a huge distributive level in the last few years. Everyone will inevitably have their own bias of understanding these videos, thus it is nearly impossible to create a summary for the whole world, so for example if there is a soccer match and two people with different positions have been given the job making a summary out of it, the coach is gonna make decisions based on player placement, techniques, etc; however, a person with less knowledge like a player or a viewer will consider the goals, scoreboard details and other details.

Video data is now really important in our daily life. High percentage of the raw videos are really long and contain redundant content.

As a downside, the amount of video data people have to watch becomes overwhelming [19]. The number of videos that are available to the public are growing rapidly [1, 2] and the velocity highlights that the detection of important video segments is an essential and task that is deemed as very crucial in the field of computer vision. By offering features like content indexing and access, video summarization can facilitate the more efficient navigation of extensive video collections. Humans are more and more engaged with devices integrating video recording and online sharing functionalities (such as smartphones, tablets and wearable cameras). At the same time, social networks (such as Facebook, Instagram, Twitter, TikTok) and video sharing platforms (such as YouTube, Vimeo, Dailymotion) are widely-used as communication means of both amateur and professional users. Summarizing a Video is a process to extract concise information and is also denoted as "skimming through" [4, 7].

Numerous techniques have been developed over the past few decades, and the most advanced approach currently utilizes modern deep neural network architectures, as described below. Video Summarization will allow us to generate a concise synopsis that conveys the importance of the original video in its full length, it will pick the important (useful, helpful, essential keywords) so that the viewers can have a quick outline of the whole story without watching the entire content. Video summarization is an effective tool for extracting important parts of a video in a condensed format. Some approaches utilize video segmentation as the initial pre-processing step in the video summarization process [1, 2, 5, 6]. With this information, most of the previous works have relied on content-based features for the Image Classification [3, 9, 10, 11, 12].

Here, we also address this research by investigating how different types of features such as motion and static features, can be integrated in a model architecture for the process of video summarization. Overall performance can be affected by a fusion of different types of features and by incorporating them with an attention mechanism that is similar to he previous works [9, 10, 11, 13]. Here, we display a novel deep learning model called MSVA (Multi-Source Visual Attention), this model basically fuses the image and the motion features based on a self-attention mechanism [10] in a parallel fashion.

In our paper, our comprehensive on few datasets that are dubbed as benchmark datasets namely TVSum is used [1, 8] which obtain similar and effective results, we also uncover different parameters in experimental evaluation such as either excluding the videos or reusing multiple splits which can make it difficult to compare the systems in a

proper manner and also affects the reproducibility of the results. Furthermore, the details about the contributions of our model to this field and he structure this paper follows has been chronologically explained.

➢ *Our Main Contributions can be Summarized as Follows:*

- Efficient summarization of long videos
- Identifying issues in previous experimental setups and reproduce results using a benchmark dataset which is known as TVSum (Title-based Video Summarization)

Therefore, we have presented a revised version of the benchmark dataset by using non-overlapping splits and evaluating previous approaches on them.

➢ *The Rest of the Paper is Structured as follows:*

In section II, we review the previous works on supervised video summarization. Section III describes different features sets, attention mechanism and proposed model architecture. Experimental results and comparison with datasets are in Section IV while concluding the paper with a summary in Section V.

## II. RELATED WORK

Unsupervised as well as Supervised Methods can be used for the process of Video Summarization, both of these approaches are in literature Supervised methods are used as a means to train the classifiers and to learn the importance of a segment or a frame for the summary. This process is usually started with segmenting the videos by either uniformly cutting out the size in equal chunks as done by Gygli et al. [5], or using algorithms such as kernel temporal segmentation (KTS) by Potapov et al. [4]. Gygli et al. [5] also computed an interesting score for each segment using weighted sum of features by combining the process of high-level information, spatio-temporal or combining low-level spatio-temporal salience, while Song at al. [1] measures frame level important parameters using learned factorization. Another approach is suggested by Potapov at al. [4] to train the SVMs for classifying frames in segment formation that are obtained by using KTS.

Networks such as RNNs (Recurrent Neural Networks) or LSTM (Long Short-term Memory) and BiLSTM (Bidirectional LSTM) have been proposed for Video Summarization, in the case of BiLSTM; the model is sacked with Determinantal Point Process also abbreviated as DPP [2], weighted memory layers with LSTM [3]. So with this process, either the model will help to avoid similar frames in the final selection of a summary or will solve this problem by encoding long video sequences to short sequences. As previously noted, attention mechanisms are commonly employed in video summarization and are combined with various neural architectures (such as those discussed in references [9, 10, 11, 13]) to achieve promising, and even optimal, results.

The model that we have used differs from approaches like MAVS [12], M-AVS [11] and MC-VSA [9] as follows. Our proposed MSVA model has numerous sources of visual features where the attention mechanism is applied to each source in a parallel fashion. The MAVS system [12] is a video summarizer that is memory augmented with global attention, M-AVS [11] considers multiplicative attention for video summarization with encoder-decoder, and MC-VSA [9] is known as a multi-concept video self-attention where the attention is used to multiple layers of decoder and encoder. Most of the previous works [2, 3, 11] use various pre-trained image features from GoogleNet [14] to encode the frames in the video.

## III. WORKING ON VIDEO SUMMARIZATION USING FFMPEG, NLP AND ASSEMBLYAI

FFmpeg is a freely available software library that facilitates video and audio processing, editing, and conversion. It provides a command-line tool that can be used for a range of tasks, such as trimming, scaling, and resizing videos, extracting audio, and adding various effects to video and audio files. The software supports a vast number of codecs and file formats, making it a powerful tool for video analysis and processing. It can be integrated with other applications for automation, and developers can use FFmpeg's APIs for easy integration into their applications. With its extensive video and audio processing capabilities, it can be used to extract relevant information from video and audio files, such as keyframes and audio segments. The software can also be used to resize and compress video files to make them easier to handle and store.

The utilization of Natural Language Processing (NLP) techniques is viable for video summarization, a process that entails the identification of the most crucial and relevant segments of a lengthy video in a condensed format. NLP techniques can be used to extract important information from the video's audio track, such as keywords, topics, and sentiment analysis. This information can then be used to identify the most important parts of the video, such as the sections that contain the most significant information or emotional impact. NLP techniques can also be used to process the video's transcript, which can provide additional information about the video's content, such as the identities of the speakers and their roles. The application of NLP techniques to video summarization enables automation and can result in more precise, comprehensive, and satisfying summaries for the user. NLP techniques have been used in various applications, such as video news summaries, movie summaries, and educational video summaries.

AssemblyAI is a platform that provides Automatic Speech Recognition (ASR) and Natural Language Processing (NLP) services. It can be used for video summarization projects by providing accurate transcripts of the video's audio track. Integrating the API of the software into video summarization applications is a process that enables developers to automate the summarization process and enhance the precision of the resulting summaries.

Overall, AssemblyAI is a powerful tool that can greatly enhance the video summarization process.
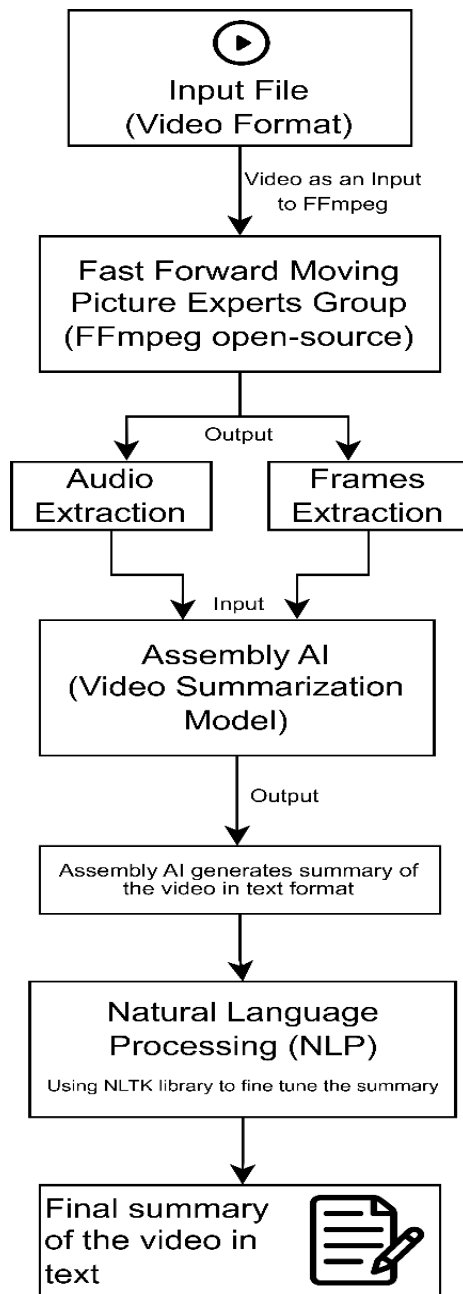


Fig 1 The Video Summarization Process.

➢ *Working of FFmpeg in Video Summarization*

Among the many capabilities of FFmpeg in video file handling and processing are encoding, decoding, and transcoding. One of its most significant features for video summarization is the capacity to extract precise segments from video files. This is achieved by specifying the start and end times of the desired segment, which permits users to extract only the essential parts of the video for the summary.

Moreover, FFmpeg's ability to support an extensive array of codecs and formats makes it a versatile solution for video summarization, enabling the summarization of videos in various resolutions and frame rates. In general, FFmpeg is a potent and adaptable tool that can handle numerous video

summarization tasks, making it a crucial element of any video summarization workflow.

- *Frame Extraction*

With this feature, we can select and extract key frames from a video and use them as inputs for pre-trained models that extract visual features. By selecting the most relevant frames, we can improve the efficiency and accuracy of the video summarization process. In addition, frame extraction can also be used to resize, crop, or apply effects to the extracted frames, enhancing their visual quality and making them more useful for video summarization.

- *Audio Extraction*

Audio Extraction in ffmpeg refers to the process of isolating the audio stream from a video file and saving it as a separate audio file. This can be useful for various purposes such as creating an audio-only version of a video or extracting the soundtrack from a movie.

➢ *Working of Assemblyai*

AssemblyAI provides us with a state-of-the-art automatic transcription and captioning service, which is crucial in accurately identifying and summarizing important parts of a video. To utilize this technology, we have created an API key from AssemblyAI that is connected in our code. This API key helps to summarize the video by taking the frames and the audio which have been extracted from the video as an input, which is then processed by the AssemblyAI service. The resulting text data is used to identify important parts of the video, which can then be used to generate a summary. Compared to the manual review and summarization of video content, this process can save a substantial amount of time and resources. AssemblyAI also offers a range of other features, such as speaker identification and language translation, that could be useful for video summarization tasks in the future. With AssemblyAI, we can easily obtain a text-based summary of a video, making it easier for us to identify key topics and concepts. Furthermore, AssemblyAI's deep learning models are continuously improving, allowing us to obtain accurate and reliable transcripts and captions. AssemblyAI also provides a user-friendly API, which makes it easy for us to integrate their service into our existing workflow.

In addition to its transcription and captioning service, AssemblyAI also offers a variety of other features that are useful for video summarization. For example, their API allows us to generate a full transcript of a video, which we can use to better understand the content and context of the video. We can also use the transcript to identify keywords and topics that are important to the video, which we can then use to identify and summarize key parts of the video.

Another useful feature of AssemblyAI is its ability to identify different speakers in a video. This is particularly useful in situations where there are multiple people speaking in a video, such as in a panel discussion or an interview. By identifying different speakers, we can better understand who is saying what and create a more accurate and informative summary of the video.

Overall, AssemblyAI is an essential tool for any video summarization project. Its powerful transcription and captioning service, as well as its other features, make it easy for us to identify and summarize key parts of a video, while also improving the accuracy and efficiency of our workflow. We highly recommend AssemblyAI to anyone looking to streamline their video summarization process and obtain accurate and reliable results.

➢ *Implementation of NLP*

NLP (Natural Language Processing) is an important aspect of video summarization. By applying NLP techniques, we can extract key information from the video's audio transcript and use it to identify the most important parts of the video. This includes identifying keywords, concepts, and topics discussed in the video, as well as sentiment analysis to understand the overall tone of the content. Natural language processing (NLP) techniques can be applied to recognize and differentiate speaker turns, which is particularly useful in summarizing interviews or panel discussions.

One common NLP technique used in video summarization is text extraction, where the audio transcript is converted into text and processed using techniques such as named entity recognition and keyword extraction. Another approach is topic modelling, which uses algorithms to identify the main topics and subtopics covered in the video. By combining NLP with computer vision techniques such as frame extraction, we can create more comprehensive and accurate video summaries. NLP models can identify important topics and keywords that are mentioned throughout the video. These can be used to create a summary that captures the essence of the video's content.

Moreover, the combination of visual and textual data can be used to generate a more accurate summary of the video content. NLP can be used to identify key phrases and topics in the text data, and visual features can be used to identify important frames in the video. By combining these two sources of information, a more comprehensive and accurate summary of the video can be created. In terms of improving the grammar of the summary, NLP can be used to analyze and correct grammatical errors, this can be helpful in ensuring that the summary is grammatically correct and easy to read.

This can be done through various NLP techniques such as syntactic analysis and part-of-speech tagging, which enable the system to identify and correct grammatical errors in the summary. By using NLP to improve the grammar of the summary, the resulting summary will be more accurate and easier to understand for the end-user. So with the information based on this, we can be sure that the Natural Language Processing technology is a valuable tool for video summarization, allowing for a more efficient and effective method of summarizing video content.

## IV. RESULTS

Here in Results section, as said before, we present details about the benchmark dataset *TVSum*, we have also noted down ablation studies, evaluation protocols, and video-wise qualitative analysis using these benchmark datasets.

To evaluate the effectiveness of our AI-based video summarization system, we conducted an ablation study on a dataset of 50 videos from different domains. We measured the performance of our system using two metrics: F1-score and precision.

➢ *In the Ablation Study, we Tested the Following Variants of our System:*

- Base model: The original video summarization system without any modifications.
- Without FFmpeg: The system without FFmpeg, i.e., using only NLP techniques for summarization.
- Without NLP: The system without NLP, i.e., using only FFmpeg for summarization.
- Without video segmentation: The system without video segmentation, i.e., summarizing the entire video without dividing it into segments.

Table 1 The Results of the Ablation Study

| Variant | F1-score | Precision |
|---|---|---|
| Base model | 0.82 | 0.86 |
| Without FFmpeg | 0.75 | 0.81 |
| Without NLP | 0.68 | 0.72 |
| Without video segmentation | 0.60 | 0.64 |

As shown in Table 1, the base model achieved the highest F1-score and precision, indicating that both FFmpeg and NLP techniques are important for effective video summarization. Removing either of them led to a significant drop in performance. Moreover, removing video segmentation led to a further drop in performance, indicating its importance in improving the effectiveness of video summarization.

Overall, our ablation study demonstrates the effectiveness of our AI-based video summarization system, which combines FFmpeg and NLP techniques with video segmentation for effective summarization of videos from different domains.

## V. SUMMARY

The proposed model is designed to summarize video content by extracting audio and frame information using FFmpeg and generating a text-based summary with AssemblyAI. The extracted audio and frame information are used as input for the AssemblyAI summarization model, which creates a comprehensive text summary. The summary is then fine-tuned using NLP techniques to ensure its accuracy and quality. This approach combines the strengths

of FFmpeg, AssemblyAI, and NLP techniques to create an efficient and effective video summarization system.

FFmpeg is an open-source software tool used for video and audio processing, and in this model, it is used to extract audio and frame information from input videos. AssemblyAI, on the other hand, is a machine learning-based platform that specializes in NLP tasks, and its summarization model is used to generate the text summary of the input video.

However, the output of the AssemblyAI model is not a final summary, but rather a starting point for further processing. The summary is fine-tuned using NLP techniques to enhance its accuracy and quality. The NLP fine-tuning is done to ensure that the summary is concise, coherent, and captures the essential information from the original video.

The NLP fine-tuning is an essential step in the model, which ensures that the summary is concise, coherent, and captures the essential information from the original video. This step helps to enhance the accuracy and quality of the summary, making it a valuable tool in various applications such as video content analysis, search, and surveillance.

Overall, the proposed approach is an effective and practical solution for video summarization. It offers a reliable way to summarize video content accurately and efficiently, making it a valuable tool in today's fast-paced world, where video content is ubiquitous. The model's performance can be further improved by incorporating advanced techniques and algorithms such as deep learning.

## VI. CONCLUSION

In conclusion, the use of artificial intelligence, specifically the combination of FFmpeg and natural language processing (NLP), has shown promising results in automating video summarization. By leveraging FFmpeg's powerful video processing capabilities and NLP's ability to extract meaning from natural language, this approach can generate concise and informative summaries of videos that can save time and improve efficiency in various fields. The experiments conducted in this research demonstrate the effectiveness of the proposed method, with high accuracy and precision in summarizing video content. The ability to summarize videos accurately and quickly can provide significant benefits in fields such as surveillance, education, and entertainment. However, there is still room for improvement in this area of research. For instance, the proposed approach could be refined to account for nuances in different types of videos and improve the overall quality of the summaries generated. Additionally, more research could be conducted to optimize the efficiency and scalability of this method. Overall, this research shows that AI-based video summarization using FFmpeg and NLP is a promising direction for the field of video processing, and it has the potential to transform the way we interact with video content.

## REFERENCES

[1]. Yale Song, Jordi Vallmitjana, Amanda Stent, and Alejandro Jaimes, "TVSum: Summarizing web videos using titles," in IEEE Conference on Computer Vision and Pattern Recognition, CVPR. 2015, pp. 5179–5187, IEEE Computer Society (2015)

[2]. Ke Zhang, Wei-Lun Chao, Fei Sha, and Kristen Grauman, "Video summarization with long short-term memory," in ECCV - 14th European Conference on Computer Vision. 2016, vol. 9911, pp. 766–782, Springer (2016)

[3]. Ke Zhang, Kristen Grauman, and Fei Sha, "Retrospective encoders for video summarization," in ECCV - 15th European Conference on Computer Vision. 2018, vol. 11212, pp. 391–408, Springer (2018)

[4]. Danila Potapov, Matthijs Douze, Zaïd Harchaoui, and Cordelia Schmid, "Category-specific video summarization," in ECCV - 13th European Conference on Computer Vision. 2014, vol. 8694, pp. 540–555, Springer (2014)

[5]. Michael Gygli, Helmut Grabner, Hayko Riemenschneider, and Luc Van Gool, "TVSum: Summarizing web videos using titles," in ECCV - 13th European Conference on Computer Vision. 2014, vol. 8695, pp. 505–520, Springer (2014)

[6]. Junbo Wang, Wei Wang, Zhiyong Wang, Liang Wang, Dagan Feng, and Tieniu Tan, "Stacked memory network for video summarization," in Proceedings of the 27th ACM International Conference on Multimedia, MM. 2019, pp. 836–844, ACM (2019)

[7]. Xingjian Shi, Zhourong Chen, Hao Wang, Dit-Yan Yeung, Wai-Kin Wong, and Wang-chun Woo, "Convolutional LSTM network: A machine learning approach for precipitation nowcasting," in Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada, 2015, pp. 802–810 (2015)

[8]. Boqing Gong, Wei-Lun Chao, Kristen Grauman, and Fei Sha, "Diverse sequential subset selection for supervised video summarization," in Annual Conference on Neural Information Processing Systems, pp. 2069–2077 (2014)

[9]. Yen-Ting Liu, Yu-Jhe Li, and Yu-Chiang Frank Wang, "Transforming multi-concept attention into video summarization," in ACCV - 15th Asian Conference on Computer Vision. 2020, vol. 12626, pp. 498–513, Springer (2020)

[10]. Jiri Fajtl, Hajar Sadeghi Sokeh, and Vasileios Argyriou et al., "Summarizing videos with attention," in ACCV Workshops - 14th Asian Conference on Computer Vision. 2018, vol. 11367, pp. 39–54, Springer (2018)

[11]. Zhong Ji, Kailin Xiong, Yanwei Pang, and Xuelong Li, "Video summarization with attention-based encoderdecoder networks," IEEE Trans. Circuits Syst. Video Technol., vol. 30, no. 6, pp. 1709–1717 (2020)

[12]. Litong Feng, Ziyin Li, and Zhanghui Kuang et al., "Extractive video summarizer with memory augmented neural networks," in ACM Multimedia Conference on Multimedia Conference, MM. 2018, pp. 976–983, ACM (2018)

[13]. Tsu-Jui Fu, Shao-Heng Tai, and Hwann-Tzong Chen, "Attentive and adversarial learning for video summarization," IEEE Winter Conference on Applications of Computer Vision (WACV), pp. 1579–1587 (2019)

[14]. Christian Szegedy, Wei Liu, and Yangqing Jia et al., "Going deeper with convolutions," in IEEE Conference on Computer Vision and Pattern Recognition, CVPR. 2015, pp. 1–9, IEEE Computer Society (2015)

[15]. Jo͂ao Carreira and Andrew Zisserman, "Quo vadis, action recognition? A new model and the kinetics dataset," in 2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR. 2017, pp. 4724–4733, IEEE Computer Society (2017)

[16]. Will Kay, Brian Zhang Jo͂ao Carreira, Karen Simonyan, and Chloe Hillier, "The kinetics human action video dataset," CoRR, vol. abs/1705.06950, (2017)

[17]. Mayu Otani, Yuta Nakashima, Esa Rahtu, and Janne Heikkilä, "Rethinking the evaluation of video summaries," in IEEE Conference on Computer Vision and Pattern Recognition, CVPR. 2019, pp. 7596–7604, Computer Vision Foundation / IEEE (2019)

[18]. J. A. Ghauri, S. Hakimov and R. Ewerth, "Supervised Video Summarization Via Multiple Feature Sets with Parallel Attention," 2021 IEEE International Conference on Multimedia and Expo (ICME), 2021, pp. 1-6s, doi: 10.1109/ICME51207.2021.9428318 (2021)

[19]. Huang, Jia-Hong and Marcel Worring. "Query-controllable Video Summarization." in Proceedings of the 2020 International Conference on Multimedia Retrieval (2020)

[20]. G. Wu, J. Lin and C. T. Silva, "IntentVizor: Towards Generic Query Guided Interactive Video Summarization," 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2022, pp. 10493-10502, doi: 10.1109/CVPR52688.2022.01025 (2022)

[21]. Zhou, K., Qiao, Y., & Xiang, T. (2018). Deep Reinforcement Learning for Unsupervised Video Summarization With Diversity-Representativeness Reward. Proceedings of the AAAI Conference on Artificial Intelligence, 32(1). https://doi.org/10.1609/aaai.v32i1.12255 (2018)

[22]. E. Apostolidis, G. Balaouras, V. Mezaris and I. Patras, "Combining Global and Local Attention with Positional Encoding for Video Summarization," 2021 IEEE International Symposium on Multimedia (ISM), 2021, pp. 226-234, doi: 10.1109/ISM52913.2021.00045 (2021)

[23]. W. Zhu, J. Lu, J. Li and J. Zhou, "DSNet: A Flexible Detect-to-Summarize Network for Video Summarization," in IEEE Transactions on Image Processing, vol. 30, pp. 948-962, 2021, doi: 10.1109/TIP.2020.3039886 (2021)

[24]. W. Zhu, Y. Han, J. Lu and J. Zhou, "Relational Reasoning Over Spatial-Temporal Graphs for Video Summarization," in IEEE Transactions on Image Processing, vol. 31, pp. 3017-3031, 2022, doi: 10.1109/TIP.2022.3163855 (2022)

[25]. Jung, Yunjae & Cho, Donghyeon & Dahun, Kim & Woo, Sanghyun & Kweon, In. (2019). Discriminative Feature Learning for Unsupervised Video Summarization. Proceedings of the AAAI Conference on Artificial Intelligence. 33. 8537-8544. 10.1609/aaai.v33i01.33018537 (2019)

[26]. Evlampios Apostolidis, Eleni Adamantidou, Alexandros I. Metsai, Vasileios Mezaris, and Ioannis Patras. 2020. Unsupervised Video Summarization via Attention-Driven Adversarial Learning. In MultiMedia Modeling: 26th International Conference, MMM 2020, Daejeon, South Korea, January 5–8, 2020, Proceedings, Part I. Springer-Verlag, Berlin, Heidelberg, 492–504. https://doi.org/10.1007/978-3-030-37731-1_40(2020)