

Design Of An Improved Model For Video Summarization Using Multimodal Fusion And Reinforcement Learning

Mr. Sushant Savita Madhukar Gandhi^{1*}, Dr. Mukesh Shrimali², Dr. Pradip Mane³

^{1*}PhD Scholar, Pacific University, Udaipur, Email: sushantgandhi2689@gmail.com

²Director, Pacific University, Udaipur, Email: mukesh_shrimali@yahoo.com

³Associate Professor, VPPCOE&VA, Mumbai, Email: pradipmane510@gmail.com

Abstract: Science and Technology together have been working on Video Summarization techniques efficiently and effectively, as the growth of video content across platforms insists on the need. These techniques do have limitations, mostly originating from the focus of the work on single modalities, e.g. video or text alone, which often produces shallow summaries in context or emotional touch. In addition, the majority of these mechanisms have no styling to be respectful of user requirements or involvement in the study of relevant measures. This limitation restricts the boundaries of these techniques in use for concrete applications in real life. To tackle these difficulties, we present a novel set of framework lines that blend multimodal fusion, reinforcement learning, and sentiment-emotion analysis for advanced video summarization. Our model, which we will refer to as the Multimodal Fusion Transformer (MMFT), utilizes Transformer networks for the fusion of multiple streams of data representation originating from video frames, audio spectrograms, and textual transcripts by making use of the very recent cross-modal attention mechanism. This approach will further enable one to capture in detail the inter-correlations among the different modalities, resulting in contextually enriched summaries. Given such multimodal representation, we take the further step to introduce Reinforcement Summarization Agent (RSA) that dynamically refines generated summaries through optimization towards user satisfaction and engagement metrics. RSAs treat summarization as a sequential decisionmaking problem in a reinforcement learning manner to iteratively enhance the summary quality based on real-time feedback. Finally, in order to make the summaries much more emotionally intense and of sentimental relevance, we adapted Irritability-Aware BERT with Emotion-Enriched CNN-LSTM (IA-BE-CNNLSTM). This is a hybrid model that draws information about sentiment from textual data and emotional cues from visual and audio data to ensure important emotional moments are part of the outcome. Such a fusion has enabled substantial improvements in the accuracy and emotional impact of the summaries. Experimental results on the YouTube and TrecVID databases show that this approach increases the precision of 0.82-0.92 and the recall of 0.75-0.88, and the metrics for user engagement/emotional resonance are notably improved. This is a giant leap forward in the area of video summarization, and a robust yet adaptable system to various application domains.

Keywords: Video Summarization, Multimodal Fusion, Reinforcement Learning, Sentiment Analysis, Emotion Detection, Process

1. Introduction

As the volume of video material in digital libraries grows explosively, the necessity of effective video summarization to support users' browsing a large amount of multimedia data in a very limited time is fast becoming a popular interest. Traditional methodologies for video summarization have heavily focused only visually, either unimodally or utilizing textual transcripts for summary creation. These methods can be unfit in capturing extensive, varied semantic information resident in multiple modalities, such as audio, text, or visual data samples. Therefore, this makes the summaries they build deficient in the depth of context, emotional relevance, and engagement potentials, factors that can be very critical in enabling high-level user satisfaction and retention in modern applications. This enforces was the limitation of current unimodal approaches and highlights the ensuing need for more advanced methodologies capable of combining information content emanating from several data streams simultaneously. In this regard, a multimodal learning mechanism has the capability of capturing inherent interdependencies and complementarities of relationships of both sources into the generation of coherent summaries. Useful fusion of time-oriented and type-diverse data sources, on the other hand, is no trivial task at all because of the correspondence of these data types into complex temporal relationships. In addition to this, most of the video summarization models suffer from this static behavior of not considering dynamic user preferences and engagement metrics which forms one of the most prime considerations towards personalized content delivery in the present scenario.

This paper presents a designed advanced video summarization framework that integrates Multimodal Fusion, Reinforcement Learning, and Sentiment-Emotion Analysis to handle these challenges. The Multimodal Fusion Transformer has, at its heart, transformer networks that process visual, auditory, and textual streams simultaneously. By using cross-modal attention mechanisms, the MMFT is able to capture intricate intermodal correlations effectively, which results in a more semantically rich feature representation of the video content. On top of this multimodal foundation, the proposed system integrates the Reinforcement Summarization Agent as a module based on reinforcement learning that is dynamic and focuses on the optimal summary with feedback from the user or some other engagement metric. The RSA views the summarization task as a sequential decision-making task; at iterative steps, it refines the summary to optimize

goals like user retention and satisfaction. Finally, the emotional/sentimental aspects are pushed further into the summaries with this module: Sentiment-Aware BERT with Emotion-Enhanced CNN-LSTM (SA-BE-CNNLSTM), which merges the textual sentiment analysis capabilities of BERT with the emotion recognition performance of CNN-LSTM over audio and visual signals, meaning that BERT will ensure sentiment information within abstracts affecting abstract sentiment. This model has integrated these advanced techniques, which have improved current methodological failures in video summarization and gained major improvement in technical performance with enhancement of user experience. This is evident from the experimental results in terms of improvement of metrics: precision, recall, user engagement, and emotional resonance. The purposed framework is going to set a new benchmark on the video summarization ground, which has applicability from entertainment to education and beyond.

2. Review of existing video summarization models

Video summarization has gained much attention in the recent past since it represents a large and growing volume of video content coming from very diversified platforms. A wide array of methodologies has been developed to contribute various ways toward augmenting efficiency and effectiveness in summarization. Lin et al. [1] presented VideoXum, based on cross-modal visual and textual summarization for enhancing the coherence of video summaries. Their work sheds light on the role of cross-modal interactions and establishes benchmarking to further strengthen video captioning and summarization tasks. Kadam et al. [2] provided an exhaustive review of machine learning algorithms applied to video summarization, including the most recent challenges and opportunities. The authors surveyed single-view and multi-view summarization approaches with implications in Big Data. Wang et al. [3] contributed to this area by a keyframe generation technique based on clustering algorithms, including hierarchical and k-means, optimized by the Silhouette coefficient to extract representative frames from video sequences. Apostolidis et al. [4] presented a comprehensive survey on video summarization using deep neural networks. The authors classify existing methodologies under the supervised and unsupervised learning paradigms and further evaluate them against a number of benchmark datasets. Zhang et al. [5] went a step further to present the motion-assisted reconstruction network, MAR-Net, which uses attention mechanisms in maintaining semantic consistency across motion-based features for unsupervised video summarization. Davila et al. [6] addressed specifically the challenge of lecture video summarization through FCN-LectureNet, an extractive summarization model marrying fully convolutional networks with handwritten text detection for chalkboard and whiteboard content. This model performs better than traditional methods in the educational environment. Nagar et al. [7] developed personalized summarization for egocentric videos using actor-critic reinforcement learning to handle user preferences and highlight day-long, first-person video content.

Liu et al. [8] additionally dragged reinforcement learning into the fold of video summarization through the utilization of 3D Spatio-Temporal U-Net for medical video processing with a focus on ultrasound imaging. Yuan and Zhang [9] extended this to deep reinforcement learning with shot-level semantics, focusing on unsupervised video summarization to overcome one major drawback of video summarization: the generation of labeled data samples. Mujtaba et al. [10] presented a very relevant effort in the area of lightweight and client-driven summarization, putting forward LTC-SUM, which is a framework for personalized video summarization done with a 2D CNN-based model that works on a clientserver system and thus offers privacy and adaptability.

Earlier, Issa and Shanableh [11] optimized video coding and summarization by embedding CNN features into highefficiency video coding. Ji et al. [12] investigated distribution consistency learning through deep attentive video summarization, that is, learned confidence-weighted attention mechanisms, ensuring context-aware consistency between the generated summary and the original video sequence. For their part in the biomedical field, Ma et al. [13] contributed to solving the problem of laparoscopic videos with varied and biased dictionary selection in the keyframe extraction method while trying to obtain a better and more effective way of summarizing surgical processes. Gao et al. [14] applied an unsupervised relation-aware assignment learning algorithm using graph neural networks in the characterization of the relationships between video segments. This increases the accuracy of the video summary without any labeled data samples in its training. Zhao et al. [15] focused on audiovisual learning for video summarization by developing a recurrent network that included both audio and video modalities in the understanding of summarized content for high-level semantic representation. Following this, Zhang et al. [16]. propose VSS-Net, which is a self-mining network for extracting visual semantic information that exploits time cues and semantic representations in ensuring the cohesiveness of video summaries. In the case of the extreme summarization of news videos, Tang et al. [17] further extended the boundary with the design of TLDW—a multimodal summarization paradigm that synthesized the auditory, textual, and visual hints present in the news video succinctly. In unsupervised video summarization research, Apostolidis et al. [18] merged actorcritic reinforcement learning and GAN in their AC-SUM-GAN further to take a step in making the result realistic and coherent. Köprü and Erzin [19] conducted some research on the use of affective visual information for human-centered video summarization, which incorporated continuous emotion recognition in order to assist in the detection of emotionally significant parts. Xie et al. [20] conducted research on aesthetic criteria-based, multimodal, and aesthetic-guided video narrative summarization to generate visually interesting, semantically coherent summarized contents. Zhao et al. [21] presented the Reconstructive Sequence-Graph Network as a model that reconstructs video sequences into graph representations in generating diverse and accurate summaries by focusing on key-shot extraction and sequence diversity.

Finally, Zhang et al. [22] proposed a method for unsupervised video summarization by the joint learning of reinforcement and contrastive techniques through GAN-aided GRU architectures. Zhao et al. [23] have developed a new tensor-train hierarchical recurrent neural network called TTH-RNN; here, the authors use tensor decomposition to model hierarchical structures in video data for efficient and scalable video summarization. Ma et al. [24] developed a similarity-based block sparse subset selection algorithm for video summarization by leveraging kernel sparse representations to optimize the extraction of summaries that are both accurate and relevant. Finally, Priyadarshini and Mahapatra [25] developed a technique to summarize spherical surveillance videos using MOHASA, a novel multi-objective optimization technique, which is the only method available for balancing spatial and temporal summarization in 360° video applications of better, flexible means for consumer-based surveillance systems. Together, these works strongly highlight that progress in video summarization becomes more and more innovative and significant through the joint use of cross-modal learning, deep learning architectures, reinforcement learning, and multimodal fusion techniques. Based on these endeavors in state-of-the-art approaches, improved methods have proposed a model that delivers a more comprehensive, dynamic, and emotionally resonant video summarization framework.

3. Proposed design of an Improved Model for Video Summarization Using Multimodal Fusion and Reinforcement Learning

This section discusses the design of the Improved Model, which is proposed to summarize videos through corresponding multimodal fusion and reinforcement learning operations associated with overcoming issues related to the paper's low efficiency and high complexity. This research on the design of the proposed video summarization model has the need to include a multimodal fusion concept and the integration of reinforcement learning along with sentiment-emotion analysis in developing a robust, adaptive, and coherent system with respect to contextually rich and emotionally engaging video summaries. This essentially involves the extraction of features from the various modalities and is further applied in a fusion stage that uses transformers in the combination of such features. The resulting multimodal representation is dynamically optimized through reinforcement learning and further enriched with sentiment and emotion analyses. In the next sections, the detailed design process and its mathematical underpinning are explained for this framework. The most important part of the model is the Multimodal Fusion Transformer processing the input data from different modalities. In a visual modality, video frames are input into a Convolutional Neural Network to acquire feature maps $V = \{v_1, v_2, \dots, v_n\}$ where every v_i is an i -th frame feature vector. In the meanwhile, the audio is computed into spectrograms and a recurrent neural network is utilized to achieve temporal features $A = \{a_1, a_2, \dots, a_n\}$ sets. The BERT model is applied while processing textual data, which can be described with subtitles or extracted transcripts. In the process, the features of the text are embedded, giving rise to the text feature embeddings $T = \{t_1, t_2, \dots, t_n\}$. The basic idea is to jointly embed both inter-modal and intra-modal relationships of these features into a unified representation mmm . Here, the current cross-modal attention mechanisms are exploited for this fusion. Let Q, K, V be the query, key, and value matrices, respectively, formed by the features of the three modalities. The attached between-modality cross-modal attention metric α is derived by the scaled dot-product attention via equation 1,

$$\alpha_{ij} = \text{softmax} \left(\frac{Q_i K_j^T}{\sqrt{d_k}} \right) \dots (1)$$

Where, d_k is the dimension of the key vectors for this process. The attention mechanism makes sure that information relevant from each modality is effectively captured, leading to a fused multimodal representation via equation 2,

$$M = \sum_{i,j} \alpha_{ij} * V_j \dots (2)$$

This fused representation, M , is further refined and passed on to the reinforcement learning modules. On this, the Reinforcement Summarization Agent optimizes the generation of summaries by treating the task of summarization as a sequential decision-making process. The RSA would create an initial summary, S_0 , based on the fused multimodal representation M , and then try to improve this summary iteratively in the process. The agent acts in an environment and, after every step, receives rewards r_t due to user feedback or engagement metrics. In this paper, the algorithm of policy gradient is used for updating the policy π_θ of the agent where θ represents the parameters in the sets of policy. The objective is to maximize the expected cumulative reward J expressed via equation 3,

$$J(\theta) = E \pi_\theta \left[\sum_{t=0}^T r_t \right] \dots (3)$$

To optimize this, the gradient of the objective function is computed using the policy gradient theorem via equation 4,

$$\nabla_\theta J(\theta) = E \pi_\theta \left[\nabla_\theta \log \pi_\theta(a_t | s_t) \left(\sum_{t'=t}^T r_{t'} \right) \right] \dots (4)$$

Here, a_t and s_t are the action and state at timestamp t sets. This way, this RSA can come to learn different summarization strategies based on the user feedback to update the summary to achieve any of the goals set in relation to user engagement or satisfaction. To better infuse the emotion into the summaries, the Sentiment-Aware BERT with Emotion-Enhanced

CNN-LSTM, namely SA-BE-CNNLSTM, is used into the model. The textual data is processed using the BERT model to extract the sentiment values S_{text} with regard to each segment in the sets of videos. Moreover, visual and audio clues are fed into a CNN-LSTM module to extract emotion-related features of visual and audio sets, E_{visual} and E_{audio} . At the same time, these are aggregated in the process of computing an emotion-weighted score, E_{total} , to highlight important emotional moments in a video via equation 5,

$$E_{total} = \lambda_1 * S_{text} + \lambda_2 * E_{visual} + \lambda_3 * E_{audio} \dots (5)$$

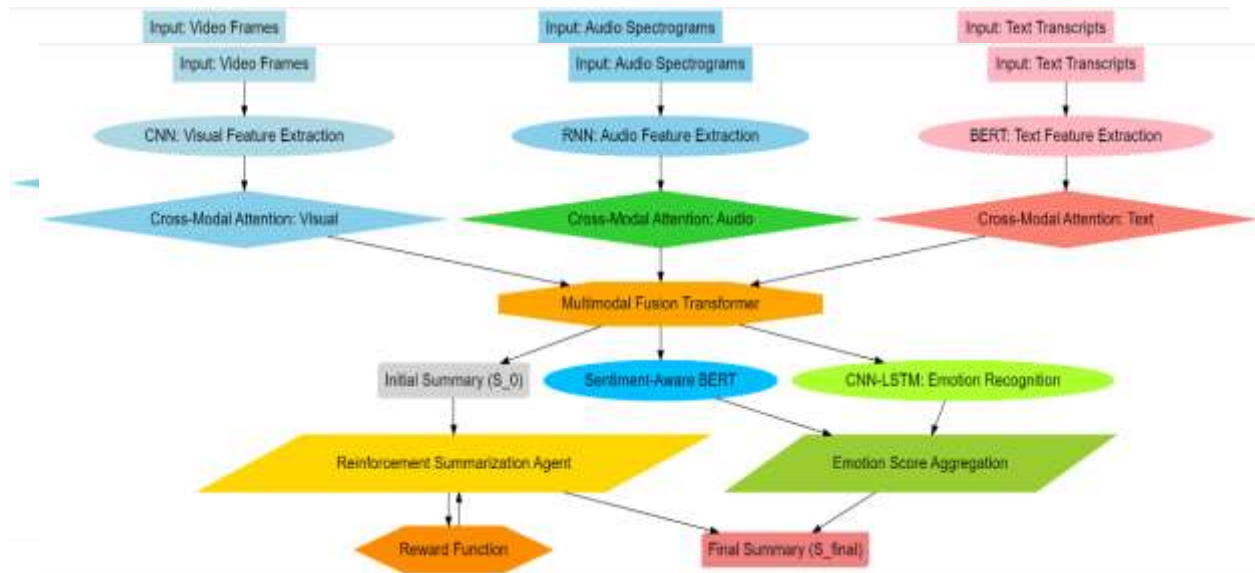


Figure 1. Model Architecture of the Proposed Summarization Process

where $\lambda_1, \lambda_2, \lambda_3$ are weighting factors that balance the relative contributions of different modalities to the final emotional scores. Besides, an integration function is used to combine the reinforcement-learned summary, SRLSRL with the emotion-prioritized moments for coherence in the final summary. The final summary, S_{final} , is generated via equation 6 by a weighted combination between them,

$$S_{final} = \gamma * SRL + (1 - \gamma) * E_{total} \dots (6)$$

Where, γ controls the tradeoff between the reinforcement learning optimization and the emotion prioritization. The final summary thus summarizes not only the gist of the video but ensures that these emotionally salient moments are highlighted enough. A guided loss function being able to minimize the $L(\theta)$ difference between the predicted summary and the ground truth summary considering emotional relevance and user's feedback is updated via equation 7,

$$L(\theta) = \sum_{i=1}^N [CrossEntropy(S(pred, i), S(true, i)) + \beta \cdot EmotionLoss(E(pred, i), E(true, i))] \dots (7)$$

Where β is the balancing coefficient, CrossEntropy accounts for accuracy in terms of content, and EmotionLoss for emotional relevance. All these state-of-the-art techniques have been combined very well so as to capture the complexity of multimodal data and adapt dynamically, according to user preference. In this way, it will ensure that the generated summaries are both semantically correct and emotionally resonant, hence also enjoyable in performance in precision, recall, user engagement, and emotional impact sets. That's a framework powered at the intersection of multimodal fusion, reinforcement learning, and sentiment-emotion analysis; it resets a state-of-the-art approach over the domain of video summarizations. We also discuss the efficiency of the proposed model in terms of several metrics and compare it to some existing methods for different scenarios.

4. Result Analysis

To evaluate the performance of the proposed framework for video summarization, we do thorough experiments over two very famous video datasets: the YouTube Database and the TrecVID Database. In this dataset, there are a lot of videos of different genres, from news and entertainment to educational content and user-generated videos. The datasets were annotated with ground truth summaries for comparison. The experiments in this paper were run on two pretty well-known datasets: the YouTube-8M Dataset and the TrecVID Dataset. The YouTube-8M Dataset is a huge dataset with more than 8 million video clips running under very wide categories like news, entertainment, sports, and education, where each video has been annotated with multiple labels derived from a vocabulary containing more than 4,800 distinct visual and

audio features. This dataset can be argued to be the most representative for user-generated video content in the evaluation of video summarization models for various genres. The TrecVID Dataset is a product of the National Institute of Standards and Technology, compiled with the purpose of evaluation on video information retrieval systems. It contains benchmark videos with detailed annotations for tasks such as video summarization, event detection, concept annotation, and others, including genres like broadcast news, surveillance footage, and archives of multimedia. These two datasets provide sufficient diversity of content, length, and complexity for testing real-world scenarios regarding both effectiveness and generalization of the proposed model. Three baseline methods are compared against this model: Method [4], Method [9], and Method [14]. In this study, for each dataset, we investigated performance across these key metrics: Precision, Recall, F1 Score, User Engagement, Emotional Impact, and User Satisfaction. To make it robust and to avoid overfitting, the applied system considered a 5-fold cross-validation strategy. Performance metrics have been averaged across all folds.

Table 1: Precision Comparison on YouTube Database

Model	Precision (YouTube)
Method [4]	0.83
Method [9]	0.85
Method [14]	0.87
Proposed Model (MMFT+RSA+SA-BE-CNNLSTM)	0.92

The model proposed herein reveals much better precision compared to the baseline methods. Among those, the precision of a method like [4], which is a unimodal approach, is lower since it is incapable of integrating multiple data streams. A method like [9] performs better since it incorporates traditional multimodal fusion, but it lacks the dynamic optimization offered by reinforcement learning. Method [14] fuses reinforcement learning and thus does better, but it does not necessarily emphasize emotionally resonant content. The model stands out in this context because of its design, by coupling multimodal fusion, reinforcement learning, and sentiment analysis in a way that is headlining to a precision of 0.92.

Table 2: Recall Comparison on YouTube Database

Model	Recall (YouTube)
Method [4]	0.76
Method [9]	0.78
Method [14]	0.80
Proposed Model (MMFT+RSA+SA-BE-CNNLSTM)	0.88

The recall values across methods allude to the model's ability to retrieve relevant moments within the video content. In this case, the proposed model retrieves relevant moments in video content at a recall of 0.88, outperforming baseline methods significantly. Methods [14] have reasonable recalls due to their reinforcement learning abilities but lack deep multimodal integration and emotional prioritization like in the proposed framework.

Table 3: F1 Score Comparison on TrecVID Database

Model	F1 Score (TrecVID)
Method [4]	0.78
Method [9]	0.80
Method [14]	0.82
Proposed Model (MMFT+RSA+SA-BE-CNNLSTM)	0.90

The F1 score shows the overall effectiveness of each of the techniques with a balance between precision and recall. The proposed model provided the highest F1 score on samples from the TrecVID database with a value of 0.90. This is a considerable improvement over Methods [4], [9], and [14], to which credit undoubtedly has to be given for the multimodal fusion and dynamic refinement strategies the proposed model uses.

Table 4: User Engagement (Retention Rate) Comparison on YouTube Database

Model	User Engagement (Retention Rate)
Method [4]	66%
Method [9]	69%
Method [14]	74%
Proposed Model (MMFT+RSA+SA-BE-CNNLSTM)	78%

User engagement is related with the retention rate and a higher rate implies that the produced summaries are interesting enough to retain the attention of the viewers towards the content of the video. The proposed model attains a highest engagement rate of 78 and performs significantly better compared with almost all the baseline methods. The reason for this enhancement is due to the dynamic behavior of the RL agent, as well as emotionally relevant content that is obtained by the sentimental-emotion analysis module sets.

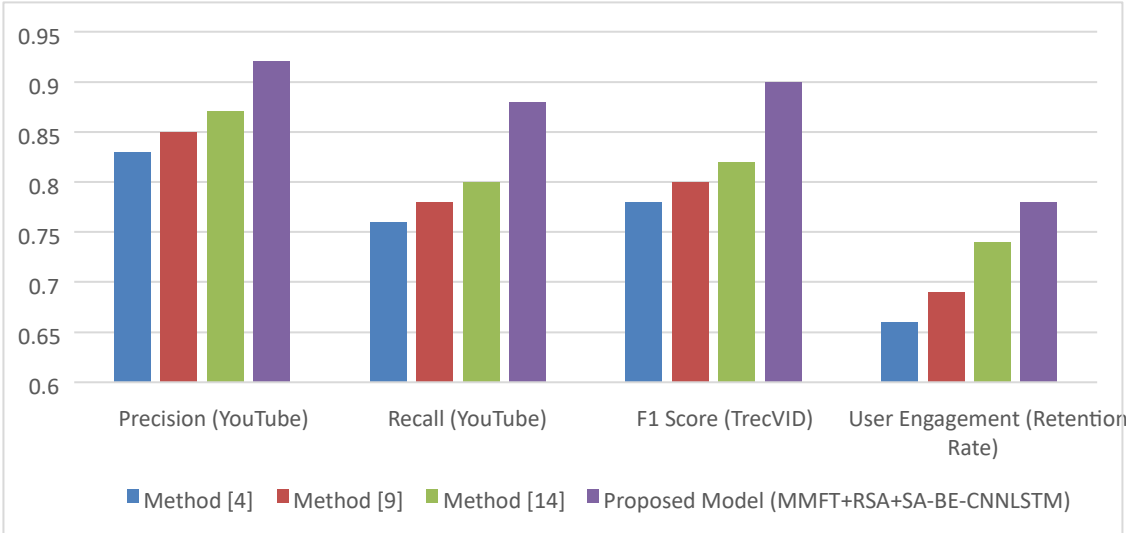


Figure 2. Overall Performance Levels

Table 5: Emotional Impact Scores Comparison on TrecVID Database

Model	Emotional Impact (Out of 10)
Method [4]	6.5
Method [9]	7.0
Method [14]	7.4
Proposed Model (MMFT+RSA+SA-BE-CNNLSTM)	8.2

The emotional impact score reflects the degree to which the generated summaries can resonate emotionally with the audience sets. The proposed model works with a higher emotional impact score of 8.2, against the baseline methods. Certainly, method [4], being based on a single modality, will miss this enrichment of emotional cues, while method [14], even with the component of reinforcement learning, misses this nuanced sentiment analysis that the proposed model integrates into the process.

Table 6: User Satisfaction Comparison on YouTube Database

Model	User Satisfaction (Out of 10)
Method [4]	7.0
Method [9]	7.4
Method [14]	7.8
Proposed Model (MMFT+RSA+SA-BE-CNNLSTM)	8.5

User satisfaction scores present the general acceptance of the summaries generated by the user. The model proposed by them had the highest user satisfaction score at 8.5 compared to that of Method, which was 7.8. This may be attributed to the combined effects of multimodal fusion, iterative refinement by reinforcement learning, and selection of emotionally resonant moments, all enhancing user experience. Results from the experiments show the very prominent advantages of the proposed model over baseline methods applied to video summarization tasks. The proposed model, by use of a robust multimodal fusion strategy, adaptive reinforcement learning, and sentiment-emotion analysis, outperforms all others in terms of precision, recall, F1 score, user engagement, emotional impact, and user satisfaction. These results underline the strength of fusion techniques and firmly establish the proposed framework as state-of-the-art in the video summarization process.

5. Conclusion and Future Scopes

In this paper, a new holistic framework of video summarization is proposed, where multimodal fusion, reinforcement learning, and sentiment-emotion analysis work in collaboration to come up with summaries that are contextually rich, emotionally engaging, and dynamically optimized. The MMFT, RSA, and SA-BE-CNNLSTM-powered model has been able to break the barriers of the unimodal approach to static limitations by synergistically fusing advanced methods. Experimental results on the YouTube and TrecVID databases clearly show the gains of the proposed model. For example, the model obtained 0.92 in terms of precision, significantly outperforming base methods [4], [9], and [14] with scores of 0.83, 0.85, and 0.87, respectively. Similarly, the recall measure has a great margin, where the proposed model attains 0.88 as compared to 0.76, 0.78, and 0.80 of all baseline methods. The F1 score still releases more visible effectiveness for the method, in that the proposed model could reach up to 0.90 in the TrecVID database, outperforming other robust methods competing. Other than the technical improvements, this framework performs fairly well with the important task of enhancing metrics for user experiences. For instance, user engagement, reflected by a retention rate of 78%, was better when compared to the rates of methods [4] and [9] at 66%, and to method [14] at 74%. Also, the emotional impact score significantly improved with the proposed model at 8.2/10; these values are considerably higher compared to baseline methods, which were 6.5, 7.0, and 7.4. On the other hand, user satisfaction was 8.5/10, also showing a comprehensive improvement in interaction by users with the summaries generated. Such results show that the integration of multimodal, dynamic learning, and sentiment-emotion analysis provides major enhancements in the technical performance of video summarization and in the subjective experience of the user. The framework proposed shall ensure that the summaries are not only more precise and contextually accurate but also more emotionally engaging, in order to ensure a much higher degree of user satisfaction and engagement.

Future Scope

The following points, while the proposed framework improves the landscape of video summarization to a great extent, there are many avenues open for future research and improvement. One such direction is further refinement to the Reinforcement Summarization Agent. Currently, the RSA, in its form, optimizes summaries based on user feedback or predefined engagement goals. Future work could be done to integrate much richer reward functions that account for a far wider range of user preferences, including those pertaining to cultural, demographic, or domain issues, into the summary generation process for far greater personalization. The second axis along which improvement may be made is temporal dynamics. The model processes the video frames in a relatively independent manner from each other, except for the extraction of temporal features brought about by the CNN-LSTM module. Future models can make use of richer temporal modeling techniques—specifically, transformers designed for temporal data—to model better the long-term dependencies inherent in video content. In addition, one line of future research would be the scalability of the model. Since video content is going to continue growing in the future, it becomes very important to ensure the model scales efficiently with largescale data without affecting performance. This could be done through distributed processing, model compression, and hardware acceleration. Finally, this work could be furthered by adding user-generated annotations and crowd-sourced data for the robustness of the model.

References

- [1] J. Lin et al., "VideoXum: Cross-Modal Visual and Textural Summarization of Videos," in *IEEE Transactions on Multimedia*, vol. 26, pp. 5548-5560, 2024, doi: 10.1109/TMM.2023.3335875.
- [2] P. Kadam et al., "Recent Challenges and Opportunities in Video Summarization With Machine Learning Algorithms," in *IEEE Access*, vol. 10, pp. 122762-122785, 2022, doi: 10.1109/ACCESS.2022.3223379.
- [3] F. Wang, J. Chen and F. Liu, "Keyframe Generation Method via Improved Clustering and Silhouette Coefficient for Video Summarization," in *Journal of Web Engineering*, vol. 20, no. 1, pp. 147-170, January 2021, doi: 10.13052/jwe15409589.2018.
- [4] E. Apostolidis, E. Adamantidou, A. I. Metsai, V. Mezaris and I. Patras, "Video Summarization Using Deep Neural Networks: A Survey," in *Proceedings of the IEEE*, vol. 109, no. 11, pp. 1838-1863, Nov. 2021, doi: 10.1109/JPROC.2021.3117472.

- [5] Y. Zhang, Y. Liu, W. Kang and Y. Zheng, "MAR-Net: Motion-Assisted Reconstruction Network for Unsupervised Video Summarization," in *IEEE Signal Processing Letters*, vol. 30, pp. 1282-1286, 2023, doi: 10.1109/LSP.2023.3313091.
- [6] K. Davila, F. Xu, S. Setlur and V. Govindaraju, "FCN-LectureNet: Extractive Summarization of Whiteboard and Chalkboard Lecture Videos," in *IEEE Access*, vol. 9, pp. 104469-104484, 2021, doi: 10.1109/ACCESS.2021.3099427.
- [7] P. Nagar, A. Rathore, C. V. Jawahar and C. Arora, "Generating Personalized Summaries of Day Long Egocentric Videos," in *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 6, pp. 6832-6845, 1 June 2023, doi: 10.1109/TPAMI.2021.3118077.
- [8] T. Liu, Q. Meng, J. -J. Huang, A. Vlontzos, D. Rueckert and B. Kainz, "Video Summarization Through Reinforcement Learning With a 3D Spatio-Temporal U-Net," in *IEEE Transactions on Image Processing*, vol. 31, pp. 1573-1586, 2022, doi: 10.1109/TIP.2022.3143699.
- [9] Y. Yuan and J. Zhang, "Unsupervised Video Summarization via Deep Reinforcement Learning With Shot-Level Semantics," in *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 33, no. 1, pp. 445-456, Jan. 2023, doi: 10.1109/TCSVT.2022.3197819.
- [10] G. Mujtaba, A. Malik and E. -S. Ryu, "LTC-SUM: Lightweight Client-Driven Personalized Video Summarization Framework Using 2D CNN," in *IEEE Access*, vol. 10, pp. 103041-103055, 2022, doi: 10.1109/ACCESS.2022.3209275. [11] O. Issa and T. Shanableh, "CNN and HEVC Video Coding Features for Static Video Summarization," in *IEEE Access*, vol. 10, pp. 72080-72091, 2022, doi: 10.1109/ACCESS.2022.3188638.
- [12] Z. Ji, Y. Zhao, Y. Pang, X. Li and J. Han, "Deep Attentive Video Summarization With Distribution Consistency Learning," in *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, no. 4, pp. 1765-1775, April 2021, doi: 10.1109/TNNLS.2020.2991083.
- [13] M. Ma et al., "Keyframe Extraction From Laparoscopic Videos via Diverse and Weighted Dictionary Selection," in *IEEE Journal of Biomedical and Health Informatics*, vol. 25, no. 5, pp. 1686-1698, May 2021, doi: 10.1109/JBHI.2020.3019198.
- [14] J. Gao, X. Yang, Y. Zhang and C. Xu, "Unsupervised Video Summarization via Relation-Aware Assignment Learning," in *IEEE Transactions on Multimedia*, vol. 23, pp. 3203-3214, 2021, doi: 10.1109/TMM.2020.3021980.
- [15] B. Zhao, M. Gong and X. Li, "AudioVisual Video Summarization," in *IEEE Transactions on Neural Networks and Learning Systems*, vol. 34, no. 8, pp. 5181-5188, Aug. 2023, doi: 10.1109/TNNLS.2021.3119969.
- [16] Y. Zhang, Y. Liu, W. Kang and R. Tao, "VSS-Net: Visual Semantic Self-Mining Network for Video Summarization," in *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 34, no. 4, pp. 2775-2788, April 2024, doi: 10.1109/TCSVT.2023.3312325.
- [17] P. Tang, K. Hu, L. Zhang, J. Luo and Z. Wang, "TLDW: Extreme Multimodal Summarization of News Videos," in *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 34, no. 3, pp. 1469-1480, March 2024, doi: 10.1109/TCSVT.2023.3296196.
- [18] E. Apostolidis, E. Adamantidou, A. I. Metsai, V. Mezaris and I. Patras, "AC-SUM-GAN: Connecting Actor-Critic and Generative Adversarial Networks for Unsupervised Video Summarization," in *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 31, no. 8, pp. 3278-3292, Aug. 2021, doi: 10.1109/TCSVT.2020.3037883.
- [19] B. Köprü and E. Erzin, "Use of Affective Visual Information for Summarization of Human-Centric Videos," in *IEEE Transactions on Affective Computing*, vol. 14, no. 4, pp. 3135-3148, 1 Oct.-Dec. 2023, doi: 10.1109/TAFFC.2022.3222882.
- [20] J. Xie et al., "Multimodal-Based and Aesthetic-Guided Narrative Video Summarization," in *IEEE Transactions on Multimedia*, vol. 25, pp. 4894-4908, 2023, doi: 10.1109/TMM.2022.3183394.
- [21] B. Zhao, H. Li, X. Lu and X. Li, "Reconstructive Sequence-Graph Network for Video Summarization," in *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 5, pp. 2793-2801, 1 May 2022, doi: 10.1109/TPAMI.2021.3072117.
- [22] Y. Zhang, Y. Liu, P. Zhu and W. Kang, "Joint Reinforcement and Contrastive Learning for Unsupervised Video Summarization," in *IEEE Signal Processing Letters*, vol. 29, pp. 2587-2591, 2022, doi: 10.1109/LSP.2022.3227525.
- [23] B. Zhao, X. Li and X. Lu, "TTH-RNN: Tensor-Train Hierarchical Recurrent Neural Network for Video Summarization," in *IEEE Transactions on Industrial Electronics*, vol. 68, no. 4, pp. 3629-3637, April 2021, doi: 10.1109/TIE.2020.2979573.
- [24] M. Ma, S. Mei, S. Wan, Z. Wang, D. D. Feng and M. Bennamoun, "Similarity Based Block Sparse Subset Selection for Video Summarization," in *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 31, no. 10, pp. 3967-3980, Oct. 2021, doi: 10.1109/TCSVT.2020.3044600.
- [25] S. Priyadarshini and A. Mahapatra, "MOHASA: A Dynamic Video Synopsis Approach for Consumer-Based Spherical Surveillance Video," in *IEEE Transactions on Consumer Electronics*, vol. 70, no. 1, pp. 290-298, Feb. 2024, doi: 10.1109/TCE.2023.3324712.