

Unsupervised Video Summarization via Multi-source Features

Hussain Kanafani², Junaid Ahmed Ghauri¹, Sherzod Hakimov¹, Ralph Ewerth^{1,2}

¹TIB – Leibniz Information Centre for Science and Technology

²L3S Research Center, Leibniz University Hannover

Hannover, Germany

hussainkanafani@gmail.com, {junaid.ghauri, sherzod.hakimov, ralph.ewerth}@tib.eu

ABSTRACT

Video summarization aims at generating a compact yet representative visual summary that conveys the essence of the original video. The advantage of unsupervised approaches is that they do not require human annotations to learn the summarization capability and generalize to a wider range of domains. Previous work relies on the same type of deep features, typically based on a model pre-trained on ImageNet data. Therefore, we propose the incorporation of multiple feature sources with chunk and stride fusion to provide more information about the visual content. For a comprehensive evaluation on the two benchmarks *TVSum* and *SumMe*, we compare our method with four state-of-the-art approaches. Two of these approaches were implemented by ourselves to reproduce the reported results. Our evaluation shows that we obtain state-of-the-art results on both datasets, while also highlighting the shortcomings of previous work with regard to the evaluation methodology. Finally, we perform error analysis on videos for the two benchmark datasets to summarize and spot the factors that lead to misclassifications.

KEYWORDS

unsupervised video summarization, multi-source combination, multi-source fusion, deep learning, video analysis

1 INTRODUCTION

Driven by the rapid growth of visual content in recent years, videos have become the dominant form of information exchange on the Web. According to Cisco Visual Networking Index [1], the video content grows annually with a rate of 33%, and it will be responsible for 80% of the global Internet traffic by 2022. However, it is time-consuming to browse long videos, and it is beneficial and preferable to watch a short and concise summary that conveys the main content of the original video. Therefore, automatic video summarization methods are required to view, search, and retrieve video content efficiently. The development of such models in a supervised fashion requires ground-truth summaries for training. However, the acquisition of a large number of ground-truth summaries is difficult, time-consuming, and expensive. Furthermore, the training data might introduce a domain or data bias. For these reasons, researchers focused on unsupervised methods that do not require human supervision and yet being able to generalize on a wider range of domains. Many approaches have tackled the task of unsupervised video summarization using different methods. Earlier approaches developed static video summaries by applying clustering algorithms to long videos [4, 6]. Deep learning approaches used different Generative Adversarial Networks (GAN) variations along with attention mechanisms [2, 3, 8–10, 15]. Other methods are

trained in a reinforcement learning-based framework with reward functions [14, 16].

In this paper, we propose a deep learning model for unsupervised video summarization called Multi-Source Chunk and Stride Fusion (MCSF), which investigates the impact of multiple visual representations extracted about visual objects and scene (i.e., places) content. It also uses two temporal constellations of the video features which give the model different perspectives of the video, similar to Jung et al [8]. Consequently, three fusion strategies are suggested and evaluated. Comprehensive experiments are conducted to compare our approach with three state-of-the-art unsupervised [2, 3, 8] methods as well as a reinforcement learning method [14] in a fair manner. We have discovered issues in the evaluation methodology used by these methods with regard to k-fold cross-validation. Some videos were excluded from the test data splits, whereas other videos were repeated multiple times in the test splits to perform k-fold cross-validation. We provide a new evaluation scheme that solves these problems in the datasets and allows for a fair comparison on the two benchmark datasets: *TVSum* [12] and *SumMe* [6]. The proposed solution yields better performance than previous state-of-the-art methods on both datasets. We share the source code for the evaluation scheme, the proposed model architecture, the re-implemented methods, and the newly generated splits of both datasets with the research community¹.

The remainder of the paper is structured as follows. The previous work in the domain of unsupervised video summarization is discussed in Section 2. In Section 3, we describe the unsupervised deep learning architecture that harnesses multi-source feature embeddings along with fusion techniques. Experimental results and the comparison with four state-of-the-art methods are presented in Section 4, while Section 5 concludes the paper.

2 RELATED WORK

In recent years, numerous works have approached video summarization from both supervised and unsupervised learning perspectives. Unsupervised methods learn the video summarization capability without the need for ground-truth summaries of videos and aim to generalize on different domains. The *Video SUMMARization* (VSUMM) approach [4] is one of the earliest methods for static video summarization, it extracts color features from the video and then performs k-medoid clustering to acquire keyframes. Gygli et al. [6] proposed a segmentation-based method that attempts to choose a subset of segments that maximizes the sum of the interestingness of segments. The interestingness score for each segment is computed using a combination of low-level features and high-level information.

¹<https://github.com/TIBHannover/UnsupervisedVideoSummarization>

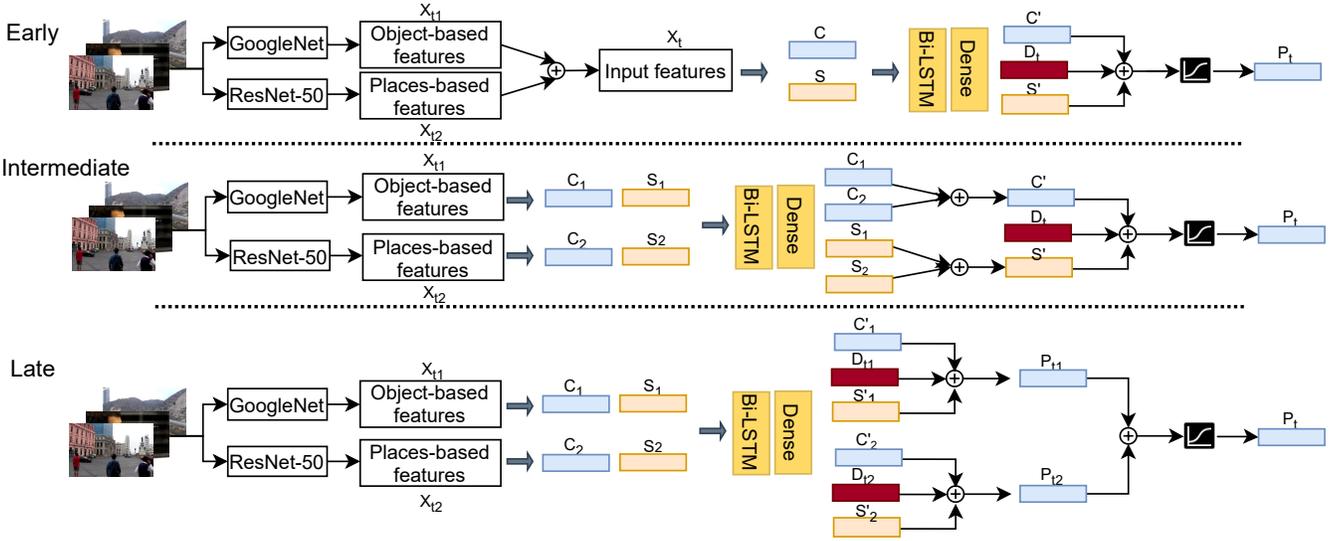


Figure 1: An overview of the Multi-Source Chunk and Stride Fusion (MCSF) architecture with different fusion techniques. X_{t1} and X_{t2} are the extracted features by GoogleNet and ResNet-50 respectively, C_1 and C_2 are the chunks, S_1 and S_2 are strides computed similar to [8] and C'_1 , S'_1 , C'_2 , and S'_2 are the chunks and strides learned by the summarizer. D_t is the difference attention, which is computed for each type of features separately and then summed.

Most deep learning methods are trained in a generative adversarial manner to overcome the absence of ground-truth video summaries. One group of methods uses generative adversarial networks (GANs) along with variational autoencoder and Long Short-term Memory (VAE-LSTM) [3, 8, 9, 15]. Among these methods *SUM-GAN* [9] represents the base that other state-of-the-art adversarial models build upon. *Chunk and Stride Network (CSNet)* [8] uses *SUM-GAN* as a baseline and tackles the problem of gradient decay when dealing with long videos as well as the ineffective feature learning due to flat distribution of frames importance scores. In CSNet, the input features are divided differently into two smaller sequences forming local chunks and global stride (view) of input features. Besides, an attention module is utilized to compute the differences of features of adjacent frames and use them as an indicator of the significance of a specific frame. Another group of methods replace VAE with a deterministic autoencoder (AE) and enhance it with an attention mechanism [2, 10]. Another approach formulated it as a decision process where the video summary algorithm is trained using reinforcement learning to find a visual summary that satisfies certain conditions such as uniformity and representativeness. However, only few publications adopt this approach [14, 18]. However, the generalizability of the aforementioned methods is not tested, since the reported experiments of most of them were conducted using the canonical settings, i.e., the data come from the same dataset and are split into 80% for training and 20% for testing.

3 UNSUPERVISED VIDEO SUMMARIZATION WITH MULTI-SOURCE FEATURES

In this section, we describe the overall architecture of the Multi-Source Chunk and Stride Fusion (MCSF) model. The overall model

architecture is shown in Figure 1. The frames in videos are sub-sampled at the rate of two frames per second where each frame’s features are extracted using the respective visual encoder model. Next, the input features are fed through different layers and finally fused in an early, intermediate, or late stage of the pipeline. As depicted in the proposed architecture, chunks and strides are fed to a Bidirectional Long-short Term Memory (Bi-LSTM) and linear layers and the output is summed up with the difference attention. The end results P_t are probabilities to select frames in the summary. The following fusion techniques [5] have been used to combine multi-source features.

Early fusion: The fusion is applied at the feature-level where features coming from different sources (X_{t1} , X_{t2}) are summed first and then fed to the next layers, which continue to compute P_t based on the fused features.

Intermediate fusion: The fusion is applied after the different chunks and strides are fed through Bi-LSTM and linear layer. The output strides are chunks from the different streams are summed accordingly with the respective difference attention, and passed through the pipeline to compute P_t .

Late fusion: The fusion is applied after both resulting importance scores from each computed. The output importance scores P_{t1} and P_{t2} are summed and passed through a sigmoid layer to produce P_t .

The intuition behind the proposed architecture is that compositional approaches can enhance the performance and enable the model to have a better understanding of the video summarization task. Current state-of-the-art methods use feature representations from a single source, mainly object-based features from pre-trained GoogleNet [13] on ImageNet data [11]. Contrary to prior work, we exploit features from two sources to enhance the representation of visual information in frames. Our approach extends previous

work [8] in order to allow for multi-source input features. These features can be fused in different network layers (early, intermediate, late) to obtain frame-level importance scores for summarization. In addition to the object-based features used in previous work, we incorporate features about (visual) scenery and place content in frames using pre-trained ResNet-50 [7] on Places365 [17] dataset. Such a fusion of multiple features allows the model to recognize changes in scene and objects, which are important indicators for summarization. For instance, different *SumMe* videos about a car crossing a railroad (*video*₆) and in a desert *video*₂₃, need to have different representations as their scene information is different, and not just rely on the visual representation about object categories present according to an ImageNet model. Thus, our proposed method uses object-based features in combination with scenery and place-related features to capture a richer representation for summarization of videos in an unsupervised fashion.

4 EXPERIMENTAL SETUP AND RESULTS

In this section, we present details about the used benchmark datasets, evaluation metrics, comparison with state of the art, discussion of the results, and qualitative video-wise error analysis.

4.1 Datasets and Evaluation Metrics

The following two benchmark datasets were used for all experiments:

- **SumMe** [6] consists of 25 videos, ranging from one to six minutes, annotated by 15 to 18 users².
- **TVSum** [12] consists of 50 videos, ranging from two to 10 minutes, annotated by 20 users³.

Methodological issues in evaluations of previous work: The evaluation of previous approaches on the two datasets is based on k -fold cross-validation where $k = 5$. When analyzing the structure of the splits used for evaluation, we found that a non-trivial number of the videos are not part of any test set of the five data splits. More specifically, some videos appear in the test data multiple times, and models are evaluated on them, whereas the model’s performance is never tested against some other videos. In total, 28% and 32% of the videos in *SumMe* and *TVSum* datasets respectively are not in the test data of the evaluated data splits provided by [3]. Therefore, this work uses new randomly generated *non-overlapping* splits to ensure that each video is contained only once in each split and all videos are included for the evaluation.

Evaluation metrics: This work assesses the performance of compared methods using F_1 scores on different splits of datasets.

4.2 Results

Ablation Study: We evaluated the variations of *MCSF* method and measured the effect of fusing multi-source features. Similar to previous work, we performed the evaluation on *SumMe* dataset by comparing the model predictions with the ground-truth summaries provided by each user and selecting the maximum value (Max). For *TVSum*, the average value (Avg) between ground-truth summaries and model predictions is used. As depicted in Table 1,

Table 1: Ablation study of MCSF model using different types of features, namely: Object-based (O) and Places and scenery-based features (P). The results are based on two different versions of splits: F'_1 on the re-evaluated provided by [3], (F_1^*) on *non-overlapping* splits proposed in this paper.

Dataset	Fusion	Features	F'_1	F_1^*
SumMe	-	O	48.1	41.5
	-	P	46.5	40.5
	Early	$O + P$	46.9	39.6
	Intermediate	$O + P$	46.0	43.3
	Late	$O + P$	47.9	40.3
TVSum	-	O	56.4	53.4
	-	P	54.2	53.6
	Early	$O + P$	54.9	54.3
	Intermediate	$O + P$	55.7	53.8
	Late	$O + P$	59.1	56.5

Table 2: The performance of different unsupervised methods with the re-evaluated (F'_1) and non-overlapping (F_1^*) splits compared to the reported results (F_1).

Dataset	Method	F_1		F'_1		F_1^*	
		Avg	Max	Avg	Max	Avg	Max
SumMe	SUM-Ind _{LU} [14]	-	51.9	22.1	46.0	18.1	42.3
	CSNet [8]	-	51.3	22.7	48.1	18.0	41.5
	SUM-GAN-AAE [2]	-	48.9	22.8	47.1	19.0	39.8
	SUM-GAN-sl [3]	-	47.3	20.4	44.6	17.7	37.7
	MCSF (ours)	-	-	21	46.0	20.1	43.3
TVSum	SUM-Ind _{LU} [14]	61.5	-	58.7	80.7	55.9	77.5
	CSNet [8]	58.8	-	56.4	77.7	53.4	76.2
	SUM-GAN-AAE [2]	58.3	-	57.7	81.6	55.1	77.8
	SUM-GAN-sl [3]	58.0	-	57.4	81.1	54.5	77.4
	MCSF (ours)	-	-	59.1	81.2	56.5	77.6

the incorporation of different types of features enhances the overall F_1 score. The highest performance on *SumMe* is achieved through intermediate fusion, whereas the late fusion performs the best on *TVSum*. Basically, we experimented with different fusion operations such as summation and averaging. However, reported results in Table 1 refer to the summation of different types of features as this operation achieved the highest performance. We also tested the model separately using only the object (O) or places and scenery (P) features without any fusion and fusing both features yielded better performance.

Overall Comparison: The evaluation results for the different state-of-the-art methods are listed in Table 2 and compared with our method. The implementations for SUM-GAN-sl [3] and SUM-GAN-AAE [2] methods are provided. We implemented both CSNet [8] and SUM-Ind_{LU} [14] since the authors did not provide their source code. The compared methods are evaluated on two versions of the five-fold data splits of both benchmark datasets using average F_1 scores using cross-validation. We computed the F_1 scores using both Average and Maximum approaches for both datasets. We also included the reported results (F_1), re-implemented and evaluation results on splits provided by [3] (F'_1), and results obtained using

²<https://gyglim.github.io/me/vsum/index.html>

³<https://github.com/yalesong/tvsum>

the *non-overlapping* splits proposed in this paper (F_1^*). We have included the best combination of our model that uses both visual features with intermediate and late fusion techniques for *SumMe* and *TVSum* datasets, respectively. It can be seen that the *MCSF* achieved higher performance on the *non-overlapping* splits of both datasets, compared to other methods. Moreover, we analyzed the results obtained and observed a drop in performance for all methods when the non-overlapping splits are used.

4.3 Error Analysis

In the following, an error analysis is performed for the following video IDs in *SumMe* dataset: 2, 5, 6, 7, 10, 12, 14, 17, 21, 24. Our goal is to find out what causes the predicted summaries to be not meaningful and fail to convey the essence of the videos. We grouped the issues under five categories as follows. Sample videos with their predictions compared with ground-truth summaries are shown in Figure 2.

Abrupt visual changes: Video segments with considerable visual changes, including shaking camera, are taken more into account in the end summary, even though these parts may not represent a relevant action related to the story of the video.

Activities with inconsiderable visual changes: In contrast to the previous point, video segments with stagnation or trivial changes in visual features are being discarded from the generated summary, although these segments may contain an essential part of the video’s story or content.

Long-temporal activities: Activities that are distributed over multiple shots are hard to capture by all state-of-the-art approaches since the methods can not infer which part of the activity is more representative and what its boundaries are.

Moving objects in the background: Current models are confused by scenes where the camera is filming a moving object.

Unrecognized certain actions: Current state-of-the-art models fail to detect essential actions such as jumping, car crash, and landing. These types of actions can be fundamental for the entire video, and discarding them makes the summary incomplete.

4.4 Discussion

Overall, results obtained from unsupervised methods on the original splits were close to the reported ones. Yet, the many videos that were from the test splits (and other videos evaluated twice) led to an unfair evaluation. The error analysis determined that existing methods have difficulties with videos filmed using moving camera settings. These difficulties can be attributed to two main reasons. First, the evaluated methods are basically trained using only object-based features that process only frame-level information. Second, those methods create a representative summary that has a similar distribution to the original video without considering the relationships between video segments. Our approach addressed the first issue and presents a corresponding solution.

5 CONCLUSIONS

In this paper, we evaluated the state-of-the-art unsupervised video summarization methods and proposed a solution to bridge the existing gaps. Therefore, we have proposed multi-source features with chunk and stride fusion to provide more information about the

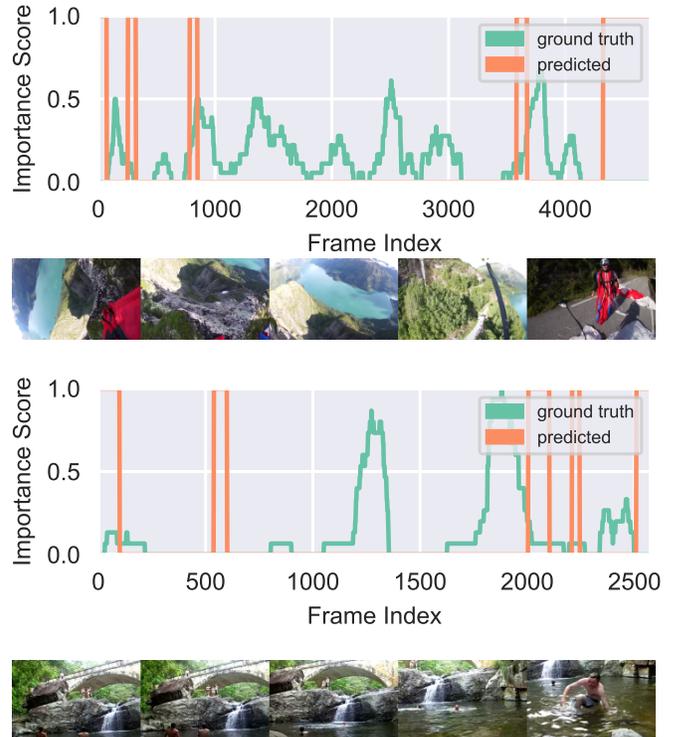


Figure 2: The upper illustrations demonstrate the generated summary of $video_2$ (Base Jumping) and $video_{24}$ (Paluma Jump) using the *MCSF* model compared to the mean of the reference summaries. The lower images show the centers of the parts selected by the generated summary. In the upper illustration, the model is confused by the abrupt visual changes at the end of $video_2$. In the lower illustration, the model does not capture the moving object in the background of $video_{24}$. In both videos, the model fails to capture the jumping event.

visual content. For the evaluation task, we re-implemented two methods and reproduced their reported results. Furthermore, all the methods are fairly compared using two different evaluation metrics and different kinds of splits. By applying the late fusion variation on *TVSum*, our approach achieved better results than the state-of-the-art methods when using the fair evaluation scheme with re-organized data splits. For *SumMe*, there was an improvement on the non-overlapping splits with an intermediate fusion. Eventually, observations of the existing methods were made based on the obtained results, and a detailed video-wise qualitative analysis on the causes for the current shortcomings was conducted. In the future, we will explore the incorporation of features from actions and other modalities to enhance model performance.

REFERENCES

- [1] Cisco and/or its affiliates. 2018. Global - 2022 Forecast Highlights. https://www.cisco.com/c/dam/m/en_us/solutions/service-provider/vni-forecast-highlights/pdf/Global_2022_Forecast_Highlights.pdf.
- [2] Evlampios Apostolidis, Eleni Adamantidou, Alexandros I. Metsai, Vasileios Mezaris, and Ioannis Patras. 2020. Unsupervised Video Summarization via Attention-Driven Adversarial Learning. In *MultiMedia Modeling*, Yong Man Ro, Wen-Huang Cheng, Junmo Kim, Wei-Ta Chu, Peng Cui, Jung-Woo Choi, Min-Chun Hu, and Wesley de Neve (Eds.). Lecture Notes in Computer Science, Vol. 11961. Springer International Publishing, Cham, 492–504. https://doi.org/10.1007/978-3-030-37731-1_40
- [3] Evlampios Apostolidis, Alexandros I. Metsai, Eleni Adamantidou, Vasileios Mezaris, and Ioannis Patras. 2019. A Stepwise, Label-based Approach for Improving the Adversarial Training in Unsupervised Video Summarization. In *Proceedings of the 1st International Workshop on AI for Smart TV Content Production, Access and Delivery - AI4TV '19*, Raphaël Troncy, Jorma Laaksonen, Hamed R. Tavakoli, Lyndon Nixon, and Vasileios Mezaris (Eds.). ACM Press, New York, New York, USA, 17–25. <https://doi.org/10.1145/3347449.3357482>
- [4] Sandra Eliza Fontes de Avila, Ana Paula Brandão Lopes, Antonio da Luz Jr., and Arnaldo de Albuquerque Araújo. 2011. VSUMM: A mechanism designed to produce static video summaries and a novel evaluation method. *Pattern Recognit. Lett.* 32, 1 (2011), 56–68. <https://doi.org/10.1016/j.patrec.2010.08.004>
- [5] Sidney K. D’Mello and Jacqueline M. Kory. 2015. A Review and Meta-Analysis of Multimodal Affect Detection Systems. *ACM Comput. Surv.* 47, 3 (2015), 43:1–43:36. <https://doi.org/10.1145/2682899>
- [6] Michael Gygli, Helmut Grabner, Hayko Riemenschneider, and Luc Van Gool. 2014. Creating Summaries from User Videos. In *Computer Vision - ECCV 2014 - 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part VII (Lecture Notes in Computer Science, Vol. 8695)*, David J. Fleet, Tomás Pajdla, Bernt Schiele, and Tinne Tuytelaars (Eds.). Springer, 505–520. https://doi.org/10.1007/978-3-319-10584-0_33
- [7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep Residual Learning for Image Recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*. IEEE Computer Society, 770–778. <https://doi.org/10.1109/CVPR.2016.90>
- [8] Yunjae Jung, Donghyeon Cho, Dahun Kim, Sanghyun Woo, and In So Kweon. 2019. Discriminative Feature Learning for Unsupervised Video Summarization. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019*. AAAI Press, 8537–8544. <https://doi.org/10.1609/aaai.v33i01.33018537>
- [9] Behrooz Mahasseni, Michael Lam, and Sinisa Todorovic. 2017. Unsupervised Video Summarization with Adversarial LSTM Networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*. IEEE Computer Society, 2982–2991. <https://doi.org/10.1109/CVPR.2017.318>
- [10] Mrigank Rochan and Yang Wang. 2019. Video Summarization by Learning From Unpaired Data. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*. Computer Vision Foundation / IEEE, 7902–7911. <https://doi.org/10.1109/CVPR.2019.00809>
- [11] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael S. Bernstein, Alexander C. Berg, and Fei-Fei Li. 2015. ImageNet Large Scale Visual Recognition Challenge. *Int. J. Comput. Vis.* 115, 3 (2015), 211–252. <https://doi.org/10.1007/s11263-015-0816-y>
- [12] Yale Song, Jordi Vallmitjana, Amanda Stent, and Alejandro Jaimes. 2015. TVSum: Summarizing web videos using titles. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*. IEEE Computer Society, 5179–5187. <https://doi.org/10.1109/CVPR.2015.7299154>
- [13] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott E. Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. 2015. Going deeper with convolutions. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*. IEEE Computer Society, 1–9. <https://doi.org/10.1109/CVPR.2015.7298594>
- [14] Gökhan Yalınz. 2019. *Unsupervised Video Summarization With Independently Recurrent Neural Networks And Multiple Rewards*. Master’s thesis. Fen Bilimleri Enstitüsü.
- [15] Li Yuan, Francis E. H. Tay, Ping Li, Li Zhou, and Jiashi Feng. 2019. Cycle-SUM: Cycle-Consistent Adversarial LSTM Networks for Unsupervised Video Summarization. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019*. AAAI Press, 9143–9150. <https://doi.org/10.1609/aaai.v33i01.33019143>
- [16] Bin Zhao, Xuelong Li, and Xiaoqiang Lu. 2018. HSA-RNN: Hierarchical Structure-Adaptive RNN for Video Summarization. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*. IEEE Computer Society, 7405–7414. <https://doi.org/10.1109/CVPR.2018.00773>
- [17] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. 2017. Places: A 10 million Image Database for Scene Recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2017).
- [18] Kaiyang Zhou, Yu Qiao, and Tao Xiang. 2018. Deep Reinforcement Learning for Unsupervised Video Summarization With Diversity-Representativeness Reward. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th Innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018*, Sheila A. McIlraith and Kilian Q. Weinberger (Eds.). AAAI Press, 7582–7589. <https://www.aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/16395>