

Facilitating the Production of Well-tailored Video Summaries for Sharing on Social Media^{*}

Evlampios Apostolidis^[0000-0001-5376-7158], Konstantinos Apostolidis^[0000-0002-9470-6332], and Vasileios Mezaris^[0000-0002-0121-4364]

Information Technologies Institute - Centre for Research and Technologies, Hellas,
6th km Charilaou - Thermi Road, 57001, Thessaloniki, Greece
{apostolid,kapost,bmezaris}@iti.gr

Abstract. This paper presents a web-based tool that facilitates the production of tailored summaries for online sharing on social media. Through an interactive user interface, it supports a “one-click” video summarization process. Based on the integrated AI models for video summarization and aspect ratio transformation, it facilitates the generation of multiple summaries of a full-length video according to the needs of target platforms with regard to the video’s length and aspect ratio.

Keywords: Video summarization · Video aspect ratio transformation · Saliency prediction · Artificial Intelligence · Social media.

1 Introduction

Social media users crave short videos that attract the viewers’ attention and can be ingested quickly. Therefore, for sharing on social media platforms, video creators often need a trimmed-down version of their original full-length video. However, different platforms impose different restrictions on the duration and aspect ratio of the video that they accept, e.g., on Facebook’s feed videos up to 2 min. appear in a 16:9 ratio, whereas Instagram and Facebook stories usually allow for 20 sec. and are shown in a 9:16 ratio. This makes the generation of tailored versions of video content for sharing on multiple platforms a tedious task. In this paper, we introduce a web-based tool that harnesses the power of AI (Artificial Intelligence) to automatically generate video summaries that encapsulate the flow of the story and the essential parts of the full-length video and are already adapted to the needs of different social media platforms in terms of video length and aspect ratio.

2 Related Work

Several video summarization tools can be found online, that are based on AI models. However, most of them produce a textual summary of the video, by analyzing the available [6,7] or automatically-extracted transcripts [8,1] using NLP

^{*} This work was supported by the EU Horizon 2020 programme under grant agreement H2020-951911 AI4Media.

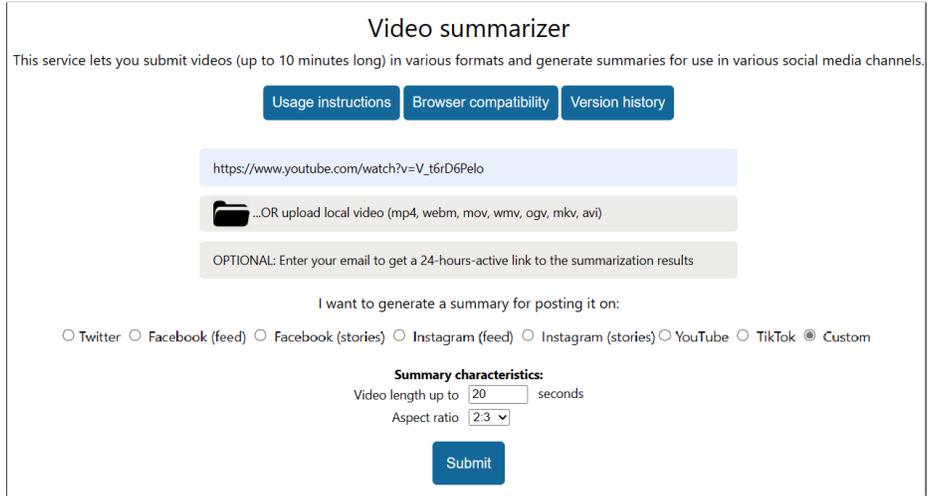
(Natural Language Processing) models. Similar solutions are currently provided by various browser plugins [4,9]. Going one step further, another tool creates a video summary by stitching in chronological order the parts of the video that correspond to the selected transcripts for inclusion in the textual summary [5]. Focusing on the visual content, Cloudinary released a tool that allows users to upload a video and receive a summarized version of it based on a user-defined summary length (max 30 sec.) [2]. In addition, Cognitive Mill released a paid AI-based platform which, among other media content management tasks, supports the semi-automatic production of movie trailers and video summaries [3]. The proposed solution in this paper is most closely related to the tool of Cloudinary. However, contrary to this tool, our solution offers various options for the video summary duration and applies video aspect ratio transformation techniques to fully meet the specifications of the target video-sharing platform.

3 Proposed Solution

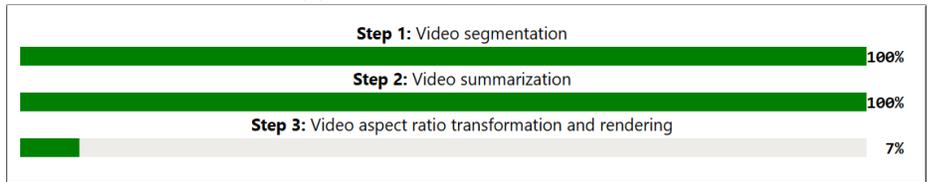
The proposed solution (available at <https://idt.iti.gr/summarizer>) is an extension of the web-based service for video summarization, presented in [16]. It is composed of a front-end user interface (UI) that allows interaction with the user (presented in Sec. 3.1), and a back-end component that analyses the video and produces the video summary (discussed in Sec. 3.2). The front-end and back-end communication is carried out via REST calls that initiate the analysis, periodically request its status, and, after completion, retrieve the analysis results (i.e., the video summary) for presentation to the user. Our solution extends [16] by: i) using an advanced AI-based method for video summarization, ii) integrating an AI-based approach for spatially cropping the video given a target aspect ratio, and iii) supporting customized values for the target duration and aspect ratio of the generated video summary.

3.1 Front-end UI

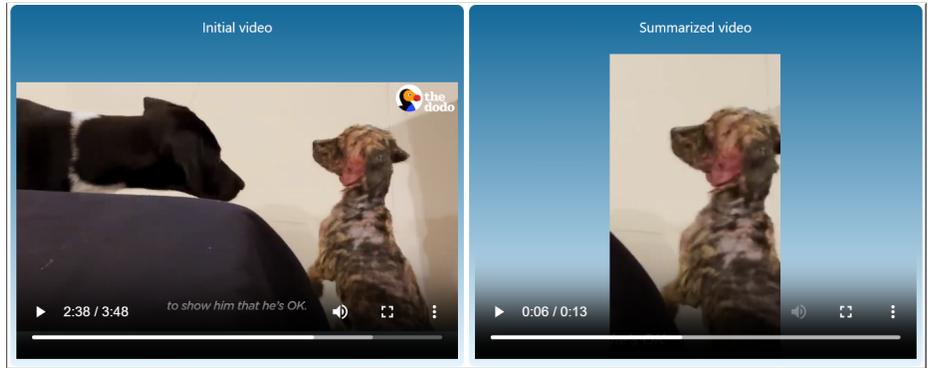
The UI of the proposed solution (see the top part of Fig. 1) allows the user to submit a video (that is either online-available or locally-stored in the user's device) for summarization, and choose the duration and aspect ratio of the produced summary. This choice can be made either by selecting among presets for various social media channels, or in a fully-custom manner. After initiating the analysis, the user can monitor its progress (see the middle part of Fig. 1) and submit additional requests while the previous ones are being analyzed. When the analysis is completed, the original video and the produced summary are shown to the user through an interactive page containing two video players that support all standard functionalities (see the bottom part of Fig. 1); through the same page, the user is able to download the produced video summary. Further details about the supported online sources, the permitted file types, and the management of the submitted and produced data, can be found in [16].



(a) The landing page of the UI.



(b) The progress-reporting bars.



(c) The video players of the page showing the analysis results.

Fig. 1: Instances of the updated and extended UI.

3.2 Back-end component

The submitted video is initially fragmented to shots using a pre-trained model of the method from [28], which exhibits 11% improved performance on the RAI dataset [15] compared to the previous (used in [16]) approach. Following, **video**

Table 1: Performance (F-Score (%)) of SUM-GAN-AEE and AC-SUM-GAN on SumMe and TVSum; the last row reports AC-SUM-GAN’s performance for augmented training data.

| Method | SumMe | TVSum |
|--------------------------------------|-------------|-------------|
| SUM-GAN-AAE [10] (used in [16]) | 48.9 | 58.3 |
| AC-SUM-GAN [11] | 50.8 | 60.6 |
| AC-SUM-GAN _{aug} (used now) | 52.0 | 61.0 |

Table 2: Used datasets for training and evaluating the AC-SUM-GAN model.

| Dataset | Videos | Duration (min) | Content |
|--------------|--------|----------------|--|
| SumMe [18] | 25 | 1-6 | holidays, events, sports, travelling, cooking |
| TVSum [27] | 50 | 1 - 11 | news, how-to’s, user-generated, documentaries |
| OVP [17] | 50 | 1 - 4 | documentary, educational, ephemeral, historical, lecture |
| YouTube [17] | 50 | 1 - 10 | cartoons, sports, tv-shows, commercial, home videos |

summarization is performed using a pre-trained model of AC-SUM-GAN [11], a top-performing unsupervised video summarization method [12]. This method embeds an Actor-Critic model into a Generative Adversarial Network and formulates the selection of important video fragments as a sequence generation task. At training time, the Actor-Critic model utilizes the Discriminator’s feedback as a reward, to progressively explore a space of actions and learn a value function (Critic) and a policy (Actor) for key-fragment selection. As shown in Table 1, AC-SUM-GAN performs much better than SUM-GAN-AAE [10] (used in [16]), on the SumMe [18] and TVSum [27] benchmark datasets for video summarization. Both methods learn the task with the help of a summary-to-video reconstruction mechanism and using the received feedback from an adversarially-trained Discriminator. We argue that the advanced performance of AC-SUM-GAN relates to the use of this feedback as a reward for training an Actor-Critic model and learning a good policy for key-fragment selection, rather than using it as part of a loss function to train a bi-directional LSTM for frame importance estimation.

The proposed solution uses a model of AC-SUM-GAN that has been trained using augmented data. Following the typical approach in the literature [12], we extended the pool of training samples of the SumMe and TVSum datasets, by including videos of the OVP and YouTube [17] datasets. As presented in Table 2, the utilized data include videos from various categories and thus facilitate the training of a general-purpose video summarization model. Nevertheless, we anticipate a better summarization performance on videos from the different categories found in the used datasets, such as tutorials, “how-to”, product demos, and event videos (e.g., birthday parties) that are commonly shared on social me-

Table 3: Performance comparison (F-Score (%)) with SotA unsupervised approaches after using augmented training data. The reported scores for the listed methods are from the corresponding papers.

| Method | SumMe | TVSum |
|-------------------------------|-------------|-------------|
| ACGAN [19] | 47.0 | 58.9 |
| RSGN _{unsup} [29] | 43.6 | 59.1 |
| 3DST-UNet [23] | 49.5 | 58.4 |
| DSR-RL-GRU [25] | 48.5 | 59.2 |
| ST-LSTM [24] | 52.0 | 58.1 |
| CAAN [22] | 50.9 | 59.8 |
| SUM-GDA _{unsup} [21] | 50.2 | 60.5 |
| SUM-FCN _{unsup} [26] | 51.1 | 59.2 |
| AC-SUM-GAN _{aug} | 52.0 | 61.0 |

Table 4: Video aspect ratio transformation performance (IoU (%)) on the RetargetVid dataset.

| | Method | Worst | Best | Mean |
|-------------------------|-------------------------------|-------------|-------------|-------------|
| 1:3 target aspect ratio | SVC (used in [14]) | 51.7 | 53.8 | 52.9 |
| | SVC _{ext} (used now) | 53.8 | 57.6 | 55.6 |
| 3:1 target aspect ratio | SVC (used in [14]) | 74.4 | 77.0 | 75.3 |
| | SVC _{ext} (used now) | 76.3 | 78.0 | 77.6 |

dia, compared to the expected performance on videos from completely unseen categories, such as movies, football games and music shows. This data augmentation process resulted in improvements on both benchmarking datasets (see the last row of Table 1) and to a very competitive performance against several state-of-the-art (SotA) unsupervised methods from the literature that have been assessed under the same evaluation settings (see Table 3).

To minimize the possibility of losing semantically-important visual content or resulting in visually-unpleasant results during **video aspect ratio transformation** (that would be highly possible when using naive approaches, such as fixed cropping of a central area of the video frames, or padding of black borders to reach the target aspect ratio) the proposed solution integrates an extension of the smart video cropping (SVC) method of [14]. The latter starts by computing the saliency map for each frame that was chosen for inclusion in the video summary. Then, to select the main part of the viewers’ focus, the integrated method applies a filtering-through-clustering procedure on the pixel values of each predicted saliency map. Finally, it infers a single point as the center of the viewer’s attention and computes a crop window for each frame based on the displacement of this point. The applied extension on [14], relates to the use of a SotA method for saliency prediction [20], that resulted in improved performance on the RetargetVid dataset [13]. As shown in Table 4, the averaged Intersection-over-Union (IoU) scores for all video frames have been increased by over 2 percentage points.

4 Conclusions

In this paper, we presented a web-based tool that facilitates the generation of video summaries that are tailored to the specifications of various social media platforms, in terms of video length and aspect ratio. After reporting on the applied extensions to a previous instance of the tool, we provided more details about the front-end user interface and the back-end component of this technology, and we documented the advanced performance of the newly integrated AI models for video summarization and aspect ratio transformation. This tool will be demonstrated in MMM2024, while the participants in the relevant demonstration session will have the opportunity to test our tool in real-time.

References

1. Brevify: Video Summarizer. <https://devpost.com/software/brevify-video-summarizer>, accessed: 2023-09-29
2. Cloudinary: Easily create engaging video summaries. <https://smart-ai-transformations.cloudinary.com>, accessed: 2023-09-29
3. Cognitive Mill: Cognitive Computing Cloud Platform For Media And Entertainment. <https://cognitivemill.com>, accessed: 2023-09-29
4. Eightify: Youtube Summary with ChatGPT. <https://chrome.google.com/webstore/detail/eightify-youtube-summary/cdcpabkolgalpgeingbdcebojebfelgb>, accessed: 2023-09-29
5. Pictory: Automatically summarize long videos. <https://pictory.ai/pictory-features/auto-summarize-long-videos>, accessed: 2023-09-29
6. summarize.tech: AI-powered video summaries. <https://www.summarize.tech>, accessed: 2023-09-29
7. Video Highlight: the fastest way to summarize and take notes from videos. <https://videohighlight.com>, accessed: 2023-09-29
8. Video Summarizer - Summarize Youtube Videos. <https://mindgrasp.ai/video-summarizer>, accessed: 2023-09-29
9. VidSummize - AI YouTube Summary with Chat GPT. <https://chrome.google.com/webstore/detail/vidsummize-ai-youtube-sum/gidcfccogfdmkfdmfhdfmfnibafoopic>, accessed: 2023-09-29
10. Apostolidis, E., Adamantidou, E., Metsai, A.I., Mezaris, V., Patras, I.: Unsupervised video summarization via attention-driven adversarial learning. In: Proc. 26th Int. Conf. MultiMedia Modeling (MMM), Part I 26. pp. 492–504. Springer (2020)
11. Apostolidis, E., Adamantidou, E., Metsai, A.I., Mezaris, V., Patras, I.: AC-SUM-GAN: Connecting actor-critic and generative adversarial networks for unsupervised video summarization. *IEEE Trans. on Circuits and Systems for Video Technology* **31**(8), 3278–3292 (2021). <https://doi.org/10.1109/TCSVT.2020.3037883>
12. Apostolidis, E., Adamantidou, E., Metsai, A.I., Mezaris, V., Patras, I.: Video summarization using deep neural networks: A survey. *Proceedings of the IEEE* **109**(11), 1838–1863 (2021). <https://doi.org/10.1109/JPROC.2021.3117472>
13. Apostolidis, K., Mezaris, V.: A fast smart-cropping method and dataset for video retargeting. In: Proc. IEEE Int. Conf. on Image Processing (ICIP). pp. 1956–1960 (2021)
14. Apostolidis, K., Mezaris, V.: A web service for video smart-cropping. In: 2021 IEEE Int. Symposium on Multimedia (ISM). pp. 25–26. IEEE (2021)

15. Baraldi, L., Grana, C., Cucchiara, R.: Shot and scene detection via hierarchical clustering for re-using broadcast video. In: Proc. 16th Int. Conf. Computer Analysis of Images and Patterns (CAIP), Part I 16. pp. 801–811. Springer (2015)
16. Collyda, C., Apostolidis, K., Apostolidis, E., Adamantidou, E., Metsai, A.I., Mezaris, V.: A web service for video summarization. In: ACM Int. Conf. on Interactive Media Experiences (IMX). pp. 148–153 (2020)
17. De Avila, S.E.F., Lopes, A.P.B., da Luz Jr, A., de Albuquerque Araújo, A.: VSUMM: A mechanism designed to produce static video summaries and a novel evaluation method. *Pattern recognition letters* **32**(1), 56–68 (2011)
18. Gygli, M., Grabner, H., Riemenschneider, H., Van Gool, L.: Creating summaries from user videos. In: 13th European Conf. on Computer Vision (ECCV), Zurich, Switzerland, 2014, Proc., Part VII 13. pp. 505–520. Springer (2014)
19. He, X., Hua, Y., Song, T., Zhang, Z., Xue, Z., Ma, R., Robertson, N., Guan, H.: Unsupervised video summarization with attentive conditional generative adversarial networks. In: Proc. of the 27th ACM Int. Conf. on Multimedia (MM '19). pp. 2296–2304. ACM, New York, NY, USA (2019)
20. Hu, F., Palazzo, S., Salanitri, F.P., Bellitto, G., Moradi, M., Spampinato, C., McGuinness, K.: Tinyhd: Efficient video saliency prediction with heterogeneous decoders using hierarchical maps distillation. In: Proc. of the IEEE/CVF Winter Conf. on Applications of Computer Vision. pp. 2051–2060 (2023)
21. Li, P., Ye, Q., Zhang, L., Yuan, L., Xu, X., Shao, L.: Exploring global diverse attention via pairwise temporal relation for video summarization. *Pattern Recognition* **111**, 107677 (2021)
22. Liang, G., Lv, Y., Li, S., Zhang, S., Zhang, Y.: Video summarization with a convolutional attentive adversarial network. *Pattern Recognition* **131**, 108840 (2022)
23. Liu, T., Meng, Q., Huang, J.J., Vlontzos, A., Rueckert, D., Kainz, B.: Video summarization through reinforcement learning with a 3d spatio-temporal u-net. *Trans. Img. Proc.* **31**, 1573–1586 (jan 2022)
24. Min, H., Ruimin, H., Zhongyuan, W., Zixiang, X., Rui, Z.: Spatiotemporal two-stream lstm network for unsupervised video summarization. *Multimedia Tools and Applications* **81**, 40489–40510 (2022)
25. Phaphuangwittayakul, A., Guo, Y., Ying, F., Xu, W., Zheng, Z.: Self-attention recurrent summarization network with reinforcement learning for video summarization task. In: Proc. of the 2021 IEEE Int. Conf. on Multimedia and Expo (ICME). pp. 1–6 (2021). <https://doi.org/10.1109/ICME51207.2021.9428142>
26. Rochan, M., Ye, L., Wang, Y.: Video summarization using fully convolutional sequence networks. In: Ferrari, V., Hebert, M., Sminchisescu, C., Weiss, Y. (eds.) *Europ. Conf. on Computer Vision (ECCV) 2018*. pp. 358–374. Springer International Publishing, Cham (2018)
27. Song, Y., Vallmitjana, J., Stent, A., Jaimes, A.: TVSum: Summarizing web videos using titles. In: IEEE Conf. on Computer Vision and Pattern Recognition (CVPR). pp. 5179–5187 (2015). <https://doi.org/10.1109/CVPR.2015.7299154>
28. Souček, T., Lokoč, J.: Transnet V2: an effective deep network architecture for fast shot transition detection. *arXiv preprint arXiv:2008.04838* (2020)
29. Zhao, B., Li, H., Lu, X., Li, X.: Reconstructive sequence-graph network for video summarization. *IEEE Trans. on Pattern Analysis and Machine Intelligence* pp. 1–1 (2021). <https://doi.org/10.1109/TPAMI.2021.3072117>