
This is an electronic reprint of the original article.
This reprint may differ from the original in pagination and typographic detail.

Harrando, Ismail; Reboud, Alison; Lisena, Pasquale; Troncy, Raphaël; Laaksonen, Jorma; Virkkunen, Anja; Kurimo, Mikko

Using Fan-Made Content, Subtitles and Face Recognition for Character-Centric Video Summarization

Published in:
Proceedings of the TRECVID 2020 Workshop

Published: 08/12/2020

Document Version
Peer-reviewed accepted author manuscript, also known as Final accepted manuscript or Post-print

Please cite the original version:
Harrando, I., Reboud, A., Lisena, P., Troncy, R., Laaksonen, J., Virkkunen, A., & Kurimo, M. (2020). Using Fan-Made Content, Subtitles and Face Recognition for Character-Centric Video Summarization. In *Proceedings of the TRECVID 2020 Workshop* National Institute of Standards and Technology (NIST).

This material is protected by copyright and other intellectual property rights, and duplication or sale of all or part of any of the repository collections is not permitted, except that material may be duplicated by you for your research use or educational purposes in electronic or print form. You must obtain permission for any other use. Electronic or print copies may not be offered, whether for sale or otherwise to anyone who is not an authorised user.

Using Fan-Made Content, Subtitles and Face Recognition for Character-Centric Video Summarization

ISMAIL HARRANDO*, ALISON REBOUD*, PASQUALE LISENA, and RAPHAËL TRONCY, EURECOM, France
JORMA LAAKSONEN, ANJA VIRKKUNEN, and MIKKO KURIMO, Aalto University, Finland

This paper describes a fan-driven and character-centered approach proposed by the MeMAD team for the 2020 TRECVID [Awad et al. 2020] Video Summarization Task. In terms of data, besides the provided videos, scripts and master shot boundaries, our approach relies on fan-made content, more precisely on the BBC EastEnders episode synopses from its Fandom Wiki¹. We also use BBC EastEnders characters' images crawled from a search engine to train a face recognition system. All our runs use the same method, but with varying constraints of the number of shots and the maximum duration. The shots included in the summaries are the ones whose transcripts and visual content have the highest similarity with sentences from the synopsis. The runs submitted are as follows:

- MeMAD1:5 shots with highest similarity scores and the total duration is < 150 sec
- MeMAD2:10 shots with highest similarity scores and the total duration is < 300 sec
- MeMAD3: 15 shots with highest similarity scores and the total duration is < 450 sec
- MeMAD4: 20 shots with highest similarity scores and the total duration is < 600 sec

Surprisingly, the scores obtained for each run were very similar for the questions answering part. Only for the character Ryan, one question more was answered by choosing 15 shots rather than less. For all our runs, the redundancy score improved with the number of shots included in the summary while the relation with the scores for tempo and contextuality seem more varying. The scores were lower for the question answering evaluation part. This is rather unsurprising to us as we realized while deciding on a similarity measure score that it was rather challenging for humans too to choose between two potentially interesting moments without knowing beforehand the questions included in the evaluation set. Overall, we consider that the results obtained speak in favour of using fan-made content as a starting point for such a task. As we did not try to optimize for tempo and contextuality, we believe there is some margin for improvement here, however the task of answering unknown questions remains challenging.

ACM Reference Format:

Ismail Harrando, Alison Reboud, Pasquale Lisena, Raphaël Troncy, Jorma Laaksonen, Anja Virkkunen, and Mikko Kurimo. 2020. Using Fan-Made

*Both authors contributed equally to this research.

¹https://eastenders.fandom.com/wiki/EastEnders_Wiki

Authors' addresses: Ismail Harrando, ismail.harrando@eurecom.fr; Alison Reboud, alison.reboud@eurecom.fr; Pasquale Lisena, pasquale.lisena@eurecom.fr; Raphaël Troncy, raphael.troncy@eurecom.fr, EURECOM, Sophia Antipolis, France; Jorma Laaksonen, jorma.laaksonen@aalto.fi; Anja Virkkunen, anja.virkkunen@aalto.fi; Mikko Kurimo, mikko.kurimo@aalto.fi, Aalto University, Espoo, Finland.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2020 Association for Computing Machinery.

XXXX-XXXX/2020/11-ART \$15.00

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

Content, Subtitles and Face Recognition for Character-Centric Video Summarization. 1, 1 (November 2020), 3 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

Considering video summarization as an important task for digital content retrieval and reuse, the TRECVID [Awad et al. 2020] Video Summarization Task (VSUM) 2020 aims at fostering the research in the field by asking its participants to automatically summarize "the major life events of specific characters over a number of weeks of programming on the BBC EastEnders TV series"². More precisely, for three different characters of the series, the participants had to submit 4 summaries with 5, 10, 15 and 20 automatically selected shots. These generated summaries are evaluated by the assessors according to their tempo, contextuality and redundancy as well as with regards to how well they answered a set of questions unknown to the participants before submission. In addition to the videos, the episodes transcripts are provided by the organizers. We propose a character centered content summary approach based on fan-written synopses. The approach relies on scraping the Fandom EastEnders Wiki content for the episode synopsis and casting, in order to align them with the corresponding episodes. We include the shots that obtain the best similarity score with a sentence from the synopsis in our runs.

2 APPROACH

Our fan-driven and character centered approach is presented in Figure 1.

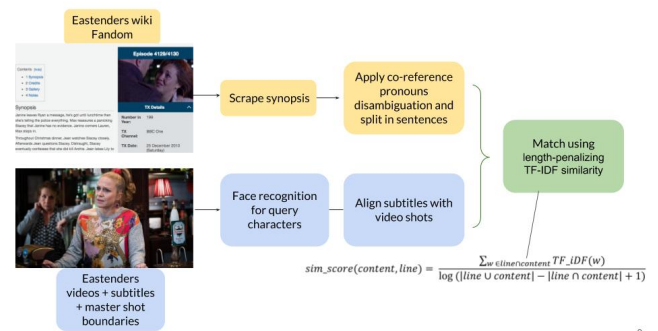


Fig. 1. Fan-driven and character centered approach

²<https://www-nlpir.nist.gov/projects/tv2020/vsum.html>

2.1 Corpus creation: Scraping synopses from the Fandom Wiki and shots selection

The first step of our approach consists in scraping synopses found on the Fandom EastEnders Wiki³. Our main hypothesis is that every sentence (ending with a period) represents an important event to be added to the final video summary. We scrape the Synopsis and the Cast sections for each episode broadcasted between the dates of the provided episodes. The mapping between the episodes and their dates is in `eastenders.collection.xml` provided by the challenge organizers.

In parallel, we extract the shots in which the three characters of interest appear from the video. To do that, we run Face Celebrity Recognition⁴. This system relies on pictures crawled from search engines using the actor’s name as search keyword (we add "East-Enders" to their names to avoid picking pictures of different people with the same name). For each picture, faces are detected using the MTCNN algorithm and the FaceNet model is applied to obtain face embeddings. Following the assumption that the majority of faces are actually representing the searched actor, other faces – e.g. person portrayed together with the actor – are automatically filtered out by removing outliers until the cosine similarity of face embeddings has a standard deviation below a threshold of 0.24 (empirically defined). The remaining faces are used to train a multi-class SVM classifier, which is used to label the faces detected on frames. For more consistent results between frames, the Simple Online and Realtime Tracking algorithm (SORT) has been included, returning groups of detections of the same person in consecutive frames. We select the shots displaying any of the the three characters of interests, keeping only those detections having a confidence score greater than 0.5. We also tried to use speaker diarization to corroborate the visual information about the characters. However, given the limitations of the current technologies in terms of number of characters and the difficulty of identifying the character corresponding to each voice, we could not pursue the idea further.

2.2 Synopses and transcript alignment and preprocessing

Using `eastenders.collection.xml` and `eastenders.episodeDescriptions.xml`, a synopsis for each episode was created. Since these were “EastEnders Omnibus” episodes, they correspond to multiple actual weekday episodes. We use the dates and the continuation to generate one synopsis for each “long” episode (typically made of 4 episodes). We then split the synopses into sentences and performed coreference resolution on the synopses to explicit character mentions⁵. In parallel, the provided XML transcripts were also converted into timestamped text and aligned with the given shot segmentation. Finally, both the synopses sentences and shot transcripts were lower cased, stop words removed and lemmatized. We also produced automatically-generated visual captions following the method presented by the PicSOM Group of Aalto University’s submissions for the TRECVID2018 VTT task [Sjöberg et al. 2018]. The hypothesis was that by describing the visual information of a shot, visual

Table 1. Average score for each run and team

TeamRun	Percentage
MeMAD1	31%
MeMAD2	31%
MeMAD3	35%
MeMAD4	32%
NIIUIT1	9%
NIIUIT2	8%
NIIUIT3	8%
NIIUIT4	6%

captions could complement well the dialog transcript and therefore allow for a better matching between the shots and synopses sentences.

2.3 Matching and runs generation

We perform a synopsis sentence / shot transcript pairwise comparison by generating a similarity score. We define similarity between two sentences as the sum of TF-IDF weights (computed on the transcript) for each word appearing in both of them, divided by the log length of the concatenation of both sentences (thus penalizing long sentences that match with many transcript lines). We order the shot by similarity score, picking only the best match for each shot (but not the other way around). This gives us scenes we are sure to appear in the summary, but not necessarily any guarantee about how important these scenes are. We also performed the pairwise comparison adding the automatically generated captions. A qualitative assessment revealed, however, that the captions were too noisy to complement well the transcript. We also make sure that if a line of dialog runs through the next shot, we include the next shot as well (to improve the smoothness of the viewing). This however was only relevant for the longest run (20 shots). Each run is made by picking up the N most matching shots out of the top, in chronological order.

3 RESULTS AND ANALYSIS

The average results for every team and all runs are presented on Table 1 and the MeMAD’s detailed scores are presented in Table 2. MeMAD obtained the best scores for every run. The mean scores (range 1 - 7. High is best) for tempo, contextuality and redundancy are all above average (respectively 4.75,4.75,4.1). However in terms of question answering, the results show that the shots selected did not allow to answer more than two (at best) of the five questions.

4 DISCUSSION AND OUTLOOK

This paper describes a character centered video summarization method based on fan-made content, subtitles and face recognition. One of the key contribution of this paper is to have demonstrated that despite some noise from Face Detection and matching, this method captured multiple important plot points for all three query characters. We also conclude that adding more shots to the summaries did, quite surprisingly, not always allow to answer more key moments related questions.

³<https://eastenders.fandom.com/wiki/EastEndersWiki>

⁴<https://github.com/D2KLab/Face-Celebrity-Recognition>

⁵<https://github.com/huggingface/neuralcoref>

Table 2. Detailed score for MeMAD's team

Query	Tempo	Contextuality	Redundancy	Q1	Q2	Q3	Q4	Q5
Janine1	6	4	5	No	No	No	No	Yes
Janine2	5	5	6	No	No	No	No	Yes
Janine3	5	5	6	No	No	No	No	Yes
Janine4	5	5	7	No	No	No	No	Yes
Ryan1	4	5	3	No	No	No	No	Yes
Ryan2	5	5	3	No	No	No	No	Yes
Ryan3	3	4	5	No	No	No	Yes	Yes
Ryan4	2	3	5	No	No	No	Yes	Yes
Stacey1	6	5	2	No	Yes	No	No	No
Stacey2	6	5	2	No	Yes	No	No	No
Stacey3	6	6	2	No	Yes	No	No	No
Stacey4	4	5	4	No	Yes	No	No	No

REFERENCES

- George Awad, Asad A. Butt, Keith Curtis, Yooyoung Lee, Jonathan Fiscus, Afzal Godil, Andrew Delgado, Jesse Zhang, Eliot Godard, Lukas Diduch, Jeffrey Liu, Alan F. Smeaton, Yvette Graham, Gareth J. F. Jones, Wessel Kraaij, and Georges Quénot. 2020. TRECVID 2020: comprehensive campaign for evaluating video retrieval tasks across multiple application domains. In *Proceedings of TRECVID 2020*. NIST, USA.
- Mats Sjöberg, Hamed R. Tavakoli, Zhicun Xu, Héctor Laria Mantecón, and Jorma Laaksonen. 2018. PicSOM Experiments in TRECVID 2018. In *Proceedings of the TRECVID 2018 Workshop*. Gaithersburg, MD, USA.