

MediaEval 2017 Predicting Media Interestingness Task

Claire-Hélène Demarty¹, Mats Sjöberg², Bogdan Ionescu³, Thanh-Toan Do⁴,
Michael Gygli⁵, Ngoc Q. K. Duong¹

¹Technicolor, Rennes, France

²Dept. of Computer Science and Helsinki Institute for Information Technology HIIT, University of Helsinki, Finland

³LAPI, University Politehnica of Bucharest, Romania

⁴University of Science, Vietnam, University of Adelaide, Australia

⁵ETH Zurich, Switzerland & Gifs.com, US

ABSTRACT

NOT CHANGED - copy-past from last year This paper provides an overview of the Predicting Media Interestingness task that is organized as part of the MediaEval 2017 Benchmarking Initiative for Multimedia Evaluation. The task, which is running for the second year, expects participants to create systems that automatically select images and video segments that are considered to be the most interesting for a common viewer. In this paper, we present the task use case and challenges, the proposed data set and ground truth, the required participant runs and the evaluation metrics.

1 INTRODUCTION

Bogdan, could you help for this?

2 TASK DESCRIPTION

NOT REALLY CHANGED - This comes from the wiki. Ngoc, if you want to rephrase, please help.

The Predicting Media Interestingness Task was proposed for the first time last year. This year's edition is a follow-up which builds incrementally upon the previous experience. The task requires participants to automatically select images and/or video segments which are considered to be the most interesting for a common viewer. Interestingness of media is to be judged based on visual appearance, audio information and text accompanying the data, including movie metadata. To solve the task, participants are strongly encouraged to deploy multimodal approaches.

As in 2016, interestingness should be assessed according to a practical use case in industry, in particular at Technicolor¹, which involves helping professionals to illustrate a Video on Demand (VOD) web site by selecting some interesting frames and/or video excerpts for the movies. The frames and excerpts should be suitable in terms of helping a user to make his/her decision about whether he/she is interested in watching the whole movie. Once again, two subtasks will be offered to participants, which correspond to two types of available media content, namely images and videos. Participants are encouraged to submit to both subtasks. In both cases, the task is a binary classification task and prediction will be carried out on a per movie basis. The two tasks are:

- *Predicting Image Interestingness* Given a set of key-frames extracted from a certain movie, the task involves automatically identifying those images that viewers report to be interesting. To solve the task, participants can make use of visual content as well as accompanying metadata, e.g., Internet data about the movie, social media information, etc.
- *Predicting Video Interestingness* Given a set of video segments extracted from a certain movie, the task involves automatically identifying the segments that viewers report to be interesting. To solve the task, participants can make use of visual and audio data as well as accompanying metadata, e.g., subtitles, Internet data about the movie, etc.

3 DATA DESCRIPTION

The data is extracted from Creative Commons licensed Hollywood-like videos: 103 movie trailers and 4 continuous extracts of ca. 15min from full-length movies.

For the video interestingness subtask, the data consists of video segments obtained after a manual segmentation. These segments correspond to shots (video shots are the continuous frame sequences recorded between a camera turn on and off) for all videos but four. Their average duration is of one second. The four last videos which correspond to the full-length movie extracts cited above, were manually segmented into longer segments with an average duration of *******, to better take into account the a certain unity of meaning and the audio information of the resulting segments. For the image subtask, the data consists of collections of key-frames extracted from the video segments used for the video subtask (one key-frame per segment). This will allow comparing results from both subtasks. The extracted key-frame corresponds to the frame in the middle of each video segment. In total, 7,396 video segments and 7,396 key-frames are released in the development set, whereas the test set consists of ******* video segments and the same number of key-frames.

***** Here we could just copy-past what was said for the provided features last year or rewrite everything. Should we talk about the bug that was discovered in the face-tracking feature and the advice not to take into account negative values ? I am not sure of this. I am copying text from last year below. If we agree that this needs rephrasing, Michael, can you help for this part?*****

To facilitate participation from various communities, we also provide some pre-computed content descriptors, namely: *low level features* — *dense SIFT* (Scale Invariant Feature Transform) which are computed following the original work in [8], except that the local

¹<http://www.technicolor.com>

frame patches are densely sampled instead of using interest point detectors. A codebook of 300 codewords is used in the quantization process with a spatial pyramid of three layers [6]; *HoG descriptors* (Histograms of Oriented Gradients) [2] are computed over densely sampled patches. Following [12], HoG descriptors in a 2×2 neighborhood are concatenated to form a descriptor of higher dimension; *LBP* (Local Binary Patterns) [9]; *GIST* are computed based on the output energy of several Gabor-like filters (8 orientations and 4 scales) over a dense frame grid like in [10]; *color histogram* computed in the HSV space (Hue-Saturation-Value); *MFCC* (Mel-Frequency Cepstral Coefficients) computed over 32ms time-windows with 50% overlap. The cepstral vectors are concatenated with their first and second derivatives; *fc7 layer* (4,096 dimensions) and *prob layer* (1,000 dimensions) of AlexNet [5]; *mid level face detection and tracking related features*² — obtained by face tracking-by-detection in each video shot with a HoG detector [2] and the correlation tracker proposed in [3].

In addition to these frame-based features, we provide C3D [11] features, which were extracted from *fc6 layer* (4,096 dimensions) and averaged on a segment level.

4 GROUND TRUTH

Both video and image data was manually annotated in terms of interestingness by human assessors. The annotation process was performed separately for the video and image subtasks, to allow us to study the correlation between the two. A dedicated web-based annotation tool was developed by the organising team for the previous edition of the task [4]. This year we made some improvements to the web tool, and some minor changes to the annotation process were introduced. The tool has been released as free and open source software, so that others can benefit from it and contribute improvements³. Overall, more than 202 annotators participated in the annotation for the video data and 144 for the images. The cultural distribution is over 20 different countries in the world.

As in last year's setup we use a pair-wise comparison protocol [1] where annotators are provided with a pair of images/shots at a time and asked to tag which one in the pair is the more interesting for them. As a change from last year, we now ask the question in a way more directly connected to the commercial application: "Which image/video makes you more interested in watching the whole movie?". We felt this would make the annotators' decision clearer, as otherwise "interestingness" can be interpreted in many ways.

As an exhaustive annotation of all possible pairs is practically impossible due to the required human resources, a boosting selection was used instead. In particular, we used a modified version of the adaptive square design method [7], in which several annotators participated in each iteration. In this method the number of comparisons for each iteration is reduced from all possible pairs $n(n-1)/2 \sim O(n^2)$ to a subset of pairs $n(\sqrt{n}-1) \sim O(n^{3/2})$, where n is the number of shots or images. For the development set, we started from iteration 6, as we could reuse the annotations done last year. To achieve the ranking used as the basis for the next round, the pair-based annotations are aggregated with the Bradley-Terry-Luce (BTL) model computation [1] resulting in an interestingness

degree for each image/shot. Previously the same procedure was also used to get the final interestingness values. This year we used an alternative method, which took into account all pair comparisons from all rounds done this year into a single large BTL calculation. This was done mainly because we discovered afterwards that some annotations from earlier rounds had to be discarded. In the iterative approach, we could only discard annotations from the most recent round, as it would be based on the previous round's BTL output.

The final binary decisions are obtained using a thresholding scheme that tries to detect the boundary where interestingness values make the "jump" between the underlying distributions of the non interesting and interesting populations. See last year's overview paper for a more detailed description [4].

A new issue this year was that we discovered that some annotators were cheating. These annotators occasionally switched to cheating, where they simply always selected the first, or the second item as the most interesting one without actually assessing the media contents. We added some heuristic anti-cheating measures to the system, for example if the same position has been selected repeatedly for a given number of times in a given period of time that user would be logged out and presented with a notice about cheating having been detected. Votes matching these criteria were also marked as "invalid" and not included in the final BTL calculation. In the development set as much as 10% of the annotations were marked as invalid according to a particular definition. The high number can be explained by the fact that a cheater can annotate much faster than an honest annotator. Detecting cheating is ultimately a very difficult task, as malicious annotators may easily create a new pattern in response to anti-cheating measures. Furthermore, it is not possible to unambiguously detect cheat votes by just analysing the voting patterns. For example, an honest annotator may indeed get a surprisingly long sequence of votes for the same position just by chance, or a cheater may develop a method to select the items randomly. A possible idea for future versions of the system would be to occasionally present pairs for which the system already has a lot of consistent evidence for to test the annotator's honesty. If the annotator fails these several times it might be a good indication that he or she is cheating. On the other hand, it is possible that annotators cheat only on very difficult cases. In fact the difficulty of many of the judgements might be the reason why some people grow tired of annotating honestly. **This tells the story of events in the annotation process, rather than being useful for getting an overview. It is too long and needs serious rewriting. Also, a better check, in my opinion is to show the same pair again at a later time (and switched) and see if they still prefer the same one**

5 RUN DESCRIPTION

Every team can submit up to 10 runs, 5 per subtask. For each subtask, a required run is defined: *Image subtask - required run*: classification is to be achieved with the use of the visual information. External data is allowed. *Video subtask - required run*: classification is to be achieved with the use of *both* audio and visual information. External data is allowed.

Apart from these required runs, any additional run for each subtask will be considered as a general run, i.e., anything is allowed, both from the method point of view and the information sources.

²<http://multimediaeval.org/mediaeval2016/persondiscovery/>

³<https://github.com/mvsjober/pair-annotate>

6 EVALUATION

For both subtasks, the official evaluation metric will be the mean average precision at 10 (MAP@10) computed over all videos, and over the top 10 best ranked images/video shots. MAP@10 is selected because it reflects the VOD use case, where the goal is to select a small set of the most interesting images or video segments for each movie. To provide a large overview of the systems' performances, other common metrics will also be provided. All metrics will be computed by using the `trec_eval` tool from NIST⁴.

7 CONCLUSIONS

NOT CHANGED - copy-past from last year - Dont think we really need to change it. If yes, Bogdan can you do it? The 2017 Predicting Media Interestingness task provides participants with a comparative and collaborative evaluation framework for predicting content interestingness with explicit focus on multimedia approaches. Details on the methods and results of each individual participant team can be found in the working note papers of the MediaEval 2017 workshop proceedings.

ACKNOWLEDGMENTS

We would like to thank Yu-Gang Jiang and Baohan Xu from the Fudan University, China, Hervé Bredin, from LIMSI, France, and Michael Gygli for providing the features that accompany the released data. Part of the task was funded under research grant PN-III-P2-2.1-PED-2016-1065, agreement 30PED/2017, project SPOTTER.

REFERENCES

- [1] R. A. Bradley and M. E. Terry. 1952. Rank Analysis of Incomplete Block Designs: the method of paired comparisons. *Biometrika* 39 (3-4) (1952), 324–345.
- [2] N. Dalal and B. Triggs. 2005. Histograms of oriented gradients for human detection. In *IEEE CVPR Conference on Computer Vision and Pattern Recognition*.
- [3] Martin Danelljan, Gustav Hager, Fahad Shahbaz Khan, and Michael Felsberg. 2014. Accurate scale estimation for robust visual tracking. In *British Machine Vision Conference*.
- [4] Claire-Hélène Demarty, Mats Sjöberg, Bogdan Ionescu, Thanh-Toan Do, Hanli Wang, Ngoc Q.K. Duong, and Frédéric Lefebvre. 2016. MediaEval 2016 Predicting Media Interestingness Task. In *Proceedings of the MediaEval 2016 Workshop*. Hilversum, Netherlands.
- [5] Yu-Gang Jiang, Qi Dai, Tao Mei, Yong Rui, and Shih-Fu Chang. 2015. Super Fast Event Recognition in Internet Videos. *IEEE Transactions on Multimedia* 17, 8 (2015), 1–13.
- [6] S. Lazebnik, C. Schmid, and J. Ponce. 2006. Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories. In *IEEE CVPR Conference on Computer Vision and Pattern Recognition*. 2169–2178.
- [7] Jing Li, Marcus Barkowsky, and Patrick Le Callet. 2013. Boosting paired comparison methodology in measuring visual discomfort of 3DTV: performances of three different designs. In *SPIE Electronic Imaging, Stereoscopic Displays and Applications*, Vol. 8648.
- [8] D. Lowe. 2004. Distinctive image features from scale-invariant keypoints. *International Journal on Computer Vision* 60 (2004), 91–110.
- [9] T. Ojala, M. Pietikainen, and T. Maenpää. 2002. Multiresolution gray-scale and rotation invariant texture classification with local binary patterns. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 24(7) (2002), 971–987.
- [10] A. Oliva and A. Torralba. 2001. Modeling the shape of the scene: a holistic representation of the spatial envelope. *International Journal of Computer Vision* 42 (2001), 145–175.
- [11] Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. 2015. Learning spatiotemporal features with 3d convolutional networks. In *ICCV*.
- [12] J. Xiao, J. Hays, K. Ehinger, A. Oliva, and A. Torralba. 2010. SUN database: Large-scale scene recognition from abbey to zoo. In *IEEE CVPR Conference on Computer Vision and Pattern Recognition*. 3485–3492.

⁴http://trec.nist.gov/trec_eval/