

Annotating, understanding, and predicting long-term video memorability*

Extended Abstract[†]

Romain Cohendet

Technicolor

Rennes, France

romain.cohendet@technicolor.com

Ngoc Khanh Duong

Technicolor

Rennes, France

Quang-Khanh-Ngoc.Duong@technicolor.com

Karthik Yadati

University of Delft

Delft

N.K.Yadati@tudelft.nl

Claire-Helene Demarty

Technicolor

Rennes, France

Claire-Helene.Demarty@technicolor.com

ABSTRACT

Following on from the study of image memorability prediction, which draw continuous attention over the past six years, the computational understanding of video memorability's search field has recently hatched. The growing number of shared videos foster us to find new ways to make their occurrence the most relevant in our everyday lives. There is no available dataset of videos annotated in terms of memorability; such dataset may probably result in a launching of the field, as it had been the case for images. The first goal of the pioneers of this search should be to succeed to propose a protocol to constitute such a dataset. In this article, we propose a protocol of 700 videos annotated with memorability scores, constitute a dataset, study the quality of the collected data, and computationally modeling memorability to predict it.

We finally propose a deep-learning model base to predict video memorability, comparing several methods and class of features.

KEYWORDS

Video memorability, Long-term memory, Measurement protocol, Memorability scores, Deep learning, Multimedia information retrieval

ACM Reference Format:

Romain Cohendet, Karthik Yadati, Ngoc Khanh Duong, and Claire-Helene Demarty. 2018. Annotating, understanding, and predicting long-term video memorability: Extended Abstract. In *Proceedings of ACM Yokohama conference (ICMR 2018)*. ACM, New York, NY, USA, Article 4, 6 pages. https://doi.org/10.475/123_4

*Produces the permission block, and copyright information

[†]The full version of the author's guide is available as `acmart.pdf` document

Permission to make digital or hard copies of part or all 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 third-party components of this work must be honored. For all other uses, contact the owner/author(s).

ICMR 2018, June 2018, Yokohama, Japan

© 2018 Copyright held by the owner/author(s).

ACM ISBN 123-4567-24-567/08/06...\$15.00

https://doi.org/10.475/123_4

1 INTRODUCTION

Enhancing the relevance of multimedia occurrences in our everyday life requires to imagine new ways to organize – in particular, to retrieve – contents. Like other metrics of video "importance", such as aesthetics or interestingness, memorability can be regarded as a particularly relevant element to help us picking a video among several ones, with the advantage of a possibility to be measured less subjectively.

The computational understanding of video memorability's search field follows on from the study of image memorability prediction which has attracted increasing attention from the seminal work of Isola *et al.* released in 2001 [4]. Recently, models of image memorability prediction achieved very good results with the introduction of deep learning to address the challenge of image memorability prediction [1, 5, 7]. This success resulted on the extension of the challenge to videos. Video memorability prediction is however a new field, and there are only two studies, to our knowledge, that addressed this issue [2, 6].

At least two important problems could explained this scarcity of studies. Firstly, protocols used in the existing studies are not smart enough, but are questionable. Secondly, authors didn't make their data accessible for community, so no database exist. We think that the existence of such a database could have the importance for video memorability that the ones of [4] and [5] had on image memorability. This article wants to takes his part to develop such a protocol to collect high quality data, and such a dataset.

No existing datasets. => No Dataset non available

The fact that there exist no dataset to modeling video memorability is the most serious obstacle which prevent video memorability prediction to take off. This should be the very first goal of pioneers of this new research field: to provide data to community. For images, the pioneers created and made downloadable such a database [4], which enable to made image memorability research flourishing. The second released database, far larger, launched new possibility for deep learning and enable to obtained the aforementioned really good results [5]. Nous pouvons parier qu'il en sera de même pour le champ de la mémorabilité des vidéos.

The establishment of a dataset, in particular of the first dataset of this kind, may not be self-evident.

Elle nécessite de définir les points cardinaux qui vont guider et contraindre les chercheurs qui utiliseront ces premières données disponibles.

Pour preuve, l'influence considérable qu'Isola et al. (puis Khosla et al., de la même équipe, avec leur base de données beaucoup plus conséquente rendue disponible en 2015) ont eu sur les études réalisées dans leur sillage : jusqu'à aujourd'hui, la quasi-totalité des études ayant porté sur la mémorabilité des images ont adopté avec les bases de données de ces auteurs leurs définitions de la mémorabilité.

Or, la manière dont les auteurs ont mesuré la mémorabilité des images, qui n'est pratiquement jamais rediscutée, n'est qu'une des manières de procéder.

Aussi, si les modalités de prédiction de la mémorabilité des images ont obtenu d'excellents résultats (très proches de ceux qu'on obtiendrait en faisant des résultats d'un groupe d'observateur les résultats d'un autre groupe d'observateurs), qu'est-ce que ces modalités prédisent si bien, en fait ?

Videos are more complex than images. Contrary to images, videos do not constitute clearly defined units, but have supplementary dimensions – sound and movement – that makes difficult the definition of what a video is.

Aussi cette étude veut participer à réfléchir aux protocoles qui pourraient être employés pour continuer à l'avenir à un protocole de plus grande envergure.

1.1 Previous work

On comprend les auteurs d'avoir agi ainsi pour constituer une base de grand ampleur (Khosla et al.) It is long-term memory a proprement parlé because. but memories evolve in long mem => Ebbinghaus curve So memorabilities of souvenirs can change after a certain delays

To our knowledge, the first attempt at predicting video memorability had been proposed by [2].

The authors partially adapted the protocol proposed by [4] to measure image memorability for videos. In contrast to the "memory game" proposed by Isola *et al.* to collect memorability data, their protocol is however much heavier. They used a classical recognition task to measure memory for videos, which consists in two steps: a free viewing task, followed two days later by a recall task.

We can infer from the sparse information authors gave that the task duration, for each of the 20 participants, was about 24 hours (by taking an average video duration half-way between 15 and 30 sec), spread over 10 sessions (5 free viewing and 5 recall tasks) of about 2 hours each.

Authors used the same proportion of fillers in the free viewing and recall tasks (4/5 of fillers and 1/5 of targets), to "to guarantee that viewers are unaware of targets". If it was totally justified in [4] for which encoding and recall tasks were interlaced, there is here a way to alleviate the task without impacting its quality; indeed, reducing the number of fillers (never reasked videos) in free viewing task would have very weak impact on memorability scores (généralement on se passe même de fillers dans les tâches d'apprentissage du matériel).

Still, as pointed by [6], the long time span of the experiment makes difficult the generalization of this protocol, in particular to build an extensive dataset.

Regarding the prediction part, we also agree with these authors to say that "the method used fMRI measurements for predicting memorability, which would be difficult to generalize".

Another earlier approach was the one of [6]. The participants passed in crowdsourcing a free viewing task followed by a recall task, for which they had to answer question (such as: "Did you see a man juggling?", "Did you see a car on road?").

The major drawback of the study comes from this use of questions instead of a recognition task. The "memorability" measured may be due not to memorability of the videos but (also) to the differences between the questions rather than a difference of memorability of the videos. This bias could be of several kinds: difficulty, comprehension (non anglophone people in AMT), easiness of imaging), indexage plus facile... C'est d'autant plus problématique que ces différences pourraient avoir conduit à une haute inter-human consistency not due to memorability attribué à tort à un phénomène de mémoire. C'est d'autant plus problématique que le temps de réponse des participants entre dans le calcul de la mémorabilité, ce qui ne peut qu'accentuer ces biais, et est probablement responsable de la high human response consistency they observe.

Le fait que les questions doivent être créées à la main, et que les choix des questions sont assez limitatifs (par ex avec la question "Did you see a car on road?", il ne peut y avoir de toute la tâche qu'une seule fois une voiture) makes difficult the generalization of this protocol to build an extensive dataset. Il a été manuellement assuré que deux questions textuelles ou deux vidéos dans une séquence étaient similaires en contenu.

Les critiques exposées montrent à quel point il est difficile de construire la première base de données pour mesurer la mémorabilité des vidéos. Isola et al. Sont bien sortis, en se cantonnant à une mesure purement objective de la mémoire à tâche de reconnaissance classique, ce qui leur a permis de lancer le champ de recherche de la mémorabilité des images. Il faudrait en faire de même pour la mémorabilité des vidéos. Cette étude s'inscrit dans cette volonté à non pas de créer un nouveau dataset, mais d'essayer un protocole pour voir s'il est efficace.

To sum up: - Set up a protocol, which could be extended to create a large dataset - very long-term memorability

1.2 Goals of our study

These two reasons are linked. They give us our priority: build a dataset and define a protocol to do that. This dataset – and a fortiori the protocol – should avoid the drawbacks of previous work in video and image memorability.

=> To propose a new protocol to collect data/for annotation for VERY long-term memorability => To constitute a dataset to try prediction on => The dataset will have "lasting" long-term memorability annotations We want our work would benefit to video but also to image memorability.

We also study around the memorability concept psycho things, to ensure the quality of the obtained data and potentially shed light on interesting factors involved in memorability of videos.

Il semble des derni res  tudes qu'il serait int ressant de prendre en compte le temps dans le calcul de la m morabilit  [LaMem + papier sur la VM o   ils prennent en compte le temps dans leurs scores]

The semantization that can affect long-term memorability of videos

The following of the paper is organized like that: 1/ Construction of the dataset, 2/ Study of the memorability scores, and 3/ Video memorability deep-learning prediction

2 CONSTRUCTION OF THE DATASET

In this section, we describe the protocol we propose to measure the *long-term* memorability of videos. The major particularity of the protocol is that, contrary to memory test used in the previous studies in image and video memorability, there is no learning step intended to encode the material before the recall, as we measure the participants' memory established prior the experiment. Indeed, we measure the memory of participant for movies they have seen before. The major focus of this protocol is to measure "lasting" long-term memorability, instead of memorability measured a few seconds to minutes after encoding task. The memory task involves three materials: videos sequences, questionnaire to collect demographics and some data on the visualization (because encoding is not controlled) and participants.

2.1 Memory task

The task designed to measure the long-term memorability consists in two parts: 1/ a questionnaire, and 2/ a recognition task.

The 10-sec video sequences are displayed one after another, separated by an inter-stimuli interval (ISI).

Les vid es cibles sont m  l es aux vid es de remplissage, et le participant doit dire, pour chaque vid e qui lui est pr sent e, s'il l'a d j  vue ou non. To answer, he must press the spacebar during the displaying of the video sequences (than answer when they arrived during the 1-s inter-stimuli interval), movies they have seen (the targets) and never seen (the fillers). Even if the participant answers, the video continues to run until its end (to avoid the temptation of answering to faster finish the experiment).

"Respond ONLY if you remember the sequence AND NOT THE movie from which it comes".

2.2 The video sequences

A list of 100 movies (available with the downloadable data) was first established, for which we mixed their popularity and their genres. The movies quality was HD 720-1080, except for ??? dv-drip movies. original-english version (without subtitles) (occidental movies, except for Slumdog millionaire)

To constitute our dataset, we manually selected 7 video sequences of 10 seconds from each movies (for a total of 700 videos). The sequences had to meet several criteria to be selected. First, not to be part of the official trailers, to be sure that, if a participant reports not having seen the movie, he won't have seen the sequence presented for the recognition in a trailer. Second, we wanted our sequences to constitute "logical" units, that is to say to "protect the sense" of the videos. Indeed, semantics is linked to memorability of images [3], from which we can imagine that it is also the case

for videos. But, contrary to images that are clearly cut units, sequences memorability depends on the cut. To maintain a relative semantic consistency intra-sequence, the cutters were provided with the instruction to avoid to cut in the middle of a sentence, or to agglomerate very different plans together. The idea was to increase the coherence of the sequences to diminish the difficulty of the computational models to learn without impairing their capacity to generalize.

During the manual selection of the sequence, we chose 127 sequences (among the 700) we called "generic" by contrast to "typical" ones. Generic sequences are parts of the movie which do not contain any elements that would enable someone to easily guess that the sequences belong to the particular movie. The list of undesirable includes but is not limited to: scenes with recognizable famous actors (e.g. in Kill Bill, avoid the scenes with Uma Thurman, Lucy Liu, David Carradine  ); scenes famous for their music, gesture... (e.g. in Kill Bill, avoid the scenes with music that enable to easily guess that the sequence comes from Kill Bill); scenes with typical elements (e.g. a spaceship or an alien from which you can easily guess that you are watching a Star Wars movie). The cutters were provided with the instructions to ask themselves the following question when choosing a generic sequence in a movie: "Could I easily guess this sequence belongs to the movie I extracted it from?" – if the answer was yes, so don't select this sequence. In some movies (e.g. Kill Bill 1, Star wars, 2001 A Space Odyssey, etc.), just a few 10-sec sequences of no sequence at all were found to meet these criteria (e.g. The lord of the ring); for these movies, we don't have generic sequences. By contrast, non-generic sequences were namely two 10-sec parts of the movie that you think are very representative of the movie (you can choose any part of the movie without restriction). These two types of sequences should be considered in link with the instruction of the task written above (even if you know that a sequence comes from a movie, just answer if you recognize THE sequence): it is just a supplementary control (qui joue un r le symbolique au taux de FA pour la prise de risque).

In Isola et al., 2014, the authors showed that memorability directly correlated with color. To explain this result, they advanced the hypothesis that the difference in memorability between indoor and outdoor scenes was due to the fact that the latter tend to be less memorable than the former, and that the color tends to be linked with differences in colors, could explain this result. To validate this hypothesis, we manually annotated our 700 sequences as consisting of an indoor or outdoor scene.

2.3 Questionnaire

The targets and fillers that constitute the items of the recognition task were different for each participant. Indeed, the video sequences viewed by a participant were selected based on response to a questionnaire presented in the figure ??? Questionnaire with the 100 movies.

Le nom du film, suivi de l'affiche, du nom du r  alisateur et des acteurs principaux    qui permettent de s  assurer que le film que le participant a vu est bien celui qui lui est pr sent      est pr sent   au participant.

For each movie, the participant had first to answer the following question: "Do you remember watching this movie?" – with the fact that was made clear that the film must have been seen ENTIER.

If he answered "yes" Please let us know how confident you are of having seen this movie (I am not/slightly/50%/considerably/100% confident) Please let us know when did you last see this movie (<1 month; <1 year; <5 years; < 10 years; > 10 years) Please let us know how many times did you watch this movie (once; 2-4 times; 5-9 times; 10-19 times; >20 times)

If he answered "no" Please let us know how confident you are of NOT having seen this movie (I am not/slightly/50%/considerably/100% confident)

The goal of the confidence question was to be sure that we select video sequences from movie, to increase our certitude that target/filler were really targets/fillers.

The other two questions were used for the analysis of the results to assess factors a priori susceptible to influence video memorability. The question about the number of times the movie had been seen for us to evaluate the effect of repetition on memorization, and the question about the last time the movie had been seen to evaluate the effect of time passage on long-term memorability, that continue to be affected in longterm memory

2.4 Participants, facilities and apparatus

105 participants (22 – 58 years of age; $mean = 37.1$; $SD = 10.4$; 26% of them female) employees of [removed for double-blind review] participated in the experiment on a volunteer basis. All participants have either normal or corrected-to-normal visual acuity.

The images were displayed on a 40 inch monitor (TV SONY Bravia ???) with a display resolution of $1,920 \times 1,080$. The participants were comfortably seated at a distance of 150 centimeters from the screen (three times the screen height). The $1,024 \times 768$ images were centered on a black background; at a viewing distance of 150 cm, the stimuli subtended 18.85 degrees of vertical visual angle.

2.5 Procedure

Participants first answered the questionnaire. During the whole time, the experimenter was near the participant in case he doubt about one movie.

Based on the answers to the questionnaires, the algorithm selected the movies associated with the maximum certitude level (i.e. 5) to constitute targets and fillers sets. // At the end, the algorithm selected 80 targets and 40 fillers for the participants, and place them in a random order for presentation // The criteria of selection was the number after the response to the questionnaire the nb of annotations of the sequences, in order to harmonize the nb of annotations per video. The sequences used as targets were selected among the sequences corresponded to the movies the participants were sure (i.e. rate of 9 in the scale) they have already seen. The algorithm selected the less annotated sequences to try maintaining a number of annotation equivalent between sequences. This was to maximize the probability that targets really are targets, and fillers really are fillers, according to they number of annotations.

Then the participants was provided with the instructions about how to complete the memory task. In particular, the experimenter

insisted on the fact that we did not want the participant guess that the sequence was il the film they have already seen, but that we want they recognize.

// During about 24 minutes, you will see a series of 10-seconds video scenes. Press the SPACEBAR during the 10-seconds video scene if you RECOGNIZE/REMEMBER seeing it before. Be careful!! You only press the spacebar when you RECOGNIZE/REMEMBER the video scene, not the movie from which the video scene had been extracted. We do not want you guess that the video scene was in a movie you have seen: just press the spacebar if you remember that you saw a particular video scene.

Then the experimental phase was then launched, corresponding to the memory task.

The total time of the experiment was about 50 minutes.

3 STUDY OF THE MEMORABILITY SCORES

In this section, we describe in detail the ground truth data collected thanks to the protocol described in the previous section. We conclude by calculating the memorability scores finally used to feed our prediction models.

3.1

3.1.1 Sequences' memorabilities. A sequence's memorability is defined as the percentage of correct detections by participants. On average, each video was viewed as target by 9.3 participants. Average sequence memorability was 48.25% (SD of 28.21%). On average, each video was viewed as filler by 9.2 participants. Average false alarm rate was 5.95% (SD of 14.47%).

3.2 Number of annotations per sequence

The histogram presented in the figure 1 corresponds to the number of sequences for each possible number of annotations.

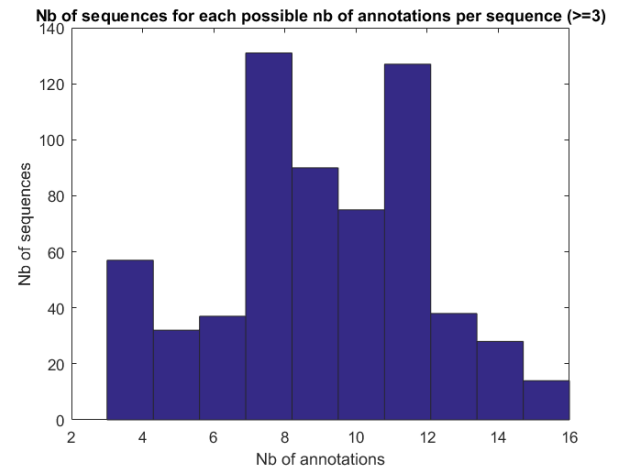


Figure 1: Nb of sequences for each possible nb of annotations per sequence (≥ 3).

3.3 Consistency analysis

We implemented the method proposed in [3] to measure the human consistency. It answers the question: "Are the videos that are more memorable (or forgettable) for a group of observers also more likely to be remembered (or forgotten) by a different group of observers?" The figure 2 presents the curve of the human consistency for each video sequence. It answers the question: "How many annotations for each do we need for the human consistency stop varying?"

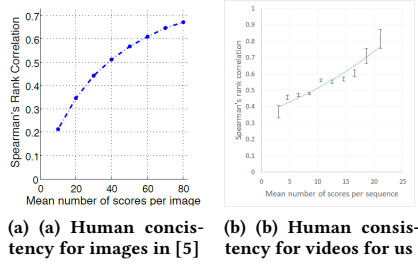


Figure 2: Human consistency averaged over 25 random splits (right) obtained using the method proposed by [3] for sequences with at least 5 annotations (with standard error).

Individual and contextual differences, besides mandatory random variability, explain the $1 - .65$ part of the memorability that is not universally derivable from the intrinsic informations of the videos.

If we compare the curve obtained by [5] (2(a)) and the curve we obtained (b), we can note that we attain human consistency far more quickly than Khosla *et al.*. At least three arguments go in the way of this results. First, we work with videos and they work with images; maybe video memorability is more universal than image memorability because of the higher length of this kind of stimulus (but we can also note that the maximum human consistency is close for images and videos!!). Second, we obtain a measure of real long term memorability, and not a memorability measured some minutes after encoding step; maybe this measure more representative of what is really memorable. Finally, we can advance the fact that in-lab experiment enable to obtain better memorability measure than crowdsourcing one, resulting to attain maximum human consistency earlier.

3.4 Generic vs. Typical sequences

3.5 Quality of the movies

dvdrip vs. HD 720-1080 => Difference of memorability? => Interesting

3.6 Response time

Previous authors have integrated response time in their score of memorability. We pose here the following question: "Is the degree of memorability related to response time?" to answer the question: "Should we exploit response time to correct memorability score?" We hypothesized that the most memorable the sequences, the faster the participant will answer (no instructions to answer faster: just to answer during the 10-sec display).

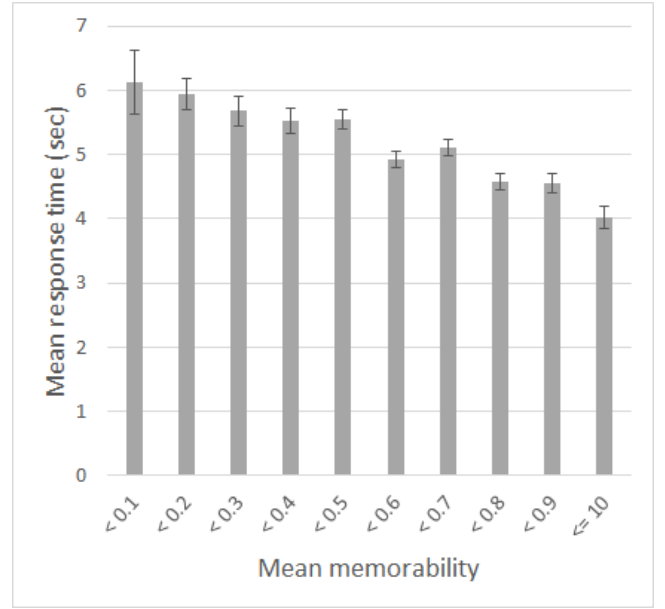


Figure 3: Mean response time for each memorability degree (error bars correspond to SEM).

We observed a Person's correlation of $-0.36 (p < .0001)$ between the response time on the target and their memorability scores. This is consistent with the figure 3. In other words, the higher the memorability, the lower the response time tends to be, suggesting that a sequence most memorable is also a sequence for which we rapidly detect that it is memorable. This is consistent with what we expected: empirically, we observe during passation that participant tend to answer as soon as they recognized the sequence.

It is interesting to link it to the response time for answer on targets (i.e. correct detections) vs. on fillers (i.e. false alarms). The global mean response time was 4.87sec on targets and 5.96sec on fillers. A Student t-test for different sample size show a significant difference ($t(2836) = -5.34, p < .0001$). This means that the participants globally answered more rapidly for targets (i.e. correct detections) than for fillers (i.e. false alarm). One explanation is that participants hesitated more for fillers they answered on, increasing their response time.

It tells us something on memory, that the memories are not evident but more blurred. This is an important point, that explained why, in our discussion, we propose to turn to a "totally objective" measure of memory to constitute a dataset for video memorability prediction: if the participants hesitated, so they choice to answer or not was probably impacted by the way they answered the instructions and their own sensibility to the risk.

3.7 Evolution of the memorability along time

3.8 New manner to compute memorability scores: take into account time and FA

3.8.1 Participants' performance. The average memory performance was the following: the average percentage of correct detection was 48.2% (SD of 14.1%) and the average percentage of false alarms was 4.78% (SD of 5.63%).

3.9 Logistic regression vs. SVM to personalize prediction model

3.10 Film genre and IMDB ratings

3.11 Context

3.12 Features linked to memorability

3.13 Indoor vs outdoor scenes

3.14 Memorability score calculation

According to what was presented before in this section, this is how we finally decided to compute our memorability scores...

4 MEMORABILITY PREDICTION

5 DISCUSSION

=> Nulle part les auteurs ne parlent de mémoire à long terme, à court terme => Ils appréhendent un problème par nature interdisciplinaire d'un point de vue purement computer science => Or la qualité des données est ce qui détermine ce que les modèles vont finalement prédire.

6 CONCLUSIONS

APPENDIX

ACKNOWLEDGMENT

REFERENCES

- [1] Yoann Baveye, Romain Cohendet, Matthieu Perreira Da Silva, and Patrick Le Callet. 2016. Deep Learning for Image Memorability Prediction: the Emotional Bias. In *Proceedings of the 2016 ACM on Multimedia Conference*. ACM, 491–495.
- [2] Junwei Han, Changyuan Chen, Ling Shao, Xintao Hu, Jungong Han, and Tianming Liu. 2015. Learning computational models of video memorability from fMRI brain imaging. *IEEE transactions on cybernetics* 45, 8 (2015), 1692–1703.
- [3] Phillip Isola, Jianxiong Xiao, Devi Parikh, Antonio Torralba, and Aude Oliva. 2014. What makes a photograph memorable? *IEEE Transactions on Pattern Analysis and Machine Intelligence* 36, 7 (2014), 1469–1482.
- [4] Phillip Isola, Jianxiong Xiao, Antonio Torralba, and Aude Oliva. 2011. What makes an image memorable?. In *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*. IEEE, 145–152.
- [5] Aditya Khosla, Akhil S Raju, Antonio Torralba, and Aude Oliva. 2015. Understanding and predicting image memorability at a large scale. In *Proceedings of the IEEE International Conference on Computer Vision*. 2390–2398.
- [6] Sumit Shekhar, Dhruv Singal, Harvinder Singh, Manav Kedia, and Akhil Shetty. 2017. Show and Recall: Learning What Makes Videos Memorable. *arXiv preprint arXiv:1707.05357* (2017).
- [7] Hammad Squalli-Houssaini, Ngoc Duong, Marquant Gwenaëlle, and Claire-Hélène Demarty. 2017. Deep learning for predicting image memorability. (2017).