

What makes an object memorable?

Anonymous ICCV submission

Paper ID ****

Abstract

Recent work by Isola et. al. (2011) has demonstrated that memorability is an intrinsic property of images that is consistent across viewers and can be predicted accurately with current computer vision techniques. Despite progress, a clear understanding of the specific components of an image that drive memorability are still unknown. While previous studies such as Khosla et. al. (2012) have tried to investigate computationally the memorability of image regions within individual images, no behavioral study has systematically explored which memorability of image regions. Here we study which region from an image is memorable or forgettable. Using a large image database, we obtained the memorability scores of the different visual regions present in every image. In our task, participants viewed a series of images, each of which were displayed for 1.4 seconds. After the sequence was complete, participants similarly viewed a series of image regions and were asked to indicate whether each region was seen in the earlier sequence of full images.

1. Introduction

Consider the image and its corresponding objects in Figure 1. Even though the person on the right is comparable in size to the left person, he is remembered far less by humans (indicated by their memorability scores of 0.18 and 0.64 respectively). People tend to remember the fish in the center and the person on the left, even after 30 minutes have passed (memorability score = 0.64). Interestingly, despite vibrant colors and considerable size, the boat is also remembered far less by humans (memorability = 0.18).

Just like aesthetics, interestingness, and other metrics of image importance, memorability quantifies something about the utility of a photograph toward our everyday lives. For many practical tasks, memorability is an especially desirable property to maximize. For example, this may be the case when creating educational materials, logos, advertisements, book covers, websites, and much more. Understanding memorability, and being able to automatically predict it,



Figure 1: Memorability of different objects. Memorability scores of objects for the image in the top row obtained from our psychophysics experiment.

lends itself to a wide variety of applications in each of these areas. **rewrite this to draw attention of reviewer to importance of image memorability.** Due to this, automatic prediction of intrinsic memorability of images using computer vision and machine learning techniques has received considerable attention in the recent years [4], [5], [3], [1], [6]. While these studies have shed light on what distinguishes the memorability of different images and the intrinsic and extrinsic properties that make those images memorable, the above example raises an interesting question: what exactly about an image is remembered? Despite progress in the computer vision literature on image memorability, a clear understanding of the memorability of the specific components of an image is still unknown. For example, not all objects in an image will be equally remembered by people and as the figure 1 seems to suggest, there exists significant and interesting differences in memorability of objects in an image. Furthermore, the memorability of complex images may be principally driven by the memorability of its objects. Can specific objects inside images be memorable to

all us and how can we better understand what makes those objects more memorable?

In this paper, we systematically explore the memorability of objects within individual images and shed light on the various factors and properties that drive object memorability by augmenting both the images and object segmentations in the 850 existing images from PASCAL 2010 [2] dataset with memorability scores and class labels. By exploring the connection between object memorability, saliency, and image memorability, our paper makes several important contributions.

Firstly, we show that just like image memorability, object memorability is a property that is shared across subjects and objects remembered by one person are also likely to be remembered by others and vice versa. Secondly, we show that there exists a strong correlation between visual saliency and object memorability and demonstrate insights when can visual saliency directly predict object memorability and when does it fail to do so. While there have been a few studies that explore the connection between image memorability and visual saliency [1], [8], our work is the first to explore the connection between object memorability and visual saliency. Third, we explore the connection between image memorability and object memorability and show that the most memorable object inside an image can be a strong predictor of image memorability in certain cases. Studying these questions, help not only understand visual saliency, image and object memorability in more detail, but it can also have important contributions to computer vision. For example, understanding which regions and objects in an image are memorable would enable us to modify the memorability of images which can have applications in advertising, user interface design etc. With this in mind, as shown in the section 4, our proposed dataset serves as a benchmark for evaluating object memorability model algorithms and can help usher in future algorithms that try to predict memorability maps.

1.1. Related works

Image Memorability: Describe Isola’s first paper n some insights that have been raised on image memorability thus far. Also describe Khosla’s comp model but we are the first work to actually describe what humans actually remember and don’t

Visual Saliency: Talk about visual attention and models that have been proposed. Also, talk about Pascal-S and how it has helped reduce dataset bias

Saliency and memorability: discuss some results related to saliency and image memorability.

and talk about our work plans on connecting and shedding light on all these phenomena together.

2. Measuring Object Memorability

As a first step towards understanding memorability of objects, we built an image database containing a variety of objects from a diverse range of categories, and measured the probability that every object in each image will be remembered by a large group of subjects after a single viewing. This helps provide ground truth memorability scores for the objects inside the images and allows for a precise analysis of the memorable elements within an image. For this task, we utilized the PASCAL-S dataset [7], a fully segmented dataset built on the validation set of the PASCAL VOC 2010 [2] segmentation challenge. For improved segmentation purposes, we manually cleaned up and refined the segmentations from this dataset. While building the improved ground-truth of full segmentation, we removed all homogeneous non-object or background segments such as ground, grass, floor, sky etc, as well as imperceptible object fragments and excessively blurred regions. All remaining object segmentations were tested for memorability. In the end, our final dataset consisted of 850 images and 3414 object segmentations i.e. on average each image consisted of approximately 4 segments for which we gathered the ground truth memorability on.

2.1. Memory Game: Measuring Object Memorability

To measure the memorability of individual objects from our dataset, we created an alternate version of the Visual Memory Game following the basic design in [4], with the exception of a few key differences. We administered the game and collected data through Amazon Mechanical Turk. In our game, participants first viewed a sequence of images one at a time, with a 1.5 second gap in between image presentations. Subjects were asked to remember the contents and objects inside those images as much as they could. To ensure that subjects would not just only look at the salient or center objects, subjects had unlimited time to freely view the images. Once they were done viewing an image, they could press any key to advance to the next image. Following the initial image sequence, participants then viewed a sequence of objects, their task then being to indicate through a key press which of those objects was present in one of the previously shown images. Each object was displayed for 1.5 second, with a 1.5 second gap in between the object sequences. Pairs of corresponding image and object sequences were broken up into 10 blocks. Each block consisted of 80 total stimuli (35 images and 45 objects), and lasted approximately 3 minutes. At the end of each block, the subject could take a short break. Overall, the experiment took approximately took 30 minutes to complete. A total of 2000 workers from Mechanical Turk (¿ 95% approval rate in Amazons system) performed the game.

Unknown to the subjects, each sequence of images was

pseudo-random and consisted of 3 'target' images taken from the Pascal-S dataset whose objects the participants were to later identify. The remaining images in the sequence consisted of 16 'filler' images and 16 'familiar' images. The 'filler' images were randomly selected from the DUT-OMRON dataset [9] and the 'familiar' objects were taken from the MSRA dataset proposed in . Successive target images were spaced 70-79 images apart, and familiar images were spaced 35-79 image apart. The corresponding sequence of objects contained 3 segmented target objects, one for each of the target images shown previously, along with 16 filler objects, 16 familiar objects, and 10 control objects. Target objects appeared only once, and participants were tested on only one object from each target image. Each target object was spaced X-X objects apart. Objects were centered within their parent frame and non-object pixels were set to grey. Participants were required to complete the entire task, which lasted approximately 30 minutes, and could not participate in the experiment a second time. A total of 2000 workers completed the task, meaning around 20 participants were given the opportunity to remember each object.

In order to measure image memorability, we presented workers on Amazon Mechanical Turk with a Visual Memory Game. In the game, participants viewed a sequence of images, each of which was displayed for 1 second, with a 1.4 second gap in between image presentations (Figure 3). Their task was to press the space bar whenever they saw an identical repeat of an image at any time in the sequence [1] [12]. Participants received feedback whenever they pressed a key (a green symbol shown at the center of the screen for correct detection, and a gray X for an error). Image sequences were broken up into levels that consisted of 120 images each. Each level took 4.8 minutes to perform. At the end of each level, the participant saw his or her correct response average score for that level, and was allowed to take a short break. Participants could complete at most 30 levels, and were able to exit the game at any time. A total of 665 workers from Mechanical Turk (≈ 95 [24]. All images were scaled and cropped about their centers to be 256x256 pixels. The role of the fillers was twofold: first, they provided spacing between the first and second repetition of a target; second, responses on repeated fillers constituted a vigilance task that allowed us to continuously check that participants were attentive to the task [1, 12]. Repeats occurred on the fillers with a spacing of 1-7 images, and on the targets with a spacing of 91-109 images. Each target was sequenced to repeat exactly once, and each filler was presented at most once, unless it was a vigilance task filler, in which case it was sequenced to repeat exactly once. Stringent criteria were used to continuously screen worker performance once they entered the game. First, the game automatically ended whenever a participant fell below a 50 or above a 50 a partic-

ipant failed any of the vigilance criteria, they were flagged. After receiving three such flags they were blocked from further participation in the experiment. Otherwise, participants were able to restart the game as many times as they wished until completing the max 30 levels. Upon each restart, the sequence was reset so that the participant would never see an image they had seen in a previous session. Finally, a qualification and training demo preceeded the actual memory game levels. After collecting the data, we assigned a memorability score to each image, defined as the percentage of correct detections by participants. On average, each image was scored by 78 participants. The average memorability score was 67.5

To measure the memorability of individual objects from our dataset, we created an alternate version of the Visual Memory Game following the basic design of [Oliva, CVPR 2011], with the exception of a few key differences. We administered the game and collected data through Amazon Mechanical Turk. In our game, participants first viewed a sequence of images one at a time. Following the initial sequence, participants then viewed a sequence of objects, their task then being to indicate through a key press which of those objects was present in one of the previously shown images. Pairs of corresponding image and object sequences were broken up into 10 blocks. Each block consisted of 80 total images, and lasted approximately 3 minutes. Each sequence of images was pseudo-random and contained of 3 target images which contained a number of objects that participants were to later identify. The remaining images in the sequence consisted of 16 filler images and 16 familiar images. Successive target images were spaced 70-79 images apart, and familiar images were spaced 35-79 image apart. The corresponding sequence of objects contained 3 segmented target objects, one for each of the target images shown previously, along with 16 filler objects, 16 familiar objects, and 10 control objects. Target objects appeared only once, and participants were tested on only one object from each target image. Each target object was spaced X-X objects apart. Objects were centered within their parent frame and non-object pixels were set to grey. Participants were required to complete the entire task, which lasted approximately 30 minutes, and could not participate in the experiment a second time. A total of 2000 workers completed the task, meaning around 20 participants were given the opportunity to remember each object.

2.2. Is object memorability a shared property across subjects?

Previous work on image memorability has found human consistency to be fairly high. That is, people tend to remember the same images, and exhibit similar performance in doing so. Despite variability due to individual differences and other sources of noise, this level of consistency

provides evidence that memorability is an intrinsic property of images that can be predicted. In contrast to full images, this paper focuses primarily on the memorability of individual objects in an image, which may or may not exhibit the same level of human consistency as full images, which often contain complex arrangements of several objects. High consistency in object memorability would indicate that, like full images, objects can potentially be predicted with high accuracy. To assess human consistency in remembering objects, we repeatedly divided our entire subject pool into two equal halves and quantified the degree to which memorability scores for the two sets of subjects were in agreement using Spearmans rank correlation (). We computed the average correlation over 25 of these random split iterations, yielding a final value of 0.76. Such a result confirms that human consistency in remembering particular objects is at least as strong as that of images.

References

- [1] Z. Bylinskii, P. Isola, C. Bainbridge, A. Torralba, and A. Oliva. Intrinsic and extrinsic effects on image memorability. *Vision research*, 2015. 1, 2
- [2] M. Everingham and J. Winn. The pascal visual object classes challenge 2010 (voc2010) development kit, 2010. 2
- [3] P. Isola, J. Xiao, D. Parikh, A. Torralba, and A. Oliva. What makes a photograph memorable? *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 36(7):1469–1482, 2014. 1
- [4] P. Isola, J. Xiao, A. Torralba, and A. Oliva. What makes an image memorable? In *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*, pages 145–152. IEEE, 2011. 1, 2
- [5] A. Khosla, J. Xiao, A. Torralba, and A. Oliva. Memorability of image regions. In *Advances in Neural Information Processing Systems*, pages 305–313, 2012. 1
- [6] J. Kim, S. Yoon, and V. Pavlovic. Relative spatial features for image memorability. In *Proceedings of the 21st ACM international conference on Multimedia*, pages 761–764. ACM, 2013. 1
- [7] Y. Li, X. Hou, C. Koch, J. M. Rehg, and A. L. Yuille. The secrets of salient object segmentation. In *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*, pages 280–287. IEEE, 2014. 2
- [8] M. Mancas and O. Le Meur. Memorability of natural scenes: the role of attention. In *ICIP*, 2013. 2
- [9] C. Yang, L. Zhang, H. Lu, X. Ruan, and M.-H. Yang. Saliency detection via graph-based manifold ranking. In *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*, pages 3166–3173. IEEE, 2013. 3