

MemCat: A new category-based image set quantified on memorability

Lore Goetschalckx^{Corresp., 1}, Johan Wagemans¹

¹ Brain & Cognition, Katholieke Universiteit Leuven, Leuven, Belgium

Corresponding Author: Lore Goetschalckx
Email address: lore.goetschalckx@kuleuven.be

Images differ in their memorability in consistent ways across observers. What makes an image memorable is not fully understood to date. Most of the current insight is in terms of high-level semantic aspects, related to the content. However, research still shows consistent differences within semantic categories, suggesting a role for factors at other levels of processing in the visual hierarchy. To aid investigations into this role as well as contributions to the understanding of image memorability more generally, we present MemCat. MemCat is a category-based image set, consisting of 10K images representing five broader, memorability-relevant categories (animal, food, landscape, sports, and vehicle) and further divided into subcategories (e.g., bear). They were sampled from existing source image sets that offer bounding box annotations or more detailed segmentation masks. We collected memorability scores for all 10K images, each score based on the responses of on average 99 participants in a repeat-detection memory task. Replicating previous research, the collected memorability scores show high levels of consistency across observers. Currently, MemCat is the second largest memorability image set and the largest offering a category-based structure. MemCat can be used to study the factors underlying the variability in image memorability, including the variability within semantic categories. In addition, it offers a new benchmark dataset for the automatic prediction of memorability scores (e.g., with convolutional neural networks). Finally, MemCat allows to study neural and behavioral correlates of memorability while controlling for semantic category.

MemCat: A new category-based image set quantified on memorability

Lore Goetschalckx¹, Johan Wagemans¹

¹Brain & Cognition, KU Leuven, Leuven, Vlaams-Brabant, Belgium

Corresponding author:

Lore Goetschalckx¹

Tiensestraat 102, Leuven, Vlaams Brabant, 3000, Belgium

Email address: lore.goetschalckx@kuleuven.be

Abstract

Images differ in their memorability in consistent ways across observers. What makes an image memorable is not fully understood to date. Most of the current insight is in terms of high-level semantic aspects, related to the content. However, research still shows consistent differences within semantic categories, suggesting a role for factors at other levels of processing in the visual hierarchy. To aid investigations into this role as well as contributions to the understanding of image memorability more generally, we present MemCat. MemCat is a category-based image set, consisting of 10K images representing five broader, memorability-relevant categories (animal, food, landscape, sports, and vehicle) and further divided into subcategories (e.g., bear). They were sampled from existing source image sets that offer bounding box annotations or more detailed segmentation masks. We collected memorability scores for all 10K images, each score based on the responses of on average 99 participants in a repeat-detection memory task. Replicating previous research, the collected memorability scores show high levels of consistency across observers. Currently, MemCat is the second largest memorability image set and the largest offering a category-based structure. MemCat can be used to study the factors underlying the variability in image memorability, including the variability within semantic categories. In addition, it offers a new benchmark dataset for the automatic prediction of memorability scores (e.g., with convolutional neural networks). Finally, MemCat allows to study neural and behavioral correlates of memorability while controlling for semantic category.

Keywords: image memorability, data set, memorability scores, categories, ground truth, visual memory, recognition memory

Introduction

Photography is booming. Estimates say that over 1.2 trillion images have been captured in 2017 (Lee, 2018). That is more than 38K images per second! With a great variability of images out there, a fair amount of research has been devoted to quantifying and understanding how images differ in meaningful ways. Among the studied image properties are, for example, aesthetic appreciation (e.g., Kong, Shen, Lin, Mech, & Fowlkes, 2016), popularity on social media (e.g., McParlane, Moshfeghi, & Jose, 2014), amount of visual clutter (e.g., Rosenholtz, Li, & Nakano, 2007), evoked fixations (e.g., Judd, Ehinger, Durand, & Torralba, 2009), etc.

It has been suggested that images also differ systematically in the likelihood of being remembered and recognized later, a property that is referred to as image memorability (Isola, Xiao, Parikh, Torralba, & Oliva, 2014). Using a repeat-detection memory task, in which participants watch a sequence of images and respond whenever they see a repeat of a previously shown image, Isola et al. (2014) assigned 2222 scene images a memorability score based on the proportion of participants recognizing the image upon its repeat. They found that memorability rank scores were highly consistent across participants. In other words, there was a lot of agreement as to which images were remembered and recognized, and which ones were easily forgotten. This suggests that memorability can indeed be considered an intrinsic image property and that whether you will remember a certain image does not only depend on you as an individual, but also on the image itself. The result has also been replicated with a more traditional long-term visual memory task, with a separate study and test phase (Goetschalckx, Moors, & Wagemans, 2018). Moreover, image memorability rankings have been shown to be stable across time (Goetschalckx, Moors, & Wagemans, 2018; Isola et al., 2014), across contexts (Bylinskii, Isola, Bainbridge, Torralba, & Oliva, 2015), and across encoding types (intentional versus incidental; Goetschalckx, Moors, & Wagemans, in press). Finally, while they might be related to some extent, image memorability does not simply boil down to popularity (Khosla, Raju, Torralba, & Oliva, 2015), aesthetics (Isola et al., 2014; Khosla et al., 2015), interestingness (Isola et al., 2014), or the ability of an image to capture attention (Bainbridge, 2017).

The findings spurred new research aimed at understanding and predicting memorability. When it comes to merely predicting, the best results so far are achieved using convolutional neural networks (CNNs; e.g., Khosla et al., 2015). However, when it comes to truly understanding, CNNs have often been critiqued to be black boxes. Nonetheless, Khosla et al.'s (2015) analyses

provided some further insight, mostly pointing at differences between broader image categories and content types. For examples, units in the network with the highest correlation with memorability seemed to respond mostly to humans, faces, and objects, while those with the lowest correlation seemed to respond to landscapes and open scenes. Furthermore, the most memorable regions of an image, according to the CNN, often capture people, animals or text. These findings are in line with earlier work (Isola et al., 2014), which showed that the predictive performance of a model trained on mere object statistics (e.g., number of objects) was boosted considerably when the object labels were taken into account (Isola et al., 2014). In addition, a model trained on the overall scene labels alone, already predicted memorability scores with a Spearman's rank correlation of .37 to the ground truth. Memorable images often had labels referring to people, interiors, foregrounds, and human-scaled objects, while labels referring to exteriors, wide-angle vistas, backgrounds, and natural scenes were associated with low image memorability scores. While differences in image categories (or labels) seem to play a large role, they do not explain all the variability in memorability. Interestingly, consistent differences in memorability scores remain even *within* image categories (Bylinskii et al., 2015). Goetschalckx, Moors, Vanmarcke, and Wagemans (2018) for example, have argued that part of that variability might be due to differences in how well the image is organized. Nonetheless, the concept of memorability and its correlates is not yet fully understood to date and further research is required to paint a clearer picture.

The current work presents a novel dataset of images quantified on memorability, that can be used in research to achieve this goal (example images in Figure 1). To the best of our knowledge, there were previously five large image sets with memorability scores, three consisting of regular photographs, which is also the focus here: Isola et al. (2014), FIGRIM (Bylinskii et al., 2015), and LaMem (Khosla et al., 2015), and two more specialized sets, which we will not further discuss here: Bainbridge, Isola, and Oliva (2013; face images), and Borkin et al. (2013; data visualizations). For completeness, we also mention a smaller set (850 images) that was used to study which objects in an image are memorable (Dubey, Peterson, Khosla, Yang, & Ghanem, 2015). Table 1 compares MemCat to the other large datasets. The comparison is discussed in more detail below.

A first feature of the current dataset is that it was designed to be representative for five different broad natural categories and to allow to study memorability differences within semantic

categories. The set is characterized by a hierarchy of five broader, but memorability relevant categories, further divided into more fine-grained subcategories. Only the FIGRIM set also offers a category structure, but the number of exemplar images per category was lower: 59-157, compared to 2,000 in the current set. We opted for broad categories to ensure that the whole was still varied and representative enough, while containing a large number of exemplar images per category at the same time. The chosen categories: animal, food, landscape, sports, and vehicle, are memorability relevant in the sense that previous research has shown that they generally differ in memorability. For example, knowing that the presence of people in an image is predictive for memorability (see above), we chose one designated people category, the sports category, and avoided including images with people in other categories. We included an animal category as a non-human animate category, food and vehicle as more object-based categories, and landscape to represent the wide exteriors that are often associated with lower memorability scores.

Second, we aimed for a large set, such that it would be suitable for machine learning approaches. With a total of 10,000 images quantified on memorability, the current set is the second largest memorability dataset, after LaMem.

Third, we sampled images from existing datasets, such that the image annotations collected there, would also be available for researchers studying memorability. In particular, we searched for images annotated with segmentation masks or at least bounding boxes, reasoning that they may hold some indications of how the image is organized (e.g., where is the subject located), which might be of particular interest when studying memorability within categories and factors other than semantics.

In summary, the unique combination of features of MemCat, together with its richness in data, make it a valuable addition to the memorability literature. We review possible uses of this set in the Discussion.

Materials & Methods

Participants. There were 249 undergraduate psychology students (KU Leuven) who participated in this study in exchange for course credits (216 female, 32 male, 1 other). Four students did not disclose their age and the remainder were aged between 18 and 27 years old ($M = 19.24$, $SD = 0.94$). The majority of the participants, however, were recruited through Amazon's Mechanical Turk (AMT) and received a monetary compensation (see further for details). The settings on

AMT were chosen such that only workers who indicated to be at least 18 years old and living in the USA could participate. Further eligibility criteria were that the worker had to have an approval rate of at least 95% on previous human intelligence tasks (HITs) and a total number of previously approved HITs of at least 100. A total of 2162 AMT-workers participated in this study (1139 female, 917 male, 4 other, and 102 who did not disclose this information). For the 1851 workers who disclosed their age, the reported ages ranged between 18 and 82 years old ($M = 37.14$, $SD = 11.89$). The AMT data collection took place from April 2018 till July 2018. Data collected through AMT has been shown to come from participant samples that are more diverse than student samples and to be comparable in quality and reliability to those collected in the lab (e.g., Buhrmester, Kwang, & Gosling, 2011).

Materials. MemCat consists of 10,000 images sampled from four previously existing image sets: ImageNet (Deng et al., 2009), COCO (Lin et al., 2014), SUN (Xiao, Hays, Ehinger, Oliva, & Torralba, 2010), and The Open Images Dataset V4 (Kuznetsova et al., 2018). The four source sets were chosen because of their large size (i.e., number of images), the availability of semantic annotations, and the availability of bounding box annotations or more complete segmentation masks for at least a subset of their images. The images selected from the source sets to be included in MemCat belonged to the five broader semantic categories outlined in the Introduction: animal (2,000 images), food (2,000 images), landscape (2,000 images), sports (2,000 images), and vehicle (2,000 images). We explain the different steps in the selection procedure in more detail below.

As a first step, we listed at least 20 subcategories for each broader category. The goal was to obtain 2,000 images per category, without including more than 100 exemplar images per subcategory. This was to ensure a reasonable level of variability and to avoid high levels of false alarms in the memory task (see further). The subcategories were then translated to semantic annotations from the source dataset. For example, for the subcategory “bear” (animal), we used COCO images annotated with a “bear” tag and ImageNet images from nodes “American black bear”, “brown bear”, and “grizzly”. An overview of our hierarchy of categories and subcategories, can be found in Figure 2.

The second step consisted of automatically sampling exemplar images from the listed subcategories, while satisfying a number of shape restrictions. To avoid that images would stand out because of an extreme aspect ratio, we only sampled images with aspect ratios between 1:2

and 2:1. Furthermore, the minimum resolution was set to 62,500 pixels. Finally, to ensure that the images would fit comfortably on most computer monitors, we adopted a maximum height of 500 pixels and a maximum width of 800 pixels. However, for SUN and The Open Images Dataset, only a low number of images satisfied the latter two restrictions (they were often too big), which is why we opted to resize (using Hamming interpolation) images from those two source datasets to meet the restrictions. Apart from the shape restriction, we also restricted the sampling to images for which bounding box annotations or more complete segmentation masks were available from the source datasets. Finally, we sampled more images than the target number (2,000 images per broad category), anticipating exclusions in the next step.

The third step constituted manual selection work, carried out by the first author, assisted by two student-interns. We manually went through the exemplar images sampled in the previous step, and eliminated images following a number of exclusion rules. The exclusion rules can roughly be divided into two kinds. A first kind of exclusion rule touches upon the quality of the image. We excluded images of poor image quality (e.g., very dark, very much overexposed, blurry, etc.), images that did not convincingly belong to the subcategory they were assigned to,¹ images in greyscale or looking like they were the result of another color filter, images that were not real photographs (e.g., drawings, digitally manipulated images, computer generated images), and collages. A second set of rules concerns factors that could affect the memorability of an image, but were not of interest for the purpose of MemCat. One such factor is text. We excluded images containing large, readable text or text not belonging to the image itself (e.g., date of capture). Furthermore, images depicting remarkably odd scenes (e.g., dog wearing Santa Clause costume) were also excluded. Similarly, we avoided images depicting famous places or people, and images of dead, wounded or fighting animals. Finally, there was one designated “people” category, the sports category, and the presence of people in an image was avoided in all other categories (but we allowed anonymous people in the background in the vehicle category or the presence of a hand in images of the food category). In addition to these exclusion rules, we also tried, to the best of our ability, not to include (near) duplicate images. If the target number of images was not obtained after Step 3, we reverted back to Step 2, if there were still images to sample from, or to Step 1 if we needed to include additional subcategories.

Finally, for those categories for which more than the target number of images survived Step 3, there was a fourth step to randomly down-sample the selection to the target number, assigning higher sampling probabilities to images annotated with segmentation masks.

In addition to MemCat, we collected 10,000 filler images, that were not quantified on memorability themselves, but were needed in the memory task used to quantify the other images. The filler images were sampled randomly from The Open Images Dataset, but from a different subset to avoid overlap.² As these images would function only as filler images, there were fewer restrictions. For example, the images could be of any category, they were allowed to contain text, etc. However, the same shape restrictions were still applied.

Procedure. Having carefully collected 10,000 images for MemCat, the next step was to quantify them on memorability. Following previous work, this was achieved by presenting the images in an online repeat-detection memory game (Isola et al., 2014; Khosla et al., 2015), in which participants watch a sequence of images and are asked to respond when they recognize a repeat of a previously shown image. Students participating for course credits played the game in the university's computer labs, hosting about 20 students at a time. AMT workers played the game from the comfort of their homes (or whichever location they preferred). Prior to starting the game, all participants were prompted to read through an informed consent page explaining the aims of the study and their rights as participants. They could give their consent by actively ticking a box. The study was approved by SMEC, the Ethical Committee of the Division of Humanities and Social Sciences, KU Leuven, Belgium (approval number: G-2015 08 298).

For the task design of the game, we closely followed Khosla et al. (2015), as their version of the game was designed to quantify large numbers of images. We divided the game into blocks of 200 trials. On each trial, an image was presented at the center of the browser window for a duration of 600 ms, with an intertrial interval of 800 ms. During this interval, a fixation cross was shown. Sixty-six images were target images, sampled randomly from MemCat, and repeated after 19 to 149 intervening images. Forty-four images were random filler images that were never repeated. Finally, there were 12 additional random filler images that were repeated after 0 to 6 intervening images to keep participants attentive and motivated. They are referred to as vigilance trials. Participants could indicate that they recognized a repeat by pressing the space-bar. They did not receive trial-by-trial feedback, but were shown their hit rate as well as number of false alarms at the end of the block. Figure 3 presents a schematic of the game.

Each block lasted a little less than 5 minutes. Care was taken to ensure that an image was never repeated more than once and never across blocks. Students were asked to complete as many blocks as they could in one hour, with one bigger, collective break of roughly ten minutes after half an hour, and smaller self-timed breaks between the remainder of the blocks. Most students could complete 8 blocks, but for some groups, slow data uploads at the end of a block resulted in lower numbers. AMT workers could complete one to 16 blocks, were allotted 48h to submit their completed blocks (so, they were allowed to spread the blocks over time), and were paid \$0.40 per block. To ensure a good quality of the AMT data and to avoid random or disingenuous responses, AMT workers were blocked from playing anymore blocks after two with a d' prime lower than 1.5 on the vigilance trials. They were warned the first time this happened.

Memorability Measures. Following Khosla et al. (2015), we deduced two image memorability measures from the data collected through the repeat-detection memory game. One measure was calculated as the proportion of participants recognizing the image when shown to them for the second time (i.e., the hit rate across participants), as was also done in Isola et al. (2014). The other measure was based on the same principle, but penalized for false alarms (i.e., when participants press the space-bar for the first presentation of the image) in the way proposed in Bainbridge et al. (2013). Rather than H/N_{resp} (first measure), the formula was the following: $(H-F)/N_{\text{resp}}$, where H is the number of participants recognizing the image, F is the number of participants making a false alarm when the image is presented for the first time, and N_{resp} is the total number of participants having been presented with the image. On average, an image's N_{resp} was 99 (after exclusions). Note that the memorability scores have an upper bound of 1 and a lower bound of 0. In theory, $(H-F)/N_{\text{resp}}$ could result in a negative score, but in practice it is highly unlikely that there would be more participants making a false alarm for the image than there are participants making a hit.

Results

Participant Performance. As mentioned, the performance on the easier vigilance trials was taken as an indication of whether participants were playing the memory game in a genuine way. If in a certain block, a participant did not distinguish vigilance repeats from non-repeat trials with a d' prime of at least 1.5 (preset performance threshold), that block was excluded from further

analyses. The exclusion rate amounted to 3% of all played blocks. Recall, however, that AMT workers were not allowed to play more blocks after two excluded ones. After exclusion, the mean dprime across participants was 2.77 ($SD = 0.56$) for the vigilance repeats, and 2.47 ($SD = 0.50$) for the target repeats. Table 2 summarizes participants' overall performance, collapsing over vigilance and target repeats. Participants generally performed well on the task.

Memorability Scores. Participants' high performance was also reflected in the average image memorability scores. Figure 4 displays the mean for each of the two memorability measures as a horizontal line ($M_{H/N_{resp}} = .76$, SD ; $M_{(H-F)/N_{resp}} = .70$). It is comparable to the mean observed in Khosla et al. (2015). In addition, Figure 4 visualizes the distribution of the collected image memorability scores for each of the five broad main categories separately. A simple linear regression revealed that the category explains 43% of the variance in the H/N_{resp} scores and 44% of the variance in the $(H-F)/N_{resp}$. In line with previous research, the landscape images were on average the least memorable ($M_{H/N_{resp}} = .60$; $M_{(H-F)/N_{resp}} = .53$). They were followed by the vehicle images ($M_{H/N_{resp}} = .76$; $M_{(H-F)/N_{resp}} = .70$). Somewhat surprisingly, the food images generally came out on top of the ranking ($M_{H/N_{resp}} = .85$; $M_{(H-F)/N_{resp}} = .80$), topping the animal $M_{H/N_{resp}} = .80$; $M_{(H-F)/N_{resp}} = .73$) and sports $M_{H/N_{resp}} = .78$; $M_{(H-F)/N_{resp}} = .71$) categories. However, there is still a large degree of variability that is not explained by differences in broad image categories. Indeed, memorability varied considerably within categories as well, with SD s of: .09 (animal; $SD_{(H-F)/N_{resp}} = .09$), .08 (food; $SD_{(H-F)/N_{resp}} = .08$), .13 (landscape; $SD_{(H-F)/N_{resp}} = .14$), .09 (sports; $SD_{(H-F)/N_{resp}} = .10$), and .09 (vehicle; $SD_{(H-F)/N_{resp}} = .09$).

Having observed that images from the same broader category indeed still differed in memorability, the next question was whether these differences are consistent across participants. This question taps into the reliability of the memorability measures. Following previous memorability work (e.g., Isola et al., 2014), the consistency was assessed by randomly splitting the participant pool in half, computing the memorability scores for each half separately and determining the Spearman's rank correlation between the two sets of scores. This was repeated for 1,000 splits and the Spearman's rank correlation was averaged across the splits. Figure 5 shows the results in function of the mean N_{resp} for each category as well as for the total image set. We first discuss the results for H/N_{resp} (see Figure 5, left panel). When collapsing over all five categories, the observed mean split-half Spearman's rank correlation with all available responses

($N_{\text{resp}} = 99$, on average) amounted to .78. In comparison, Khosla et al. (2015) reported a mean split-half Spearman's rank correlation of .67 for their LaMem dataset. However, they only collected 80 responses per image. After randomly down-sampling our data to an N_{resp} of 80, we still found a split-half consistency of .73. With the exception of the landscape category, for which we observed a total (i.e., without down-sampling) split-half consistency of .77, the total per category split-half consistency estimates were lower, ranging between .59 and .67. This is possibly due to smaller ranges of memorability scores within those categories (see Figure 4). Note, however, that the split-half consistencies are an underestimate of the reliability of the memorability scores calculated based on the full participant pool. The latter can be estimated from the split-half consistency by means of the Spearman-Brown formula (Brown, 1910; Spearman, 1910). Applying this formula, we found the following final reliabilities for the H/N_{resp} memorability scores : .87 (all), .80 (animal), .75 (food), .87 (landscape), .75 (sports), .78 (vehicle).

For the $(H-F)/N_{\text{resp}}$ memorability scores, we confine the discussion to pointing out that the pattern of results is qualitatively similar, although the final reliabilities are somewhat lower: .86 (all), .74 (animal), .71 (food), .85 (landscape), .77 (sports), .71 (vehicle).

Finally, after finding that the two image memorability measures were both acceptably reliable, we asked how they compared to each other. In the current dataset, they were highly intercorrelated, as evidenced by a Pearson correlation of .93 when collapsing over all five categories. The per category correlations were: .82 (animal), .90 (food), .91 (landscape), .85 (sports), .88 (vehicle).

Discussion

We presented a new dataset, MemCat, consisting of a total of 10,000 images, each quantified on memorability using a repeat-detection memory task (first introduced by Isola et al., 2014, version used here based on Khosla et al., 2015). MemCat is the second largest image memorability dataset available, and the largest that is based on a category structure. That is, it is divided into five broader, memorability relevant semantic categories: animal, food, landscape, sports, and vehicle, each with 2,000 exemplar images, which are further divided into subcategories (e.g., bear, cat, cow). Furthermore, the images were sampled from popular, existing datasets such that additional annotations available there (e.g., segmentations masks or bounding boxes) would also

be available to researchers wishing to use MemCat for research aimed at investigating specific factors underlying memorability.

Replicating previous research, we found that images differ considerably in memorability and that these differences are highly consistent across participants. Part but not all of this variability can be explained by differences between the five broader semantic categories. Note, however, that this result is correlational in nature, and that one should be cautious drawing causal conclusions. In line with Bylinskii et al. (2015), considerable variability in memorability remained even within the categories. However, the consistency there was somewhat lower, probably because the variance was also lower. When the differences between images become smaller, it becomes harder to reliably and consistently distinguish them. Nevertheless, the consistency estimates per category were still high, indicating that we obtained reliable memorability scores. Finally, we reported results for two different methods to compute memorability scores. One is to compute the hit rate across participants: H/N_{resp} . This was the method used in the original work by Isola et al. (2014). However, in principle, it is possible that some images elicit more key presses not because they are truly recognized, but for some other reason (e.g., they seem familiar). That is why Bainbridge et al. (2013) suggested to correct for false alarms (i.e., when participants press the key for the first presentation of an image, when it is not a repeat) by computing $(H-FA)/N_{\text{resp}}$. We report both measures for comparison, but note that they lead to a highly similar pattern of results and are also strongly intercorrelated.

In what follows, we discuss possible uses of MemCat. Most of what we learned from previous studies about what makes an image memorable is specified in terms of semantic categories or content types (e.g., images of people are more memorable than landscapes). However, a considerable amount of variability was still left unexplained. A primary use of the current dataset is in studies aiming to better understand the factors underlying image memorability. In particular, with 2,000 images for each of five broader categories, it allows to zoom in on variability within categories. This variability is of more interest to practical applications (e.g., advertising, education), because the semantic category or the content type (e.g., a certain product) will often be predefined and it will be a matter of choosing or creating a more memorable depiction of it. In addition to dividing the set into broad semantic categories, we also avoided variability due to other factors already discovered in previous studies (e.g., we excluded images depicting oddities,

images containing text or recognizable places or faces), thus creating a set designed to help understand the previously unexplained variability in image memorability. Second, MemCat is also useful as a benchmark for machine learning approaches to automatically predict memorability. Currently, LaMem (Khosla et al., 2015) is most often used, but models can now also be trained and tested on the current dataset. When taking Khosla et al.'s (2015) MemNet-CNN (without retraining), we found that its predictions show a rank correlation of 0.68 with the $(H-F)/N_{\text{resp}}$ memorability scores in the current set, suggesting that there is room for improvement. Given the category structure in MemCat, one could explore, for the first time, memorability models with one or more layers that are specific to a category. Indeed, it is possible that what makes landscape images memorable is different from what makes animal images memorable. Finally, a third possible use is in neuroscientific studies or psychophysical studies examining effects of memorability. The current set offers a large number of quantified images to choose from. Moreover, it facilitates matching memorability conditions (e.g., high versus low) on semantic category, something that is often done in neuroscientific studies (e.g., Bainbridge, Dilks, & Oliva, 2017; Khaligh-Razavi, Bainbridge, Pantazis, & Oliva, 2016).

Usage. On the MemCat project page, <http://gestaltrevision.be/projects/memcat/>, we provide a link to the collection of 10K images as well as links to two data files. One file describes the images that were used and contains columns indicating the image filename in its source dataset, the name of its source dataset, the category (e.g., animal) and subcategory (e.g., bear) we assigned it to, the label that was used to sample it from its source dataset (e.g., American black bear), the current width, the current height, the factor by which it was resized (both the original width and height were multiplied by this factor), the number of hits (H), the number of false alarms (FA), the number of participants it was presented to (N_{resp}), and the two memorability scores. The other file contains the data collected in the repeat-detection memory game. Its columns indicate the participant ID (anonymized), the participant's age, the participant's gender, whether or not they participated through AMT, the block number, the trial number, the image shown, the trial type (target, target repeat, filler, vigilance, vigilance repeat), the participant's response (hit, correct rejection, miss, false alarm), the screen width, and the screen height.

Acknowledgments

The authors would like to thank the student-interns who have assisted in the image selection and data collection: Justine Aeyels and Joran Geeraerts, and Christophe Bossens, who is responsible for the technical support in the lab and contributed greatly to the implementation of the repeat-detection memory task on Amazon’s Mechanical Turk.

373 **Endnotes**

¹ This could happen, for example, with images from the COCO source dataset. COCO images do not come with a single, overall scene label, but instead come with multiple semantic tags describing what is in the image. For this reason, an image annotated with the tag “cat,” for instance, could be more of a living room image that just happens to have cat sleeping somewhere in a corner in the background.

² The source set is presented in three different subsets: train, validation, and test. We sampled from the validation and test subsets for MemCat, and from the train subset for the fillers images used in the memory task.

References

- Bainbridge, W. A. (2017). The resiliency of memorability: A predictor of memory separate from attention and priming. *ArXiv, 1703.07738 [q-bio]*. Retrieved from <http://arxiv.org/abs/1703.07738>
- Bainbridge, W. A., Dilks, D. D., & Oliva, A. (2017). Memorability: A stimulus-driven perceptual neural signature distinctive from memory. *NeuroImage, 149*, 141–152. <https://doi.org/10.1016/j.neuroimage.2017.01.063>
- Bainbridge, W. A., Isola, P., & Oliva, A. (2013). The intrinsic memorability of face photographs. *Journal of Experimental Psychology: General, 142*(4), 1323–1334. <https://doi.org/10.1037/a0033872>
- Borkin, M. A., Vo, A. A., Bylinskii, Z., Isola, P., Sunkavalli, S., Oliva, A., & Pfister, H. (2013). What makes a visualization memorable? *IEEE Transactions on Visualization and Computer Graphics, 19*(12), 2306–2315. <https://doi.org/10.1109/TVCG.2013.234>
- Brown, W. (1910). Some experimental results in the correlation of mental abilities. *British Journal of Psychology, 1904-1920, 3*(3), 296–322. <https://doi.org/10.1111/j.2044-8295.1910.tb00207.x>
- Buhrmester, M., Kwang, T., & Gosling, S. D. (2011). Amazon's Mechanical Turk: A new source of inexpensive, yet high-quality, data? *Perspectives on Psychological Science, 6*(1), 3–5.
- Bylinskii, Z., Isola, P., Bainbridge, C., Torralba, A., & Oliva, A. (2015). Intrinsic and extrinsic effects on image memorability. *Vision Research, 116*(Part B), 165–178. <https://doi.org/10.1016/j.visres.2015.03.005>
- Deng, J., Dong, W., Socher, R., Li, L. J., Li, K., & Fei-Fei, L. (2009). ImageNet: A large-scale hierarchical image database. *IEEE Conference on Computer Vision and Pattern Recognition, 2009. CVPR 2009*, 248–255. <https://doi.org/10.1109/CVPR.2009.5206848>
- Dubey, R., Peterson, J., Khosla, A., Yang, M. H., & Ghanem, B. (2015). What makes an object memorable? *2015 IEEE International Conference on Computer Vision (ICCV)*, 1089–1097. <https://doi.org/10.1109/ICCV.2015.130>
- Goetschalckx, L., Moors, J., & Wagemans, J. (in press). Incidental image memorability. *Memory*.

- Goetschalckx, L., Moors, P., Vanmarcke, S., & Wagemans, J. (2018). Get the picture? Goodness of image organization contributes to image memorability. *PsyArXiv*. <https://doi.org/10.17605/OSF.IO/8ATJ9>
- Goetschalckx, L., Moors, P., & Wagemans, J. (2018). Image memorability across longer time intervals. *Memory*, 26(5), 581–588. <https://doi.org/10.1080/09658211.2017.1383435>
- Isola, P., Xiao, J., Parikh, D., Torralba, A., & Oliva, A. (2014). What makes a photograph memorable? *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(7), 1469–1482. <https://doi.org/10.1109/TPAMI.2013.200>
- Judd, T., Ehinger, K., Durand, F., & Torralba, A. (2009). Learning to predict where humans look. *2009 IEEE 12th International Conference on Computer Vision*, 2106–2113. <https://doi.org/10.1109/ICCV.2009.5459462>
- Khaligh-Razavi, S.-M., Bainbridge, W. A., Pantazis, D., & Oliva, A. (2016). From what we perceive to what we remember: Characterizing representational dynamics of visual memorability. *BioRxiv*, 049700. <https://doi.org/10.1101/049700>
- Khosla, A., Raju, A. S., Torralba, A., & Oliva, A. (2015). Understanding and predicting image memorability at a large scale. *2015 IEEE International Conference on Computer Vision (ICCV)*, 2390–2398. <https://doi.org/10.1109/ICCV.2015.275>
- Kong, S., Shen, X., Lin, Z., Mech, R., & Fowlkes, C. (2016). Photo aesthetics ranking network with attributes and content adaptation. *ArXiv CS*, 1606.01621. Retrieved from <http://arxiv.org/abs/1606.01621>
- Kuznetsova, A., Rom, H., Alldrin, N., Uijlings, J., Krasin, I., Pont-Tuset, J., ... Ferrari, V. (2018). The Open Images Dataset V4: Unified image classification, object detection, and visual relationship detection at scale. *ArXiv:1811.00982*.
- Lee, E. (2018, September 12). Our best photos deserve to be printed [Blog post]. Retrieved from InfoBlog - Insight from InfoTrends website: <http://blog.infotrends.com/our-best-photos-deserve-to-be-printed/#more-24362>
- Lin, T.-Y., Maire, M., Belongie, S., Bourdev, L., Girshick, R., Hays, J., ... Dollár, P. (2014). Microsoft COCO: Common Objects in Context. *ArXiv:1405.0312 [Cs]*. Retrieved from <http://arxiv.org/abs/1405.0312>
- Macmillan, N. A., & Creelman, C. D. (2005). *Detection theory: A user's guide* (2nd ed). Mahwah, NJ: Lawrence Erlbaum Associates.
- McParlane, P. J., Moshfeghi, Y., & Jose, J. M. (2014). “Nobody comes here anymore, It’s too crowded”: Predicting image popularity on Flickr. *Proceedings of International*

- 438 *Conference on Multimedia Retrieval*, 385–391.
439 <https://doi.org/10.1145/2578726.2578776>
- 440 Rosenholtz, R., Li, Y., & Nakano, L. (2007). Measuring visual clutter. *Journal of Vision*,
441 7(2), 17–17. <https://doi.org/10.1167/7.2.17>
- 442 Spearman, C. (1910). Correlation calculated from faulty data. *British Journal of Psychology*,
443 1904-1920, 3(3), 271–295. <https://doi.org/10.1111/j.2044-8295.1910.tb00206.x>
- 444 Xiao, J., Hays, J., Ehinger, K. A., Oliva, A., & Torralba, A. (2010). SUN database: Large-
445 scale scene recognition from abbey to zoo. *2010 IEEE Conference on Computer*
446 *Vision and Pattern Recognition (CVPR)*, 3485–3492.
447 <https://doi.org/10.1109/CVPR.2010.5539970>

Figure 1(on next page)

Example images of MemCat.

The memorability score, calculated as the hit rate across participants (H/N_{resp}) is indicated in the bottom right corner.



Figure 2(on next page)

Category hierarchy of MemCat.

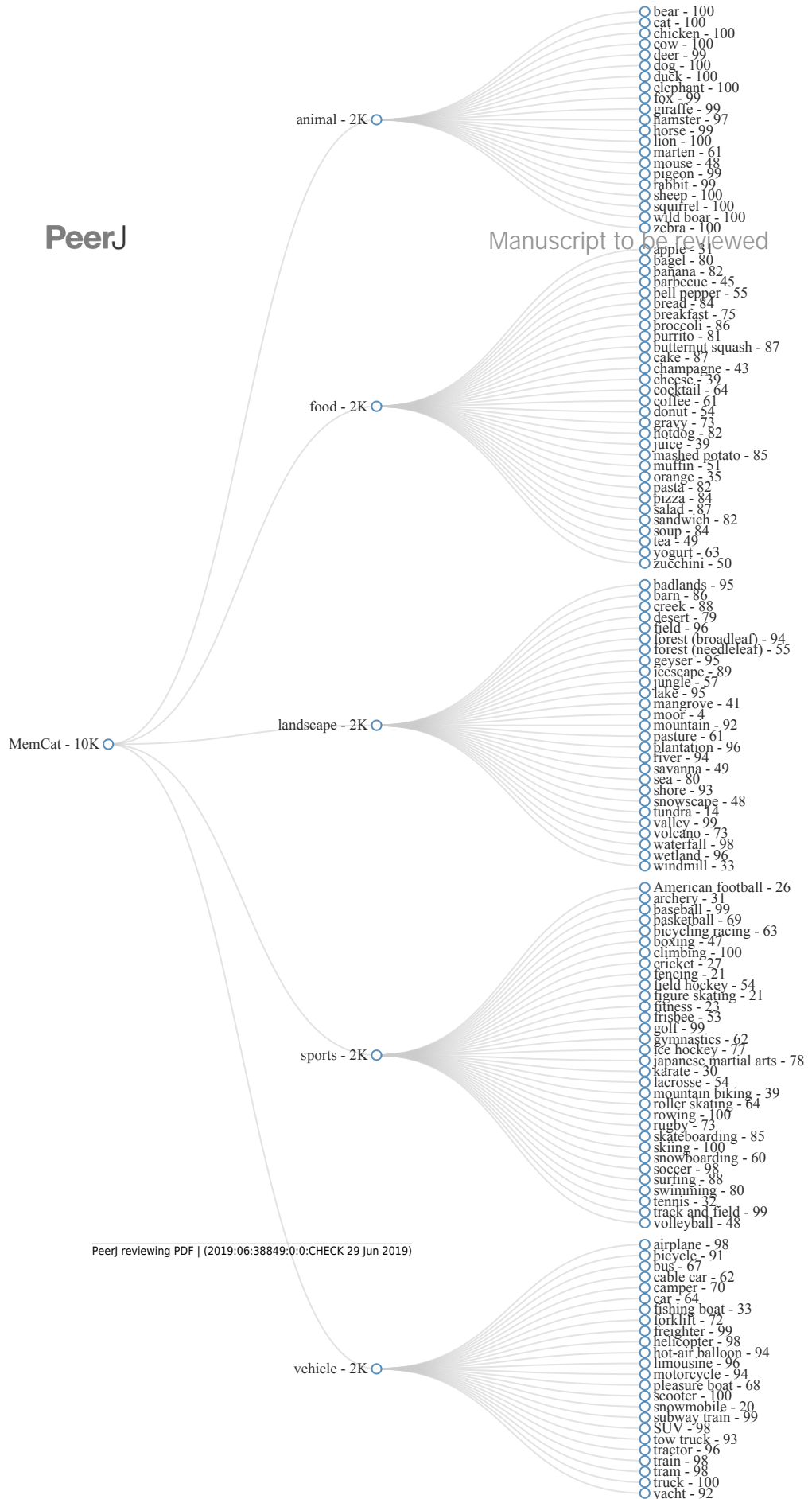


Figure 3(on next page)

Schematic of our implementation of the repeat-detection memory game first introduced by Isola et al. (2014).

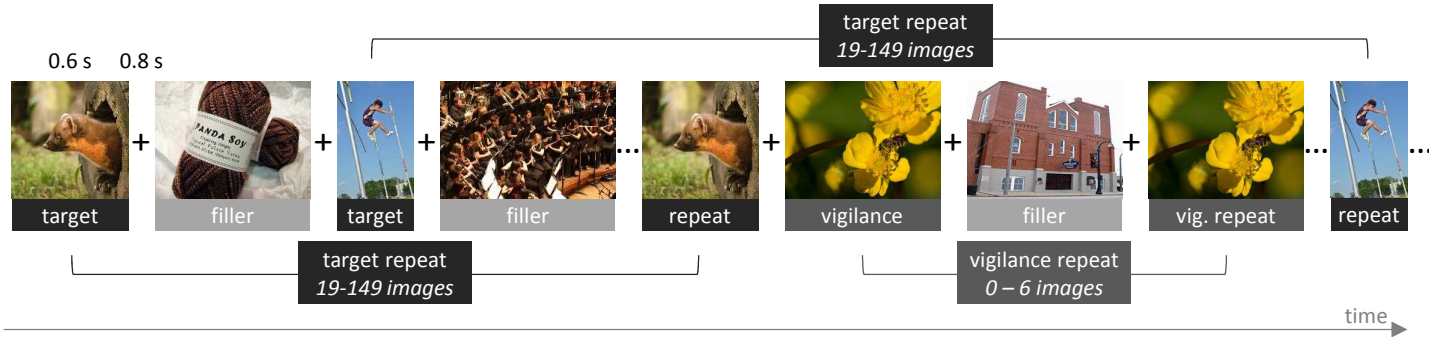


Figure 4(on next page)

Distribution of the collected memorability measures.

The left panel represents memorability scores computed as the hit rate across participants. The right panel represents scores corrected for false alarms. The horizontal lines indicate the global mean memorability scores. The asterisks represent the mean per category. Each category contains 2,000 quantified images.

H/N_{resp}

$(H - F)/N_{\text{resp}}$

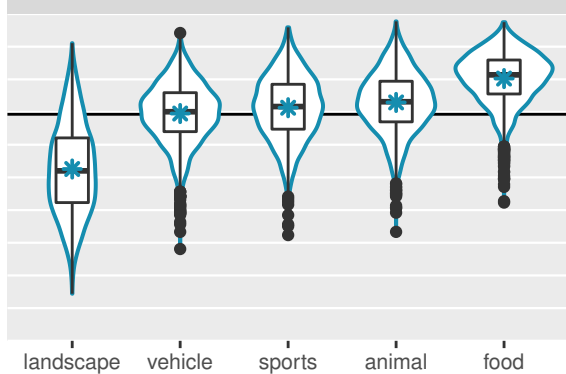
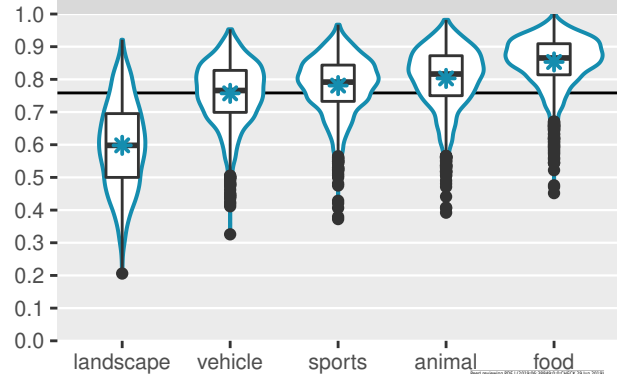


Figure 5(on next page)

Split-half consistency across participants in function of N_{resp} .

Estimates are based on 1,000 random splits. N_{resp} corresponds to the total number of data points for an image, not to the number that goes into one half during the split-half procedure. The dashed line represents predicted consistencies based on the Spearman-Brown formula (Brown, 1910; Spearman, 1910) applied to the observed consistency when N_{resp} is the maximum number of available data points.

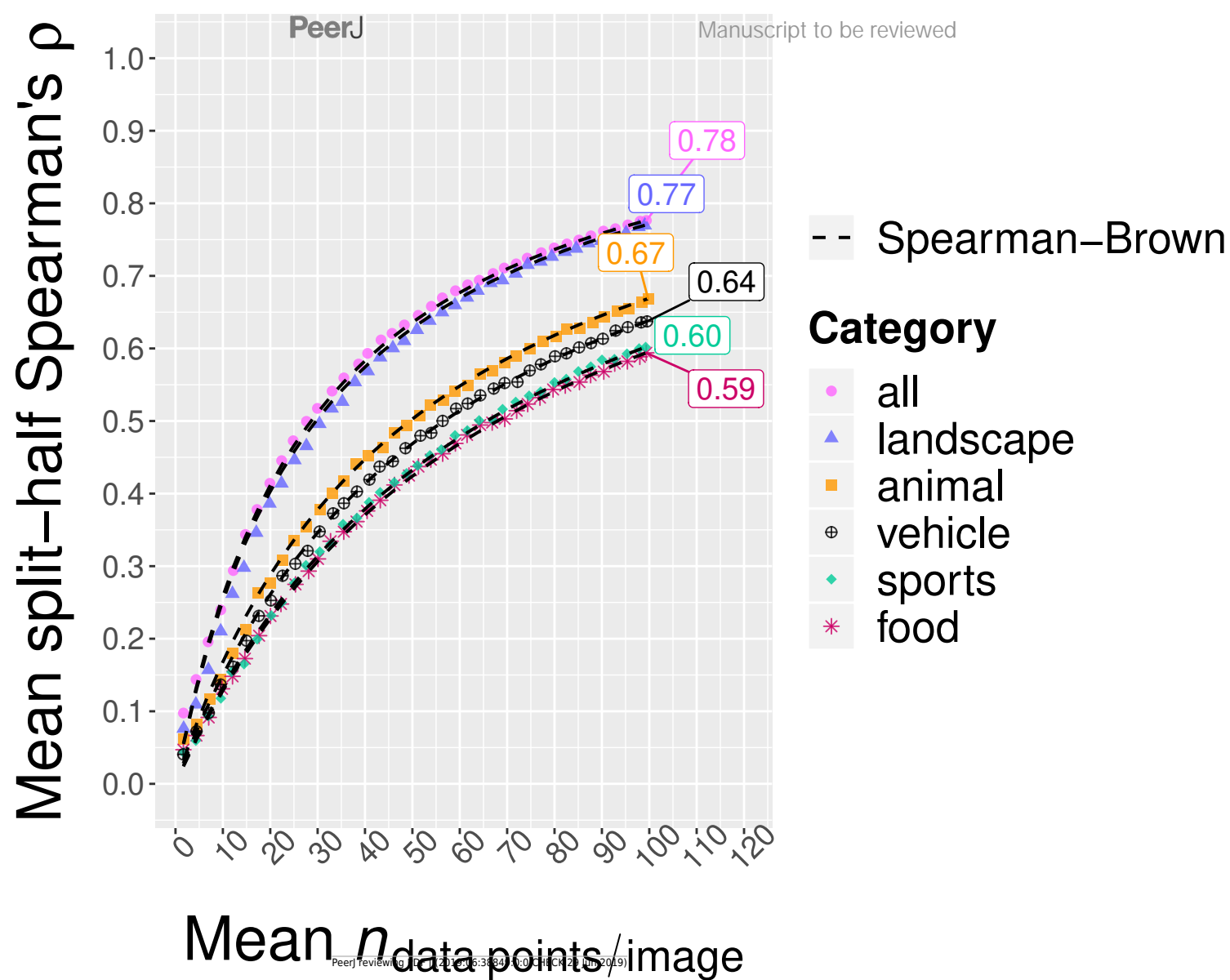


Table 1 (on next page)

Comparison MemCat to other memorability datasets.

	Isola et al. (2014)	FIGRIM	LaMem	MemCat
Category-based	no	yes	no	yes
Number of quantified images	2222	1754	~60K	10K
Bounding boxes or segmentation data	yes	yes	no	yes

Table 2(on next page)

Recognition memory performance.

The table presents descriptive statistics across participants ($n = 2291$) for five Signal Detection Theory measures. See Macmillan and Creelman (2005) for an explanation of these measures.

					False	Prop.
	d'	β	Hit rate	alarm rate	correct	
Mean	2.50	4.43	.76	.05	.87	
Median	2.48	3.00	.79	.04	.88	
SD	0.49	5.48	.14	.04	.05	
Min	0.69	0.09	.03	.00	.60	
Max	4.46	98.26	1.00	.49	.98	