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

Lore Goetschalckx [Corresp., 1] , Johan Wagemans [1]

[1] Brain & Cognition, KU 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 the study of neural and behavioral correlates of memorability while controlling for semantic category.

1 **MemCat: A new category-based image set quantified**
2 **on memorability**
3
4

5 Lore Goetschalckx[1], Johan Wagemans[1]
6
7 [1]Brain & Cognition, KU Leuven, Leuven, Vlaams-Brabant, Belgium
8
9 Corresponding author:
10 Lore Goetschalckx[1]
11 Tiensestraat 102, Leuven, Vlaams Brabant, 3000, Belgium
12 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 10,000 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 10,000 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 the study of neural and behavioral correlates of memorability while controlling for semantic category.

## Introduction

A large body of work within the visual memory field has been devoted to questions about its capacity and fidelity (for a review, see Brady, Konkle, & Alvarez, 2011). Often, these studies make abstraction of the to-be-remembered stimuli and potential differences between them. Yet, work by Isola, Xiao, Parikh, Torralba, and Oliva (2014), using everyday images, showed that they do not all share the same baseline likelihood of being remembered and recognized later. Instead, images differ in "memorability" in ways that are consistent across participants and this can be measured reliably (Isola et al., 2014).

To assess memorability, Isola et al. (2014) used 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. The researchers assigned a memorability score to 2222 scene images 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 furthermore 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, 2019). 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 (Gygli, Grabner, Riemenschneider, Nater, & Van Gool, 2013; 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 memorability score of an image, the best results so far are achieved using convolutional neural networks (CNNs; e.g., Khosla et al., 2015). When it comes to truly understanding, on the other hand, CNNs have often been critiqued to be black boxes (however, see Benitez, Castro, and Requena, 1997 for counterarguments). It is not always clear

63    to us humans why a CNN predicts a certain score for one image and not another. Nonetheless,

64    Khosla et al.'s (2015) analyses provided some further insight, mostly pointing at differences

65    between broader image categories and content types. For example, units in the network

66    displaying the highest correlation with memorability seemed to respond mostly to humans, faces,

67    and objects, while those with the lowest correlation seemed to respond to landscapes and open

68    scenes. Furthermore, the most memorable regions of an image, according to the CNN, often

69    capture people, animals or text. These findings are in line with earlier work, which also

70    predominantly revealed high-level semantic attributes. Isola et al. (2014), for example, showed

71    that the predictive performance of a model trained on mere object statistics (e.g., number of

72    objects) was boosted considerably when the object labels were taken into account. In addition, a

73    model trained on the overall scene labels alone, already predicted memorability scores with a

74    Spearman's rank correlation of .37 to the ground truth. Memorable images often had labels

75    referring to people, interiors, foregrounds, and human-scaled objects, while labels referring to

76    exteriors, wide-angle vistas, backgrounds, and natural scenes were associated with low image

77    memorability scores.

78    Together, these findings suggest that a fair share of the variability in memorability resides in

79    differences between semantic categories. Perhaps this is not surprising considering the central

80    position occupied by categories in the broader cognitive system. It has been said that carving up

81    the world around us into meaningful categories of stimuli that can be considered equivalent is a

82    core function of all organisms (Rosch, Mervis, Gray, Johnson, & Boyes-Braem, 1976). It helps

83    us understand novel events and make predictions about it (Medin & Coley, 1998). According to

84    Rosch et al. (1976), categories are represented hierarchically and organized into a taxonomy of

85    different levels of abstraction. The basic level is the best compromise between providing enough

86    information and being cognitively inexpensive. It is also the preferred naming level (e.g., "cat").

87    Other levels can be superordinate (e.g., "feline" or "mammal") or subordinate (e.g., "tabby").

88    Recently, Akagunduz, Bors, and Evans (2019) pointed out that categories are also used to

89    organize memory. More specifically, instead of encoding an image as a mere collection of pixels,

90    we extract visual memory schemas associated with its category (i.e., key regions and objects and

91    their interrelations), along with an image's idiosyncrasies. To map these visual memory schemas,

92    they had participants indicate which image regions made them recognize the image. The

93    resulting maps showed high consistency across participants, suggesting that visual memory

94    schemas partly determine what participants find memorable. Moreover, humans can visually

95    categorize an object depicted in an image very rapidly and accurately (e.g., Bacon-Macé, Macé,

96    Fabre-Thorpe, & Thorpe, 2005; Fei-Fei, Iyer, Koch, & Perona, 2007; Greene & Oliva, 2009), as

97    well as categorize the image at the level of the whole scene (e.g., Delorme, Richard, & Fabre-

98    Thorpe, 2000; VanRullen & Thorpe, 2001; Xu, Kankanhalli, & Zhao, 2019). Often a single

99    glance suffices. Interestingly, Broers, Potter, and Nieuwenstein (2018) observed enhanced

100   recognition performance for memorable versus non-memorable images in an ultra-rapid serial

101   visual presentation task. Finally, there is also evidence for a category hierarchy in the

102   representations in high-level human visual cortex, with for example clusters for animacy and

103   subclusters for faces and body parts (Carlson, Tovar, Alink, & Kriegeskorte, 2013; Cichy,

104   Pantazis, & Oliva, 2014; Kriegeskorte et al., 2008).

105   While semantic categories (or labels) seem to play a large role in image memorability, they do

106   not explain all the observed variability. Interestingly, consistent differences in memorability

107   scores remain even *within* image categories (Bylinskii et al., 2015). Goetschalckx, Moors,

108   Vanmarcke, and Wagemans (2019) for example, have argued that part of that variability might

109   be due to differences in how well the image is organized. Nonetheless, the concept of

110   memorability and its correlates is not yet fully understood to date and further research is required

111   to paint a clearer picture. The current work presents a novel, category-based dataset of images

112   quantified on memorability, designed for research to achieve this goal (example images in Figure

113   1).

114   To the best of our knowledge, there were previously five large image sets with memorability

115   scores, three consisting of regular photographs, which is also the focus here: Isola et al. (2014),

116   FIGRIM (Bylinskii et al., 2015), and LaMem (Khosla et al., 2015), and two more specialized

117   sets, which we will not further discuss here: Bainbridge, Isola, and Oliva (2013; face images),

118   and Borkin et al. (2013; data visualizations). For completeness, we also mention a smaller set

119   (850 images) that was used to study which objects in an image are memorable (Dubey, Peterson,

120   Khosla, Yang, & Ghanem, 2015). Table 1 compares MemCat to the other large datasets. The

121   comparison is discussed in more detail below.

122   A first feature of the current dataset is its hierarchical category structure. It was designed to be

123   representative for five different broad natural categories and to allow the study of memorability

124   differences within semantic categories. The set is characterized by a hierarchy of five broader

125   categories, further divided into more fine-grained subcategories. Only the FIGRIM set also offers

126   a category structure, but the number of exemplar images per category was lower: 59-157,

127   compared to 2000 in the current set. We opted for broad categories to ensure that the whole was

128   still varied and representative enough, while containing a large number of exemplar images per

129   category at the same time. Moreover, our final choice of categories: animal, food, landscape,

130   sports, and vehicle, was motivated by their relevance for memorability, meaning that they (or

131   related categories) have been observed to differ in their overall memorability in previous

132   research (Isola, Xiao, Parikh, Torralba, & Oliva, 2011; Isola et al., 2014; Aditya Khosla, Raju,

133   Torralba, & Oliva, 2015). For example, knowing that the presence of people in an image is

134   predictive for memorability (see above), we chose one category of images depicting  people as

135   the main subject and avoided including images with people in other categories. For this one

136   category, we chose "sports", because "people" in itself constitutes a category that was too broad

137   in comparison to the other categories and did not lend itself well for a division into

138   subcategories. Furthermore, we included an animal category as a non-human animate category,

139   food and vehicle as more object-based categories, and landscape to represent the wide exteriors

140   that are often associated with lower memorability scores (Isola et al., 2011, 2014; Khosla et al.,

141   2015).

142   Second, we aimed for a large set, such that it would be suitable for machine learning approaches.

143   With a total of 10,000 images quantified on memorability, the current set is the second largest

144   memorability dataset, after LaMem.

145   Third, we sampled images from existing datasets, such that the image annotations collected there

146   would also be available for researchers studying memorability. In particular, we searched for

147   images annotated with segmentation masks or at least bounding boxes, reasoning that they may

148   hold some indications of how the image is organized (e.g., where is the subject located), which

149   might be of particular interest when studying memorability within categories and factors other

150   than semantics.

151   In summary, the unique combination of features of MemCat, together with its richness in data,

152   make it a valuable addition to the memorability. Among the possible uses by memorability

153   researchers are the study of what makes an image memorable beyond its category, a benchmark

154   for machine learning approaches, and a semantically controlled stimulus set for psychophysical

155   or neuroscientific studies about the correlates of memorability (elaborated in the Discussion

156     section). However, given that categorization is a core function of the human mind, MemCat

157     would also appeal to a much broader range of cognitive (neuro)scientists.

## Materials & Methods

159     **Participants.** There were 249 undergraduate psychology students (KU Leuven) who participated

160     in this study in exchange for course credits (216 female, 32 male, 1 other). Four students did not

161     disclose their age and the remainder were aged between 18 and 27 years old ($M = 19.24$, $SD =$

162     0.94). The majority of the participants, however, were recruited through Amazon's Mechanical

163     Turk (AMT) and received a monetary compensation (see further for details). The settings on

164     AMT were chosen such that only workers who indicated to be at least 18 years old and living in

165     the USA could participate. Further eligibility criteria were that the worker had to have an

166     approval rate of at least 95% on previous human intelligence tasks (HITs) and a total number of

167     previously approved HITs of at least 100. A total of 2162 AMT-workers participated in this

168     study (1139 female, 917 male, 4 other, and 102 who did not disclose this information). For the

169     1851 workers who disclosed their age, the reported ages ranged between 18 and 82 years old ($M$

170     $= 37.14$, $SD = 11.89$). The AMT data collection took place from April 2018 till July 2018. Data

171     collected through AMT has been shown to come from participant samples that are more diverse

172     than student samples and to be comparable in quality and reliability to those collected in the lab

173     (e.g., Buhrmester, Kwang, & Gosling, 2011).

174     **Materials.** MemCat consists of 10,000 images sampled from four previously existing image sets:

175     ImageNet (Deng et al., 2009), COCO (Lin et al., 2014), SUN (Xiao, Hays, Ehinger, Oliva, &

176     Torralba, 2010), and The Open Images Dataset V4 (Kuznetsova et al., 2018). The four source

177     sets were chosen because of their large size (i.e., number of images), the availability of semantic

178     annotations, and the availability of bounding box annotations or more complete segmentation

179     masks for at least a subset of their images. The images selected from the source sets to be

180     included in MemCat belonged to the five broader semantic categories outlined in the

181     Introduction: animal (2000 images), food (2000 images), landscape (2000 images), sports (2000

182     images), and vehicle (2000 images). We explain the different steps in the selection procedure in

183     more detail below.

184     As a first step, we listed at least 20 subcategories for each broader category. The goal was to

185     obtain 2000 images per category, without including more than 100 exemplar images per

186     subcategory. This was to ensure a reasonable level of variability and to avoid high levels of false

187     alarms in the memory task (see further). The subcategories were then translated to semantic

188     annotations from the source dataset. For example, for the subcategory "bear" (animal), we used

189     COCO images annotated with a "bear" tag and ImageNet images from nodes "American black

190     bear". "brown bear", and "grizzly". An overview of our hierarchy of categories and

191     subcategories, can be found in Figure 2.

192     The second step consisted of automatically sampling exemplar images from the listed

193     subcategories, while satisfying a number of shape restrictions. To avoid that images would stand

194     out because of an extreme aspect ratio, we only sampled images with aspect ratios between 1:2

195     and 2:1. Furthermore, the minimum resolution was set to 62,500 pixels. Finally, to ensure that

196     the images would fit comfortably on most computer monitors, we adopted a maximum height of

197     500 pixels and a maximum width of 800 pixels. However, for SUN and The Open Images

198     Dataset, only a low number of images satisfied the latter two restrictions (they were often too

199     big), which is why we opted to resize (using Hamming interpolation) images from those two

200     source datasets to meet the restrictions. Apart from the shape restriction, we also restricted the

201     sampling to images for which bounding box annotations or more complete segmentation masks

202     were available from the source datasets. Finally, we sampled more images than the target number

203     (2000 images per broad category), anticipating exclusions in the next step.

204     The third step constituted manual selection work, carried out by the first author, assisted by two

205     student-interns. We manually went through the exemplar images sampled in the previous step,

206     and eliminated images following a number of exclusion rules. The exclusion rules can roughly

207     be divided into two kinds. A first kind of exclusion rule touches upon the quality of the image.

208     We excluded images of poor image quality (e.g., very dark, very much overexposed, blurry,

209     etc.), images that did not convincingly belong to the subcategory they were assigned to,[1] images

210     in greyscale or looking like they were the result of another color filter, images that were not real

211     photographs (e.g., drawings, digitally manipulated images, computer generated images), and

212     collages. A second set of rules concerns factors that could affect the memorability of an image,

213     but were not of interest for the purpose of MemCat. One such factor is text. We excluded images

214     containing large, readable text or text not belonging to the image itself (e.g., date of capture).

215     Another factor was the presence of people in the image. There was one designated "people"

216     category, the sports category, meaning that every included exemplar image depicted one or more

217    people practicing sports. However, the presence of people was avoided in all other categories

218    (but we allowed anonymous people in the background in the vehicle category or the presence of

219    a hand in images of the food category). Furthermore, images depicting remarkably odd scenes

220    (e.g., dog wearing Santa Clause costume) were also excluded. Similarly, we avoided images

221    depicting famous places or people (e.g., Roger Federer or Cristiano Ronaldo in the sports

222    category), and images of dead, wounded or fighting animals. In addition to these exclusion rules,

223    we also tried, to the best of our ability, not to include (near) duplicate images. If the target

224    number of images was not obtained after Step 3, we reverted back to Step 2, if there were still

225    images to sample from, or to Step 1 if we needed to include additional subcategories.

226    Finally, for those categories for which more than the target number of images survived Step 3,

227    there was a fourth step to randomly down-sample the selection to the target number, assigning

228    higher sampling probabilities to images annotated with segmentation masks.

229    In addition to MemCat, we collected 10,000 filler images, that were not quantified on

230    memorability themselves, but were needed in the memory task used to quantify the other images.

231    The filler images were sampled randomly from The Open Images Dataset, but from a different

232    subset to avoid overlap.[2] As these images would function only as filler images, there were fewer

233    restrictions. For example, the images could be of any category, they were allowed to contain text,

234    etc. However, the same shape restrictions were still applied.

235    **Procedure.** Having carefully collected 10,000 images for MemCat, the next step was to quantify

236    them on memorability. Following previous work, this was achieved by presenting the images in

237    an online repeat-detection memory game (Isola et al., 2014; Khosla et al., 2015), in which

238    participants watch a sequence of images and are asked to respond when they recognize a repeat

239    of a previously shown image. Students participating for course credits played the game in the

240    university's computer labs, hosting about 20 students at a time. AMT workers played the game

241    from the comfort of their homes (or whichever location they preferred). Prior to starting the

242    game, all participants were prompted to read through an informed consent page explaining the

243    aims of the study and their rights as participants. They could give their consent by actively

244    ticking a box. The study was approved by SMEC, the Ethical Committee of the Division of

245    Humanities and Social Sciences, KU Leuven, Belgium (approval number: G-2015 08 298).

246    For the task design of the game, we closely followed Khosla et al. (2015), as their version of the

247    game was designed to quantify large numbers of images. We divided the game into blocks of

248    200 trials. On each trial, an image was presented at the center of the browser window for a

249    duration of 600 ms, with an intertrial interval of 800 ms. During this interval, a fixation cross

250    was shown. Sixty-six images were target images, sampled randomly from MemCat, and repeated

251    after 19 to 149 intervening images. Forty-four images were random filler images that were never

252    repeated. Finally, there were 12 additional random filler images that were repeated after 0 to 6

253    intervening images to keep participants attentive and motivated. They are referred to as vigilance

254    trials. Participants could indicate that they recognized a repeat by pressing the space-bar. They

255    did not receive trial-by-trial feedback, but were shown their hit rate as well as number of false

256    alarms at the end of the block. Figure 3 presents a schematic of the game.

257    Each block lasted a little less than 5 min. Care was taken to ensure that an image was never

258    repeated more than once and never across blocks. Students were asked to complete as many

259    blocks as they could in one hour, with one bigger, collective break of roughly 10 min after half

260    an hour, and smaller self-timed breaks between the remainder of the blocks. Most students could

261    complete eight blocks, but for some groups, slow data uploads at the end of a block resulted in

262    lower numbers. AMT workers could complete one to 16 blocks, were allotted 48h to submit their

263    completed blocks (so, they were allowed to spread the blocks over time), and were paid $0.40

264    per block. To ensure a good quality of the AMT data and to avoid random or disingenuous

265    responses, AMT workers were blocked from playing anymore blocks after two with a $d'$ lower

266    than 1.5 on the vigilance trials. They were warned the first time this happened.

267    **Memorability Measures.** We computed two different, but related measures of memorability

268    from the data collected through the repeat-detection memory game. These were the same two

269    measures as used in LaMem, the largest available image memorability dataset yet. As mentioned

270    in the Introduction, one measure is simply the proportion of participants recognizing the image

271    when shown to them for the second time (i.e., the hit rate across participants). This is the

272    "original" memorability measure, as introduced by Isola et al. (2014), also adopted in many other

273    memorability studies (e.g., Bainbridge, Isola, & Oliva, 2013; Bylinskii, Isola, Bainbridge,

274    Torralba, & Oliva, 2015; Khosla, Raju, Torralba, & Oliva, 2015). The other memorability

275    measure computed for the LaMem images  was based on the same principle, but penalized for

276    false alarms (i.e., when participants press the space-bar for the first presentation of the image) in

277    the way proposed by Khosla, Bainbridge, Torralba, and Oliva (2013), who applied it to a dataset

278    of face images (Bainbridge et al., 2013). Rather than $H/N_{resp}$ (first measure), their formula was

279 the following: $(H-F)/N_{resp}$, where H is the number of participants recognizing the image, F is the

280 number of participants making a false alarm when the image is presented for the first time, and

281 $N_{resp}$ is the total number of participants having been presented with the image. Here, an image's

282 $N_{resp}$ was 99 (after exclusions) on average. Note that the memorability scores have an upper

283 bound of 1 and a lower bound of 0. In theory, $(H-F)/N_{resp}$ could result in a negative score, but in

284 practice it is highly unlikely that there would be more participants making a false alarm for the

285 image than there are participants making a hit.

## Results

287 **Participant Performance.** As mentioned, the performance on the easier vigilance trials was

288 taken as an indication of whether participants were playing the memory game in a genuine way.

289 If in a certain block, a participant did not distinguish vigilance repeats from non-repeat trials with

290 a $d'$ of at least 1.5 (preset performance threshold), that block was excluded from further analyses.

291 The exclusion rate amounted to 3% of all played blocks. Recall, however, that AMT workers

292 were not allowed to play more blocks after two excluded ones.

293 After exclusion, the mean $d'$ across participants was 2.77 ($SD = 0.56$) for the vigilance repeats,

294 and 2.47 ($SD = 0.50$) for the target repeats. Table 2 summarizes participants' overall

295 performance, collapsing over vigilance and target repeats. Participants generally performed well

296 on the task.

297 **Memorability Scores.** Participants' high performance was also reflected in the average image

298 memorability scores. Figure 4 displays the mean for each of the two memorability measures as a

299 horizontal line ($M_{H/Nresp} = .76$, SD; $M_{(H-F)/Nresp} = .70$). It is comparable to the mean observed in

300 Khosla et al. (2015). In addition, Figure 4 visualizes the distribution of the collected image

301 memorability scores for each of the five broad main categories separately. A simple linear

302 regression revealed that the category explains 43% of the variance in the $H/N_{resp}$ scores and 44%

303 of the variance in the $(H-F)/N_{resp.}$ In line with previous research, the landscape images were on

304 average the least memorable ($M_{H/Nresp} = .60$; $M_{(H-F)/Nresp} = .53$). They were followed by the vehicle

305 images ($M_{H/Nresp} = .76$; $M_{(H-F)/Nresp} = .70$). Somewhat surprisingly, the food images generally came

306 out on top of the ranking ($M_{H/Nresp} = .85$; $M_{(H-F)/Nresp} = .80$), topping the animal $M_{H/Nresp} = .80$; $M_{(H-F)/Nresp} = .73$) and sports $M_{H/Nresp} = .78$; $M_{(H-F)/Nresp} = .71$) categories. However, there is still a large

308 degree of variability that is not explained by differences in broad image categories. Indeed,

309    memorability varied considerably within categories as well, with *SDs* of: .09 (animal; $SD_{(H-}$

310    $_{F)/Nresp}$ = .09), .08 (food; $SD_{(H-F)/Nresp}$ = .08), .13 (landscape; $SD_{(H-F)/Nresp}$ = .14), .09 (sports; $SD_{(H-}$

311    $_{F)/Nresp}$ = .10), and .09 (vehicle; $SD_{(H-F)/Nresp}$ = .09).

312    Having observed that images from the same broader category indeed still differed in

313    memorability, the next question was whether these differences are consistent across participants.

314    This question taps into the reliability of the memorability measures. Following previous

315    memorability work (e.g., Isola et al., 2014), the consistency was assessed by randomly splitting

316    the participant pool in half, computing the memorability scores for each half separately and

317    determining the Spearman's rank correlation between the two sets of scores. This was repeated

318    for 1000 splits and the Spearman's rank correlation was averaged across the splits. Figure 5

319    shows the results in function of the mean $N_{resp}$ for each category as well as for the total image set.

320    We first discuss the results for H/ $N_{resp}$ (see Figure 5, left panel). When collapsing over all five

321    categories, the observed mean split-half Spearman's rank correlation with all available responses

322    ($N_{resp}$ = 99, on average) amounted to .78. In comparison, Khosla et al. (2015) reported a mean

323    split-half Spearman's rank correlation of .67 for their LaMem dataset. However, they only

324    collected 80 responses per image. After randomly down-sampling our data to an $N_{resp}$ of 80, we

325    still found a split-half consistency of .73. With the exception of the landscape category, for

326    which we observed a total (i.e., without down-sampling) split-half consistency of .77, the total

327    per category split-half consistency estimates were lower, ranging between .59 and .67. This is

328    possibly due to smaller ranges of memorability scores within those categories (see Figure 4).

329    Note, however, that the split-half consistencies are an underestimate of the reliability of the

330    memorability scores calculated based on the full participant pool. The latter can be estimated

331    from the split-half consistency by means of the Spearman-Brown formula (Brown, 1910;

332    Spearman, 1910). Applying this formula, we found the following final reliabilities for the H/$N_{resp}$

333    memorability scores : .87 (all), .80 (animal), .75 (food), .87 (landscape), .75 (sports), .78

334    (vehicle).

335    For the (H-F)/$N_{resp}$ memorability scores, we confine the discussion to pointing out that the pattern

336    of results is qualitatively similar, although the final reliabilities are somewhat lower: .86 (all), .74

337    (animal), .71 (food), .85 (landscape), .77 (sports), .71 (vehicle).

338    Finally, after finding that the two image memorability measures were both acceptably reliable,

339    we asked how they compared to each other. In the current dataset, they were highly

340    intercorrelated, as evidenced by a Pearson correlation of .93 when collapsing over all five

341    categories. The per category correlations were: .82 (animal), .90 (food), .91 (landscape), .85

342    (sports), .88 (vehicle).

## Discussion

344    We presented a new dataset, MemCat, consisting of a total of 10,000 images, each quantified on

345    memorability using a repeat-detection memory task (first introduced by Isola et al., 2014, version

346    used here based on Khosla et al., 2015). MemCat is the second largest image memorability

347    dataset available, and the largest that is based on a category structure. That is, it is divided into

348    five broader, memorability-relevant semantic categories: animal, food, landscape, sports, and

349    vehicle, each with 2000 exemplar images, which are further divided into subcategories (e.g.,

350    bear, cat, cow). Furthermore, the images were sampled from popular, existing datasets such that

351    additional annotations available there (e.g., segmentations masks or bounding boxes) would also

352    be available to researchers wishing to use MemCat for research aimed at investigating specific

353    factors underlying memorability.

354    Replicating previous research, we found that images differ considerably in memorability and that

355    these differences are highly consistent across participants. Part but not all of this variability can

356    be explained by differences between the five broader semantic categories. Note, however, that

357    this result is correlational in nature, and that one should be cautious drawing causal conclusions.

358    In line with Bylinskii et al. (2015), considerable variability in memorability remained even

359    within the categories. However, the consistency there was somewhat lower, probably because the

360    variance was also lower. When the differences between images become smaller, it becomes

361    harder to reliably and consistently distinguish them. Nevertheless, the consistency estimates per

362    category were still high, indicating that we obtained reliable memorability scores. Finally, we

363    reported results for two different methods to compute memorability scores. One is to compute

364    the hit rate across participants: $H/N_{resp}$. This was the method used in the original work by Isola et

365    al. (2014). However, in principle, it possible that some images elicit more key presses not

366    because they are truly recognized, but for some other reason (e.g., they seem familiar). That is

367    why Bainbridge et al. (2013) suggested to correct for false alarms (i.e., when participants press

368    the key for the first presentation of an image, when it is not a repeat) by computing $(H-F)/N_{resp}$.

369    We report both measures for comparison, but note that they lead to a highly similar pattern of

370    results and are also strongly intercorrelated. In what follows, we discuss possible uses of

371    MemCat.

372    Most of what we learned from previous studies about what makes an image memorable is

373    specified in terms of semantic categories or content types (e.g., images of people are more

374    memorable than landscapes). However, a considerable amount of variability was still left

375    unexplained. A primary use of the current dataset is in studies aiming to better understand the

376    factors underlying image memorability. In particular, with 2000 images for each of five broader

377    categories, it allows to zoom in on variability within categories. This variability is of more

378    interest to practical applications (e.g., advertising, education), because the semantic category or

379    the content type (e.g., a certain product) will often be predefined and it will be a matter of

380    choosing or creating a more memorable depiction of it. In addition to dividing the set into broad

381    semantic categories, we also avoided variability due to other factors already discovered in

382    previous studies (e.g., we excluded images depicting oddities, images containing text or

383    recognizable places or faces), thus creating a set designed to help understand the previously

384    unexplained variability in image memorability.

385    Second, MemCat is also useful as a benchmark for machine learning approaches to automatically

386    predict memorability. Currently, LaMem (Khosla et al., 2015) is most often used, but models can

387    now also be trained and tested on the current dataset. When taking Khosla et al.'s (2015)

388    MemNet-CNN (without retraining), we found that its predictions show a rank correlation of .68

389    with the $(H-F)/N_{resp}$ memorability scores in the current set, suggesting that there is room for

390    improvement. Given the category structure in MemCat, one could explore, for the first time,

391    memorability models with one or more layers that are specific to a category. Indeed, it is possible

392    that what makes landscape images memorable is different from what makes animal images

393    memorable.

394    Finally, a third possible use is in neuroscientific studies or psychophysical studies examining

395    effects of memorability. The current set offers a large number of quantified images to choose

396    from. Moreover, it facilitates matching memorability conditions (e.g., high versus low) on

397    semantic category, something that is often done in neuroscientific studies (e.g., Bainbridge,

398    Dilks, & Oliva, 2017; Khaligh-Razavi, Bainbridge, Pantazis, & Oliva, 2016; Mohsenzadeh,

399    Mullin, Oliva, & Pantazis, 2019).

400   **Usage.** On the MemCat project page, http://gestaltrevision.be/projects/memcat/, we provide a
401   link to the collection of 10,000 images as well as links to two data files, all hosted on OSF (also
402   see Additional Information). One file describes the images that were used and contains columns
403   indicating the image filename in its source dataset, the name of its source dataset, the category
404   (e.g., animal) and subcategory (e.g., bear) we assigned it to, the label that was used to sample it
405   from its source dataset (e.g., American black bear), the current width, the current height, the
406   factor by which it was resized (both the original width and height were multiplied by this factor),
407   the number of hits (H), the number of false alarms (FA), the number of participants it was
408   presented to ($N_{resp}$), and the two memorability scores. The other file contains the data collected in
409   the repeat-detection memory game. Its columns indicate the participant ID (anonymized), the
410   participant's age, the participant's gender, whether or not they participated through AMT, the
411   block number, the trial number, the image shown, the trial type (target, target repeat, filler,
412   vigilance, vigilance repeat), the participant's response (hit, correct rejection, miss, false alarm),
413   the screen width, and the screen height.

## Conclusions

415   With MemCat, we present a large new dataset of 10,000 images fully annotated with ground
416   truth memorability scores collected through an online repeat-detection memory task. It is the
417   second largest memorability dataset to date and the largest with a hierarchical category structure.
418   The results showed that images differ in memorability in ways that are consistent across
419   participants, even within semantic categories. Among other things, MemCat allows the study of
420   which factors might underlie such differences. Its richness in data and unique combination of
421   features will appeal to a broad range of researchers in cognitive science and beyond (e.g.,
422   computer vision).

423

## 424 **Acknowledgments**

425 The authors would like to thank the student-interns who have assisted in the image selection and

426 data collection: Justine Aeyels and Joran Geeraerts, and Christophe Bossens, who is responsible

427 for the technical support in the lab and contributed greatly to the implementation of the repeat-

428 detection memory task on Amazon's Mechanical Turk.

## Endnotes

429

430  [1] This could happen, for example, with images from the COCO source dataset. COCO images do

431  not come with a single, overall scene label, but instead come with multiple semantic tags

432  describing what is in the image. For this reason, an image annotated with the tag "cat," for

433  instance, could be more of a living room image that just happens to have a cat sleeping

434  somewhere in a corner in the background.

435  [2] The source set is presented in three different subsets: train, validation, and test. We sampled

436  from the validation and test subsets for MemCat, and from the train subset for the fillers images

437  used in the memory task.

## References

Akagunduz, E., Bors, A., & Evans, K. (2019). Defining image memorability using the visual memory schema. *IEEE Transactions on Pattern Analysis and Machine Intelligence*. Advance online publication. doi:10.1109/tpami.2019.2914392

Bacon-Macé, N., Macé, M. J.-M., Fabre-Thorpe, M., & Thorpe, S. J. (2005). The time course of visual processing: Backward masking and natural scene categorisation. *Vision Research, 45*(11), 1459–1469. doi:10.1016/j.visres.2005.01.004

Bainbridge, W. A. (2017). The resiliency of memorability: A predictor of memory separate from attention and priming. *arXiv e-prints*. arXiv: 1703.07738 [q-bio.NC]

Bainbridge, W. A., Dilks, D. D., & Oliva, A. (2017). Memorability: A stimulus-driven perceptual neural signature distinctive from memory. *NeuroImage, 149*, 141–152. doi: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. doi:10.1037/a0033872

Benitez, J. M., Castro, J. L., & Requena, I. (1997). Are artificial neural networks black boxes? *IEEE Transactions on Neural Networks, 8*(5), 1156–1164. doi:10.1109/72.623216

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. doi:10.1109/TVCG.2013.234

Brady, T. F., Konkle, T., & Alvarez, G. A. (2011). A review of visual memory capacity: Beyond individual items and toward structured representations. *Journal of Vision, 11*(5:4), 1–34. doi:10.1167/11.5.4

Broers, N., Potter, M. C., & Nieuwenstein, M. R. (2017). Enhanced recognition of memorable pictures in ultra-fast RSVP. *Psychonomic Bulletin & Review, 25*(3), 1080–1086. doi: 10.3758/s13423-017-1295-7

Brown, W. (1910). Some experimental results in the correlation of mental abilities. *British Journal of Psychology, 1904-1920, 3*(3), 296–322. doi:10.1111/j.2044-8295.1910.tb00207.x

467    Buhrmester, M., Kwang, T., & Gosling, S. D. (2011). Amazon's mechanical turk: A new source
468        of inexpensive, yet high-quality, data? *Perspectives on Psychological Science*, *6*(1), 3–5.
469        doi:10.1177/1745691610393980

470    Bylinskii, Z., Isola, P., Bainbridge, C., Torralba, A., & Oliva, A. (2015). Intrinsic and extrinsic
471        effects on image memorability. *Vision Research*, *116*(Part B), 165–178.
472        doi:10.1016/j.visres.2015.03.005

473    Carlson, T., Tovar, D. A., Alink, A., & Kriegeskorte, N. (2013). Representational dynamics of
474        object vision: The first 1000 ms. *Journal of Vision*, *13*(10), 1–1. doi:10.1167/13.10.1

475    Cichy, R. M., Pantazis, D., & Oliva, A. (2014). Resolving human object recognition in space and
476        time. *Nature Neuroscience*, *17*(3), 455–462. doi:10.1038/nn.3635

477    Delorme, A., Richard, G., & Fabre-Thorpe, M. (2000). Ultra-rapid categorisation of natural
478        scenes does not rely on colour cues: A study in monkeys and humans. *Vision Research*,
479        *40*(16), 2187–2200. doi:10.1016/S0042-6989(00)00083-3

480    Deng, J., Dong, W., Socher, R., Li, L. J., Li, K., & Fei-Fei, L. (2009). ImageNet: A large-scale
481        hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern
482        Recognition (CVPR)*, 248–255. Red Hook, NY: Curran Associates.
483        doi:10.1109/CVPR.2009.5206848

484    Dubey, R., Peterson, J., Khosla, A., Yang, M. H., & Ghanem, B. (2015). What makes an object
485        memorable? *2015 IEEE International Conference on Computer Vision (ICCV)*, 1089–
486        1097. doi:10.1109/ICCV.2015.130

487    Fei-Fei, L., Iyer, A., Koch, C., & Perona, P. (2007). What do we perceive in a glance of a real-
488        world scene? *Journal of Vision*, *7*(1), 10. doi:10.1167/7.1.10

489    Goetschalckx, L., Moors, J., & Wagemans, J. (2019). Incidental image memorability. *Memory,
490        27*(9), 1273-1282. doi:10.1080/09658211.2019.1652328

491    Goetschalckx, L., Moors, P., Vanmarcke, S., & Wagemans, J. (2019). Get the picture? Goodness
492        of image organization contributes to image memorability. *Journal of Cognition, 2*(1).
493        doi:10.5334/joc.80

494    Goetschalckx, L., Moors, P., & Wagemans, J. (2018). Image memorability across longer time
495        intervals. *Memory*, *26*(5), 581–588. https://doi.org/10.1080/09658211.2017.1383435

496     Greene, M., & Oliva, A. (2009). Recognition of natural scenes from global properties: Seeing the

497          forest without representing the trees. *Cognitive Psychology*, *58*(2), 137–176.

498          doi:10.1016/j.cogpsych.2008.06.001

499     Gygli, M., Grabner, H., Riemenschneider, H., Nater, F., & Van Gool, L. (2013). The

500          interestingness of images. In *2013 IEEE International Conference on Computer* Vision

501          (ICCV). Red Hook, NY: Curran Associates. doi:10.1109/iccv.2013.205

502     Isola, P., Parikh, D., Torralba, A., & Oliva, A. (2011). Understanding the intrinsic memorability

503          of images. In J. Shawe-Taylor, R. S. Zemel, P. L. Bartlett, F. Pereira, & K. Q.

504          Weinberger (Eds.), *Advances in Neural Information Processing Systems (NIPS) 24* (pp.

505          2429-2437). Red Hook, NY: Curran Associates. doi:10.21236/ada554133

506     Isola, P., Xiao, J., Parikh, D., Torralba, A., & Oliva, A. (2014). What makes a photograph

507          memorable? *IEEE Transactions on Pattern Analysis and Machine Intelligence*, *36*(7),

508          1469–1482. doi:10.1109/TPAMI.2013.200

509     Khaligh-Razavi, S.-M., Bainbridge, W. A., Pantazis, D., & Oliva, A. (2016). From what we

510          perceive to what we remember: Characterizing representational dynamics of visual

511          memorability. *bioRxiv*. doi:10.1101/049700. eprint:

512          https://www.biorxiv.org/content/early/2016/04/22/049700.full.pdf

513     Khosla, A., Bainbridge, W. A., Torralba, A., & Oliva, A. (2013). Modifying the memorability of

514          face photographs. In 2*013 IEEE International Conference on Computer Vision (ICCV)*.

515          Red Hook, NY: Curran Associates. doi:10.1109/iccv.2013.397

516     Khosla, A., Raju, A. S., Torralba, A., & Oliva, A. (2015). Understanding and predicting image

517          memorability at a large scale. In *2015 IEEE International Conference on Computer*

518          *Vision (ICCV)*, 2390–2398. Red Hook, NY: Curran Associates.

519          doi:10.1109/ICCV.2015.275

520     Kriegeskorte, N., Mur, M., Ruff, D. A., Kiani, R., Bodurka, J., Esteky, H., … Bandettini, P. A.

521          (2008). Matching categorical object representations in inferior temporal cortex of man

522          and monkey. *Neuron*, *60*(6), 1126–1141. doi:10.1016/j.neuron.2008.10.043

523     Kuznetsova, A., Rom, H., Alldrin, N., Uijlings, J., Krasin, I., Pont-Tuset, J., … Ferrari, V.

524          (2018). The open images dataset V4: Unified image classification, object detection, and

525          visual relationship detection at scale. *CoRR*. arXiv:1811.00982.

526 Lin, T.-Y., Maire, M., Belongie, S., Bourdev, L., Girshick, R., Hays, J., … Dollár, P. (2014).

527         Microsoft COCO: Common objects in context. In D. Fleet, T. Pajdla, B. Schiele, & T.

528         Tuytelaars (Eds.), *2014 European Conference on Computer Vision (ECCV)* (pp. 740–

529         755). Cham: Springer International Publishing. doi:10.1007/978-3-319-10602-1_48

530 Medin, D. L., & Coley, J. D. (1998). Concepts and categorization. In J. Hochberg (Ed.),

531         *Perception and cognition at century's end* (pp. 403–439). San Diego, CA: Academic

532         Press. doi:10.1016/B978-012301160-2/50015-0

533 Mohsenzadeh, Y., Mullin, C., Oliva, A., & Pantazis, D. (2019). The perceptual neural trace of

534         memorable unseen scenes. *Scientific Reports, 9*(1), 6033. doi:10.1038/s41598-019-

535         42429-x

536 Rosch, E., Mervis, C. B., Gray, W. D., Johnson, D. M., & Boyes-Braem, P. (1976). Basic objects

537         in natural categories. *Cognitive Psychology*, *8*(3), 382–439. https://doi.org/10.1016/0010-

538         0285(76)90013-X

539 Spearman, C. (1910). Correlation calculated from faulty data. *British Journal of Psychology,*

540         *1904-1920*, *3*(3), 271–295. https://doi.org/10.1111/j.2044-8295.1910.tb00206.x

541 VanRullen, R., & Thorpe, S. J. (2001). Is it a bird? Is it a plane? Ultra-rapid visual categorisation

542         of natural and artifactual objects. *Perception*, *30*(6), 655–668.

543         https://doi.org/10.1068/p3029

544 Xiao, J., Hays, J., Ehinger, K. A., Oliva, A., & Torralba, A. (2010). SUN database: Large-scale

545         scene recognition from abbey to zoo. In *2010 IEEE Conference on Computer Vision and*

546         *Pattern Recognition (CVPR)* (pp. 3485–3492). Red Hook, NY: Curran Associates.

547         doi:10.1109/CVPR.2010.5539970

548 Xu, B., Kankanhalli, M. S., & Zhao, Q. (2019). Ultra-rapid object categorization in real-world

549         scenes with top-down manipulations. *PLOS ONE*, *14*(4), e0214444.

550         https://doi.org/10.1371/journal.pone.0214444

# Figure 1

Example images of MemCat.

*The memorability score, calculated as the hit rate across participants (H/Nresp) is indicated in the bottom right corner. In line with previous research, images differed consistently in their memorability score, even within semantic categories. MemCat represents five broader semantic categories: animal, food, landscape, sports and vehicle. Each row (A–C) displays exemplar images in that category order.*
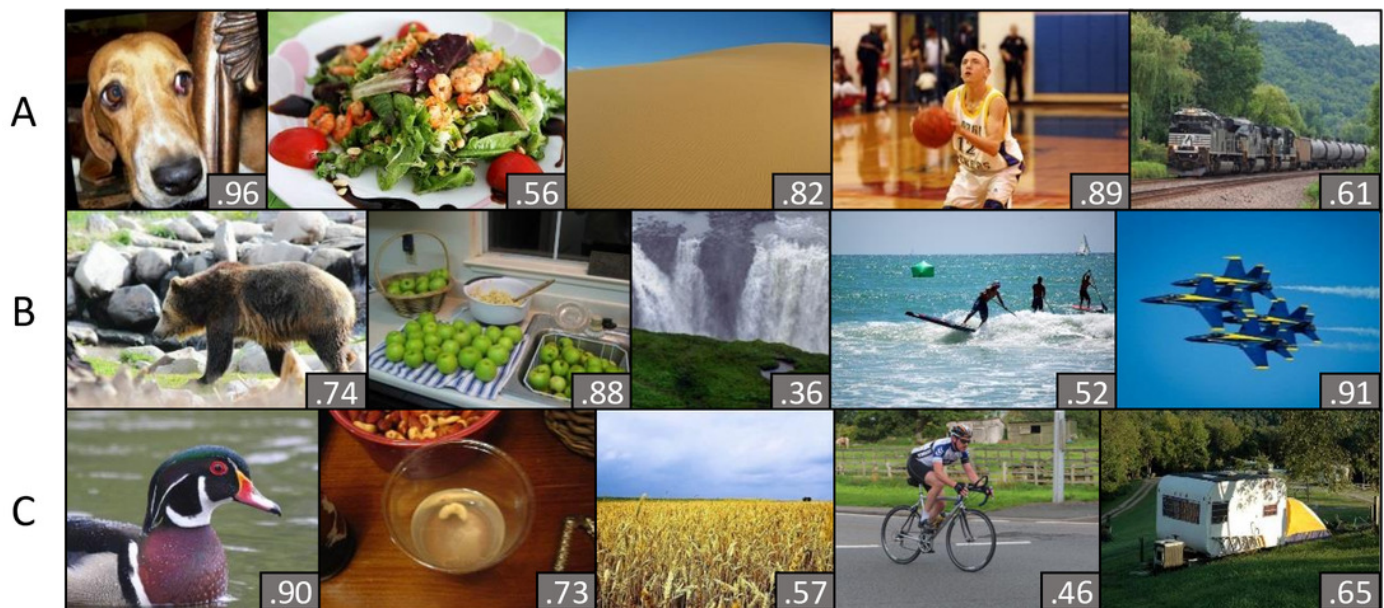
# Figure 2

Category hierarchy of MemCat.

animal - 2K
- bear - 100
- cat - 100
- chicken - 100
- cow - 100
- deer - 99
- dog - 100
- duck - 100
- elephant - 100
- fox - 99
- giraffe - 99
- hamster - 97
- horse - 99
- lion - 100
- marten - 61
- mouse - 48
- pigeon - 99
- rabbit - 99
- sheep - 100
- squirrel - 100
- wild boar - 100
- zebra - 100

food - 2K
- apple - 31
- bagel - 80
- banana - 82
- barbecue - 45
- bell pepper - 55
- bread - 84
- breakfast - 75
- broccoli - 86
- burrito - 81
- butternut squash - 87
- cake - 87
- champagne - 43
- cheese - 39
- cocktail - 64
- coffee - 61
- donut - 54
- gravy - 73
- hotdog - 82
- juice - 39
- mashed potato - 85
- muffin - 51
- orange - 35
- pasta - 82
- pizza - 84
- salad - 87
- sandwich - 82
- soup - 84
- tea - 49
- yogurt - 63
- zucchini - 50

landscape - 2K
- badlands - 95
- barn - 86
- creek - 88
- desert - 79
- field - 96
- forest (broadleaf) - 94
- forest (needleleaf) - 55
- geyser - 95
- icescape - 89
- jungle - 57
- lake - 95
- mangrove - 41
- moor - 4
- mountain - 92
- pasture - 61
- plantation - 96
- river - 94
- savanna - 49
- sea - 80
- shore - 93
- snowscape - 48
- tundra - 14
- valley - 99
- volcano - 73
- waterfall - 98
- wetland - 96
- windmill - 33

MemCat - 10K

sports - 2K
- American football - 26
- archery - 31
- baseball - 99
- basketball - 69
- bicycling racing - 63
- boxing - 47
- climbing - 100
- cricket - 27
- fencing - 21
- field hockey - 54
- figure skating - 21
- fitness - 23
- frisbee - 53
- golf - 99
- gymnastics - 62
- ice hockey - 77
- japanese martial arts - 78
- karate - 30
- lacrosse - 54
- mountain biking - 39
- roller skating - 64
- rowing - 100
- rugby - 73
- skateboarding - 85
- skiing - 100
- snowboarding - 60
- soccer - 98
- surfing - 88
- swimming - 80
- tennis - 32
- track and field - 99
- volleyball - 48

vehicle - 2K
- airplane - 98
- bicycle - 91
- bus - 67
- cable car - 62
- camper - 70
- car - 64
- fishing boat - 33
- forklift - 72
- freighter - 99
- helicopter - 98
- hot-air balloon - 94
- limousine - 96
- motorcycle - 94
- pleasure boat - 68
- scooter - 100
- snowmobile - 20
- subway train - 99
- SUV - 98
- tow truck - 93
- tractor - 96
- train - 98
- tram - 98
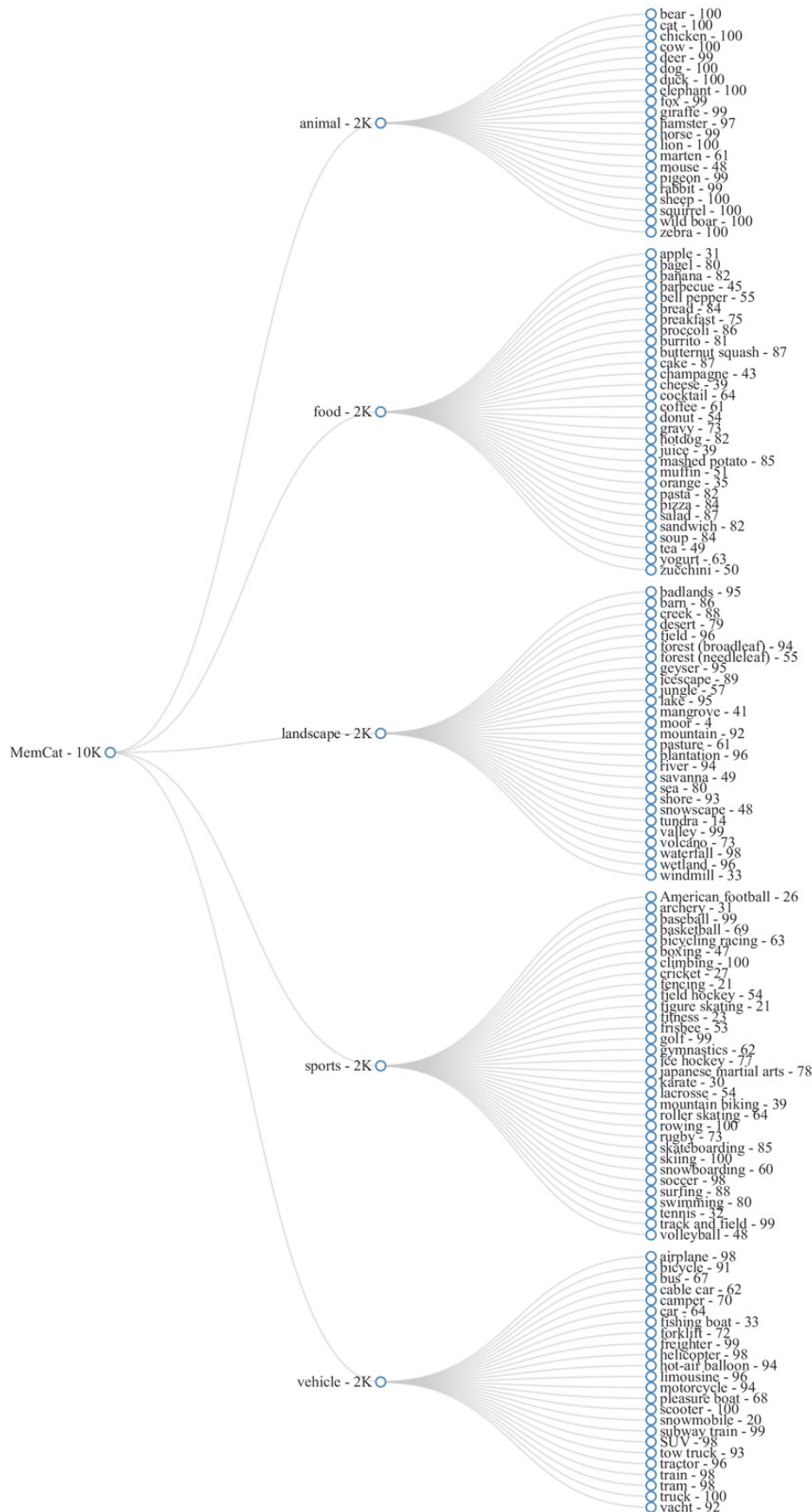- truck - 100
- yacht - 92

# Figure 3

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

Each image is presented for 600 ms, with an intertrial interval of 800 ms. Participants are instructed to press the space-bar whenever they recognize a repeat of a previously shown image.
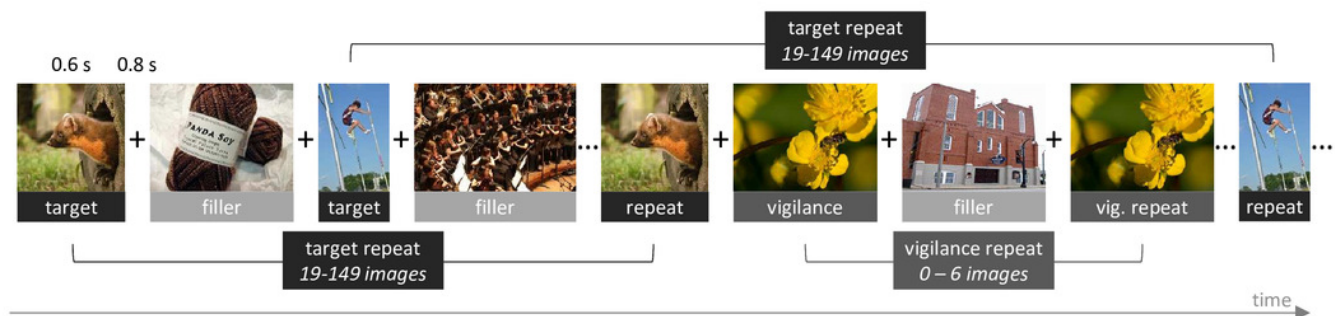
# Figure 4

Distribution of the collected memorability measures.

Panel A represents memorability scores computed as the hit rate across participants. Panel B 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 2000 quantified images. In addition to overall differences across categories, we observed considerable variability in memorability within categories too.
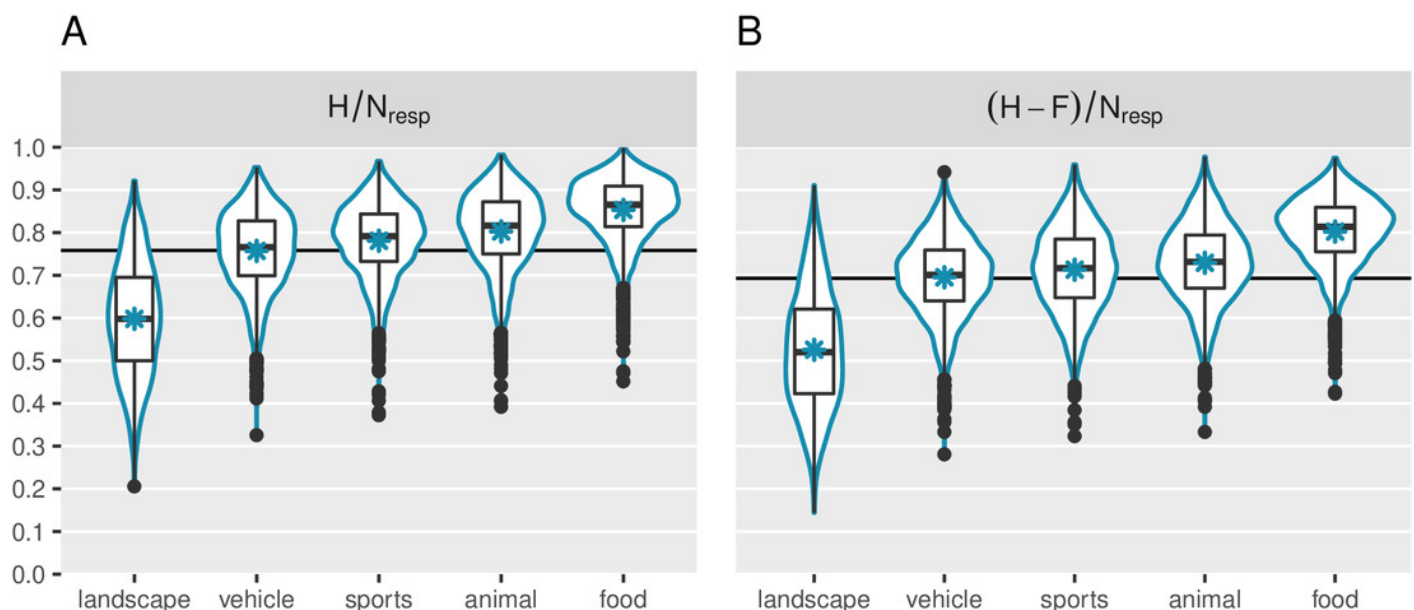
# Figure 5

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

Estimates are based on 1000 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. Even though the consistency was lower when zooming in on a single category compared to the whole set at once (possibly due to a smaller range of scores), images still showed highly consistent differences in memorability within categories.
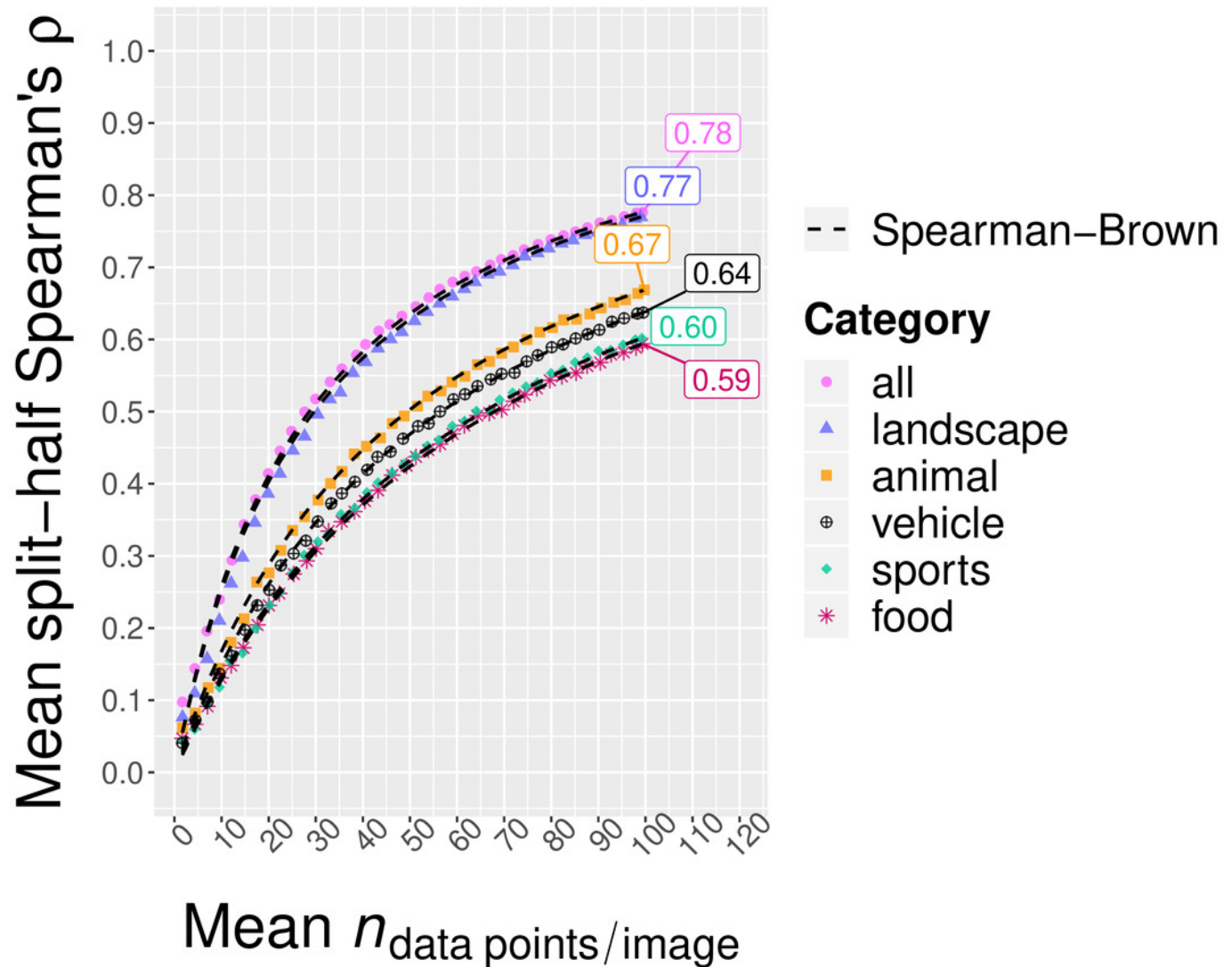
**Table 1**(on next page)

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 |

1

**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.

| | d′ | β | Hit rate | False alarm rate | Prop. 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 |

1