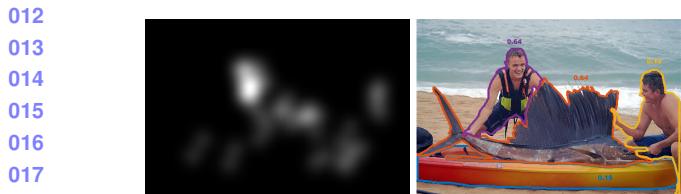


000
001
002
003
004
005
006
007
008
009
010
011

012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053

Figure 1: Not all objects are equally remembered. Memorability scores of objects for the image in the top row obtained from our psychophysics experiment.

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053

Anonymous ICCV submission

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053

Paper ID ****

Abstract

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

1. Introduction

Consider the image and its corresponding objects in Figure 1. Even though the person on the right is comparable in size to person on the left, he is remembered far less by human subjects (indicated by their respective memorability scores of 0.18 and 0.64). Moreover, people tend to remember the person on the left and the fish in the center, even after 3 minutes and more than 70 additional stimuli have

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

intervened (memorability score = 0.64). Interestingly, despite vibrant colors and considerable size, the boat is far less memorable (score = 0.18).

Considerable effort has recently been aimed at understanding and predicting image memorability, and with good reason: such a property of images is useful to understand as it can be leveraged in a variety of fascinating applications [cite papers here including modifying face memorability]. However, such an analysis does not provide direct information about just what exactly in those images are the memorable parts. Consider again the example in Figure 1. What is it about the fish and the person on the left that makes them much more memorable than the other objects in the image? Moreover, how does this object level memorability influence the overall memorability of the photo? Is it by virtue of a single object or a combination/configuration of multiple objects? Inferring such detailed information from image-level knowledge alone is an extremely difficult task.. An understanding of memorability at an object level provides relevant information at a more granular level and may also eventually produce bottom up explanation of image memorability. Such a complete picture about memorability can provide vision systems with information that humans consider meaningful which can lead to a wide variety of applications. As an analogy to what humans remember and consider meaningful in an image, we draw attention to the large body of work on saliency that has sought to determine what is important in an image. Saliency, like memorability and aesthetic value, is an image property that machines must leverage in order to be utilized effectively by humans. What is more, these properties may influence each other as well. For example, while memorability does not appear to be increased by aesthetic value, saliency may actively interact with it (cite zoya and icip paper here). To what extent do gaze patterns correlate with memorability, and more generally, to what extent are visual saliency and memorability related?

In this paper, we systematically explore the memorability of objects within individual images and shed light on the various factors that drive object memorability. In exploring the connection between object memorability, saliency,

108 and image memorability, our paper makes several important
109 contributions.

110 Firstly, we show that just like image memorability, object
111 memorability is a property that is shared across subjects: objects remembered by one person are also likely to
112 be remembered by others and vice versa. Secondly, we
113 uncover the relationship between visual saliency and object
114 memorability, and demonstrate those instances where
115 visual saliency directly predicts object memorability and
116 when/why it fails. While there have been a few studies
117 that explore the connection between image memorability
118 and visual saliency [7], [28], our work is the first to ex-
119 plore the connection between object memorability and visual
120 saliency. Third, we make the initial leaps in disam-
121 biguating the link between image memorability and object
122 memorability, and show that in many cases, the memorability
123 of an image is primarily driven by the memorability
124 of its most memorable object. Studying these questions,
125 help not only understand visual saliency, image and object
126 memorability in more detail, but it can also have impor-
127 tant contributions to computer vision. For example, under-
128 standing which regions and objects in an image are memo-
129 rable would enable us to modify the memorability of images
130 which can have applications in advertising, user interface
131 design etc. With this in mind, we show in section 4 that
132 our proposed dataset can serve as a benchmark for eval-
133 uating object memorability algorithms and encourage future
134 object and region memorability prediction schemes. Taken
135 together, our efforts offer a deeper understanding of mem-
136 orability in general and broaden the contribution of human-
137 level insight to improve machine prediction

140 1.1. Related works

141 In this section, we briefly discuss existing work related
142 to visual memory and image memorability. We also review
143 research related to visual saliency prediction and discuss the
144 relationship of memorability and visual attention.

145 **Image Memorability:** Describe Isola’s first paper n
146 some insights that have been raised on image memorabil-
147 ity thus far. Also describe Khosla’s comp model but we are
148 the first work to actually describe what humans actually re-
149 member and don’t

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

153 **Saliency and memorability:** discuss some results re-
154 lated to saliency and image memorability.

155 and talk about our work plans on connecting and shed-
156 ding light on all these phenomena together.

157 2. Measuring Object Memorability

158 As a first step towards understanding memorability of
159 objects, we built an image database containing a variety of

160 objects from a diverse range of categories, and measured the
161 probability that every object in each image will be remem-
162 bered by a large group of subjects after a single viewing.
163 This helps provide ground truth memorability scores for the
164 objects inside the images and allows for a precise analysis
165 of the memorable elements within an image. For this task,
166 we utilized the PASCAL-S dataset [24], a fully segmented
167 dataset built on the validation set of the PASCAL VOC 2010
168 [11] segmentation challenge. For improved segmentation
169 purposes, we manually cleaned up and refined the segmen-
170 tations from this dataset. We removed all homogenous non-
171 object or background segments such as ground, grass, floor,
172 sky etc, as well as imperceptible object fragments and ex-
173 cessively blurred regions. All remaining object segmenta-
174 tions were tested for memorability. In the end, our final
175 dataset consisted of 850 images and 3412 object segmenta-
176 tions i.e. on average each image consisted of approximately
177 4 object segments for which we gathered the ground truth
178 memorability on.

179 2.1. Object Memory Game

180 To measure the memorability of individual objects from
181 our dataset, we created an alternate version of the Visual
182 Memory Game through Amazon Mechanical Turk follow-
183 ing the basic design in [16], with the exception of a few
184 key differences. In our game, participants first viewed a se-
185 quence of images one at a time, with a 1.5 second gap in be-
186 tween image presentations. Subjects were asked to remem-
187 ber the contents and objects inside those images as much as
188 they could. To ensure that subjects would not just only look
189 at the salient or center objects, subjects had unlimited time
190 to freely view the images. Once they were done viewing an
191 image, they could press any key to advance to the next im-
192 age. Following the initial image sequence, participants then
193 viewed a sequence of objects, their task then being to indi-
194 cate through a key press which of those objects was present
195 in one of the previously shown images. Each object was
196 displayed for 1.5 second, with a 1.5 second gap in between
197 the object sequences. Pairs of corresponding image and ob-
198 ject sequences were broken up into 10 blocks. Each block
199 consisted of 80 total stimuli (35 images and 45 objects), and
200 lasted approximately 3 minutes. At the end of each block,
201 the subject could take a short break. Overall, the experiment
202 took approximately took 30 minutes to complete.

203 Unknown to the subjects, inside each block, each se-
204 quence of images was pseudo-random and consisted of 3
205 ‘target’ images taken from the Pascal-S dataset whose ob-
206 jects the participants were to later identify. The remaining
207 images in the sequence consisted of 16 ‘filler’ images and
208 16 ‘familiar’ images. The ‘filler’ images were randomly se-
209 lected from the DUT-OMRON dataset [32] and the ‘famili-
210 ar’ images were randomly sampled from the MSRA dataset
211 proposed in [26]. Similarly, the object sequence was also

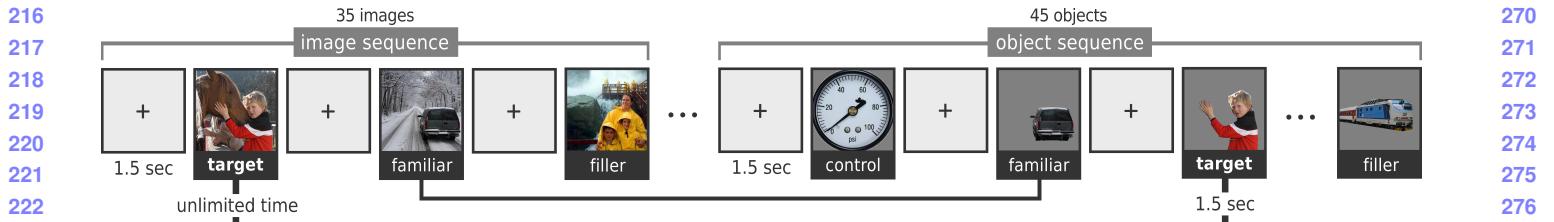


Figure 2: Main task. add-in later.

pseudo-random and consisted of 3 'target' objects (1 object taken randomly from each previously shown target image). The remaining objects in the sequence consisted of 10 'control' objects, 16 'filler' objects, and 16 'familiar' objects. The 'filler' objects were taken randomly from the 80 different object categories in the Microsoft COCO dataset [25] and the 'familiar' objects were the objects taken from the previously displayed 'familiar' images in the image sequence. The fillers and familiars helped provide spacing between the target images and target objects, whereas the control objects allowed us to check if the subjects were paying attention to the task [5], [16]. While the fillers and familiars (both the images and objects) were taken from datasets resembling real world scenes and objects, the 'control' objects were artificial stimuli randomly sampled from the dataset proposed in [5] and helped serve as a control to test the attentiveness of the subjects. The target images and the respective target objects were spaced 70 – 79 stimuli apart, and familiar images and their respective objects were spaced 1 – 79 stimuli apart. All images and objects appeared only once, and each subject was tested on only one object from each target image. Objects were centered within their parent frame and non-object pixels were set to grey. Participants were required to complete the entire task, which included 10 blocks (overall time approximately 30 minutes), and could not participate in the experiment a second time. After collecting the data, we assigned a 'memorability score' to each target object in our dataset, defined as the percentage of correct detections by subjects. In all our analysis, we removed all subjects whose accuracy on the control objects was below 70%. In the end, our analysis was performed on a total of 1823 workers from Mechanical Turk ($> 95\%$ approval rate in Amazons system). The memorability score of an object corresponded to the number of subjects that correctly detected the repetition of that object. On average, each object was scored by 16 subjects and the average memorability score was 33% ($SD = 28\%$).

2.2. Consistency Analysis

To assess human consistency in remembering objects, we repeatedly divided our entire subject pool into two equal halves and quantified the degree to which memorability scores for the two sets of subjects were in agreement using

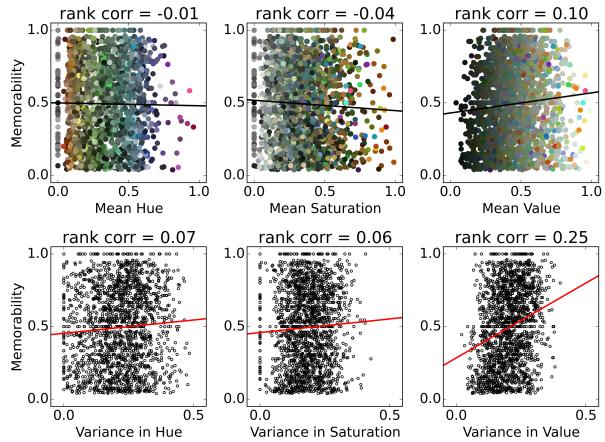


Figure 3: Correlations between simple color features and object memorability. Color features were computed from HSV representations of the objects.

Spearman's rank correlation (ρ). We computed the average correlation over 25 of these random split iterations, yielding a final value of 0.76. This high consistency in object memorability indicates that, like full images, object memorability is a shared property across subjects. People tend to remember (and forget) the same objects in images, and exhibit similar performance in doing so. Thus memorability of objects in images can potentially be predicted with high accuracy. In the next section, we study the various factors that possibly drive object memorability in images.

3. Understanding Object Memorability

In this section, we aim to better understand object memorability and the factors that make an object more memorable or forgettable to humans. We first investigate the role that simple color features play in determining object memorability.

3.1. Can simple features explain memorability?

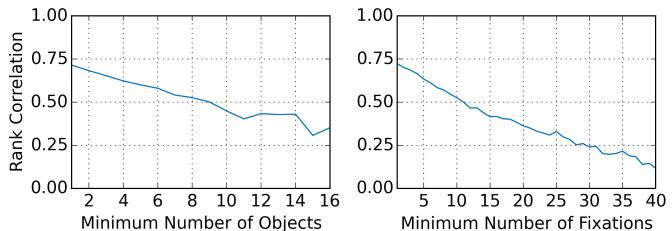
While simple image features are traditionally poor predictors of memorability in full images [16], and with good reason [21], do they play any role in determining object memorability? We decomposed each image into its hue, saturation, and value components and calculated the mean

324 and standard deviation of each channel. Mean value ($\rho = 0.1$) and variance in value ($\rho = 0.25$) were weakly correlated with object memorability suggesting that brighter and higher contrast objects may be more memorable (Figure 3).
 325 On the other hand, essentially no relationship was found
 326 between memorability and either hue or saturation (Figure
 327 3). This deviates slightly from the findings in [16] that
 328 showed mean hue to be weakly predictive of image mem-
 329 orability. However, this makes sense since the effect was
 330 speculated to be due to the blue and green outdoor land-
 331 scapes being less memorable than warmly colored human
 332 faces and indoor scenes. While our dataset contained plenty
 333 of indoor objects and people, outdoor scene-related image
 334 regions such as sky and ground were not included as ob-
 335 jects. Taken together, these results show that, like image
 336 memorability, basic pixel statistics do not play a significant
 337 role in determining the memorability of objects in images.
 338

342 3.2. What is the role of saliency in memorability?

343 Intuitively, the regions within an image that are most
 344 salient are likely to have a higher probability of being re-
 345 membered, since they will draw the attention of viewers and
 346 a majority of a viewers eye fixations will be spent looking
 347 at those regions.. On the other hand, it is conceivable that
 348 some visually appealing regions will not be memorable, es-
 349 pecially since aesthetic images are known to be less mem-
 350 orable [16], [15]. When can visual saliency predict object
 351 memorability and what are the possible differences between
 352 these two phenomena? Quantifying the precise relationship
 353 between saliency and memorability will be paramount to-
 354 wards understanding object memorability in greater depth.
 355

356 To this aim, we utilized the eye fixation dataset made
 357 available for the Pascal-S dataset in [24]. With this dataset
 358 in hand, we first calculated the number of unique fixation
 359 points within the area of each object and computed the cor-
 360 relation between this metric and the objects memorability
 361 score (Figure 5 a). We found this correlation to be positive
 362 and considerably high ($\rho = 0.71$), suggesting that fixation count
 363 and visual saliency may drive object memorability
 364 considerably. However, the large concentration of points
 365 on the bottom left part of scatter plot in Figure 5 a suggests
 366 that part of the reason for this high correlation is that ob-
 367 jects that have not been viewed at all have essentially no
 368 memorability. Indeed, only objects that have been seen can
 369 be remembered. In addition, the points toward the top left
 370 appear to decrease in trend. Looking deeper, Figure 4 plots
 371 the change in correlation between object memorability and
 372 fixations as the minimum number of fixations inside ob-
 373 jects increases. The downward monotonic trend indicates
 374 that as the number of fixations inside an object increases,
 375 the predictive ability diminishes significantly. In addition,
 376 Figure 4 plots the correlation between object memorability
 377 and number of fixations as a function of total number of



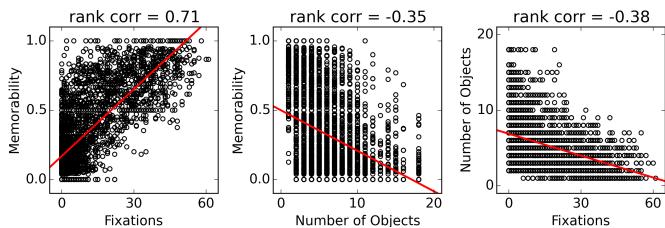
378
 379
 380
 381
 382
 383
 384
 385
 386
 387
 388
 389
 390
 391
 392
 393
 394
 395
 396
 397
 398
 399
 400
 401
 402
 403
 404
 405
 406
 407
 408
 409
 410
 411
 412
 413
 414
 415
 416
 417
 418
 419
 420
 421
 422
 423
 424
 425
 426
 427
 428
 429
 430
 431

Figure 4: Correlation between object memorability and number of fixations. add-in later.

objects in an image. Similar to the previous trend, as the number of objects in an image increases, the correlation between saliency i.e. number of fixations and memorability score decreases sharply. This finding is in agreement with the two remaining scatter plots in Figure 5 b (shows that the memorability of an object decreases in the presence of many other objects) and Figure 5 c (shows that number of fixations decreases with the number of objects). This makes intuitive sense since people have more to look at in an image when more objects are present, and so they may look less at any one object, especially if they compete for saliency, and therefore may have a more difficult time remembering those objects.

To sum up, saliency is a surprisingly good index of object memorability in simple contexts where there are few objects in the image, or when an object has few interesting points, but it is a much weaker predictor of object memorability in complex scenes containing multiple objects that have many points of interest (Figure 7).

Center Bias: Figure 6 elucidates another example where saliency and memorability diverge. Previous studies related to visual saliency have showed that saliency is heavily influenced by center bias [19], [33], primarily due to photographer bias (also evident from the leftmost plot in Figure 6) and viewing strategy [31]. Since our data collection experiment tries to control for the viewing strategy, memorability exhibits comparatively less center bias than saliency. This is most apparent when considering the difference in the solid ellipse in the right plot (shows where 95% of fixation positions are located), and the dashed ellipse (shows where the 95% of the above-median memo-



423
 424
 425
 426
 427
 428
 429
 430
 431

Figure 5: Correlation between object memorability and number of fixations. add-in later.

Figure 7: Saliency Fail cases. add-in later.

variable objects are located).

3.3. How do object categories affect memorability?

In the previous sections, we showed that simple features have little predictive power over object memorability and explored the relationship between visual saliency and object memorability. In this section, we explore how the category of an object influences the probability that the object will be remembered.

3.3.1 Are some classes more memorable than others?

For this analysis, we first assigned three in-house annotators the task of assigning class labels to each object segmentation in our dataset. The annotators were given the original image (for reference) and the object segmentation and asked to assign a single category to the segment out of 7 possible categories: animal, building, device, furniture, nature, person, and vehicle. We choose these high-level categories such that a wide range of object classes could be covered under these categories. For example, device included object segments such as utensils, bottles, televisions, computers etc, nature included segments like trees, mountains, flowers, and vehicle contained segments like cars, bikes, buses, airplanes etc.

Figure 9 shows the distribution of the memorability scores for all 7 object classes in our dataset. This visualisation gives a sense of how the memorability changes across different object categories. Animal, person, and vehicle are all highly memorable classes each associated with an average memorability score greater than or close to 0.5. Inter-

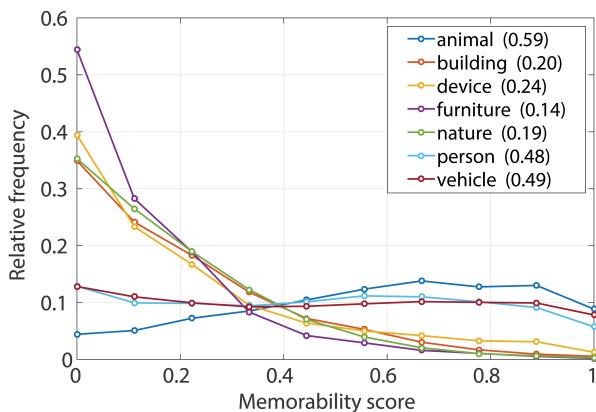


Figure 9: Average memorability per object class. Figure showing some object classes are more memorable than others.

estingly, all other object categories have an average memorability score lower than 0.25, indicating that humans do not remember objects from these categories very well. In particular, furniture is the least memorable object class with an average memorability score of only 0.14. This could be possibly due to the fact that most objects from classes like furniture, nature, and building either appear mostly in the background or are occluded which likely decreases their memorability significantly. By contrast, objects from the animal, person, and vehicle classes appear mostly in the foreground, leading to a higher memorability score on average. Interestingly, the topmost memorable objects from building, furniture, and nature tend to have an average memorability score in the range of 0.4 – 0.8, whereas the topmost memorable objects from classes person, animal and vehicle have an average memorability higher than 0.90. This is particularly interesting as these top objects are not occluded and most of them tend to appear in the foreground. While the differences in the memorability of different classes could be driven primarily due to factors like occlusion, size, background/foreground, or photographic bias, the distribution in figure 9 suggests that humans remember some object classes such as person, animal, and vehicle irrespective of external nuisance factors and these object classes are *intrinsically* more memorable than others.

3.3.2 Why are some objects not memorable in a class?

As demonstrated above, some object classes (i.e. animal, person, vehicle) are more memorable than others. However, not all objects in a class are equally memorable. The examples in Figure 8 show the most memorable, medium memorable, and least memorable objects for each object class. Across classes, non-memorable objects tend to be those that are occluded and obstructed by other objects. What other possible factors could influence the memorability of an object within a class? Among the various possible factors, we

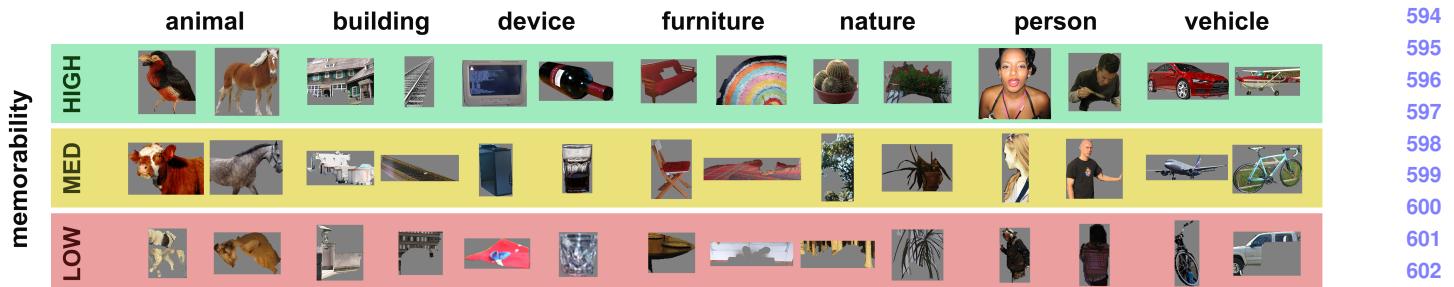


Figure 8: Qualitative results from object categories. add-in later.

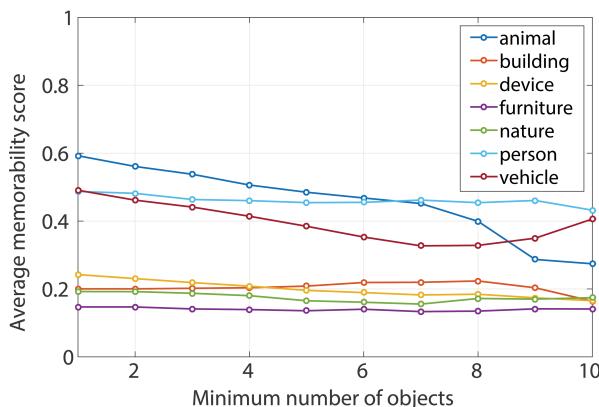


Figure 10: Correlation between object class and number of objects. add-in later.

explored how object memorability within a particular class is influenced by a) the number of objects in an image and b) the presence of other object classes.

Number of objects: We first examined how the memorability of each object class is affected by the number of objects inside an image. Figure 10 shows the change in average memorability of the different object classes with respect to the increase in number of objects in an image. Results indicate that the number of objects present in an image is an important factor in determining memorability. For example, as the number of objects in an image increases, the memorability of animals and vehicles decreases sharply, likely as a result of competition for attention, or decreased spotlight on a single subject of the composition. Interestingly, the memorability of the person class does not change significantly with an increase in number of objects. This suggests that people are not only one of the most memorable object classes, but are also more robust to the presence of clutter in images. This may be because single people in images steal all of the attention of the viewer, but how do they behave in the presence of other people? To answer this, we turn to the question of interclass memorability next.

Inter-class memorability: How does the presence of a particular object class influence the memorability of another object class? For all pairwise combinations of ob-

ject category, we gathered all images that contained at least one object from both categories and computed the change in the average memorability scores for the two object categories. Figure 11 plots these data and visualizes how the memorability of each object class is affected by the presence of other object classes. The first thing to note is that the values of low-memorability classes (i.e. nature, furniture, device, and building) are not greatly affected by the presence of other object categories. Instead, their memorability tends to remain low across all contexts. The memorability of the animal class remains close to its high average memorability score in presence of most classes, but drops significantly in the presence of other animals, vehicles, and people. The memorability of people also remains close to its average memorability score and tends to be unaffected by the presence of most object categories (including other people). However, the memorability of a person decreases in the presence of vehicles and buildings. This could be due to the fact that people in images containing vehicles or buildings are usually zoomed out and are usually smaller in size (also illustrated in figure 12). The memorability of the vehicle class is strongly affected by the presence of other object categories. In particular, its memorability drops significantly in the presence of another vehicle, people, and animals. Taken together, when an animal, vehicle or a person is present in the same image, the memorability of all three classes usually goes down. However, this pattern of change in memorability varies by class, leading to interesting results. For example, when a vehicle and animal are present in the same image, the animal is generally more memorable, even though the memorability of both of these classes drops significantly. When a vehicle or an animal co-occurs with a person, the person is generally more memorable (also shown in Figure 12).

3.4. How are object & image memorability related?

We now know what objects people remember and the factors that influence their memorability, but to what extent does the memorability of individual objects affect the overall memorability of a scene? Moreover, if an image is highly memorable, what can we say about the memorabil-

648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
ity of the objects inside those images (and vice versa)? To shed light on this question, we conducted a second large-scale experiment on Amazon Mechanical Turk for all images in our dataset to gather their respective image memorability scores. For this experiment, we followed the exact paradigm as the memory game experiment proposed in [16]. A series of images from our dataset and Microsoft COCO dataset [25] (i.e. the 'filler' images) were flashed for 1 second each, and subjects were instructed to press a key whenever they detected a repeat presentation of an image. A total of 350 workers participated in this experiment with each image being viewed 80 times on average. The rank correlation after averaging over 25 random split half trials was found to be 0.70, providing evidence for consistency in the image memorability scores.

Utilizing results from both experiments, we computed the correlation between the the scores of the single most memorable object in each image (from Experiment 1) and the overall memorability score of each image (from Experiment 2). We found this correlation to be moderately high ($\rho = 0.4$), suggesting that the most highly memorable object in the image plays a crucial role in determining the overall memorability of an image. To investigate this finding in relation to extreme cases only, we performed the same analysis as above on a subset of the data containing the topmost 100 memorable images and the bottommost 100 memorable images. The correlation between maximum object memorability and image memorability for this subset of the images increased significantly ($\rho = 0.62$), meaning maximum object memorability serves as a strong indicator of whether an image is *highly* memorable or non-memorable. That is, images that are highly memorable contain at least one highly memorable object, and images with low memorability usually do not contain a single highly memorable object (also shown in Figure 13).

It seems that maximum object memorability is highly explanatory, but does this behavior generalize across ob-

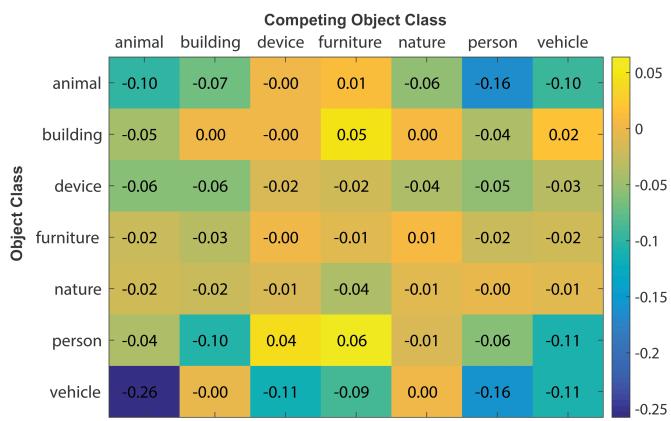


Figure 11: **inter-class object memorability relationship.** add-in later.



Figure 12: **Qual Results.** Figure showing how memorability of different classes is effected in presence of other classes. Bottom row is the memorability map

Animal	Building	Device	Furniture	Nature	Person	Vehicle	All
0.38	0.22	0.47	0.53	0.64	0.54	0.30	0.40

Table 1: **Max object memorability and image memorability.** add-in later.

ject classes? We further computed the correlation between maximum object and image memorability for each individual object class. Results shown in Table 1 show that certain object classes are more strongly correlated than others. For example, images containing animals, buildings, or vehicles as the most memorable objects tend to have varying degree of image memorability (indicated by their lower ρ values). On the other hand, classes like device, furniture, nature, and person are strongly correlated with image memorability, indicating that if an image's most memorable object belongs to one of these classes, the object memorability score is strongly predictive of the image memorability score. We can imagine scenarios in which this information would be potentially useful. For example, in the case where vision systems are tasked to predict scene memorability, a *single* object and its class can serve as a strong prior in predicting image memorability.

4. Predicting Object Memorability

Along with understanding what drives memorability of objects in a scene, our work also makes available the very first dataset containing the ground truth memorability of constituent objects from a highly diverse image set. In this section, we show that our dataset can be used to benchmark computational models and serve as a stepping stone in the direction of object memorability prediction.

Baseline models: As a first step, we propose a simple baseline model that utilizes a conv-net [22], [18] trained on the ImageNet database [10]. Since object categories play an important role in determining object memorability (??), and deep learning models have recently been shown to achieve state-of-the-art results in various recognition tasks, including object recognition and object categorization [12], [23], we believe that this simple model can serve as a good initial baseline for object memorability prediction. We first generated object segments by using MCG, a generic object proposal method proposed in [2]. Next, we trained an

756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

SVR using 6-fold cross-validation on the original segments to map deep features to memorability scores. We then used this model to predict the memorability scores for the top K ($K = 20$) segments (obtained via the ranking scores provided by the MCG algorithm) for each image. After obtaining the predicted memorability scores, the memorability maps were generated by averaging the top K segments at the pixel level. Since image features like SIFT [27] and HOG [9] have previously been shown to achieve good performance in predicting image memorability [16], [15], we built a second baseline model using these features for comparison. Training and testing of this model was performed similar to the deep-net baseline model.

Evaluation: To evaluate the accuracy of the predicted memorability maps, we computed the rank correlation between the mean predicted memorability score inside each of the original object segments and their ground truth memorability scores. From figure 14, we first note that our deep-net baseline model, DL-MCG performs considerably well ($\rho = 0.39$). In contrast, the baseline model trained using HOG and SIFT features, H+S exhibits much lower overall performance ($\rho = 0.27$). Saliency maps generated from saliency algorithms are also likely to have some degree of overlap with memorability and are therefore worth comparing to our baseline, especially given the absence of alternative memorability prediction methods¹. Thus, we also included 8 state-of-the-art-saliency methods GB [13], AIM [6], DV [14], IT [17], GC [8], PC [29], SF [30], and FT [1] to our comparison (some of the top performing methods according to benchmarks in [4], [3]). Results from figure 14) show that the H+S baseline is outperformed by most saliency methods. Thus, even though models using SIFT and HOG have previously demonstrated high predictive power for image memorability, they may not be as well suited for the task of predicting object memorability. The deep-net baseline model, DL-MCG performs better than all other saliency methods and only PC ($\rho = 0.38$), SF ($\rho = 0.37$), and GB ($\rho = 0.36$) show performance comparable to the model. A common factor between these saliency methods is that they explicitly add center bias to their implementation. Even though memorability exhibits lesser center bias when compared to eye fixations, it still tends to be biased slightly towards the center due to photographer bias (section 3.2), which could be a part of the reason for the high performance of these methods. Despite this, DL-MCG performs favorably against them and is potentially much better suited for memorability prediction on a wide range of datasets. Thus, we recommend in the future, memorability algorithms compare their methods against our

¹The only other algorithm that generates memorability maps was proposed in [20]. We contacted the authors and they said they will be releasing an updated version of their paper and codes soon. We will add it to the comparison once they release the code.

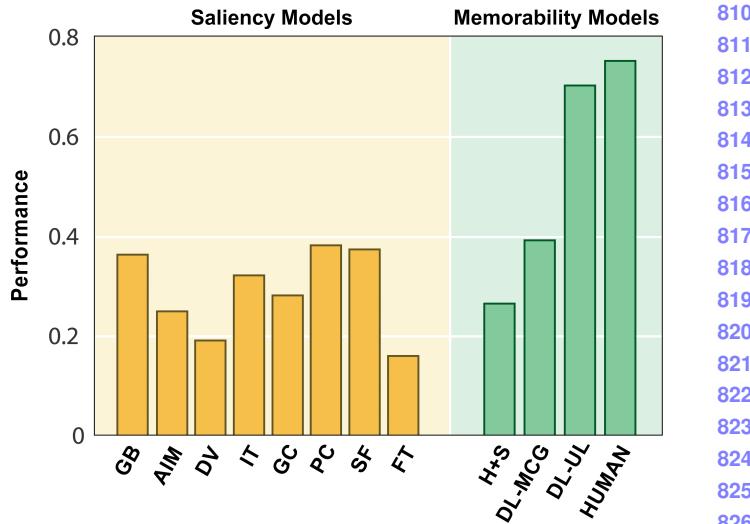


Figure 14: **Main task.** add-in later.

DL-MCG baseline. While DL-MCG performed fairly well, part of the performance of this model is dependent on the quality of the segmentations used. For this reason, we also consider the upper bound of our current predictive power by showing the results for our model containing predictions on the original segments (referred to as DL-UL in Figure 14). Interestingly, the accuracy of this model is very high and close to human performance ($\rho = 0.7$). This demonstrates that the deep-net model has high predictive ability that is suppressed most heavily by constraints of the segmentation task. The main insight of our evaluation is that deep features serve as strong predictors of memorability and selection of higher quality segments can potentially lead to improved memorability prediction algorithms.

References

- [1] R. Achanta, S. Hemami, F. Estrada, and S. Sussstrunk. Frequency-tuned salient region detection. In *Computer vision and pattern recognition, 2009. cvpr 2009. ieee conference on*, pages 1597–1604. IEEE, 2009. 8
- [2] P. Arbelaez, J. Pont-Tuset, J. Barron, F. Marques, and J. Malik. Multiscale combinatorial grouping. In *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*, pages 328–335. IEEE, 2014. 7
- [3] A. Borji, D. N. Sihite, and L. Itti. Salient object detection: A benchmark. In *Computer Vision–ECCV 2012*, pages 414–429. Springer, 2012. 8
- [4] A. Borji, D. N. Sihite, and L. Itti. Quantitative analysis of human-model agreement in visual saliency modeling: A comparative study. *Image Processing, IEEE Transactions on*, 22(1):55–69, 2013. 8
- [5] T. F. Brady, T. Konkle, G. A. Alvarez, and A. Oliva. Visual long-term memory has a massive storage capacity for object details. *Proceedings of the National Academy of Sciences*, 105(38):14325–14329, 2008. 3

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

- 864 [6] N. Bruce and J. Tsotsos. Saliency based on information maximization. In *Advances in neural information processing systems*, pages 155–162, 2005. 8
- 865 [7] Z. Bylinskii, P. Isola, C. Bainbridge, A. Torralba, and A. Oliva. Intrinsic and extrinsic effects on image memorability. *Vision research*, 2015. 2
- 866 [8] M.-M. Cheng, G.-X. Zhang, N. J. Mitra, X. Huang, and S.-M. Hu. Global contrast based salient region detection. In *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*, pages 409–416. IEEE, 2011. 8
- 867 [9] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*, volume 1, pages 886–893. IEEE, 2005. 8
- 868 [10] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pages 248–255. IEEE, 2009. 7
- 869 [11] M. Everingham and J. Winn. The pascal visual object classes challenge 2010 (voc2010) development kit, 2010. 2
- 870 [12] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*, pages 580–587. IEEE, 2014. 7
- 871 [13] J. Harel, C. Koch, and P. Perona. Graph-based visual saliency. In *Advances in neural information processing systems*, pages 545–552, 2006. 8
- 872 [14] X. Hou and L. Zhang. Dynamic visual attention: Searching for coding length increments. In *Advances in neural information processing systems*, pages 681–688, 2009. 8
- 873 [15] P. Isola, J. Xiao, D. Parikh, A. Torralba, and A. Oliva. What makes a photograph memorable? *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 36(7):1469–1482, 2014. 4, 8
- 874 [16] P. Isola, J. Xiao, A. Torralba, and A. Oliva. What makes an image memorable? In *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*, pages 145–152. IEEE, 2011. 2, 3, 4, 7, 8
- 875 [17] L. Itti, C. Koch, and E. Niebur. A model of saliency-based visual attention for rapid scene analysis. *IEEE Transactions on pattern analysis and machine intelligence*, 20(11):1254–1259, 1998. 8
- 876 [18] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. In *Proceedings of the ACM International Conference on Multimedia*, pages 675–678. ACM, 2014. 7
- 877 [19] T. Judd, K. Ehinger, F. Durand, and A. Torralba. Learning to predict where humans look. In *Computer Vision, 2009 IEEE 12th international conference on*, pages 2106–2113. IEEE, 2009. 4
- 878 [20] A. Khosla, J. Xiao, A. Torralba, and A. Oliva. Memorability of image regions. In *Advances in Neural Information Processing Systems*, pages 305–313, 2012. 8
- 879 [21] T. Konkle, T. F. Brady, G. A. Alvarez, and A. Oliva. Conceptual distinctiveness supports detailed visual long-term mem-
880 ory for real-world objects. *Journal of Experimental Psychology: General*, 139(3):558, 2010. 3
- 881 [22] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012. 7
- 882 [23] H. Lee, R. Grosse, R. Ranganath, and A. Y. Ng. Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pages 609–616. ACM, 2009. 7
- 883 [24] Y. Li, X. Hou, C. Koch, J. M. Rehg, and A. L. Yuille. The secrets of salient object segmentation. In *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*, pages 280–287. IEEE, 2014. 2, 4
- 884 [25] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014*, pages 740–755. Springer, 2014. 3, 7
- 885 [26] T. Liu, Z. Yuan, J. Sun, J. Wang, N. Zheng, X. Tang, and H.-Y. Shum. Learning to detect a salient object. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 33(2):353–367, 2011. 2
- 886 [27] D. G. Lowe. Distinctive image features from scale-invariant keypoints. *International journal of computer vision*, 60(2):91–110, 2004. 8
- 887 [28] M. Mancas and O. Le Meur. Memorability of natural scenes: the role of attention. In *ICIP*, 2013. 2
- 888 [29] R. Margolin, A. Tal, and L. Zelnik-Manor. What makes a patch distinct? In *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*, pages 1139–1146. IEEE, 2013. 8
- 889 [30] F. Perazzi, P. Krahnenbuhl, Y. Pritch, and A. Hornung. Saliency filters: Contrast based filtering for salient region detection. In *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, pages 733–740. IEEE, 2012. 8
- 890 [31] P.-H. Tseng, R. Carmi, I. G. Cameron, D. P. Munoz, and L. Itti. Quantifying center bias of observers in free viewing of dynamic natural scenes. *Journal of vision*, 9(7):4, 2009. 4
- 891 [32] C. Yang, L. Zhang, H. Lu, X. Ruan, and M.-H. Yang. Saliency detection via graph-based manifold ranking. In *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*, pages 3166–3173. IEEE, 2013. 2
- 892 [33] L. Zhang, M. H. Tong, T. K. Marks, H. Shan, and G. W. Cottrell. Sun: A bayesian framework for saliency using natural statistics. *Journal of vision*, 8(7):32, 2008. 4
- 893 [34] L. Zhang, M. H. Tong, T. K. Marks, H. Shan, and G. W. Cottrell. Sun: A bayesian framework for saliency using natural statistics. *Journal of vision*, 8(7):32, 2008. 4
- 894 [35] L. Zhang, M. H. Tong, T. K. Marks, H. Shan, and G. W. Cottrell. Sun: A bayesian framework for saliency using natural statistics. *Journal of vision*, 8(7):32, 2008. 4
- 895 [36] L. Zhang, M. H. Tong, T. K. Marks, H. Shan, and G. W. Cottrell. Sun: A bayesian framework for saliency using natural statistics. *Journal of vision*, 8(7):32, 2008. 4
- 896 [37] L. Zhang, M. H. Tong, T. K. Marks, H. Shan, and G. W. Cottrell. Sun: A bayesian framework for saliency using natural statistics. *Journal of vision*, 8(7):32, 2008. 4
- 897 [38] L. Zhang, M. H. Tong, T. K. Marks, H. Shan, and G. W. Cottrell. Sun: A bayesian framework for saliency using natural statistics. *Journal of vision*, 8(7):32, 2008. 4
- 898 [39] L. Zhang, M. H. Tong, T. K. Marks, H. Shan, and G. W. Cottrell. Sun: A bayesian framework for saliency using natural statistics. *Journal of vision*, 8(7):32, 2008. 4
- 899 [40] L. Zhang, M. H. Tong, T. K. Marks, H. Shan, and G. W. Cottrell. Sun: A bayesian framework for saliency using natural statistics. *Journal of vision*, 8(7):32, 2008. 4
- 900 [41] L. Zhang, M. H. Tong, T. K. Marks, H. Shan, and G. W. Cottrell. Sun: A bayesian framework for saliency using natural statistics. *Journal of vision*, 8(7):32, 2008. 4
- 901 [42] L. Zhang, M. H. Tong, T. K. Marks, H. Shan, and G. W. Cottrell. Sun: A bayesian framework for saliency using natural statistics. *Journal of vision*, 8(7):32, 2008. 4
- 902 [43] L. Zhang, M. H. Tong, T. K. Marks, H. Shan, and G. W. Cottrell. Sun: A bayesian framework for saliency using natural statistics. *Journal of vision*, 8(7):32, 2008. 4
- 903 [44] L. Zhang, M. H. Tong, T. K. Marks, H. Shan, and G. W. Cottrell. Sun: A bayesian framework for saliency using natural statistics. *Journal of vision*, 8(7):32, 2008. 4
- 904 [45] L. Zhang, M. H. Tong, T. K. Marks, H. Shan, and G. W. Cottrell. Sun: A bayesian framework for saliency using natural statistics. *Journal of vision*, 8(7):32, 2008. 4
- 905 [46] L. Zhang, M. H. Tong, T. K. Marks, H. Shan, and G. W. Cottrell. Sun: A bayesian framework for saliency using natural statistics. *Journal of vision*, 8(7):32, 2008. 4
- 906 [47] L. Zhang, M. H. Tong, T. K. Marks, H. Shan, and G. W. Cottrell. Sun: A bayesian framework for saliency using natural statistics. *Journal of vision*, 8(7):32, 2008. 4
- 907 [48] L. Zhang, M. H. Tong, T. K. Marks, H. Shan, and G. W. Cottrell. Sun: A bayesian framework for saliency using natural statistics. *Journal of vision*, 8(7):32, 2008. 4
- 908 [49] L. Zhang, M. H. Tong, T. K. Marks, H. Shan, and G. W. Cottrell. Sun: A bayesian framework for saliency using natural statistics. *Journal of vision*, 8(7):32, 2008. 4
- 909 [50] L. Zhang, M. H. Tong, T. K. Marks, H. Shan, and G. W. Cottrell. Sun: A bayesian framework for saliency using natural statistics. *Journal of vision*, 8(7):32, 2008. 4
- 910 [51] L. Zhang, M. H. Tong, T. K. Marks, H. Shan, and G. W. Cottrell. Sun: A bayesian framework for saliency using natural statistics. *Journal of vision*, 8(7):32, 2008. 4
- 911 [52] L. Zhang, M. H. Tong, T. K. Marks, H. Shan, and G. W. Cottrell. Sun: A bayesian framework for saliency using natural statistics. *Journal of vision*, 8(7):32, 2008. 4
- 912 [53] L. Zhang, M. H. Tong, T. K. Marks, H. Shan, and G. W. Cottrell. Sun: A bayesian framework for saliency using natural statistics. *Journal of vision*, 8(7):32, 2008. 4
- 913 [54] L. Zhang, M. H. Tong, T. K. Marks, H. Shan, and G. W. Cottrell. Sun: A bayesian framework for saliency using natural statistics. *Journal of vision*, 8(7):32, 2008. 4
- 914 [55] L. Zhang, M. H. Tong, T. K. Marks, H. Shan, and G. W. Cottrell. Sun: A bayesian framework for saliency using natural statistics. *Journal of vision*, 8(7):32, 2008. 4
- 915 [56] L. Zhang, M. H. Tong, T. K. Marks, H. Shan, and G. W. Cottrell. Sun: A bayesian framework for saliency using natural statistics. *Journal of vision*, 8(7):32, 2008. 4
- 916 [57] L. Zhang, M. H. Tong, T. K. Marks, H. Shan, and G. W. Cottrell. Sun: A bayesian framework for saliency using natural statistics. *Journal of vision*, 8(7):32, 2008. 4
- 917 [58] L. Zhang, M. H. Tong, T. K. Marks, H. Shan, and G. W. Cottrell. Sun: A bayesian framework for saliency using natural statistics. *Journal of vision*, 8(7):32, 2008. 4

972													1026
973													1027
974													1028
975													1029
976													1030
977													1031
978													1032
979													1033
980													1034
981													1035
982													1036
983													1037
984													1038
985													1039
986													1040
987													1041
988													1042
989													1043
990													1044
991													1045
992													1046
993													1047
994		0.94	1.00	0.93	0.75	0.92	0.80	0.92	0.94	0.92	0.80	0.91	0.76
995	Most memorable												
996													
997													
998													
999	Least memorable	0.47	0.56	0.48	0.38	0.50	0.38	0.50	0.61	0.51	0.50	0.52	0.39
1000													
1001													
1002													
1003													
1004													
1005													
1006													
1007													
1008													
1009													
1010													
1011													
1012													
1013													
1014													
1015													
1016													
1017													
1018													
1019													
1020													
1021													
1022													
1023													
1024													
1025													

Figure 13: Qual image-object results. add-in later.