

000
001
002
003
004
005
006
007
008
009
010
011054
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

Salicon: Saliency in Context

Anonymous CVPR submission

Paper ID 234

Abstract

This paper presents a new method to collect large-scale human data during natural explorations on images. While current datasets present a rich set of images and task-specific annotations such as category labels and object segments, this work focuses on recording and logging how humans shift their attention during visual exploration. The goal is to offer new possibilities to (1) complement task-specific annotations to advance the ultimate goal in visual understanding, and (2) understand visual attention and learn saliency models, all with human attentional data at a much larger scale.

We designed a mouse-contingent multi-resolutinal paradigm based on neurophysiological and psychophysical studies of peripheral vision, to simulate the natural viewing behavior of humans. The new paradigm allowed using a general-purpose mouse instead of an eye tracker to record viewing behaviors, thus enabling large-scale data collection. The paradigm was validated with controlled laboratory as well as large-scale online data. We report in this paper a proof-of-concept Saliency in Context (Salicon) dataset of human “free-viewing” data on 10,000 images from the Microsoft COCO (MS COCO) dataset with rich contextual information. Finally we evaluate the use of the collected data in the context of saliency prediction, and demonstrate them a good source as ground truth for the evaluation of saliency algorithms.

1. Introduction

Motivation One of the ultimate goals in computer vision is to describe the contents of an image. Humans are known to perform better than their machine counterparts in telling a story from an image, and we aim to leverage human intelligence and computer vision algorithms to bridge the gap between humans and computers in visual understanding.

In the recent years, there are several datasets with unprecedented numbers of images and annotations [31, 5, 33, 18], which enable breakthroughs in visual scene understanding, especially goal-specific tasks like object clas-

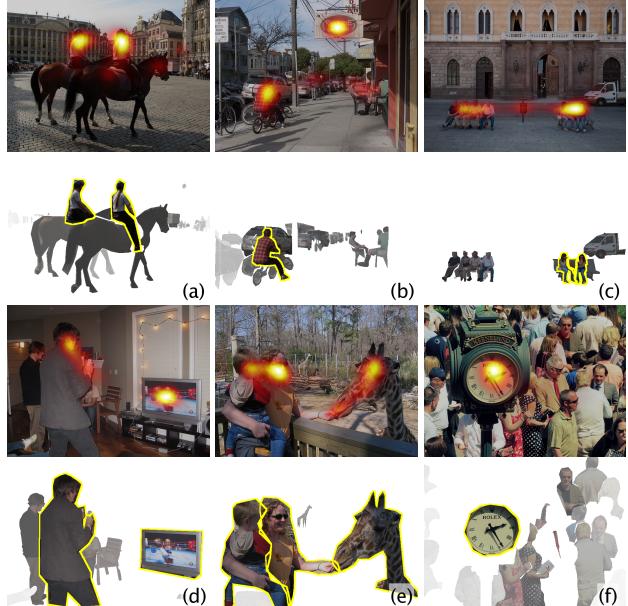


Figure 1. Contextual information is crucial in image understanding (image examples from MS COCO). We propose a new method to collect large-scale attentional data (Saliency in Context, 1st row) to provide additional data in visual understanding. For example, with annotated object segments, it naturally highlights key components in an image (ranked object segments in the 2nd row, with key objects outlined in yellow) to (a) rank object categories, (b) suggest new categories important to characterize a scene (text in this example), (c-e) convey social cues, and (f) direct to places designed for attention in advertisement.

sification and segmentation. In the recently published MS COCO dataset [18], non-iconic images and objects in context are emphasized to understand natural scenes. On top of annotations for the conventional computer vision tasks, it also includes sentences to describe an image, a big step toward the Turning test in the visual domain.

Complementary to all the existing big datasets, in this work we focus on how people direct their gaze when inspecting a visual scene. Humans and other primates shift their gaze to allocate processing resources to the most important subset of the visual input. Understanding and em-

108 ulating the way that human observers free-view a natural
109 scene to respond rapidly and adaptively has both scientific
110 and economic impact. The logging of human viewing data
111 during the assumption-free exploration also provides in-
112 sights to other vision tasks and complement them to better
113 understand and describe image contents (see Figure 1). For
114 example, it naturally ranks labeled object categories, and
115 suggests new categories for current classification datasets.
116 By highlighting important objects by humans, it leverages
117 human intelligence in visual understanding.

118 To collect large-scale human behavioral data in scene ex-
119 ploration, we first propose a novel psychophysical paradigm
120 to record mouse-movement data that mimic the ways hu-
121 mans view a scene [30]. The designed stimuli encode the
122 visual acuity drop-off as a function of retinal eccentricity.
123 The mouse-contingent paradigm motivates mouse move-
124 ments, to reveal interesting objects in the periphery with
125 high resolution, similarly as humans shift their gazes to
126 bring objects-of-interest to the fovea. Rather than record-
127 ing the task-specific end outcomes by human annotators,
128 we record the natural viewing patterns during the explo-
129 ration. Therefore, our method is general and task-free. We
130 then propose a crowdsourcing mechanism to collect large-
131 scale mouse-tracking data through Amazon Mechanic Turk
132 (AMT).

133 **Challenges** To record where humans look, eye-tracking
134 experiments are commonly conducted, where subjects sit
135 still in front of a screen with their eye movements recorded
136 by a camera. Normally an infrared illuminator is necessary
137 to help acquire high-quality data. There are several chal-
138 lenges particular to data collecting and usage.

139 First, large-scale data collection is prohibitive. An
140 eye tracker used in laboratories generally costs between
141 \$30,000 - \$80,000. Despite recent advances in gaze and eye
142 modeling and detection (e.g., [9]), accurate eye-tracking ex-
143 periments are still difficult without customized eye-tracking
144 hardware. Data collection with general-purpose webcams
145 is not yet possible, especially in uncontrolled settings such
146 as through the AMT platform. This greatly limits the data
147 collection process. As a result, the sizes of the current eye-
148 tracking datasets are at the order of hundreds images and
149 tens subjects, much smaller than those for object detection,
150 object classification, scene categorization, or segmentation.

151 Second, eye-tracking data are not sufficiently general.
152 Datasets collected from different labs are quite different in
153 nature due to various image selection criteria, experimental
154 setup, and instructions. Thus datasets cannot be directly
155 combined, nor models learned from one dataset directly
156 generalize to another [37].

157 **Objectives** This paper focuses on two major objectives:

158 1. We propose a novel psychophysical paradigm as an
159 alternative to eye-tracking, to provide approximation of
160 human gaze in natural exploration. We design a gaze-

162 contingent multi-resolutional mechanism where subjects
163 can move the mouse to direct the high-resolutional fovea to
164 where they find interesting in the image stimuli. The mouse
165 trajectories from multiple subjects are aggregated to indi-
166 cate where people look most in the images.

167 2. We propose a crowdsourcing platform to collect large-
168 scale mouse-tracking data. We first sample 10,000 images
169 from the MS COCO dataset with rich contextual informa-
170 tion, and collect mouse-movement data using AMT. The
171 “free-viewing” dataset is by far the largest one in both scale
172 and context variability. We would like to point out that, with
173 the crowdsourcing platform, it allows us to easily collect
174 and compare various data with different top-down instruc-
175 tions, for example, to investigate the attention shifts during
176 story-telling vs. category labeling.

2. Related work

177 **Eye-tracking datasets** There is a growing interest in the
178 cognitive science and computer science disciplines to un-
179 derstand how humans and other animals shift their gazes
180 to interact with the complex visual scenes. Several eye-
181 tracking datasets have been recently constructed and shared
182 in the community to understand visual attention and to build
183 computational saliency models.

184 An eye-tracking dataset includes natural images as the
185 visual stimuli and eye movement data recorded using eye-
186 tracking devices. A typical dataset contains hundreds or a
187 thousand images, viewed by tens of subjects while the lo-
188 cations of their eyes in image coordinates are tracked over
189 time. Even if POET, the largest dataset we know by far, con-
190 tains 6,270 images and is only viewed by 5 subjects [20].
191 While instructions are known to affect eye movement pat-
192 terns, most common in eye-tracking dataset is the use of a
193 so-called “free-viewing” task [4, 17, 3, 26] due to its task-
194 free nature.

195 Most datasets have their own distinguishing features in
196 image selection. For example, most images in the FIFA
197 dataset [4] contain faces, and the NUSeF dataset [26] fo-
198 cuses on semantically affective objects/scenes. Compared
199 with FIFA and NUSeF, the widely used Toronto dataset
200 has less noticeably salient objects in the scenes. The MIT
201 dataset [17] is more general due to its relatively large size,
202 *i.e.*, 1003 images, and the generality of the image source.
203 Quite a few images in these datasets are with dominant ob-
204 jects in the center. To facilitate object and semantic saliency,
205 the OSIE dataset [34] features in multiple dominant ob-
206 jects in an image. Besides general purpose images, there
207 are also recent datasets in focused domains like the MIT
208 Low Resolution dataset [16] for saliency in low resolution,
209 EyeCrowd [15] for saliency in crowd, and FiWI [28] for
210 web page saliency. Human labeling such as object bounding
211 boxes [15], contours [26, 34], and social attributes [34, 15]
212 are available in certain datasets as ground truth data for
213 214 215

216
217
218
219
220
221
222
223
224
225

learning and analysis of problems of interest.

The scale of the current datasets is inherently limited by the experimental requirements. We envision that the collection of a larger-scale eye-tracking dataset would not only improve saliency prediction with big ground truth data, but driving new research directions in visual attention studies as well as complementing current efforts in computer vision datasets and annotations for more ambitious tasks in visual understanding.

Crowdsourcing Manual labeling to obtain ground truth human data is important for computer vision applications. Human knowledge and experience in this way is leveraged to train better computer models. Services like Amazon Mechanical Turk (AMT) has been extensively used to distribute the labeling task to many people, allowing the collection of large-scale labeling data. Recent works [31, 5, 32, 33, 6, 18] mainly focused on crowdsourcing image classification, object detection, and segmentation using this marketplace. Some of the most successful datasets along the line include Tiny Images [31], ImageNet [5] SUN [33], and MS COCO [18]. These datasets include hundreds thousands to millions of images containing hundreds or thousands of categories of interest, aiming at capturing general objects, scenes, or context in the visual world.

Current crowdsourcing tasks focus on the end output from humans (e.g., a category label, an object segment), while our method records the procedure during which humans explore the scene in a real-time manner. We expect that the viewing patterns reveal cognitive process and can be leveraged for intelligent visual understanding. Our current experiments use task-free scenarios, and it could work with any other task-specific annotation procedure to log how humans explore the scene to complete a certain task.

Mouse tracking Mouse tracking and eye-mouse coordination have been studied in the human-computer interaction literature. For example, one of the most popular application of mouse-tracking is web page analysis [12, 19]. Huang *et al.* [12] studied mouse behaviors in web searching tasks, suggesting the plausibility of using mouse positions to predict user behavior and gaze positions. Navalpakkam *et al.* [19] integrated the mouse position on web pages with task relevance, and developed computational models to predict eye movement from mouse activity. Web pages contain domain-specific contents that motivate users to move their mouse to click links and to navigate. In natural images, however, to motivate users to move their mouse as one shifts attention requires specific design of the visual stimuli.

3. Mouse-contingent free-viewing paradigm

To verify the feasibility of replacing eye-tracking data collection with mouse-tracking, and to investigate the correlations between the two modalities, we designed a novel mouse-contingent paradigm with multi-resolitional images

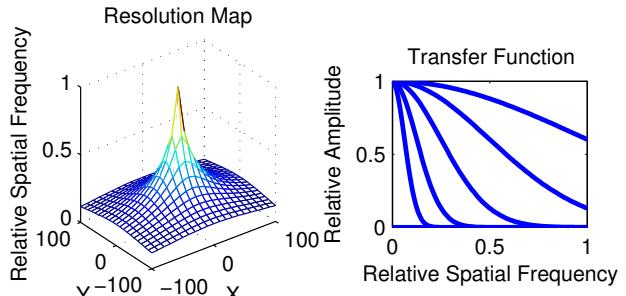
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

Figure 2. The resolution map and transfer functions.

generated in real-time. The dataset we used in this project was OSIE dataset which contains 700 images with the resolution of 800×600 . We collected the mouse-tracking data in a controlled laboratory environment, with similar hardware and software configurations as reported in [34].

3.1. Stimuli

To simulate the free-viewing patterns of human visual attention with mouse-tracking, we created an interactive paradigm by producing multi-resolitional images in real-time, based on the simulation method proposed by Perry and Geisler [23]. Gaze-contingent and mouse-contingent stimuli have been used in a variety of psychophysical studies, such as reading [27] and visual search [25]. The production of multi-resolitional images is based on neurophysiological and psychophysical studies of peripheral vision. Human visual system shows a well-defined contrast sensitivity by retinal eccentricity relationship. Specifically, contrast sensitivity to higher spatial frequencies drops off as a function of retinal eccentricity (e.g., [21, 24]). Therefore, we first generated a resolution map to simulate the sensitivity drop-off in peripheral vision (see Figure 2). [13] It is defined as a function $R : \Theta \rightarrow [0, 1]$, where Θ is the set of viewing angles θ with respect to the retinal eccentricity, and $[0, 1]$ represents the set of relative spatial frequency. The resolution map approximates a normal adult's vision with the exclusion of the blind spot. A higher $R(\theta)$ indicates a higher resolution at the visual eccentricity θ . Specifically, the resolution map is formulated as

$$R(x, y) = \frac{\alpha}{\alpha + \theta(x, y)}, \quad (1)$$

where $\alpha = 2.5^\circ$ is the half-height angle, which means that when $\theta(x, y) = \alpha$ the image will become only half the resolution of the center of attention ($\theta(x, y) = 0$). In our experiments, we set $\alpha = 2.5$ to approximate the actual acuity of human retina. The image coordinates were mapped to the visual angles by the following function:

$$\theta(x, y) = \frac{1}{p} \sqrt{(x - x_c)^2 + (y - y_c)^2}, \quad (2)$$

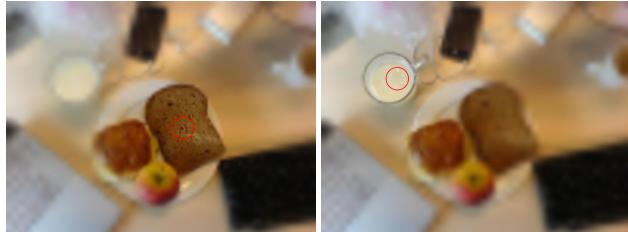
324
325
326
327
328
329
330
331

Figure 3. An example of the mouse-contingent stimuli. The red circles indicate the movement of mouse cursor from one object to another.

where θ is the visual angle, x and y are pixel coordinates, and (x_c, y_c) is the center of attention. The parameter p represents the number of pixels a person can see in a degree of visual angle, which can be changed to simulate different viewing distances. Generally, the closer the distance is, the less can be seen in the high-resolutonal fovea. We found that $p = 7.5$ led to a more comfortable and natural experience, according to the subjects' performances and feedbacks in pilot experiments. An example of the produced multi-resolutonal image is shown in Figure 3. To compute the multi-resolutonal image in real-time, we applied a fast approximation with a 6-level Gaussian pyramid from A_1 to A_6 . A_1 was the original image and A_i was down-sampled to A_{i+1} with a factor of 2 in both dimensions. The standard deviation of the Gaussian distribution was set to $\sigma = 0.248$ pixel. After that, all the down-sampled images (A_2 to A_6) were then interpolated to the original image size. We then computed six matrices of blending coefficients, $M_1 \dots M_6$. We used transfer function $T(f)$ (see Function 3 and Figure 2) and blending function $B(x, y)$ (see Equation 1 in [23]) to calculate these blending coefficients. The transfer function maps relative spatial frequency $f = R(x, y)$ to relative amplitude $T(f)$ in the Gaussian pyramid:

$$T_i(f) = \begin{cases} e^{1/2 \times (-2^{i-3} f / \sigma)^2}, & i = 1, \dots, 5 \\ 0, & i = 6, \end{cases} \quad (3)$$

The blending function $B(x, y)$ calculates the blending coefficients of each pixel (x, y) :

$$B(x, y) = \frac{0.5 - T_i(x, y)}{T_{i-1}(x, y) - T_i(x, y)}, \quad (4)$$

where i is the layer number of (x, y) . To calculate the layer number, we first determined six bandwidths $w_i, i = 1 \dots 6$ such that $T_i(w_i) = 0.5, i = 1 \dots 5$ and $w_6 = 0$. Then we normalized all w_i to $[0, 1]$. The layer number of pixel (x, y) is i such that $w_{i-1} \geq R(x, y) \geq w_i$. Next we calculated entries of $M_1 \dots M_6$. For each pair of indices (x, y) , we considered it as a pair of coordinates of a pixel and we

calculated its layer number i_0 , then

$$M_i(x, y) = \begin{cases} B(x, y), & i = i_0 - 1 \\ 1 - B(x, y), & i = i_0 \\ 0, & \text{otherwise} \end{cases} \quad (5)$$

for $i = 1 \dots 6$. Finally, the multi-resolutonal stimulus was a linear combination of M_i and A_i for $i = 1 \dots 6$.

3.2. Subjects and procedure

Sixteen subjects (10 male and 6 female) aged between 19 and 28 participated in the mouse-tracking experiment. All participants had normal or corrected-to-normal vision, and normal color vision as assessed by Ishihara plates. All subjects had not participated in any eye-tracking experiment or seen the OSIE images before. The images were presented to the subjects in 700 trials at random order. Each trial consists of a 5-second image presentation followed by a 2-second waiting interval. The mouse cursor was displayed as a red circle with a radius of 2 degrees of visual field that is sufficiently large not to block the high-resolutonal region of focus, and automatically moved to the image center when the image onset. The subjects were instructed to explore the image freely by moving the mouse cursor to anywhere they wanted to look. No further instructions were given on how to move the mouse or where they should look in the images. Whenever they moved the mouse, the mouse-contingent stimuli was updated by shifting the center of the resolution map to the mouse position. In the meantime, the mouse position and the timestamp were recorded. Each block contains 50 trials, and the subject can take a short break between blocks.

Presentation of stimuli and recording of mouse position were implemented in Matlab® (Mathworks, Natick, MA) using the Psychophysics Toolbox [1, 22] under Linux 14.04 LTS. The experiment PC was a Dell Precision™ T5610 with Intel® Xeon® E5-2609@2.5GHz (CPU), 32GB RAM and Nvidia® Quadro® K600 (GPU). The mouse speed and acceleration were adjusted to the maximum in the system settings. There was a practice session for the subjects to get familiar with the mouse-contingent paradigm and the mouse configuration, which consists of 10 other images from the Internet with the same resolution as the OSIE images. The practice trials were identical to the formal trials in terms of all parameters.

4. Large-scale attentional data collection by AMT deployment

The motivation for the mouse-tracking paradigm is for large-scale data collection. In this section, we report implementation and design issues to deploy the mouse-tracking experiments on the paid AMT crowdsourcing marketplace.

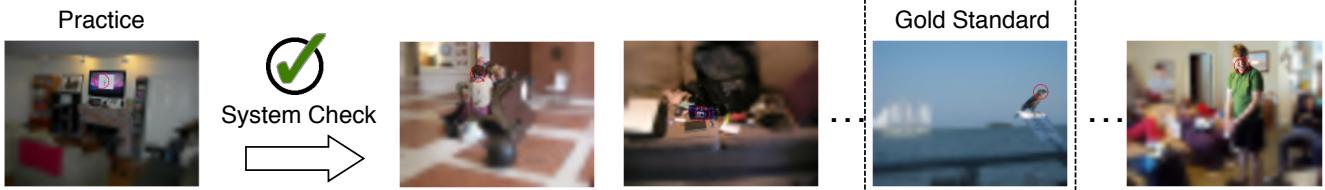
432
433
434
435
436
437
438
439
440486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

Figure 4. The procedure of an AMT task.

We employed the same paradigm and parameter configurations as described in Section 3, while making a few minor adjustments to the procedure to accommodate the more uncontrolled online situations. Figure 4 illustrates the online experiment procedure on AMT.

Our task required real-time rendering of the mouse-contingent stimuli, *i.e.*, the image rendering was triggered by the mouse events in the browser. Therefore it was important to perform a system check to ensure a smooth rendering during visual exploration. The system check was conducted at the practice stage of an AMT task, which detected failures due to a variety of reasons such as unsupported browser features, unfriendly browser plug-ins, and low memory capacity. To ensure that our paradigm was shown smoothly without noticeable lag at the browser side, we evaluated the synchronization quality of the display and the mouse activity, by measuring the distances between the mouse positions and the rendered centers of attention. Only workers who passed the system check could continue the task.

We deployed the experiment on AMT using 10,000 MS COCO training images with 640×480 pixels and 700 OSIE images (scaled to 640×480 pixels). The OSIE images were added as “gold standard”, where the eye tracking data in OSIE can be used as a baseline to evaluate the performance of workers. Currently in each task, a worker viewed 40 images, including 36 images from the MS COCO dataset and 4 images from the OSIE dataset. With the large-scale data collection, we created a Saliency in Context (Salicon) dataset, with 10,000 MS COCO images viewed by 60 observers each. Details of the mouse-tracking results and statistics of the experiments are reported in Section 5.

5. Statistics and results

In this section, we report the mouse-tracking statistics of the two datasets – OSIE and Salicon. For OSIE images, we compare three sets of data: eye-tracking, mouse-tracking in lab, and mouse-tracking with AMT. For Salicon, we report the mouse-tracking statistics in terms of the MS COCO object categories.

5.1. Data preprocessing

Due to the differences in hardware and software settings, the mouse-tracking data have a large variety of sample rates. In the lab experiments, the mean sample rate was 285.61

Hz, across all subjects. While in the AMT data, due to the event system of the browser environments, the sampling was not triggered until the mouse moved. Therefore, the mean sample rate was 69.42 Hz. We discarded the data with sample rate lower than 12 Hz, and resampled the data with a shape-preserving piecewise cubic interpolation that matched the data in position, velocity and acceleration. This was to equalize the number of samples across all observers. The normalized mouse samples had a uniform sample rate at 100 Hz. We added a simple pre-processing step by excluding half samples with high mouse-moving velocity (*i.e.* saccades) for each observer, while keeping the fixations. All pre-processed mouse samples for the same image were then aggregated and blurred with a Gaussian filter to generate a saliency map, same as the common practice to generate the fixation maps from eye-tracking data [34].

5.2. Center bias

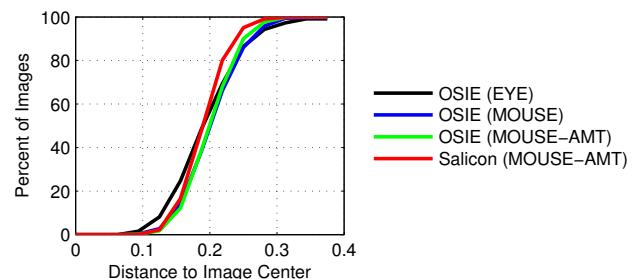


Figure 5. Cumulative distributions of mean distances from fixations / samples to the image center. Distances are normalized with the image width.

In almost all eye-tracking datasets, there exists a spatial prior that pixels near the image center attract more fixations, known as the center bias [29]. The main reasons of the center bias include photographer bias, experimental configuration, and viewing strategy. Similarly, our mouse-tracking data are also biased towards the image center. The cumulative distribution of the mean distance from sample points to the image center is shown in Figure 5. We normalized the distance to center by the image width, and did not observe significant differences in the average distance to center between the AMT and controlled mouse-tracking data or between mouse-tracking and eye-tracking data.

540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559

5.3. Evaluating mouse maps with eye fixations

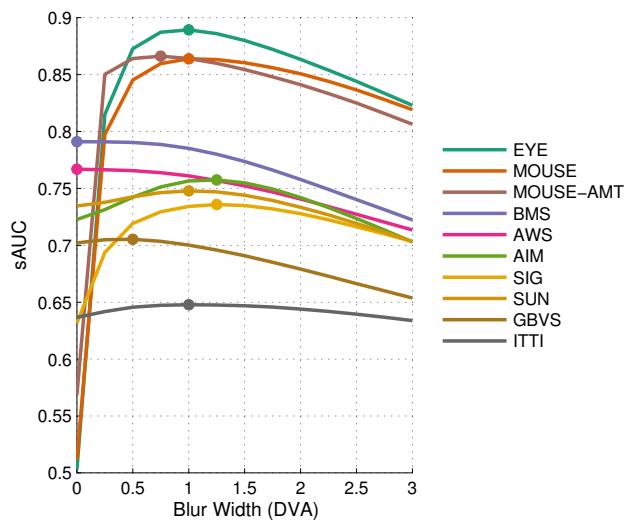


Figure 6. Eye fixation prediction performance with mouse-tracking and the highly referred/state-of-the-art computational saliency models: eye tracking (EYE), mouse map in lab (MOUSE), mouse map on AMT (MOUSE-AMT) the Itti & Koch model (ITTI) [14], the information maximization model (AIM) [2], the graph-based saliency (GBVS) [10], the saliency using natural statistics (SUN) [36], the image signature (SIG) [11], the adaptive whitening saliency (AWS) [8], and the boolean map saliency (BMS) [35].

We evaluated the similarities of the mouse maps and the eye fixation maps, using the most commonly used evaluation metric – the shuffled AUC (sAUC) [36]. The sAUC computes the area under the receiver operating characteristic (ROC) curve, taking positive samples from the fixations of a test image, and negative samples from all fixations in other images. This way it discounts the global center bias in the dataset. We compared the performances of the mouse maps with the inter-observer performance of eye-tracking (computed by aggregating fixations from other subjects than each tested subject, used as a baseline). We also included the highly referred and the state-of-the-art saliency algorithms in the comparison [14, 10, 36, 2, 8, 11, 35]. All saliency maps were blurred by a Gaussian kernel with σ from 0 (no blurring) to 3 degrees of visual angle (DVA; 24 pixels according to the eye-tracking configuration), and the optimal blur width was chosen for each model.

As shown in Figure 6, the lab and AMT mouse models scored closely in sAUC (~ 0.86). They are much closer to the human performance (~ 0.89) in eye-tracking than the computational models. Figure 7 presents the images with high and low sAUC scores in mouse-tracking (with AMT). While the consistency is high in most images, we noticed that disagreement sometimes happens in images containing small text. This may be caused by the relatively small

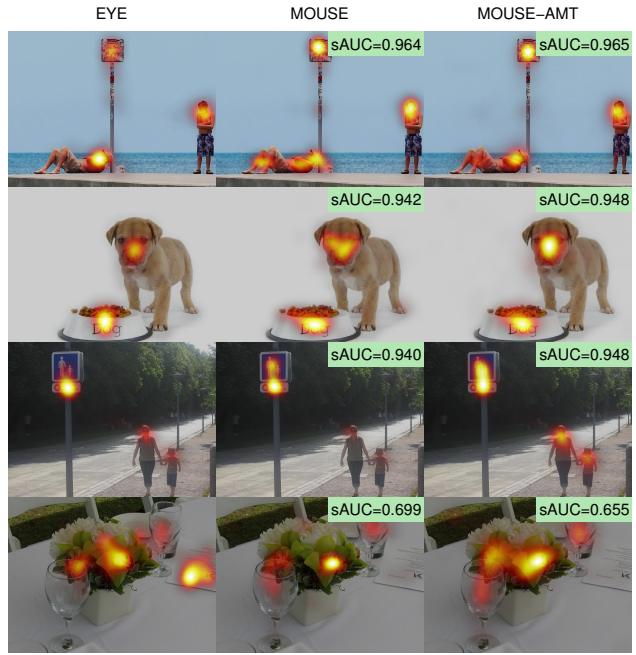


Figure 7. Image examples with high and low eye-mouse similarities evaluated with sAUC. Eye fixation maps and mouse maps are overlaid.

visual angle we use (7.5 pixels per degree) in the mouse-contingent paradigm. Text far from the gaze center may not be easily observed due to the relatively low peripheral resolution. As described in Section 3, the free parameter p corresponds to the visual angle to the scene, ecologically valid in natural vision. While the conventional eye tracking experiments mostly fix this parameter, the proposed paradigm allows the change of this parameter to mimic scenarios with varying distances to the stimuli.

5.4. Categorical analysis

For the Salicon dataset, we sampled 10,000 images from the currently released MS COCO training set, which contains 80 of the 91 categories. The subset was selected from a total of 17,797 images with the resolution of 640×480 . The selection was based on the number of categories in each image. Figure 8 reports the statistics of the dataset in comparison with the MS COCO training images. Our selected images have more instances and categories per image, and are in general richer in contextual information. The rich context is helpful to compare the relative importance of each category in visual exploration. The instance sizes in the Salicon is not significantly different from the full MS COCO training set.

With the mouse maps from the aggregated AMT data, we computed the “maximum object saliency” as the maximum of the map values inside each object’s outline, as it does not scale with the object size [7]. This way we rank the objects

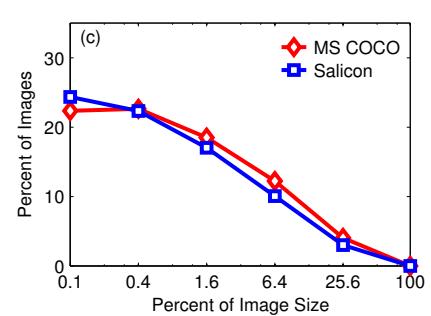
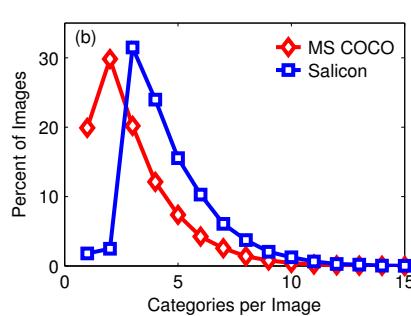
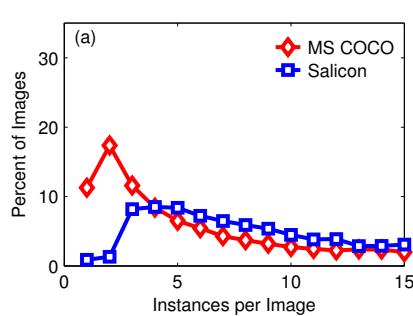


Figure 8. Distributions of (a) number of instances per image, (b) number of categories per image, and (c) instance area.

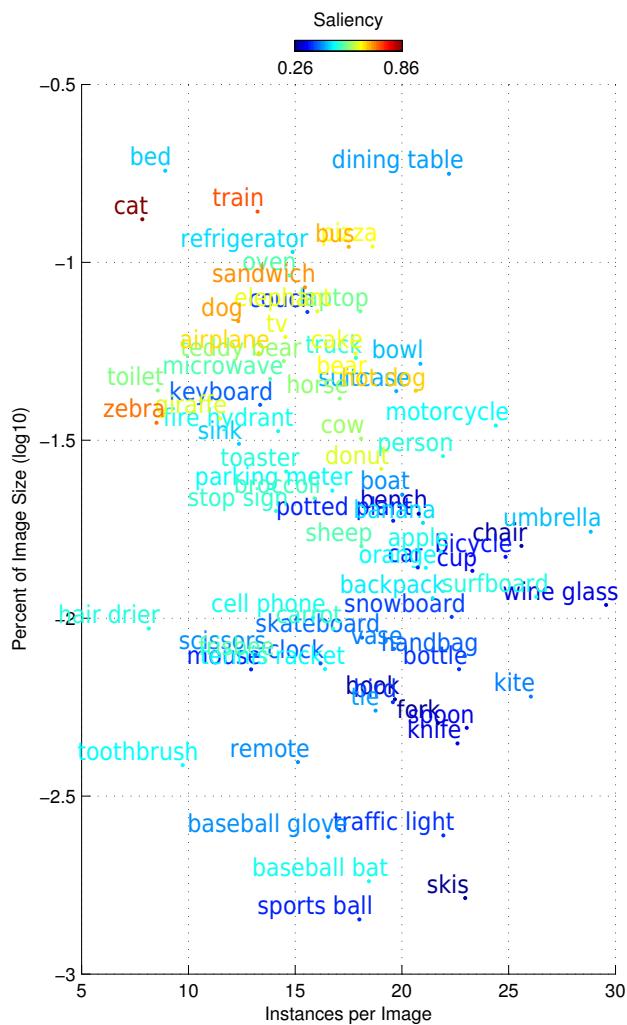


Figure 9. Average saliency values for each of the 80 object categories in Salicon.

in the same image by these values to decide their relative importance.

To quantify the importance of categories with the attentional data, for each category of instances, we computed the



Figure 10. Examples of salient but missed object categories, including face, text, picture, food, door and window, etc. Segmented object instances are masked with colors indicating the categories.

mean instance size, the average number of total instances in the scene which has instances of the particular category, and average saliency value. Figure 9 shows the average saliency values for all the 80 object categories in our dataset. As observed, the importance of a category correlates with its average size and number of instances in the same scene. For the most salient categories, objects appear relatively large in images and are with fewer distractors. Examples include animals, food and train. In comparison, furniture like bed, dining table and refrigerator are relatively less salient, although large in size. Small objects are mostly less salient, except categories that are interactive with humans such as surfboard, baseball bat, and tennis racket.

We further explored the collected attentional data as a natural way to suggest new categories for object annotations and segmentation. The MS COCO has selected 91 categories leveraging domain references, children experiences, and mutual agreement from co-authors. Human attentional data provide yet another complementary source that iden-

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
tify objects that humans look at frequently and rapidly during natural exploration. Figure 10 illustrates examples of typical scenarios where fixations land on unlabeled objects, and suggests several categories be added to the MS COCO dataset to improve its contextual richness. For example, faces attract attention consistently and strongly. Since it is not defined as a category but subregion of ‘person’, we observe that (1) most fixations land on faces though the entire persons are annotated, and (2) some faces are missed if the objects do not belong to the existing category (e.g., toy face, animal face in the first row in Figure 10). Text and pictures also attract attention consistently, but not explicitly defined category in MS COCO (second and third rows). As illustrated in the fourth row, food is frequently missed as only certain types of food are defined (e.g., broccoli, sandwich). Doors and windows attract considerable gaze (fifth row) mostly due to their contextual importance. Detecting these objects would help to understand the context of the scene. These examples demonstrate a potential application of the proposed work in complementing other annotations for visual understanding.

5.5. Mouse tracking as an evaluation benchmark

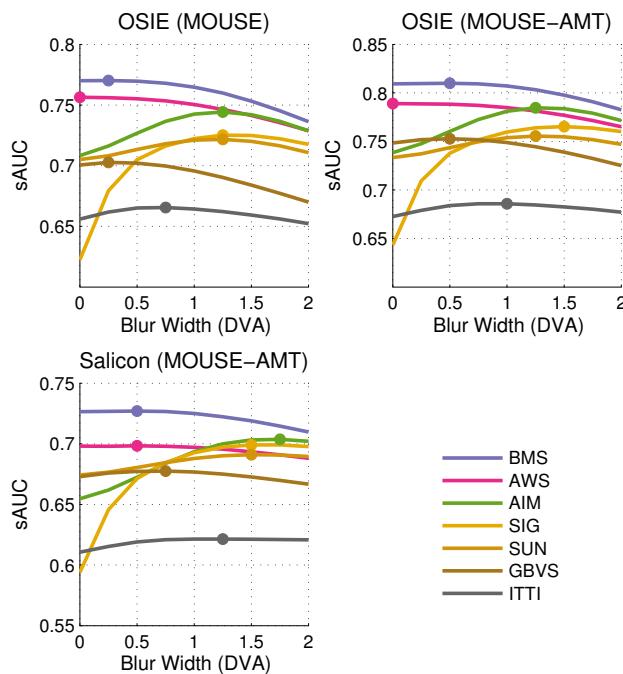


Figure 11. Evaluation of saliency algorithms against mouse-tracking data.

Since the mouse-tracking and eye-tracking data were qualitatively and quantitatively similar, we further exploited the mouse-tracking as a benchmark to evaluate computation saliency algorithms. We tested the state-of-the-art saliency algorithms on the OSIE dataset and randomly selected

2,000 images from the Salicon dataset. We used the pre-processed mouse samples as positive samples in the sAUC computation. For the AMT mouse-tracking data (OSIE and Salicon), in order to reduce the computational cost in the evaluation, we filtered the mouse samples by only keeping the pixels viewed by at least two observers. The comparative results are shown in Figure 11. From the comparison we observe that (1) on OSIE, the sAUC scores for both mouse-tracking data (laboratory and AMT) are close to the eye-tracking ones (see Figure 6), and their ranks are basically preserved. The results show that mouse-tracking is a good replacement of eye-tracking in model evaluation. (2) Comparing the saliency algorithm performance on Salicon vs. on OSIE, similar patterns are observed too. The difference in score reflects dataset difference in image properties.

6. Conclusion

This paper presents a new paradigm to collect human attentional data. Our paradigm enables large-scale data collection by using a general-purpose mouse instead of an expensive eye tracker to record viewing behaviors. With the proposed method, a large mouse-tracking dataset for saliency in context (Salicon) was created on 10,000 images from MS COCO. Salicon is by far the largest attention dataset in both scale and context variability, and data collection on more images is ongoing with the same protocol. With the visual attentional data collected from mouse-tracking, the Salicon dataset complements existing task-specific annotations with natural behavior of visual exploration in task-free situations. The paradigm can also be easily generalized to various types of tasks with top-down instructions. We also envision Salicon to be a good source for learning and benchmarking saliency algorithms with more data.

References

- [1] D. H. Brainard. The Psychophysics Toolbox. *Spat. Vis.*, 10(4):433–436, 1997. 4
- [2] N. Bruce and J. Tsotsos. Saliency Based on Information Maximization. In *NIPS*, pages 155–162, 2005. 6
- [3] N. D. B. Bruce and J. K. Tsotsos. Saliency, attention, and visual search: an information theoretic approach. *J. Vis.*, 9(3):5.1–24, 2009. 2
- [4] M. Cerf, E. P. Frady, and C. Koch. Faces and text attract gaze independent of the task: Experimental data and computer model. *J. Vis.*, 9(12):10.1–15, 2009. 2
- [5] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A large-scale hierarchical image database. In *CVPR*, pages 248–255. IEEE, 2009. 1, 3
- [6] J. Deng, J. Krause, and F.-F. Li. Fine-Grained Crowdsourcing for Fine-Grained Recognition. In *CVPR*, pages 580–587. IEEE, 2013. 3

- 864 [7] W. Einhäuser, M. Spain, and P. Perona. Objects predict fixations better than early saliency. *J. Vis.*, 8(14):18.1–26, 2008. 6
- 865 [8] A. Garcia-Diaz, X. R. Fdez-Vidal, X. M. Pardo, and R. Dosil. 918
- 866 Saliency from hierarchical adaptation through decorrelation 919
- 867 and variance normalization. *Image Vis. Comput.*, 30(1):51– 920
- 868 64, 2012. 6
- 869 [9] D. W. Hansen and Q. Ji. In the eye of the beholder: a survey 921
- 870 of models for eyes and gaze. *TPAMI*, 32(3):478–500, 2010. 922
- 871 2
- 872 [10] J. Harel, C. Koch, and P. Perona. Graph-based visual 923
- 873 saliency. In *NIPS*, pages 545–552, 2006. 6
- 874 [11] X. Hou, J. Harel, and C. Koch. Image signature: Highlighting 924
- 875 sparse salient regions. *TPAMI*, 34:194–201, 2012. 6
- 876 [12] J. Huang, R. W. White, G. Buscher, and K. Wang. Improving 925
- 877 searcher models using mouse cursor activity. In *SIGIR*, page 926
- 878 195, New York, USA, 2012. ACM Press. 3
- 879 [13] H.-W. Hunziker. *Im Auge des Lesers: foveale und periphere 927*
- 880 Wahrnehmung - vom Buchstabieren zur Lesefreude
- 881 Transmedia Verlag, 2006. 3
- 882 [14] L. Itti, C. Koch, and E. Niebur. A model of saliency-based 928
- 883 visual attention for rapid scene analysis. *TPAMI*, 20(11):1254– 929
- 884 1259, 1998. 6
- 885 [15] M. Jiang, J. Xu, and Q. Zhao. Saliency in Crowd. In *ECCV*, 930
- 886 2014. 2
- 887 [16] T. Judd, F. Durand, and A. Torralba. Fixations on low- 931
- 888 resolution images. *J. Vis.*, 11:1–20, 2011. 2
- 889 [17] T. Judd, K. Ehinger, F. Durand, and A. Torralba. Learning 932
- 890 to predict where humans look. In *ICCV*, pages 2106–2113. 933
- 891 IEEE, 2009. 2
- 892 [18] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. 934
- 893 Ramanan, P. Dollár, and C. L. Zitnick. Microsoft COCO: 935
- 894 Common Objects in Context. In *ECCV*, volume cs.CV, pages 936
- 895 740–755, 2014. 1, 3
- 896 [19] V. Navalpakkam, L. L. Jentzsch, R. Sayres, S. Ravi, 937
- 897 A. Ahmed, and A. Smola. Measurement and modeling of 938
- 898 eye-mouse behavior in the presence of nonlinear page 939
- 899 layouts. In *ICWWW*, pages 953–964. International World Wide 940
- 900 Web Conferences Steering Committee, 2013. 3
- 901 [20] D. P. Papadopoulos, A. D. F. Clarke, F. Keller, and V. 941
- 902 Ferrari. Training object class detectors from eye tracking data. 942
- 903 In *Computer Vision - ECCV 2014 - 13th European Conference, 943*
- 904 Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V
- 905
- 906 pages 361–376, 2014. 2
- 907 [21] E. Peli, J. Yang, and R. B. Goldstein. Image invariance with 944
- 908 changes in size: the role of peripheral contrast thresholds. *J. 945*
- 909 Opt. Soc. Am. A
- 910 8(11):1762, 1991. 3
- 911 [22] D. G. Pelli. The VideoToolbox software for visual 946
- 912 psychophysics: transforming numbers into movies. *Spat. Vis.*, 947
- 913 10:437–442, 1997. 4
- 914 [23] J. S. Perry and W. S. Geisler. Gaze-contingent real-time 948
- 915 simulation of arbitrary visual fields. In *SPIE*, volume 4662, 949
- 916 pages 57–69, 2002. 3, 4
- 917 [24] J. S. Pointer and R. F. Hess. The contrast sensitivity gradient 950
- 918 across the human visual field: with emphasis on the low 951
- 919 spatial frequency range. *Vision Res.*, 29(9):1133–1151, 1989. 952
- 920 3
- 921 [25] M. Pomplun, E. M. Reingold, and J. Shen. Investigating 953
- 922 the visual span in comparative search: The effects of task 954
- 923 difficulty and divided attention. *Cognition*, 81, 2001. 3
- 924 [26] S. Ramanathan, H. Katti, N. Sebe, M. Kankanhalli, and T. S. 955
- 925 Chua. An eye fixation database for saliency detection in 956
- 926 images. In *ECCV*, volume 6314 LNCS, pages 30–43, 2010. 957
- 927 2
- 928 [27] K. Rayner and G. W. McConkie. What guides a reader’s eye 958
- 929 movements? *Vision Res.*, 16(8):829–837, 1976. 3
- 930 [28] C. Shen and Q. Zhao. Webpage Saliency. In *ECCV*, num- 959
- 931 ber 65, pages 1–15, 2014. 2
- 932 [29] B. W. Tatler. The central fixation bias in scene viewing: se- 960
- 933 lecting an optimal viewing position independently of motor 961
- 934 biases and image feature distributions. *J. Vis.*, 7(14):4.1–17, 962
- 935 2007. 5
- 936 [30] L. N. Thibos. Acuity perimetry and the sampling theory of 963
- 937 visual resolution. *Optom. Vis. Sci.*, 75:399–406, 1998. 2
- 938 [31] A. Torralba, R. Fergus, and W. T. Freeman. 80 million tiny 964
- 939 images: A large data set for nonparametric object and scene 965
- 940 recognition. *TPAMI*, 30:1958–1970, 2008. 1, 3
- 941 [32] P. Welinder, S. Branson, P. Perona, and S. Belongie. The 966
- 942 Multidimensional Wisdom of Crowds. In *NIPS*, volume 6, 967
- 943 pages 2424–2432, 2010. 3
- 944 [33] J. Xiao, J. Hays, K. A. Ehinger, A. Oliva, and A. Torralba. 968
- 945 SUN database: Large-scale scene recognition from abbey to 969
- 946 zoo. In *CVPR*, pages 3485–3492. IEEE, 2010. 1, 3
- 947 [34] J. Xu, M. Jiang, S. Wang, M. S. Kankanhalli, and Q. Zhao. 970
- 948 Predicting human gaze beyond pixels. *J. Vis.*, 14(1):1–20, 971
- 949 2014. 2, 3, 5
- 950 [35] J. Zhang and S. Sclaroff. Saliency Detection: A Boolean 972
- 951 Map Approach. In *ICCV*, pages 153–160, 2013. 6
- 952 [36] L. Zhang, M. H. Tong, T. K. Marks, H. Shan, and G. W. 973
- 953 Cottrell. SUN: A Bayesian framework for saliency using 974
- 954 natural statistics. *J. Vis.*, 8(7):32.1–20, 2008. 6
- 955 [37] Q. Zhao and C. Koch. Learning a saliency map using fixated 975
- 956 locations in natural scenes. *J. Vis.*, 11(3):9, 2011. 2
- 957
- 958
- 959
- 960
- 961
- 962
- 963
- 964
- 965
- 966
- 967
- 968
- 969
- 970
- 971