# What is a salient object? A dataset and a baseline model for salient object detection

Ali Borji, *Member, IEEE*

*Abstract*—Salient object detection or salient region detection models, diverging from fixation prediction models, have traditionally been dealing with locating and segmenting the most salient object or region in a scene. While the notion of most salient object is sensible when multiple objects exist in a scene, current datasets for evaluation of saliency detection approaches often have scenes with only one single object. We introduce three main contributions in this paper: First, we take an in-depth look at the problem of salient object detection by studying the relationship between where people look in scenes and what they choose as the most salient object when they are explicitly asked. Based on the agreement between fixations and saliency judgments, we then suggest that the most salient object is the one that attracts the highest fraction of fixations. Second, we provide two new less biased benchmark datasets containing scenes with multiple objects that challenge existing saliency models. Indeed, we observed a severe drop in performance of 8 state-of-the-art models on our datasets (40% to 70%). Third, we propose a very simple yet powerful model based on superpixels to be used as a baseline for model evaluation and comparison. While on par with the best models on MSRA-5K dataset, our model wins over other models on our data highlighting a serious drawback of existing models, which is convoluting the processes of locating the most salient object and its segmentation. We also provide a review and statistical analysis of some labeled scene datasets that can be used for evaluating salient object detection models. We believe that our work can greatly help remedy the over-fitting of models to existing biased datasets and opens new venues for future research in this fast-evolving field.

*Index Terms*—Salient object detection, explicit saliency, bottom-up attention, regions of interest, eye movements

## I. INTRODUCTION

**P**LEASE take a look at the images in the top row of Fig. 1. Which object stands out the most (i.e., is the most salient one) in each of these scenes? The answer is trivial. There is only one object, thus it is the most salient one. Now, look at the images in the third row. These scenes are much more complex and contain several objects, thus it is more challenging for a vision system to select the most salient object.

This problem, known as *salient object detection (and segmentation)*, has recently attracted a great deal of interest in computer vision community. The goal is to simulate the astonishing capability of human attention in prioritizing objects for high-level processing. Such a capability has several applications in recognition (e.g., [68]–[70]), image and video compression (e.g., [71], [72]), video summarization (e.g., [73], [74], media re-targeting and photo collage (e.g., [75], [76]),

A. Borji is with the Computer Science Department, University of Wisconsin, Milwaukee, WI, 53211. E-mail: borji@uwm.edu.

image quality assessment (e.g., [77], [78]), image segmentation (e.g., [79]), content-based image retrieval and image collection browsing (e.g., [80], [81]), image editing and manipulating (e.g., [83], [84]), visual tracking (e.g., [82], [86], [87]), object discovery (e.g., [88], [90]), and human-robot interaction (e.g., [91]).

A large number of saliency detection methods have been proposed in the past 7 years (since [17]). In general, a salient object detection model involves two steps: 1) *selecting objects to process* (i.e., determining saliency order of objects), and 2) *segmenting the object area* (i.e., isolating the object and its boundary). So far, models have bypassed the first challenge by focusing on scenes with single objects (See Fig. 1). They do a decent job on the second step as witnessed by very high performances on existing biased datasets (e.g., on ASD dataset [10]) which contain low-clutter images with often a single object at the center. However, it is unclear how current models perform on complex cluttered scenes with several objects. Despite the volume of past research, this trend has not been yet fully pursued, mainly due to the lack of two ingredients: 1) suitable benchmark datasets for scaling up models and model development, and 2) a widely-agreed objective definition of the most salient object. In this paper, we strive to provide solutions for these problems. Further, we aim to discover which component might be the weakest link in the possible failure of models when migrating to complex scenes.

Some related topics, closely or remotely, to visual saliency modeling and salient object detection include: object importance [46], [47], object proposal generation [27], memorability [49], scene clutter [50], image interestingness [52]–[54], video interestingness [51], surprise [55], image quality assessment [56], scene typicality [57], [58], aesthetic [54], and attributes [59], [60].

## II. RELATED WORK

One of the earliest models, which generated the *first wave* of interest in image saliency in computer vision and neuroscience communities, was proposed by Itti *et al.* [2]. This model was an implementation of earlier general computational frameworks and psychological theories of bottom-up attention based on center-surround mechanisms. In [2], Itti *et al.* showed examples where their model was able to detect spatial discontinuities in scenes. Subsequent behavioral (e.g., [41]) and computational studies (e.g., [3]) started to predict fixations with saliency models to verify models and to understand human visual attention. A *second wave* of
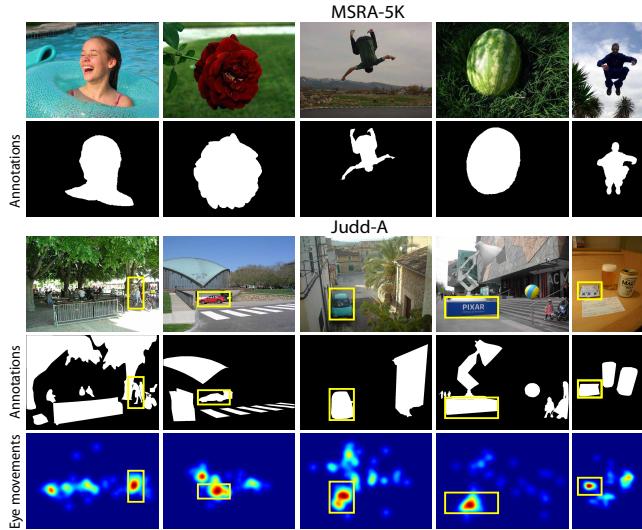
Fig. 1: Top: sample images from the MSRA-5K dataset [15]. Bottom: sample images from Judd-A (Judd-Annotation) dataset [1] with fixations and our annotations. Yellow boxes contain peaks of fixation maps. While scenes in MSRA-5K [16] (and other salient object datasets; Table 1) have close-up views of single objects in uniform and structured backgrounds, scenes in Judd-A dataset are more cluttered and have larger fields of view.

interest appeared with works of Liu *et al.* [17] and Achanta *et al.* [10] who treated saliency detection as a binary segmentation problem with 1 for a foreground pixel and 0 for a pixel of the background region. Since then it has been less clear where this new definition stands as it shares many concepts with other well-established computer vision areas such as general segmentation algorithms (e.g., [34], [35]), category independent object proposals (e.g., [27]), fixation prediction saliency models (e.g. [1], [3], [4]), and general object detection methods. This is partly because current datasets have shaped a definition for this problem, which might not totally reflect full potential of models to *select and segment salient objects in an image with an arbitrary level of complexity*.

Reviewing all saliency detection models goes beyond the scope of this paper (See [26], [30], [39], [89]). Some breakthrough efforts are as follows. Liu *et al.* [17] introduced a conditional random field (CRF) framework to combine multiscale contrast and local contrast based on surrounding, context, and color spatial distributions for binary saliency estimation. Achanta *et al.* [10] proposed subtracting the average color from the low-pass filtered input for saliency detection. Goferman *et al.* [9] used a patch-based approach to incorporate global context, aiming to detect image regions that represent the scene. Cheng *et al.* [14] proposed a region contrast-based method to measure global contrast in the Lab color space. In [11], Wang *et al.* estimated local saliency, leveraging a dictionary learned from other images, and global saliency using a dictionary learned from other patches of the same image. Perazzi *et al.* [29] observed that decomposing an image into perceptually uniform regions, which abstracts away unnecessary details, is important for high quality saliency detection.

In [15], Jiang *et al.* utilized the difference between the color histogram of a region and its immediately neighboring regions for measuring saliency. Feng *et al.* [12] defined a measure of saliency as the cost of composing an image window using the remaining parts of the image, and tested it on PASCAL VOC dataset [85]. This method, in its essence, follows the same goal as in [27]. Chang *et al.* [14] proposed a graphical model for fusing generic objectness [27] and visual saliency for salient object detection. Shen and Wu [38] modeled an image as a low-rank matrix (background) plus sparse noises (salient regions) in a feature space. More recently, Margolin *et al.* [37] integrated pattern and color distinctnesses with high-level cues to measure saliency of an image patch.

Some studies have considered the relationship between fixations and saliency judgments similar to [65]. For example, Xu et al. [43] investigated the role of high-level semantic knowledge (e.g., object operability, watchability, gaze direction) and object information (e.g., object center-bias) for fixation prediction in free viewing of natural scenes. They constructed a large dataset[1] called "Object and Semantic Images and Eye-tracking (OSIE)". Indeed they found an added value for this information for fixation prediction and proposed a regression model (to find combination weights for different cues) that improves fixation prediction performance. Koehler et al. [44] collected a dataset known as the UCSB dataset[2]. This dataset contains 800 images. One hundred observers performed an explicit saliency judgment task, 22 observers performed a free viewing task, 20 observers performed a saliency search task, and 38 observers performed a cued object search task. Observers completing the free viewing task were instructed to freely view the images. In the explicit saliency judgment task, observers were instructed to view a picture on a computer monitor and click on the object or area in the image that was most salient to them. *Salient* was explained to observers as something that stood out or caught their eye (similar to [65]). Observers in the saliency search task were instructed to determine whether or not the most salient object or location in an image was on the left or right half of the scene. Finally, observers who performed the cued object search task were asked to determine whether or not a target object was present in the image. Then, they conducted a benchmark and introduced models that perform the best on each of these tasks.

A similar line of work to ours in this paper has been proposed by Mishra *et al.* [23] where they combined monocular cues (color, intensity, and texture) with stereo and motion features to segment a region given an initial user-specified seed point, practically ignoring the first stage in saliency detection (which we address here by automatically generating a seed point). Ultimately, our attempt in this work is to bridge the interactive segmentation algorithms (e.g., [23], [24]) and saliency detection models and help transcend their applicability.

Perhaps the most similar work to ours has been published by Li *et al.* [42]. In their work, they offer two contributions. *First*, they collect an eye movement dataset using annotated images

---

[1]http://www.ece.nus.edu.sg/stfpage/eleqiz/predicting.html

[2]https://labs.psych.ucsb.edu/eckstein/miguel/research_pages/saliencydata.html

| Dataset | Ref | Total Scenes | Num. of Objects | Annt. | Scene Resolution | Num Annt. | Eye Data |
|---------|-----|-------------|-----------------|-------|-----------------|-----------|----------|
| ASD | [10], [17] | 1000 | ~1 | BD | 400 × 300 | 1 | - |
| MSRA-A | [17] | 20K | ~1 | BB | 400 × 300 | 3 | - |
| MSRA-B | [17] | 5K | ~1 | BB | 400 × 300 | 9 | - |
| MSRA-5K | [16], [17] | 5K | ~1 | BD | 400 × 300 | 1 | - |
| SED-1 | [22] | 100 | 1 | BD | ~300 × 225 | 3 | - |
| SED-2 | [22] | 100 | 2 | BD | ~300 × 225 | 3 | - |
| SOD | [20], [34] | 300 | ~3 | BD | 481 × 321 | 7 | - |
| CSSD | [19] | 200 | ~1 | BD | ~400 × 300 | 1 | - |
| ECSSD | [19] | 10K | ~1 | BD | ~400 × 300 | 1 | - |
| ImgSal | [18] | 235 | ~ 2 | BD | 640 × 480 | 19 | 50 |
| THUR10K | [8], [17] | 10K | ~1 | BD | 400 × 300 | 1 | - |
| THUR15K | [8] | 15K | ~1 | BD | 400 × 300 | 1 | - |
| iCoseg | [40] | 643 | ~1 | BD | ~500 × 400 | 1 | - |
| DUT-OMRON | [45] | 5,172 | ~5 | BD | 400 × 400 | 5 | 5 |
| PASCAL-S | [42], [85] | 850 | ~5 | BD | Variable | 12 | 8 |
| UCSB | [44] | 700 | ~5 | BD | 405 × 405 | 100 | 22 |
| OSIE | [43] | 700 | ~7 | BD | 800 × 600 | 1 | 15 |
| Bruce-A | [3] | 120 | ~4 | BD | 681 × 511 | 70 | 20 |
| Judd-A | [1] | 900 | ~5 | BD | 1024 × 768 | 2 | 15 |

TABLE I: Overview of popular salient object datasets. The last two proposed here (A stands for "Annotation") avoid the dreaded entry of "1" in the number of objects. Compared with other datasets, scenes in Judd-A and Bruce-A datasets have more variety and are less structured. BB and BD stand for bounding box and boundary (i.e., pixel accuracy), respectively. Last column shows the number of eye tracking subjects. Datasets derived from the MSRA carry its problems, which are images with single objects, low clutter, and high degree of center-bias. iCoSeg is a co-segmentation dataset. The ASD dataset is also known as MSRA1000.

from the PASCAL dataset [85] and call their dataset PASCAL-S. *Second*, they propose a model that outperforms other state-of-the-art salient object detection models on this dataset (as well as four other benchmark datasets). Their model decouples the salient object detection problem into two processes: 1) *a segment generation process*, followed by 2) *a saliency scoring mechanism* using fixation prediction. Here, similar to Li *et al.*, we also take advantage of eye movements to measure object saliency but instead of first fully segmenting the scene, we perform a shallow segmentation using superpixels. We then only focus on segmenting the object that is most likely to attract attention. In other words, the two steps are similar to Li *et al.* but are performed in the reverse order. This can potentially lead to better efficiency as the first expensive segmentation part is now only an approximation.

We also offer another dataset which is complimentary to Li *et al.*'s dataset and together both datasets (and models) could hopefully lead to a paradigm shift in the salient object detection field to avoid using simple biased datasets. Further, we situate this field among other similar fields such as general object detection and segmentation, objectness proposal generation models, and saliency models for fixation prediction.

Several salient object detection datasets have been created as more models have been introduced in the literature to extend capabilities of models to more complex scenes. Table I lists properties of 19 popular salient object detection datasets. Although these datasets suffer from some biases (e.g., low scene clutter, center-bias, uniform backgrounds, and non-

ambiguous objects), they have been very influential for the past progress. Unfortunately, recent efforts to extend existing datasets have only increased the number of images without really addressing core issues specifically background clutter and number of objects. Majority of datasets (in particular large scale ones such as those derived from the MSRA dataset) have scenes with often one object which is usually located at the image center. This has made model evaluation challenging since some high-performing models that emphasize image center fail in detecting and segmenting the most salient off-center object [30]. We believe that now is the time to move on to more versatile datasets and remedy biases in salient object datasets.

## III. WHAT IS A SALIENT OBJECT?

In this section, we briefly explain how salient object detection models differ from fixation prediction models, what people consider the most salient object when they are explicitly asked to choose one, what are the relationships between these judgments and eye movements, and what salient object detection models actually predict.

We investigate properties of salient objects from humans' point of view when they are explicitly asked to choose such objects. We then study whether (and to what extent) saliency judgments agree with eye movements. While it has been assumed that eye movements are indicators of salient objects, so far few studies (e.g., [64], [65]) have directly and quantitatively confirmed this assumption. Moreover, the level of agreement and cases of disagreement between fixations and saliency judgments have not been fully explored. Some studies (e.g., [28]), have shown that human observers choose to annotate salient objects or regions first but they have not asked humans explicitly (LabelMe data was analyzed in [28]) and they have ignored eye movements. Knowing which objects humans consider as salient is specially crucial when outputs of a model are going to be interpreted by humans.

### A. Salient object detection vs. fixation prediction

There are two major differences between models defining saliency as "where people look" and models defining saliency as "which objects stand out". *First*, the former models aim to predict points that people look in free-viewing of natural scenes usually for 3 to 5 seconds while the latter aim to detect and segment salient objects (by drawing pixel-accurate silhouettes around them). In principle a model that scores well on one problem should not score very well on the other. An optimal model for fixation prediction should only highlight those points that a viewer will look at (few points inside an object and not the whole object region). Since salient object detection models aim to segment the whole object region they will generate a lot of false positives (these points belong to the object but viewers may not fixate at them) when it comes to fixation prediction. On the contrary, a fixation prediction model will miss a lot of points inside the object (i.e., false negatives) when it comes to segmentation.

*Second*, due to noise in eye tracking or observers' saccade landing (typically around 1 degrees and $\sim$ 30 pixels), highly
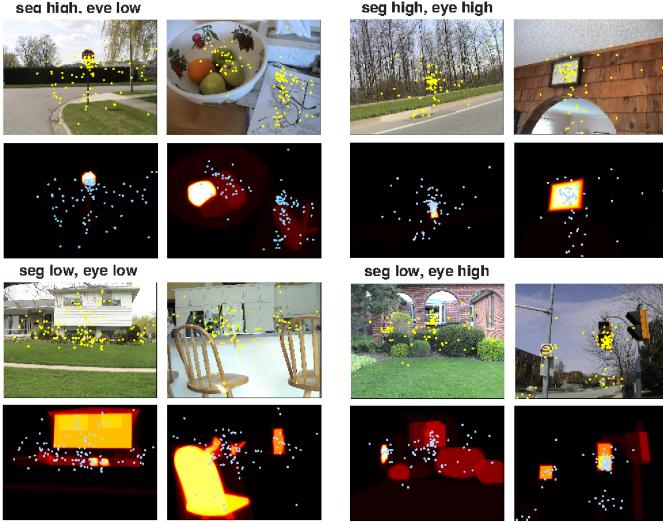
Fig. 2: Sample images from Bruce-A dataset [3] and annotation heat maps embedded with fixations. eye-low/eye-high indicate low/high agreement in eye movements measured in termed of shuffled AUC score [39]. seg-low/seg-high indicate low/high agreement in saliency annotations using the equation in section III.B.

accurate pixel-level prediction maps are less desired. In fact, due to these noises, sometimes blurring prediction maps increases the scores [25], [42], [61]. On the contrary, producing salient object detection maps that can accurately distinguish object boundaries are highly desirable specially in applications. Due to these, different evaluation and benchmarks have been developed for comparing models in these two categories.

In practice, models, whether they address segmentation or fixation prediction, are applicable interchangeably as both entail generating similar saliency maps. For example, several researches have been thresholding saliency maps of their models, originally designed to predict fixations, to detect and segment salient proto-objects (e.g., [62], [63]).

### B. Human explicit saliency judgments

In our previous study [65], we addressed what people consider as the most outstanding (i.e., salient) object in a scene. While in [65] we studied the explicit saliency problem from a behavioral perspective, here we are mainly interested in constructing computational models for automatic salient object detection in arbitrary visual scenes. A total of 70 students (13 male, 57 female) undergraduate USC students with normal or corrected-to-normal vision in the age range between 18 and 23 (mean = 19.7, std = 1.4) were asked to draw a polygon around the object that stood out the most. Participants' annotations were supposed not to be too loose (general) or too tight (specific) around the object. They were shown an illustrative example for this purpose. Participants were able to relocate their drawn polygon from one object to another or modify its outline. We were concerned with the case of selection of the single most salient object in an image. Stimuli were the images

from the dataset by Bruce and Tsotsos (2005) [3][3]. See Fig. 2 for sample images from this dataset.

We first measured the degree to which annotations of participants agree with each other using the following quantitative measure:

$$r_k = \frac{2}{n(n-1)} \sum_{i=1}^{n-1} \sum_{j=i+1}^{n=70} |s_{ik} \cap s_{jk}|/|s_{ik} \cup s_{jk}| \quad (1)$$

where $s_{ik}$ and $s_{jk}$ are annotations of $i$-th and $j$-th participants, respectively (out of $n = 70$ participants) over the $k$-th image. Above measure has the well-defined lower-bound of 0, when there is no overlap in segmentations of users, and the upper-bound of 1, when segmentations have perfect overlap. Fig. 3.left shows histogram of $r$ values. Participants had moderate agreement with each other (mean $r = 0.37$; std 0.17; significantly above chance). Inspection of images with lowest $r$ values shows that these scenes had several foreground objects while images with highest annotation agreement had often one visually distinct salient object (e.g., a sign, a person, or an animal; see Fig. 2).

### C. Relationship between saliency judgments and fixations

We also investigated the relationship between explicit saliency judgments and freeviewing fixations as two indicators of visual attention. Here we used shuffled AUC (sAUC) score to tackle center-bias in eye movement data [6], [39]. For each of 120 images, we showed that a map built from annotations of 70 participants explains fixations of free viewing observers significantly above chance (sAUC of $0.62 \pm 0.07$, chance $0.50$, $t$-test $p < 0.05$; Fig. 3.right). The prediction power of this map was as good as the ITTI98 model [2]. Hence, we concluded that explicit saliency judgments agree with fixations. Fig. 2 shows high- and low-agreement cases between fixations and annotations.

Here, we merge annotations of all 70 participants on each image, normalize the resultant map to [0 1], and threshold it at 0.7 to build our first benchmark saliency detection dataset (called Bruce-A). Prevalent objects in Bruce and Tsotsos dataset are man-made home supplies in indoor scenes (see [3] for more details on this dataset).

Similar results, to link fixations with salient objects, have been reported by Koehler *et al.* [44]. As in [64], they asked observers to click on salient locations in natural scenes. They showed high correlation between clicked locations and observers' eye movements (from a different group of subjects) in free-viewing. While the most salient [65], important [46], or interesting [28], [51], [64] object may tell us a lot about a scene, eventually there is a subset of objects that can minimally describe a scene. This has been addressed in the past somewhat indirectly in the contexts of saliency [9], language and attention [66], and phrasal recognition [59], [67].

---

[3]This dataset contains eye movements over 120 color photographs of indoor and outdoor environments with the resolution of $511 \times 681$ pixels. Images in this dataset have been presented at random to 20 observers (in a free-viewing task) for 4 sec each, with 2 sec of delay (a gray mask) in between.
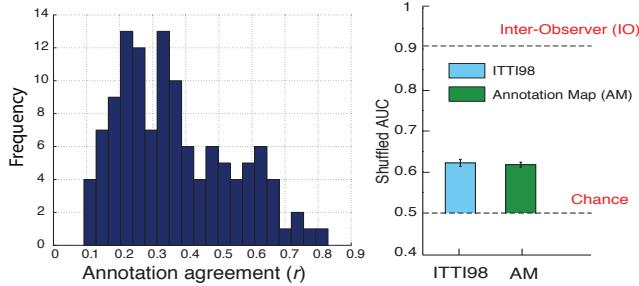
Fig. 3: Left: Histogram of agreements in our explicit saliency judgment task, Right: Prediction power of the annotation map (AM) and the saliency model by Itti *et al.* [2] for explaining eye movements over Bruce and Tsotsos dataset. Chance level is the accuracy of a random map with the value of each pixel drawn uniformly random between 0 and 1. Inter-observer (IO) model is a map build from fixations of other observers over the same image and is then smoothed with a small Gaussian kernel.
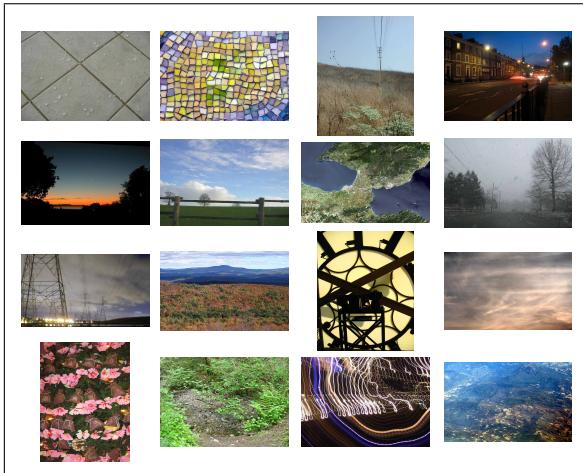


Fig. 4: Sixteen sample images from the Judd dataset that are not included in Judd-A dataset. These images do not have well-defined salient objects in them, contain lots of background clutter and are often boring!.

## IV. OUR NEW LARGE-SCALE DATASET: JUDD-A

Based on results from our saliency judgment experiment [65], we then decided to annotate scenes of the dataset by Judd *et al.* [1]. The reason for choosing this dataset is because it is currently the most popular dataset for benchmarking fixation prediction models [1], [39]. It contains eye movements of 15 observers freely viewing 1003 scenes from variety of topics. Thus, using fixations we can easily determine which object, out of several annotated objects, is the most salient one. We only used 900 images from the Judd dataset and discarded images without well-defined objects (e.g., mosaic tiles, flames) or images with very cluttered backgrounds (e.g., nature scenes). Figure 4 shows examples of discarded scenes.

We asked 2 observers to manually outline objects using the LabelMe [7] open annotation tool (http://new-labelme.csail.mit.edu/). Observers were instructed to accurately segment as many objects as possible following three rules: 1) discard reflection of objects in mirrors, 2) segment objects that are not separable as one (e.g., apples in a basket), and 3) interpolate the boundary of occluded objects only if doing otherwise may create several parts for an occluded object. These cases, however, did happen rarely. Observers were also told that their outline should be good enough for somebody to recognize the object just by seeing the drawn polygon. Observers were paid for their effort. Fig. 5 shows sample images and their annotated objects. To determine which object is the most salient one, we selected the object at the peak of the human fixation map.

### A. Dataset statistics

Here we explore some summary statistics of our data. On average, 36.93% of an image pixels was annotated by the 1st observer with a std of 29.33% (44.52%, std=29.36% for the 2nd observer). 27.33% of images had more than 50% of their pixels segmented by the 1st observer (34.18% for the 2nd). The number of annotated objects in a scene ranged from 1 to 31 with median of 3 for the 1st observer (1 to 24 for the 2nd observer with median of 4; Fig. 6.left). The median object size was 10% of the total image area for the 1st observer (9% for the 2nd observer). Fig. 6.left (inset) shows the average annotation map for each observer over all images. It indicates that either more objects were present at the image center and/or observers tended to annotate central objects more. Overall, our data suggests that both observers agree to a good extent with each other. Finally, in order to create one ground truth segmentation map per image, we asked 5 other observers to choose the best of two annotations (criteria based on selection of annotated objects and boundary accuracy). The best annotation was the one with max number of votes (611 images with 4 to 1 votes).

Next, we quantitatively analyzed the relationship between fixations and annotations (Note that we explicitly define the most salient object as the one with the highest fraction of fixations on it). We first looked into the relationship between the object annotation order and the fraction of fixations on objects. Fig. 6.middle shows fraction of fixations as a function of object annotation order. In alignment with previous findings [28], [65] we observe that observers chose to annotate objects that attract more fixations. But here, unlike [28] which used saliency models to demonstrate that observers prioritize annotating interesting and salient objects, we used actual eye movement data. We also quantized the fraction of fixations that fall on scene objects over the Judd-A annotations, and observed that in about 55% of images, the most salient object attracts more than 50% of fixations (mean fixation ratio of 0.54; image background=0.45; Fig. 6.right).

The most salient object ranged in size from 0.1% to 90.2% of the image size (median=10.17%). The min and max aspect ratio (W/H) of bounding boxes fitted to the most salient object were 0.04 and 13.7, respectively (median=0.94).

Judd dataset is known to be highly center-biased [1], [39], in terms of eye movements [6], due to two factors: 1) the tendency of observers to start viewing the image from the center (a.k.a viewing strategy), and 2) tendency of photographers to frame interesting objects at the image center (a.k.a

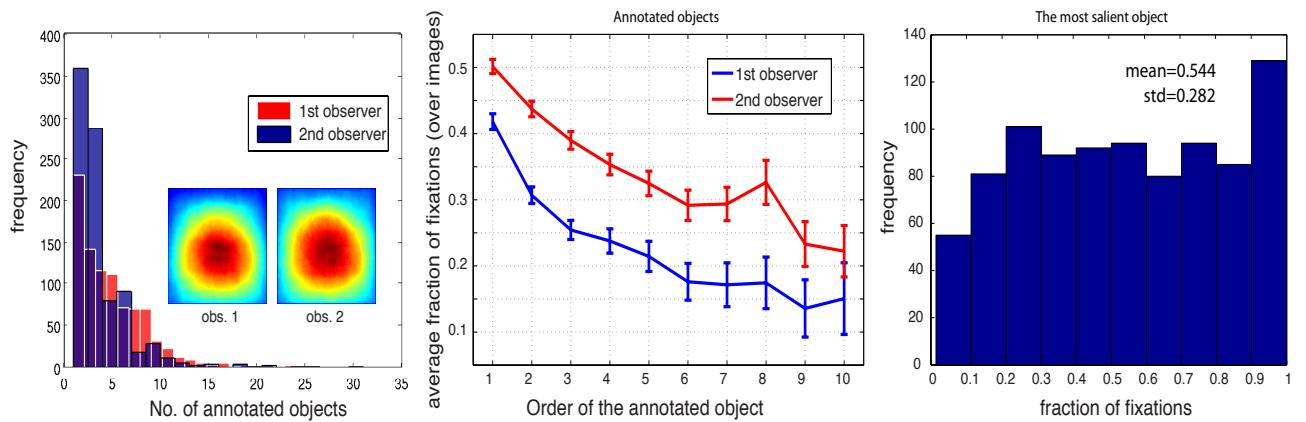Fig. 5: Sample scenes from our dataset and their corresponding annotations.

Fig. 6: Left: The histogram of number of annotated objects by two observers, Middle: Average fraction of fixations as a function of annotation order. Error bars show the standard error of the mean (s.e.m) over 900 images. Right: Histogram of fixation ratios on the most salient object of the Judd-A dataset.

photographer bias). Here we verify the second factor by showing the average annotation map of the most salient object in Fig. 7. Our datasets seem to have relatively less center-bias compared to MSRA-5K and CSSD datasets. Note that other datasets mentioned in Table I are also highly center-biased. To count the number of images with salient objects at the image center, we defined the following criterion. An image is on-centered if its most salient object overlaps with a normalized (to [0 1]) central Gaussian filter with $\sigma = 50$. This Gaussian filter is resized to the image size and is then truncated above 0.95. Utilizing this criterion, we selected 667 and 223 on-centered and off-centered scenes, respectively. Partitioning data in this manner helps scrutinize performance of models and tackle the problem of center-bias.

To further explore the amount of center-bias in Bruce-A and Judd-A datasets, we first calculated the Euclidean distance from center of bounding boxes, fitted to object masks, to the image center. We then normalized this distance to the half of the image diagonal (i.e., image corner to image center). Fig. 8.left shows the distribution of normalized object distances. As opposed to MSRA-5K and CSSD datasets that show an unusual peak around the image center, objects in our datasets are further apart from the image center.
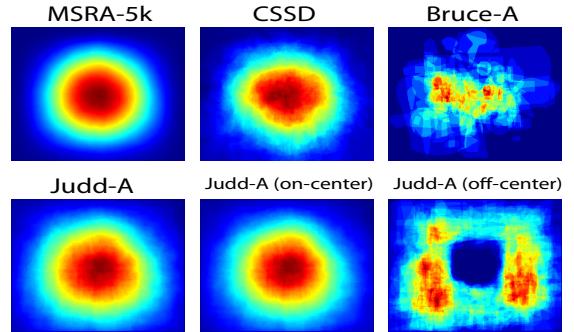
Fig. 7: Average annotation maps of two salient object datasets and our datasets. Distributions of salient objects for on-center and off-center scenes are also shown.

Fig. 8.right shows distributions of normalized object sizes. A majority of salient objects in Bruce-A and Judd-A datasets occupy less than 10% of the image. On average, objects in our datasets are smaller than MSRA-5K and CSSD making salient object detection more challenging.

We also analyzed complexity of scenes on four datasets. To this end, we first used the popular graph-based superpixel seg-
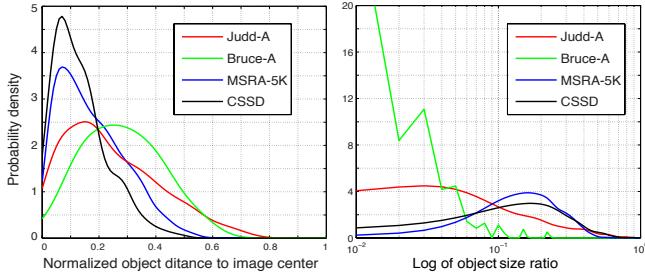
Fig. 8: Left: distribution of normalized object distance (*object center - image center* over *half of the image diagonal*). Right: distribution of salient object size ratio (*object size* over *image area*) in log scale across four datasets.

mentation algorithm by Felzenszwalb and Huttenlocher [21] to segment an image into contiguous regions larger than 60 pixels each (parameter settings: $\sigma = 1$, segmentation coefficient $K = 300$). The basic idea is that the more superpixels an image contains, the more complex and cluttered it is [13]. By analogy to scenes, an object with several superpixels is less homogeneous, and hence is more complex (e.g., a person vs. a ball ). Fig. 9 shows distributions of number of superpixels on the most salient object, the background, and the entire scene. If a superpixel overlapped with the salient object and background, we counted it for both. In general, complexities of backgrounds and whole scenes in our datasets, represented by blue and red curves, are much higher than in the other two datasets. The most salient object in Judd-A dataset on average contains more superpixels than salient objects in MSRA-5K and CSSD datasets, even with smaller objects. The reason why number of superpixels is low on the Bruce-A dataset is because of its very small salient objects (See Fig. 9.right).

Further, we inspected types of objects in Judd-A images. We found that 45% of images have at least one person in them and 27.2% have more than two people. On average each scene has 1.56 persons (std = 3.2). In about 27% of images, annotators chose a person as the most salient object. We also found that 280 out of 900 images (31.1%) had one or more text in them. Other frequent objects were animals, cars, faces, flowers, and signs.
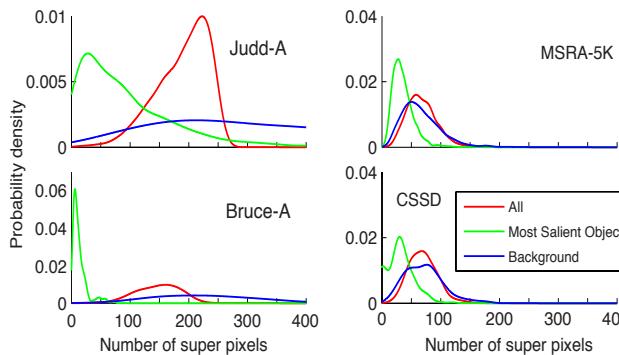


Fig. 9: Distributions of the number of superpixels for the most salient object, the background, and entire scene (All) across four datasets.

## V. OUR BASELINE SALIENCY MODEL: SALBASE

In general, it is agreed that for good saliency detection, a model should meet the following three criteria: 1) *High detection rate*. There should be a low probability of failing to detect real salient regions, and low probability of falsely detecting background regions as salient regions, 2) *High resolution*. Saliency maps should have high or full resolution to accurately locate salient objects and retain original image information as much as possible, and 3) *High computational efficiency*. Saliency models with low processing time are preferred. Here, we analyze these factors by proposing a simple baseline salient object detection model.

We propose a straightforward model to serve two purposes: 1) *to assess the degree to which our data can be explained by a simple model*. This way our model can be used for measuring bias and complexity of a saliency dataset, and 2) *to gauge progress and performance of the state of the art models*. By comparing performance of best models relative to this baseline model over existing datasets and our datasets, we can judge how powerful and scalable these models are. Note that we deliberately keep the model simple to achieve above goals.

Our model involves the following two steps:

***Step 1***: Given an input image, we compute a saliency map and an over-segmented region map. For the former, we use a fixation prediction model (traditional saliency models) to find spatial outliers in scenes that attract human eye movements and visual attention. Here, we use two models for this purpose: AWS [5] and HouNIPS [4], which have been shown to perform very well in recent benchmarks and to be computationally efficient [39]. As controls, we also use the generic *objectness* measure by Alexe *et al.* [27], as well as the human fixation map to determine the upper-bound performance. The reason for using fixation saliency models is to obtain an quick initial estimation of locations where people may look in the hope of finding the most salient object. These regions are then fed to the segmentation component in the next step. It is critical to first limit the subsequent expensive processes onto the right region. For the latter, as in the previous section, we use the fast and robust algorithm by Felzenszwalb and Huttenlocher [21][4] with same parameters as in section IV-A.

***Step 2***: The saliency map is first normalized to [0 1] and is then thresholded at $\beta$ (here $\beta = 0.7$). Then all unique image superpixels that spatially overlap with the truncated saliency map are included. Here we discarded those superpixels that touch the image boundary because they are highly likely to be part of the background. Finally, after this process, the holes inside the selected region will be considered as part of the salient object (e.g., filling in operation). Fig. 10 illustrates the process of segmentation and shows outputs of our model for some images from MSRA-5K, Bruce-A, and Judd-A datasets.

The essential feature of our simple model is dissociating saliency detection from segmentation, such that now it is possible to pinpoint what might be the cause of mistakes or low performance of a model, i.e., detecting the wrong object or faulty segmentation. This is particularly important since

---

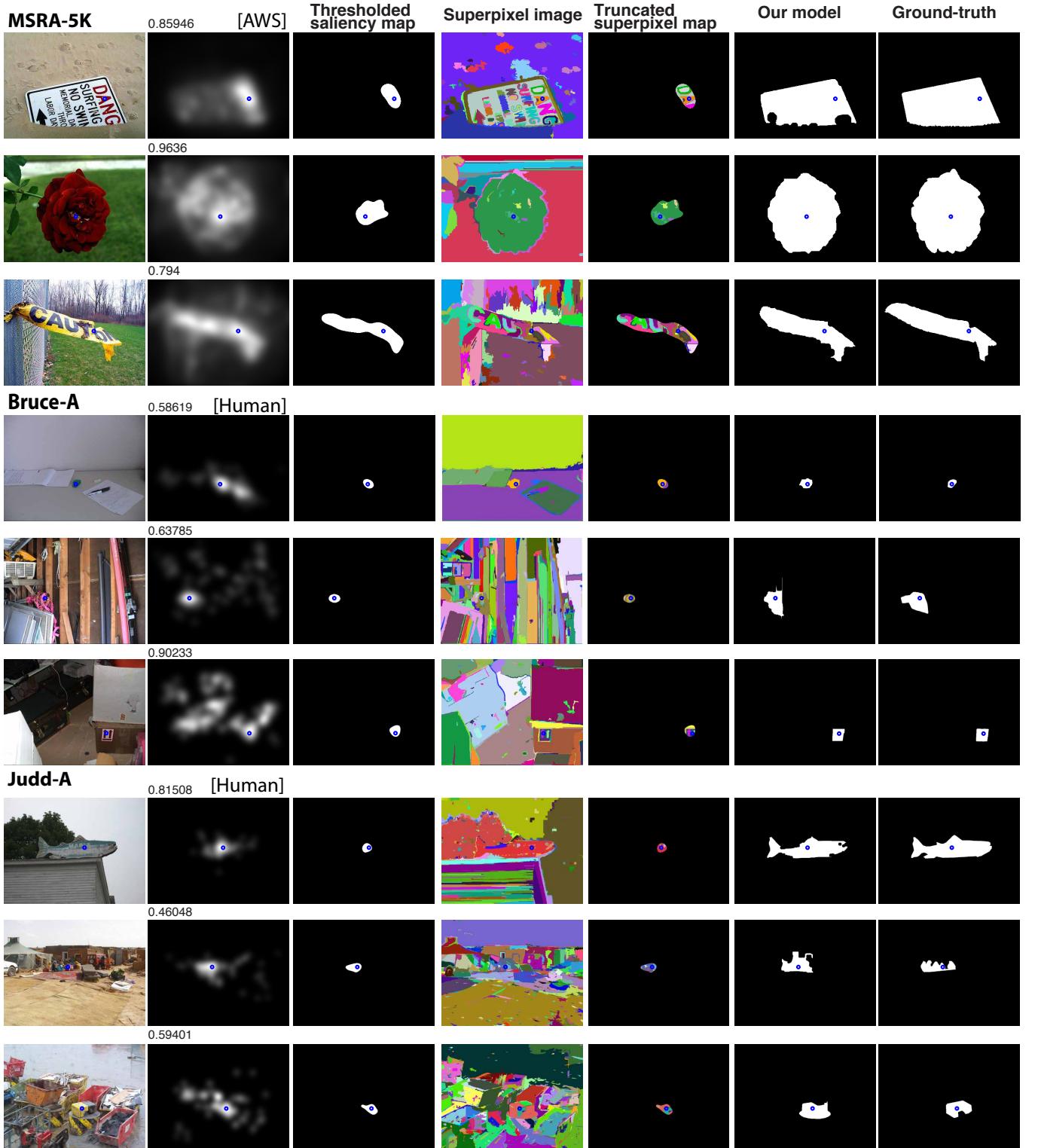[4]We achieved lower accuracy using the normalized cut algorithm [31].

Fig. 10: Columns from left to right: original image, saliency map (human or a model), top 70% of the saliency map ($\beta = 0.7$), graph-based superpixel segmentation [21], truncated superpixel image, our model prediction, and ground-truth. Small blue dot represents the location of saliency map maximum. Numbers above images in the second column show the PASCAL criterion $\Omega = |\omega \cap o|/|\omega \cup o|$ where $\omega$ is our segmentation and $o$ is the ground-truth annotation, so the higher the $\Omega$ the better.
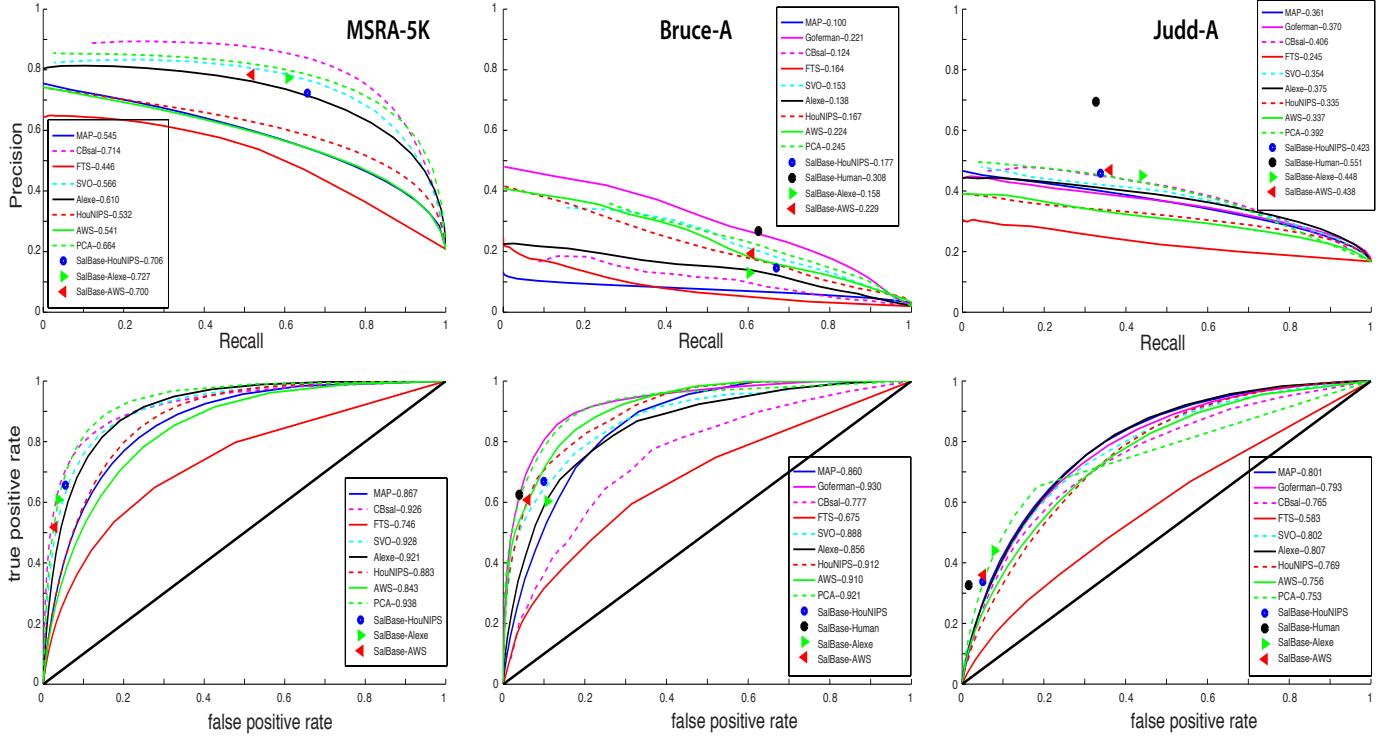
Fig. 11: PR and ROC curves of our model versus 8 other models. Note how drastically models are degraded on our dataset, which contain scenes with multiple objects. We did not run Goferman model [9] over MSRA-5K due to its slow computation. Alexe model [27] is used with 1000 windows. MAP stands for the mean annotation map (Fig. 7). Also note that we expect similar performance for even better recent salient object detection models (e.g.,DRFI [16]).

almost all models have confused these two steps and have faded the boundary.

Note that currently there is no training stage in our model and it is manually constructed with fixed parameters. The second stage in our model is where more modeling contribution can be made, for example by devising more elaborate ways to include or discard superpixels in the final segmentation. One strategy is to learn model parameters from data. Some features to include in a learning method are size and position of a superpixel, a measure of elongatedness, a measure of concavity or convexity, distance between feature distributions of a superpixel and its neighbors, etc. To some extent, some of these these features have already been utilized in previous models [8], [30]. Another direction will be expanding our model to multi scale (similar to [19]).

## VI. MODEL EVALUATION AND COMPARISON

We exhaustively compared our model to 8 state of the art methods which have been shown to perform very well on previous benchmarks [30]. These models come from 3 categories allowing us to perform cross-category comparison: 1) *salient object detection models* including CBsal [15], SVO [14], PCA [37], Goferman [9], and FTS [10], 2) *generic objectness measure* by Alexe *et al.* [27], and 3) *fixation prediction models* including AWS [5] and HouNIPS [4].

We use two widely adopted metrics:

- **Precision-recall (PR) curve:** For a saliency map $S$ normalized to $[0, 255]$, we convert it to a binary mask $M$ with a threshold $T_f$. $Precision$ and $Recall$ are then computed as follows given the ground truth mask $G$:

$$Precision = \frac{|M \cap G|}{|M|}, \quad Recall = \frac{|M \cap G|}{|G|} \quad (2)$$

To measure the quality of saliency maps produced by several algorithms, we vary the threshold $T_f$ from 0 to 255. On each threshold, $Precision$ and $Recall$ values are computed. Finally, we can get a precision-recall (PR) curve to describe the performance of different algorithms. We also report the F-Measure defined as:

$$F_\alpha = \frac{(1 + \alpha)Precision \times Recall}{\alpha \times Precision + Recall} \quad (3)$$

Here, as in [10] and [8], we set $\alpha = 0.3$ to weigh precision more than recall.

- **Receiver operating characteristics (ROC) curve:** We also report the false positive rate ($FPR$) and true positive rate ($TPR$) during the thresholding a saliency map:

$$TPR = \frac{|M \cap G|}{|G|}, \quad FPR = \frac{|M \cap G|}{|M \cap G| + |\bar{M} \cap \bar{G}|} \quad (4)$$

where $\bar{M}$ and $\bar{G}$ denote the opposite (complement) of the binary mask $M$ and ground-truth, respectively. The ROC curve is the plot of $TPR$ versus $FPR$ by varying the threshold $T_f$.

| Model | Dataset | | |
|---|---|---|---|
| | MSRA - 5K | Bruce-A | Judd-A |
| MEP | 0.545 | 0.10 | 0.361 |
| Goferman | - | 0.221 | 0.370 |
| CBsal | 0.714 | 0.124 | 0.406 |
| FTS | 0.446 | 0.164 | 0.245 |
| SVO | 0.566 | 0.153 | 0.354 |
| Alexe | 0.610 | 0.138 | 0.375 |
| HouNIPS | 0.532 | 0.167 | 0.335 |
| AWS | 0.541 | 0.224 | 0.337 |
| PCA | 0.664 | 0.245 | 0.392 |
| SalBase-HouNIPS | 0.706 | 0.177 | 0.423 |
| SalBase-Alexe | **0.727** | 0.158 | 0.448 |
| SalBase-AWS | 0.700 | 0.229 | 0.438 |
| SalBase-Human | - | **0.308** | **0.551** |

TABLE II: F-measure accuracy of models. Performance of the best model is highlighted in boldface font.

| Model | Dataset | | |
|---|---|---|---|
| | MSRA - 5K | Bruce-A | Judd-A |
| MEP | 0.867 | 0.860 | 0.801 |
| Goferman | - | **0.930** | 0.793 |
| CBsal | 0.926 | 0.777 | 0.765 |
| FTS | 0.746 | 0.675 | 0.583 |
| SVO | 0.928 | 0.888 | 0.802 |
| Alexe | 0.921 | 0.856 | **0.807** |
| HouNIPS | 0.883 | 0.912 | 0.769 |
| AWS | 0.843 | 0.910 | 0.756 |
| PCA | **0.938** | 0.921 | 0.753 |
| SalBase-HouNIPS | 0.781 | 0.751 | 0.633 |
| SalBase-Alexe | 0.773 | 0.714 | 0.662 |
| SalBase-AWS | 0.737 | 0.756 | 0.644 |
| SalBase-Human | - | 0.780 | 0.654 |

TABLE III: AUC accuracy of models. Performance of the best model is highlighted in boldface font.

### A. Quantitative evaluation

Results are shown in Fig. 11. Consistent with previous reports over the MSRA-5K dataset [16], CBsal, PCA, SVO, and Alexe models rank on the top (with F-measures above 0.55 and AUCs above 0.90). Fixation prediction models perform lower at the level of the MAP. FTS model ranked on the bottom again in alignment with previous results. Our models work on par with the best models on this dataset with all F-measures above 0.70 (max with Alexe model about 0.73). Moving from this simple dataset (because our simple models ranked on the top; see also the analysis in section IV-A) to more complex datasets (middle column in Fig. 11) we observed a dramatic drop in performance of all models. The best performance now is 0.24 belonging to the PCA model. We observed about 72% drop in performance averaged over 5 models (CBsal, FTS, SVO, PCA, and Alexe) from MSRA-5K to Bruce-A dataset. Note in particular how MAP model is severely degraded here (poorest with F measure of 0.1) since objects are now less at the center. Our best model on this dataset is the SalBase-Human (F-measure about 0.31). Surprisingly, AUC results are still high on this dataset since objects are small thus true positive rate is high at all levels of false positive rate (See also performance of MAP). Patterns of results over Judd-A dataset are similar to those over Bruce-A with all of our models performing higher than others. The lowest performance here belongs to FTS followed by the two fixation prediction models. Our SalBase-Human model scores the best with the F-measure about 0.55. Among our models that used a model to pick the most salient location, SalBase-AWS scores higher over Bruce-A and Judd-A datasets possibly because AWS is better able to find the most salient location. The average drop from MSRA-5K to Judd-A dataset is $\sim 41\%$ (for 5 saliency detection models). Fig. 12 shows that these findings are robust to F-measure parameterization. Tables II and III summarize the F-measure and AUC of models.

### B. Analysis of saliency map thresholding

To study the dependency of results on saliency map thresholding (i.e., how many superpixels to include), we varied the saliency threshold $\beta$ and calculated F-measure for SalBase-Human and SalBase-AWS models (See Fig. 13). We observed

that even higher scores are achievable using different parameters. For example, since objects in the Judd-A dataset are larger, a lower threshold yields a better accuracy. The opposite holds over the Bruce-A dataset.

### C. Analysis of superpixel segmentation parameters

To investigate the dependency of results on segmentation parameters, we varied the parameters of the segmentation algorithm from too fine ($\sigma = 1$, K = 100, min = 20; many segments; over-segmenting) to too coarse ($\sigma = 1$, K = 1000, min = 800; fewer segments; under-segmenting). Both of these settings yielded lower performances than results in Fig. 11. Results with another parameter setting with $\sigma = 1$, K = 500, and min = 50 are shown in Fig. 14. Scores and trends are similar to those shown in Fig. 11, with SalBase-Human and SalBase-AWS being the top contenders.
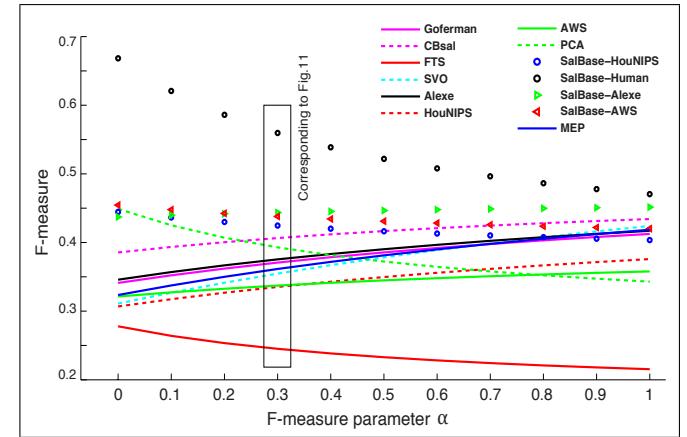


Fig. 12: F-measure as a function of parameter $\alpha$ on Judd-A dataset.

### D. Analysis of model failure cases

Analysis of cases where our model fails, shown in Fig. 15, reveals four reasons: *First*, on Bruce-A dataset when humans look at an object more but annotators chose a different object. *Second*, when a segment that touches the image border is
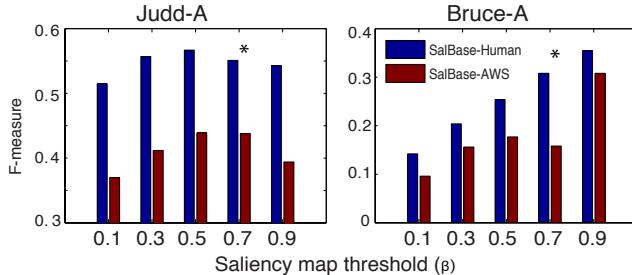
Fig. 13: F-measure as a function of saliency threshold $\beta$. Stars correspond to points shown in Fig. 11. Note that even higher accuracies are possible with different thresholds over our datasets. Corresponding F-measure values over MSRA-5K for SalBase-AWS model are: 0.62, 0.67, 0.70, 0.73, and 0.62.

part of the salient object. *Third*, when the object segment falls outside the thresholded saliency map (or a wrong one is included). *Fourth*, when the first stage (i.e., fixation prediction model) pick the wrong object as the most salient one (See Fig. 16, first column). Regarding the first problem, care must be taken in assuming what people look is what they choose as the most salient object. Although this assumption is correct in a majority of cases (Fig. 3), it does not hold in some cases. With respect to the second and third problems, future modeling effort is needed to decide which superpixels to include/discard to determine the extent of an object. The fourth problem points toward shortcomings of fixation prediction models. Indeed, in several scenes where our model failed, people and text were the most salient objects. Person and text detectors were not utilized in the saliency models employed here.

### E. Qualitative comparison of models

Fig. 16 shows a visual comparison of models over 12 scenes from the Judd-A dataset. CBsal and SVO generate more visually pleasant maps. Goferman highlights object boundaries more than object interiors. PCA generates center-biased maps. Some models (e.g., Goferman, FTS) generate sparse saliency maps while some others generate smoother ones (e.g., SVO, CBSal). AWS and HouNIPS models generate pointy maps to better account for fixation locations.

### VII. DISCUSSION AND CONCLUSION

In this work, we showed that: 1) explicit human saliency judgments agree with free-viewing fixations (thus extending our previous results in [65]), 2) our new benchmark datasets challenge existing state-of-the-art salient object detection models (in alignment with Li *et al.*'s dataset [42]), and 3) a conceptually simple and computationally efficient model ($\sim$0.2 s for $400 \times 300$ saliency and segmentation maps on a PC with a 3.2 GHz Intel i7 CPU and 6GB RAM using Matlab) wins over the state of the art models and can be used as a baseline in the future. We also highlighted a limitation of models which is the main reason behind their failure on complex scenes. They often segment the wrong object as the most salient one.

Previous modeling effort has been mainly concentrated on biased datasets with images containing objects at the center.

Here, we focused on this shortcoming and described how unbiased salient object detection datasets can be constructed. We also reviewed datasets that can be used for saliency model evaluation (in addition to datasets in Table I) and measured their statistics. No dataset exists so far that has all of object annotations, eye movements, and explicit saliency judgments. Bruce-A has fixations, and only explicit saliency judgments but not all object labels. Judd-A, OSIE, and PASCAL-S datasets have annotations and fixations but not explicit saliency judgments. Here, we chose the object that falls at the peak of the fixation map as the most salient one. UCSB dataset lacks object annotations but it has fixations and saliency judgments using clicks (as opposed to object boundaries in Bruce-A). Future research by collecting all information on a large scale dataset will benefit salient object detection research.

Here we suggested that the most salient object in a scene is the one that attracts the majority of fixations (similar to [42]). One can argue that the most salient object is the one that observers look at first. While in general, these two definitions may choose different objects, given the short presentation times in our datasets (3 sec on Judd, 4 sec on Bruce) we suspect that both suggestions will yield to similar results.

Our model separates detection from segmentation. A benefit of this way of modeling is that it can be utilized for other purposes (e.g., segmenting interesting or important objects) by replacing the first component of our model. Further, augmented with a top-down fixation selection strategy, our model can be used as an active observer (e.g., [94]).

Our analysis suggests two main reasons for model performance drop over the Judd-A dataset: The *first reason* that the literature has focused so far is to avoid incorrectly segmenting the object region (i.e., increasing true positives and reducing false positives). Therefore, low performance is partially due to inaccurately highlighting (segmenting) the salient object. The *second reason* that we attempted to highlight in this paper (we believe is the main problem causing performance drop as models performed poorly on Judd-A compared to MSRA-5K) is segmenting the wrong object (i.e., not the most salient object). Note that although here we did not consider the latest proposed salient object detection models in our model comparison (e.g., [16], [29], [45], [92], [93], [95]), we believe that our results are likely to generalize compared to newer models. The rationale is that even recent models have also used the ASD dataset [10] (which is highly center-biased) for model development and testing. Nonetheless, we encourage future works to use our model (as well as Li et al.'s model) as a baseline for model benchmarking.

Two types of cues can be utilized for segmenting an object: appearance [21], [31] (i.e., grouping contiguous areas based on surface similarities) and boundary [35], [36] (i.e., cut regions based on observed pixel boundaries). Here we mainly focused on the appearance features. Taking advantage of both region appearance and contour information (similar to [23], [34]) for saliency detection (e.g., growing the foreground salient region until reaching the object boundary) is an interesting future direction. In this regard, it will be helpful to design suitable measures for evaluating accuracy of models for detecting boundary (e.g., [20]).

Fig. 14: Left: F-measure with superpixel parameters as $\sigma = 1$, K = 500, and min = 50. Right: A sample image and its corresponding superpixel segmentations.
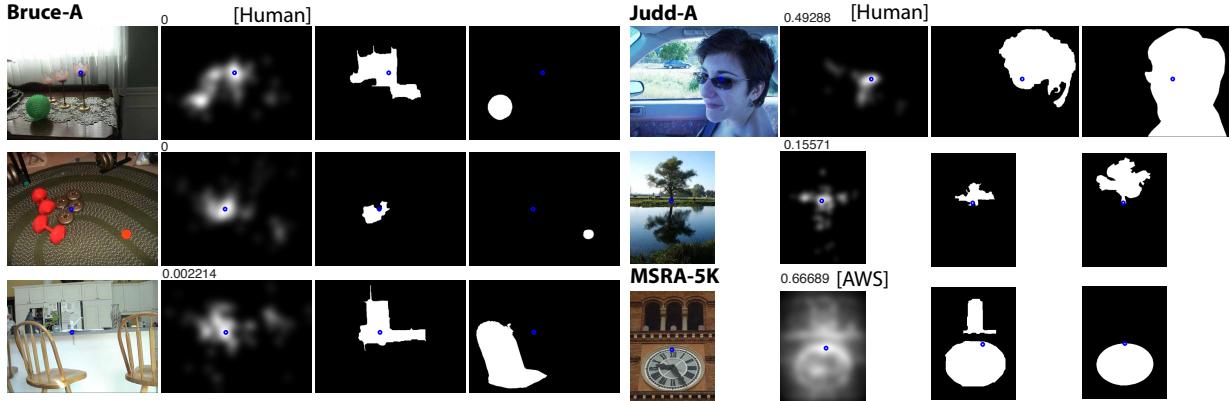


Fig. 15: Sample images where our model (3rd column) fails. Fourth columns show the ground truth. Left: images from the Bruce dataset where our model fails because human saliency judgments and fixations do not agree. Right: failure cases where some superpixels are mistakenly discarded or included.

Our datasets allow more elaborate analysis of the interplay between saliency detection, fixation prediction, and object proposal generation. Obviously, these models depend on the other. On one hand, it is critical to correctly predict where people look to know which object is the most salient one. On the other hand, labeled objects in scenes can help us study how objects guide attention and eye movements. For example, by verifying the hypotheses that some parts of objects (e.g., object center [33]) or semantically similar objects [32]) attract fixations more, better fixation prediction models become feasible.

## REFERENCES

[1] T. Judd, K. Ehinger, F. Durand and, A. Torralba. Learning to predict where humans look. *International Conference on Computer Vision*, 2009.

[2] L. Itti, C. Koch, and E. Niebur. A model of saliency-based visual attention for rapid scene analysis. *IEEE Trans. PAMI*, 1998.

[3] N. D. B. Bruce and J. K. Tsotsos. Saliency based on information maximization. *NIPS*, 2005.

[4] X. Hou and L. Zhang. Dynamic attention: Searching for coding length increments. *NIPS*, 2008.

[5] A. Garcia-Diaz, X. R. Fdez-Vidal, X. M. Pardo, and R. Dosil. Decorrelation and distinctiveness provide with human-like saliency. *ACIVS*, 2009.

[6] B. W. Tatler. The central fixation bias in scene viewing: selecting an optimal viewing position independently of motor bases and image feature distributions. *J. Vision*, 14(7): 2007.

[7] B.C. Russell, A. Torralba, K.P. Murphy, and W.T. Freeman. LabelMe: A database and web-based tool for image annotation. *International Journal of Computer Vision*, 2008.

[8] M. M. Cheng, G. X. Zhang, N. J. Mitra, X. Huang, and S. M. Hu. Global contrast based salient region detection. *IEEE Conference on Computer vision and Pattern Recognition*, 2011.

[9] S. Goferman, L. Zelnik-Manor, and A. Tal. Context-aware saliency detection. *Computer vision and Pattern Recognition (CVPR)*, 2010.

[10] R. Achanta, S. Hemami, F. Estrada, and S. Süsstrunk. Frequency-tuned salient region detection. *Computer vision and Pattern Recognition (CVPR)*, 2009.

[11] M. Wang, J. Konrad, P. Ishwar, Y. Jing, and H. Rowley. Image saliency: from intrinsic to extrinsic context. *IEEE Conference on Computer vision and Pattern Recognition*, 2011.

[12] J. Feng, Y. Wei, L. Tao, C. Zhang, and J. Sun. Salient object detection by composition. *International Conference on Computer Vision*, 2011.

[13] M. Bravo and H. Farid. A scale invariant measure of clutter. *Journal of Vision*, 2008.

[14] K.Y Chang, T.L Liu, H.T Chen, and S. Lai. Fusing generic objectness and visual saliency for salient object detection. *International Conference on Computer Vision*, 2011.

[15] H. Jiang, J. Wang, Z. Yuan, T. Liu, N. Zheng, and S. Li. Automatic salient object segmentation based on context and shape prior. *British Machine Vision Conference*, 2011.

[16] H. Jiang, J. Wang, Z. Yuan, Y. Wu, N. Zheng, and S. Li. Salient object detection: A discriminative regional feature integration approach. *IEEE Conference on Computer vision and Pattern Recognition*, 2013.

[17] T. Liu, J. Sun, N Zheng, X. Tang, and H. Shum. Learning to detect a salient object. *Computer vision and Pattern Recognition (CVPR)*, 2007.

[18] J. Li, M.D. Levine, X. An, and H. He. Saliency detection based on frequency and spatial domain analysis. *BMVC*, 2011.

[19] Q. Yan, L. Xu, J. Shi, and J. Jia. Hierarchical saliency detection. *IEEE Conference on Computer vision and Pattern Recognition*, 2013.

[20] V. Movahedi and J. H. Elder, Design and perceptual validation of performance measures for salient object segmentation. *POCV*, 2010.

[21] P. Felzenszwalb and D. Huttenlocher. Efficient graph-based image segmentation, *International Journal of Computer Vision*, 2004.
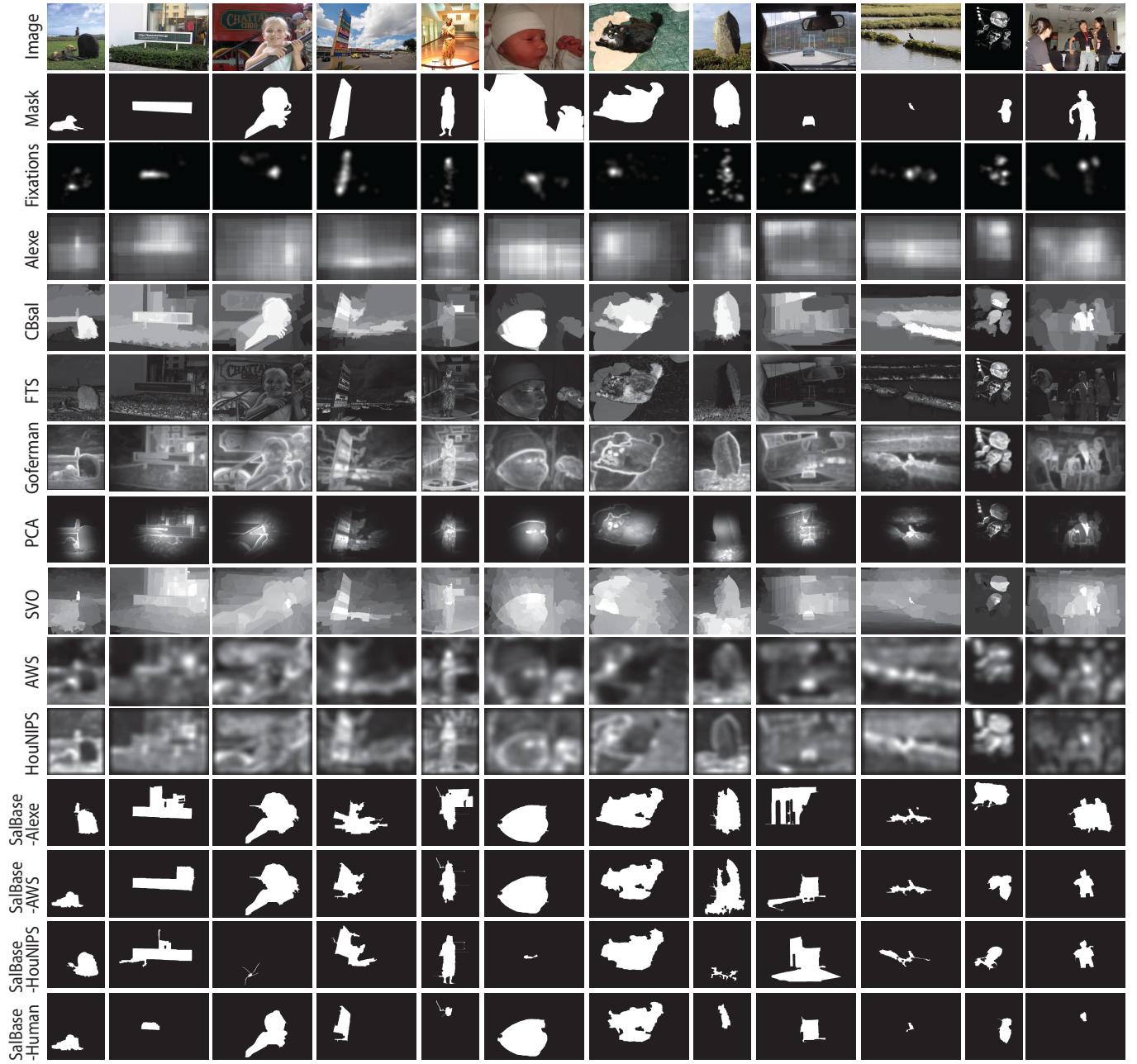
Fig. 16: Sample images from Judd-A dataset along with saliency detection maps of models.

[22] S. Alpert, M. Galun, R. Basri, A. Brandt. Image segmentation by probabilistic bottom-up aggregation and cue integration. *IEEE Conference on Computer vision and Pattern Recognition*, 2007.

[23] A. K. Mishra, Y. Aloimonos, and C. L. Fah. Active segmentation with Fixation. *International Conference on Computer Vision*, 2011.

[24] C. Rother, V. Kolmogorov, and A. Blake. "grabcut": interactive foreground extraction using iterated graph cuts. *ACM Trans. Graph.*, 2004.

[25] A. Borji, and L. Itti. Exploiting local and global patch rarities for saliency detection. *CVPR*, 2012.

[26] A. Borji, M. M., Cheng, H., Jiang, and J., Li. Salient Object Detection: A Survey. *arXiv preprint arXiv:1411.5878*, 2014.

[27] B. Alexe, T. Deselares, and V. Ferrari. What is an object? *IEEE Conference on Computer vision and Pattern Recognition*, 2010.

[28] L. Elazary and L. Itti. Interesting objects are visually salient. *Journal of Vision*, 2008.

[29] F. Perazzi, P. Krähenbühl, Y. Pritch, A. Hornung. Saliency filters: Contrast based filtering for salient region detection. *IEEE Conference on Computer vision and Pattern Recognition*, 2012.

[30] A. Borji, D.N. Sihite, and L. Itti. Salient object detection: A benchamrk. *ECCV*, 2012.

[31] J. Shi and J. Malik. Normalized cuts and image segmentation. *IEEE. Trans. PAMI*, 2000.

[32] A.D. Hwang, H.C. Wang, and M. Pomplun. Semantic guidance of eye movements in real-world scenes. *Vision Research*, 2011.

[33] A. Nuthman and J.M. Henderson. Object-based attentional selection in scene viewing. *Journal of Vision*, 2010.

[34] P. Arbelaez, M. Maire, C. Fowlkes, and J. Malik. Contour detection and hierarchical image segmentation. *IEEE T-PAMI*, 2011.

[35] D. Martin, C. Fowlkes, J. Malik. Learning to detect natural image boundaries using local brightness, color and texture cues. *PAMI*, 2004.

[36] X. Ren and C. C. Fowlkes and J. Malik. Scale-invariant contour completion using conditional random field. *International Conference on Computer Vision*, 2003.

[37] R. Margolin, L. Zelnik-Manor, and A. Tal. What Makes a Patch Distinct? *IEEE Conference on Computer vision and Pattern Recognition*, 2013.

[38] X. Shen and Y. Wu. A unified approach to salient object detection via

low rank matrix recovery. *IEEE Conference on Computer vision and Pattern Recognition*, 2012.

[39] A. Borji, D.N. Sihite, and L. Itti. Quantitative analysis of human-model agreement in visual saliency modeling: a comparative study. *IEEE Trans. Image Processing*, 2012.

[40] D. Batra, A. Kowdle, D. Parikh, J. Luo, T. Chen. icoseg: Interactive co-segmentation with intelligent scribble guidance. *IEEE Conference on Computer vision and Pattern Recognition*, 2010.

[41] D. Parkhurst, K. Law, and E. Niebur. Modeling the role of salience in the allocation of overt visual attention. *Vision Research*, 2002.

[42] Y. Li, X. Hou, C. Koch, J. M. Rehg, and A. L. Yuille. The Secrets of Salient Object Segmentation. *IEEE Conference on Computer vision and Pattern Recognition*, 2014.

[43] J. Xu, M. Jiang, S. Wang, M. S. Kankanhalli, Q. Zhao. Predicting human gaze beyond pixels. *Journal of Vision*, vol. 14, no. 1, pp: 1:20, 2014.

[44] K. Koehler, F. Guo, S. Zhang, M. P. Eckstein. What do saliency models predict?. *Journal of Vision*, vol. 14, no. 3, 2014.

[45] Z. Yang, D. Li, J. Wang, and X. Li. Saliency detection based on manifold learning. *Eighth International Symposium on Multispectral Image Processing and Pattern Recognition*, 2013.

[46] M. Spain and P. Perona, Measuring and predicting object importance, International Journal of Computer Vision, vol. 91, no. 1, pp. 59-76, 2011.

[47] A. C. Berg, T. L. Berg, H. Daume, J. Dodge, A. Goyal, X. Han, A. Mensch, M. Mitchell, A. Sood, K. Stratos et al., Understanding and predicting importance in images, in IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2012, pp. 3562-3569.

[48] X. Han, A. Mensch, M. Mitchell, A. Sood, K. Stratos et al., Understanding and predicting importance in images, in IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2012, pp. 3562-3569.

[49] P. Isola, J. Xiao, A. Torralba, and A. Oliva, What makes an image memorable? in IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2011, pp. 145-152.

[50] R. Rosenholtz, Y. Li, and L. Nakano, Measuring visual clutter, Journal of Vision, vol. 7, no. 2, 2007.

[51] Y.-G. Jiang, Y. Wang, R. Feng, X. Xue, Y. Zheng, and H. Yang, Understanding and predicting interestingness of videos, AAAI, 2013.

[52] H. Katti, K. Y. Bin, T. S. Chua, and M. Kankanhalli, Preattentive discrimination of interestingness in images, in IEEE International Conference on Multimedia and Expo (ICME), 2008, pp. 1433-1436.

[53] M. Gygli, H. Grabner, H. Riemenschneider, F. Nater, and L. Van Gool, The interestingness of images, IEEE International Conference on Computer Vision (ICCV), 2013.

[54] S. Dhar, V. Ordonez, and T. L. Berg, High level describable attributes for predicting aesthetics and interestingness, in IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2011, pp. 1657-1664.

[55] L. Itti and P. Baldi, Bayesian surprise attracts human attention, in NIPS, 2005, pp. 547-554.

[56] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, Image quality assessment: From error visibility to structural similarity, IEEE Transactions on Image Processing, vol. 13, no. 4, pp. 600-612, 2004.

[57] K. A. Ehinger, J. Xiao, A. Torralba, and A. Oliva, Estimating scene typicality from human ratings and image features, in Proceedings of the 33rd Annual Cognitive Science Conference, 2011.

[58] J. Vogel and B. Schiele, A semantic typicality measure for natural scene categorization, *Pattern Recognition*, vol. 3175, pp. 195-203, 2004.

[59] A. Farhadi, I. Endres, D. Hoiem, and D. Forsyth, Describing objects by their attributes, in IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2009, pp. 1778-1785.

[60] D. Parikh and K. Grauman, Relative attributes. *In IEEE International Conference on Computer Vision (ICCV)*, pp. 503-510, 2011.

[61] X. Hou, J. Harel, and C. Koch, Image signature: Highlighting sparse salient regions, IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 34, no. 1, pp. 194-201, 2012.

[62] H. J. Seo and P. Milanfar, Static and space-time visual saliency detection by self-resemblance, Journal of Vision, vol. 9, no. 12, pp. 15, 1-27, 2009.

[63] E. Erdem and A. Erdem, Visual saliency estimation by nonlinearly integrating features using region covariances, Journal of Vision, vol. 13, no. 4, pp. 11, 1-20, 2013.

[64] C. Masciocchi, S. Mihalas, D. Parkhurst, and E. Niebur, Everyone knows what is interesting: Salient locations which should be fixated. Journal of Vision, vol. 9, pp. 1-22, 2009.

[65] A. Borji, D. N. Sihite, and L. Itti, What stands out in a scene? a study of human explicit saliency judgment, Vision Research, vol. 91, no. 0, pp. 62-77, 2014.

[66] L. Itti and M. A. Arbib, Attention and the minimal subscene, Action to language via the mirror neuron system, pp. 289-346, 2006.

[67] G. Kulkarni, V. Premraj, S. Dhar, S. Li, Y. Choi, A. C. Berg, and T. L. Berg, Baby talk: Understanding and generating simple image descriptions, in IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2011, pp. 1601-1608.

[68] U. Rutishauser, D. Walther, C. Koch, and P. Perona, Is bottom-up attention useful for object recognition? in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2004, vol. 2, 2004, pp. II-37.

[69] J. Uijlings, K. van de Sande, T. Gevers, and A. Smeulders, Selective search for object recognition, International Journal of Computer Vision, vol. 104, no. 2, pp. 154-171, 2013.

[70] A. Borji and L. Itti, Scene classification with a sparse set of salient regions, in IEEE International Conference on Robotics and Automation (ICRA), 2011, pp. 1902-1908.

[71] C. Guo and L. Zhang, A novel multiresolution spatiotemporal saliency detection model and its applications in image and video compression, IEEE Trans. on Image Processing, vol. 19, no. 1, pp. 185-198, 2010.

[72] L. Itti, Automatic foveation for video compression using a neurobiological model of visual attention, IEEE Transactions on Image Processing, vol. 13, no. 10, pp. 1304-1318, 2004.

[73] Y.-F. Ma, X.-S. Hua, L. Lu, and H.-J. Zhang, A generic framework of user attention model and its application in video summarization, IEEE Transactions on Multimedia, vol. 7, no. 5, pp. 907-919, 2005.

[74] P. Bodesheim, Spectral clustering of rois for object discovery, in DAGM-Symposium, ser. Lecture Notes in Computer Science, R. Mester and M. Felsberg, Eds., vol. 6835, 2011, pp. 450-455.

[75] S. Goferman, A. Tal, and L. Zelnik-Manor, Puzzle-like collage, Computer Graphics Forum, vol. 29, no. 2, pp. 459-468, 2010.

[76] J. Wang, L. Quan, J. Sun, X. Tang, and H.-Y. Shum, Picture collage, in IEEE Conference on Computer Vision and Pattern Recognition (CVPR), vol. 1, 2006, pp. 347-354.

[77] A. Ninassi, O. Le Meur, P. Le Callet, and D. Barbba, Does where you gaze on an image affect your perception of quality? applying visual attention to image quality metric, in IEEE Conference on Image Processing (ICIP), vol. 2, 2007, pp. II-169.

[78] H. Liu and I. Heynderickx, Studying the added value of visual attention in objective image quality metrics based on eye movement data, ICIP, pp. 3097-3100, 2009.

[79] M. Donoser, M. Urschler, M. Hirzer, and H. Bischof, Saliency driven total variation segmentation, in IEEE International Conference on Computer Vision (ICCV), 2009, pp. 817-824.

[80] T. Chen, M.-M. Cheng, P. Tan, A. Shamir, and S.-M. Hu, Sketch2photo: internet image montage, ACM Trans. on Graphics, vol. 28, no. 5, p. 124, 2009.

[81] S. Feng, D. Xu, and X. Yang, Attention-driven salient edge (s) and region (s) extraction with application to cbir, Signal Processing, vol. 90, no. 1, pp. 1-15, 2010.

[82] A. Borji, S. Frintrop, D. N. Sihite, and L. Itti. Adaptive object tracking by learning background context. *In Computer Vision and Pattern Recognition Workshops (CVPRW), CVPR (pp. 23-30).* 2012.

[83] A. Y.-S. Chia, S. Zhuo, R. K. Gupta, Y.-W. Tai, S.-Y. Cho, P. Tan, and S. Lin, Semantic colorization with internet images, ACM Transactions on Graphics, vol. 30, no. 6, p. 156, 2011.

[84] H. Liu, L. Zhang, and H. Huang, Web-image driven best views of 3d shapes, The Visual Computer, vol. 28, no. 3, pp. 279-287, 2012

[85] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, A. Zisserman, The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88(2), 303-338, 2010.

[86] S. Stalder, H. Grabner, and L. Van Gool, Dynamic objectness for adaptive tracking, ACCV, pp. 43-56, 2013.

[87] J. Li, M. Levine, X. An, X. Xu, and H. He, Visual saliency based on scale-space analysis in the frequency domain, IEEE Trans. on Pattern Analysis and Machine Intelligence, vol. 35, no. 4, pp. 996-1010, 2013.

[88] A. Karpathy, S. Miller, and L. Fei-Fei, Object discovery in 3d scenes via shape analysis, ICRA, 2013, pp. 2088-2095.

[89] A., Borji, H. R., Tavakoli, D. N., Sihite, and L. Itti. Analysis of scores, datasets, and models in visual saliency prediction. *ICCV*, pp. 921-928. 2013.

[90] S. Frintrop, G. M. Garca, and A. B. Cremers, A cognitive approach for object discovery, ICPR, 2014.

[91] D. Meger, P.-E. Forssen, K. Lai, S. Helmer, S. McCann, T. Southey, M. Baumann, J. J. Little, and D. G. Lowe, Curious george: An attentive semantic robot, Robotics and Autonomous Systems, vol. 56, no. 6, pp. 503-511, 2008.

[92] B. Jiang, L. Zhang, H. Lu, C. Yang, and M.-H. Yang, Saliency detection via absorbing markov chain, in ICCV, 2013.

[93] P. Jiang, H. Ling, J. Yu, and J. Peng, Salient region detection by ufo: Uniqueness, focusness and objectness, in ICCV, 2013.

[94] A. Borji, D. N. Sihite, and L. Itti. What/where to look next? Modeling top-down visual attention in complex interactive environments, *IEEE Trans. SMC, Part A*.

[95] J. Kim, D. Han, Y.-W. Tai, and J. Kim, Salient region detection via high-dimensional color transform, in CVPR, 2014.

**Ali Borji** received his BS and MS degrees in computer engineering from Petroleum University of Technology, Tehran, Iran, 2001 and Shiraz University, Shiraz, Iran, 2004, respectively. He did his Ph.D. in cognitive neurosciences at Institute for Studies in Fundamental Sciences (IPM) in Tehran, Iran, 2009 and spent four years as a postdoctoral scholar at iLab, University of Southern California from 2010 to 2014. He is currently an assistant professor at University of Wisconsin, Milwaukee. His research interests include visual attention, active learning, object and scene recognition, and cognitive and computational neurosciences.