





Saliency Benchmarking Made Easy: Separating Models, Maps and Metrics

Matthias Kümmerer¹(✉) , Thomas S. A. Wallis^{1,2} , and Matthias Bethge¹

¹ Werner Reichardt Centre for Integrative Neuroscience, University of Tübingen,
Tübingen, Germany

{matthias.kuemmerer,tom.wallis,matthias}@bethgelab.org

² Wilhelm-Schickard Institute for Computer Science (Informatik),
University of Tübingen, Tübingen, Germany

Abstract. Dozens of new models on fixation prediction are published every year and compared on open benchmarks such as MIT300 and LSUN. However, progress in the field can be difficult to judge because models are compared using a variety of inconsistent metrics. Here we show that no single saliency map can perform well under all metrics. Instead, we propose a principled approach to solve the benchmarking problem by separating the notions of saliency models, maps and metrics. Inspired by Bayesian decision theory, we define a saliency model to be a probabilistic model of fixation density prediction and a saliency map to be a metric-specific prediction derived from the model density which maximizes the expected performance on that metric given the model density. We derive these optimal saliency maps for the most commonly used saliency metrics (AUC, sAUC, NSS, CC, SIM, KL-Div) and show that they can be computed analytically or approximated with high precision. We show that this leads to consistent rankings in all metrics and avoids the penalties of using one saliency map for all metrics. Our method allows researchers to have their model compete on many different metrics with state-of-the-art in those metrics: “good” models will perform well in all metrics.

Keywords: Saliency · Benchmarking · Metrics · Fixations
Bayesian decision theory · Model comparison

1 Introduction

Humans have a foveated visual system: only a small central part of the retina has high receptor density allowing the perception of the details of a scene. Therefore humans make eye movements to place the high resolution fovea on things they want to see. Understanding where they choose to look is therefore an important component of understanding behaviour.

Electronic supplementary material The online version of this chapter (https://doi.org/10.1007/978-3-030-01270-0_47) contains supplementary material, which is available to authorized users.

A long-standing account of bottom-up attentional guidance posits the existence of a “saliency map” (or maps) in the human brain [26, 48]. Here, a saliency map represents spatial importance, usually defined to be local contrast in low-level features such as luminance, color or orientation. Since Itti and Koch formulated this concept into their seminal image-based model [17], a large number of models have been proposed for predicting fixations from image features, e.g. [1, 6, 15, 24, 25, 55, 56] and more recently many models based on deep learning, e.g. [16, 28, 30, 31, 36, 49]; see [4, 19] for extensive reviews of the literature. New models are published on a regular basis with contributions coming mainly from the communities of computer vision and psychology. It has been extensively discussed which effects are important for fixation prediction, from low and high-level influences [3, 7, 12, 14, 18, 31, 50] to biases [8, 44–46], tasks [27, 41, 43] and semantic effects [11]. Over time, the concept of a saliency map has moved away from its origins in low-level feature integration, and can now refer more generally to “a map that predicts fixations”. In practice, saliency maps are now synonymous with saliency models.

The large number of models created the need for quantitative metrics to assess progress in the field and compare models. Many different metrics have been proposed. The AUC-type metrics [45] used to be most common while the last years have seen a shift towards metrics like CC [22], NSS [37] and SIM [23], and recently the information gain metric has been proposed [32]. For an overview of the different metrics in use see e.g. [4, 23]. The community uses these metrics in benchmarks to keep track of the progress: the MIT Saliency Benchmark [9, 23] and the LSUN Challenge [21, 52–54].

The most widely accepted MIT benchmark evaluates submissions in eight different metrics. Depending on which metric one chooses, the model rankings and performances change dramatically. This fact has led to substantial research analyzing the differences between metrics and giving recommendations in which situation to use which metric [10, 33, 38–40, 51]. Other authors have instead proposed new approaches to modeling and evaluation: Modeling as point processes [2, 42], other loss functions [20] and GLMMs [35].

The general conclusion in the field is that the metrics measure qualitatively different things [10, 40, 51], and that it is even conceptually impossible to determine a best model independent of the different metrics. Recently, Kümmerer et al. [32] tried to argue for a unique ranking between different models by showing that much of the disagreement between different metrics can be removed via postprocessing of the saliency maps by optimizing the saliency scale and smoothing kernel for information gain (IG, essentially log-likelihood).

However, this does not seem to be a satisfactory solution: For one, this approach requires access to all models one wants to compare to and needs tedious postprocessing for each of them. In addition to this practical barrier the approach also suffers from the major conceptual shortcoming that optimizing for IG cannot be optimal for all metrics. In fact, we show below that the log densities proposed in [32] perform suboptimally on most metrics and can still produce inconsistent rankings. Ideally one would like a model to be able to compete in all metrics on the metric’s original scale with other models, even with models that are directly

optimized for that metric and where only the metric performances are known. This is not possible when evaluating on log densities as proposed in [32].

In fact, we show in this paper that even with knowledge of the true fixation distribution, no single saliency map can perform well in all metrics. In practice however, researchers must still decide on a particular saliency map to submit to the benchmark. Therefore, their model cannot compete with state-of-the-art models in all metrics – not because the model is intrinsically bad on those metrics, but because different metrics require the saliency maps to look different, independent of the encoded information about fixation placement (see Fig. 1). As long as one evaluates all saliency metrics on the *same* saliency maps, it is impossible to solve the benchmarking problem.

Here, we argue that the fundamental problem is that saliency models and saliency maps are considered to be the same. A major insight from Bayesian decision theory is that the derivation of optimal decisions can be decomposed into a task-independent probability distribution over possible outcomes of an experiment and a task-dependent error metric. In the saliency setting, one decides on a saliency map to submit to a certain metric. Correspondingly, saliency *models* should be defined as *metric-independent* probability densities over possible fixations and subsequently many different *metric-dependent* saliency *maps* can be derived from the same density for different error metrics.

We show that saliency maps for the most influential metrics AUC, sAUC, NSS, CC, SIM, and KL-Div can be derived from fixation densities in a principled way. We demonstrate the validity of our approach on real models and real data. By decoupling the notions of saliency models and saliency maps, saliency models can be meaningfully compared on all metrics *in their original scale*, and the MIT saliency benchmark will implement our suggested approach.

2 Theory

Motivated by the line of thoughts presented above we here propose to use the following definitions:

1. a *saliency model* predicts a fixation probability density $p(x, y | I)$ given an image I .
2. a *saliency metric* is a performance measure for a saliency map on ground truth data.
3. a *saliency map* $s_{p, \text{metric}}(x, y, I)$ is a metric-specific prediction derived from the model density.

It has been argued before that formulating saliency models as probabilistic models is advantageous (e.g. [2, 32]). In this definition, a saliency model predicts a fixation probability density, that is, the probability $p(x, y | I)$ of observing a fixation at a given pixel in a given image¹. The three definitions we propose

¹ Note that we use the fixation probability density for single fixations (as in [32]) whereas [2] define a point process density for a whole scanpath.

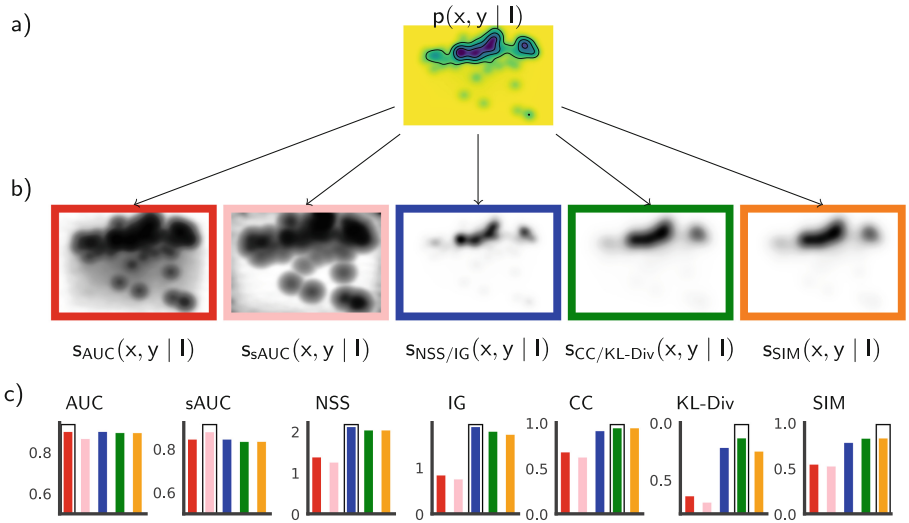


Fig. 1. No single saliency map can perform best in all metrics even when the true fixation distribution is known. This problem can be solved by separating saliency models from saliency maps. **(a)** Fixations are distributed according to a ground truth fixation density $p(x, y | I)$ for some stimulus I (see supplementary material for details on the visualization). **(b)** This ground truth density predicts different saliency maps depending on the intended metric. The saliency maps differ dramatically due to the different properties of the metrics but always reflect the same underlying model. Note that the maps for the NSS and IG metrics are the same, as are those for CC and KL-Div. **(c)** Performances of the saliency maps from (b) under seven saliency metrics on a large number of fixations sampled from the model distribution in (a). Colors of the bars correspond to the frame colors in (b). The predicted saliency map for the specific metric (framed bar) yields best performance in all cases.

above follow the rationale of Bayesian decision theory: the saliency model is a posterior density over all possible events and the saliency metric is a utility function. Based on the posterior density and the utility function, a saliency map is then chosen to maximize the expected utility.

2.1 Predicting Saliency Maps from Saliency Models

From the predicted fixation density of a model, one can use expected utility maximization to derive the saliency map which the model expects to yield highest performance in some metric².

Evaluating a saliency metric involves a saliency map $s(x, y | I)$ for a stimulus I and ground truth fixation data (x_i, y_i) . Therefore, we can phrase a metric

² Note that the term “metric” is a slight abuse of notation: strictly speaking, a metric measures the distance between two objects and is usually desired to be minimal. However, in saliency, the term “metric” denotes the performance that one wants to maximize (with a few exceptions, e. g., KL-Div and earth mover’s distance).

as a function $M[s(x, y | I); (x_1, y_1), \dots, (x_n, y_n)]$. Note that some metrics as CC or SIM use an empirical saliency map instead of ground truth fixations (*distribution-based metrics, r1cheSaliency2013*). However, the empirical saliency map is always constructed from ground truth fixations, usually by convolving them with a Gaussian. This can be taken to be part of the metric evaluation, as we will demonstrate below. Simplifying notation with $D = (x_1, y_1), \dots, (x_n, y_n)$, the metric evaluation can be written as

$$M[s(x, y | I); D].$$

Assuming that the fixations are distributed according to some distribution $(x_i, y_i) \sim p(x, y | I)$ and therefore $D \sim \prod_1^n p(x, y)$, the expected performance of the metric on a saliency map is $\mathbb{E}_D M[s(x, y | I); D]$. One should choose the saliency map which is expected to yield highest performance for the metric M : that is, the solution of

$$\max_{s(x, y | I)} \mathbb{E}_D M[D, s(x, y | I)]$$

Solving this optimization problem for a fixation distribution p given by a model of interest essentially answers the following question: if we assume that the unknown fixations, on which the saliency map later will be evaluated, come from the model density p (and therefore $D = \prod_i^n p$), what would be the best saliency map to use for metric M ? For a metric M the solutions to the optimization problem give rise to a transformation $p(x, y | I) \mapsto s_M(x, y | I)$ from fixation densities to derived metric-specific saliency maps. While the optimization problem might be hard in general, for most commonly-used saliency metrics it can be solved exactly or approximately, as we show below. Importantly, the methods we outline here are deterministic transformations depending only on the model’s density prediction. No optimization using ground truth data is necessary.

In the following we give exact or approximate solutions for six of the most widely used metrics, including three metrics which operate directly on ground truth fixations (AUC, sAUC, and NSS) and three distribution-based metrics which first convert the ground truth fixations into a empirical saliency map (CC, SIM, KL-Div). Additionally we include the IG metric introduced in [32] since we use this metric for converting existing saliency map models to probabilistic models.

AUC, sAUC. The AUC-type metrics (“Area Under the Curve”, [45]) measure the model performance in a 2AFC (2 alternative forced choice) task where the model has to decide which one of two locations has been fixated: in a 2AFC task, a system is presented with one signal and one noise stimulus and chooses which stimulus is the “signal”. In the case of the AUC in saliency, signal and noise correspond to fixated and non-fixated image locations respectively (See supplementary material for a proof of the equivalence between the ROC curve and the 2AFC task). Denoting the model’s fixation distribution $p_{\text{fix}}(x, y)$, the nonfixation distribution $p_{\text{nonfix}}(x, y)$ (which is uniform for AUC and the image independent center bias for sAUC) and denote the two locations by (x_1, y_1) resp. (x_2, y_2) . The 2AFC

task reduces to deciding whether these points are sampled from $p_{\text{fix}} \times p_{\text{nonfix}}$ or from $p_{\text{nonfix}} \times p_{\text{fix}}$. The likelihoods of the two points given these two distributions are $p_{\text{fix}}(x_1, y_1)p_{\text{nonfix}}(x_2, y_2)$ resp. $p_{\text{nonfix}}(x_1, y_1)p_{\text{fix}}(x_2, y_2)$. The model expects optimal performance by choosing the distribution which has higher likelihood, or equivalently, the point for which $p_{\text{fix}}(x, y)/p_{\text{nonfix}}(x, y)$ has the higher value. Therefore the model should expect the saliency map $p_{\text{fix}}(x, y)/p_{\text{nonfix}}(x, y)$ to yield highest performance. In the special case of the standard AUC metric, p_{nonfix} is constant and the saliency map boils down to p_{fix} . An additional practical consideration is that the MIT benchmark currently only accepts submissions as JPEG images. To compensate for this limited precision and possible JPEG-artefacts, one should additionally histogram-equalize the saliency map (see Supplementary Material).

NSS. The *Normalized Scanpath Saliency* (NSS, [37]) performance of a saliency map model is defined to be the average saliency value of fixated pixels in the normalized (zero mean, unit variance) saliency maps (i.e., the average z-score of the fixated saliency values).

We can show analytically that one should expect the highest NSS score from the predicted fixation density itself: given an image with N pixels let the probability for a single fixation falling onto pixel i be p_i . Then the expected NSS of a saliency map $q = (q_1, \dots, q_N)$ with $\frac{1}{N} \sum_i q_i = \bar{q} = 0$, $\|q\|_2^2 = 1$ is $\sum_i^N p_i \cdot q_i = \langle p, q \rangle$. Finding the saliency map with the best possible NSS is equivalent to finding the solution of the problem

$$\max \langle p, q \rangle \quad \text{s.t.} \quad \bar{q} = 0, \|q\|^2 = 1$$

Since $q \mapsto q' = \bar{p} + \alpha q$ with $\alpha = \sqrt{\|p\|^2 - 1/N}$ induces a maximum-preserving bijection between $\{q \mid \bar{q} = 0, \|q\|^2 = 1\}$ and $\{q' \mid \bar{q}' = \bar{p}/N, \|q'\|^2 = \|p\|^2\}$, we can look for the maximum of $\langle p, q' \rangle \quad \text{s.t.} \quad \bar{q}' = \bar{p}, \|q'\|^2 = \|p\|^2$ instead (and normalize q afterwards to get the normalized saliency map). Because of $\langle x, y \rangle = \frac{1}{2}(\|x\|^2 + \|y\|^2 - \|x - y\|^2)$, the maximum under these conditions is identical with the minimum of $\|p - q\|^2$, which is p .

Therefore, the best possible saliency map with respect to NSS is the density of the fixation distribution.

IG. The *information gain* (IG, [32]) metric requires the saliency map to be a probability distribution and compares the average log-probability of fixated pixels to that given by a baseline model (usually the centerbias or a uniform model). The optimal saliency map for IG depends on how the metric interprets saliency maps as probability densities. We normalize the saliency maps to be probability vectors (nonnegative, unit sum) and in this case the predicted density itself yields the highest expected performance: Let $p = (p_1, \dots, p_N)$ with $p \geq 0$, $\sum_i p_i = 1$ denote the predicted probabilities for each pixel and q with $q \geq 0$, $\sum_i q_i = 1$ a saliency map. Let $p_{bl} = (p_{bl,1}, \dots, p_{bl,N})$ be the pixel probabilities of the baseline model. Then the expected IG of q is $\mathbb{E}_p IG(q) = \sum_i p_i (\log q_i - \log p_{bl,i})$ and its maximum is $\operatorname{argmax}_q \mathbb{E}_p IG(q) = \operatorname{argmax}_q \sum_i p_i (\log q_i - \log p_{bl,i}) = \operatorname{argmax}_q \sum_i p_i \log q_i = \operatorname{argmax}_q \sum_i p_i (\log q_i - \log p_i) = \operatorname{argmin}_q \sum_i p_i (\log p_i - \log q_i) = \operatorname{argmin}_q KL[p, q] = p$.

CC. The *correlation coefficient* (CC, [22]) measures the correlation between model saliency map and empirical saliency map after normalizing both saliency maps to have zero mean and unit variance. This is equivalent to measuring the euclidean distance between the predicted saliency map and the normalized empirical saliency map. The expected euclidean distance to a random variable is minimized by its expectation value. Therefore the optimal saliency map with respect to CC is the expected normalized empirical saliency map.

This shows that predicting the optimal saliency map for CC crucially depends on how the empirical saliency maps are computed. Empirical saliency maps are typically computed by blurring observed fixation positions from eye movement data with a Gaussian kernel of a certain size. In this case the expected empirical saliency map would be $\mathbb{E}_{x_i \sim p} \frac{1}{N} \sum_i G_\sigma(x) = \frac{1}{N} \sum_i \mathbb{E}_{x \sim p} G_\sigma(x) = \frac{1}{N} \sum_i G_\sigma * p = G_\sigma * p$, that is, the density blurred with a Gaussian kernel of size σ .

Unfortunately, the expected empirical saliency map is not the expected normalized empirical saliency map which was earlier shown to be optimal for CC. Normalization involves subtracting the mean and dividing by the standard deviation, and the latter is nonlinear. Effectively, normalizing the variance just changes the weight by which the different empirical saliency maps are averaged in the expectation value. As long as the variances of the different empirical saliency maps don't differ too much, this won't have much of an effect and our simulations suggest that this is the case (Supplementary Material). Therefore, as an approximation to the expected normalized empirical saliency map, we use the expected saliency map in this paper, which is computed by convolving the expected density by a Gaussian.

Obviously, if more involved techniques are used to compute the empirical saliency maps (e.g. cross validation of the kernel size as in [32]), then the expected empirical saliency map is harder or impossible to calculate analytically. However, one can still approximate it numerically by sampling normalized empirical saliency maps from the expected fixation distribution and averaging them.

KL-Div. The KL-Div metric computes the *Kullback-Leibler divergence* between the empirical saliency maps and the model saliency maps after converting both of them into probability distributions (by making them nonnegative and normalizing them to have unit sum) Therefore, unlike for most other metrics, in KL-Div lower values are better.

We can show that for the KL-Div metric, the expected empirical saliency map expects the best performance: let $e = (e_1, \dots, e_N)$ with $e \geq 0$, $\sum_i e_i = 1$ denote the random variable which represents the empirical saliency map and q with $q \geq 0$, $\sum_i q_i = 1$ the model saliency map. Then we are looking for the q which minimizes $\mathbb{E}_p KL[e, q]$. Since $\mathbb{E}_p [KL[e, q]] = \mathbb{E}_p \left[\sum_i e_i \frac{\log e_i}{\log q_i} \right] = \mathbb{E}_p [\sum_i e_i \log e_i] - \sum_i \mathbb{E}_p [e_i] \log q_i$, this is equivalent to finding the maximum of $\sum_i \mathbb{E}_p [e_i] \log q_i$, which is again equivalent to finding the minimum of $\sum_i \mathbb{E}_p [e_i] \log \mathbb{E}_p [e_i] - \sum_i \mathbb{E}_p [e_i] \log q_i = KL[\mathbb{E}_p [e], q]$. This is obviously minimized by $q = \mathbb{E}_p [e]$, the expected empirical saliency map. As for CC, this is the density blurred by the same kernel size as used for the empirical saliency map.

SIM. The *Similarity* (SIM, [23]) metric normalizes the model saliency map and the empirical saliency map to be probability vectors (in the same way as KL-Div) and sums the pixelwise minimum of two saliency maps. As opposed to the CC-metric, which can be interpreted as measuring the l_2 -distance between normalized saliency maps, this effectively measures the l_1 -distance between saliency maps ($\sum_i \min(p_i, q_i) = \sum_i \frac{1}{2} (p_i + q_i - |p_i - q_i|) = 1 - \frac{1}{2} \|p - q\|_1$.) This optimization problem cannot be solved analytically in general. Instead we solve it numerically: we perform a constrained stochastic gradient descend on sets of fixations sampled from the probability density (see Sect. 3 for details). Note that the optimal saliency map for SIM, unlike all other saliency maps presented here, depends on the number of fixations per image (see the Supplement for details on this effect).

3 Experiments and Results

We use the pysaliency toolbox [29] to compute saliency metrics (see Supplement for details). From a probability density over an image we compute five types of saliency maps: **AUC saliency maps** are created by equalizing the probability density to yield a uniform histogram over all pixels. **sAUC saliency maps** are created by dividing the probability density by the center bias density and again equalizing the saliency map to yield a uniform histogram over all pixels. The center bias density was estimated using a Gaussian kernel density estimate over all fixations from the MIT1003 dataset and crossvalidated across images. **NSS/IG saliency maps** are simply the probability density. **CC/KL-Div saliency maps** are calculated by convolving the probability density with a Gaussian kernel with $\sigma = 35px$ (corresponding to 1dva, as commonly used on the MIT1003 dataset). **SIM saliency maps:** We divide the CC saliency map by its sum to normalize it. Starting from there, we perform constrained (nonnegative, unit sum) stochastic gradient descend on fixations sampled from the predicted density to maximize the expected SIM performance (see Supplementary Material for implementation details).

3.1 No Saliency Map to Rule Them All

Here we illustrate using simulated data that even if the true fixation density is known, no single saliency map can win in all saliency metrics. From a fictional fixation density (Fig. 1a) we compute the saliency maps that we predict to be optimal for the seven saliency metrics AUC, sAUC, NSS/IG, CC/KL-Div and SIM (Fig. 1b). We sample 1000 sets of 100 fixations from the fixation density and evaluate all five saliency maps using the seven different saliency metrics on this dataset (Fig. 1c, raw data in the Supplement).

Although the saliency maps in Fig. 1b all are predicted by the same model, they appear visually different: while the AUC saliency map is essentially just the normalized density, the sAUC saliency map removes the center bias contribution (see above). The NSS/IG saliency map is exactly the density and shows

large areas with very low values. The CC/KL-Div saliency map, being a blurred version of the density, is much smoother than the NSS saliency map. The SIM saliency map looks mostly like the CC/KL-Div saliency map but is slightly more sparse.

The ranking of the five saliency maps is highly inconsistent across metrics (Fig. 1c): even with knowledge of the real fixation distribution, no saliency map can be optimal for all saliency metrics. However, each saliency map is optimal for exactly those metrics for which it has been predicted to be optimal (framed bars). This illustrates our main result: By deriving metric-specific saliency maps in a principled way from fixation densities, one model can perform optimally in all metrics. Notice that in current practice, the situation faced by an individual research team is rather to pick from one of the maps in Fig. 1b and be penalized accordingly on other metrics in Fig. 1c.

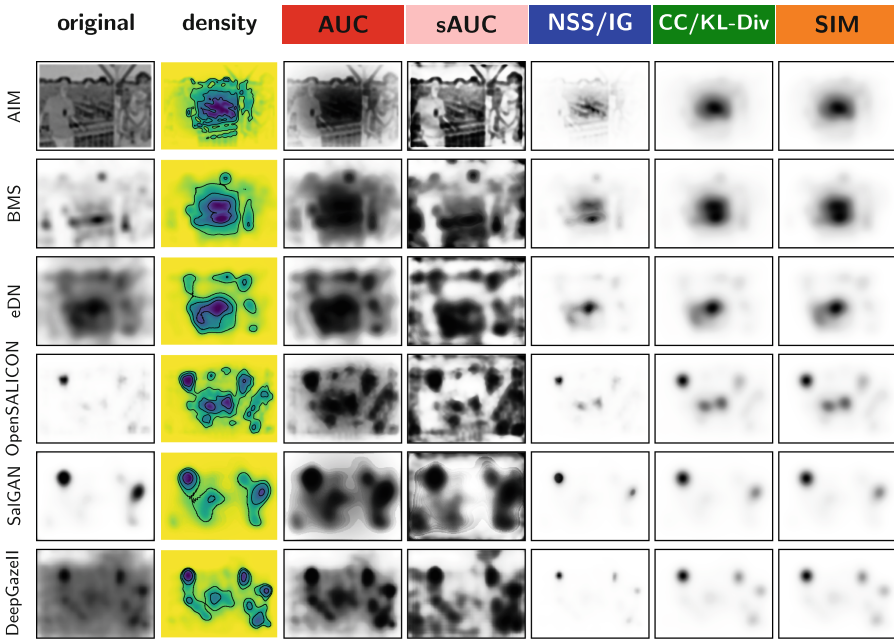


Fig. 2. The predicted saliency map for various metrics according to different models, for the same stimulus. For six models (rows) we show their original saliency map (first column), the probability distribution after converting the model into a probabilistic model (second column) and the saliency maps predicted for seven different metrics (columns three through seven). The predictions of different models for the same metric (column) appear more similar than the predictions of the same model for different metrics (row). In particular, note the inconsistency of the original models (what are typically compared on the benchmark) relative to the per-metric saliency maps. It is therefore difficult to visually compare original model predictions, which have been formulated for different metrics.

3.2 MIT1003

In our main experiment, we use our approach to evaluate six saliency models on the popular benchmarking dataset MIT1003 (freeviewing fixations of 15 subjects on 1003 images, [24]). For all evaluated models, the original source code and default parameters have been used. The included models are **AIM** [6], Boolean Map-based Saliency (**BMS**) [55], the Ensemble of Deep Networks (**eDN**) [49], **OpenSALICON** [47], **SalGAN** [36] and **DeepGaze II** [31].

Converting existing models that produce arbitrary saliency maps into probabilistic models is not straightforward [32]. We used the method described in [32] and implemented in the pysaliency toolbox as `optimize_for_information_gain`: we fitted a pixelwise monotone nonlinearity and a center bias for each model to yield maximum information gain for the MIT1003 dataset (see supplementary material for details). Unlike [32] we did not optimize an additional Gaussian convolution to smooth the predictions. Since DeepGaze II is already formulated as a probabilistic model, there was no need to convert this model. For showing the “original saliency map” we use the log density in this case.

Example saliency maps. In Fig. 2, we show the probability distribution and the predicted saliency maps (columns) for the saliency models (rows) for one example stimulus. Comparing the saliency maps within and between columns, i.e. metrics, one notices that the process of predicting saliency maps for certain metrics has a strong effect on the shape of the saliency maps that is consistent across models. It influences the visual appearance of the saliency map to a larger degree than the actual model does: the AUC and sAUC maps are very high contrast, while the NSS and CC saliency maps have large areas of very little saliency. The CC and SIM saliency maps are much smoother than all other saliency maps. It is a quite common technique in the field to compare the saliency maps of different models visually (e.g., see [13], Figure 6; [5], Figure 6; [4], Figure 9). Figure 2 shows that this technique can be very misleading unless the saliency maps are of the same type (i.e. intended for the same saliency metric).

Comparing model performance. In Fig. 3 we evaluate the saliency maps of the saliency models (AIM, BMS, eDN, OpenSALICON, SalGAN, DeepGaze II; x-axis) on the seven saliency metrics (subplots, raw data in the Supplement). Each line indicates the models’ performances in the evaluated metric when using a specific type of saliency map. The dashed lines indicate performance using the models’ original saliency maps (i.e. not transformed into true probability densities). The performances are very inconsistent between the different metrics on the original saliency maps. The solid lines indicate the metric performances on the five types of derived saliency maps (red: AUC, pink: sAUC, blue: NSS and IG, green: CC and KL-Div, orange: SIM). Additionally, we included log-density saliency maps as proposed in [32] (purple dotted lines).

For each metric, the saliency map predicted for that metric (thick line in each sub plot) yields highest performance for all models. Conversely, saliency maps derived for other metrics often incur severe penalties (except for very few borderline cases, see below). While the model rankings given by the different

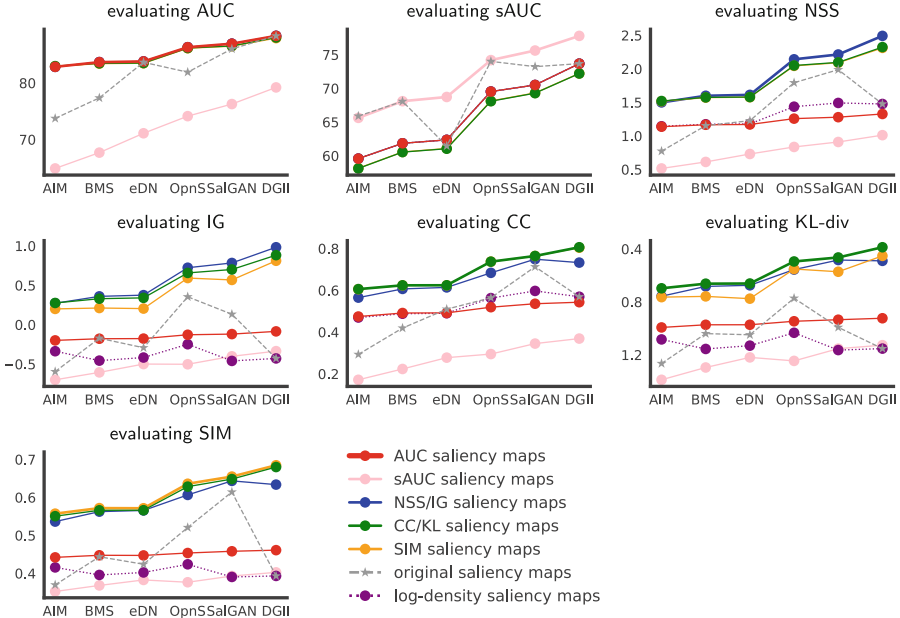


Fig. 3. We reformulated several saliency models in terms of fixation densities and evaluated AUC, sAUC, NSS, IG, CC, KL-Div and SIM on the original saliency maps (dashed line) and the saliency maps derived from the probabilistic model for the different saliency metrics (solid lines) on the MIT1003 dataset. Saliency maps derived for a given metric always yield the highest performance for that metric(thick line), and for each metric the model ranking is consistent when using the correct saliency maps – unlike for the original saliency maps and some other derived saliency maps. Note that AUC metrics yield identical results on AUC saliency maps, NSS saliency maps and log-density saliency maps, therefore the blue and purple lines are hidden by the red line in the AUC and sAUC plots. Also, the CC metric yields only slightly worse results on the SIM saliency map than on the CC saliency map, therefore the orange line is hidden by the green line in the CC plot. OpnS=OpenSALICON, DGII=DeepGaze II.

metrics on each saliency map type are much more consistent than on the original saliency maps, there is still disagreement between metrics left when evaluating all metrics on the same saliency map type.

Interestingly, the AIM model reaches better NSS performance with the CC saliency map than with the NSS saliency map. This is easy to explain: the AIM model’s predicted density improves after blurring. For the better models this effect vanishes. For example, DeepGaze II reaches significantly higher NSS scores with the NSS saliency map than with the CC saliency map and vice versa for the CC metric. The SIM metric seems to show only slightly better performance on the SIM saliency map than on the CC saliency map, with the average difference being just 0.006. However, the best five models with respect to SIM in the MIT

Saliency Benchmark perform within a range of less than 0.02. A difference of 0.006 could easily change a model’s ranking by multiple places.

Figure 3 also serves to illustrate a key difference between the metric unification proposed in [32] and our method of predicting saliency maps from fixation densities: the metric results presented in [32] correspond to the purple dotted log-density lines for AUC, sAUC, NSS and to the blue density lines for IG and KL-Div (in our implementation taking the logarithm of the density is part of the metric itself). As reported in [32], the model rankings are more consistent for those lines than for the original saliency maps. However, except for AUC and IG, in all other metrics the models are penalized when evaluated like this and additionally for the best models even the agreement between metric rankings is lost (SalGAN vs DeepGaze II, AUC/sAUC/IG vs NSS/CC/KL-Div). This shows that the method proposed in [32], while managing to remove a significant amount of the disagreement between metrics, is not perfect.

To summarize, Fig. 3 illustrates the main result of this paper: No matter what saliency map type you decide for, even state-of-the-art models will perform suboptimally in some metrics and rankings will still be inconsistent. Only by using the right saliency map for each metric given the model density, every model performs as well as it can theoretically and all model rankings agree. Consequently, our evaluation yields a unique winner of the benchmark: from all included models, DeepGaze II performs best in all considered metrics.

4 Discussion

Despite much progress in fixation prediction in recent years, comparing saliency models to each other can be confusing due to the large number of benchmarking metrics, giving inconsistent model rankings. Here we argue that benchmarking can be simplified by considering *saliency models* to be probability density predictors, *saliency metrics* to be performance measures that assess saliency maps against ground truth fixations, and subsequently *saliency maps* to be metric-specific predictions derived from the model’s density. We have shown that probabilistic models can predict good saliency maps for the most common saliency metrics: “good models” perform well in many metrics.

Importantly, this metric-specific prediction reflects the same underlying model. It is not the case that the model is being re-trained for each metric. Rather, the saliency maps we show are derived deterministically from the fixation density predicted by a model. In this way it is possible to obtain optimal predictions from a given saliency density for arbitrary metrics without retraining. The saliency model density captures all necessary information in the training data and represents it in a way that it can readily be used in combination with arbitrary error metrics. Information gain (equivalently, log-likelihood) is an ideal optimization metric because it reflects all information in the structure of the fixation density, independent of any particular metric. Therefore, it should lead to good results in all metrics.

The fact that metrics impose strong constraints on saliency maps means that it is misleading to visually compare saliency maps intended for different metrics

(see Fig. 2)—but this is commonly done in the field ([4, 5, 13]) For example, the optimal saliency maps for distribution-based metrics like CC, SIM and KL-Div require blurring unlike those for NSS and IG.

Another consequence of the present work is that the eight metrics available on the MIT benchmark can now be seen as a benefit rather than a possible source of confusion. Since each metric assesses different aspects of the fixation prediction, the benchmark would now allow fair comparison over a number of tasks of interest, which may be more or less relevant for certain applications. For example, sAUC is most relevant when one is interested in a model’s predictive performance once the center bias is excluded (e.g., in applying to a setting with a different center bias from the MIT1003 training data).

While the saliency maps we have derived give the optimal metric-specific saliency map for a given fixation density, it is nevertheless still possible that a given model could do better on a metric with a saliency map not intended for that metric, rather than the metric-specific saliency map itself. If the model’s density is not the correct one (i.e. does not reflect the data-generating density), then the derived saliency maps can be suboptimal. If the model’s density is especially bad, some metrics might even perform better on saliency maps not predicted for this metric than on the one predicted for this metric. For example: if a model’s density prediction is too sparse, the AUC metric will perform better on the smoothed CC saliency map than it will perform on the actual AUC saliency map. Therefore, actually optimizing model predictions for each specific metric may yield insights into the differences between the metrics (by comparing the underlying densities). Indeed, this could in practice produce better performance on the training metric than an information gain optimized density. The fact that we don’t observe this effect on the original saliency maps (which *were* trained in the case of eDN, OpenSALICON, SalGAN and DeepGaze II: Fig. 3, dashed lines) suggests any improvement is likely small, and can come at the price of performing substantially worse in other metrics.

Finally, we would like to note that the distinction between saliency models and saliency maps we draw here does not contradict ideas that a “saliency map” or maps may be instantiated in the human brain, as a corollary of bottom-up attentional guidance or an importance map for (e.g.) choosing the next place to fixate in a scene [26, 34, 48]. Our nomenclature is rather independent and intended for saliency model benchmarking.

The code for evaluating saliency models as demonstrated in this work has been released as part of the `pysaliency` python library (available at <https://github.com/matthias-k/pysaliency>).

Conclusion. Our work solves the problem that one saliency model cannot reach state-of-the-art performance in all relevant saliency metrics. Our key theoretical contribution is to decouple the notions of saliency models and saliency maps. For benchmarking practice, this means that saliency models can be meaningfully compared on all metrics *in their original scale*. Therefore, our method allows comparing to traditional models that do not use this method; it works even if only metric scores of other models are known (as for example in cases where

metric scores are published in a paper). Practically, this means that there is no need to revise an existing benchmark: researchers who submit model densities can have their performance fairly evaluated, but existing models can remain in the table. The MIT saliency benchmark will implement this option.

Acknowledgements. This study is part of Matthias Kümmerer’s thesis work at the International Max Planck Research School for Intelligent Systems (IMPRS-IS). The research has been funded by the German Science Foundation (DFG; Collaborative Research Centre 1233) and the German Excellency Initiative (EXC307).

References

1. Adeli, H., Vitu, F., Zelinsky, G.J.: A model of the superior colliculus predicts fixation locations during scene viewing and visual search. *J. Neurosci.* **37**(6), 1453–1467 (2016). <https://doi.org/10.1523/jneurosci.0825-16.2016>
2. Barthelme, S., Trukenbrod, H., Engbert, R., Wichmann, F.: Modeling fixation locations using spatial point processes. *J. Vis.* **13**(12), 1–1 (2013). <https://doi.org/10.1167/13.12.1>
3. Borji, A., Sihite, D.N., Itti, L.: Objects do not predict fixations better than early saliency: a re-analysis of einhauser et al.’s data. *J. Vis.* **13**(10), 18–18 (2013). <https://doi.org/10.1167/13.10.18>
4. Borji, A., Itti, L.: State-of-the-art in visual attention modeling. *IEEE Trans. Pattern Anal. Mach. Intell.* **35**(1), 185–207 (2013). <https://doi.org/10.1109/tpami.2012.89>
5. Borji, A., Sihite, D.N., Itti, L.: Quantitative analysis of human-model agreement in visual saliency modeling: a comparative study. *IEEE Trans. Image Process.* **22**(1), 55–69 (2013). <https://doi.org/10.1109/tip.2012.2210727>
6. Bruce, N.D.B., Tsotsos, J.K.: Saliency, attention, and visual search: an information theoretic approach. *J. Vis.* **9**(3), 5–5 (2009). <https://doi.org/10.1167/9.3.5>
7. Bruce, N.D.B., Catton, C., Janjic, S.: A deeper look at saliency: Feature contrast, semantics, and beyond. In: 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). IEEE (2016). <https://doi.org/10.1109/cvpr.2016.62>
8. Bruce, N.D., Wloka, C., Frosst, N., Rahman, S., Tsotsos, J.K.: On computational modeling of visual saliency: examining what’s right, and what’s left. *Vis. Res.* **116**, 95–112 (2015). <https://doi.org/10.1016/j.visres.2015.01.010>
9. Bylinskii, Z., Judd, T., Durand, F., Oliva, A., Torralba, A.: MIT saliency benchmark. <http://saliency.mit.edu/>
10. Bylinskii, Z., Judd, T., Oliva, A., Torralba, A., Durand, F.: What do different evaluation metrics tell us about saliency models? [cs] (2016), [arXiv:1604.03605](https://arxiv.org/abs/1604.03605)
11. Bylinskii, Z., Recasens, A., Borji, A., Oliva, A., Torralba, A., Durand, F.: Where should saliency models look next? In: Leibe, B., Matas, J., Sebe, N., Welling, M. (eds.) ECCV 2016. LNCS, vol. 9909, pp. 809–824. Springer, Cham (2016). https://doi.org/10.1007/978-3-319-46454-1_49
12. Cerf, M., Harel, J., Huth, A., Einhäuser, W., Koch, C.: Decoding what people see from where they look: predicting visual stimuli from scanpaths. In: Paletta, L., Tsotsos, J.K. (eds.) WAPCV 2008. LNCS (LNAI), vol. 5395, pp. 15–26. Springer, Heidelberg (2009). https://doi.org/10.1007/978-3-642-00582-4_2
13. Cornia, M., Baraldi, L., Serra, G., Cucchiara, R.: Predicting human eye fixations via an LSTM-based saliency attentive model. [cs] (2016), [arXiv:1611.09571](https://arxiv.org/abs/1611.09571)

14. Einhauser, W., Spain, M., Perona, P.: Objects predict fixations better than early saliency. *J. Vis.* **8**(14), 18–18 (2008). <https://doi.org/10.1167/8.14.18>
15. Harel, J., Koch, C., Perona, P.: Graph-based visual saliency. In: *Advances in neural information processing systems*, pp. 545–552 (2006)
16. Huang, X., Shen, C., Boix, X., Zhao, Q.: SALICON: Reducing the semantic gap in saliency prediction by adapting deep neural networks. In: *2015 IEEE International Conference on Computer Vision (ICCV)*. IEEE (2015). <https://doi.org/10.1109/iccv.2015.38>
17. Itti, L., Koch, C., Niebur, E.: A model of saliency-based visual attention for rapid scene analysis. *IEEE Trans. Pattern Anal. Mach. Intell.* **20**(11), 1254–1259 (1998). <https://doi.org/10.1109/34.730558>
18. Itti, L.: Quantifying the contribution of low-level saliency to human eye movements in dynamic scenes. *Vis. Cogn.* **12**(6), 1093–1123 (2005). <https://doi.org/10.1080/13506280444000661>
19. Itti, L., Borji, A.: *Computational models: Bottom-up and top-down aspects*. The Oxford Handbook of Attention. Oxford University Press, Oxford (2014)
20. Jetley, S., Murray, N., Vig, E.: End-to-end saliency mapping via probability distribution prediction. In: *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE (2016). <https://doi.org/10.1109/cvpr.2016.620>
21. Jiang, M., Huang, S., Duan, J., Zhao, Q.: SALICON: saliency in context. In: *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE (2015). <https://doi.org/10.1109/cvpr.2015.7298710>
22. Jost, T., Ouerhani, N., Wartburg, R.V., Müri, R., Hügli, H.: Assessing the contribution of color in visual attention. *Comput. Vis. Image Underst.* **100**(1–2), 107–123 (2005). <https://doi.org/10.1016/j.cviu.2004.10.009>
23. Judd, T., Durand, F.d., Torralba, A.: A benchmark of computational models of saliency to predict human fixations. *CSAIL Technical reports* (2012). 1721.1/68590
24. Judd, T., Ehinger, K., Durand, F., Torralba, A.: Learning to predict where humans look. In: *2009 IEEE 12th International Conference on Computer Vision*. IEEE (2009). <https://doi.org/10.1109/iccv.2009.5459462>
25. Kienzle, W., Franz, M.O., Scholkopf, B., Wichmann, F.A.: Center-surround patterns emerge as optimal predictors for human saccade targets. *J. Vis.* **9**(5), 7–7 (2009). <https://doi.org/10.1167/9.5.7>
26. Koch, C., Ullman, S.: Shifts in selective visual attention: towards the underlying neural circuitry. *Hum. Neurobiol.* **4**, 219–227 (1985). <https://cseweb.ucsd.edu/classes/fa09/cse258a/papers/koch-ullman-1985.pdf>
27. Koehler, K., Guo, F., Zhang, S., Eckstein, M.P.: What do saliency models predict? *J. Vis.* **14**(3), 14–14 (2014). <https://doi.org/10.1167/14.3.14>
28. Kruthiventi, S.S.S., Ayush, K., Babu, R.V.: DeepFix: a fully convolutional neural network for predicting human eye fixations. *IEEE Trans. Image Process.* **26**(9), 4446–4456 (2017). <https://doi.org/10.1109/tip.2017.2710620>
29. Kümmerer, M.: *pysaliency*. <https://github.com/matthias-k/pysaliency>
30. Kümmerer, M., Theis, L., Bethge, M.: Deep gaze i: boosting saliency prediction with feature maps trained on ImageNet. In: *2015 International Conference on Learning Representations - Workshop Track (ICLR)* (2015), [arXiv:1411.1045](https://arxiv.org/abs/1411.1045)
31. Kümmerer, M., Wallis, T.S.A., Gatys, L.A., Bethge, M.: Understanding low- and high-level contributions to fixation prediction. In: *The IEEE International Conference on Computer Vision (ICCV)*. IEEE (2017)
32. Kümmerer, M., Wallis, T.S.A., Bethge, M.: Information-theoretic model comparison unifies saliency metrics. *Proc. Natl. Acad. Sci. USA* **112**(52), 16054–16059 (2015). <https://doi.org/10.1073/pnas.1510393112>

33. Le Meur, O., Baccino, T.: Methods for comparing scanpaths and saliency maps: strengths and weaknesses. *Behav. Res.* **45**(1), 251–266 (2012). <https://doi.org/10.3758/s13428-012-0226-9>
34. Li, Z.: A saliency map in primary visual cortex. *Trends Cogn. Sci.* **6**(1), 9–16 (2002). [https://doi.org/10.1016/s1364-6613\(00\)01817-9](https://doi.org/10.1016/s1364-6613(00)01817-9)
35. Nuthmann, A., Einhäuser, W., Schütz, I.: How well can saliency models predict fixation selection in scenes beyond central bias? a new approach to model evaluation using generalized linear mixed models. *Front. Hum. Neurosci.* **11**, 491 (2017). <https://doi.org/10.3389/fnhum.2017.00491>
36. Pan, J., et al.: SalGAN: visual saliency prediction with generative adversarial networks. [cs] (2017), [arXiv:1701.01081](https://arxiv.org/abs/1701.01081)
37. Peters, R.J., Iyer, A., Itti, L., Koch, C.: Components of bottom-up gaze allocation in natural images. *Vis. Res.* **45**(18), 2397–2416 (2005). <https://doi.org/10.1016/j.visres.2005.03.019>
38. Riche, N.: Metrics for saliency model validation. From Human Attention to Computational Attention, pp. 209–225. Springer, New York (2016). https://doi.org/10.1007/978-1-4939-3435-5_12
39. Riche, N.: Saliency model evaluation. From Human Attention to Computational Attention, pp. 245–267. Springer, New York (2016). https://doi.org/10.1007/978-1-4939-3435-5_14
40. Riche, N., Duvinage, M., Mancas, M., Gosselin, B., Dutoit, T.: Saliency and human fixations: state-of-the-art and study of comparison metrics. In: 2013 IEEE International Conference on Computer Vision. IEEE (2013). <https://doi.org/10.1109/iccv.2013.147>
41. Rothkopf, C.A., Ballard, D.H., Hayhoe, M.M.: Task and context determine where you look. *J. Vis.* **7**(14), 16 (2016). <https://doi.org/10.1167/7.14.16>
42. Schütt, H.H., Rothkegel, L.O.M., Trukenbrod, H.A., Reich, S., Wichmann, F.A., Engbert, R.: Likelihood-based parameter estimation and comparison of dynamical cognitive models. *Psychol. Rev.* **124**(4), 505–524 (2017). <https://doi.org/10.1037/rev0000068>
43. Tatler, B.W., Hayhoe, M.M., Land, M.F., Ballard, D.H.: Eye guidance in natural vision: reinterpreting salience. *J. Vis.* **11**(5), 5–5 (2011). <https://doi.org/10.1167/11.5.5>
44. Tatler, B.W.: The central fixation bias in scene viewing: selecting an optimal viewing position independently of motor biases and image feature distributions. *J. Vis.* **7**(14), 4 (2007). <https://doi.org/10.1167/7.14.4>
45. Tatler, B.W., Baddeley, R.J., Gilchrist, I.D.: Visual correlates of fixation selection: effects of scale and time. *Vis. Res.* **45**(5), 643–659 (2005). <https://doi.org/10.1016/j.visres.2004.09.017>
46. Tatler, B.W., Vincent, B.T.: Systematic tendencies in scene viewing. *J. Eye Mov. Res.* **2**(2), 1–18 (2008). http://csi.ufs.ac.za/resres/files/tatler_2008_jemr.pdf
47. Thomas, C.: OpenSalicon: an open source implementation of the salicon saliency model. CoRR abs/1606.00110 (2016), [arXiv:1606.00110](https://arxiv.org/abs/1606.00110)
48. Treisman, A.M., Gelade, G.: A feature-integration theory of attention. *Cogn. Psychol.* **12**(1), 97–136 (1980). [https://doi.org/10.1016/0010-0285\(80\)90005-5](https://doi.org/10.1016/0010-0285(80)90005-5)
49. Vig, E., Dorr, M., Cox, D.: Large-scale optimization of hierarchical features for saliency prediction in natural images. In: 2014 IEEE Conference on Computer Vision and Pattern Recognition. IEEE (2014). <https://doi.org/10.1109/cvpr.2014.358>

50. Vincent, B.T., Baddeley, R., Correani, A., Troscianko, T., Leonards, U.: Do we look at lights? using mixture modelling to distinguish between low- and high-level factors in natural image viewing. *Vis. Cogn.* **17**(6–7), 856–879 (2009). <https://doi.org/10.1080/13506280902916691>
51. Wilming, N., Betz, T., Kietzmann, T.C., König, P.: Measures and limits of models of fixation selection. *PLoS ONE* **6**(9), e24038 (2011). <https://doi.org/10.1371/journal.pone.0024038>
52. Xiao, J., Xu, P., Zhang, Y., Ehinger, K., Finkelstein, A., Kulkarni, S.: What can we learn from eye tracking data on 20,000 images? *J. Vis.* **15**(12), 790 (2015). <https://doi.org/10.1167/15.12.790>
53. Yu, F., et al.: Large-scale scene understanding challenge. <http://sun.cs.princeton.edu/2017/>
54. Yu, F., et al.: SALICON saliency prediction challenge. <http://salicon.net/challenge-2017/>
55. Zhang, J., Sclaroff, S.: Saliency detection: a Boolean map approach. In: 2013 IEEE International Conference on Computer Vision. IEEE (2013). <https://doi.org/10.1109/iccv.2013.26>
56. Zhang, L., Tong, M.H., Marks, T.K., Shan, H., Cottrell, G.W.: SUN: a Bayesian framework for saliency using natural statistics. *J. Vis.* **8**(7), 32 (2008). <https://doi.org/10.1167/8.7.32>