

# Exploiting inter-image similarity and ensemble of extreme learners for fixation prediction using deep features

Hamed R.-Tavakoli<sup>a,\*</sup>, Ali Borji<sup>b</sup>, Jorma Laaksonen<sup>a</sup>, Esa Rahtu<sup>c</sup>

<sup>a</sup>*Department of computer science, Aalto University, Espoo, Finland*

<sup>b</sup>*Center for Research in Computer Vision, University of Central Florida, USA*

<sup>c</sup>*Center for machine vision research, University of Oulu, Finland*

---

## Abstract

This paper presents a novel fixation prediction and saliency modeling framework based on inter-image similarities and ensemble of Extreme Learning Machines (ELM). The proposed framework is inspired by two observations, 1) the contextual information of a scene along with low-level visual cues modulates attention, 2) the influence of scene memorability on eye movement patterns caused by the resemblance of a scene to a former visual experience. Motivated by such observations, we develop a framework that estimates the saliency of a given image using an ensemble of extreme learners, each trained on an image similar to the input image. That is, after retrieving a set of similar images for a given image, a saliency predictor is learnt from each of the images in the retrieved image set using an ELM, resulting in an ensemble. The saliency of the given image is then measured in terms of the mean of predicted saliency value by the ensemble's members.

*Keywords:* Visual attention, saliency prediction, fixation prediction, inter-image similarity, extreme learning machines

---

## 1. Introduction

The fixation prediction, also known as saliency modeling, is associated with the estimation of a saliency map, the probability map of the locations an observer will be looking at for a long enough period of time meanwhile viewing a scene. It is part of the computational perspective of visual attention [1], the process of narrowing down the available visual information upon which to focus for enhanced processing.

---

\*Corresponding author

*Email address:* `hamed.r-tavakoli@aalto.fi` (Hamed R.-Tavakoli)

Computer vision community has been investigating the fixation prediction and saliency modeling extensively because of its wide range of applications, including, recognition [2, 3, 4, 5, 6], detection [7, 8, 9, 10, 11, 12], compression [13, 14, 15, 16], tracking [17, 18, 19, 20], segmentation [21, 22, 23], super-resolution [24], advertisement [25], perceptual designing [26], image quality assessment [27, 28], motion detection and background subtraction [29, 30, 31], scene memorability [32] and visual search [33, 34]. In many of these applications, a saliency map can facilitate the selection of a subset of regions in a scene for elaborate analysis which reduces the computation complexity and improves energy efficiency [35].

From a human centric point of view, the formation of a saliency map is not a pure bottom-up process and is influenced by several factors such as the assigned task, level of expertise of the observer, scene familiarity, and memory. It is shown that human relies on the prior knowledge about the scene and long-term memory as crucial components for construction and maintenance of scene representation [36]. In a similar vein, [37] suggests that an abstract visual representation can be retained in memory upon a short exposure to a scene and this representation influences eye movements later.

The study of the role of scene memory in guiding eye movements in a natural experience entailing prolonged immersion in three-dimensional environments [38] suggests that observers learn the location of objects over time and use a spatial-memory-guided search scheme to locate them. These findings have been the basis of research for measuring memorability of scenes from pure observer eye movements [39, 32], that is similar images have alike eye movement patterns and statistics. Inspired by the findings of [37, 36, 38] and scene memorability research, we incorporate the similarity of images as an influencing factor in fixation prediction.

Besides the fact that similar images may induce similar eye movement patterns due to memory recall, it is well agreed that the interaction of low-level visual cues (e.g., edges, color, etc.) affect saliency formation [40] and contextual information of a scene can modulate the saliency map [41, 42]. Imagine that you are watching two pairs of images, a pair of street scene and a pair of nature beach images, meanwhile having your eye movements recorded. It is not surprising to find similar salient regions for the images of alike scenes because similar low-level cues and contextual data are mostly present in each pair. Figure 1 depicts examples of such a scenario. In the case of street scene, the observers tend to converge to the traffic signs, while they tend to spot low-level structural information in beach images. This further motivates us to exploit learning saliency from inter-image similarities.

This paper presents a novel fixation prediction algorithm based on inter-image similarities and an ensemble of saliency learners using features from deep convolutional neural networks. To meet this end, we first investigate the benefits from inter-image similarities for fixation prediction. Then, we introduce 1) an image similarity metric using *gist* descriptor [41] and *classemes* [43], 2) a fixation prediction algorithm, using an ensemble of extreme learning machines, where for a given image, each member of the ensemble is trained with an image similar

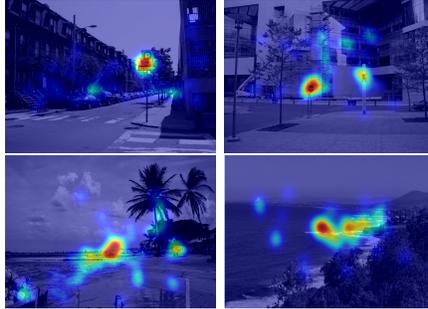


Figure 1: Images with similar contextual information with observers’ fixation density maps overlaid. Top: pair of street images, down: pair of natural beach images.

to the input image. We report the performance of the proposed framework on MIT saliency benchmarks [44], both MIT300 and CAT2000 databases<sup>1</sup>, along with evaluations on databases with publicly available ground-truth.

In the rest of this paper, we briefly review the related work. Afterwards, using a toy problem, we demonstrate the benefit from inter-image similarity. In section 4, we explain the proposed model. We then continue with the experiments to assess the performance of the model. The paper ends with discussion and conclusion remarks.

## 2. Related work

The field of computer vision is replete with a numerous variety of saliency models. A widely recognized group of models apply the feature integration theory [40] and consider a center-surround interaction of features [45, 2, 46, 47, 48, 49, 50, 51, 52, 53, 54]. There are models which consider the information theoretic foundations [55, 56, 57, 58, 59, 60], frequency domain aspect [61, 62, 16, 63, 64, 65, 66, 67, 68], diffusion and random walk techniques [69, 70, 71], and etc. Investigating the extent of saliency modeling approaches is beyond the scope of this article and readers are advised to consult relevant surveys [72, 73]. We, however, briefly review some of the most relevant techniques.

Learning-based techniques are a large group of methods which are establishing a relation between a feature space and human fixations. For example, [74] uses a nonlinear transformation to associate image patches with human eye movement statistics. In [75], a linear SVM classifier is used to establish a relation between three channels of low- (intensity, color, etc), mid- (horizon line)

---

<sup>1</sup>These databases have their ground-truth unavailable to public in order to provide a fair model evaluation. Thus, the scores are computed by MIT saliency research team using our submitted maps.

and high-level (faces and people) features and human eye movements in order to produce a saliency map. In a similar vein, [76] employs multiple-instance learning. By learning a classifier, [77, 78] estimate the optimal weights for fusing several conspicuity maps from observers' eye movement data. These approaches often learn a probabilistic classifier to determine the probability of a feature being salient. Then, they employ the estimated saliency probability in order to build a saliency map.

The recent saliency modeling methods, akin to other computer vision techniques, are revolutionized and advanced significantly by applying deep Convolutional Neural Networks (CNN). There exists significant number of models that employ CNNs, of which many are relevant to the proposed model.

Ensembles of Deep Networks (eDN) [79] adopts the neural filters learned during image classification task by deep neural networks and learns a classifier to perform fixation prediction. eDN can be considered an extension to [75] in which the features are obtained from layers of a deep neural network. For each layer of the deep neural network, eDN first learns the optimal blend of the neural responses of all the previous layers and the current layer by a guided hyperparameter search. Then, it concatenates the optimal blend of all the layers to form a feature vector for learning a linear SVM classifier.

Deep Gaze I [80] utilizes CNNs for the fixation prediction task by treating saliency prediction as point processing. Despite this model is justified differently than [79] and [75], in practice, it boils down to the same framework. Nonetheless, the objective function to be minimized is slightly different due to the explicit incorporation of the center-bias factor and the imposed sparsity constraint in the framework. SalNet [81] is another technique that employs a CNN-based architecture, where the last layer is a deconvolution. The first convolution layers are initialized by the VGG16 [82] and the deconvolution is learnt by fine-tuning the architecture for fixation prediction.

Multiresolution CNN (Mr-CNN) [83] designs a deep CNN-based technique to discriminate image patches centered on fixations from non-fixated image patches at multiple resolutions. It hence trains a convolutional neural network at each scale, which results in three parallel networks. The outputs of these networks are connected together through a common classification layer in order to learn the best resolution combination.

SALICON [84] develops a model by fine-tuning the convolutional neural network, trained on ImageNet, using saliency evaluation metrics as objective functions. It feeds an image into a CNN architecture at two resolutions, coarse and fine. Then, the response of the last convolution layer is obtained for each scale. These responses are then concatenated together and are fed into a linear integration scheme, optimizing the Kullback-Leibler divergence between the network output and the ground-truth fixation maps in a regression setup. The error is back-propagated to the convolution layers for fine-tuning the network.

The proposed method can be considered a learning-based approach. While many of the learning-based techniques are essentially solving a classification problem, the proposed model has a regression ideology in mind. It is thus closer to the recent deep learning approaches that treat the problem as estimation of

a probability map in terms of a regression problem [81, 84, 85]. Nonetheless, it exploits an ensemble of extreme learning machines.

### 3. Saliency benefits from inter-image similarity

The main motivation behind the proposed model is that people may have similar fixation patterns in exposure to alike images. In other words, inter-image saliency benefits saliency prediction. In order to investigate such an assertion, we build a toy problem to tell *how well the saliency map of an image predicts saliency in a similar image*.

We choose a common saliency database [75] and computed the gist [41] of the scene for each image. Afterwards, the most similar image pairs and the most dissimilar pairs were identified. For each image pair, we use the fixation density map of one as the predicted saliency map of the other. The assessment reveals that such a fixation prediction scheme produces significantly different ( $p \leq 0.05$ ) shuffled AUC scores [86] where the score of prediction using similar pairs is 0.54 and the score of prediction by dissimilar image pairs is 0.5. The results indicate that while there is a degree of prediction for similar pairs, the dissimilar pairs are not doing better than chance. We observe the same performance difference for other metrics such as correlation score (0.175 vs. 0.115) and normalized scanpath score (0.86 vs. 0.59). Given the above observation, we lay the foundation of our saliency model for fixation prediction.

### 4. Saliency Model

A high-level conceptual schematic of our proposed model is depicted in Figure 2. The framework components include: 1) an image feature transform, 2) a similar image retrieval engine and a scene repository bank, and 3) an ensemble of neural saliency (fixation) predictors. The image feature transform performs the feature extraction and produces a pool of features used by the other units in the system. The similar image retrieval finds the top most similar images, stored in the scene bank, corresponding to a given image. It then retrieves the predictors trained using those images in order to facilitate the formation of the ensemble of saliency predictors. In the rest of this section, we explained the details of the mentioned components.

#### 4.1. Image feature transform

The image feature transform unit extracts several features from the image and feeds them forward to the other units. There has been a recent surge in the application of features learnt from image statistics and deep convolutional neural networks (CNNs) in a wide range of computer vision related applications. In this work, we adopt a filter-bank approach to the use of CNNs [87] for saliency prediction. We, thus, build an image pyramid and compute the CNNs' responses over each scale using the architecture of VGG16 [82]. To combine the convolution responses of each scale, we employ an upsampling procedure

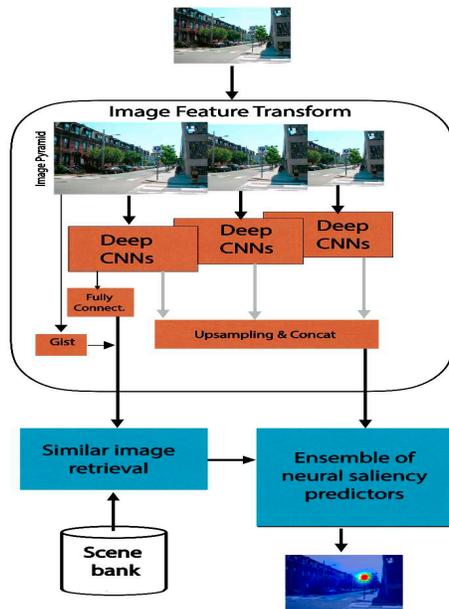


Figure 2: General architecture of the model.

and concat the features from the last convolution layer of each scale in order to build a feature map.

Furthermore, we compute the classemes [43] from deep pipeline, that is, the probability of each of the one thousand classes of ImageNet [88] is computed using the fully-connected layers of the VGG16. The classemes are complemented by the low-level scene representation to make the gist of the scene [9]. The classemes and low-level scene features of [41] build a spatial representation of the outside world that is rich enough to convey the meaning of a scene as envisioned in [89]. The feature vector obtained by concatenating classemes and gist features is used for the recognition and retrieval of similar images.

#### 4.2. Similar image retrieval & scene bank

The similar image retrieval unit fetches the information required for building an ensemble of neural predictors from the scene bank. The scene bank holds a set of images in terms of scene representation feature vector, consisting of classemes feature and the gist descriptor, and a neural fixation predictor unit for each image.

Given the scene representation vector of an input image, denoted as  $v^q$ , the retrieval method fetches the most  $n$  similar images from the set of scene vectors,  $\mathbf{V} = \{v_1, \dots, v_{n'}\}$ , using the Euclidean distance, that is,  $dist_i = \|v^q - v_i\|$ . It then fetches the neural fixation predictor units corresponding to the  $n$



Figure 3: Image retrieval examples. The input (query) image is on the left and its closest match is on the right. The query images are from [75] and the closest match is from [90]. The observers’ fixation density map is overlaid.

images with the smallest  $dist_i$  in order to form the ensemble of neural fixation predictors, to be discussed in Section 4.3.

Figure 3 demonstrates the results of retrieval system. It visualizes a query image and its corresponding most similar retrieved image between two different databases with the observer gaze information overlaid. Interestingly, the retrieved images not only share similar objects and bottom-up structures, but can also have similar attention grabbing regions. It is worth noting that the closest scene is not necessarily of the same scene category, however, it often contains similar low-level and/or high-level perceptual elements.

#### 4.3. Saliency prediction

We define the saliency of an image in terms of features and locations, that is,  $\mathbf{Sal} = p(\mathbf{y}|\mathbf{x}, \mathbf{m})$ , where  $\mathbf{y}$  corresponds to pixel level saliency,  $\mathbf{x}$  represents image features and  $\mathbf{m}$  is the location. Under the independence assumption, the saliency formulation boils down to the following:

$$\mathbf{Sal} = p(\mathbf{y}|\mathbf{x})p(\mathbf{y}|\mathbf{m}). \quad (1)$$

The  $p(\mathbf{y}|\mathbf{x})$  corresponds to saliency prediction from image features and  $p(\mathbf{y}|\mathbf{m})$  represents a spatial prior. We estimate  $p(\mathbf{y}|\mathbf{x})$  using an ensemble of neural predictors and  $p(\mathbf{y}|\mathbf{m})$  is learnt from human gaze information.

Figure 4 depicts the ensemble of neural saliency predictors. The ensemble of neural predictors consists of several neural units with equal contributions. In training phase, we train one neural unit for each image in the training set and store them in the scene bank. In the test phase, the retrieval unit fetches several neural units, corresponding to the  $n$  images most similar to the input image. The ensemble, then, computes the responses of each of the units and aggregates them in order to produce an estimate of  $p(\mathbf{y}|\mathbf{x})$ , as follows:

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{Q} \left( \left( \sum_j \mathcal{C}(\tanh(\mathbf{y}_j)) \right)^\alpha \right), \quad (2)$$

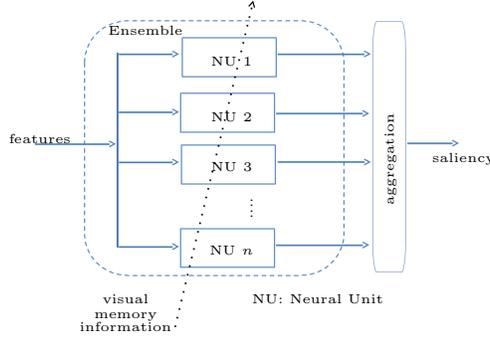


Figure 4: Ensemble of neural saliency predictors.

$$\mathcal{C}(x) = \begin{cases} x & x > 0 \\ 0 & x \leq 0 \end{cases}, \quad (3)$$

where  $\mathcal{Q}(\cdot)$  resizes image or saliency data to the size of preference (the size of input image),  $\alpha$  is an attenuation factor to emphasize more salient areas, and  $\mathbf{y}_j$  is the output of the  $j$ th unit of the ensemble.

#### 4.3.1. Neural units

The neural saliency predictor utilizes randomly-weighted single-layer feed-forward networks in order to establish a mapping from the feature space to the saliency space. The idea of randomly-weighted single-hidden-layer feedforward networks (SLFNs) can be traced back to the Gamba perceptron [91] followed by others like [92, 93]. In the neural saliency predictor, we adopt the recent implementation of Extreme Learning Machines (ELM) [94]. The theory of ELM facilitates the implementation of a neural network architecture such that the hidden layer weights can be chosen randomly meanwhile the output layer weights are determined analytically [95]. Motivated by better function approximation properties of ELMs [96, 97], we employ them as the primary entity of the neural saliency prediction.

Having a set of training samples  $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N \subset \mathbb{R}^k \times \mathbb{R}^m$ , the image features  $\mathbf{x}_i$  and the corresponding fixation density value  $\mathbf{y}_i$  are associated using a SLFNs with  $L$  hidden nodes defined as

$$\mathbf{y}_i = \sum_{j=1}^L \boldsymbol{\gamma}_j f(\boldsymbol{\omega}_j \cdot \mathbf{x}_i + b_j), \quad (4)$$

where  $f(\cdot)$  is a nonlinear activation function,  $\boldsymbol{\gamma}_j \in \mathbb{R}^m$  is the output weight vector,  $\boldsymbol{\omega}_j \in \mathbb{R}^k$  is the input weight vector, and  $b_j$  is the bias of the  $j$ th hidden node. The conventional solution to (4) is gradient-based, which is a slow iterative process that requires to tune all the parameters like  $\boldsymbol{\gamma}_j$ ,  $\boldsymbol{\omega}_j$  and  $b_j$ . The iterative scheme is prone to divergence, local minima, and overfitting. The

ELM tries to soften such problems and avoid them by random selection of the hidden layer parameters ( $\omega_j$  and  $b_j$ ) and the estimation of output weights. To this end, (4) can be rewritten as

$$\mathbf{Y} = \mathbf{H}\mathbf{\Gamma}, \quad (5)$$

where  $\mathbf{Y} = [\mathbf{y}_1 \ \mathbf{y}_2 \ \dots \ \mathbf{y}_N]^T \in \mathbb{R}^{N \times m}$ ,  $\mathbf{\Gamma} = [\gamma_1 \ \gamma_2 \ \dots \ \gamma_L]^T \in \mathbb{R}^{L \times m}$ , and

$$\mathbf{H} = \begin{bmatrix} f(\omega_1 \cdot \mathbf{x}_1 + b_1) & \cdots & f(\omega_L \cdot \mathbf{x}_1 + b_L) \\ \vdots & \ddots & \vdots \\ f(\omega_1 \cdot \mathbf{x}_N + b_1) & \cdots & f(\omega_L \cdot \mathbf{x}_N + b_L) \end{bmatrix}_{N \times L}, \quad (6)$$

which is the hidden layer matrix of the neural network. Once the matrix  $\mathbf{H}$  is decided by random selection of input weights and biases, the solution of (5) can be approximated as  $\mathbf{\Gamma} = \mathbf{H}^\dagger \mathbf{Y}$ , where  $\mathbf{H}^\dagger$  is the *Moore-Penrose pseudoinverse* of matrix  $\mathbf{H}$ .

#### 4.3.2. Learning spatial prior

In order to learn the spatial prior,  $p(\mathbf{y}|\mathbf{m})$ , we fit a mixture of Gaussian over the eye fixation data. We learn the spatial prior using the gaze data of [90], where the number of kernels corresponds to the number of fixation points. The spatial prior puts more weight on the regions that are more agreed by observers. As demonstrated in many saliency research papers, the spatial prior introduces a center-bias effect [98]. The same phenomenon is observed in Figure 5, depicting the spatial prior. While there exist arguments on getting advantage of location priors, we address the issue by selecting proper evaluation metrics and benchmarks. It is also worth noting that we are not using summation prior integration, which generally boosts all the regions in the center of the image equally.

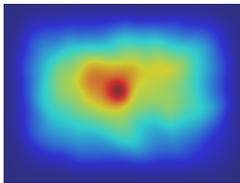


Figure 5: Spatial prior learnt from [90].

## 5. Experiments

We conduct several experiments in order to evaluate the model. The test databases include MIT [75], MIT300 [99], and CAT2000 [100]. The MIT database consists of 1003 images of indoor and outdoor scenes with eye movements of 15 observers. MIT300 consists of 300 natural indoor and outdoor scenes and CAT2000 consists of 4000 images divided into two sets of train and test, with

2000 images in each set. CAT2000 includes 20 categories of images, including, action, affective, art, black & white, cartoon, fractal, indoor, outdoor, inverted, jumbled, line drawings, low resolution, noisy, object, outdoor man made, outdoor natural, pattern, random, satellite, sketch, and social.

MIT300 and CAT2000 (test set) do not allow the ground-truth access in order to provide a fair comparison. At the moment, they are the widely accepted benchmarks and the results presented are provided by the MIT saliency benchmark team using our submitted maps. The results of the proposed model are also accessible on the benchmark website<sup>2</sup> under the acronym “iSEEL”.

We learn two ensembles,  $ensemble_{OSIE}$  and  $ensemble_{CAT2k}$ . The first is trained on the OSIE database [90] and the latter is trained using the training set of CAT2000. We employ  $ensemble_{CAT2k}$  in predicting the CAT2000 test images. The system parameters are optimized for each ensemble.

In this section, we first explain the system parameters. We then evaluate the performance generalization of the proposed model in comparison with a baseline model using the MIT database. We continue with the Benchmark results on the MIT300 and the CAT2000 databases.

### 5.1. System parameters

The system parameters are the number of neural units in each ensemble, denoted  $n$ , the number of hidden layers in each unit,  $L$ , and the attenuation factor,  $\alpha$ . We furthermore learn a post processing smoothing Gaussian kernel, denoted as  $\sigma$ , which is used to smooth the model’s maps. All the parameters, except the number of hidden nodes are learnt. For each of the ensembles, the number of hidden nodes of each neural unit is fixed and equal to 20. The rest of the parameters of the system are optimized on Toronto database [56]. The tuning cost function minimizes the KL-divergence between the maps of the model and the ground-truth fixation density maps.

Figure 6 depicts the effect of the number of neural units in conjunction with the value of the attenuation factor  $\alpha$  on the ensemble performance. Based on our observations, an ensemble of size 10 is required to obtain an acceptable result. The optimization of parameters, however, recommend the following parameters for each ensemble,  $ensemble_{OSIE}$  :  $[n = 697, \alpha = 6, \sigma = 13]$  and  $ensemble_{CAT2k}$ :  $[n = 1710, \alpha = 9, \sigma = 13]$ , where  $L = 20$  has been fixed.

### 5.2. Performance generalization

To test the generalization of the model, we evaluate its performance using the MIT database [75]. We choose the ensemble of deep neural networks (eDN) [79] as a baseline model because of the use of deep features and SVM classifiers. The proposed model, however, utilizes an ensemble of ELM regression units. We also evaluate several models including, AIM [55], GBVS [69], AWS [101], Judd [75], and FES [51] for the sake of comparison with traditional models.

---

<sup>2</sup>MIT saliency benchmark website:[http://saliency.mit.edu/results\\_mit300.html](http://saliency.mit.edu/results_mit300.html)

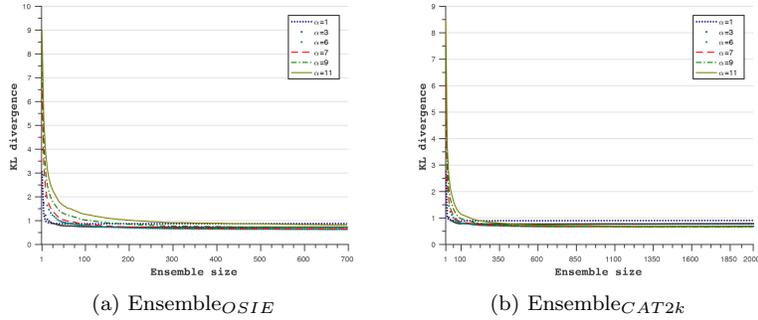


Figure 6: System parameters: the effect of ensemble size and  $\alpha$  on the saliency model.

In order to ease the interpretation of evaluation, we choose a subset of scores that complement each other. We employ shuffled AUC (sAUC, an AUC metric that is robust towards center bias), similarity metric (SIM, a metric indicating how two distributions resemble each other [44]), and normalized scanpath saliency (NSS, a metric to measure consistency with human fixation locations). NSS and sAUC scores are utilized in [86], which we borrow part of the scores from, and complement them with the SIM score.

Figure 7 reports the results. As depicted, the proposed model outperforms all other models on two metrics and outperforms the eDN on all the three metrics. The highest gain compared to the eDN is on the NSS score, indicating a high consistency with human fixation locations which explains the high SIM score as well. To summarize, the proposed model generalizes well and has the edge over traditional models. We later compare the proposed model with the recent state-of-the-art models on well-established benchmarks.

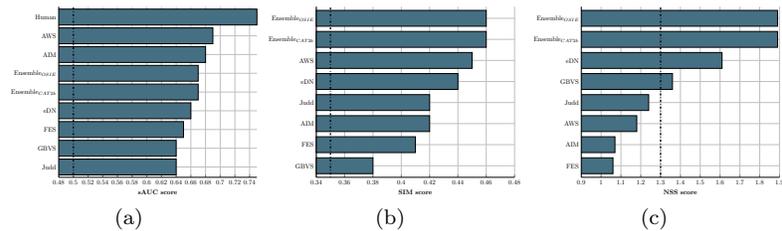


Figure 7: Performance generalization: the performance of the proposed model compared to traditional models and eDN [79] as a baseline model. The dashed vertical line indicates the performance of a Gaussian dummy model. The human score for SIM and NSS are 1 and 3.1, respectively.

Table 1: MIT300 Benchmark results, sorted using NSS.

Model	AUC-based metrics			Similarity-based metrics				NSS
	Judd	Borji	shuffled	SIM	EMD	CC	KL	
<b>Infinite Human</b>	0.92	0.88	0.81	1.00	0	1	0	3.29
<b>SALICON [84]</b>	0.87	0.85	0.74	0.60	2.62	0.74	0.54	2.12
<b>PDP [85]</b>	0.85	0.80	0.73	0.60	2.58	0.70	0.92	2.05
<b>ML-Net [102]</b>	0.85	0.75	0.70	0.59	2.63	0.67	1.10	2.05
<b>ensemble<sub>OSIE</sub> (iSEEL)</b>	0.84	0.81	0.68	0.57	2.72	0.65	0.65	1.78
<b>Mean One Human</b>	0.80	0.66	0.63	0.38	3.48	0.52	6.19	1.65
<b>SalNet [81]</b>	0.83	0.82	0.69	0.52	3.31	0.58	0.81	1.51
<b>BMS [103]</b>	0.83	0.82	0.65	0.51	3.35	0.55	0.81	1.41
<b>Mr-CNN [83]</b>	0.79	0.75	0.69	0.48	3.71	0.48	1.08	1.37
<b>eDN [79]</b>	0.82	0.81	0.62	0.41	4.56	0.45	1.14	1.14

### 5.3. Benchmark

Many of the recent deep saliency models have their codes and maps unavailable to public, making comparisons difficult. We, hence, rely on available benchmarks. We report the performance using all the metrics and published works, reported on the MIT benchmark. For brevity, the focus will be on recent top-performing models. The results also include the performance of “Infinite Human” and “Mean One Human” to indicate how well a model performs in comparison with mean eye position of several human (upper-bound performance) and the on average performance of one human, respectively.

*Results on MIT300.* Table 1 summarizes the performance comparison, where the proposed model is 4th among published works on this benchmark on the basis of NSS. MIT300 is the largest benchmark with over 60 models at the time of this writing. We, however, report the best performing models and the most recent state-of-the-art ones. The comparison indicates that the models are becoming powerful enough to capture fixation location. It is, hence, difficult to distinguish them from each other on many metrics. NSS, however, seems to be the most informative metric that determines the models’ performance well, particularly for top-performing models that judging AUC-based metrics and Similarity-based metrics are difficult.

*Results on CAT2000.* Table 2 contains the performance comparison on the CAT2000 database. 19 models, which are mostly traditional ones, are evaluated on this database. The proposed model, ensemble<sub>CAT2k</sub>, ranks similarly with BMS [103] at the top of the ranking. Both models produce the highest NSS score among models and on average have indistinguishable values for the AUC-based and the Similarity-based metrics.

We also evaluate ensemble<sub>OSIE</sub> along with ensemble<sub>CAT2k</sub> in order to further investigate the improvements caused by incorporating similar images in the training phase. Backing the hypothesis, the ensemble trained on CAT2000 outperforms the ensemble that is learnt from only indoor and outdoor images of OSIE in terms of the overall scores.

We look into the performance of the models in each of the twenty class categories of CAT2000 database. To be concise, we investigate ensemble<sub>CAT2k</sub>,

Table 2: CAT2000 Benchmark results, sorted using NSS.

Model	AUC-based metrics			Similarity-based metrics				NSS
	Judd	Borji	shuffled	SIM	EMD	CC	KL	
<b>Infinite Human</b>	0.90	0.84	0.62	1.00	0	1	0	2.85
<b>ensemble<sub>CAT2k</sub> (iSEEL)</b>	0.84	0.81	0.59	0.62	1.78	0.66	0.92	1.67
<b>BMS [103]</b>	0.85	0.84	0.59	0.61	1.95	0.67	0.83	1.67
<b>ensemble<sub>OSIE</sub></b>	0.83	0.81	0.59	0.59	2.24	0.64	0.67	1.62
<b>FES [103]</b>	0.82	0.76	0.54	0.57	2.24	0.64	2.10	1.61
<b>Mean One Human</b>	0.76	0.67	0.56	0.43	2.51	0.56	7.77	1.54
<b>Judd [75]</b>	0.84	0.84	0.56	0.46	3.60	0.54	0.94	1.30
<b>eDN [79]</b>	0.85	0.84	0.55	0.52	2.64	0.54	0.97	1.30

ensemble<sub>OSIE</sub>, and BMS, which are the top three best performing models, using the three metrics of shuffled AUC (sAUC), SIM, and NSS. The results are summarized in Figure 8. The proposed model, both ensemble<sub>CAT2k</sub> and ensemble<sub>OSIE</sub>, are outperforming the BMS on low resolution, noisy, outdoor, black & white, action, affective and social categories. The BMS seems performing better when there is no particular contextual information and more low-level feature interactions matter, e.g., fractal category, and pattern. The other categories are, however, more difficult to judge. Overall, it seems the three models can complement each other in the areas where one falls behind the others.

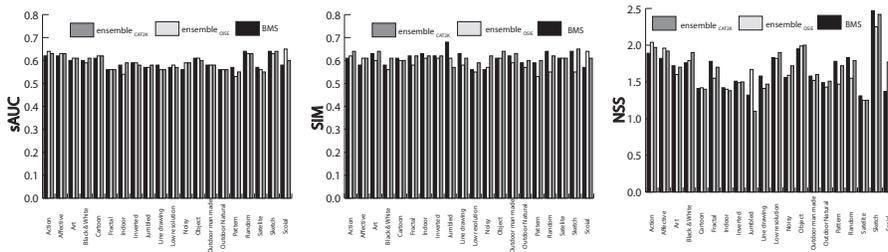


Figure 8: Performance on categories of CAT2000: the performance of the proposed model compared to BMS.

## 6. Discussion & conclusion

We demonstrated the usefulness of scene similarity in predicting the saliency motivated by the effect of the familiarity of a scene on the observer’s eye movements. The idea can, however, be easily extended to the utilization of observers’ eye movements in task-specific models, where a model is trained for a specific task and experts’ eye movements are incorporated. An expert approach for solving a specific task is different from that of a naive observer. Thus, we can consider the encoding of expert observers’ eye movements as an implicit expert knowledge utilization, which can be handy in scenarios of scene analysis such as spotting object-specific anomalies from saliency maps in order to reduce the search time.

We introduced a saliency model with the motive of exploiting the effect of immediate scene recall on the human perception. The proposed model uses randomly-weighted neural networks as an ensemble architecture. It establishes a mapping from a feature space, consisting of deep features, to the saliency space. The saliency prediction relies only on the neural units corresponding to the images that are similar to the input image. The neural units are pre-trained and stored in a scene bank from a handful of images. For each neural unit, the scene bank also stores a scene descriptor, consisting of classemes and gist descriptor. To find the similar images from scene bank, the proposed model employs the distance between the scene descriptor of the input image and neural units.

The proposed model was evaluated on several databases. The results were reported on two well-established benchmark databases by the MIT benchmark team, namely MIT300 and CAT2000. Among the published methods and on the basis of NSS, consistency with the locations of human fixation, the proposed method was ranked 4th and 1st (in conjunction with BMS) on MIT300 and CAT2000, respectively. The results indicate benefit from learning saliency from images similar to the input image. The code for the proposed model is available at: <http://github.com/hrtavakoli/iseel>.

## Acknowledgement

Hamed R.-Tavakoli and Jorma Laaksonen were supported by the Finnish Center of Excellence in Computational Inference Research (COIN). The authors would like to thank the MIT saliency benchmark team, particularly Zoya Bylinskii, for their quick response on benchmark request.

## References

- [1] J. K. Tsotsos, *A Computational Perspective on Visual Attention*, The MIT Press, 2011.
- [2] S. Frintrop, *Vocus: A visual attention system for object detection and goal-directed search*, Ph.D. thesis (2006).
- [3] A. Salah, E. Alpaydin, L. Akarun, A selective attention-based method for visual pattern recognition with application to handwritten digit recognition and face recognition, *IEEE Trans. Pattern Anal. Mach. Intell.* 24 (3) (2002) 420–425.
- [4] C. Siagian, L. Itti, Rapid biologically-inspired scene classification using features shared with visual attention, *IEEE Trans. Pattern Anal. Mach. Intell.* 29 (2) (2007) 300–312.
- [5] D. Gao, S. Han, N. Vasconcelos, Discriminant saliency, the detection of suspicious coincidences, and applications to visual recognition, *IEEE Trans. Pattern Anal. Mach. Intell.* 31 (6) (2009) 989 – 1005.

- [6] C. Kanan, G. Cottrell, Robust classification of objects, faces, and flowers using natural image statistics, in: CVPR, 2010.
- [7] C. Papageorgiou, T. Poggio, A trainable system for object detection, *Int. J. Comput. Vision* 38 (1) (2000) 15–33.
- [8] G. Bouchard, B. Triggs, Hierarchical part-based visual object categorization, in: CVPR, 2005.
- [9] A. Oliva, A. Torralba, Modeling the shape of the scene: A holistic representation of the spatial envelope, *Int. J. Comput. Vision* 42 (2001) 145–175.
- [10] A. Torralba, Contextual priming for object detection, *Int. J. Comput. Vision* 53 (2) (2003) 169–191.
- [11] K. Ehinger, B. Hidalgo-Sotelo, A. Torralba, A. Oliva, Modeling search for people in 900 scenes, *Vis. Cogn.* 17 (2009) 945–978.
- [12] G. Fritz, C. Seifert, L. Paletta, H. Bischof, Attentive object detection using an information theoretic saliency measure, in: WAPCV, 2005.
- [13] M. Kunt, A. Ikonomopoulos, M. Kocher, Second-generation image-coding techniques, *Proc. IEEE* 73 (4) (1985) 549–574.
- [14] N. Ouerhani, J. Bracamonte, H. Hugli, M. Ansorge, F. Pellandini, Adaptive color image compression based on visual attention, in: ICIP, 2001.
- [15] N. Dhavale, L. Itti, Saliency-based multifoventated mpeg compression, in: ISSPA, 2003.
- [16] C. Guo, L. Zhang, A novel multiresolution spatiotemporal saliency detection model and its applications in image and video compression, *IEEE Trans. Img. Proc.* 19 (1) (2010) 185–198.
- [17] V. Mahadevan, N. Vasconcelos, Saliency-based discriminant tracking, in: CVPR, 2009.
- [18] S. Frintrop, General object tracking with a component-based target descriptor, in: ICRA, 2010.
- [19] A. Borji, S. Frintrop, D. Sihite, L. Itti, Adaptive object tracking by learning background context, in: CVPRW, 2012.
- [20] H. R.-Tavakoli, M. Shahram Moin, J. Heikkilä, Local similarity number and its application to object tracking, *Int. J. Adv. Robot. Syst.* 10 (184). doi:10.5772/55337.
- [21] A. Mishra, Y. Aloimonos, C. L. Fah, Active segmentation with fixation, in: CVPR, 2009.

- [22] Y. Fu, J. Cheng, Z. Li, H. Lu, Saliency cuts: An automatic approach to object segmentation, in: ICPR, 2008.
- [23] V. Yanulevskaya, J. Uijlings, J.-M. Geusebroek, Salient object detection: From pixels to segments, *Image and Vision Comput.* 31 (1) (2013) 31 – 42.
- [24] N. G. Sadaka, L. J. Karam, Efficient perceptual attentive super-resolution, in: ICIP, 2009.
- [25] H. Liu, S. Jiang, Q. Huang, C. Xu, A generic virtual content insertion system based on visual attention analysis, in: ACM MM, 2008.
- [26] R. Rosenholtz, A. Dorai, R. Freeman, Do predictions of visual perception aid design?, *ACM Trans. Appl. Percept.* 8 (2) (2011) 12:1–12:20.
- [27] A. Ninassi, O. L. Meur, P. L. Callet, D. Barba, Does where you gaze on an image affect your perception of quality? applying visual attention to image quality metric, in: ICIP, 2007.
- [28] Q. Ma, L. Zhang, Saliency-based image quality assessment criterion, in: ICIC, 2008.
- [29] V. Mahadevan, N. Vasconcelos, Spatiotemporal saliency in dynamic scenes, *IEEE Trans. Pattern Anal. Mach. Intell.* 32 (1) (2010) 171 –177.
- [30] H. R.-Tavakoli, E. Rahtu, J. Heikkilä, Saliency detection using joint temporal and spatial decorrelation, in: SCIA, 2013.
- [31] H. R.-Tavakoli, E. Rahtu, J. Heikkil, Temporal saliency for fast motion detection, in: ACCV workshops, 2013.
- [32] M. Mancas, O. L. Meur, Memorability of natural scenes: The role of attention, in: ICIP, 2013.
- [33] N. Butko, J. R. Movellan, Optimal scanning for faster object detection, in: CVPR, 2009.
- [34] M.-M. Cheng, Z. Zhang, W.-Y. Lin, P. H. S. Torr, BING: Binarized normed gradients for objectness estimation at 300fps, in: CVPR, 2014.
- [35] R. Kasturi, D. Goldgof, R. Ekambaram, R. Sharma, G. Pratt, M. Anderson, M. Peot, M. Aguilar, E. Krotkov, D. Hackett, D. Khosla, Y. Chen, K. Kim, Y. Ran, Q. Zheng, L. Elazary, R. Voorhies, D. Parks, L. Itti, Performance evaluation of neuromorphic-vision object recognition algorithms, in: ICPR, 2014.
- [36] A. Hollingworth, J. M. Henderson, Accurate visual memory for previously attended objects in natural scenes, *J. Exp. Psychol. Hum. Percept. Perform.* 28 (2002) 113–136.

- [37] M. Castelhana, J. Henderson, Initial scene representations facilitate eye movement guidance in visual search., *J. Exp. Psychol. Hum. Percept. Perform.* 33 (4) (2007) 753–63.
- [38] D. Kit, L. Katz, B. Sullivan, K. Snyder, D. Ballard, M. Hayhoe, Eye movements, visual search and scene memory, in an immersive virtual environment, *PLoS ONE* 9 (4).
- [39] A. Bulling, D. Roggen, Recognition of visual memory recall processes using eye movement analysis, in: *UbiComp*, 2011.
- [40] A. M. Treisman, G. Gelade, A feature-integration theory of attention, *Cognitive Psychol* 12 (1) (1980) 97 – 136.
- [41] A. Torralba, A. Oliva, M. Castelhana, J. Henderson, Contextual guidance of eye movements and attention in real-world scenes: the role of global features in object search, *Psychol. Rev.* 113 (4) (2006) 766–86.
- [42] A. Oliva, A. Torralba, The role of context in object recognition., *Trends Cogn Sci.* 11 (12).
- [43] L. Torresani, M. Szummer, A. Fitzgibbon, Efficient object category recognition using classemes, in: *ECCV*, 2010.
- [44] T. Judd, F. Durand, A. Torralba, A benchmark of computational models of saliency to predict human fixations, *Tech. Rep. MIT-CSAIL-TR-2012-001*, Massachusetts institute of technology (2012).
- [45] L. Itti, C. Koch, E. Niebur, A model of saliency-based visual attention for rapid scene analysis, *IEEE Trans. Pattern Anal. Mach. Intell.* 20 (11) (1998) 1254 –1259.
- [46] S.-B. Choi, B.-S. Jung, S.-W. Ban, H. Niitsuma, M. Lee, Biologically motivated vergence control system using human-like selective attention model, *Neurocomputing* 69 (2006) 537 – 558.
- [47] O. Le Meur, P. Le Callet, D. Barba, D. Thoreau, A coherent computational approach to model bottom-up visual attention, *IEEE Trans. Pattern Anal. Mach. Intell.* 28 (5) (2006) 802 –817.
- [48] N. Murray, M. Vanrell, X. Otazu, C. Parraga, Saliency estimation using a non-parametric low-level vision model, in: *CVPR*, 2011.
- [49] D. Gao, V. Mahadevan, N. Vasconcelos, The discriminant center-surround hypothesis for bottom-up saliency, in: *NIPS*, 2007.
- [50] H. J. Seo, P. Milanfar, Nonparametric bottom-up saliency detection by self-resemblance, in: *CVPR*, 2009.
- [51] H. R.-Tavakoli, E. Rahtu, J. Heikkil, Fast and efficient saliency detection using sparse sampling and kernel density estimation, in: *SCIA*, 2011.

- [52] E. Erdem, A. Erdem, Visual saliency estimation by nonlinearly integrating features using region covariances, *J. Vis.* 13 (4).
- [53] Q. Wang, Y. Yuan, P. Yan, Visual saliency by selective contrast, *IEEE Transactions on Circuits and Systems for Video Technology* 23 (7) (2013) 1150–1155.
- [54] H. R. Tavakoli, E. Rahtu, J. Heikkilä, Stochastic bottomup fixation prediction and saccade generation, *Image and Vision Computing* 31 (9) (2013) 686 – 693.
- [55] N. D. B. Bruce, J. K. Tsotsos, Saliency based on information maximization, in: *NIPS*, 2006.
- [56] N. D. B. Bruce, J. K. Tsotsos, Saliency, attention, and visual search: An information theoretic approach, *J. Vis.* 9 (3).
- [57] M. Mancas, Computational attention: Towards attentive computers, Ph.D. thesis, CIACO University (2007).
- [58] X. Hou, L. Zhang, Dynamic visual attention: searching for coding length increments, in: *NIPS*, 2008.
- [59] Y. Li, Y. Zhou, J. Yan, Z. Niu, J. Yang, Visual saliency based on conditional entropy, in: *ACCV*, 2010.
- [60] Y. Li, Y. Zhou, L. Xu, X. Yang, J. Yang, Incremental sparse saliency detection, in: *ICIP*, 2009.
- [61] X. Hou, L. Zhang, Saliency detection: A spectral residual approach, in: *CVPR*, 2007.
- [62] C. Guo, Q. Ma, L. Zhang, Spatio-temporal saliency detection using phase spectrum of quaternion fourier transform, in: *CVPR*, 2008.
- [63] P. Bian, L. Zhang, Biological plausibility of spectral domain approach for spatiotemporal visual saliency, in: *ICONIP*, 2008.
- [64] P. Bian, L. Zhang, Visual saliency: a biologically plausible contourlet-like frequency domain approach, *Cogn. Neurodyn.* 4 (3) (2010) 189–198.
- [65] J. Li, M. Levine, X. An, H. He, Saliency detection based on frequency and spatial domain analyses, in: *BMVC*, 2011.
- [66] J. Li, M. Levine, X. An, X. Xu, H. He, Visual saliency based on scale-space analysis in the frequency domain, *IEEE Trans. Pattern Anal. Mach. Intell.* 35 (4) (2013) 996–1010.
- [67] X. Hou, J. Harel, C. Koch, Image signature: Highlighting sparse salient regions, *IEEE Trans. Pattern Anal. Mach. Intell.* 34 (1).

- [68] B. Schauerte, R. Stiefelhagen, Predicting human gaze using quaternion dct image signature saliency and face detection, in: WACV, 2012.
- [69] J. Harel, C. Koch, P. Perona, Graph-based visual saliency, in: NIPS, 2007.
- [70] V. Gopalakrishnan, Y. Hu, D. Rajan, Random walks on graphs to model saliency in images, in: CVPR, 2009.
- [71] W. Wang, Y. Wang, Q. Huang, W. Gao, Measuring visual saliency by site entropy rate, in: CVPR, 2010.
- [72] A. Toet, Computational versus psychophysical bottom-up image saliency: A comparative evaluation study, *IEEE Trans. Pattern Anal. Mach. Intell.* 33 (11) (2011) 2131–2146.
- [73] A. Borji, L. Itti, State-of-the-art in visual attention modeling, *IEEE Trans. Pattern Anal. Mach. Intell.* 35 (PrePrints) (2013) 185–207.
- [74] W. Kienzle, M. O. Franz, B. Schlkopf, F. A. Wichmann, Center-surround patterns emerge as optimal predictors for human saccade targets, *J. Vis.* 9 (5).
- [75] T. Judd, K. Ehinger, F. Durand, A. Torralba, Learning to predict where humans look, in: ICCV, 2009.
- [76] Q. Wang, Y. Yuan, P. Yan, X. Li, Saliency detection by multiple-instance learning, *IEEE Trans. Cybern.* 43 (2) (2013) 660–672.
- [77] Q. Zhao, C. Koch, Learning a saliency map using fixated locations in natural scenes, *J. Vis.* 11 (3).
- [78] Q. Zhao, C. Koch, Learning visual saliency by combining feature maps in a nonlinear manner using adaboost, *J. Vis.* 12 (6).
- [79] E. Vig, M. Dorr, D. Cox, Large-scale optimization of hierarchical features for saliency prediction in natural images, in: CVPR, 2014.
- [80] M. Kümmerer, L. Theis, M. Bethge, Deep gaze i: Boosting saliency prediction with feature maps trained on imagenet, in: ICLR Workshop, 2015.
- [81] J. Pan, K. McGuinness, E. Sayrol, N. O’Connor, X. Giro-i Nieto, Shallow and deep convolutional networks for saliency prediction, in: CVPR, 2016.
- [82] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition, in: ICLR, 2015.
- [83] N. Liu, J. Han, D. Zhang, S. Wen, T. Liu, Predicting eye fixations using convolutional neural networks, in: CVPR, 2015.
- [84] X. Huang, C. Shen, X. Boix, Q. Zhao, Salicon: Reducing the semantic gap in saliency prediction by adapting deep neural networks, in: ICCV, 2015.

- [85] S. Jetley, N. Murray, E. Vig, End-to-end saliency mapping via probability distribution prediction, in: CVPR, 2016.
- [86] A. Borji, H. R.-Tavakoli, D. N. Sihite, L. Itti, Analysis of scores, datasets, and models in visual saliency prediction, in: ICCV, 2013.
- [87] M. Cimpoi, S. Maji, A. Vedaldi, Deep filter banks for texture recognition and segmentation, in: CVPR, 2015.
- [88] A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with deep convolutional neural networks, in: NIPS, 2012.
- [89] A. Oliva, Gist of the scene, *Neurobiol. Atten.* 696 (64) (2005) 251–258.
- [90] J. Xu, M. Jiang, S. Wang, M. S. Kankanhalli, Q. Zhao, Predicting human gaze beyond pixels, *J. Vis* 14 (1) (2014) 1–20.
- [91] M. Minsky, S. Papert, *Perceptrons: an introduction to computational geometry*, MIT Press, 1969.
- [92] W. Schmidt, M. Kraaijveld, R. Duin, Feedforward neural networks with random weights, in: ICPR, 1992.
- [93] Y.-H. Pao, G.-H. Park, D. J. Sobajic, Learning and generalization characteristics of the random vector functional-link net, *Neurocomputing* 6 (2) (1994) 163 – 180.
- [94] G.-B. Huang, Q.-Y. Zhu, C.-K. Siew, Extreme learning machine: a new learning scheme of feedforward neural networks, in: IJCNN, 2004.
- [95] G.-B. Huang, Q.-Y. Zhu, C.-K. Siew, Extreme learning machine: Theory and applications, *Neurocomput.* 70 (2006) 489–501.
- [96] G.-B. Huang, Q.-Y. Zhu, C.-K. Siew, Real-time learning capability of neural networks, *IEEE Trans. Neural Netw.* 17 (4) (2006) 863–878.
- [97] G. Huang, G.-B. Huang, S. Song, K. You, Trends in extreme learning machines: A review, *Neural Netw.* 61 (2015) 32–48.
- [98] B. W. Tatler, B. T. Vincent, The prominence of behavioural biases in eye guidance, *Vis. Cogn.* 17 (6-7) (2009) 1029–1054.
- [99] Z. Bylinskii, T. Judd, A. Borji, L. Itti, F. Durand, A. Oliva, A. Torralba, Mit saliency benchmark (Jul 2016).
- [100] A. Borji, L. Itti, Cat2000: A large scale fixation dataset for boosting saliency research, in: CVPR workshops, 2015.
- [101] A. Garcia-Diaz, X. Fdez-Vidal, X. Pardo, R. Dosil, Decorrelation and distinctiveness provide with human-like saliency, in: ACIVS, 2009.

- [102] M. Cornia, L. Baraldi, G. Serra, R. Cucchiara, A deep multi-level network for saliency prediction, in: ICPR, 2016.
- [103] J. Zhang, S. Sclaroff, saliency detection: a Boolean map approach, in: ICCV, 2013.