

Contextual Encoder-Decoder Network for Visual Saliency Prediction

Alexander Kroner^{a,b,*}, Mario Senden^{a,b}, Kurt Driessens^c, Rainer Goebel^{a,b,d}

^a*Department of Cognitive Neuroscience, Faculty of Psychology and Neuroscience,
Maastricht University, Maastricht, The Netherlands*

^b*Maastricht Brain Imaging Centre, Faculty of Psychology and Neuroscience,
Maastricht University, Maastricht, The Netherlands*

^c*Department of Data Science and Knowledge Engineering, Faculty of Science and Engineering,
Maastricht University, Maastricht, The Netherlands*

^d*Department of Neuroimaging and Neuromodeling, Netherlands Institute for Neuroscience,
Royal Netherlands Academy of Arts and Sciences (KNAW), Amsterdam, The Netherlands*

Abstract

Predicting salient regions in natural images requires the detection of objects that are present in a scene. To develop robust representations for this challenging task, high-level visual features at multiple spatial scales must be extracted and augmented with contextual information. However, existing models aimed at explaining human fixation maps do not incorporate such a mechanism explicitly. Here we propose an approach based on a convolutional neural network pre-trained on a large-scale image classification task. The architecture forms an encoder-decoder structure and includes a module with multiple convolutional layers at different dilation rates to capture multi-scale features in parallel. Moreover, we combine the resulting representations with global scene information for accurately predicting visual saliency. Our model achieves competitive and consistent results across multiple evaluation metrics on two public saliency benchmarks and we demonstrate the effectiveness of the suggested approach on five datasets and selected examples. Compared to state of the art approaches, the network is based on a lightweight image classification backbone and hence presents a suitable choice for applications with limited computational resources, such as (virtual) robotic systems, to estimate human fixations across complex natural scenes. Our TensorFlow implementation is openly available at <https://github.com/alexanderkroner/saliency>.

1. Introduction

Humans demonstrate a remarkable ability to obtain relevant information from complex visual scenes [1, 2]. Overt attention is the mechanism that governs the processing of stimuli by directing gaze towards a spatial location within the visual field [3]. This sequential selection ensures that the eyes sample prioritized aspects from all available information to reduce the cost of cortical computation [4]. In addition, only a small central region of the retina, known as the fovea, transforms incoming light into neural responses with high spatial resolution, whereas acuity

*Corresponding author. *Email address:* kroner.contact@gmail.com



Figure 1: A visualization of four natural images with the corresponding empirical fixation maps, our model predictions, and estimated maps based on the work by Itti et al. [9]. The network proposed in this study was not trained on the stimuli shown here and thus exhibits its generalization ability to unseen instances. All image examples demonstrate a qualitative agreement of our model with the ground truth data, assigning high saliency to regions that contain semantic information, such as a door (a), flower (b), face (c), or text (d). On the contrary, the approach by Itti et al. [9] detected low-level feature contrasts and wrongly predicted high values at object boundaries rather than their center.

decreases rapidly towards the periphery [5, 6]. Given the limited number of photoreceptors in the eye, this arrangement allows to optimally process visual signals from its environment [7]. The function of fixations is thus to resolve the trade-off between coverage and sampling resolution of the whole visual field [8].

The spatial allocation of attention when viewing natural images is commonly represented in the form of topographic saliency maps that depict which parts of a scene attract fixations reliably. Identifying the underlying properties of these regions would allow us to predict human fixation patterns and gain a deeper understanding of the processes that lead to the observed behavior. In computer vision, this challenging problem has originally been approached using models rooted in *Feature Integration Theory* [10]. The theory suggests that early visual features must first be registered in parallel before serial shifts of overt attention combine them into unitary object-based representations. This two-stage account of visual processing has emphasized the role of stimulus properties for explaining human gaze. In consequence, the development of feature-driven models has been considered sufficient to enable the prediction of fixation patterns under task-free viewing conditions. Koch and Ullman [11] have introduced the notion of a central saliency map which integrates low-level information and serves as the basis for eye movements. This has resulted in a first model implementation by Itti et al. [9] that influenced later work on biologically-inspired architectures.

With the advent of deep neural network solutions for visual tasks such as image classification [12], saliency modeling has also undergone a paradigm shift from manual feature engineering towards automatic representation learning. In this work, we leveraged the capability of

convolutional neural networks (CNNs) to extract relevant features from raw images and decode them towards a distribution of saliency across arbitrary scenes. Compared to the seminal work by Itti et al. [9], this approach allows predictions to be based on semantic information instead of low-level feature contrasts (see Figure 1). This choice was motivated by studies demonstrating the importance of high-level image content for attentional selection in natural images [13, 14].

Furthermore, it is expected that complex representations at multiple spatial scales are necessary for accurate predictions of human fixation patterns. We therefore incorporated a contextual module that samples multi-scale information and augments it with global scene features. The contribution of the contextual module to the overall performance was assessed and final results were compared to previous work on two public saliency benchmarks. We achieved predictive accuracy on unseen test instances at the level of current state of the art approaches, while utilizing a computationally less expensive network backbone with roughly one order of magnitude fewer processing layers. This makes our model suitable for applications in (virtual) robotic environments, as demonstrated by Bornet et al. [15], and we developed a webcam-based interface for saliency prediction in the browser with only moderate hardware requirements (see <https://storage.googleapis.com/msi-net/demo/index.html>).

2. Related Work

Early approaches towards computational models of visual attention were defined in terms of different theoretical frameworks, including Bayesian [16] and graph-based formulations [17]. The former was based on the notion of self-information derived from a probability distribution over linear visual features as acquired from natural scenes. The latter framed saliency as the dissimilarity between nodes in a fully-connected directed graph that represents all image locations in a feature map. Hou and Zhang [18] have instead proposed an approach where images were transformed to the log spectrum and saliency emerged from the spectral residual after removing statistically redundant components. A mechanism inspired more by biological than mathematical principles was first implemented and described in the seminal work by Itti et al. [9]. Their model captures center-surround differences at multiple spatial scales with respect to three basic feature channels: color, intensity, and orientation. After normalization of activity levels, the output is fed into a common saliency map depicting local conspicuity in static scenes. This standard cognitive architecture has since been augmented with additional feature channels that capture semantic image content, such as faces and text [19].

With the large-scale acquisition of eye tracking measurements under natural viewing conditions, data-driven machine learning techniques became more practicable. Judd et al. [20] introduced a model based on support vector machines to estimate fixation densities from a set of low-, mid-, and high-level visual features. While this approach still relied on a hypothesis specifying which image properties would successfully contribute to the prediction of saliency, it marked the beginning of a progression from manual engineering to automatic learning of features. This development has ultimately led to applying deep neural networks with emergent representations for the estimation of human fixation patterns. Vig et al. [21] were the first to train an ensemble of shallow CNNs to derive saliency maps from natural images in an end-to-end fashion, but failed to capture object information due to limited network depth.

Later attempts addressed that shortcoming by taking advantage of classification architectures pre-trained on the *ImageNet* database [22]. This choice was motivated by the finding that features extracted from CNNs generalize well to other visual tasks [23]. Consequently, *DeepGaze I* [24] and *II* [25] employed a pre-trained classification model to read out salient image locations

from a small subset of encoding layers. This is similar to the network by Cornia et al. [26] which utilizes the output at three stages of the hierarchy. Oyama and Yamanaka [27] demonstrated that classification performance of pre-trained architectures strongly correlates with the accuracy of saliency predictions, highlighting the importance of object information. Related approaches also focused on the potential benefits of incorporating activation from both coarse and fine image resolutions [28], and recurrent connections to capture long-range spatial dependencies in convolutional feature maps [29, 30]. Our model explicitly combines semantic representations at multiple spatial scales to include contextual information in the predictive process. For a more complete account of existing saliency architectures, we refer the interested reader to a comprehensive review by Borji [31].

3. Methods

We propose a new CNN architecture with modules adapted from the semantic segmentation literature to predict fixation density maps of the same image resolution as the input. Our approach is based on a large body of research regarding saliency models that leverage object-specific features and functionally replicate human behavior under free-viewing conditions. In the following sections, we describe our contributions to this challenging task.

3.1. Architecture

Image-to-image learning problems require the preservation of spatial features throughout the whole processing stream. As a consequence, our network does not include any fully-connected layers and reduces the number of downsampling operations inherent to classification models. We adapted the popular *VGG16* architecture [32] as an image encoder by reusing the pre-trained convolutional layers to extract increasingly complex features along its hierarchy. Striding in the last two pooling layers was removed, which yields spatial representations at $\frac{1}{8}$ of their original input size. All subsequent convolutional encoding layers were then dilated at a rate of 2 by expanding their kernel, and thereby increased the receptive field to compensate for the higher resolution [33]. This modification still allowed us to initialize the model with pre-trained weights since the number of trainable parameters remained unchanged. Prior work has shown the effectiveness of this approach in the context of saliency prediction problems [29, 30].

For related visual tasks such as semantic segmentation, information distributed over convolutional layers at different levels of the hierarchy can aid the preservation of fine spatial details [34, 35]. The prediction of fixation density maps does not require accurate class boundaries but still benefits from combined mid- to high-level feature responses [24, 25, 26]. Hence, we adapted the multi-level design proposed by Cornia et al. [26] and concatenated the output from layers 10, 14, and 18 into a common tensor with 1,280 activation maps.

This representation constitutes the input to an *Atrous Spatial Pyramid Pooling* (ASPP) module [36]. It utilizes several convolutional layers with different dilation factors in parallel to capture multi-scale image information. Additionally, we incorporated scene content via global average pooling over the final encoder output, as motivated by the study of Torralba et al. [37] who stated that contextual information plays an important role for the allocation of attention. Our implementation of the ASPP architecture thus closely follows the modifications proposed by Chen et al. [38]. These authors augmented multi-scale information with global context and demonstrated performance improvements on semantic segmentation tasks.

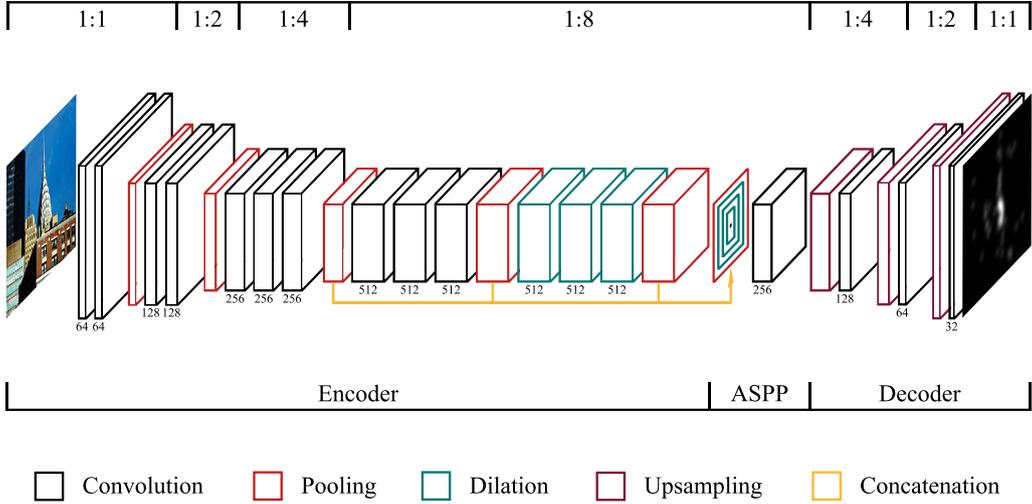


Figure 2: An illustration of the modules that constitute our encoder-decoder architecture. The VGG16 backbone was modified to account for the requirements of dense prediction tasks by omitting feature downsampling in the last two max-pooling layers. Multi-level activations were then forwarded to the ASPP module, which captured information at different spatial scales in parallel. Finally, the input image dimensions were restored via the decoder network. Subscripts beneath convolutional layers denote the corresponding number of feature maps.

In this work, we laid out three convolutional layers with kernel sizes of 3×3 and dilation rates of 4, 8, and 12 in parallel, together with a 1×1 convolutional layer that could not learn new spatial dependencies but nonlinearly combined existing feature maps. Image-level context was represented as the output after global average pooling (i.e. after averaging the entries of a tensor across both spatial dimensions to a single value) and then brought to the same resolution as all other representations via bilinear upsampling, followed by another point-wise convolutional operation. Each of the five branches in the module contains 256 filters, which resulted in an aggregated tensor of 1,280 feature maps. Finally, the combined output was forwarded to a 1×1 convolutional layer with 256 channels that contained the resulting multi-scale responses.

To restore the original image resolution, extracted features were processed by a series of convolutional and upsampling layers. Previous work on saliency prediction has commonly utilized bilinear interpolation for that task [29, 30], but we argue that a carefully chosen decoder architecture, similar to the model by Pan et al. [39], results in better approximations. Here we employed three upsampling blocks consisting of a bilinear scaling operation, which doubled the number of rows and columns, and a subsequent convolutional layer with kernel size 3×3 . This setup has previously been shown to prevent checkerboard artifacts in the upsampled image space in contrast to deconvolution [40]. Besides an increase of resolution throughout the decoder, the amount of channels was halved in each block to yield 32 feature maps. Our last network layer transformed activations into a continuous saliency distribution by applying a final 3×3 convolution. The outputs of all but the last linear layer were modified via rectified linear units. Figure 2 visualizes the overall architecture design as described in this section.

3.2. Training

Weight values from the ASPP module and decoder were initialized according to the *Xavier* method by Glorot and Bengio [41]. It specifies parameter values as samples drawn from a uniform distribution with zero mean and a variance depending on the total number of incoming and outgoing connections. Such initialization schemes are demonstrably important for training deep neural networks successfully from scratch [42]. The encoding layers were based on the VGG16 architecture pre-trained on both *ImageNet* [22] and *Places2* [43] data towards object and scene classification respectively.

We normalized the model output such that all values are non-negative with unit sum. The estimation of saliency maps can hence be regarded as a *probability distribution prediction* task as formulated by Jetley et al. [44]. To determine the difference between an estimated and a target distribution, the *Kullback-Leibler* (KL) divergence is an appropriate measure rooted in information theory to quantify the statistical distance D . This can be defined as follows:

$$D_{KL}(P \parallel Q) = \sum_i Q_i \ln\left(\epsilon + \frac{Q_i}{\epsilon + P_i}\right) \quad (1)$$

Here, Q represents the target distribution, P its approximation, i each pixel index, and ϵ a regularization constant. Equation (1) served as the loss function which was gradually minimized via the *Adam* optimization algorithm [45]. We defined an upper learning rate of 10^{-6} and modified the weights in an online fashion due to a general inefficiency of batch training according to Wilson and Martinez [46]. Based on this general setup, we trained our network for 10 epochs and used the best-performing checkpoint for inference.

4. Experiments

The proposed encoder-decoder model was evaluated on five publicly available eye tracking datasets that yielded qualitative and quantitative results. First, we provide a brief description of the images and empirical measurements utilized in this study. Second, the different metrics commonly used to assess the predictive performance of saliency models are summarized. Finally, we report the contribution of our architecture design choices and benchmark the overall results against baselines and related work in computer vision.

4.1. Datasets

A prerequisite for the successful application of deep learning techniques is a wealth of annotated data. Fortunately, the growing interest in developing and evaluating fixation models has led to the release of large-scale eye tracking datasets such as *MIT1003* [20], *CAT2000* [47], *DUT-OMRON* [48], *PASCAL-S* [49], and *OSIE* [50]. The costly acquisition of measurements, however, is a limiting factor for the number of stimuli. New data collection methodologies have emerged that leverage webcam-based eye movements [51] or mouse movements [52] instead via crowdsourcing platforms. The latter approach resulted in the *SALICON* dataset, which consists of 10,000 training and 5,000 validation instances serving as a proxy for empirical gaze measurements. Due to its large size, we first trained our model on *SALICON* before fine-tuning the learned weights towards fixation predictions on either of the other datasets with the same optimization parameters. This widely adopted procedure has been shown to improve the accuracy of eye movement estimations despite some disagreement between data originating from gaze and mouse tracking experiments [53].

The images presented during the acquisition of saliency maps in all aforementioned datasets are largely based on natural scenes. Stimuli of CAT2000 additionally fall into predefined categories such as *Action*, *Fractal*, *Object*, or *Social*. Together with the corresponding fixation patterns, they constituted the input and desired output to our network architecture. In detail, we rescaled and padded all images from the SALICON and OSIE datasets to 240×320 pixels, the MIT1003, DUT-OMRON, and PASCAL-S datasets to 360×360 pixels, and the CAT2000 dataset to 216×384 pixels, such that the original aspect ratios were preserved. For the latter five eye tracking sets we defined 80% of the samples as training data and the remainder as validation examples with a minimum of 200 instances. The correct saliency distributions on test set images of MIT1003 and CAT2000 are held out and predictions must hence be submitted online for evaluation.

4.2. Metrics

Various measures are used in the literature and by benchmarks to evaluate the performance of fixation models. In practice, results are typically reported for all of them to include different notions about saliency and allow a fair comparison of model predictions [54, 55]. A set of nine metrics is commonly selected: *Kullback-Leibler divergence* (KLD), *Pearson’s correlation coefficient* (CC), *histogram intersection* (SIM), *Earth Mover’s distance* (EMD), *information gain* (IG), *normalized scanpath saliency* (NSS), and three variants of *area under ROC curve* (AUC-Judd, AUC-Borji, shuffled AUC). The former four are location-based metrics, which require ground truth maps as binary fixation matrices. By contrast, the remaining metrics quantify saliency approximations after convolving gaze locations with a Gaussian kernel and representing the target output as a probability distribution. We refer readers to an overview by Bylinskii et al. [56] for more information regarding the implementation details and properties of the stated measures.

In this work, we adopted KLD as an objective function and produced fixation density maps as output from our proposed network. This training setup is particularly sensitive to false negative predictions and thus the appropriate choice for applications aimed at salient target detection [56]. Defining the problem of saliency prediction in a probabilistic framework also enables fair model ranking on public benchmarks for the MIT1003, CAT2000, and SALICON datasets [54]. As a consequence, we evaluated our estimated gaze distributions without applying any metric-specific postprocessing methods.

4.3. Results

A quantitative comparison of results on independent test datasets was carried out to characterize how well our proposed network generalizes to unseen images. Here, we were mainly interested in estimating human eye movements and regarded mouse tracking measurements merely as a substitute for attention. The final outcome for the 2017 release of the SALICON dataset is therefore not reported in this work but our model results can be viewed on the public leaderboard¹ under the user name akroner.

To assess the predictive performance for eye tracking measurements, the MIT saliency benchmark [59] is commonly used to compare model results on two test datasets with respect to prior work. Final scores can then be submitted on a public leaderboard to allow fair model ranking on eight evaluation metrics. Table 1 summarizes our results on the test dataset of MIT1003, namely *MIT300* [60], in the context of previous approaches. The evaluation shows that our model only

¹<https://competitions.codalab.org/competitions/17136>

	<i>AUC-J</i> ↑	<i>SIM</i> ↑	<i>EMD</i> ↓	<i>AUC-B</i> ↑	<i>sAUC</i> ↑	<i>CC</i> ↑	<i>NSS</i> ↑	<i>KLD</i> ↓
DenseSal [27]	0.87	0.67	1.99	0.81	0.72	0.79	2.25	0.48
DPNSal [27]	0.87	0.69	2.05	0.80	0.74	0.82	2.41	0.91
SALICON [28] [†]	0.87	0.60	2.62	0.85	0.74	0.74	2.12	0.54
DSCLRCN [30]	0.87	0.68	2.17	0.79	0.72	0.80	2.35	0.95
DeepFix [57] [†]	0.87	0.67	2.04	0.80	0.71	0.78	2.26	0.63
EML-NET [58]	0.88	0.68	1.84	0.77	0.70	0.79	2.47	0.84
DeepGaze II [25]	0.88	0.46	3.98	0.86	0.72	0.52	1.29	0.96
SAM-VGG [29] [†]	0.87	0.67	2.14	0.78	0.71	0.77	2.30	1.13
ML-Net [26] [†]	0.85	0.59	2.63	0.75	0.70	0.67	2.05	1.10
SAM-ResNet [29]	0.87	0.68	2.15	0.78	0.70	0.78	2.34	1.27
DeepGaze I [24]	0.84	0.39	4.97	0.83	0.66	0.48	1.22	1.23
Judd [20]	0.81	0.42	4.45	0.80	0.60	0.47	1.18	1.12
eDN [21]	0.82	0.41	4.56	0.81	0.62	0.45	1.14	1.14
GBVS [17]	0.81	0.48	3.51	0.80	0.63	0.48	1.24	0.87
Itti [9]	0.75	0.44	4.26	0.74	0.63	0.37	0.97	1.03
SUN [16]	0.67	0.38	5.10	0.66	0.61	0.25	0.68	1.27
Ours [†]	0.87	0.68	1.99	0.82	0.72	0.79	2.27	0.66

Table 1: Quantitative results of our model for the MIT300 test set in the context of prior work. The first line separates deep learning approaches with architectures pre-trained on image classification (the superscript [†] represents models with a VGG16 backbone) from shallow networks and other machine learning methods. Entries between the second and the third line are models based on theoretical considerations and define a baseline rather than competitive performance. Arrows indicate whether the metrics assess similarity ↑ or dissimilarity ↓ between predictions and targets. The best results are marked in bold and models are sorted in descending order of their cumulative rank across a subset of weakly correlated evaluation measures within each group.

	<i>AUC-J</i> ↑	<i>SIM</i> ↑	<i>EMD</i> ↓	<i>AUC-B</i> ↑	<i>sAUC</i> ↑	<i>CC</i> ↑	<i>NSS</i> ↑	<i>KLD</i> ↓
SAM-VGG [29] [†]	0.88	0.76	1.07	0.79	0.58	0.89	2.38	0.54
SAM-ResNet [29]	0.88	0.77	1.04	0.80	0.58	0.89	2.38	0.56
DeepFix [57] [†]	0.87	0.74	1.15	0.81	0.58	0.87	2.28	0.37
EML-NET [58]	0.87	0.75	1.05	0.79	0.59	0.88	2.38	0.96
Judd [20]	0.84	0.46	3.60	0.84	0.56	0.54	1.30	0.94
eDN [21]	0.85	0.52	2.64	0.84	0.55	0.54	1.30	0.97
Itti [9]	0.77	0.48	3.44	0.76	0.59	0.42	1.06	0.92
GBVS [17]	0.80	0.51	2.99	0.79	0.58	0.50	1.23	0.80
SUN [16]	0.70	0.43	3.42	0.69	0.57	0.30	0.77	2.22
Ours [†]	0.88	0.75	1.07	0.82	0.59	0.87	2.30	0.36

Table 2: Quantitative results of our model for the CAT2000 test set in the context of prior work. The first line separates deep learning approaches with architectures pre-trained on image classification (the superscript [†] represents models with a VGG16 backbone) from shallow networks and other machine learning methods. Entries between the second and third lines are models based on theoretical considerations and define a baseline rather than competitive performance. Arrows indicate whether the metrics assess similarity ↑ or dissimilarity ↓ between predictions and targets. The best results are marked in bold and models are sorted in descending order of their cumulative rank across a subset of weakly correlated evaluation measures within each group.

marginally failed to achieve state-of-the-art performance on any of the individual metrics. When computing the cumulative rank (i.e. the sum of ranks according to the standard competition ranking procedure) on a subset of weakly correlated measures (*sAUC*, *CC*, *KLD*) [55, 56], we ranked third behind the two architectures *DenseSal* and *DPNSal* from Oyama and Yamanaka [27]. However, their approaches were based on a pre-trained *Densely Connected Convolutional Network*

<i>Parameters</i>	
DeepGaze I [24]	3,750,913
ML-Net [26] [†]	15,452,145
DeepGaze II [25]	20,065,973
DenseSal [27]	26,474,209
SALICON [28] [†]	29,429,889
DSCLRCN [30]	30,338,437
DeepFix [57] [†]	35,455,617
SAM-VGG [29] [†]	51,835,841
SAM-ResNet [29]	70,093,441
DPNSal [27]	76,536,513
EML-NET [58]	84,728,569
Ours	24,934,209

Table 3: The number of trainable parameters for all deep learning models listed in Table 1 that are competing in the MIT300 saliency benchmark. Entries of prior work are sorted according to increasing network complexity and the superscript [†] represents pre-trained models with a VGG16 backbone.

	<i>Sizes</i>	<i>Speed</i>	<i>Memory</i>	<i>Computations</i>
MIT1003	360 × 360 px	43 FPS	194 MB	75 GFLOPS
CAT2000	216 × 384 px	56 FPS	175 MB	48 GFLOPS
DUT-OMRON	360 × 360 px	43 FPS	194 MB	75 GFLOPS
PASCAL-S	360 × 360 px	43 FPS	194 MB	75 GFLOPS
OSIE	240 × 320 px	58 FPS	173 MB	44 GFLOPS

Table 4: The results after evaluating our model with respect to its computational efficiency. We tested five versions trained on different eye tracking datasets, each receiving input images of their preferred sizes in *pixels* (px). After running each network on 10,000 test set instances from the ImageNet database for 10 times, we averaged the inference speed and described the results in *frames per second* (FPS). All settings demonstrated consistent outcomes with a standard deviation of less than 1 FPS. The minimal GPU memory utilization was measured with TensorFlow in *megabytes* (MB) and included the requirements for initializing a testing session. Finally, we estimated the *floating point operations per second* (FLOPS) at a scale of 9 orders of magnitude.

<i>Hardware specifications</i>		<i>Software specifications</i>	
GPU	NVIDIA TITAN Xp	TensorFlow	1.14.0
CPU	Intel Xeon E5-1650	CUDA	10.0
RAM	32 GB DDR4	cuDNN	7.5
HDD	256 GB SSD	GPU driver	418.74

Table 5: Details regarding the hardware and software specifications used throughout our evaluation of computational efficiency. The system ran under the Debian 9 operating system and we minimized usage of the computer during the experiments to avoid interference with measurements of inference speed.

with 161 layers [61] and *Dual Path Network* with 131 layers [62] respectively, both of which are computationally far more expensive than the VGG16 model used in this work (see Table 5 by Oyama and Yamanaka [27] for a comparison of the computational efficiency). Furthermore, DenseSal and DPNSal implemented a multi-path design where two images of different resolutions are simultaneously fed to the network, which substantially reduces the execution speed compared to single-stream architectures. Among all entries of the MIT300 benchmark with a VGG16 backbone [26, 28, 29, 57], our model clearly achieved the highest performance.

			<i>AUC-J</i> \uparrow	<i>SIM</i> \uparrow	<i>EMD</i> \downarrow	<i>AUC-B</i> \uparrow	<i>sAUC</i> \uparrow	<i>CC</i> \uparrow	<i>NSS</i> \uparrow	<i>KLD</i> \downarrow
MIT1003	\oplus ASPP	μ	0.899*	0.602*	2.430*	0.832*	0.713	0.741*	2.663*	0.818*
		σ	0.001	0.004	0.033	0.005	0.003	0.003	0.012	0.051
	\ominus ASPP	μ	0.890	0.573	2.645	0.827	0.720	0.700	2.540	0.867
		σ	0.001	0.005	0.033	0.006	0.003	0.004	0.016	0.042
CAT2000	\oplus ASPP	μ	0.882*	0.734*	2.553*	0.812	0.582	0.854*	2.359*	0.430*
		σ	0.000	0.002	0.025	0.005	0.003	0.003	0.007	0.010
	\ominus ASPP	μ	0.873	0.683	2.975	0.824	0.591	0.770	2.092	0.501
		σ	0.000	0.002	0.033	0.004	0.002	0.002	0.006	0.013
DUT-OMRON	\oplus ASPP	μ	0.921*	0.649*	1.155*	0.864*	0.775	0.776*	2.873*	0.634*
		σ	0.001	0.002	0.016	0.005	0.003	0.002	0.008	0.023
	\ominus ASPP	μ	0.915	0.627	1.237	0.855	0.777	0.748	2.789	0.695
		σ	0.001	0.003	0.022	0.004	0.003	0.002	0.011	0.031
PASCAL-S	\oplus ASPP	μ	0.914*	0.667*	1.015*	0.853*	0.701	0.818*	2.610*	0.645*
		σ	0.000	0.003	0.010	0.004	0.003	0.002	0.008	0.044
	\ominus ASPP	μ	0.901	0.610	1.195	0.831	0.715	0.759	2.420	0.720
		σ	0.001	0.004	0.013	0.004	0.003	0.004	0.015	0.022
OSIE	\oplus ASPP	μ	0.918*	0.648*	1.647*	0.816*	0.788*	0.808*	3.010	0.749
		σ	0.001	0.002	0.017	0.004	0.004	0.002	0.010	0.031
	\ominus ASPP	μ	0.916	0.641	1.733	0.808	0.781	0.804	3.000	0.767
		σ	0.001	0.002	0.025	0.008	0.006	0.002	0.010	0.046

Table 6: A summary of the quantitative results for the models with \oplus and without \ominus an ASPP module. The evaluation was carried out on five eye tracking datasets respectively. Each network was independently trained 10 times resulting in a distribution of values characterized by the mean μ and standard deviation σ . The star * denotes a significant increase of performance between the two conditions according to a one sided paired t-test. Arrows indicate whether the metrics assess similarity \uparrow or dissimilarity \downarrow between predictions and targets. The best results are marked in bold.

We further evaluated the model complexity of all relevant deep learning approaches listed in Table 1. The number of trainable parameters was computed based on either the official code repository or a replication of the described architectures. In case a reimplementation was not possible, we faithfully estimated a lower bound given the pre-trained classification network. Table 3 summarizes the findings and shows that our model compares favorably to the best-performing approaches. While the number of parameters provides an indication about the computational efficiency of an algorithm, more measures are needed. Therefore, we recorded the inference speed and GPU memory consumption of our model and calculated the number of computations (see Table 4) for our given hardware and software specifications (see Table 5). The results highlight that our approach achieves fast inference speed combined with a low GPU memory footprint, and thus enables applications to systems constrained by computational resources.

Table 2 demonstrates that we obtained state-of-the-art scores for the CAT2000 test dataset regarding the AUC-J, sAUC, and KLD evaluation metrics, and competitive results on the remaining measures. The cumulative rank (as computed above) suggests that our model outperformed all previous approaches, including the ones based on a pre-trained VGG16 classification network [29, 57]. Our final evaluation results for both the MIT300 and CAT2000 datasets can be viewed on the MIT saliency benchmark under the model name *MSI-Net*, representing our multi-scale information network. Qualitatively, the proposed architecture successfully captures semantically meaningful image features such as faces and text towards the prediction of saliency, as can be seen in Figure 1. Unfortunately, a visual comparison with the results from prior work was not possible since most models are not openly available.

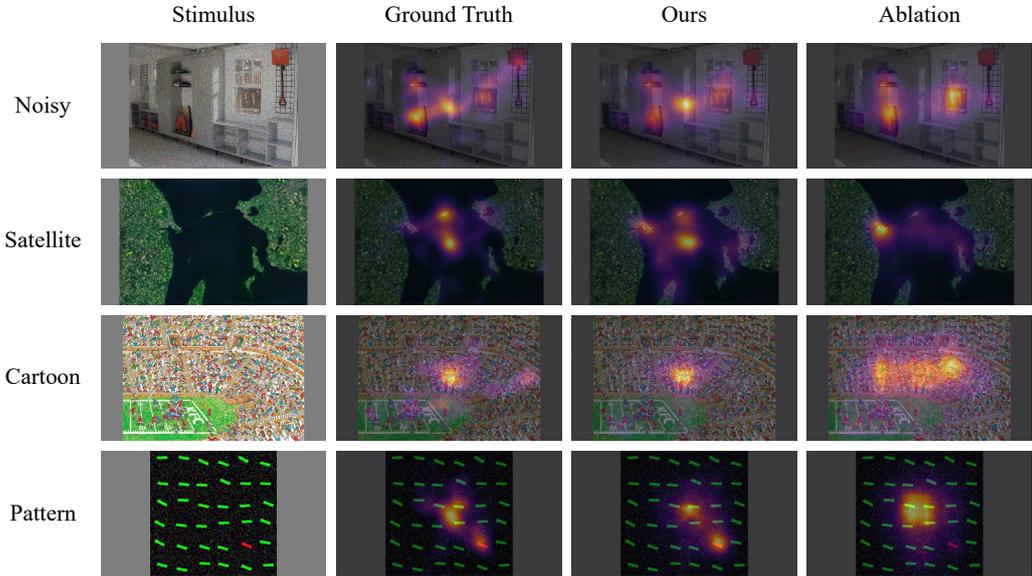


Figure 3: A visualization of four example images from the CAT2000 validation set with the corresponding fixation heat maps, our best model predictions, and estimated maps based on the ablated network. The qualitative results indicate that multi-scale information augmented with global context enables a more accurate estimation of salient image regions.

	<i>AUC-J</i> \uparrow	<i>SIM</i> \uparrow	<i>EMD</i> \downarrow	<i>AUC-B</i> \uparrow	<i>sAUC</i> \uparrow	<i>CC</i> \uparrow	<i>NSS</i> \uparrow	<i>KLD</i> \downarrow
Noisy	+0.010	+0.073	-0.506	-0.015	-0.009	+0.122	+0.395	-0.099
Satellite	+0.015	+0.060	-0.663	-0.012	-0.007	+0.137	+0.362	-0.100
Cartoon	+0.015	+0.066	-0.652	-0.010	-0.004	+0.125	+0.349	-0.121
Pattern	+0.011	+0.050	-0.437	-0.003	+0.001	+0.078	+0.277	-0.065

Table 7: A list of the four image categories from the CAT2000 validation set that showed the largest average improvement by the ASPP architecture based on the cumulative rank across a subset of weakly correlated evaluation measures. Arrows indicate whether the metrics assess similarity \uparrow or dissimilarity \downarrow between predictions and targets. Results that improved on the respective metric are marked in green, whereas results that impaired performance are marked in red.

To quantify the contribution of multi-scale contextual information to the overall performance, we conducted a model ablation analysis. A baseline architecture without the ASPP module was constructed by replacing the five parallel convolutional layers with a single 3×3 convolutional operation that resulted in 1,280 activation maps. This representation was then forwarded to a 1×1 convolutional layer with 256 channels. While the total number of feature maps stayed constant, the amount of trainable parameters increased in this ablation setting. Table 6 summarizes the results according to validation instances of five eye tracking datasets for the model with and without an ASPP module. It can be seen that our multi-scale architecture reached significantly higher performance (one tailed paired t-test) on most metrics and is therefore able to leverage the information captured by convolutional layers with different receptive field sizes. An ablation analysis of the multi-level component adapted from Cornia et al. [26] can be viewed in the Appendix A.

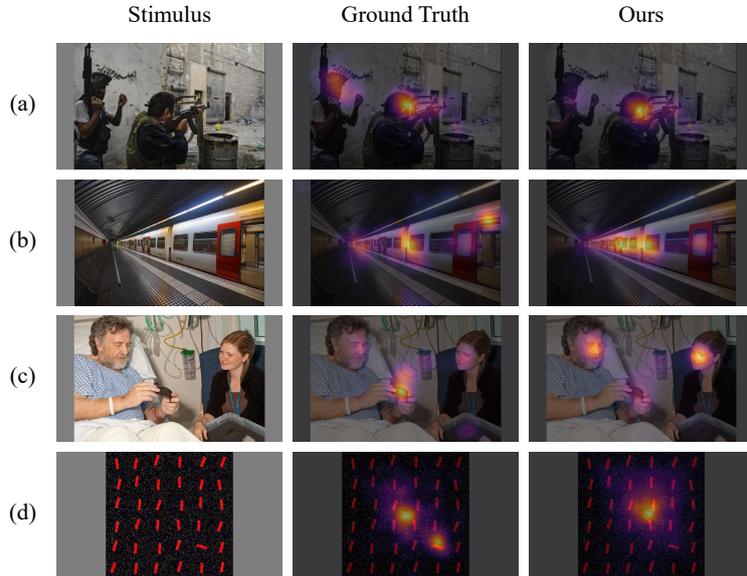


Figure 4: A visualization of four example images from the CAT2000 validation set with the corresponding eye movement patterns and our model predictions. The stimuli demonstrate cases with a qualitative disagreement between the estimated saliency maps and ground truth data. Here, our model failed to capture an occluded face (a), small text (b), direction of gaze (c), and low-level feature contrast (d).

The categorical organization of the CAT2000 database also allowed us to quantify the improvements by the ASPP module with respect to individual image classes. Table 7 lists the four categories that benefited the most from multi-scale information across the subset of evaluation metrics on the validation set: *Noisy*, *Satellite*, *Cartoon*, *Pattern*. To understand the measured changes in predictive performance, it is instructive to inspect qualitative results of one representative example for each image category (see Figure 3). The visualizations demonstrate that large receptive fields allow the reweighting of relative importance assigned to image locations (*Noisy*, *Satellite*, *Cartoon*), detection of a central fixation bias (*Noisy*, *Satellite*, *Cartoon*), and allocation of saliency to a low-level color contrast that pops out from an array of distractors (*Pattern*).

5. Discussion

Our proposed encoder-decoder model clearly demonstrated competitive performance on two datasets towards visual saliency prediction. The ASPP module incorporated multi-scale information and global context based on semantic feature representations, which significantly improved the results both qualitatively and quantitatively on five eye tracking datasets. This suggests that convolutional layers with large receptive fields at different dilation factors can enable a more holistic estimation of salient image regions in complex scenes. Moreover, our approach is computationally lightweight compared to prior state-of-the-art approaches and could thus be implemented in (virtual) robotic systems that require computational efficiency. It also outperformed all other networks defined with a pre-trained VGG16 backbone as calculated by the cumulative rank on a subset of evaluation metrics to resolve some of the inconsistencies in ranking models by a single measure or a set of correlated ones [55, 56].

Further improvements of benchmark results could potentially be achieved by a number of additions to the processing pipeline. Our model demonstrates a learned preference for predicting fixations in central regions of images, but we expect performance gains from modeling the central bias in scene viewing explicitly [24, 25, 26, 29, 57]. Additionally, Bylinskii et al. [59] summarized open problems for correctly assigning saliency in natural images, such as robustness in detecting semantic features, implied gaze and motion, and importance weighting of multiple salient regions. While the latter was addressed in this study, Figure 4 indicates that the remaining obstacles still persist for our proposed model.

Overcoming these issues requires a higher-level scene understanding that models object interactions and predicts implicit gaze and motion cues from static images. Robust object recognition could however be achieved through more recent classification networks as feature extractors [27] at the cost of added computational complexity. However, this study does not investigate whether the benefits of the proposed modifications generalize to other pre-trained architectures. That would constitute an interesting avenue for future research. To detect salient items in search array stimuli (see Figure 4d), a mechanism that additionally captures low-level feature contrasts might explain the empirical data better. Besides architectural changes, data augmentation in the context of saliency prediction tasks demonstrated its efficiency to improve the robustness of deep neural networks according to Che et al. [63]. These authors stated that visual transformations such as mirroring or inversion revealed a low impact on human gaze during scene viewing and could hence form an addition to future work on saliency models.

Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgement

This study has received funding from the European Union’s Horizon 2020 Framework Programme for Research and Innovation under the Specific Grant Agreement Nos. 720270 (Human Brain Project SGA1) and 785907 (Human Brain Project SGA2). Furthermore, we gratefully acknowledge the support of NVIDIA Corporation with the donation of a Titan X Pascal GPU used for this research.

References

- [1] J. Jonides, D. E. Irwin, S. Yantis, Integrating visual information from successive fixations, *Science* 215 (1982) 192–194.
- [2] D. E. Irwin, Information integration across saccadic eye movements, *Cognitive Psychology* 23 (1991) 420–456.
- [3] M. I. Posner, Orienting of attention, *Quarterly Journal of Experimental Psychology* 32 (1980) 3–25.
- [4] P. Lennie, The cost of cortical computation, *Current Biology* 13 (2003) 493–497.
- [5] A. Cowey, E. Rolls, Human cortical magnification factor and its relation to visual acuity, *Experimental Brain Research* 21 (1974) 447–454.
- [6] M. A. Berkley, F. Kitterle, D. W. Watkins, Grating visibility as a function of orientation and retinal eccentricity, *Vision Research* 15 (1975) 239–244.
- [7] B. Cheung, E. Weiss, B. Olshausen, Emergence of foveal image sampling from learning to attend in visual scenes, *arXiv preprint arXiv:1611.09430* (2016).
- [8] K. R. Gegenfurtner, The interaction between vision and eye movements, *Perception* 45 (2016) 1333–1357.

- [9] L. Itti, C. Koch, E. Niebur, A model of saliency-based visual attention for rapid scene analysis, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 20 (1998) 1254–1259.
- [10] A. M. Treisman, G. Gelade, A feature-integration theory of attention, *Cognitive Psychology* 12 (1980) 97–136.
- [11] C. Koch, S. Ullman, Shifts in selective visual attention: Towards the underlying neural circuitry, *Human Neurobiology* 4 (1985) 219–227.
- [12] A. Krizhevsky, I. Sutskever, G. E. Hinton, ImageNet classification with deep convolutional neural networks, *Advances in Neural Information Processing Systems* 25 (2012) 1097–1105.
- [13] W. Einhäuser, M. Spain, P. Perona, Objects predict fixations better than early saliency, *Journal of Vision* 8 (2008) 18.
- [14] A. Nuthmann, J. M. Henderson, Object-based attentional selection in scene viewing, *Journal of Vision* 10 (2010) 20.
- [15] A. Bornet, J. Kaiser, A. Kroner, E. Falotico, A. Ambrosano, K. Cantero, M. H. Herzog, G. Francis, Running large-scale simulations on the Neurorobotics Platform to understand vision – the case of visual crowding, *Frontiers in Neurobotics* 13 (2019) 33.
- [16] L. Zhang, M. H. Tong, T. K. Marks, H. Shan, G. W. Cottrell, SUN: A Bayesian framework for saliency using natural statistics, *Journal of Vision* 8 (2008) 32.
- [17] J. Harel, C. Koch, P. Perona, Graph-based visual saliency, *Advances in Neural Information Processing Systems* 19 (2006) 545–552.
- [18] X. Hou, L. Zhang, Saliency detection: A spectral residual approach, *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2007) 1–8.
- [19] M. Cerf, E. P. Frady, C. Koch, Faces and text attract gaze independent of the task: Experimental data and computer model, *Journal of Vision* 9 (2009) 10.
- [20] T. Judd, K. Ehinger, F. Durand, A. Torralba, Learning to predict where humans look, *Proceedings of the International Conference on Computer Vision* (2009) 2106–2113.
- [21] E. Vig, M. Dorr, D. Cox, Large-scale optimization of hierarchical features for saliency prediction in natural images, *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2014) 2798–2805.
- [22] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, L. Fei-Fei, ImageNet: A large-scale hierarchical image database, *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2009) 248–255.
- [23] J. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng, T. Darrell, DeCAF: A deep convolutional activation feature for generic visual recognition, *Proceedings of the International Conference on Machine Learning* (2014) 647–655.
- [24] M. Kümmerer, L. Theis, M. Bethge, DeepGaze I: Boosting saliency prediction with feature maps trained on ImageNet, *arXiv preprint arXiv:1411.1045* (2014).
- [25] M. Kümmerer, T. S. Wallis, M. Bethge, DeepGaze II: Reading fixations from deep features trained on object recognition, *arXiv preprint arXiv:1610.01563* (2016).
- [26] M. Cornia, L. Baraldi, G. Serra, R. Cucchiara, A deep multi-level network for saliency prediction, *Proceedings of the International Conference on Pattern Recognition* (2016) 3488–3493.
- [27] T. Oyama, T. Yamanaka, Influence of image classification accuracy on saliency map estimation, *arXiv preprint arXiv:1807.10657* (2018).
- [28] X. Huang, C. Shen, X. Boix, Q. Zhao, SALICON: Reducing the semantic gap in saliency prediction by adapting deep neural networks, *Proceedings of the International Conference on Computer Vision* (2015) 262–270.
- [29] M. Cornia, L. Baraldi, G. Serra, R. Cucchiara, Predicting human eye fixations via an LSTM-based saliency attentive model, *IEEE Transactions on Image Processing* 27 (2018) 5142–5154.
- [30] N. Liu, J. Han, A deep spatial contextual long-term recurrent convolutional network for saliency detection, *IEEE Transactions on Image Processing* 27 (2018) 3264–3274.
- [31] A. Borji, Saliency prediction in the deep learning era: An empirical investigation, *arXiv preprint arXiv:1810.03716* (2018).
- [32] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition, *arXiv preprint arXiv:1409.1556* (2014).
- [33] F. Yu, V. Koltun, Multi-scale context aggregation by dilated convolutions, *arXiv preprint arXiv:1511.07122* (2015).
- [34] B. Hariharan, P. Arbeláez, R. Girshick, J. Malik, Hypercolumns for object segmentation and fine-grained localization, *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2015) 447–456.
- [35] J. Long, E. Shelhamer, T. Darrell, Fully convolutional networks for semantic segmentation, *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2015) 3431–3440.
- [36] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, A. L. Yuille, DeepLab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected CRFs, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 40 (2018) 834–848.
- [37] A. Torralba, A. Oliva, M. S. Castelano, J. M. Henderson, Contextual guidance of eye movements and attention in real-world scenes: The role of global features in object search, *Psychological Review* 113 (2006) 766.

- [38] L.-C. Chen, G. Papandreou, F. Schroff, H. Adam, Rethinking atrous convolution for semantic image segmentation, arXiv preprint arXiv:1706.05587 (2017).
- [39] J. Pan, C. C. Ferrer, K. McGuinness, N. E. O’Connor, J. Torres, E. Sayrol, X. Giro-i Nieto, SalGAN: Visual saliency prediction with generative adversarial networks, arXiv preprint arXiv:1701.01081 (2017).
- [40] A. Odena, V. Dumoulin, C. Olah, Deconvolution and checkerboard artifacts, *Distill* 1 (2016) e3.
- [41] X. Glorot, Y. Bengio, Understanding the difficulty of training deep feedforward neural networks, *Proceedings of the International Conference on Artificial Intelligence and Statistics* (2010) 249–256.
- [42] I. Sutskever, J. Martens, G. Dahl, G. Hinton, On the importance of initialization and momentum in deep learning, *Proceedings of the International Conference on Machine Learning* (2013) 1139–1147.
- [43] B. Zhou, A. Lapedriza, A. Khosla, A. Oliva, A. Torralba, Places: A 10 million image database for scene recognition, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 40 (2017) 1452–1464.
- [44] S. Jetley, N. Murray, E. Vig, End-to-end saliency mapping via probability distribution prediction, *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2016) 5753–5761.
- [45] D. P. Kingma, J. Ba, Adam: A method for stochastic optimization, arXiv preprint arXiv:1412.6980 (2014).
- [46] D. R. Wilson, T. R. Martinez, The general inefficiency of batch training for gradient descent learning, *Neural Networks* 16 (2003) 1429–1451.
- [47] A. Borji, L. Itti, CAT2000: A large scale fixation dataset for boosting saliency research, arXiv preprint arXiv:1505.03581 (2015).
- [48] C. Yang, L. Zhang, H. Lu, X. Ruan, M.-H. Yang, Saliency detection via graph-based manifold ranking, *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2013) 3166–3173.
- [49] Y. Li, X. Hou, C. Koch, J. M. Rehg, A. L. Yuille, The secrets of salient object segmentation, *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2014) 280–287.
- [50] J. Xu, M. Jiang, S. Wang, M. S. Kankanhalli, Q. Zhao, Predicting human gaze beyond pixels, *Journal of Vision* 14 (2014) 28.
- [51] P. Xu, K. A. Ehinger, Y. Zhang, A. Finkelstein, S. R. Kulkarni, J. Xiao, TurkerGaze: Crowdsourcing saliency with webcam based eye tracking, arXiv preprint arXiv:1504.06755 (2015).
- [52] M. Jiang, S. Huang, J. Duan, Q. Zhao, SALICON: Saliency in context, *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2015) 1072–1080.
- [53] H. R. Tavakoli, F. Ahmed, A. Borji, J. Laaksonen, Saliency revisited: Analysis of mouse movements versus fixations, *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2017) 6354–6362.
- [54] M. Kümmerer, T. Wallis, M. Bethge, Saliency benchmarking made easy: Separating models, maps and metrics, *Proceedings of the European Conference on Computer Vision* (2018) 770–787.
- [55] N. Riche, M. Duvinage, M. Mancas, B. Gosselin, T. Dutoit, Saliency and human fixations: State-of-the-art and study of comparison metrics, *Proceedings of the International Conference on Computer Vision* (2013) 1153–1160.
- [56] Z. Bylinskii, T. Judd, A. Oliva, A. Torralba, F. Durand, What do different evaluation metrics tell us about saliency models?, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 41 (2018) 740–757.
- [57] S. S. Kruthiventi, K. Ayush, R. V. Babu, DeepFix: A fully convolutional neural network for predicting human eye fixations, *IEEE Transactions on Image Processing* 26 (2017) 4446–4456.
- [58] S. Jia, EML-NET: An expandable multi-layer network for saliency prediction, arXiv preprint arXiv:1805.01047 (2018).
- [59] Z. Bylinskii, T. Judd, A. Borji, L. Itti, F. Durand, A. Oliva, A. Torralba, MIT saliency benchmark, <http://saliency.mit.edu/>, 2015.
- [60] T. Judd, F. Durand, A. Torralba, A benchmark of computational models of saliency to predict human fixations, 2012.
- [61] G. Huang, Z. Liu, L. Van Der Maaten, K. Q. Weinberger, Densely connected convolutional networks, *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2017) 2261–2269.
- [62] Y. Chen, J. Li, H. Xiao, X. Jin, S. Yan, J. Feng, Dual path networks, *Advances in Neural Information Processing Systems* 30 (2017) 4467–4475.
- [63] Z. Che, A. Borji, G. Zhai, X. Min, Invariance analysis of saliency models versus human gaze during scene free viewing, arXiv preprint arXiv:1810.04456 (2018).

Appendix A. Feature Concatenation Ablation Analysis

			$AUC-J \uparrow$	$SIM \uparrow$	$EMD \downarrow$	$AUC-B \uparrow$	$sAUC \uparrow$	$CC \uparrow$	$NSS \uparrow$	$KLD \downarrow$
MIT1003	\oplus CONCAT	μ	0.899	0.602	2.430	0.832	0.713	0.741	2.663	0.818
		σ	0.001	0.004	0.033	0.005	0.003	0.003	0.012	0.051
	\ominus CONCAT	μ	0.898	0.599	2.445	0.837	0.715	0.740	2.649	0.794
		σ	0.001	0.007	0.048	0.008	0.003	0.004	0.022	0.039
CAT2000	\oplus CONCAT	μ	0.882	0.734	2.553	0.812	0.582	0.854	2.359	0.430
		σ	0.000	0.002	0.025	0.005	0.003	0.003	0.007	0.010
	\ominus CONCAT	μ	0.881	0.734	2.545	0.814	0.586	0.855	2.354	0.436
		σ	0.000	0.003	0.032	0.003	0.003	0.002	0.012	0.013
DUT-OMRON	\oplus CONCAT	μ	0.921*	0.649*	1.155*	0.864	0.775	0.776*	2.873*	0.634*
		σ	0.001	0.002	0.016	0.005	0.003	0.002	0.008	0.023
	\ominus CONCAT	μ	0.919	0.641	1.191	0.861	0.773	0.766	2.825	0.656
		σ	0.001	0.003	0.018	0.005	0.003	0.003	0.009	0.042
PASCAL-S	\oplus CONCAT	μ	0.914*	0.667*	1.015*	0.853	0.701	0.818*	2.610*	0.645
		σ	0.000	0.003	0.010	0.004	0.003	0.002	0.008	0.044
	\ominus CONCAT	μ	0.907	0.636	1.130	0.855	0.711	0.791	2.494	0.649
		σ	0.001	0.006	0.019	0.004	0.005	0.005	0.017	0.031
OSIE	\oplus CONCAT	μ	0.918*	0.648*	1.647*	0.816	0.788	0.808*	3.010*	0.749
		σ	0.001	0.002	0.017	0.004	0.004	0.002	0.010	0.031
	\ominus CONCAT	μ	0.908	0.605	1.932	0.821	0.788	0.760	2.774	0.740
		σ	0.001	0.005	0.028	0.009	0.007	0.005	0.027	0.039

Table A.8: A summary of the quantitative results for the models with \oplus and without \ominus the concatenation of encoder features. The evaluation was carried out on five eye tracking datasets respectively. Each network was independently trained 10 times resulting in a distribution of values characterized by the mean μ and standard deviation σ . The star * denotes a significant increase of performance between the two conditions according to a one sided paired t-test. Arrows indicate whether the metrics assess similarity \uparrow or dissimilarity \downarrow between predictions and targets. The best results are marked in bold.

	$AUC-J \uparrow$	$SIM \uparrow$	$EMD \downarrow$	$AUC-B \uparrow$	$sAUC \uparrow$	$CC \uparrow$	$NSS \uparrow$	$KLD \downarrow$
Action	+0.001	+0.007	-0.062	+0.001	\pm 0.000	+0.010	+0.025	-0.020
Social	+0.004	+0.003	-0.064	+0.002	+0.002	+0.007	+0.025	-0.037
Fractal	+0.001	-0.001	+0.034	-0.017	-0.004	\pm 0.000	+0.018	+0.018
Pattern	\pm 0.000	+0.006	-0.051	-0.005	-0.004	\pm 0.000	-0.005	+0.016

Table A.9: A list of the image categories from the CAT2000 validation set that either showed the largest average improvement (first two entries) or impairment (last two entries) by the multi-level design based on the cumulative rank across a subset of weakly correlated evaluation measures. Arrows indicate whether the metrics assess similarity \uparrow or dissimilarity \downarrow between predictions and targets. Results that improved on the respective metric are marked in green, whereas results that impaired performance are marked in red.

In this experimental setting, we removed the concatenation operation from the network architecture and compared the model performance of the ablated version to the one including a multi-level design (see Table A.8). While models trained on the CAT2000 dataset did not consistently benefit from the aggregation of features at different stages of the encoder, all other cases demonstrated a mostly significant improvement according to the majority of metric scores. Table A.9 indicates that predictions on natural image categories (*Action*, *Social*) leveraged the multi-level information for better performance, whereas adverse results were achieved on artificial and simplified stimuli (*Fractal*, *Pattern*). In conclusion, the feature concatenation design

might only be recommendable for training models on datasets that mostly consist of complex natural images, such as MIT1003, DUT-OMRON, PASCAL-S, or OSIE.