

TempSAL - Uncovering Temporal Information for Deep Saliency Prediction

Bahar Aydemir, Ludo Hoffstetter, Tong Zhang, Mathieu Salzmann, Sabine Süsstrunk
School of Computer and Communication Sciences, EPFL, Switzerland

{bahar.aydemir, tong.zhang, mathieu.salzmann, sabine.susstrunk}@epfl.ch

Abstract

Deep saliency prediction algorithms complement the object recognition features, they typically rely on additional information, such as scene context, semantic relationships, gaze direction, and object dissimilarity. However, none of these models consider the temporal nature of gaze shifts during image observation. We introduce a novel saliency prediction model that learns to output saliency maps in sequential time intervals by exploiting human temporal attention patterns. Our approach locally modulates the saliency predictions by combining the learned temporal maps. Our experiments show that our method outperforms the state-of-the-art models, including a multi-duration saliency model, on the SALICON benchmark. Our code will be publicly available on GitHub¹.

1. Introduction

Humans have developed attention mechanisms that allow them to selectively focus on the important parts of a scene. Saliency prediction algorithms aim to computationally detect these regions that stand out relative to their surroundings. These predictions have numerous applications in image compression [34], image enhancement [48], image retargeting [1], rendering [40], and segmentation [26].

Since the seminal work of Itti et al. [16], many have developed solutions using both handcrafted features [5] and deep ones [7, 15, 24, 31, 43, 45]. Nowadays, employing deep neural networks is preferred in saliency prediction as they outperform bottom-up models. These methods typically depend on pretrained object recognition networks to extract features from the input image [28]. In addition to these features, scene context [44], object co-occurrence [47], and dissimilarity [2] have been exploited to improve the saliency prediction. However, while these approaches model the scene context and objects, they fail to consider that humans dynamically observe scenes [46]. In neuroscience, the inhibition of return paradigm states that a suppression mecha-

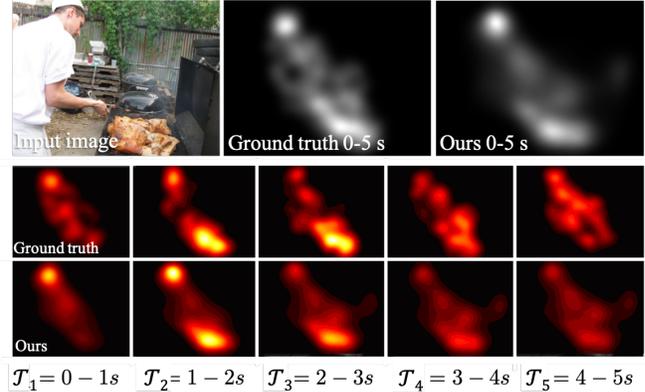


Figure 1. An example of how human attention shifts over time. We show the input image and the corresponding image saliency ground truth from the SALICON [18] dataset. Notice that in \mathcal{T}_1 , the cook is salient, while in \mathcal{T}_2 and \mathcal{T}_3 , the food on the barbeque becomes the most salient region in this scene. We can predict saliency maps in these sequential time intervals as well as combine them into a refined single image saliency map for the whole observation duration.

nism reduces visual attention towards recently attended objects [36] and encourages selective attention to novel regions. Motivated by this principle, we develop a saliency prediction model that incorporates temporal information.

Fosco et al. [12] also exploit temporal information in saliency prediction, but they consider snapshots containing observations up to 0.5, 3, and 5 seconds, thus not leveraging saliency trajectory but rather saliency accumulation. By contrast, here, we model consecutive time slices, connecting our approach more directly with the human gaze and thus opening the door to automated visual appeal assessment in applications such as website design [29], advertisement [33] and infographics [11].

To achieve this, we show that when viewing images, human attention yields temporally evolving patterns, and introduce a network capable of exploiting this temporal information for saliency prediction. Specifically, our model learns time-specific predictions and is able to combine them with a conventional image saliency map to obtain a tempo-

¹<https://baharay.github.io/tempsal/>

rally modulated image saliency prediction. As evidenced by our experiments, this consistently boosts the accuracy of the baseline network, enabling us to outperform the state-of-the-art models on the SALICON saliency benchmark. In particular, we outperform image saliency prediction models by 2.9% and 1.8% in the standard IG and KLD metrics, respectively, and the multi-duration model of [12] by 35.7% and 68% in IG and KLD, respectively.

We summarize our contributions as follows:

- We evidence the presence of temporally evolving patterns in human attention.
- We show that temporal information in the form of a saliency trajectory is important for saliency prediction in natural images, providing an investigation of the SALICON dataset for temporal attention shifts.
- We introduce a novel, saliency prediction model, namely TempSAL, capable of simultaneously predicting conventional image saliency and temporal saliency trajectories.
- We propose a spatiotemporal mixing module that learns time dependent patterns from temporal saliency maps. Our approach outperforms the state-of-the-art image saliency models that either do not consider temporal information or encode it in a cumulative manner.

2. Related work

2.1. Saliency prediction for natural images

Early saliency prediction methods were biologically inspired and bottom up. In particular, Itti et al. used color, intensity, and orientation contrast [16]. Goferman employed global and local contrast as contextual cues [13]. Judd et al. [21] further incorporated mid-level and high-level semantic features, using horizon, face, person, and car detectors. Later, Vig et al. [43] showed that deep neural networks can be applied to saliency prediction. Yet, saliency prediction lacks the large scale annotated datasets that are available for image classification tasks [9], which prevents training robust models. To overcome this, Kummerer et al. [23] showed that using pretrained object recognition networks significantly improves saliency predictions. Subsequent state-of-the-art models such as EML-Net [17], DeepGaze2 [24], and SALICON [15] similarly use pretrained convolutional neural network (VGG [39]) encoders. Recent works utilize additional sources of information such as scene context [22], external knowledge [47] and object dissimilarity [2] to improve saliency prediction. Yet, none of these methods take into account the temporal evolution of human gaze, which occurs even when the image stimuli are static. In our work, we make use of these temporal patterns as an

additional source of information to boost conventional image saliency prediction.

2.2. Multi-duration saliency

Existing image saliency ground-truth maps include all fixations made throughout the observation period. Aggregating all these fixations that have different timestamps into a single ground-truth map results in the loss of temporal information. Representing the fixations as scanpaths retains the temporal clues by encoding the change of gaze of an individual over time. However, merging numerous scanpaths is challenging [14]. Fosco et al. [12] proposed *multi-duration saliency* to characterize the attention of a group of individuals while taking into account time-dependent attention shifts. The temporal maps they rely on, however, encode the attention distribution of many observers across overlapping time periods of increasing durations. While this is a convenient way of capturing a population’s attention patterns, it does not reflect the saliency trajectory over time. Similarly, [32] use order of fixations as a sequential metadata for deep supervision but they do not model evolution of attention through time. In our work, we model multi-duration saliency to analyze underlying attention patterns by using *mutually exclusive* time slices. We provide this temporal information to our spatiotemporal mixing module to refine the initial image saliency prediction with temporal information. Moreover, this lets us predict temporal saliency maps for each second of attention.

3. Temporal saliency data analysis

3.1. Dataset

SALICON [18] is the largest human attention dataset on natural images. It was created via a crowdsourced mouse tracking experiment, which was shown to be similar to eye-tracking [18] and widely used in the saliency prediction literature. SALICON consists of 10000 training, 5000 validation and 5000 test images from the MS-COCO dataset [27]. The SALICON dataset provides saliency maps, fixations, and gaze points for each image and observer. A gaze point is a raw data point recorded by a tracking device. It describes the spatial coordinates of the eye/mouse on the associated stimuli at a given timestamp. Conversely, fixations describe the coordinates of the long pause when the eyes are fixated on an image detail. Following common practice in eye tracking experiments, Jiang et al. [18] grouped spatially and temporally close gaze points to create fixations. Since the fixations were created by grouping multiple gaze points, they do not have an associated timestamp. To address this, SALICON-MD [12] assumes that the fixations are uniformly distributed across the total viewing time. We provide a finer approximation for recovering the fixations’

timestamps, by minimizing the spatial and temporal distance between a fixation and the nearest gaze point. We refer the reader to the supplementary material for the details of this approximation process.

3.2. Temporal patterns in the dataset

In this section, we examine how temporality evolves during human visual attention. To observe the evolution of attention over time, we slice the data into five slices, one for each second of observation. In particular, we inspect the dissimilarity between slices, the agreement with the average, and the distribution of fixations in the time-saliency space.

Average maps. Viewing patterns in image saliency experiments tend to show a concentration towards the image center [41], which is known as the photographer’s bias or center bias. We observe a similar spatial bias for each temporal slice, shown in Figure 2, where we plot the average heat maps for each temporal slice. Note that the gaze tends to converge to the center of the image as time passes. This means that the observers revisit the previously seen important center regions [46].

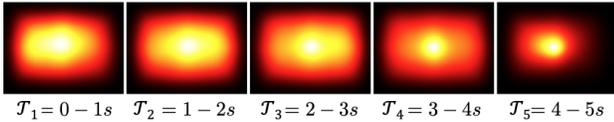


Figure 2. Average heat maps for each one second interval. Note that a center-bias occurs, similar to image saliency prediction’s average ground-truth maps.

We plot the differences of the consecutive average temporal slices in Figure 3 to illustrate attention shifts. Light blue indicates the regions with reduced attention, whereas (light) red indicates increased attention. We observe that attention shifts from left to right, with a subsequent dispersion from the center towards the corners. Then, attention increases at the center of the image, slightly skewed to the left. Interestingly, the trend (especially in $T_2 - T_1$) coincides with the western left-to-right reading direction [6].

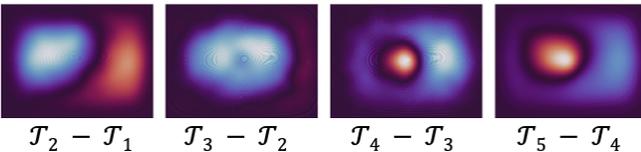


Figure 3. Differences of the consecutive average temporal slices shown in Fig. 2. Red indicates regions of increased attention whereas blue indicates decreased attention.

CC	\mathcal{T}_1	\mathcal{T}_2	\mathcal{T}_3	\mathcal{T}_4	\mathcal{T}_5
\mathcal{T}_1	1.00	<u>0.70</u>	0.54	0.50	0.54
\mathcal{T}_2	0.70	1.00	<u>0.73</u>	0.66	0.65
\mathcal{T}_3	0.54	0.73	1.00	<u>0.75</u>	0.70
\mathcal{T}_4	0.50	0.66	<u>0.75</u>	1.00	0.73
\mathcal{T}_5	0.54	0.65	0.70	<u>0.73</u>	1.00
Average	0.66	0.75	0.74	0.73	0.72

Table 1. Correlation scores of the temporal slices with each other in a single image, averaged over all images. All slices show more similarity to their direct temporal neighbors. The last row shows the average similarity of a slice with the other slices, \mathcal{T}_1 being the most dissimilar one.

Inter-slice similarity across time. We expect temporal saliency slices to be more similar to their closer-in-time slices than to the ones further away since human attention is continuous over time. Table 1 contains the correlation coefficients between each pair of saliency slices in a single image, averaged over all images. We calculate the correlation coefficient between slices \mathcal{T}_j and \mathcal{T}_k as

$$CC(\mathcal{T}_j, \mathcal{T}_k) = \frac{1}{N} \sum_{n=i}^N CC(\mathcal{T}_{ij}, \mathcal{T}_{ik}), \quad j, k \in \{1, \dots, 5\}, \quad (1)$$

where N is the total number of images, and \mathcal{T}_{ij} and \mathcal{T}_{ik} denote the j^{th} and k^{th} slice of the i^{th} image.

By calculating t-test scores on the pairwise comparisons, we observe that all of the pairwise differences except $\mathcal{T}_1, \mathcal{T}_3$ and $\mathcal{T}_1, \mathcal{T}_5$ are statistically significant ($p < 0.01$). Thus, the attention residuals between different time intervals in one image are significantly different. We provide more details in the supplementary material.

Intra-slice similarity across images. We also investigate the deviation of each slice from its respective average slices. The average slices are depicted in Figure 2. Table 2 shows the deviation of a slice from the average time slices per image. We compute CC scores between a single slice and the corresponding average slice as

$$CC(\mathcal{T}_j, A_j) = \frac{1}{N} \sum_{n=i}^N CC(\mathcal{T}_{ij}, A_j), \quad j \in \{1, \dots, 5\}, \quad (2)$$

where A_j denotes the j^{th} average slice.

We average the scores across all images. Higher values of CC indicate more agreement with the average whereas lower values of CC indicate more deviation from the average. Note that the similarity with the average across images decreases with time, except for the last slice. This can be explained by the more prominent center bias in $\mathcal{T}_5 - \mathcal{T}_4$ as seen in Figure 2. \mathcal{T}_1 has the least deviation from the average by a significant margin ($p \ll 0.01$). This shows

	\mathcal{T}_1	\mathcal{T}_2	\mathcal{T}_3	\mathcal{T}_4	\mathcal{T}_5
CC	0.574	0.433	0.431	0.426	0.447

Table 2. Correlation scores of each time slice in a single image with the average maps presented in Figure 2, averaged over all images. The similarity between slices across images decreases with time, with the exception of the last slice.

that humans tend to look at similar places at first, and then their attention scatters around to less important regions.

Early and late fixations versus saliency. Lastly, we investigate the relationship between the fixation timestamps and their respective saliency values. We assign a saliency value to each fixation as the normalized pixel value in the corresponding saliency map. We plot the histogram of number of fixations with their saliency and timestamp values in Figure 4. The fixation time stamps range from 0 to 5000 ms and the saliency values range from 0 to 1. Late fixations tend to have lower saliency values than earlier fixations, as indicated by the darker color towards the bottom right corner. That is, the first region we glance at in an image is more important (salient) than the following regions [5, 16].

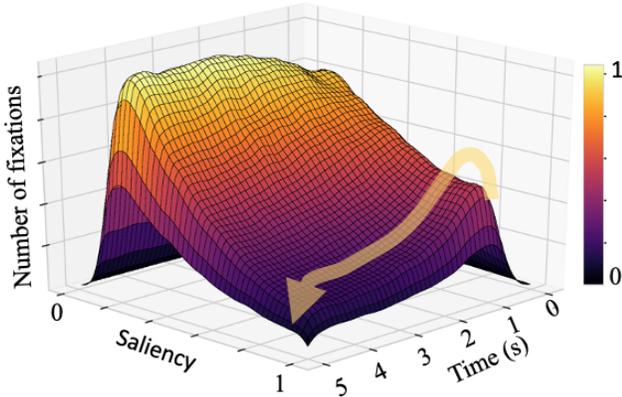


Figure 4. Number of fixations with their respective saliency values and timestamps. Lighter colors indicate higher number of occurrences while darker areas denote fewer occurrences. We see that late fixations tend to be less salient, which can be seen as the decrease in the number of salient fixations along the arrow. The most salient fixations appear at approximately 1s.

4. Methodology

4.1. Temporal slices

We aim to recover fixation timestamps to train models with this temporal information. We extract temporal slices by grouping the fixations in several time-intervals and, following common practice [3], blurring with a Gaussian kernel. We break down the fixations into time slices with two

time-slicing (grouping) alternatives, namely equal duration and equal distribution. The equal duration model outperforms the equal distribution one and is easier to interpret; we present a comparison of these two alternatives as an ablation study in Section 5.7.

4.2. Temporal saliency model

Let us now introduce our framework that exploits temporal human attention information. Our model is depicted in Figure 5. We extract image features using a pre-trained object recognition encoder [30]. Then, we decode these features by a temporal slice decoder to obtain one saliency map per time slice. These temporal saliency slices are useful in automated visual appeal assessment in applications such as website design [29], advertisement [33] and infographics [11]. In parallel, we decode the same image features into an initial image saliency prediction. Finally, we combine the temporal slices and the image saliency predictions in the spatiotemporal mixing module to produce a final image saliency map. We describe each component in detail in the following sections.

4.2.1 Image encoder and saliency decoders

Following the previous saliency prediction architectures [24, 28, 37], we first encode the input image with a pre-trained image classification network, in our case PNASNet-5 [30]. We extract encoded features at various levels for multi-level integration, similar to a U-Net structure [38]. Formally, we denote the image encoder as

$$\mathcal{E}(\mathcal{I}) = [\mathcal{E}_i], \quad i \in \{1, \dots, 5\}, \quad (3)$$

where \mathcal{I} is the input image, and \mathcal{E}_i the i^{th} encoder block. The output of $\mathcal{E}(\cdot)$ therefore is a 5D vector. The early encoder blocks extract low-level features, such as edges, color, and contrast, while the later blocks encode high-level semantics. We pass these blocks to the temporal saliency decoder, the image saliency decoder, and the spatiotemporal mixing module.

Our temporal slice decoder, namely \mathcal{D}_T , processes the encoder blocks with four 3x3 convolution layers followed by ReLU functions, integrating one encoder block after each convolution. Later, two 3x3 convolution layers with a ReLU in-between and a sigmoid function at the end produce n temporal saliency maps. Formally, we write the temporal saliency decoder as

$$\mathcal{D}_T(\mathcal{E}(\mathcal{I})) = [\mathcal{T}_n] =: \mathcal{T}, \quad n \in \{1, \dots, 5\} \quad (4)$$

where \mathcal{T}_n denotes the n^{th} temporal saliency slice. Through this branch of the network, a single image input produces n temporal saliency slices. We use this component to provide temporal predictions to the spatiotemporal mixing module.

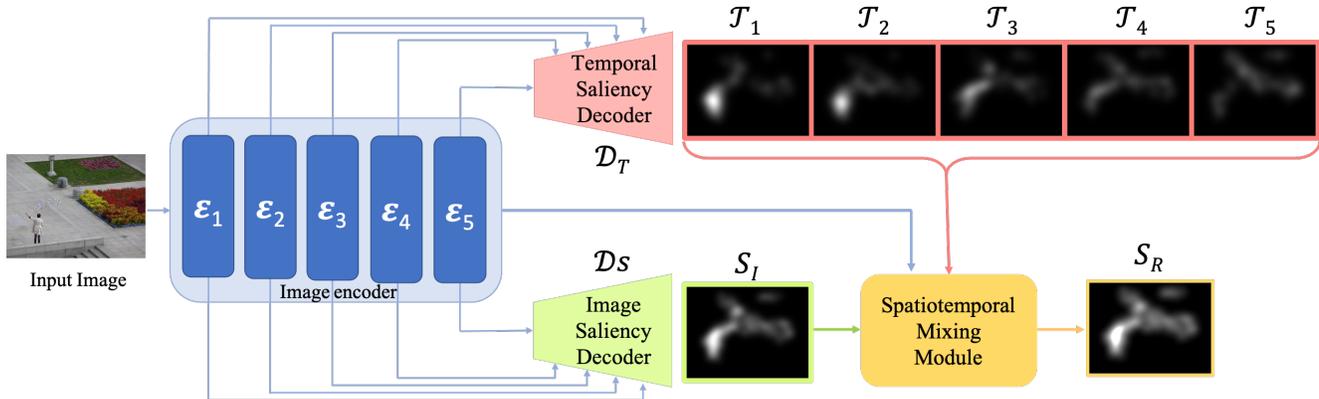


Figure 5. **Overview of the proposed architecture.** We encode image features into encoder blocks consisting of multi-level image features. We then pass these blocks to the temporal saliency decoder (shown in pink) to decode them into temporal saliency predictions, which are saliency maps in sequential time intervals. In parallel, the image saliency decoder (shown in green) decodes the encoder blocks into an image saliency prediction. We then combine (1) the temporal saliency maps, (2) the image saliency map, and (3) the encoder blocks in the spatiotemporal mixing module (shown in orange). (Best viewed in color.)

Our image saliency decoder, namely \mathcal{D}_S , has the same structure as \mathcal{D}_T , with the exception of the number of output channels. This component produces a single map \mathcal{S}_I , which corresponds to the conventional image saliency map of the input image. As such, we can write

$$\mathcal{S}_I = \mathcal{D}_S(\mathcal{E}(\mathcal{I})). \quad (5)$$

We use this module to provide cumulative saliency information to the spatiotemporal mixing module.

4.2.2 Spatiotemporal Mixing Module

To incorporate temporal information into the saliency prediction, we introduce a module that combines temporal and spatial saliency maps. This module takes temporal saliency predictions, the initial image saliency prediction, and the encoded image feature blocks as input. We write this as

$$\mathcal{S}_R = \text{SMM}(\mathcal{E}(\mathcal{I}), \mathcal{T}, \mathcal{S}_I), \quad n \in \{1, \dots, 5\}, \quad (6)$$

where \mathcal{S}_R denotes the final, refined image saliency map. We use the encoded image features in this module to benefit from both low-level and high-level saliency features by multi-level integration.

The module architecture is shown in Figure 6. It takes the last two encoder blocks $[\mathcal{E}_5, \mathcal{E}_4]$ and passes them through a 3×3 convolution. We then concatenate the other encoder blocks with the image saliency and temporal saliency maps passing through 3×3 convolution, ReLU, and linear upsampling to keep the spatial dimensions consistent. In the last step, we only add the saliency maps to output a final refined saliency map \mathcal{S}_R . This module eliminates the need for optimizing a weight parameter between the spatial and temporal maps. It can also modulate the maps within the

spatial range of convolutions, which allows the selection of different regions from different maps.

4.2.3 Loss Functions

To train our network, we use the Kullback-Leibler divergence (KL) [42] and the Correlation Coefficient (CC) [19] between the predicted and ground-truth saliency maps. First, we train the temporal branch using

$$\mathcal{L}_1(\mathcal{I}) = \lambda_1 * \text{CC}(GT_n, \mathcal{T}_n) + \beta_1 * \text{KL}(GT_n, \mathcal{T}_n), \quad (7)$$

where GT_n denotes the temporal ground truth for the n^{th} slice. We then freeze the weights in this component and train the spatiotemporal mixing module using

$$\mathcal{L}_2(\mathcal{I}) = \lambda_2 * \text{CC}(GT, \mathcal{S}_R) + \beta_2 * \text{KL}(GT, \mathcal{S}_R), \quad (8)$$

where GT is the image saliency ground truth for image I .

5. Experiments and Results

5.1. Experimental Setup

We use a batch size of 32 and an initial learning rate of $1e-4$, reduced by a factor of ten every two epochs. We train the temporal branch first and then freeze the weights. We found that 10 epochs of training on SALICON was sufficient. For SALICON, we used the provided test, train, and validation splits.

5.2. Metrics

We evaluate the obtained saliency predictions according to the following standard metrics used by the community.

Area Under the Curve (AUC) [3]: Saliency prediction can

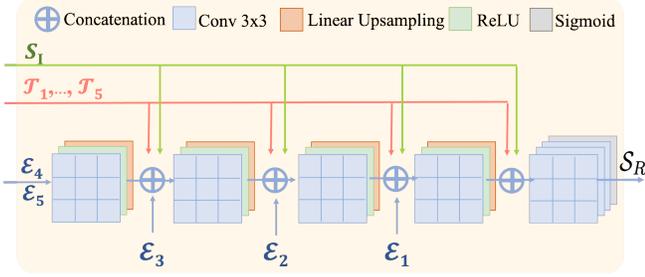


Figure 6. The spatiotemporal mixing module combines temporal saliency predictions with the conventional image saliency prediction with multi-level image feature integration. \mathcal{S}_T denotes the predicted image saliency map, $\mathcal{T}_{1,\dots,5}$ the temporal saliency predictions for n frames, and $\mathcal{E}_{1,\dots,5}$ the encoder blocks. This multi-level integration scheme provides information from earlier layers of the network to the next blocks in this module. \mathcal{S}_R denotes the temporally refined image saliency map output.

be interpreted as classifying fixation vs non-fixation points. The area under the ROC curve shows the trade-off between true positives (TP) and false positives (FP). A higher AUC score indicates less FPs.

Normalized Scanpath Saliency (NSS) [35]: This metric compares the predicted saliency values at the ground-truth fixation points to the average predicted saliency. An NSS score of one indicates that the predicted saliency values at the ground-truth fixation points are one standard deviation above the average.

Kullback - Leibler Divergence (KL) [42]: The KL measures the cumulative distance between the predicted and the ground-truth saliency maps. A KL score close to zero indicates a better approximation of the ground-truth saliency map by the predicted one.

Pearson’s correlation coefficient (CC) [19]: This metric measures the linear relationship between the predicted and ground-truth saliency maps. It ranges from -1 to 1. A CC score close to one indicates a strong linear correlation between the two maps.

Similarity (SIM) score [20]: The similarity score sums the minimum value between the predicted and the ground-truth saliency maps over all pixels. A similarity score of 1 indicates a perfect prediction since both of the maps are probability distributions summing to 1.

Information Gain (IG) score [25]: The information gain is a information-theoretic metric which measures the difference in average log-likelihood between the predicted saliency map and center-bias prior.

5.3. Quantitative Results

We compare our method with the state-of-the-art models, namely SAM-Resnet [8], MSI-Net, GazeGAN, MDNSal [37], SimpleNet [37], DeepGaze IIE [28], and MD-SEM [12], in image saliency prediction. Our model

outperforms these methods in five out of seven metrics, showing the benefit of incorporating temporal information. Moreover, our model outperforms the only other multi-duration saliency model in image saliency prediction by a significant margin. Furthermore, we compare our model with this multi-duration model and a multi-duration baseline. Our model improves the saliency prediction in two durations consisting of 0.5 and 3 seconds in two out of three metrics and in all three metrics in the five second duration.

5.4. Comparison with state-of-the-art methods

We first evaluate the performance of our model on image saliency prediction on the SALICON benchmark [18]. The ground truth of SALICON’s test set is exclusively hosted on the CodaLab website². Table 3 shows the comparison of standard evaluation metrics for different state-of-the-art saliency models alongside our model TempSAL. TempSAL outperforms all the baselines in almost all metrics. When it does not, it still yields competitive results.

5.5. Comparison with the multi-duration method

To compare our method with the only other multi-duration model [12], we modify our network to output three temporal slices. We train our network on a three slice SALICON multi-duration dataset first and then fine-tune it on the CodeCharts1k dataset [12] using the given training and validation splits. We report the results of the comparison in Table 4.

5.6. Qualitative Results

In Figure 7, we compare the temporal and image saliency maps obtained with our method with the ground truth from SALICON [18]. Our model learns time-specific predictions and is able to combine such predictions with a conventional image saliency map. We provide additional qualitative results in the supplementary material.

5.7. Ablation studies

In this section, we investigate the effect of different components in our model, and of two temporal slicing methods. We also provide a comparison with a multi-duration baseline model on the temporal SALICON dataset.

Effect of the SMM module: We evaluate the effect of the spatiotemporal mixing module (SMM) and the image saliency decoder in Table 5. The first model consists of the image encoder and temporal saliency decoder only. We take the average of the temporal slices to measure its performance by comparing with the ground-truth image saliency map. In the second row, we add the image saliency decoder to our model. Similarly, we take the average of the predicted maps \mathcal{T}_n and \mathcal{S}_T . Lastly, we add the spatiotemporal mixing

²<https://competitions.codalab.org/competitions/17136>

Model	MD	AUC \uparrow	CC \uparrow	KL \downarrow	SAUC \uparrow	IG \uparrow	NSS \uparrow	SIM \uparrow
SAM-Resnet [8]	\times	0.865	0.899	0.610	0.741	0.538	1.990	0.793
MSI-Net [22]	\times	0.865	0.899	0.307	0.736	0.793	1.931	0.784
GazeGAN [4]	\times	0.864	0.879	0.376	0.736	0.720	1.899	0.773
SimpleNet [37]	\times	0.869	0.907	0.201	0.743	0.880	1.960	0.793
MDNSal [37]	\times	0.865	0.899	0.221	0.736	0.863	1.935	0.790
UNISAL [10]	\times	0.864	0.879	0.354	0.739	0.780	1.952	0.775
DeepGaze IIE [28]	\times	0.869	0.872	0.285	0.767	0.766	1.996	0.733
MD-SEM [12]	\checkmark	0.864	0.868	0.568	0.746	0.660	2.058	0.774
TempSAL	\checkmark	0.869	0.911	0.195	0.745	0.896	1.967	0.800

Table 3. Evaluation results on the SALICON (LSUN 2017) test benchmark. We compare our model with the state-of-the-art saliency prediction models, namely SAM-Resnet [8], MSI-Net [22], GazeGAN [4], MDNSal [37], SimpleNet [37], DeepGaze IIE [28], and MD-SEM [12]. The results in bold show the best performance. Our method outperforms the state-of-the-art on conventional image saliency in five metrics. The MD column denotes the ability of the models to predict multi-duration saliency. Our model outperforms the only other multi-duration saliency model by a significant margin on six out of seven metrics.

Model	TempSAL			MD-SEM			SAM-MD		
	CC \uparrow	KL \downarrow	NSS \uparrow	CC \uparrow	KL \downarrow	NSS \uparrow	CC \uparrow	KL \downarrow	NSS \uparrow
Accuracy metrics									
Slice 1 (0-500 ms)	0.819	0.496	3.422	0.816	0.351	3.374	0.805	0.370	3.181
Slice 2 (0-3000 ms)	0.752	0.512	2.703	0.745	0.452	2.694	0.738	0.469	2.541
Slice 3 (0-5000 ms)	0.822	0.471	3.337	0.734	0.487	2.677	0.715	0.535	2.495
Average	0.797	0.493	3.154	0.765	0.430	2.915	0.753	0.458	2.739

Table 4. Results of our model, MD-SEM [12], and the baseline SAM-MD [12] across different durations on the CodeCharts1k dataset [12]. Our model improves the saliency prediction in the first two intervals, consisting of 500 ms and 3000 ms observations, in two out of three metrics. Our model benefits from both non-overlapping temporal slices and image saliency. In this comparison, the time slices are cumulative, not mutually exclusive, which diminishes the separation between different time slices. However, our model performs well in the last slice since this slice corresponds to the image saliency with 5000 ms observation duration.

Models	CC \uparrow	KL \downarrow	NSS \uparrow	SIM \uparrow
$\mathcal{D}_T(\mathcal{E}(\mathcal{I}))$	0.852	0.243	1.973	0.754
$+\mathcal{D}_S(\mathcal{E}(\mathcal{I}))$	0.857	0.252	1.943	0.760
+SMM	0.906	0.198	1.930	0.798

Table 5. Results of ablation studies on the temporal SALICON validation dataset. The first row denotes the model with only the temporal saliency decoder. In the second row, the model has both the temporal and image saliency decoders. The last row denotes the performance with the spatiotemporal mixing module (SMM). As evidenced by the improved accuracy metrics, the SMM effectively modulates the spatial and temporal saliency maps to refine the initial image saliency prediction.

module, which effectively modulates these predicted maps and combines them into a final image saliency map \mathcal{S}_R .

Comparison with a temporal baseline model: The performance of our TempSAL model with five temporal slices is provided in Table 6. Note that each saliency slice contains five times fewer samples than the original image saliency map. Therefore, individual slices contain more variation compared to conventional accumulated maps. As a baseline to our model, we compute the performance of an architec-

ture composed of five replicated SimpleNet models (5xSimpleNet) [37], each trained on one saliency slice. This baseline model uses an unshared encoder and decoder for each slice, while we share the decoder among slices. Therefore, we do not benefit from increased model capacity. We observe an accuracy decline in the baseline model, which confirms the increased discrepancy in the data.

Model	Baseline				TempSAL			
	CC \uparrow	KL \downarrow	NSS \uparrow	SIM \uparrow	CC \uparrow	KL \downarrow	NSS \uparrow	SIM \uparrow
Accuracy metrics								
\mathcal{T}_1	0.898	0.211	2.436	0.778	0.899	0.214	2.453	0.782
\mathcal{T}_2	0.870	0.219	2.159	0.765	0.877	0.215	2.211	0.776
\mathcal{T}_3	0.840	0.247	1.840	0.753	0.843	0.247	1.878	0.758
\mathcal{T}_4	0.820	0.273	1.729	0.743	0.825	0.264	1.740	0.749
\mathcal{T}_5	0.811	0.275	1.646	0.738	0.813	0.276	1.654	0.743
Average	0.848	0.245	1.962	0.756	0.852	0.243	1.987	0.761

Table 6. Results of the baseline model (left) and our TempSAL model (right) across different time slices. In 18 out of 20 comparisons, our model consistently outperforms the baseline. Note that both models perform best in the first slice, in which the intra-slice agreement is more prominent than in the other slices, as mentioned in Section 3.2.

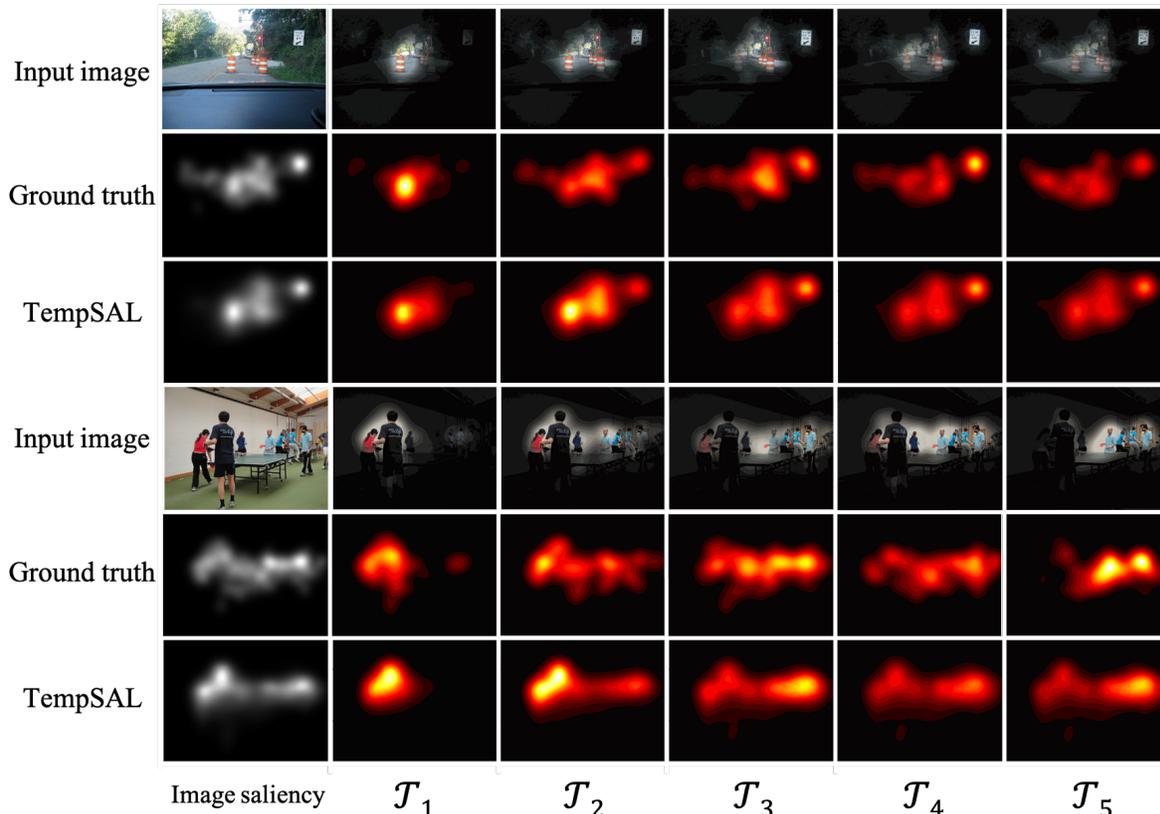


Figure 7. Temporal saliency predictions of our model and ground truth temporal saliency maps. Black and white maps are image saliency maps for the whole observation duration. Red-yellow maps are temporal saliency maps for one second intervals. The first row shows the input image with the ground-truth saliency overlaid. The second row shows the ground truth saliency and the third row shows our temporal saliency predictions. Our approach captures the attention shifts in sequential temporal maps. Moreover, our model is able to produce accurate image saliency predictions which are close to the ground truth maps.

Equal duration versus equal distribution: We break down the fixations into time slices with two time-slicing alternatives, namely equal duration and equal distribution. The equal duration method groups the fixations based on their timestamps. Each slice has a different total number of fixations. On the other hand, the equal distribution method groups an equal number of fixations in each slice. Therefore, the duration of each slice is different from that of the other ones. We provide more details on the slicing methods in the supplementary material. We train and evaluate two models using both sampling methods. The results are presented in Table 7.

6. Conclusion

We present a saliency prediction method that can learn time-specific predictions and is also able to exploit temporal information to improve overall image saliency prediction. In particular, we show that the temporally evolving patterns in human attention play an important

Model	Equal distribution				Equal duration			
	CC \uparrow	KL \downarrow	NSS \uparrow	SIM \uparrow	CC \uparrow	KL \downarrow	NSS \uparrow	SIM \uparrow
Slice 1	0.899	0.213	2.426	0.785	0.899	0.214	2.453	0.782
Slice 2	0.857	0.252	2.152	0.760	0.877	0.215	2.211	0.776
Slice 3	0.836	0.263	1.877	0.752	0.843	0.247	1.878	0.758
Slice 4	0.821	0.278	1.750	0.746	0.825	0.264	1.740	0.749
Slice 5	0.815	0.283	1.676	0.744	0.813	0.276	1.654	0.743
Average	0.846	0.258	1.976	0.757	0.852	0.243	1.987	0.761

Table 7. Results of the equal distribution model (first column) and the equal duration one (second column) across different time slices. The equal duration model achieves better results in 13 out of 20 comparisons.

role in saliency prediction in natural images. This is evidenced by our experiments that demonstrate our TempSAL method outperforming the state-of-the-art, including a multi-duration method exploiting cumulative temporal saliency maps.

Acknowledgement. This work was supported in part by the Swiss National Science Foundation via the Sinergia grant CRSII5-180359.

References

- [1] Radhakrishna Achanta and Sabine Süssstrunk. Saliency detection for content-aware image resizing. In *IEEE International Conference on Image Processing (ICIP)*, pages 1005–1008. IEEE, 2009. 1
- [2] Bahar Aydemir, Deblina Bhattacharjee, Seungryong Kim, Tong Zhang, Mathieu Salzmann, and Sabine Süssstrunk. Modeling object dissimilarity for deep saliency prediction. *CoRR*, abs/2104.03864, 2021. 1, 2
- [3] Zoya Bylinskii, Tilke Judd, Aude Oliva, Antonio Torralba, and Fredo Durand. What do different evaluation metrics tell us about saliency models? *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(3):740, 2019. 4, 5
- [4] Zhaohui Che, Ali Borji, Guangtao Zhai, Xiongkuo Min, Guodong Guo, and Patrick Le Callet. Gazegan: A generative adversarial saliency model based on invariance analysis of human gaze during scene free viewing. *ArXiv*, abs/1905.06803, 2019. 7
- [5] Ming-Ming Cheng, Niloy J. Mitra, Xiaolei Huang, Philip H. S. Torr, and Shi-Min Hu. Global contrast based salient region detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(3):569–582, 2015. 1, 4
- [6] Hannah Faye Chua, Julie E. Boland, and Richard E. Nisbett. Cultural variation in eye movements during scene perception. *Proceedings of the National Academy of Sciences*, 102(35):12629–12633, 2005. 3
- [7] Marcella Cornia, Lorenzo Baraldi, Giuseppe Serra, and Rita Cucchiara. A deep multi-level network for saliency prediction. In *IEEE International Conference on Pattern Recognition (ICPR)*, pages 3488–3493. IEEE, 2016. 1
- [8] Marcella Cornia, Lorenzo Baraldi, Giuseppe Serra, and Rita Cucchiara. Predicting human eye fixations via an LSTM-based saliency attentive model. *IEEE Transactions on Image Processing*, 27(10):5142–5154, 2018. 6, 7
- [9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 2
- [10] Richard Droste, Jianbo Jiao, and J. Alison Noble. Unified Image and Video Saliency Modeling. In *European Conference on Computer Vision (ECCV)*, 2020. 7
- [11] Camilo Fosco, Vincent Casser, Amish Kumar Bedi, Peter O’Donovan, Aaron Hertzmann, and Zoya Bylinskii. Predicting visual importance across graphic design types. *CoRR*, abs/2008.02912, 2020. 1, 4
- [12] Camilo Fosco, Anelise Newman, Pat Sukhum, Yun Bin Zhang, Nanxuan Zhao, Aude Oliva, and Zoya Bylinskii. How much time do you have? modeling multi-duration saliency. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4473–4482, 2020. 1, 2, 6, 7
- [13] Stas Goferman, Lihi Zelnik-Manor, and Ayellet Tal. Context-aware saliency detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(10):1915–1926, 2011. 2
- [14] Joseph Goldberg and Jonathan Helfman. Scanpath clustering and aggregation. pages 227–234, 01 2010. 2
- [15] Xun Huang, Chengyao Shen, Xavier Boix, and Qi Zhao. SALICON: Reducing the semantic gap in saliency prediction by adapting deep neural networks. In *IEEE International Conference on Computer Vision (ICCV)*, pages 262–270. IEEE, 2015. 1, 2
- [16] Laurent Itti, Christof Koch, and Ernst Niebur. A model of saliency-based visual attention for rapid scene analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20(11):1254–1259, 1998. 1, 2, 4
- [17] Sen Jia and Neil D. B. Bruce. EML-NET: An expandable Multi-Layer NETWORK for saliency prediction. *Image and Vision Computing*, 95:103887, 2020. 2
- [18] Ming Jiang, Shengsheng Huang, Juanyong Duan, and Qi Zhao. SALICON: Saliency in context. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2015. 1, 2, 6
- [19] Timothée Jost, Nabil Ouerhani, Roman von Wartburg, René Müri, and Heinz Hügli. Assessing the contribution of color in visual attention. *Computer Vision and Image Understanding*, 100(1-2):107–123, 2005. 5, 6
- [20] Tilke Judd, Fredo Durand, and Antonio Torralba. A benchmark of computational models of saliency to predict human fixations. *MIT Technical Report*, 2012. 6
- [21] Tilke Judd, Krista Ehinger, Fredo Durand, and Antonio Torralba. Learning to predict where humans look. In *IEEE International Conference on Computer Vision (ICCV)*. IEEE, 2009. 2
- [22] Alexander Kroner, Mario Senden, Kurt Driessens, and Rainer Goebel. Contextual encoder–decoder network for visual saliency prediction. *Neural Networks*, 129:261 – 270, 2020. 2, 7
- [23] Matthias Kümmerer, Lucas Theis, and Matthias Bethge. Deep Gaze I: Boosting saliency prediction with feature maps trained on ImageNet. In *International Conference on Learning Representations (ICLR) Workshops*, 2015. 2
- [24] Matthias Kümmerer, Tom Wallis, and Matthias Bethge. DeepGaze II: Predicting fixations from deep features over time and tasks. *Journal of Vision*, 17(10):1147, 2017. 1, 2, 4
- [25] Matthias Kümmerer, Thomas S. A. Wallis, and Matthias Bethge. Information-theoretic model comparison unifies saliency metrics. *Proceedings of the National Academy of Sciences*, 112(52):16054–16059, 2015. 6
- [26] Qingshan Li, Yue Zhou, and Jie Yang. Saliency based image segmentation. In *2011 International Conference on Multimedia Technology*, pages 5068–5071. IEEE, 2011. 1
- [27] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft COCO: Common objects in context. In *European Conference on Computer Vision (ECCV)*, pages 740–755. Springer, 2014. 2
- [28] Akis Linardos, Matthias Kümmerer, Ori Press, and Matthias Bethge. Deepgaze iie: Calibrated prediction in and out-of-domain for state-of-the-art saliency modeling. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12919–12928, October 2021. 1, 4, 6, 7
- [29] Gitte Lindgaard, Gary Fernandes, Cathy Dudek, and J. Brown. Attention web designers: You have 50 milliseconds

- to make a good first impression! *Behaviour & Information Technology*, 25(2):115–126, 2006. 1, 4
- [30] Chenxi Liu, Barret Zoph, Maxim Neumann, Jonathon Shlens, Wei Hua, Li-Jia Li, Li Fei-Fei, Alan Yuille, Jonathan Huang, and Kevin Murphy. Progressive neural architecture search, 2018. 4
- [31] Nian Liu and Junwei Han. A deep spatial contextual long-term recurrent convolutional network for saliency detection. *IEEE Transactions on Image Processing*, 27(7):3264–3274, 2018. 1
- [32] Sandeep Mishra and Oindrila Saha. Recsal : Deep recursive supervision for visual saliency prediction, 2020. 2
- [33] Johanna Palcu, Jennifer Sudkamp, and Arnd Florack. Judgments at gaze value: Gaze cuing in banner advertisements, its effect on attention allocation and product judgments. *Frontiers in Psychology*, 8, 2017. 1, 4
- [34] Yash Patel, Srikar Appalaraju, and R. Manmatha. Saliency driven perceptual image compression, 2020. 1
- [35] Robert J. Peters, Asha Iyer, Laurent Itti, and Christof Koch. Components of bottom-up gaze allocation in natural images. *Vision Research*, 45(18):2397–2416, 2005. 6
- [36] Michael Posner, Robert Rafal, Lisa Choatec, and Jonathan Vaughan. Inhibition of return: Neural basis and function. *Cognitive Neuropsychology*, Vol. 2:211 – 228, 09 1985. 1
- [37] Navyasri Reddy, Samyak Jain, Pradeep Yarlagadda, and Vineet Gandhi. Tidying deep saliency prediction architectures. *CoRR*, abs/2003.04942, 2020. 4, 6, 7
- [38] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation, 2015. 4
- [39] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *International Conference on Learning Representations*, 2015. 2
- [40] Markus Steinberger, Bernhard Kainz, Stefan Hauswiesner, Rostislav Khlebnikov, Denis Kalkofen, and Dieter Schmalstieg. Ray prioritization using stylization and visual saliency. *Computers & Graphics*, 36(6):673–684, 2012. 1
- [41] Po-He Tseng, Ran Carmi, Ian G. M. Cameron, Douglas P. Munoz, and Laurent Itti. Quantifying center bias of observers in free viewing of dynamic natural scenes. *Journal of Vision*, 9(7):4–4, 07 2009. 3
- [42] Mathukumalli Vidyasagar. Kullback-leibler divergence rate between probability distributions on sets of different cardinalities. In *IEEE Conference on Decision and Control (CDC)*, pages 948–953. IEEE, 2010. 5, 6
- [43] Eleonora Vig, Michael Dorr, and David Cox. Large-scale optimization of hierarchical features for saliency prediction in natural images. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2798–2805. IEEE, 2014. 1, 2
- [44] Wenguan Wang, Qiuxia Lai, Huazhu Fu, Jianbing Shen, Haibin Ling, and Ruigang Yang. Salient object detection in the deep learning era: An in-depth survey. *arXiv*, page arXiv:1904.09146, 2019. 1
- [45] Sheng Yang, Guosheng Lin, Qiuping Jiang, and Weisi Lin. A dilated inception network for visual saliency prediction. *IEEE Transactions on Multimedia*, 22(8):2163–2176, 2020. 1
- [46] Alfred L. Yarbus. *Eye Movements and Vision*, volume 2. Plenum Press, 1967. 1, 3
- [47] Yifeng Zhang, Ming Jiang, and Qi Zhao. Saliency prediction with external knowledge. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 484–493, January 2021. 1, 2
- [48] Jufeng Zhao, Yueting Chen, Huajun Feng, Zhihai Xu, and Qi Li. Fast image enhancement using multi-scale saliency extraction in infrared imagery. *Optik*, 125(15):4039–4042, 2014. 1