

000
001
002
003
004
005
006
007
008
009
010
011054
055
056
057
058
059
060
061
062
063
064
065

How Many Glances? Modeling Multi-duration Saliency

Anonymous ICCV submission

Paper ID 5002

Abstract

What jumps out at you in a single glance of an image is different than what you might notice after closer inspection. Despite this, models of visual saliency have ignored the temporal aspect of visual attention and have produced prediction maps at fixed viewing durations. As a result, current applications leveraging saliency models are rigidly tailored for a fixed viewing duration, depending on the attention dataset they were trained on. To incorporate knowledge of viewing duration into saliency modeling, we first develop a “Codecharts UI” and use it to crowdsource human attention at various viewing durations. We collect the CodeCharts1K dataset, which contains viewing data at 0.5, 3, and 5 seconds on 1000 images from diverse computer vision datasets. Our analysis shows distinct differences in gaze locations at these time points and exposes recurring temporal patterns about which objects attract attention. Using insights from this analysis, we develop an LSTM-based model of saliency that simultaneously trains on data from multiple time points. Our Multi-Duration Saliency Excited Model (MD-SEM) achieves state-of-the-art performance on the LSUN 2017 Challenge with 57% fewer parameters than comparable architectures. Additionally, it provides new predictions at time points that current models can not predict. We explore applications where multi-duration saliency can be used to prioritize the visual content to keep, transmit, and render.

1. Introduction

How long an observer examines an image determines what they notice and what tasks they can complete. Despite this dependency of viewing behavior on time, most models of visual attention predict saliency at a fixed duration, e.g., by training on data collected with a viewing time of 3 or 5 seconds per image¹ [5, 18, 27, 25, 34].

In this paper, we introduce the first saliency model that outputs multiple saliency maps based on the viewing dura-

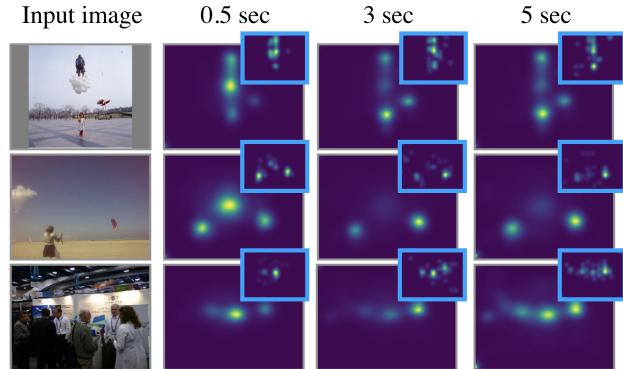


Figure 1. Predictions of our Multi-Duration Saliency Excited Model for images at viewing durations of 0.5, 3, and 5 seconds. Images here are from the Abnormal Objects [37], SALICON [27], and EyeCrowd [28] datasets (top to bottom row). Insets with blue borders contain human ground-truth gaze locations collected using our CodeCharts UI.

tion of an observer. We propose an efficient crowdsourcing methodology that allows us to collect at large scale human attention data at several viewing durations. We use it to assemble CodeCharts1K, a dataset of 1000 images with viewing patterns at three durations (0.5, 3, and 5 seconds) encompassing a variety of image genres such as actions, out-of-context objects, and memorable images. We find that humans are highly consistent with each other at a specific viewing duration but show significant differences across durations.

We then introduce an original multi-duration saliency model which, given an image as input, predicts saliency maps for three different durations. Our model achieves state-of-the art performance when evaluated at a single duration and outperforms other baseline models if they are trained to predict multiple durations. We show that the predicted saliency maps can be used as input to applications such as image cropping, compression and rendering, and captioning, to tailor them to different contexts based on time.

¹See <http://saliency.mit.edu/datasets.html>

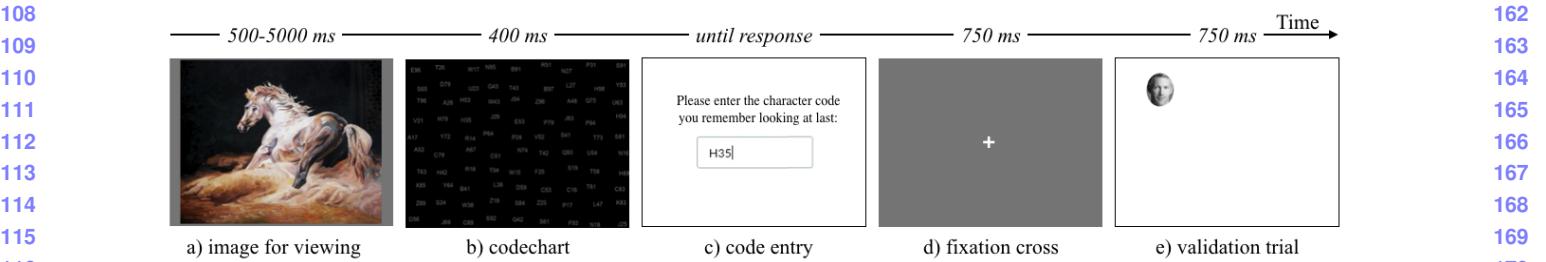


Figure 2. CodeCharts UI task flow: participants view an image for a variable amount of time (a), after which a codechart is briefly flashed on the screen (b). Participants are asked to report the character code they remember looking at last (c), as an indication of their gaze on the image. This process proceeds for a sequence of images, separated by fixation crosses (d) to re-center gaze. Validation trials (e) are interspersed among the experimental sequence to check if participants move their gaze as expected.

2. Related work

Recent efforts have engaged in collecting large-scale attention data using interfaces that can be run remotely, without an eye tracker [7, 27, 31, 32, 36]. Webcam-based approaches [32, 36] rely on good visibility and stability of the viewer in front of the web camera and raise privacy concerns. Such approaches, trained to predict gaze location from images of the eyes, are sufficient for coarse estimates of gaze location but have error rates that do not permit good localization. On the other hand, moving-window approaches like SALICON [27] and BubbleView [31] allow people to use their cursors to inspect small, bubble-shaped regions of blurred images. Data collected is pixel-accurate but has two downsides: (1) blurring images distorts the visual content and interferes with feature sizes in an image, and (2) moving the mouse cursor requires a different process than moving the eyes, which can impact which image regions are explored [31]. The CodeCharts UI is a response to these concerns. This approach does not require webcams nor does it distort the image, and it allows for fine-grained control of image presentation time. We exploit this latter property to capture multi-duration attention data.

The large-scale attention data captured using the SALICON [27] and BubbleView [31] methods has enabled training neural network models of saliency (e.g., [7, 18, 25, 34]). The top performers on the MIT Saliency Benchmark [5] have been trained on SALICON data, and have opened a wide performance gap to the previous, non neural network based models of saliency [8]. Driven by such improvements in efficiency and accuracy, saliency models have found wide use in applications like image cropping, retargeting, and view-finding for improved composition [3, 11, 20, 40].

3. Dataset of multi-duration attention

We introduce an approach to capture human gaze data with a user interface that does not require explicit eye tracking but successfully captures the same attention patterns. Importantly, this interface allows easily manipulating the

image presentation time for large-scale collection of multi-duration attention data.

3.1. CodeCharts UI

Task flow: We revisit the user interface of Rudoy et al. [39], which introduced the idea of using a character chart to validate self-reported gaze locations. In our task, participants view an image for 500-5000 milliseconds (ms), followed by a jittered grid of random, three-character codes (“codechart”). They then self-report the first three-character code they see when the image vanishes (Fig. 2a-c). Based on pilot experiments, the codechart duration is set as short as possible to maintain accurate performance. By construction, participants report the region of the image they were looking at last. We repeat the steps in Fig. 2a-c for dozens of images, shown in sequence in a single experiment. A fixation cross between trials helps to re-center gaze (Fig. 2d).

Validation: To ensure good quality data, we check that the three-character codes entered by participants are “valid”, i.e., they occur in the corresponding codechart. We also insert validation trials throughout the image sequence: cropped faces [2] randomly placed on the image canvas. A code is marked “correct” if it overlaps with some portion of the face in the corresponding codechart. A pre-test phase is used to screen participants based on their performance entering valid codes for regular images and correct codes for face images. Our validation procedure has a significant impact on collected data quality (see Supplement).

3.2. Data collection

Pilot experiments: We ran an initial experiment with images from the CAT2000 dataset [4] to determine the optimal task flow (timing of task components, validation procedure, number of participants, etc.). We then sampled 50 images from the OSIE dataset [41] and collected the gaze locations of 50 participants per image for each of 6 image durations: 0.5, 1, 2, 3, 4, and 5 seconds. We found that the gaze maps obtained at 0.5, 3, and 5 seconds were the most different from each other (Fig. 3a), which motivated our de-

216 cision to collect further data at these 3 durations. The gaze
 217 locations collected with our CodeCharts UI at 3 seconds
 218 most closely matched the ground-truth OSIE data, which
 219 was originally also collected at 3 seconds, with a Pearson’s
 220 Correlation Coefficient (CC score) of .71. This validates the
 221 ability of CodeCharts data to model natural human gaze.
 222

223 **CodeCharts1K:** For our data collection, we sampled
 224 a variety of image types to better understand how attention
 225 evolves with time. We used 500 images from SALI-
 226 CON [27], 130 from LaMem [30], 120 from CAT2000 [4]²,
 227 100 from EyeCrowd [28], 100 from a mix of Abnormal
 228 Objects [37] and Out-of-Context Objects [15], and 50 from the
 229 Stanford 40K Actions dataset [43]³. Images were padded to
 230 have the same aspect ratio, and were resized in-browser to
 231 fit in an image window of 700×1000 pixels. The task se-
 232 quence included 6 “tutorial” images to test participant atten-
 233 tiveness, 50 dataset images, and 5 validation trials of faces,
 234 spaced throughout the sequence. We used Amazon’s Me-
 235 chanical Turk and paid participants at an hourly rate of \$10.
 236 Data collection cost \$4.90 per image for 150 unique gaze
 237 points (50 participants each at 0.5, 3, and 5 seconds).
 238

4. Data analysis

4.1. Do people look in the same places?

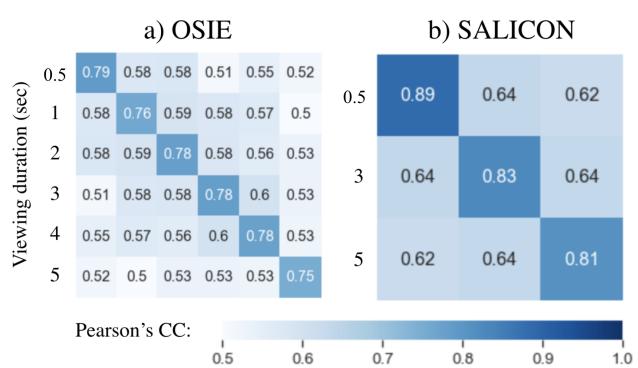
241 To determine whether participants look in the same lo-
 242 cations, we perform a split-half consistency analysis. We
 243 compare the gaze points of one half of the participants to
 244 the gaze points of the other half, by converting both sets of
 245 gaze points to heatmaps and computing a Pearson’s Corre-
 246 lation Coefficient (CC score) between them. We perform
 247 this analysis per viewing duration. By performing this anal-
 248 ysis also *across* viewing durations, we measure whether the
 249 gaze patterns are different for different viewing durations.
 250

251 From the diagonal entries in Fig. 3, we see that the split-
 252 half consistency between participants is very high across all
 253 durations. This finding replicates across all the image sets
 254 we tested (more results in the Supplement). Importantly,
 255 however, **consistency does not degrade over time**: the dif-
 256 ferences between the diagonal entries are not statistically
 257 significant according to a two-sided T-test ($p > 0.05$ for all
 258 pairwise comparisons). In other words, **saliency is equally
 259 predictable across all the viewing durations**.

260 From the off-diagonal entries in Fig. 3, we see that
 261 the gaze locations of participants viewing an image at one
 262 duration are different from the gaze locations of partici-
 263 pants viewing an image at a different duration. These
 264 differences are statistically significant for all pairs of du-
 265 rations ($p < 0.001$, Bonferroni-corrected), with a replication
 266

²We used 100 images from the ‘Action’ [43] and 20 from the ‘Low Resolution’ [29] categories.

³We used action classes that explicitly contained an interaction of a person and an object, by selecting 10 images each of: shooting an arrow, throwing a frisby, walking the dog, writing on a board, writing on a book.



270
 271
 272
 273
 274
 275
 276
 277
 278
 279
 280
 281
 282
 283
 284
 285
 286
 287
 288
 289
 290
 291
 292
 293
 294
 295
 296
 297
 298
 299
 300
 301
 302
 303
 304
 305
 306
 307
 308
 309
 310
 311
 312
 313
 314
 315
 316
 317
 318
 319
 320
 321
 322
 323

Figure 3. Split-half consistency of viewers across datasets and viewing durations (using Pearson’s Correlation Coefficient). Participants are remarkably consistent in where they look. This consistency persists across different viewing durations, and replicates across datasets. Consistency is much higher between participants viewing images at the same duration, than across durations, indicating that the saliency patterns are different across durations.

302 across datasets. In other words, **different things are salient at different viewing durations**.

303 These analyses indicate that gaze data collected using
 304 our CodeCharts UI (i) contains a consistent signal at each
 305 of the viewing durations, (ii) this signal is equally strong at
 306 all the viewing durations, and (iii) the signal is significantly
 307 different between viewing durations. These findings set the
 308 stage for the computational model presented in Sec. 5.

4.2. What is salient at what time?

309 **Things and stuff:** We used COCO segmentation
 310 maps [9] of the SALICON images to compute gaze counts
 311 per image segment across time (full analysis in Supple-
 312 ment). From 0.5 sec to 3 sec, gaze frequently moves away
 313 from people and is pulled towards objects and furniture
 314 (e.g., paper, bottle, table). From 3 to 5 sec, there is an in-
 315 crease of attention on “stuff” like grass, carpet, and road that
 316 may contain other objects. At these longer durations people
 317 gaze more at small and distant objects in an image.

318 **Faces:** It is known that gaze is attracted by faces [8, 10].
 319 For a finer-grained analysis over time, we ran a face detec-
 320 tion network [21] over all the images in CodeCharts1K.
 321 For the 303 images on which faces were detected, we com-
 322 puted face saliency by aggregating gaze counts per face
 323 region and then normalizing both by the number of gaze
 324 points per image and across all 3 durations. This ensures
 325 that for a given image across time, face saliency ranges be-
 326 tween 0 and 1. Fig. 4 plots these results as one line per
 327 image, where thick lines are averages to visualize the dom-
 328 inant patterns. Across our CodeCharts1K dataset, we find a
 329 dominant “boomerang” pattern: people start out by looking
 330 at faces at 0.5 sec, their gaze shifts elsewhere at 3 sec, and
 331

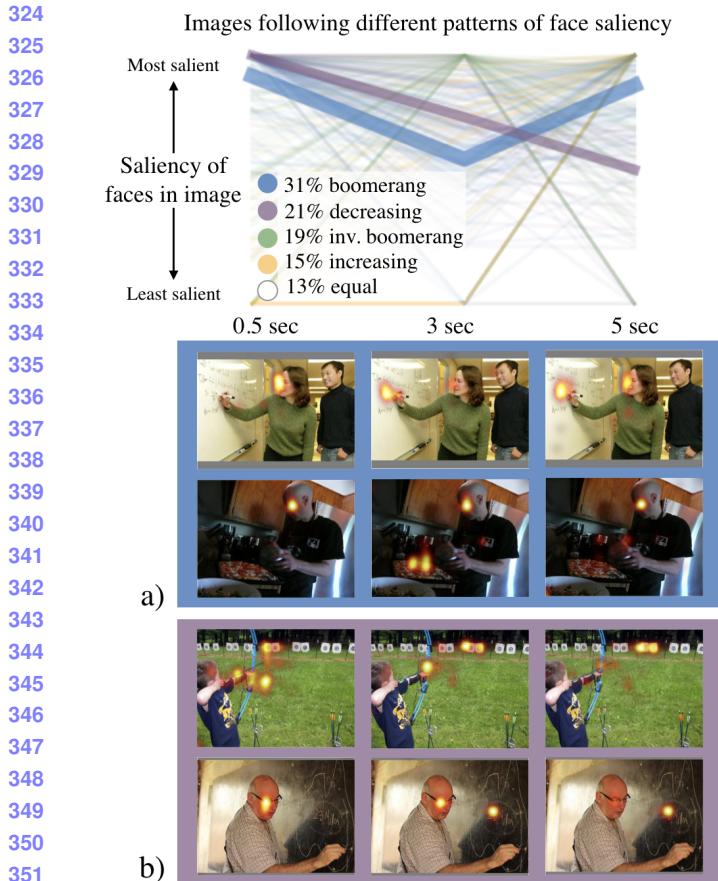


Figure 4. Dominant patterns of human gaze on faces across time: (a) face saliency decreases from 0.5 to 3 sec, and increases again from 3 to 5 sec (“boomerang”), and (b) face saliency decreases from 0.5 to 5 sec. Each line in the graph above represents how the saliency of faces within a single image varies over time, while the thicker lines represent an average over images. We normalize the saliency of faces across time on a per-image basis in this plot.

returns to faces at 5 sec. The second most prevalent pattern is a decrease in gaze on faces over time. Other patterns, like an increase in face saliency over time, were in the minority.

Qualitatively, human gaze frequently moves from the actor (at 0.5 sec) to the action (at 3 and 5 sec). Sometimes this shift in attention is gradual: saliency at 3 sec is a combination of saliency at 0.5 and 5 sec (Fig. 4b). In other cases, saliency at 5 sec is a combination of saliency at the shorter durations: we hypothesize that in these cases, people’s attention is initially grabbed by faces, before they explore an image and return to the most salient regions (Fig. 4a).

5. Multi-duration saliency model

In order to predict these changes in human attention over time, we introduce the Multi-Duration Saliency Excited Model (MD-SEM), a new architecture adapted to

multi-duration saliency. Our model is capable of producing saliency maps at T different viewing durations (where $T=3$ in our implementation). The architecture is based on two concepts: (1) a powerful encoder-decoder architecture built on an ImageNet-pretrained Xception network [16]; (2) a Temporal Excitation Module, a novel block that applies a time-based re-weighting to saliency feature maps with a minimal increase in parameters. Our model achieves first place on the LSUN 2017 SALICON saliency challenge, and is the current state-of-the-art on CodeCharts1K, our multi-duration saliency dataset.

5.1. Architecture motivation

The current state-of-the-art saliency models are built on convolutional neural networks, exploiting their ability to understand semantic components of images to produce accurate saliency maps. Current architectures tend to be bulky, with large numbers of parameters and specialized modules. One of the current top performers, EML-NET [26], contains both a DenseNet and NASNet. Another top performer, the Saliency Attentive Model (SAM) [18], uses an Attentive Convolutional LSTM for refinement as well as a Learned Prior module with several 5×5 convolutions, bringing the total model size to more than 70 million parameters. We present a model of reduced size and complexity that nevertheless achieves an increase in accuracy.

We simplify saliency prediction by distilling the required components to the bare minimum: (1) a strong encoder, (2) a processing module on the compressed representation, and (3) a regularized decoder. Our multi-duration saliency network obeys these three constraints: we use a parameter-efficient encoder, a concise temporal re-weighting module (which explicitly performs computations on vectors with reduced dimensionality), and a regularized decoder composed of 3 sets of dilated convolutions with Dropout.

5.2. Convolutional encoder-decoder

Convolutional encoder-decoder architectures have been extensively used, being particularly effective for image-to-image tasks like segmentation [1, 12, 33], saliency prediction [26, 34] and importance prediction [7]. Encoding the image input through a convolutional encoder allows for rich feature extraction and reduces the dimensionality of the input to values that facilitate feature manipulation. Decoding the generated representation through convolutional layers can be thought of as learned upsampling. Given the necessity of deep semantic understanding, we use a state-of-the-art network in image recognition as our encoder: the Xception network [16], pretrained on ImageNet. The Xception network is lightweight and accurate (0.790 top-1 accuracy on ImageNet with only 22M parameters), and has shown success in semantic segmentation [13], a related task. It proposes a complete decoupling of cross-channel and spatial

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

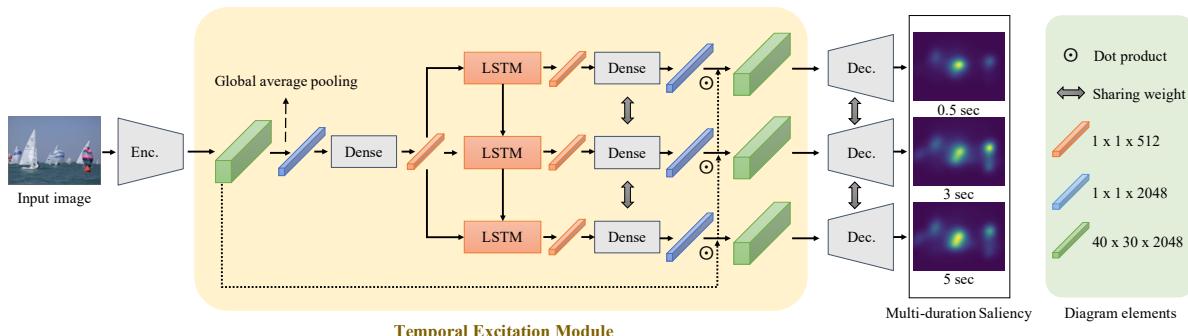


Figure 5. Multi-Duration Saliency Excited Model (MD-SEM) architecture.

correlations of feature maps. Its convolutional layers perform depthwise-separable convolutions, operating on each channel separately before merging the transformed channels through a pointwise (1×1) convolution.

5.3. Temporal Excitation module

In order to predict saliency across durations, we introduce a module capable of recursively manipulating the feature representation generated by the encoder to adapt it to each timestep. Our module utilizes a Long Short Term Memory (LSTM) network to generate scaling vectors that re-weight the feature maps differently for each timestep. Feature map re-weighting is a concept that has been explored before in the form of Squeeze and Excitation Modules [24], but has not been exploited as a temporal modification tool. A previous top-performing saliency model, SAM [18], also contains an LSTM, but only for internal map refinement, rather than for capturing the temporal aspect of saliency. The main advantage is simplicity: whereas SAM utilizes a parameter-heavy Attentive Convolutional LSTM that operates on full feature maps, we approach the problem with a simple LSTM that works on a squeezed low-dimensional vector obtained from pooling those maps. In Sec. 6.3, we show that our architecture performs better than SAM on predicting multi-duration saliency.

The architecture of the Temporal Excitation Module is included in Fig. 5. First, the feature maps generated by the encoder are pooled through global average pooling, generating a $1 \times 1 \times C$ vector. The vector is then passed through a dense (fully connected) layer, which reduces the dimensionality of the vector and aids in generalization. As explained in [24], applying this type of operation is key to modeling channel-wise relationships. The output of the dense layer is replicated T (in our case, 3) times and fed as a sequence to the LSTM. This is the key of the Temporal Excitation module: the LSTM outputs T vectors, which contain information specific to each timestep and will scale each feature map differently. A ReLU non-linearity is applied before passing the vectors through a dimensionality-

increasing dense layer. Finally, a sigmoid non-linearity ensures that the scaling weights s_c remain in a sensible range, following insights by [24]. The block outputs a set of T feature maps, which are rescaled according to each one of the weight vectors generated. Formally, the block outputs T sets of feature maps $F^{(t)}$, where each feature map $f_c^{(t)}$ is computed as:

$$f_c^{(t)} = I_c * s_c^{(t)},$$

where I_c are the c input feature maps, $f_c^{(t)}$ are the T sets of c feature maps, and $s_c^{(t)}$ are the scaling weights generated by the previously defined operations.

The Temporal Excitation module introduces few new parameters to the model. With the choice of hyper-parameters described in Sec. 5.4, we obtain an architecture with 30 million parameters, 57% smaller than SAM [18]. This allows us to further reduce model complexity and overfitting.

5.4. Implementation details

Architecture: Our network's blocks are implemented as follows. The Xception decoder is instantiated with its original definition, removing the last fully connected layer, and obtaining a decoder that outputs an $H \times W \times 2048$ tensor of feature maps. The Temporal Excitation Module is instantiated with a 512-wide fully connected layer, followed by an LSTM with 512 cells, a ReLU non-linearity, and a sigmoid-activated fully connected layer with 2048 parameters to transform the scaling vector back to its input size. Finally, the decoder is composed of 3 sets of convolutional blocks. The first block contains two convolutional layers with 256 filters each, followed by a 2×2 upsampling layer and a dropout operator with probability 0.3. The second and third blocks are identical, but use 128 and 64 filters each on their convolutions. Finally, a 1×1 convolution with 1 filter is used to get the final set of feature maps to a single-channel saliency heatmap. We set T=3 to match the data available in the CodeCharts1K dataset. Note that the same decoder is applied to each one of the T outputs of the Temporal Excitation Module, thus concentrating time information exclu-

540 sively in that module, and reducing model complexity. A
 541 description of the architecture is in the Supplement.
 542

543 **Core loss:** Following previous works, the network’s
 544 core loss is defined as a weighted combination of Kull-
 545 back Leibler divergence (KL), Linear Correlation Coeffi-
 546 cient (CC) and Normalized Scanpath Saliency (NSS) (see
 547 [6] for details on the formulas). As discussed in [6], NSS
 548 is more robust than other metrics at measuring the quality
 549 of saliency predictions. Thus, we defined the weights of the
 550 network with a higher emphasis on NSS. After evaluation,
 551 we defined the weights to be 10 for KL, -5 for CC and -10
 552 for NSS during SALICON-MD training (see subsection 6.1
 553 for details on the dataset), but changed the NSS weight to -1
 554 for CodeCharts training to account for the reduction in the
 555 number of fixations per image.

556 **CC Match:** To ensure that our network is correctly mod-
 557 eling the differences across durations, we introduce a novel
 558 training loss called Correlation Coefficient Match (CCM).
 559 This loss forces the network to output saliency maps that are
 560 similar to the ground-truth saliency maps at each duration
 561 yet different across durations. Given a set of T durations
 562 for which we want to predict saliency maps, we calculate
 563 the CCM loss by computing Pearson’s Correlation Coeffi-
 564 cient (CC) on pairs of saliency maps at adjacent durations,
 565 then computing the difference between the ground-truth and
 566 predicted scores. Let $y^{(j)}$ be the heatmap corresponding to
 567 duration j where CC is defined as:

$$CC(y_1, y_2) = \frac{\sigma(y_1, y_2)}{\sigma(y_1) \cdot \sigma(y_2)}, \quad (1)$$

571 where $\sigma(y_1, y_2)$ is the covariance of y_1 and y_2 . Our CCM
 572 loss is:

$$L_{CCM}(y_g, y_p) = \frac{1}{T-1} \sum_{j=0}^{T-1} \left| CC\left(y_g^{(j)}, y_g^{(j+1)}\right) - CC\left(y_p^{(j)}, y_p^{(j+1)}\right) \right| \quad (2)$$

579 where $y_g^{(t)}$ and $y_p^{(t)}$ are the ground-truth and predicted
 580 saliency maps for time t . Our experiments showed that a
 581 weight of 3 is appropriate for this loss. We show in Ta-
 582 ble 1 that this novel loss significantly boosts our per-
 583 formance when evaluating multi-duration predictions.

585 6. Evaluation

586 6.1. Datasets

589 For training, we rely on the SALICON-MD (Multi-
 590 Duration) and CodeCharts1K datasets. SALICON-MD is
 591 derived from the original SALICON dataset by bucketing
 592 a participant’s fixations based on when they occurred. We
 593 assume that fixations occur at a constant interval across the

Model	NSS ↑	CC ↑	KL ↓	SIM ↑	CCM ↑	
SAM-MD w/o CCM	2.700	0.744	0.434	0.616	0.041	594
SAM-MD w/ CCM	2.765	0.778	0.401	0.641	0.061	595
MD-SEM w/o CCM	2.778	0.754	0.565	0.598	0.073	596
MD-SEM w/ CCM	2.875	0.773	0.392	0.633	0.091	597

598 Table 1. MD-SEM results on CodeCharts1K with and without
 599 CCM loss. We report performance on NSS, CC, KL, SIM and our
 600 custom CC Match loss. NSS and CC are effective saliency metrics
 601 (they are symmetric to false positives and false negatives [6]). A
 602 lower value is better for KL.

603 viewing duration (from 0 to 5 sec) and split them into 6 non-
 604 overlapping buckets. This time-separated data serves as an
 605 approximate but large pretraining dataset. For final train-
 606 ing and evaluation, we use ground-truth multi-duration data
 607 from CodeCharts1K, defined in Sec. 3.2.

608 6.2. Training details

609 Our training scheme takes advantage of both datasets
 610 to create a model that is generalizable and accurate at ev-
 611 ery duration. In order to learn from as much data as pos-
 612 sible, we pretrain on SALICON-MD. Pretraining on tem-
 613 poral data that exhibits differences across durations is im-
 614 portant so that our model learns at the outset to discrimi-
 615 nate between timesteps. We then fine-tune on ground-truth
 616 CodeCharts1K. For both datasets, we set the batch size to 8
 617 and the initial learning rate to 1e-4, which is reduced by a
 618 factor of ten every three epochs. At the beginning of train-
 619 ing we freeze the weights of the encoder for one epoch,
 620 which forces the Temporal Excitation Module to make use
 621 of the ImageNet-pretrained feature maps from the encoder
 622 instead of distorting them. We find that 10 epochs of train-
 623 ing on SALICON-MD and 5 on CodeCharts1K is sufficient.

624 For SALICON-MD, we use the pre-established test,
 625 train, and validation sets. For the code charts data, we train
 626 on 70% of the images, validate on 5%, and test on 25%.
 627 Our reported test results for CodeCharts1K are the average
 628 scores from cross-validating on four splits of the data.

629 6.3. Comparison to state-of-the-art

630 On CodeCharts1K, our final model achieves an NSS of
 631 2.875 and CC of 0.773 averaged across the three durations
 632 (Table 2). This is compared to an average ground-truth hu-
 633 man score of 0.843 on the same dataset (Fig. 3). In addition
 634 to accurately capturing gaze location at a single timestep,
 635 our high CCM score (see Table 1) indicates that our network
 636 effectively captures the variability of human gaze across du-
 637 rations. Our model reflects many of the qualitative patterns
 638 of human gaze movement that stand out in the code charts
 639 data, such as the boomerang effect and the tendency of hu-
 640 mans to focus on the object of an action (see Fig. 6).

641 Our model is first-of-its-kind in its ability to predict
 642 multi-duration saliency. Thus, we benchmark against the

648
649
650
651
652

Model	NSS \uparrow	CC \uparrow	KL \downarrow	SIM \uparrow	Params \downarrow
SAM x3	2.020	0.803	0.330	0.708	210.3M
SAM-MD	2.057	0.792	0.338	0.695	70.1M
MD-SEM	2.061	0.811	0.324	0.708	30.9M

653
654
655
656
657
Table 2. Comparison to state-of-the-art models on SALICON-MD.
Our model marginally outperforms SAM-MD (SAM adapted by
ourselves for multi-duration prediction) with 57% less parameters.
We also report improvements on SAM trained on each timestep
separately (SAM x3).

Model	NSS \uparrow	CC \uparrow	KL \downarrow	SIM \uparrow
SAM-res [18]	1.990	0.899	0.610	0.793
EML-Net [26]	2.050	0.886	0.520	0.780
SalNet [35]	1.859	0.622	-	-
CEDNS*	2.045	0.862	1.026	0.753
MD-SEM (Ours)	2.058	0.868	0.568	0.774

666
667
Table 3. Comparison to state-of-the-art on SALICON test set
(LSUN 2017 Challenge).668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
currently-available option for generating saliency at multiple durations: training T separate copies of a standard saliency network, each on one of T splits of multi-duration saliency data. Furthermore, in order to prove that our architecture is uniquely effective at the multi-duration problem, we also benchmark against a modified version of SAM. In SAM-MD, the LSTM is modified to produce a map at each timestep, and during training these maps are compared to the ground-truth at the corresponding duration. The results of these benchmarks are in 2. Not only is MD-SEM better at approximating human gaze and differentiating across durations, but it also uses significantly fewer parameters than the other baselines.683
684
685
Finally, we show that our model performs at state-of-the-art for traditional single-duration saliency by achieving first place on the LSUN 2017 challenge (Table 3).686

7. Applications

688
689
690
691
692
693
As saliency neural network models approach human-level prediction, saliency has been used for many applications including cropping, retargeting, and image captioning. Our multi-duration saliency model can bring additional context by accounting for the expected time that a viewer may have to explore an image.694
695
696
697
698
699
700
701
Cropping: Automatic cropping of images is useful for image thumbnails, view-finding for improved composition, and retargeting for different use cases [19]. Multi-duration saliency allows us to additionally take into account the expected time a viewer will spend on an image, and retarget the image accordingly (e.g., an image that is part of a passing advertisement should be simpler than if it is the main image on a page). In Fig. 7 we used our multi-duration702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
Figure 6. Saliency predictions of MD-SEM on various datasets. Insets with blue borders contain human ground-truth gaze locations collected using our CodeCharts UI. a) The “boomerang” pattern of saliency: starts out at face (0.5 sec), moves to object (3 sec), and back to face (5 sec). b) The boomerang pattern applied to objects: gaze starts at central object, moves to new salient location, then back to initial object. c) More distant objects (especially faces) become the objects of focus at later durations. d) Saliency distributes over time from faces to multiple scene elements. Difficult cases for our model: e) Our model predicts faces will be salient for all durations, but people focus on the abnormal region of the scene (woman standing on man’s head). f) Model knows to change focus of saliency from one face to another between durations, but chooses the wrong face to focus on relative to ground-truth. g) Model distributes saliency over too many objects in a scene when it’s not clear which one should be the focus.756
757
758
759
760
761
saliency maps to crop windows [14] within an image that capture 90% of the most salient image regions *at that viewing duration*. Our automatically-generated thumbnails tend to contain close-ups of content for shorter durations.762
763
Compression and rendering: Just as saliency has been



Figure 7. **Cropping.** Images cropped based on viewing time by selecting the window with 90% of the most salient image regions, as predicted by our multi-duration saliency model. Automatically-computed image crops for shorter viewing durations tend to contain close-ups of image content, focusing on fewer elements.

used as a mechanism to prioritize the visual content to preserve and render, multi-duration saliency can add a temporal aspect to these applications. For instance, if an image is expected to be viewed for shorter periods of time, fewer visual elements need to be rendered (or preserved, in the case of compression). In Fig. 8 we provide a visualization of which visual content would be prioritized at different viewing durations for such applications. To generate these visualizations, we use Mask R-CNN for instance segmentation on the image [22]. The detected instances are then mapped to the saliency predictions from our model. Rather than using the saliency directly at a given time point, we accumulate the saliency across time. Instances that have a mean saliency score in the 90-th percentile are kept, using saliency at 0.5 sec as a reference, and the rest of the image is blurred and darkened, for visualization purposes.

Captioning: Captions can be used to facilitate search and improve accessibility. Despite captioning being tied to a sequential viewing of the image, current approaches do not leverage this information to produce better captions. Some approaches have used an attention mechanism to attend to different parts of the image [42], yet machine attention does a worse job of approximating human attention than saliency models when performing captioning [23]. While some recent work attempts to produce an attention map at each timestep for captioning [17], the saliency model used does not explicitly model the temporal aspect of human attention. Our multi-duration saliency maps offer a closer approximation to how humans view images and provide an opportunity to focus . Here, we used our saliency



Figure 8. **Compression.** The more you look, the more you see. Visualized are image regions that are predicted to attract gaze at different viewing durations (accumulated over time). A better understanding of saliency across time can facilitate saliency-driven applications like image compression, transmission, and rendering to take viewing duration into account. Will the end user have less than 1 second to look at an image? Saliency can be used to prioritize the visual content to retain, transmit, and render.



Figure 9. **Captioning.** Captions generated by passing saliency-enhanced images to an image captioning model [38], using saliency at different durations to prioritize image content.

predictions to focus an image captioning model [38] on image regions that should stand out at different viewing durations. Removing the non-salient visual clutter indeed produces promising initial results (Fig. 9).

8. Conclusion

Guided by results on how time affects people’s focus of attention, we introduce an LSTM-based model of saliency which, when trained on crowdsourced gaze data from different viewing times, achieves state-of-the-art performance, while providing novel predictions at time points that other saliency approaches cannot model. Importantly, we show that the time an observer, or a system, spent on an image determine which tasks they can complete. We demonstrate the direct impact of our multi-duration saliency approach to applications which require prioritizing the visual content for image cropping, compression and rendering and captioning.

864

References

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

- [1] V. Badrinarayanan, A. Kendall, and R. Cipolla. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 39(12):2481–2495, 2017. 4
- [2] W. A. Bainbridge, P. Isola, and A. Oliva. The intrinsic memorability of face photographs. *Journal of Experimental Psychology: General*, 142(4):1323, 2013. 2
- [3] A. Borji. Saliency prediction in the deep learning era: An empirical investigation. *arXiv preprint arXiv:1810.03716*, 2018. 2
- [4] A. Borji and L. Itti. CAT2000: A large scale fixation dataset for boosting saliency research. *CVPR’15 Workshop on the Future of Datasets*, 2015. 2, 3
- [5] Z. Bylinskii, T. Judd, A. Borji, L. Itti, F. Durand, A. Oliva, and A. Torralba. Mit saliency benchmark. 1, 2
- [6] Z. Bylinskii, T. Judd, A. Oliva, A. Torralba, and F. Durand. What do different evaluation metrics tell us about saliency models? *IEEE transactions on pattern analysis and machine intelligence*, 41(3):740–757, 2019. 6
- [7] Z. Bylinskii, N. W. Kim, P. O’Donovan, S. Alsheikh, S. Madan, H. Pfister, F. Durand, B. Russell, and A. Hertzmann. Learning visual importance for graphic designs and data visualizations. In *Proceedings of the 30th Annual ACM Symposium on User Interface Software and Technology*, pages 57–69. ACM, 2017. 2, 4
- [8] Z. Bylinskii, A. Recasens, A. Borji, A. Oliva, A. Torralba, and F. Durand. Where should saliency models look next? In *European Conference on Computer Vision*, pages 809–824. Springer, 2016. 2, 3
- [9] H. Caesar, J. Uijlings, and V. Ferrari. Coco-stuff: Thing and stuff classes in context. In *Computer vision and pattern recognition (CVPR), 2018 IEEE conference on*. IEEE, 2018. 3
- [10] M. Cerf, E. P. Frady, and C. Koch. Faces and text attract gaze independent of the task: Experimental data and computer model. *Journal of vision*, 9(12):10–10, 2009. 3
- [11] J. Chen, G. Bai, S. Liang, and Z. Li. Automatic image cropping: A computational complexity study. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 507–515, 2016. 2
- [12] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and 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(4):834–848, 2018. 4
- [13] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 801–818, 2018. 4
- [14] Y.-L. Chen, T.-W. Huang, K.-H. Chang, Y.-C. Tsai, H.-T. Chen, and B.-Y. Chen. Quantitative analysis of automatic image cropping algorithms:a dataset and comparative study. In *IEEE WACV 2017*, 2017. 7

- [15] M. J. Choi, A. Torralba, and A. S. Willsky. Context models and out-of-context objects. *Pattern Recognition Letters*, 33(7):853–862, 2012. 3
- [16] F. Chollet. Xception: Deep learning with depthwise separable convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1251–1258, 2017. 4
- [17] M. Cornia, L. Baraldi, G. Serra, and R. Cucchiara. Paying more attention to saliency: Image captioning with saliency and context attention. *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, 14(2):48, 2018. 8
- [18] M. Cornia, L. Baraldi, G. Serra, and R. Cucchiara. Predicting human eye fixations via an lstm-based saliency attentive model. *IEEE Transactions on Image Processing*, 27(10):5142–5154, 2018. 1, 2, 4, 5, 7
- [19] R. England. Twitter uses smart cropping to make image previews more interesting. <https://engadget.com/2018/01/25/twitter-uses-smart-cropping-to-make-image-previews-more-interest/>. 7
- [20] S. A. Esmaeili, B. Singh, and L. S. Davis. Fast-at: Fast automatic thumbnail generation using deep neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4622–4630, 2017. 2
- [21] A. Geitgey. Face recognition. http://github.com/ageitgey/face_recognition. 3
- [22] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017. 8
- [23] S. He, H. R. Tavakoli, A. Borji, and N. Pugeault. A synchronized multi-modal attention-caption dataset and analysis. *arXiv preprint arXiv:1903.02499*, 2019. 8
- [24] J. Hu, L. Shen, and G. Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7132–7141, 2018. 5
- [25] X. Huang, C. Shen, X. Boix, and Q. Zhao. Salicon: Reducing the semantic gap in saliency prediction by adapting deep neural networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 262–270, 2015. 1, 2
- [26] S. Jia. Eml-net: An expandable multi-layer network for saliency prediction. *arXiv preprint arXiv:1805.01047*, 2018. 4, 7
- [27] M. Jiang, S. Huang, J. Duan, and Q. Zhao. Salicon: Saliency in context. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1072–1080, 2015. 1, 2, 3
- [28] M. Jiang, J. Xu, and Q. Zhao. Saliency in crowd. In *European Conference on Computer Vision*, pages 17–32. Springer, 2014. 1, 3
- [29] T. Judd, F. Durand, and A. Torralba. Fixations on low-resolution images. *Journal of Vision*, 11(4):14–14, 2011. 3
- [30] A. Khosla, A. S. Raju, A. Torralba, and A. Oliva. Understanding and predicting image memorability at a large scale. In *International Conference on Computer Vision (ICCV)*, 2015. 3
- [31] N. W. Kim, Z. Bylinskii, M. A. Borkin, K. Z. Gajos, A. Oliva, F. Durand, and H. Pfister. Bubbleview: an interface for crowdsourcing image importance maps and tracking 918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

- 972 visual attention. *ACM Transactions on Computer-Human In-* 1026
973 *teraction (TOCHI)*, 24(5):36, 2017. 2 1027
974 [32] K. Kafka, A. Khosla, P. Kellnhofer, H. Kannan, S. Bhan- 1028
975 darkar, W. Matusik, and A. Torralba. Eye tracking for 1029
976 everyone. In *Proceedings of the IEEE conference on computer 1030
977 vision and pattern recognition*, pages 2176–2184, 2016. 2 1031
978 [33] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional 1032
979 networks for semantic segmentation. In *Proceedings of the 1033
980 IEEE conference on computer vision and pattern recogni- 1034
981 tion*, pages 3431–3440, 2015. 4 1035
982 [34] J. Pan, C. Canton, K. McGuinness, N. E. O’Connor, J. Tor- 1036
983 res, E. Sayrol, and X. a. Giro-i Nieto. Salgan: Visual saliency 1037
984 prediction with generative adversarial networks. In *arXiv*, 1038
985 January 2017. 1, 2, 4 1039
986 [35] J. Pan, E. Sayrol, X. Giro-i Nieto, K. McGuinness, and N. E. 1040
987 O’Connor. Shallow and deep convolutional networks for 1041
988 saliency prediction. In *Proceedings of the IEEE Conference 1042
989 on Computer Vision and Pattern Recognition*, pages 598– 1043
990 606, 2016. 7 1044
991 [36] A. Papoutsaki, P. Sangkloy, J. Laskey, N. Daskalova, 1045
992 J. Huang, and J. Hays. Webgazer: Scalable webcam eye 1046
993 tracking using user interactions. In *Proceedings of the 1047
994 Twenty-Fifth International Joint Conference on Artificial 1048
995 Intelligence-IJCAI 2016*, 2016. 2 1049
996 [37] S. Park, W. Kim, and K. M. Lee. Abnormal object detection 1050
997 by canonical scene-based contextual model. In *European 1051
998 Conference on Computer Vision (ECCV)*, 2012. 1, 3 1052
1000 [38] S. J. Rennie, E. Marcheret, Y. Mroueh, J. Ross, and V. Goel. 1053
1001 Self-critical sequence training for image captioning. In *Pro- 1054
1002 ceedings of the IEEE Conference on Computer Vision and 1055
1003 Pattern Recognition*, pages 7008–7024, 2017. 8 1056
1004 [39] D. Rudoy, D. B. Goldman, E. Shechtman, and L. Zelnik- 1057
1005 Manor. Crowd sourcing gaze data collection. *Proceedings of 1058
1006 ACM Collective Intelligence Conference*, 2012. 2 1059
1007 [40] Z. Wei, J. Zhang, X. Shen, Z. Lin, R. Mech, M. Hoai, and 1060
1008 D. Samaras. Good view hunting: learning photo composi- 1061
1009 tion from dense view pairs. In *Proceedings of the IEEE Con- 1062
1010 ference on Computer Vision and Pattern Recognition*, pages 1063
1011 5437–5446, 2018. 2 1064
1012 [41] J. Xu, M. Jiang, S. Wang, M. S. Kankanhalli, and Q. Zhao. 1065
1013 Predicting human gaze beyond pixels. *Journal of vision*, 1066
1014 14(1):28–28, 2014. 2 1067
1015 [42] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhudinov, 1068
1016 R. Zemel, and Y. Bengio. Show, attend and tell: Neural 1069
1017 image caption generation with visual attention. In *Inter- 1070
1018 national conference on machine learning*, pages 2048–2057, 1071
1019 2015. 8 1072
1020 [43] B. Yao, X. Jiang, A. Khosla, A. L. Lin, L. Guibas, and L. Fei- 1073
1021 Fei. Human action recognition by learning bases of action 1074
1022 attributes and parts. In *2011 International Conference on 1075
1023 Computer Vision*, pages 1331–1338. IEEE, 2011. 3 1076
1024
1025