# Are All the Frames Equally Important?

**Oleksii Sidorov**
**Marius Pedersen**
The Norwegian Colour and
Visual Computing Laboratory,
NTNU
Gjøvik, Norway
oleksiis@stud.ntnu.no
marius.pedersen@ntnu.no

**Nam Wook Kim**
Harvard University
Cambridge, USA
namwkim@seas.harvard.edu

**Sumit Shekhar**
Adobe Research
San Jose, USA
sushekha@adobe.com

## Abstract

In this work, we address the problem of measuring and pre-
dicting temporal video saliency - a metric which defines the
importance of a video frame for human attention. Unlike the
conventional spatial saliency which defines the location of
the salient regions within a frame (as it is done for still im-
ages), temporal saliency considers importance of a frame
as a whole and may not exist apart from context.

The proposed interface is an interactive cursor-based al-
gorithm for collecting experimental data about temporal
saliency. We collect the first human responses and perform
their analysis. As a result, we show that qualitatively, the
produced scores have very explicit meaning of the semantic
changes in a frame, while quantitatively being highly corre-
lated between all the observers.

Apart from that, we show that the proposed tool can simul-
taneously collect fixations similar to the ones produced
by eye-tracker in a more affordable way. Further, this ap-
proach may be used for creation of first temporal saliency
datasets which will allow training computational predic-
tive algorithms. The proposed interface does not rely on
any special equipment, which allows to run it remotely and
cover a wide audience.

## Author Keywords
Attention; video; saliency; temporal saliency; eye-tracking

## CCS Concepts

•**Information systems** → **Multimedia information systems;** •**Human-centered computing** → **Human computer interaction (HCI);** *User interface toolkits;*

## Introduction

It seems obvious that some fragments of a video are more important than others. Such fragments concentrate most of the viewer's attention while others remain of no interest. The naïve examples are: a culmination scene in a movie, a screamer in a horror film, the moment of an explosion, or even a slight motion in very calm footage. We denote such fragments as groups of frames with high *temporal saliency*. Information about temporal saliency is an essential part of a video characterization which gives valuable insights about the video structure. Such information is directly applicable in video compression (frames which do not attract attention may be compressed more), video summarization (salient frames contain the most of perceived video content), indexing, memorability prediction, and others tasks. So, the reader may expect that there is a big number of algorithms and techniques aimed at measuring and predicting temporal saliency. However, this is not the case. The most, if not all, of the well-known works on video saliency are aimed at spatial saliency, *i.e.*, a prediction of spatial distribution of the observer's attention across the frame (in a similar way as if it was an individual image). We hypothesize that this is due to the absence of established methodology for measuring temporal saliency in the experiment, which is crucial for obtaining ground truth data. Conventionally, saliency data are collected using eye-tracking, which is a technique that produces a continuous temporal signal. In other words, it does not allow to differentiate between the frames as a whole, because each frame produces the same kind of output – a pair of gaze fixation coordinates with a rate defined by hardware.

In this work, we propose a new methodology for measuring temporal video saliency in the experiment – the first, to the best of our knowledge, method of this kind. For this, we develop a special interface based on mouse-contingent moving-window approach for measuring saliency maps of static images. We also show that it can simultaneously gather meaningful spatial information which can serve as an approximation of gaze fixations.

During the experiment, observers are presented with repeated blurry video-sequences which they can partially deblur using mouse click (Fig. 1). Users can deblur a circular region with a center at cursor location which approximates the confined area of focus in the human eye fovea surrounded by a blurred periphery [3]. Since the number of clicks is limited - observers are forced to use clicks only on most "interesting" frames which attract their attention. Statistical analysis of the collected clicks allows to assign the corresponding level of importance to each frame. This information can be applied directly in numerous tasks of video processing.

To summarize, unlike the conventional approaches which only try to understand *where* the observer looks, we also study *when* the observer pay the most attention.

## Related works

The straightforward method of retrieving the information about attention is based on the utilization of commercial eye-trackers (*e.g.* EyeLink, Tobii). Hardware-based eye-tracking has been used widely in various studies on human-computer interaction [6][13]. A less accurate, but much more affordable, way of measuring saliency is based on measuring the mouse cursor position which was proven to correlate strongly with gaze fixations [4][5][15]. The most successful algorithms of this type utilize a moving-window paradigm, which masks information outside of the area adjacent to the cursor and requires a user to move the cur-

sor (followed by a window around it) to make other regions visible. Such algorithms include Restricted Focus Viewer software by Jansen *et al.* [7] and more recent SALICON [8] and BubbleView [10]. These algorithms were also used in large online crowdsourcing experiments due to the native scalability of cursor-based approaches. However, they were studied only in the context of spatial saliency of static images. This is fair for static images, but for video-sequences, temporal information is commonly even more important than spatial regions. Furthermore, there are no well-known experimental datasets which can provide this kind of information[1] and be used for training of computational algorithms. For example, the popular video saliency datasets Hollywood-2 [16], UCF sports [12], SAVAM [2], DHF1K [17] only provide eye-tracking results which are constant in the temporal domain.

## Methodology

Our approach is inspired by moving-window gaze approximations methods for still images. In the proposed setup all video frames are blurred. Clicking the mouse deblurs a round window around the cursor. Users are demonstrated repeated video sequence during which they can click the mouse for short periods of time. The total number of times when the frame was deblurred defines temporal saliency score, while location of the cursor when the mouse button is pressed approximates gaze fixation location and allows to detect what caused the interest.

*Discretization*
Short fragment of a video is more likely to attract user's attention rather than a single frame, so we let the users keep the mouse button pressed instead of clicking on each frame they find interesting. However, when not forced explicitly,

observers tend to keep the mouse button pressed all the time, which is natural. Thus, to obtain variation of scores, it is crucial to restrict users artificially. Our solution is to simply limit the amount of deblurred frames (time period), after which clicking the mouse button stops working, and additionally limit the amount of deblurred frames per one continuous click. The users cannot see the limits, instead, they learn them during a test trial and then follow them intuitively. For example, a 10-second video may have up to 4 seconds of deblurred frames, but no more than 1 second at once. In the result, a user can make 4 long clicks 1 second each or a larger number of short clicks, while we are guaranteed to have at least four discrete responses after one run.

*Repetition*
The idea of repeating the videos may be used to gather more responses from one observer and have richer statistics. Moreover, if a salient event happens at the end, the observer may reach the limit before seeing it, so it is necessary to make a second round. Also, eye-motion and cognitive processing are faster than clicking the mouse, so giving the user an opportunity to predict when an event will happen is beneficial for the creation of more accurate saliency maps with a shorter delay. However, we observed that in the majority of the cases, the first run is the most informative one, and the user is able to detect most salient information without preparation. *Subsequent repeats lead to shifting the user's attention to smaller details.* Eventually, we used repetition in our experiments, but analyze different numbers of repeats in results.

*Other parameters*
Other important parameters are the blur radius and the radius of the window. Their definition requires more detailed study. The task given to an observer also influences where they look [18][10], so, this parameter depends on the par-
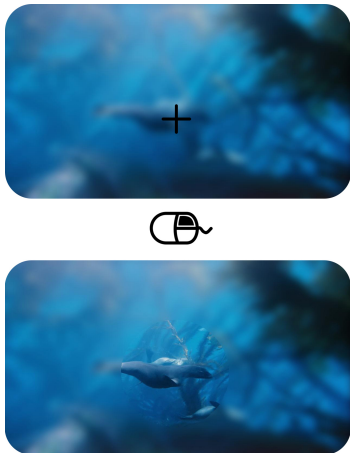


**Figure 1:** The proposed interface. A more representative video demonstration is available online: [link].

---

[1] A comprehensive list of saliency datasets: http://saliency.mit.edu/datasets.html

ticular context in which the experiment is performed. In our case, we are interested in basic watching of a video without a particular task, so we worked under a "free-view" setup.

## Experimental setup

The experiments were performed offline using a special setup in the laboratory (Fig. 2) for the sake of fully-controlled conditions (in future we are also planning to run the experiment on Amazon Mechanical Turk for gathering larger database, which would be impossible to do with an eye-tracker). The display used is 24.1" EIZO ColorEdge CG241W color-calibrated with X-Rite Eye-One Pro. The distance between the display and the observer was 50 cm.

The code is written in MatLab with Psychtoolbox-3 [11] and is publicly available by the link [2].

Videos with ground-truth eye-tracking data were taken from SAVAM dataset [2] due to their high quality and diverse content. We used eight 10-seconds long HD videos including two test videos. The content of the videos is diverse and includes: a basketball game with a score moment, a calm shot of leaves in the wind, marine animals underwater, a cinematic scene of a child coming home, a surveillance camera footage of two men meeting, a suffocating diver emerging from the water.

Interface parameters: radius of a circular window – 200 px ($6.2°$ visual angle), blur kernel – Gaussian with standard deviation of 15, video duration – 10 s, limit of deblurred frames per one round – 4 s (100 frames), limit of deblurred frames at one click – 1 s (25 frames), number of repetitions – 5, frame-rate of the videos – 25 fps, video resolution – 1280 px $\times$ 720 px ($38.2° \times 22°$ visual angle), videos are silent.

The observers were invited from the University staff and students. 30 subjects in total, 15 women and 15 men. Age: 21-42 (mean 25.6).



**Figure 2:** Experimental setup (the light is off during the session).
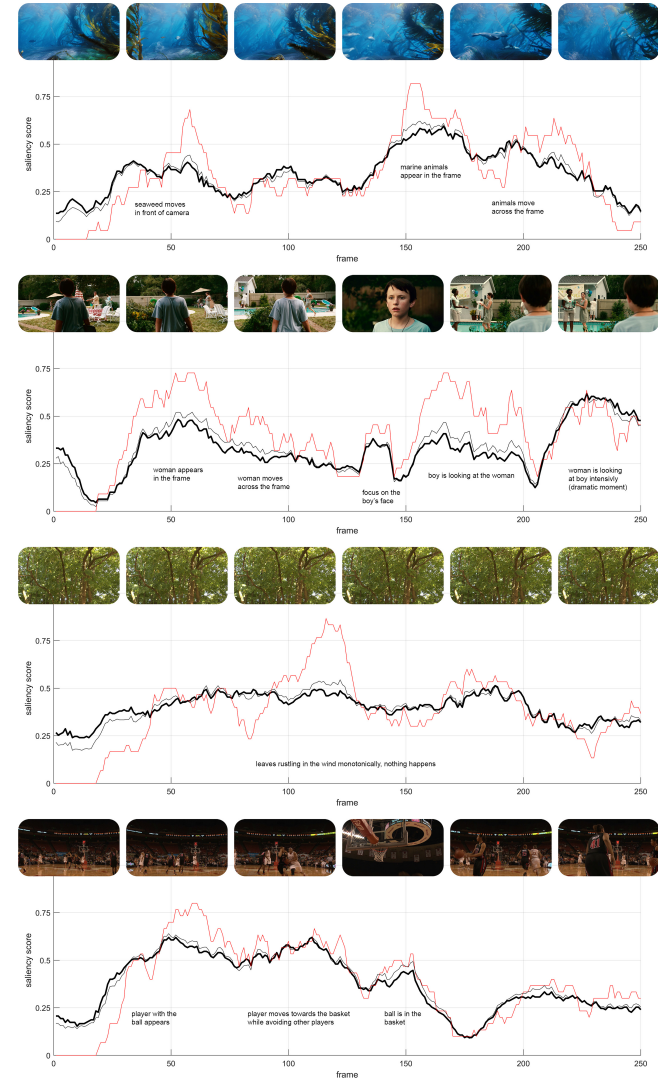
[2]https://github.com/acecreamu/temporal-saliency



**Figure 3:** The produced temporal saliency graphs. Thick black line $C_{1-5}$, red line $C_1$, thin black line $C_{1-5}^{(W)}$. Zoom is required.

| | Pearson Correlation Coefficient (mean $_{std}$) | | | | Kolmogorov-Smirnov test (mean $p$-value) | | | |
|---|---|---|---|---|---|---|---|---|
| | $C_1$ | $C_{1-2}$ | $C_{1-5}$ | $C_{1-5}^{(W)}$ | $C_1$ | $C_{1-2}$ | $C_{1-5}$ | $C_{1-5}^{(W)}$ |
| "The underwater world" | $0.663_{\,0.082}$ | $0.694_{\,0.082}$ | $0.740_{\,0.074}$ | $0.770_{\,0.064}$ | 0.119 | 0.048 | 0.011 | 0.036 |
| "Cinematic scene" | $0.615_{\,0.092}$ | $0.711_{\,0.057}$ | $0.803_{\,0.053}$ | $0.789_{\,0.051}$ | 0.164 | 0.107 | 0.033 | 0.067 |
| "Leaves in the wind" | $0.694_{\,0.068}$ | $0.563_{\,0.099}$ | $0.545_{\,0.108}$ | $0.647_{\,0.092}$ | 0.081 | 0.073 | 0.044 | 0.057 |
| "Basketball game" | $0.741_{\,0.072}$ | $0.766_{\,0.070}$ | $0.863_{\,0.050}$ | $0.845_{\,0.051}$ | 0.164 | 0.099 | 0.055 | 0.063 |
| "Diver suffocating" | $0.789_{\,0.050}$ | $0.788_{\,0.054}$ | $0.820_{\,0.057}$ | $0.834_{\,0.051}$ | 0.134 | 0.092 | 0.043 | 0.068 |
| "Meeting of the two" | $0.660_{\,0.089}$ | $0.701_{\,0.085}$ | $0.740_{\,0.069}$ | $0.753_{\,0.069}$ | 0.121 | 0.112 | 0.061 | 0.053 |

**Table 1:** Inter-observer consistency of the measured temporal saliency maps. $C_{1-N}$ denotes sum of $N$ rounds used for computation.



**Figure 4:** The comparison of spatial saliency maps. Top row in each pair – eye-tracking results, bottom – our results. Zoom is required.

## Results and discussion

The proposed interface allows measuring both temporal and spatial saliency at the same time, thus, we evaluate the accuracy of both these outputs.

*Temporal saliency results*
Considering that there are no ground truth temporal saliency data, we evaluate the output of the algorithm by analyzing the produced temporal saliency "maps" and estimating inter-observer consistency. The examples of obtained temporal saliency "maps" are illustrated in Fig. 3. The demonstration of the videos with saliency scores encoded as a color-map is available online: [link]. Figure 3 demonstrates three plots for each video which correspond to different averaging approaches: the sum of all clicks from all five video repeats ($C_{1-5}$); the sum of clicks only from the first round without repeating ($C_1$); and the *weighted* sum of clicks from 5 rounds ($C_{1-5}^{(W)} = \sum_{n=1}^{5} C_n W_n$, where $W_n = \{1, 0.8, 0.6, 0.4, 0.2\}$). All the scores are normalized by a maximum number of clicks the frame can have. Qualitative analysis shows that *most of the peaks on the temporal saliency graph correspond to the semantically meaningful salient events on the video*. This is the main achievement of the proposed interface. It can also be seen that an intentionally taken monotonic video without salient events ("leaves in the wind") has relatively flat saliency

graph without strongly pronounced peaks (which could be even flatter when the response statistics is larger). Apart from that, it may be seen that in the case of other videos, the output of the first round (red line) is very similar to the total output of all five rounds. This means that even when in the next rounds observers start exploring smaller, less salient details, they still return to the "main" events and follow a similar pattern of clicks as in the first round. Also, adding weights to the sum (thin black line) does not influence the results significantly, which again indicates the similarity of clicks from all the rounds. However, using $N$ rounds indeed allows to gather $N$ times more responses making the graph smoother and, as we show next, produces more consistent responses from each observer.

In order to estimate consistency between different groups of observers, we synthetically split observers into two groups of 15 people each. Then, we compute temporal saliency maps for each group independently and compare the results. The comparison is done using the Pearson Correlation Coefficient between the saliency maps from different groups, as well as performing the Kolmogorov-Smirnov test between two distributions and reporting the p-value. Results are averaged between 100 random splits (standard deviation is also reported for PCC). Table 1 shows that the correlation between responses from different observers is very high, up to 0.86. Increasing the number of rounds considered increases the correlation of responses significantly, with maximum values achieved when all five rounds are included.

*Spatial saliency results*

The spatial saliency maps produced by eye-tracking data versus our interface can be compared visually in Fig. 4. (fixation points are blurred with a Gaussian of sigma equal to $1°$ of visual angle (33 px)). As may be seen, the results are very similar, even though we did not use any special

|  | AUC (mean $_{std}$) | NSS (mean $_{std}$) |
|---|---|---|
| "The underwater world" | 0.617 $_{0.108}$ | 0.73 $_{0.78}$ |
| "Cinematic scene" | 0.712 $_{0.119}$ | 1.59 $_{1.05}$ |
| "Leaves in the wind" | 0.548 $_{0.055}$ | 0.18 $_{0.21}$ |
| "Basketball game" | 0.727 $_{0.114}$ | 1.52 $_{0.93}$ |
| "Diver suffocating" | 0.794 $_{0.113}$ | 2.66 $_{1.41}$ |
| "Meeting of the two" | 0.625 $_{0.060}$ | 0.95 $_{0.43}$ |

Table 2: Comparison of the measured spatial saliency maps and gaze-fixations obtained using eye-tracker.

equipment and collected spatial data additionally to the main temporal output.

Saliency maps are evaluated quantitatively using standard saliency metrics: Area under ROC Curve (AUC) [9][1] and Normalized Scanpath Saliency (NSS) [14]. Table 2 presents statistics of the scores computed per frame. Results demonstrate both good and poor performance, and differ significantly from video to video. Additionally, quality of spatial saliency can be assessed visually via the rendered videos with map overlay [link], as well as the videos with both eye-tracking (blue dots) and our results (red dots) simultaneously [link].

## Conclusions

In this work, we presented a novel mouse-contingent interface designed for measuring temporal and spatial video saliency. Temporal saliency is a novel concept which is studied incongruously less than it should in comparison to spatial saliency. Temporal video saliency allows identifying the important fragments of a video by assigning a saliency score to each frame. The analysis of the experimental study shows that the use of the proposed interface allows to accurately approximate the temporal saliency "map" as well as gaze-fixations of the observers at the same time.

## REFERENCES

[1] Ali Borji, Hamed R Tavakoli, Dicky N Sihite, and Laurent Itti. 2013. Analysis of scores, datasets, and models in visual saliency prediction. In *IEEE ICCV*. 921–928.

[2] Yury Gitman, Mikhail Erofeev, Dmitriy Vatolin, Andrey Bolshakov, and Alexey Fedorov. 2014. Semiautomatic Visual-Attention Modeling and Its Application to Video Compression. In *IEEE ICIP*. Paris, France, 1105–1109.

[3] Frédéric Gosselin and Philippe G Schyns. 2001. Bubbles: a technique to reveal the use of information in recognition tasks. *Vision research* 41, 17 (2001), 2261–2271.

[4] Qi Guo and Eugene Agichtein. 2010. Towards predicting web searcher gaze position from mouse movements. In *CHI'10 Extended Abstracts*. ACM, 3601–3606.

[5] Jeff Huang, Ryen White, and Georg Buscher. 2012. User see, user point: gaze and cursor alignment in web search. In *SIGCHI*. ACM, 1341–1350.

[6] Robert J. K. Jacob and Keith S. Karn. 2003. Eye Tracking in Human-Computer Interaction and Usability Research: Ready to Deliver the Promises. *Mind* 2, 3 (2003), 4.

[7] Anthony R Jansen, Alan F Blackwell, and KIM Marriott. 2003. A tool for tracking visual attention: The restricted focus viewer. *Behavior research methods, instruments, & computers* 35, 1 (2003), 57–69.

[8] Ming Jiang, Shengsheng Huang, Juanyong Duan, and Qi Zhao. 2015. Salicon: Saliency in context. In *IEEE CVPR*. 1072–1080.

[9] Tilke Judd, Krista Ehinger, Frédo Durand, and Antonio Torralba. 2009. Learning to predict where humans look. In *IEEE ICCV*. IEEE, 2106–2113.

[10] Nam Wook Kim, Zoya Bylinskii, Michelle A Borkin, Krzysztof Z Gajos, Aude Oliva, Fredo Durand, and Hanspeter Pfister. 2017. BubbleView: an interface for crowdsourcing image importance maps and tracking visual attention. *ACM TOCHI* 24, 5 (2017), 36.

[11] Mario Kleiner, David Brainard, Denis Pelli, Allen Ingling, Richard Murray, Christopher Broussard, and others. 2007. WhatâĂŹs new in Psychtoolbox-3. (2007).

[12] Stefan Mathe and Cristian Sminchisescu. 2015. Actions in the eye: Dynamic gaze datasets and learnt saliency models for visual recognition. *IEEE TPAMI* 37, 7 (2015), 1408–1424.

[13] Jakob Nielsen and Kara Pernice. 2010. *Eyetracking web usability*. New Riders.

[14] Robert J Peters, Asha Iyer, Laurent Itti, and Christof Koch. 2005. Components of bottom-up gaze allocation in natural images. *Vision research* 45, 18 (2005), 2397–2416.

[15] Kerry Rodden, Xin Fu, Anne Aula, and Ian Spiro. 2008. Eye-mouse Coordination Patterns on Web Search Results Pages. In *CHI'08 Extended Abstracts (CHI EA '08)*. ACM, New York, NY, USA, 2997–3002.

[16] E. Vig, M. Dorr, and D. Cox. 2014. Large-Scale Optimization of Hierarchical Features for Saliency Prediction in Natural Images. In *IEEE CVPR*. 2798–2805.

[17] Wenguan Wang, Jianbing Shen, Fang Guo, Ming-Ming Cheng, and Ali Borji. 2018. Revisiting video saliency: A large-scale benchmark and a new model. In *IEEE CVPR*. 4894–4903.

[18] Alfred L Yarbus. 1967. *Eye movements and vision*. Plenum, New York, NY, USA.