



HAL
open science

ChaboNet : Design of a deep CNN for prediction of visual saliency in natural video

Souad Chaabouni, Jenny Benois-Pineau, Chokri Ben Amar

► To cite this version:

Souad Chaabouni, Jenny Benois-Pineau, Chokri Ben Amar. ChaboNet : Design of a deep CNN for prediction of visual saliency in natural video. *Journal of Visual Communication and Image Representation*, 2019, 60, pp.79-93. <10.1016/j.jvcir.2019.02.004>. <hal-02326279>

HAL Id: hal-02326279

<https://hal.science/hal-02326279v1>

Submitted on 22 Oct 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire HAL, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons CC BY-NC 4.0 - Attribution - Non-commercial use - International License

Elsevier Editorial System(tm) for Journal of
Visual Communication and Image Representation
Manuscript Draft

Manuscript Number: JVCI-18-30R2

Title: ChaboNet : Design of a Deep CNN for Prediction of Visual Saliency
in Natural Video

Article Type: Research paper

Keywords: Visual attention prediction in video; Deep Convolutional Neural
Networks; saliency map; residual motion; dynamic content

Corresponding Author: Dr. souad chaabouni, Ph.D

Corresponding Author's Institution: Laboratoire bordelais d'informatique

First Author: souad chaabouni, Ph.D

Order of Authors: souad chaabouni, Ph.D; Jenny Benois-Pineau, Professor;
Chokri Ben Amar, Professor

Abstract: Prediction of visual saliency in images and video is needed for
video understanding, coding and other applications.

The majority of prediction models are founded only on "bottom-up"
features.

Nevertheless, the "top-down" component of human visual attention becomes
prevalent as human observers explore the visual scene. Visual saliency
which is always a mix of bottom-up and top-down cues can be predicted on
the basis of

seen data. In this paper, a model of prediction of visual saliency in
video on the basis of Deep convolutional neural networks is proposed. A
Deep CNN architecture is designed. Various input channels for a CNN
architecture

are studied: using the known sensitivity of human visual system to
residual motion, pixel colour values are completed with residual motion
map. The experiments show that the choice of the input features for the
Deep CNN depends on visual task.

Souad CHAABOUNI
LaBRI UMR 5800
351 crs de la Liberation F33405, France
Jenny BENOI
LaBRI UMR 5800
351 crs de la Liberation F33405, France
Chokri BEN AMAR
REGIM-Lab LR11ES48 ENIS BP1173, 3038 Sfax, Tunisie

12/01/2018

Dear ,

We wish to submit a new manuscript entitled "Chabonet: Design of a Deep CNN for Prediction of visual saliency in Natural Video" for consideration by the Journal of Visual Communication and Image Representation .We confirm tat this work is original. Some preliminary results were published in conference proceedings which we reference. The paper is not currently under consideration for publication elsewhere.

In this paper, we report on prediction of visual saliency using deep convolution network in natural video. The paper should be of interest to readers specialized in computer vision, visual saliency prediction and Deep learning.The novelty of the work is in the use of specific input layers in video, a transfer learning scheme and a good benchmark with classical saliency prediction models widely used by the community. Please address all correspondence concerning this manuscript tony self at :

chaabounisouad@yahoo.fr

Thank you for your consideration of this manuscript.

Yours Sincerely,

Souad CHAABOUNI

on behalf of all co-authors

We would like to thank the reviewers for their time spent evaluating our work and for their interest remarks which help us to better expose our paper.

<p>Reviewer #2: Thank you for revising the paper by considering our comments. I think that the revised manuscript is better than the previous submission. However, there still exist several other major concerns.</p>	
<p>1 The authors claim that their first contribution is that the architecture is benchmarked against known architectures AlexNet and LeNet. I think this cannot be considered as a novel contribution. LeNet is proposed in 1994, and AlexNet is in 2012. Nowadays many architectures (e.g. VGG, ResNet) outperforms these two.</p>	<p>We have proposed a specific architecture in a family of « light » deep architectures such as AlexNet. Its particularities consist in the re-enforcement of features, as human visual system is sensitive to contrasts.</p> <p>We stress that we remain in the framework of such a «light » strategy in the section 3, see : « In this section, the proposed architecture ChaboNet for the visual saliency prediction is presented. This is a relatively "light" deep architecture compared to havier popular GoogleNet ... »</p> <p>As our experiences show, the same architecture is applicable to very different datasets, subject to an adequate training and initialization, see section 4 « transfer learning ».</p> <p>Obviously, the architecture can be adapted to the new optimization strategies such as proposed in ResNet.</p>
<p>2 Fonts in Table 11, 12, 13 and 14 are too small.</p>	<p>The tables 11, 12, 13 and 14 are presented vertically, please see the page 31 .. 35</p>
<p>3 There are some typos in this paper, a thorough proofreading is needed. Here we list some examples in the following sentences. (1) On Page 5, in line 121: to avfoid overfitting -> to avoid overfitting; (2) On Page 26, in Table 9: %78.73+-0.930 -> 78.73%+-0.930</p>	<p>Typos corrected</p>

1
2
3
4
5
6
7
8
9 ChaboNet : Design of a Deep CNN for Prediction of
10 Visual Saliency in Natural Video
11

12
13 Souad CHAABOUNI^{a,b}, Jenny BENOIS-PINEAU^a, Chokri BEN AMAR^b
14

15 ^aLaBRI UMR 5800, Univ. Bordeaux, 33400 France

16 ^bDepartment of Computer science, REGIM-Lab LR11ES48, Univ. Sfax, 3038 Tunisia
17
18
19

20
21 **Abstract**

22
23 Prediction of visual saliency in images and video is needed for video under-
24 standing, search and retrieval, coding, watermarking and other applications.
25 The majority of prediction models are founded only on “bottom-up” features.
26 Nevertheless, the “top-down” component of human visual attention becomes
27 prevalent as human observers explore the visual scene. Visual saliency which is
28 always a mix of bottom-up and top-down cues can be predicted on the basis of
29 seen data. In this paper, a model of prediction of visual saliency in video on
30 the basis of Deep convolutional neural networks (CNNs) is proposed. A Deep
31 CNN architecture is designed. Various input channels for a CNN architecture
32 are studied: using the known sensitivity of human visual system to residual
33 motion, pixel colour values are completed with residual motion map. The latter
34 is a normalized energy of residual motion in video frames with regard to the
35 estimated global affine motion model. The experiments show that the choice
36 of the input features for the Deep CNN depends on visual task: for highly dy-
37 namic content, the proposed model with residual motion is more efficient and
38 gives decent results with relatively shallow Deep architecture.
39

40
41
42
43
44
45
46
47
48 *Keywords:* Visual attention prediction in video, Deep Convolutional Neural
49 Networks, saliency map, residual motion, dynamic content.
50

51
52
53
54 *URL:* souad.chaabouni@u-bordeaux.fr (Souad CHAABOUNI),
55 jenny.benois-pineau@u-bordeaux.fr (Jenny BENOIS-PINEAU),
56 chokri.benamar@ieee.org (Chokri BEN AMAR)
57

1. Introduction and related work

Prediction of visual saliency in images and video is an intensively researched topic. With the growing volumes of digital video it is of highest interest for variety of applications involving video understanding, coding, watermarking... Prediction of pixel-wise saliency means assigning to each pixel in the image plane a measure characterizing the attraction of this image locus for a human observer. Several visual low-level characteristics: luminance, color, orientation and movement provoke human gaze attraction when observing visual content. This is why a very large amount of research works was devoted to prediction of visual saliency in images and video on the basis of popular “feature integration theory” [1]. These models simulate stimuli-driven, or “bottom-up” attention. [2] proposed a saliency detection algorithm for panoramic landscape images of outdoor scenes. Hence, the background of a panoramic image was estimated using the characteristics of geodesic similarity on a graph and the spatial distribution of homogeneous background regions. In [3], the saliency was considered as a ranking loss function which is designed to rank saliency values in the descending order of their relevance. In [4], the similarity computed between the input image and each similar image is measured and used for computing adaptive fusion weights for multiple saliency maps fusion.

For saliency detection in videos, local features and global features were extracted in [5] to create a pixel-level temporal and spatial saliency map. The authors of [6] defined the primary salient object in a video using the integration of the local visual/motion saliency, the global appearance consistency, and spatio-temporal smoothness constraint on object trajectories. In [7] the properties of saliency, low-rank, connectivity and sparsity were integrated into an unified objective function to detect a moving object using the saliency map. Models that use motion features are reported in [8]. The model of Wang [9] computes the gradient flow field and energy optimization using intra-frame boundary information and the inter-frame information to build a consistent spatio-temporal

1
2
3
4
5
6
7
8
9 saliency. In [10], a video saliency model was proposed in order to detect the
10 attended regions. The latter correspond to both interesting objects and domi-
11 nant motions. The so-called “top-down” attention is driven by observation task.
12 In a task-driven visual search users search for particular objects and the goal
13 of saliency models is to predict them in the image plane such as in [11]. In a
14 free observation process of unknown video content, top-down visual attention is
15 triggered when the observer understands the scene and selects the targets to fol-
16 low. It becomes prevalent [12] when the human subject observes visual content
17 with progressively increasing observation time, which is the case in continuous
18 video scenes [13], [14]. Therefore, the areas of interest for the user cannot be
19 predicted sufficiently well by purely bottom-up models. Various new cues have
20 been proposed to enhance bottom-up models from the “top-down” perspective
21 [12]. Boujut et al. [15] proposed a bottom-up spatio-temporal saliency model,
22 enhanced by considering “face detection” as semantics in the video. The ELD
23 [16] deep saliency model extracts high level features, using the deep network
24 VGG-net, and low level features (Average RGB value, Gabor filter, ...). The
25 concatenation of both encoded features, using a fully connected neural network,
26 generates the final saliency map. The above mentioned research directly or not
27 introduces elements of top-down visual attention prediction in model building.
28 From methodological point of view frequent are the attempts to learn the in-
29 terest of users with regard to elements of visual scenes whatever the content
30 is when deploying supervised machine learning approaches [17]. Using super-
31 vised learning in the field of saliency prediction generates classifiers that can
32 predict the focus of attention on the basis of already seen data thus combin-
33 ing bottom-up and top-down visual cues. Here, deep learning has emerged as
34 an active research trend. It involves learning, at multiple levels of abstraction,
35 for mining data such as images, sound, and text [18]. In addition to the com-
36 mon design of the classifier architectures on the basis of neural networks, deep
37 learning presents a philosophy to model the complex relationships between data
38 [19], [20]. Generally, deep neural networks are multi-layer predictive networks
39 formed to maximize the probability of input data with regard to target classes
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9 [21]. Due to the increased computational capacities with Graphical Processing
10 Units (GPU) deep neural networks have outperformed all previous classification
11 models in the tasks of video understanding. Deep convolutional neural net-
12 works (CNNs) were developed in Computer Vision, first by Yann LeCun with
13
14 65 works (CNNs) were developed in Computer Vision, first by Yann LeCun with
15 the LeNet [19] architecture that was used to recognize digits. Then, AlexNet
16 [22] network has become very popular as the architecture for visual recognition
17 tasks. Now, deep learning architectures, which have recently been proposed for
18 the prediction of salient areas in images, differ essentially by the quantity of
19 convolution and pooling layers, the input data, pooling strategies, the nature
20 of the final classifiers, the loss functions to optimize and the formulation of the
21 70 convolution and pooling layers, the input data, pooling strategies, the nature
22 of the final classifiers, the loss functions to optimize and the formulation of the
23 problem. Recently a fully convolutional network “FCN” for saliency prediction
24 was proposed. This model estimates a dense saliency map of a given image using
25 a set of extreme learners, each trained on an image similar to the input image
26 [23]. A lot of works today, are devoted to saliency prediction in still images
27 using “FCN”: In [24], a global scene information that was trained on diverse
28 categories of an eye-tracking data set, was used in addition to local information.
29 [25] present an architectural extension to any CNN to fine-tune traditional 2D
30 75 [23]. A lot of works today, are devoted to saliency prediction in still images
31 using “FCN”: In [24], a global scene information that was trained on diverse
32 categories of an eye-tracking data set, was used in addition to local information.
33 [25] present an architectural extension to any CNN to fine-tune traditional 2D
34 saliency prediction to Omni-directional Images. In [26], the authors proposed a
35 deep CNN that predicts eye fixations and segments salient objects. [27] reuses
36 an existing neural network trained on the task of object recognition to predict
37 80 deep CNN that predicts eye fixations and segments salient objects. [27] reuses
38 an existing neural network trained on the task of object recognition to predict
39 eye fixations. [28] formulated the prediction of eye fixations as a minimization
40 of a loss function that measures the Euclidean distance of the predicted saliency
41 map with the provided ground truth. Despite the popularity of these models,
42 they still need a thorough study in real-life situations.
43
44 85 they still need a thorough study in real-life situations.

47 Prediction of visual attention in images reveals the binary classification prob-
48 lem for areas in images as “salient” and “non-salient”. It corresponds to the vi-
49 sual experiment in free viewing conditions, when the subjects are simply asked to
50 look at the content without any specific visual task. In [29], firstly, the learning
51 of the relevant characteristics of the saliency of natural images was performed,
52 90 of the relevant characteristics of the saliency of natural images was performed,
53 and secondly the eye fixations on objects with semantic content was predicted.
54 In Simonyan’s work [30], the subjects are asked to look for an object from a
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9 given taxonomy in the images. Therefore, the classification problem is multi-
10 class, and can be expressed as a task-dependent visual experiment. In [31], first,
11
12
13 95 a random bank of uniform filters is used to generate multiple representations of
14 input images. The second phase provides the combination of different localized
15 representations. While a significant effort has been already made for building
16 saliency prediction models from still images with deep learning approach, very
17 few have been built for video content with it [32]. Video has a supplementary
18
19
20 100 dimension: the temporality that is expressed by apparent motion in the image
21 plane. In such early works several combinations of primary features (input)
22 such as color values and residual motion [33] are used to feed a CNN.
23

24
25 The intelligent application of Deep NNs to the well-studied problem of pre-
26 diction of visual saliency in images and video requires deep understanding of
27
28
29 105 their inherent weaknesses. The first one is the need of a very large amount of
30 training data for building a well-generalizing model. Next, they are sensitive
31 to the noise in training data. The noise is quite common when automatic and
32 manual annotation of a large amount of data is required.
33

34
35 In the present paper we are interested in the prediction of saliency maps,
36
37 110 which means the degree of interest for the observer in each pixel position. We
38 do not consider the problem of a “scanpath” prediction [34], that is dynamics
39 of gaze saccades across image plane. Hence, the Deep CNNs turn to be a right
40 tool for such a map prediction. In this study, an approach with Deep CNN
41 that ensures learning salient areas in order to predict pixel-wise saliency maps
42
43
44 115 in video is proposed. We systematize our early works [32] and go further in
45 studying our network architecture for the saliency detection problem. A specific
46 attention is paid to the input layer of the proposed architecture, i.e. the data
47 we extract from video to feed the network.
48
49

50
51 We stress that we remain with relatively “light” deep architecture, compa-
52
53 120 rable to AlexNet. This is why it is benchmarked against known “light” ar-
54 chitectures AlexNet and LeNet. This is our first contribution. The second
55 contribution consists in the design of a transfer learning scheme to avoid over-
56 fitting on real-life small saliency data sets. It differs from usual approaches of
57
58

1
2
3
4
5
6
7
8
9 transfer learning in Deep CNNs classifiers, where the majority of researches use
10 the pre-trained models on ImageNet whatever the target domain and classifica-
11 125 tion problem is accordingly to the method first proposed in [21]. The method
12 is benchmarked with Bengio’s method [21].
13
14

15 Despite the appearance of fully-convolutional (FC) deep networks for saliency
16 prediction [35] which allow prediction of dense saliency maps, our method re-
17 mains ”sparse”. This means that we i)predict saliency of regularly sampled
18 130 patches in video frames with a classical Convolutional Network architecture,
19 ii)and then densify the map by interpolation. Both FC-nets and our scheme
20 require interpolation. Nevertheless, we consider that in our case, it is easier to
21 select training data which will not be corrupted by distractors and changes of
22 the focus-of-attention along the time in video. We explain this choice in section
23 135 2.
24
25
26
27
28

29 The paper is organized as follows. In section 2 data selection method is
30 described, in section 3 the designed CNN architecture is presented, section 4
31 is devoted to the proposed transfer-learning scheme. Pixel-wise computation
32 of predicted visual attention/saliency maps is then introduced in section 5. In
33 140 section 6, results and comparison with classical saliency prediction methods are
34 presented. Section 7 concludes the paper and outlines its perspectives.
35
36
37
38
39

40 2. Policy of data set creation: salient and non-salient patches

41
42 To train any CNN model selection of a training data set which would contain
43 as less noise as possible is the must. At the first step of our method : prediction
44 145 of saliency of patches in the video frames, classification problem is two-class.
45 The training set has to comprise salient and Non-salient regions in video frames.
46 The ground-truth for saliency here are the Gaze Fixation Density Maps (GFDM)
47 [36]. They are built upon gaze fixations of a cohort of subjects recorded during
48 a psycho-visual experiment.
49 150
50
51
52
53

54 For salient patches extraction the intuition is clear: we need to extract
55 patches in the video frames where the GFDM has strong values. For Non-salient
56
57
58

patches extraction, the situation is more complex. First, due to the distractors
 and visual fatigue, the areas in a given video frame which are salient can become
 Non-salient in the next frame. Next, in the bottom-up saliency mechanisms of
 human visual attention, local contrasts can invoke human gaze. But if the ob-
 servers are attracted by a semantic object in a different locus of a video frame,
 contrasted areas can have low values of GFDM. They thus become Non-salient.
 Hence, if the selection of Non-salient patches is based only on low GFDM values,
 then the resulting training data in Non-salient class of patches would contain
 noise. We illustrate such a phenomena in figure 1 (c). The focus of attention of
 subjects changes and Non-salient patches are selected even on the moving object
 (red ball). (a “Non-salient” patch in figure 1 (b) is selected on a contrasted
 background).



(a) heat map of frame #0013



(b) selected patches on frame #0013



(c) heat map of frame #0014



(d) selected patches on frame #0014

Figure 1: Extraction of Non-salient patches by random selection in the Non-salient area of a video frame (SRC07 video IRCCyN [37]).

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

165 We thus proposed a strategy based on video production rules which will
allows to reduce the noise in Non-salient training data. The method of patch
selection was explained in detail in [38]. Here we shortly remind it.

Figure 2 summarizes different steps to select salient patches. Firstly, the
GFDMs are computed, then morphological erosion is applied. The illustration
170 is given at a frame from HOLLYWOOD ¹ data set. Patches centered on local
maxima of saliency values are selected as salient. A Non-salient patch is a
squared region in the image plane which is not supposed to attract human gaze.
It should not be situated in the area-of-interest in a video frame, and must
not be already selected as salient. According to the rule of thirds in produced
175 and post-produced content, the most interesting details of the image or of a
video frame have to cover the frame center and the intersections of the three
horizontal and vertical lines that divide the image into nine equal parts [39].
Hence, we exclude the area -of-interest defined in the rule of thirds and selected
salient patches when randomly selecting Non-salient patches in each video frame
180 of training set.

3. Deep Convolutional Neural Network for visual saliency: ChaboNet

In this section, the proposed architecture ChaboNet for the visual saliency
prediction is presented. This is a relatively "light" deep architecture compared
to the popular GoogleNet [40], VGG [41] or variants of ResNet [42]. As the
185 purpose is in predicting visual saliency in video and not in static images, specific
features which are added to conventional RGB pixel values are described first.
The implementation of ChaboNet is realized on the basis of Caffe framework
[43].

3.1. A specific input data layer

190 When addressing visual attention prediction in video, the sensitivity of HVS
to motion has to be taken into account [44]. The sensitivity of HSV to mo-

¹available at <http://www.di.ens.fr/~laptev/actions/hollywood2/>

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

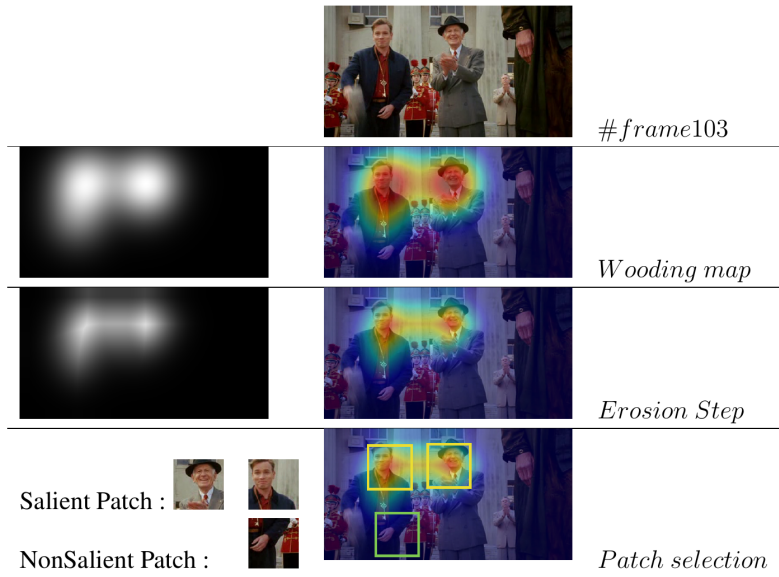


Figure 2: Policy of patch selection : example and steps, HOLLYWOOD data set

tion in a dynamic scene is modeled by residual motion [45]. Human observers accommodate to the global, i.e. ego-motion of the observer himself or camera motion in a recorded visual scene, and are attracted by specific local motions of objects. They first execute a saccade to a moving target and then continue with the “smooth pursuit” or visual tracking [46] keeping focus-of-attention on it. Local motion, of the target is expressed by residual motion relatively to the global camera motion [45]. To compute residual motion, the approach described in detail in [38] was followed. Here a pixel-wise motion field is computed by an optical flow method first. Using the dense motion field vectors as raw measures, the affine linear model of global motion is estimated by least square estimator and RANSAC algorithm [47]. Finally, the residual motion is the vector - difference between the initial motion vector and the one generated by the estimated affine model. As motion features, the squared $L2$ norm of residual motion vectors in each pixel in a video frame, normalized by its maximum in the frame, is used.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

The composition of the input layer of the CNN is illustrated in figure 3. Here for each patch the input layer is composed of three color channel values and the residual motion feature map. Due to this configuration, the model is called “4K-model” in contrast to “3K-model”, where only color channel values are used.

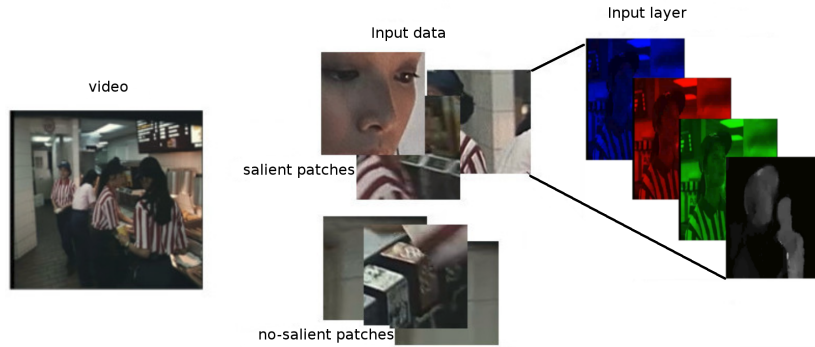


Figure 3: Input data layer : different features joined to the network.

3.2. The ChaboNet network architecture design

ChaboNet architecture was designed for the two-class classification problem: prediction of category of a patch in a given video frame as salient or non-salient. We aimed i) to preserve a reasonable deepness and ii) to remain comparable in the number of layers with a quite efficient network Alexnet [22]. The proposed ChaboNet architecture is summarized in figure 4.

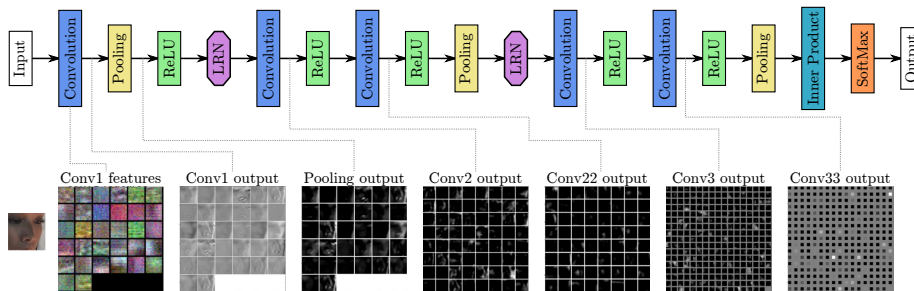


Figure 4: Architecture of video saliency convolution network ‘ChaboNet’.

As in Deep CNN architectures designed for image classification tasks [43], ChaboNet is composed of a hierarchy of patterns. Each pattern consists of a cascade of operations, followed by a normalization operation in some cases. The result of normalization operation is denoted as *Norm*. The cascading of linear and nonlinear operations successively produces high-level features. They are transmitted via a fully connected layer to the deepest layer which is a softmax classifier. It assigns the confidence for each patch to be salient or not. Due to quite a limited size of input patches three patterns are proposed in this architecture. The pattern P^1 below is a usual combination of convolution, pooling and non-linear layers, P^2 and P^3 have the same structure. The whole network can be detailed as follows.

$$\begin{aligned}
 & \text{Pattern } P^1 : \\
 & \text{Input} \xrightarrow{\text{convolution}} \text{Conv}^1 \xrightarrow{\text{pooling}} \text{Pool}^1 \xrightarrow{\text{ReLU}} R^1 \\
 & \text{Pattern } P^p : \text{ with } p \in \{2, 3\} \\
 & \text{Norm}^{p-1} \xrightarrow{\text{convolution}} \text{Conv}^p \xrightarrow{\text{ReLU}} R^p \xrightarrow{\text{convolution}} \text{Conv}^{pp} \xrightarrow{\text{ReLU}} R^{pp} \xrightarrow{\text{pooling}} \text{Pool}^p
 \end{aligned}$$

The normalization operation that leads to the output *Norm*, is added after the patterns P^1 and P^2 only, as after the pattern P^3 the features are quite sparse. The architecture of ChaboNet is depicted in figure 4. The features after convolution layers are presented for the example image from figure 3. It can be seen that the first layer of the network performs more as low-pass filters and deeper the convolution layer is more “high-pass” effect is observable.

3.2.1. Convolutional layers

In order to extract the most important information for further analysis or exploitation of image patches, the convolution with a fixed number of filters is needed. It is performed accordingly to the equation (1):

$$\mathbf{x}_j^l = \sum_{k \in \Omega_j} \mathbf{x}_k^{l-1} \odot \mathbf{w}_k^l + \mathbf{b}_j^l \quad (1)$$

with Ω_j - is the kernel support, i.e. the receptive field of j -th neuron;

l - is the network layer;
 \mathbf{x}_j^l - is the input of j -th neuron at layer l , that is feature-map vector;
 \mathbf{w}_k^l - is the weight of k -th neuron in the receptive field Ω_j ;
 \mathbf{b}_j^l - is the bias of j -th neuron at the layer l .
 \odot is Hadamard product which is a coordinate-wise operation.

Inspired by literature as [22], [29] where the size of convolution kernels is either maintained constant or is decreasing with the depth of layers, in ChaboNet network, 32 kernels were used with the size of 12×12 for the convolution layer of the first pattern P^1 . In the second pattern P^2 , 128 kernels for each convolutional layer were used. In P^2 the size of the kernels for the first convolutional layer was chosen as 6×6 and for the second convolution layer, a kernel of 3×3 was used. Finally, 288 kernels with the size of 3×3 were used for each convolution layer of the last pattern P^3 . This allows a progressive reduction of highly dimensional data before conveying them to the fully connected layers. The number of filters in the convolutional layers is growing, on the contrary, to explore the richness of the original data and to highlight structural patterns. For the filter size, several tests were made with the same values as in AlexNet [22], Shen’s network [29], LeNet [19], Cifar [48] and finally, the size of 12×12 was retained in the first layer of the pattern P^1 as it yielded the best accuracy in saliency prediction problem. The bias term was chosen as a null term.

3.2.2. Pooling layers

Pooling reduces the computational complexity for the upper layers and summarizes the outputs of neighboring groups of neurons from the same kernel map. It reduces the size of each input feature map by the acquisition of a value for each receptive field of neurons of the next layer. In our architecture max-pooling was used. For each channel of data at the l -th network layer it is expressed by equation (2):

$$\mathbf{y}^l(x, y) = \max_{x', y' \in \mathcal{N}} \mathbf{x}^{l-1}(x', y') \quad (2)$$

Here \mathcal{N} denotes the neighborhood of (x, y) . *max* operation is performed coordinate-wise. The kernel size of the pooling operation for the both patterns P^1 and P^2 was set to 3×3 . The pooling of the third pattern P^3 was done with a size of 1×1 , which means the full connection to the inner product layer.

3.2.3. Non Linear Response Layer

The non-linear transformation layer simulates the response of a neuron on excitement. In previous works on image classification the rectified linear unit (ReLU) function has been shown efficient [22]. The ReLU operation is expressed as (3)

$$\mathbf{z} = \max(\mathbf{0}, \mathbf{x}) \quad (3)$$

It is also applied in a coordinate-wise manner

Compared to usual for neural networks sigmoid function, ReLU does not suppress high frequency features. This is a good property for saliency prediction task. Indeed, HVS is sensitive to contrasts. The first pattern P^1 is designed in the manner that the ReLU operation is introduced after the pooling one. As the operations of pooling and ReLU compute the maximum, they are commutative. Cascading pooling before ReLU can reduce the execution time as pooling step reduces the number of neurons or nodes. In the two last patterns, stacking two convolutional layers before the destructive pooling layer ensures the computation of more complex features that will be more “expressive”.

3.2.4. Local response normalization layers

A local Response Normalization (LRN) layer normalizes values of feature maps which are calculated through the neurons having unbounded (due to ReLU) activation. This operation is used to detect the high-frequency characteristics with a high response of the neuron, and to scale down answers that are uniformly greater in a local area (see equation (4)).

$$\psi(\mathbf{z}(x, y)) = \frac{\mathbf{z}(x, y)}{\left(1 + \frac{\alpha}{N^2} \sum_{x'=\max(0, x-[N/2])}^{\min(S, x+[N/2]+N)} \sum_{y'=\max(0, y-[N/2])}^{\min(S, y+[N/2]+N)} (\mathbf{z}(x', y'))^2\right)^\beta} \quad (4)$$

Here $\mathbf{z}(x, y)$ represents the value of the feature map after ReLU operation at (x, y) coordinates and the sums are taken in the neighborhood of (x, y) of size $N \times N$, α and β regulate normalization strength. Normalization is also a coordinate-wise operation. Figure 5 summarizes the parameters used for each layer of the three patterns.

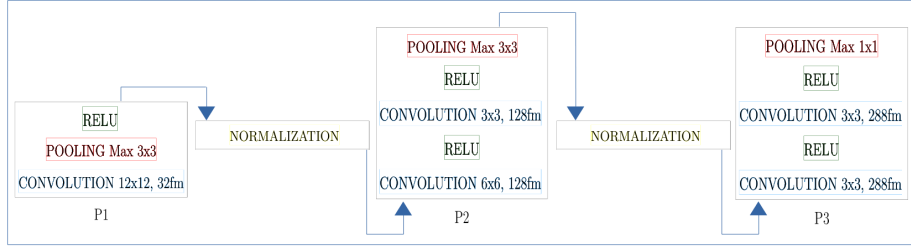


Figure 5: Detailed setting of each layer of ‘ChaboNet’ network.

3.3. Training and validation of the model

At a training step, the coefficients of convolution filters are repeatedly optimized in a forward-backward loop in the Caffe framework[43]. The optimization method used is the stochastic gradient descent ‘SGD’ with momentum. The initialization of convolution coefficients is realized randomly according to Gaussian law as proposed in [43]. The setting of the solver requires the definition of the number of iterations. It is defined accordingly to the equation (5):

$$iterations_numbers = epochs \times \frac{Total_images_number}{batch_size} \quad (5)$$

here $batch_size$ represents the number of images for each network switching, $epochs$ is the number of times the totality of the data set is switched by the network. The accuracy and loss are recorded with the interval of 2000 iterations both on training and validation sets.

4. Domain-dependent transfer learning for small databases

The generalization power of Deep CNN classifiers strongly depends on the quantity of data and on the coverage of data space in the training data set. In

1
2
3
4
5
6
7
8
9 real-life applications, e.g. prediction of benchmark models for studies of visual
10 attention of specific populations [49] or saliency prediction for visual quality
11 assessment [37] the database volumes are small. In order to predict saliency
12 in these small collections of videos, transfer learning approach was needed. It
13
14
15
16 315 presents a technique used in the field of machine learning that increases the
17 accuracy of learning either by using it in different tasks, or in the same task
18 [50]. Several studies have proven the power of this technique [51], [52]. In terms
19 of optimization method which is SGD, transfer learning means that the network
20 parameters are not initialized randomly, but their initialization corresponds to
21
22
23
24 320 a local minimum of loss function for a large data set. A small database can
25 be considered as slightly different data and starting from pre-trained parameter
26 values can bring improvement in optimization.
27

28 In saliency prediction problem, one could have supposed that as HVS is
29 sensitive to the singularities such as e.g. contrasts, the same trained saliency
30
31 325 model can be applied whatever the database is. In [49] we have shown that it is
32 not the case. Parameters trained on a large database for the same classification
33 task are not efficient when directly applied to another database, we observe
34 the so-called "over-fitting". Hence, the transfer learning or fine-tuning of the
35 parameters from a large database to a small database is necessary.
36
37
38

39 330 Transfer learning scheme consists of two steps: i) learning the whole classifi-
40 cation model on a large data set, ii) transfer on small data set. The latter means
41 initialization of parameters' values in learning process by the optimal parameter
42 values obtained on a large data set. In the present work, step i) - learning on a
43 large data set was performed from a scratch, i.e. with random initialization of
44
45
46
47 335 parameters at all layers. For the step ii) two initialization schemes were tested:
48 that one proposed by Bengio et al. [50] and ours [33] explained in the follow-
49 ing. For both steps i) and ii) the classification problem was the same : binary
50 classification of patches into salient and non-salient.
51
52

53 In the research of Bengio et al. [50] addressing object recognition problem,
54
55 340 the authors show that the first layers of a Deep CNN learn characteristics sim-
56 ilar to the responses of Gabor filters regardless of the data set or classification
57
58

task. Hence, in their transfer learning scheme just the three first convolutional layers pre-trained on a large database are used for the other database as the initialization of parameters. The coefficients on deeper layers are left free for optimization, that is initialized randomly.

In our case, saliency prediction task differs from object recognition task. Thus the proposal is to initialize all parameters in all layers of the network to train on a small data set, by the best model trained on a large data set. Equation (6) expresses the transfer of the convolutional weights, $W^l = \{\mathbf{w}_k\}^l$ for layer l , obtained from the larger database to the new smaller database. Here the Stochastic Gradient Descent with momentum is used as in [43]. In the following equation we omit any indexes except iteration number i for simplicity:

$$\begin{cases} V_{i+1} = m \cdot V_i - \gamma \cdot \epsilon \cdot W_i - \epsilon \cdot \langle \frac{\partial L}{\partial W} | W_i \rangle_{D_i} \\ W_{i+1} = W_i + V_{i+1} \quad | \quad W_0 = W' \end{cases} \quad (6)$$

With $\epsilon = 0.001$ - a fixed learning rate, $m = 0.9$ - momentum coefficient; $\gamma = 0.00004$ - weight decay and W' presents the best learned model parameters pre-trained on the large data set. The initial value of the velocity V_0 was set to zero. These parameter values are inspired by the values used in [43] with the same fixed learning rate and show the best performances on a large training data set. Hence in this section, the network for classification of patches in video frames as salient or non has been presented. In the next section, the method for generation of pixel-wise saliency maps on the basis of predicted salient patches is introduced.

5. Generation of saliency map

The saliency map of each frame I of the video is computed using the output value of the trained deep CNN model. Here we use the method from previous works [33], [38]. The soft-max classifier gives the probability for a patch of

1
2
3
4
5
6
7
8
9 belonging to the "salient" class accordingly to the equation 7.

$$\phi(\mathbf{u})_q = \frac{e^{u_q}}{\sum_r e^{u_r}}, r = 1, \dots, d \quad (7)$$

10
11
12
13
14 To build a dense predicted saliency map we classify patches of the input video
15 frames first. The patches of the same size s ($s = 100$ in our experiments) are
16 sampled with a stride of 5 pixels. The output value of the soft-max classifier
17 with regard to the salient class on each patch defines its degree of saliency. If
18 the score is assigned to the center of each patch (x_0, y_0) , a sparse saliency map
19 is obtained $M(x, y)$. Then, to densify the map, score values are interpolated
20 with Gaussian filters: in the center of each patch, a Gaussian $G(x, y)$ is applied
21 with a peak value of $\frac{A * M(x_0, y_0)}{2\pi\sigma^2}$. The A -parameter value was experimentally
22 chosen as 10. The spread parameter σ was fixed as a half-size of the patch. For
23 each pixel in the image plane, the Gaussian are summed-up. Finally the map is
24 normalized by saliency peak as in Wooding method for GFDM (see section 2).
25
26
27
28
29
30
31
32

33 **6. Experiments and results**

34 *6.1. Data sets*

35
36
37
38 To learn the model, three data sets were used, HOLLYWOOD[53] [54], the
39 well-known CRCNS [55] and IRCCyN [37].
40

41 The HOLLYWOOD database contains 823 training videos and 884 videos for
42 the validation step. The number of subjects with recorded gaze fixations varies
43 according to each video with up to 19 subjects. The spatial resolution of videos
44 varies as well. Despite the discrepancy of these parameters, we use it for model
45 building as it is the only large-scale video database with recorded gaze fixations.
46
47 The CRCNS ² data set [55] is one of the oldest and the most known data sets for
48 saliency prediction benchmarking. It contains 50 videos of 640×480 resolution
49 Gaze recordings of up to eight different subjects are available. To create the
50 training, validation and testing set, each video of CRCNS was split according to
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

²available at <https://crcns.org/data-sets/eye/eye-1>

the following scheme: one frame for testing, one frame for validation and four frames for training set were selected.

IRCCyN database is composed of 31 Standard-Definition (720x576) videos and comprises gaze fixations of 37 subjects. It is a homogeneous database designed for video quality assessment and by its good visual quality and the interest of the content, is the most interesting one. These videos contain certain categories of attention attractors such as high contrast, faces. However, videos with objects in motion are not frequent. Our model was designed for saliency prediction in video, to capture “smooth pursuit” motion of human eyes. It cannot be evaluated by using all available videos of IRCCyN data set. Hence, videos that do not contain a real object motion were eliminated. Only SRC02, SRC03, SRC04, SRC05, SRC06, SRC07, SRC10, SRC13, SRC17, SRC19, SRC23, SRC24 and SRC27 video files were thus used in experiments. This data set is referenced in the following as IRCCyN-MVT. For each chosen video of IRCCyN-MVT, the same selection scheme as for CRCNS data set was used: one frame is taken for the testing step, one frame for the validation step and four frames for the training step. The distribution of selected data between salient and non-salient classes is presented in table 1.

Table 1: Distribution of learning data: total number of salient and non-salient patches selected from each database.

data sets		training step	validation step
HOLLYWOOD	SalientPatch	222863	251294
	Non-SalientPatch	221868	250169
	<i>total</i>	444731	501463
CRCNS	SalientPatch	33370	8373
	Non-SalientPatch	30491	7730
	<i>total</i>	63861	16103
IRCCyN-MVT	SalientPatch	2013	511
	Non-SalientPatch	1985	506
	<i>total</i>	3998	1017

1
2
3
4
5
6
7
8
9 *6.2. Evaluation of the interest of residual motion*

10
11 410 To evaluate the interest of residual motion in saliency prediction, we have
12 computed the AUC metric[34] between gaze fixations and the energy-of-residual-
13 motion map. Here, we used popular data sets CRCNS [55] and IRCCyN [37]
14 that have been created and benchmarked for the task of saliency prediction in
15 natural videos.
16
17

18
19 415 Results summarized in table 2 and 4 show a correspondence between gaze
20 fixations and residual motion map especially for the “gamecube02” video of
21 CRCNS database with 0.56 value of AUC metric, and for the “SRC23” video of
22 IRCCyN database where we obtain a very interesting result ($AUC = 0.68$). In
23 table 4, 8 videos from 12 tested videos give an AUC value higher than 0.55. In
24 remaining sequences “SRC02”, “SRC07” and “SRC13’ moving objects are not
25 significant for their understanding. This experience can just encourage us to go
26 for the integration of residual motion as an input to ChaboNet architecture.
27
28 420
29
30
31

32 Table 2: The comparison of AUC metric of gaze fixations ‘GFM’ vs the energy of residual
33 motion map ‘ResidualMotion’ for 890 frames of CRCNS videos.
34

VideoName	TotFrame = 890	GFM vs ResidualMotion
beverly03	80	0.54 ± 0.119
gamecube02	303	0.56 ± 0.152
monica05	102	0.52 ± 0.110
standard02	86	0.499 ± 0.06
tv-announce01	73	0.472 ± 0.181
tv-news04	82	0.535 ± 0.186
tv-sports04	164	0.500 ± 0.147

35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50 Table 3: Frames of CRCNS videos.



Table 4: The comparison of AUC metric of gaze fixations 'GFM' vs Residual Motion map for 456 frames of IRCCyN videos.

VideoName	TotFrame = 456	GFM vs ResidualMotion
SRC02	37	0.46 ± 0.025
SRC03	28	0.55 ± 0.112
SRC04	35	0.55 ± 0.191
SRC05	35	0.57 ± 0.148
SRC06	36	0.603 ± 0.156
SRC07	36	0.48 ± 0.028
SRC10	33	0.55 ± 0.086
SRC13	35	0.59 ± 0.147
SRC17	42	0.48 ± 0.071
SRC19	33	0.64 ± 0.078
SRC23	40	0.68 ± 0.094
SRC24	33	0.51 ± 0.045
SRC27	33	0.53 ± 0.074

6.3. Evaluation of patches' saliency prediction with ChaboNet

The network was implemented using a graphic card Tesla K40m and processor (2×14 cores). A sufficiently large amount of patches, 256, was used per iteration (see the *batch_size* parameter in equation (5)). After a fixed number of training iterations, a model validation step was implemented: here the accuracy of the model at the current iteration was computed on the validation data set. It is denoted as "Test accuracy" in the figures 6, 7, 8.

To evaluate the deep network and to prove the importance of the addition of the residual motion map, two models were created with the same parameter settings and architecture of the network: the first one contains R, G and B primary pixel values in patches, denoted as *ChaboNet3k*. The *ChaboNet4k* is the model using RGB values and the normalized energy of residual motion as input data, see section 3.1.

Figure 6 illustrates the variations of the accuracy along iterations of all the models tested for the database HOLLYWOOD. Peak and mean accuracy values are presented in table 5). The results of learning experiments on HOLLYWOOD data set yield the following conclusions: i) when adding residual motion as an input feature to RGB values, the accuracy is improved by almost 2%. ii) the accuracy curve (figure 6 (a)) and the corresponding loss curve (figure 6(b)) show that the best trained model reached 80% of accuracy with the smallest

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

loss (at the iteration #8690 see table 5). Thus, it does not present an over-fitting situation. Figures 6 (c) and (d) show a better performance of the *ChaboNet4k* model in terms of speed for training and validation. Hence, the loss curve (Figure 6 (c)) of the *ChaboNet4k* model still decreasing during training. But since the 10000 iterations, the test loss curve (6 (d)) starts increasing which reflects the beginning of an over-fitting situation.

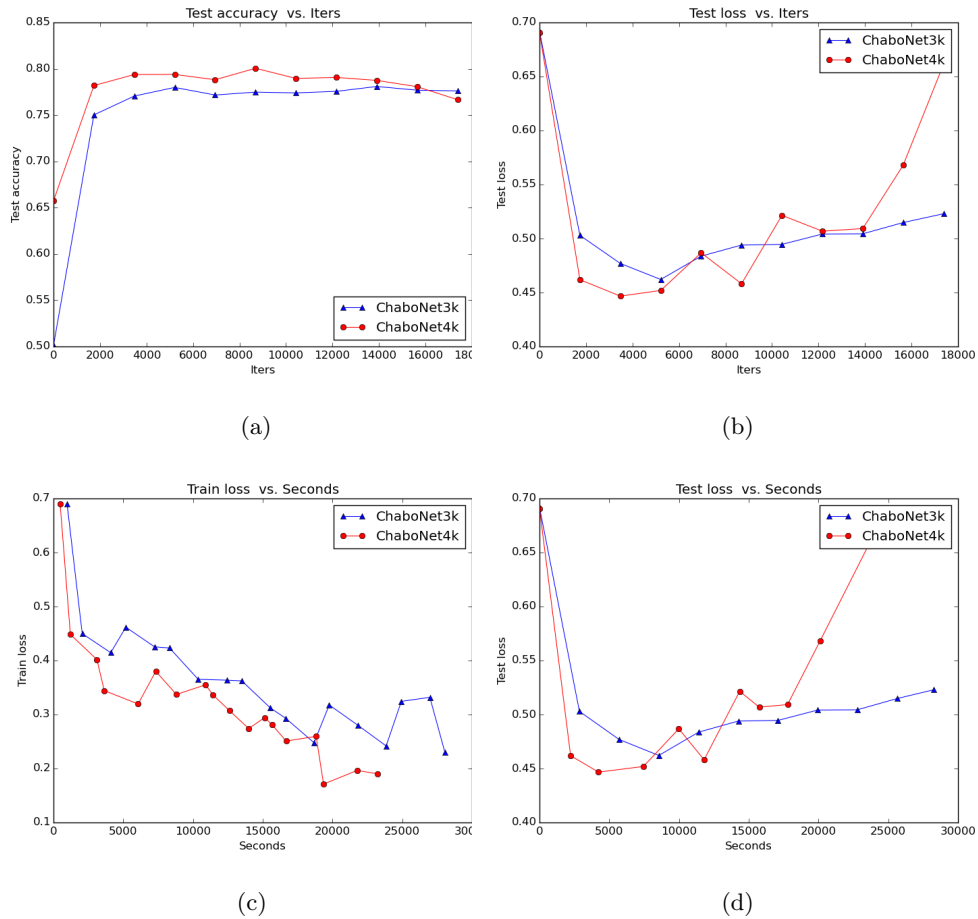


Figure 6: Training the network - Accuracy and loss vs time of *ChaboNet3k* and *ChaboNet4k* for HOLLYWOOD database : (a) Accuracy vs iterations, (b) Loss on validation data set vs iterations, (c) Train loss vs seconds, (d) Loss on validation data set vs seconds.

So, the model obtained after 8690 iterations is used to predict saliency on the validation set of this database, and to initialize the parameters when learning with transfer on other used data sets. Mean accuracy is also slightly higher. Indeed, 1.53% of mean accuracy increase is observed with merely the same stability of training which expressed by the standard deviation in the table 5.

Table 5: The accuracy results on HOLLYWOOD data set

	<i>ChaboNet3k</i>	<i>ChaboNet4k</i>
<i>training – time</i>	<i>7h47m33s</i>	<i>6h27m2s</i>
<i>min – Accuracy</i> _(#iter)	50.11% _(#0)	65.73% _(#0)
<i>max – Accuracy</i> _(#iter)	77.98% _(#5214)	80.05% _(#8690)
<i>avg – Accuracy ± std</i>	77.30% ± 0.864	78.73% ± 0.930

Figure 7 illustrates the variations of the accuracy along iterations of all models tested on IRCCyN-MVT data set. To overcome the lack of data for training, the proposed transfer learning scheme, see section 4 was used from HOLLYWOOD to IRCCyN-MVT. The best HOLLYWOOD model obtained at the iteration 8690 was used to initialize training of parameters with *ChaboNet4k* model on IRCCyN-MVT data set. For *ChaboNet3k*, training started from the iteration 5214 of the best model on HOLLYWOOD data set. Table 6 summarizes obtained results. The gain of using 4k- against 3k- data as input layer to the deep CNNs is about 1.12% in terms of mean accuracy.

Table 6: The accuracy results on IRCCyN-MVT data set.

	<i>ChaboNet3k</i>	<i>ChaboNet4k</i>
<i>training – time</i>	<i>0h4m6s</i>	<i>0h4m25s</i>
<i>min – Accuracy</i> _(#iter)	70.80% _(#5216)	77.83% _(#8848)
<i>max – Accuracy</i> _(#iter)	92.67% _(#6544)	92.77% _(#9664)
<i>avg – Accuracy ± std</i>	89.96 ± 4.159	91.08% ± 3.107

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

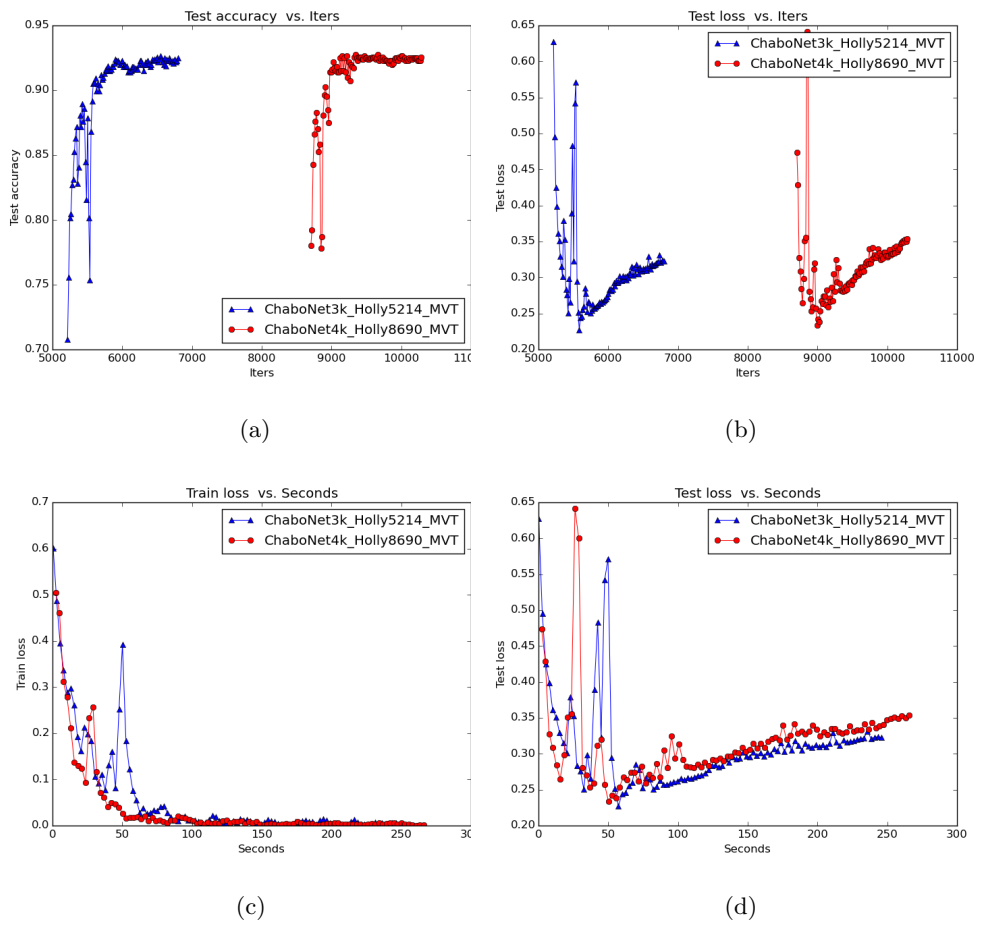


Figure 7: Accuracy and loss vs time of 3k and 4k for videos with motion from IRCCyN-MVT database : (a) Accuracy vs iterations, (b) Loss on validation data set vs iterations, (c) Train loss vs seconds, (d) Loss on validation data set vs seconds.

Figure 8 illustrates the variations of the accuracy and loss along the time expressed in number of iterations and seconds for CRCNS data set. The best model is obtained at the iteration #32500 with an accuracy of 91.66%.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

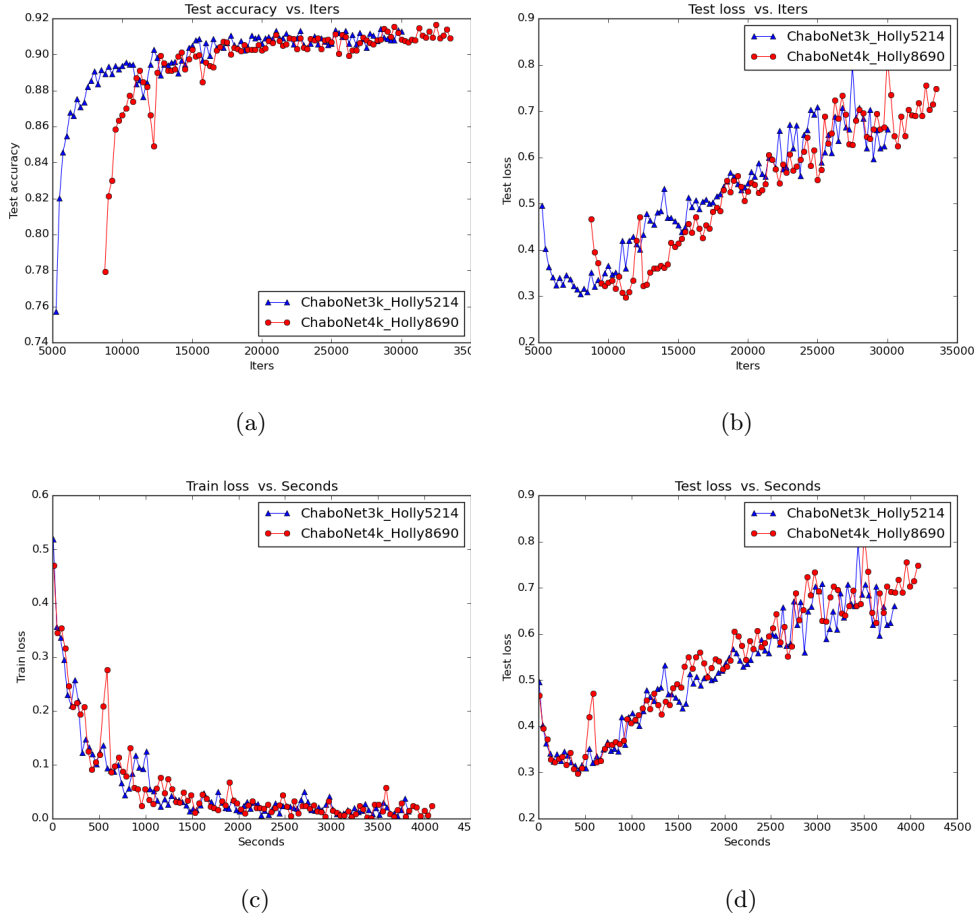


Figure 8: Accuracy and loss vs iterations of 3k and 4k for CRCNS database : a) Accuracy vs iterations, (b) Loss on validation data set vs iterations, (c) Train loss vs seconds, (d) Loss on validation data set vs seconds.

Table 7: The accuracy results on CRCNS data set

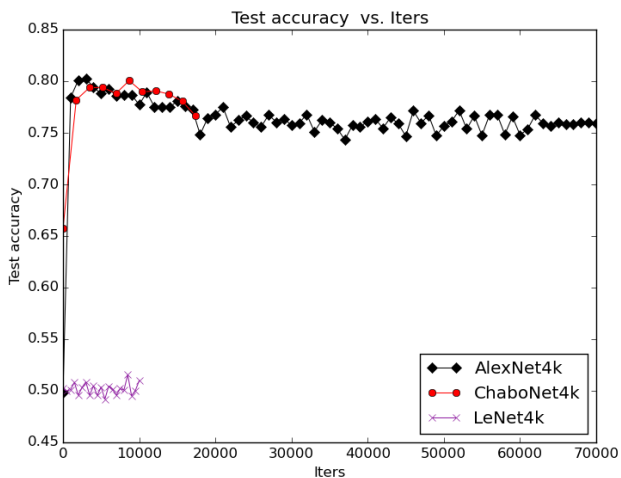
	<i>ChaboNet3k</i>	<i>ChaboNet4k</i>
<i>training – time</i>	1h3min42s	1h7min58s
<i>min – Accuracy</i> (#iter)	75.71% (#5250)	77.95% (#8750)
<i>max – Accuracy</i> (#iter)	91.45% (#28500)	91.66% (#32500)
<i>avg – Accuracy</i> \pm <i>std</i>	89.77% \pm 2.085	89.81% \pm 2.035

1
2
3
4
5
6
7
8
9
6.4. Validation of the ChaboNet architecture

10
11 To evaluate the ChaboNet *architecture* designed for saliency prediction,
12 an experiment was conducted with the HOLLYWOOD data set. The popular
13 AlexNet [22] and the original LeNet [19] network architectures were used as a
14
15
470 base-line with data patches extracted from HOLLYWOOD data.

17 For AlexNet, the network settings were taken exactly as in [22], that means
18 the same number and size of filters at all layers and the same learning parameter
19 such as the number of iterations (450.000). To better visualize, in figure 9 the
20 iterations of AlexNet were limited to 70.000. Similarly, the original settings
21 of AlexNet were limited to 70.000. Similarly, the original settings
22 of LeNet were preserved from [19]. Here the number of iterations was 10.000.
23
475 ChaboNet Network training was performed with 17.400 iterations.

26
27 Obtained results summarized in figure 9 showed that the ChaboNet network
28 outperformed the AlexNet and LeNet architectures (see table 8). In fact, with
29 17.400 iterations, ChaboNet outperformed by 2% in mean accuracy the AlexNet
30
31
480 architecture which needs 450.000 iterations. When comparing the 10.000 first
32 iterations of ChaboNet and LeNet, mean accuracy was discovered to be better
33 by more than 20%. Furthermore, the stability of training expressed by small
34 standard deviation is much stronger, see line 4 of the table 8.
35
36
37
38



56 Figure 9: ChaboNet 4k architecture vs AlexNet and LeNet on HOLLYWOOD data set.
57
58

Table 8: Accuracy results : validation of ChaboNet 4k architecture vs AlexNet and LeNet networks on HOLLYWOOD data set.

	<i>ChaboNet4k</i>	<i>AlexNet4k</i>	<i>LeNet4k</i>
$min_{(\#iter)}$	65.73% _(#0)	49,84% _(#0)	49,2% _(#5500)
$max_{(\#iter)}$	80.05% _(#8690)	80,27% _(#3000)	51,56% _(#8500)
$avg \pm std$	78.73% \pm 0,930	76,77% \pm 6,633	50,17% \pm 0,575

6.5. BN or LRN normalization for ChaboNet4k architecture

485 Recently, batch normalization ‘BN’ [56] which allows normalizing layer inputs, have shown its efficiency in designed architectures for image classification tasks. We compare and evaluate the use of batch normalization and the local response normalization for saliency prediction task. Figure 10 illustrates the variations of the accuracy and loss along the time expressed in number of iterations and seconds for the ChaboNet4k architecture using batch normalization 490 ‘BN’ and local response normalization ‘LRN’. Obtained results summarized in table 9, showed that the use of LRN layer outperforms the BN layer with 7% for saliency prediction tasks.

Table 9: The accuracy results on HOLLYWOOD data set

	<i>ChaboNet4k – LRN</i>	<i>ChaboNet4k – BN</i>
$min - Accuracy_{(\#iter)}$	65.73% _(#0)	45.71% _(#0)
$max - Accuracy_{(\#iter)}$	80.05% _(#8690)	76.82% _(#15642)
$avg - Accuracy \pm std$	78.73% \pm 0.930	70.78% \pm 7.034

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

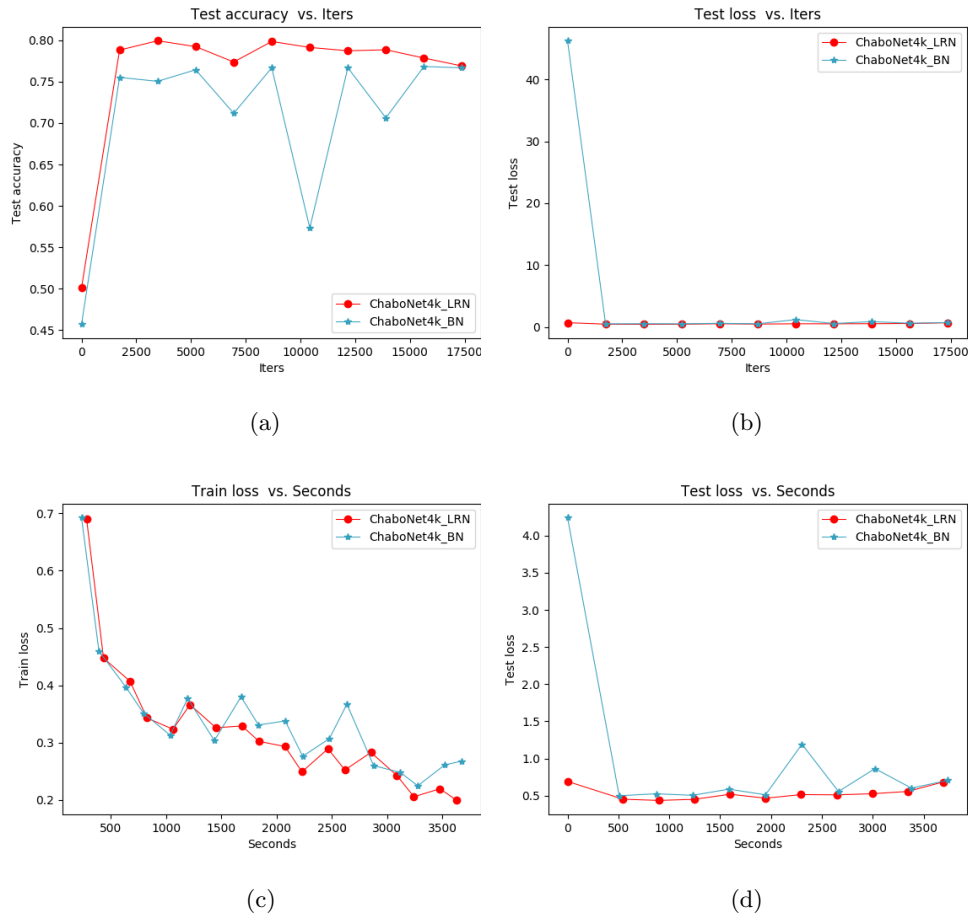


Figure 10: Training the network - Accuracy and loss vs time of *ChaboNet4k* for HOLLYWOOD database using local response normalization (LRN) and batch normalization (BN) : (a) Accuracy vs iterations, (b) Loss on validation data set vs iterations, (c) Train loss vs seconds, (d) Loss on validation data set vs seconds.

6.6. Validation of the proposed method of transfer learning

495 The previous work [33] already showed that training on a small database (IRCCyN) with transfer learning increases the accuracy by 4% in average and makes training process more stable (up to 50% of decrease of standard deviation of accuracy along iterations). In the present work, the proposed transfer

learning scheme is benchmarked with that one proposed by Bengio. Two ex-
 500 periments were conducted with the same small data set IRCCyN-MVT and
 CRCNS, and the same definition of network ChaboNet: i) Our method: start
 training of all ChaboNet layers from the best model already trained on the
 large HOLLYWOOD data set (see section 4). ii) Bengio’s method: the three
 first convolutional layers are trained on the HOLLYWOOD data set and then
 505 fine-tuned on the target data set, other layers are trained on target data set
 with random initialization.

Figure 11 illustrates the variations of the accuracy along iterations of the two
 experiments performed with the two small data sets. One can see less stable
 behavior when the transfer method of Bengio et al. is applied. The proposed
 510 method of transfer learning outperformed the Bengio’s method by almost 3.6%
 in IRCCyN-MVT data set and almost 0.44% in CRCNS data set in terms of
 mean accuracy. The gain on stability of training in our method is more than
 50%, see table 10.

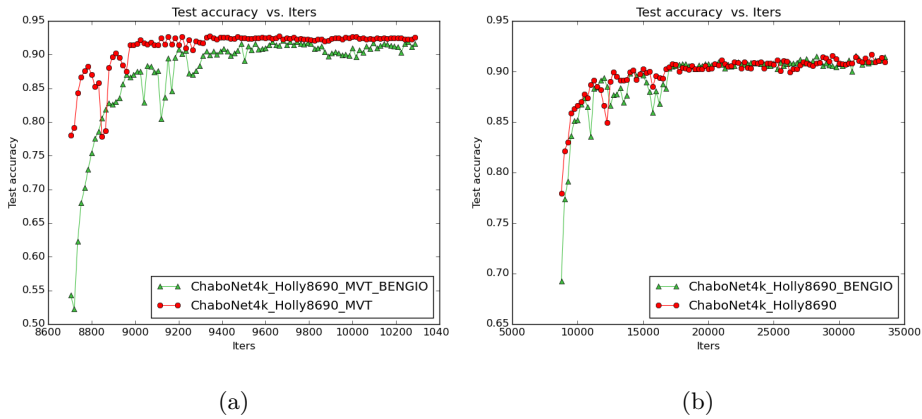


Figure 11: Evaluation and comparison of the proposed method of transfer learning :
 (a) Comparison on IRCCyN-MVT, (b) Comparison on CRCNS.

Table 10: The accuracy results on IRCCyN-MVT and CRCNS data set.

	<i>Our transfer method</i>		<i>BENGIO transfer method</i>	
	IRCCyN-MVT	CRCNS	IRCCyN-MVT	CRCNS
$max(\#iter)$	92.77% (#9664)	91.66% (#32500)	92.08% (#9680)	91.55% (#31250)
$avg \pm std$	91.08% \pm 3.107	89.81% \pm 2.035	87.48% \pm 7.243	89.37% \pm 3.099

6.7. Evaluation of predicted visual saliency maps

515 After training and validation of the model on HOLLYWOOD data set, the model obtained at the iteration #8690 having the maximum value of accuracy 80.05% was chosen. This model was used to predict the probability of a patch to be salient. For the CRCNS data set, the model obtained at the iteration #32500 with the accuracy of 91.66% is used to predict saliency. In the same
 520 manner, the model with the accuracy of 92.77% obtained at the iteration #9664 is used for the IRCCyN-MVT data set.

To evaluate the method of saliency prediction by interpolation of classifier outputs as presented in section 5, performances were compared with the most popular saliency models from the literature. Several spatial saliency models
 525 were chosen: Itti and Koch spatial model [57], Signature Sal [58] (a simple image descriptor is introduced here referred to as the “image signature”. The authors show that it performs better than Itti and Koch model) and GBVS (regularized spatial saliency model of Harel [59]). We also benchmarked our model against spatio-temporal models for saliency prediction in videos like the model proposed
 530 by Seo [60] which is built upon optical flow and the model of Wang [9] which is based on the gradient flow field and energy optimization. Finally, the last benchmark is with the ELD [16] deep saliency model which used both high level and low level features for saliency detection under an unified deep learning framework.

535 In tables 12, 13, 11 below, the comparison of Deep CNN prediction of pixel-wise saliency maps with the Gaze Fixations Density Maps (Gaze-fix) is shown.

The quality of predicted maps is compared with prediction by classical saliency models (Signature Sal, GBVS, Seo) also compared to the same reference: GFDM and the recent state-of-the-art methods of Wang and ELD. The

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

540 comparison is given in terms of the widely used AUC metric [61]. Mean value
of the metric for each saliency model compared to the GFDM is given together
with standard deviation for a sample of videos. In table 11, the *ChaboNet4k*
is compared with reference models for CRCNS data set. In table 12 the maps
built on HOLLYWOOD database with its best patch saliency prediction model
545 *Chabonet4K* are compared with three reference models and in table 13 the com-
parison is fulfilled on IRCCyN-MVT data set for both 3K and 4K ChaboNet
models.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Table 11: The comparison of AUC metric of gaze fixations 'Gaze-fix' vs predicted saliency 'GBVS', 'IttiKoch', 'Seo', 'Wang' and 'ELD' and the ChaboNet4k for 890 frames of CRCNS videos.

VideoName	<i>TotalFrame</i> = 890	Gaze-fix vs GBVS	Gaze-fix vs IttiKoch	Gaze-fix vs Seo	Gaze-fix vs Wang	Gaze-fix vs ELD	Gaze-fix vs ChaboNet4k
beverly03	80	0.78 ± 0.151	0.77 ± 0.124	0.66 ± 0.172	0.73 ± 0.155	0.59 ± 0.118	0.79 ± 0.118
gamecube02	303	0.73 ± 0.165	0.74 ± 0.180	0.61 ± 0.179	0.76 ± 0.126	0.78 ± 0.145	0.82 ± 0.126
monica05	102	0.75 ± 0.183	0.73 ± 0.158	0.54 ± 0.156	0.73 ± 0.162	0.67 ± 0.162	0.79 ± 0.133
standard02	86	0.78 ± 0.132	0.72 ± 0.141	0.61 ± 0.169	0.78 ± 0.094	0.60 ± 0.159	0.71 ± 0.181
tv-announce01	73	0.60 ± 0.217	0.64 ± 0.203	0.52 ± 0.206	0.00 ± 0.00	0.64 ± 0.209	0.63 ± 0.215
tv-news04	82	0.78 ± 0.169	0.79 ± 0.154	0.61 ± 0.162	0.00 ± 0.00	0.85 ± 0.091	0.72 ± 0.145
tv-sports04	164	0.68 ± 0.182	0.69 ± 0.162	0.56 ± 0.193	0.00 ± 0.00	0.64 ± 0.179	0.78 ± 0.172

1
2
3
4
5
6
7
8
9 The best AUC metric values are underscored. It can be stated that in general
10 spatial models (Signature Sal, GBVS or Itti) performed better in half of the
11
12 550 tested videos. This is due to the fact that these videos contain very contrasted
13 areas in the video frames, which attract human gaze. They do not contain
14 areas having an interesting residual motion. Nevertheless, the *ChaboNet4K*
15 model systematically outperforms Seo’s and Wang’s model which use motion
16 features. Our proposed deep network still remains competitive with the ELD
17
18 555 deep saliency method. This shows definitively that the use of a Deep CNN is a
19 way for prediction of visual saliency in video scenes. However, for IRCCyN-MVT
20 data set, see table 13, despite videos without any motion were left aside, the gain
21 in performance of the proposed model is not very clear due to the complexity of
22 these visual scenes, such as presence of strong contrasts and faces. Using high
23 level features in ELD method ensures an interesting results on IRCCyN-MVT
24 data set.
25
26
27
28 560
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Table 12: The comparison of AUC metric of gaze fixations 'Gaze-fix' vs predicted saliency 'GBVS', 'SignatureSal', 'Seo', 'Wang' and 'ELD') and the ChaboNet4k for the videos from HOLLYWOOD.

VideoName	$TotFrame = 2248$	Gaze-fix vs GBVS	Gaze-fix vs SignatureSal	Gaze-fix vs Seo	Gaze-fix vs Wang	Gaze-fix vs ELD	Gaze-fix vs ChaboNet4k
clipTest56	137	0.76 ± 0.115	0.75 ± 0.086	0.64 ± 0.116	0.68 ± 0.105	0.66 ± 0.105	0.77 ± 0.118
clipTest105	154	0.63 ± 0.169	0.57 ± 0.139	0.54 ± 0.123	0.77 ± 0.069	0.77 ± 0.103	0.69 ± 0.186
clipTest147	154	0.86 ± 0.093	0.90 ± 0.065	0.70 ± 0.103	0.84 ± 0.063	0.76 ± 0.063	0.81 ± 0.146
clipTest250	160	0.74 ± 0.099	0.69 ± 0.110	0.47 ± 0.101	0.81 ± 0.073	0.75 ± 0.116	0.71 ± 0.180
clipTest350	66	0.65 ± 0.166	0.68 ± 0.249	0.57 ± 0.124	0.74 ± 0.177	0.82 ± 0.083	0.72 ± 0.177
clipTest400	200	0.75 ± 0.127	0.67 ± 0.110	0.60 ± 0.106	0.74 ± 0.133	0.74 ± 0.119	0.71 ± 0.146
clipTest451	132	0.70 ± 0.104	0.59 ± 0.074	0.57 ± 0.068	0.67 ± 0.096	0.68 ± 0.070	0.63 ± 0.151
clipTest500	166	0.82 ± 0.138	0.84 ± 0.150	0.75 ± 0.152	0.88 ± 0.076	0.89 ± 0.052	0.84 ± 0.156
clipTest600	200	0.75 ± 0.131	0.678 ± 0.149	0.53 ± 0.108	0.81 ± 0.091	0.72 ± 0.136	0.71 ± 0.180
clipTest650	201	0.72 ± 0.106	0.74 ± 0.087	0.61 ± 0.092	0.65 ± 0.086	0.58 ± 0.090	0.70 ± 0.078
ClipTest700	262	0.74 ± 0.128	0.76 ± 0.099	0.50 ± 0.059	0.87 ± 0.058	0.70 ± 0.107	0.78 ± 0.092
clipTest800	200	0.70 ± 0.096	0.75 ± 0.071	0.53 ± 0.097	0.56 ± 0.118	0.71 ± 0.128	0.66 ± 0.141
ClipTest803	102	0.86 ± 0.106	0.87 ± 0.068	0.73 ± 0.148	0.85 ± 0.070	0.77 ± 0.088	0.88 ± 0.078
ClipTest849	114	0.75 ± 0.155	0.91 ± 0.070	0.55 ± 0.122	0.92 ± 0.024	0.86 ± 0.033	0.74 ± 0.132

Table 13: The comparison of AUC metric of gaze fixations 'Gaze-fix' vs predicted saliency 'GBVS', 'SignatureSal', 'Seo', 'Wang' and 'ELD') and the ChaboNet3k and ChaboNet4k for the videos from IRCCyN-MVT.

VideoName	<i>TotalFrame</i> = 1227	Gaze-fix vs GBVS	Gaze-fix vs SignatureSal	Gaze-fix vs Seo	Gaze-fix vs Wang	Gaze-fix vs ELD	Gaze-fix vs ChaboNet3k	Gaze-fix vs ChaboNet4k
src02	37	0.68 ± 0.076	0.49 ± 0.083	0.44 ± 0.017	0.43 ± 0.039	0.60 ± 0.043	0.012 ± 0.077	0.48 ± 0.073
src03	28	0.82 ± 0.088	0.87 ± 0.057	0.76 ± 0.091	0.80 ± 0.063	0.69 ± 0.099	0.00 ± 0.000	0.70 ± 0.149
src04	35	0.79 ± 0.058	0.81 ± 0.029	0.59 ± 0.057	0.76 ± 0.112	0.81 ± 0.048	0.12 ± 0.214	0.57 ± 0.135
src05	35	0.73 ± 0.101	0.67 ± 0.122	0.48 ± 0.071	0.71 ± 0.088	0.75 ± 0.113	0.39 ± 0.186	0.53 ± 0.128
src06	36	0.85 ± 0.080	0.71 ± 0.151	0.73 ± 0.148	0.81 ± 0.092	0.72 ± 0.153	0.00 ± 0.000	0.60 ± 0.180
src07	36	0.72 ± 0.070	0.73 ± 0.060	0.57 ± 0.060	0.74 ± 0.062	0.65 ± 0.090	0.34 ± 0.284	0.55 ± 0.135
src10	33	0.87 ± 0.048	0.92 ± 0.043	0.82 ± 0.101	0.90 ± 0.045	0.82 ± 0.074	0.90 ± 0.045	0.60 ± 0.173
src13	35	0.79 ± 0.103	0.75 ± 0.111	0.64 ± 0.144	0.82 ± 0.067	0.71 ± 0.116	0.36 ± 0.201	0.52 ± 0.138
src17	42	0.55 ± 0.092	0.33 ± 0.099	0.45 ± 0.033	0.45 ± 0.062	0.62 ± 0.095	0.00 ± 0.000	0.51 ± 0.098
src19	33	0.76 ± 0.094	0.68 ± 0.086	0.59 ± 0.117	0.75 ± 0.062	0.89 ± 0.020	0.46 ± 0.075	0.75 ± 0.123
src23	40	0.76 ± 0.050	0.69 ± 0.070	0.58 ± 0.067	0.77 ± 0.045	0.83 ± 0.027	0.03 ± 0.169	0.66 ± 0.105
src24	33	0.63 ± 0.071	0.58 ± 0.054	0.55 ± 0.059	0.65 ± 0.040	0.55 ± 0.046	0.23 ± 0.252	0.50 ± 0.052
src27	33	0.59 ± 0.117	0.64 ± 0.091	0.52 ± 0.057	0.77 ± 0.064	0.68 ± 0.071	0.00 ± 0.000	0.54 ± 0.106

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Table 14: Visual evaluation of different saliency models for video taken from Hollywood data set '*NameVideo#_frameNumber*'

	Human	ChaboNet4k	GBVS	Seo	Wang	ELD
clipTest56#33						
clipTest147#34						
clipTest500#163						
clipTest700#02						
clipTest803#39						

Table 15, presents the time needed for testing one patch and the creation of the saliency map across one frame with a stride of 5 pixels. We used machines with Xeon E5 processor. The computation time for one patch is quite reasonable, the overall time for one full HD video frame remains high despite parallelization of computations. This is explained by a large quantity of scanning windows, more GPU processors are need for faster computation.

Table 15: Time for testing one patch and one frame of video.

	machine $8\mu p$	machine $20\mu p$	machine $2 \times 14cores\mu p$
patch 100×100	0.015s	0.028s	0.011s
frame 720×576	42.31s	18.49s	8.56s

7. Conclusion and perspectives

This study addressed the problem of prediction of visual attention in video content with Deep CNNs. In this paper we have further extended experiments and confirmed partial results we obtained in our previous works. We hypothesized that the model could capture gaze attraction by moving objects due to the residual motion maps added to primary color pixel values. First of all, we measured the correspondence of residual motion maps with gaze fixations of observers which confirmed the interest of their incorporation into input layer of proposed Deep architecture. The performances of prediction when different kinds of features are ingested by the network -color pixel values only, color values with residual motion- were compared. As far as dynamic content is concerned, the saliency is better predicted with spatio-temporal features (RGB and residual motion) when scenes do not contain distracting contrasts. The proposed relatively shallow architecture ChaboNet was compared to similar architectures AlexNet and LeNet and showed better prediction power in terms of mean accuracy and stability of training phase. The transfer learning scheme applied to the prediction of saliency on small data sets by fine-tuning parameters pre-trained on

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

585 a large data set (HOLLYWOOD) successfully outperforms the state-of-the-art,
i.e. Bengio’s method. Finally, a method for building pixel-wise saliency maps,
using the probability of patches to be salient, was extensively tested against ref-
erence spatial, spatio-temporal and a deppe learning-based prediction models.
We come to the conclusion that the interpolation of classification results is not a
590 definite way to build dense predicted saliency maps. Indeed, accuracy in patch
classification are high, but the AUC metric values with reference GFDMs are
not systematically better than reference classical models. We hypothesize that
this is partly a distractor problem, but still the way to predict a dense map can
be further developed with FC-networks. Last but not least, Deep NNs supply
595 an interesting framework for use of temporal continuity of saliency maps in the
absence of distractors. This is the future of the present research.

Acknowledgment

This research has been supported by University of Bordeaux, University of
Sfax and the grant UNetBA.

600 **References**

[1] A. M. Treisman, G. Gelade, A feature-integration theory of attention, *Cog-
nitive Psychology* 12 (1) (1980) 97–136.

[2] B. J. Han, J. Y. Sim, Saliency detection for panoramic landscape images
of outdoor scenes, *Journal of Visual Communication and Image Representa-
tion* 49 (Supplement C) (2017) 27 – 37.

605 [3] Z. Li, C. Lang, S. Feng, T. Wang, Saliency ranker: A new salient object
detection method, *Journal of Visual Communication and Image Representa-
tion* 50 (Supplement C) (2018) 16 – 26.

[4] J. Ren, Z. Liu, X. Zhou, G. Sun, C. Bai, Saliency integration driven by sim-
610 ilar images, *Journal of Visual Communication and Image Representation*
50 (2018) 227 – 236.

- 1
2
3
4
5
6
7
8
9 [5] Z. Liu, X. Zhang, S. Luo, O. Le Meur, Superpixel-based spatiotemporal
10 saliency detection, *IEEE Transactions on Circuits and Systems for Video*
11 *Technology* 24 (9) (2014) 1522–1540. doi:10.1109/TCSVT.2014.2308642.
12
13
14
15 615 [6] J. Yang, G. Zhao, J. Yuan, X. Shen, Z. Lin, B. L. Price, J. Brandt, Discov-
16 ering primary objects in videos by saliency fusion and iterative appearance
17 estimation, *IEEE Trans. Circuits Syst. Video Techn.* 26 (6) (2016) 1070–
18 1083. doi:10.1109/TCSVT.2015.2433171.
19 URL <http://dx.doi.org/10.1109/TCSVT.2015.2433171>
20
21
22
23 620 [7] Y. Pang, L. Ye, X. Li, J. Pan, Incremental learning with saliency map
24 for moving object detection, *IEEE Transactions on Circuits and Systems*
25 *for Video Technology* PP (99) (2016) 1–1. doi:10.1109/TCSVT.2016.
26 2630731.
27
28
29
30 [8] L. Duan, T. Xi, S. Cui, H. Qi, A. C. Bovik, A spatiotemporal
31 625 weighted dissimilarity-based method for video saliency detection,
32 *Signal Processing: Image Communication* 38 (2015) 45 – 56, re-
33 cent *Advances in Saliency Models, Applications and Evaluations*.
34 doi:<http://dx.doi.org/10.1016/j.image.2015.08.005>.
35 URL [http://www.sciencedirect.com/science/article/pii/](http://www.sciencedirect.com/science/article/pii/S0923596515001307)
36 [S0923596515001307](http://www.sciencedirect.com/science/article/pii/S0923596515001307)
37
38 630
39
40
41 [9] W. Wang, J. Shen, L. Shao, Consistent video saliency using local gradi-
42 ent flow optimization and global refinement, *IEEE Transactions on Image*
43 *Processing* 24 (11) (2015) 4185–4196.
44
45
46
47 [10] S. Zhong, Y. Liu, F. Ren, J. Zhang, T. Ren, Video saliency detection via
48 635 dynamic consistent spatio-temporal attention modelling, in: *Proceedings*
49 *of the Twenty-Seventh AAAI Conference on Artificial Intelligence*, AAAI
50 Press, 2013, pp. 1063–1069.
51
52
53
54 [11] I. González-Díaz, V. Buso, J. Benois-Pineau, Perceptual modeling in the
55 problem of active object recognition in visual scenes, *Pattern Recognition*
56 56 (2016) 129–141.
57 640
58
59
60
61
62
63
64
65

- 1
2
3
4
5
6
7
8
9 [12] Y. Pinto, A. R. van der Leij, I. G. Sligte, V. F. Lamme, H. S. Scholte,
10 Bottom-up and top-down attention are independent, *Journal of Vision*
11 13 (3) (2013) 16.
12
13
14 [13] J. Shen, L. Itti, Top-down influences on visual attention during listening
15 are modulated by observer sex, *Vision Research* 65 (2012) 62–76.
16 645
17
18 [14] A. Borji, D. N. Sihite, L. Itti, What stands out in a scene? A study of
19 human explicit saliency judgment, *Vision Research* 91 (2013) 62–77.
20
21
22 [15] H. Boujut, J. Benois-Pineau, T. Ahmed, O. Hadar, P. Bonnet, No-reference
23 video quality assessment of h.264 video streams based on semantic saliency
24 maps, Vol. 8293, 2012, pp. 82930T–82930T–9.
25 650
26
27 [16] G. Lee, Y.-W. Tai, J. Kim, Deep saliency with encoded low level distance
28 map and high level features, 2016 IEEE Conference on Computer Vision
29 and Pattern Recognition (CVPR) (2016) 660–668.
30
31
32 [17] G. Sharma, F. Jurie, C. Schmid, Discriminative spatial saliency for image
33 classification, in: 2012 IEEE Conference on Computer Vision and Pattern
34 655 Recognition, 2012, pp. 3506–3513.
35
36
37 [18] L. Deng, D. Yu, DEEP LEARNING: Methods and Applications, *Founda-
38 tions and Trends in Signal Processing* 7 (MSR-TR-2014-21) (2014) 197–387.
39
40
41 [19] Y. LeCun, L. Bottou, Y. Bengio, P. Haffner, Gradient-Based Learning Ap-
42 plied to Document Recognition, no. 86(11), 1998, pp. 2278–2324.
43 660
44
45 [20] J. Bruna, S. Mallat, Invariant Scattering Convolution Networks, *IEEE
46 Trans. Pattern Anal. Mach. Intell.* 35 (8) (2013) 1872–1886.
47
48
49 [21] Y. Bengio, A. Courville, P. Vincent, Representation Learning: A Review
50 and New Perspectives, *IEEE Transactions on Pattern Analysis and Ma-
51 chine Intelligence* (35 (8)) (2014) 1798–1828.
52 665
53
54 [22] A. Krizhevsky, I. Sutskever, G. E. Hinton, ImageNet Classification with
55 Deep Convolutional Neural Networks, in: F. Pereira, C. Burges, L. Bottou,
56
57
58

1
2
3
4
5
6
7
8
9 K. Weinberger (Eds.), *Advances in Neural Information Processing Systems*
10 25, Curran Associates, Inc., 2012, pp. 1097–1105.

11
12
13 670 [23] H. R. Tavakoli, A. Borji, J. Laaksonen, E. Rahtu, Exploiting inter-image
14 similarity and ensemble of extreme learners for fixation prediction using
15 deep features, *CoRR* abs/1610.06449.

16
17 URL <http://arxiv.org/abs/1610.06449>

18
19
20 [24] S. F. Dodge, L. J. Karam, Visual saliency prediction using a mixture of
21 675 deep neural networks, *CoRR* abs/1702.00372. [arXiv:1702.00372](https://arxiv.org/abs/1702.00372).

22
23
24 [25] R. Monroy, S. Lutz, T. Chalasani, A. Smolic, Salnet360: Saliency maps
25 for omni-directional images with CNN, *CoRR* abs/1709.06505. [arXiv:](https://arxiv.org/abs/1709.06505)
26 1709.06505.

27
28 URL <http://arxiv.org/abs/1709.06505>

29
30
31 680 [26] S. S. S. Kruthiventi, V. Gudisa, J. H. Dholakiya, R. V. Babu, Saliency
32 Unified: A Deep Architecture for simultaneous Eye Fixation Prediction
33 and Salient Object Segmentation, 2016 IEEE Conference on Computer
34 Vision and Pattern Recognition (CVPR) 00 (2016) 5781–5790. [doi:10.](https://doi.org/10.1109/CVPR.2016.623)
35 1109/CVPR.2016.623.

36
37
38
39 685 [27] M. Kümmerer, L. Theis, M. Bethge, Deep Gaze I: Boosting Saliency Pre-
40 diction with Feature Maps Trained on ImageNet, *CoRR* abs/1411.1045.

41
42
43 [28] J. Pan, Giró i Nieto, X., End-to-end Convolutional Network for Saliency
44 Prediction, *CoRR* abs/1507.01422.

45
46
47 [29] C. Shen, Q. Zhao, Learning to Predict Eye Fixations for Semantic Contents
48 690 Using Multi-layer Sparse Network, *Neurocomputing* 138 (2014) 61–68.

49
50
51 [30] K. Simonyan, A. Vedaldi, A. Zisserman, Deep Inside Convolutional Net-
52 works: Visualising Image Classification Models and Saliency Maps, *CoRR*
53 abs/1312.6034.

- 1
2
3
4
5
6
7
8
9 [31] E. Vig, M. Dorr, D. Cox, Large-Scale Optimization of Hierarchical Features
10 for Saliency Prediction in Natural Images , in: Proceedings of the 2014
11 695 IEEE Conference on Computer Vision and Pattern Recognition, CVPR
12 '14, 2014, pp. 2798–2805.
13
14
15
16 [32] J. Han, L. Sun, X. Hu, J. Han, L. Shao, Spatial and temporal visual atten-
17 tion prediction in videos using eye movement data, *Neurocomputing* 145
18 (2014) 140–153.
19 700
20
21 [33] S. Chaabouni, J. Benois-Pineau, C. Ben Amar, Transfer learning with deep
22 networks for saliency prediction in natural video, in: 2016 IEEE Interna-
23 tional Conference on Image Processing (ICIP), Vol. 91, 2016, pp. 1604–
24 1608.
25
26
27
28 [34] A. Coutrot, O. Le Meur, Visual attention saccadic models: taking into
29 705 account global scene context and temporal aspects of gaze behaviour, 2016,
30 poster.
31 URL <https://hal.inria.fr/hal-01391751>
32
33
34
35 [35] M. Cornia, L. Baraldi, G. Serra, R. Cucchiara, Multi-level net: A visual
36 saliency prediction model, in: *Computer Vision - ECCV 2016 Workshops -*
37 710 *Amsterdam, The Netherlands, October 8-10 and 15-16, 2016, Proceedings,*
38 *Part II*, 2016, pp. 302–315. doi:10.1007/978-3-319-48881-3_21.
39 URL https://doi.org/10.1007/978-3-319-48881-3_21
40
41
42
43
44 [36] D. S. Wooding, Eye movements of large populations: II. Deriving regions
45 715 of interest, coverage, and similarity using fixation maps, *Behavior Research*
46 *Methods, Instruments, & Computers* 34 (4) (2002) 518–528.
47
48
49 [37] F. Boulos, W. Chen, B. Parrein, P. Le Callet, Region-of-Interest Intra Pre-
50 diction for H.264/AVC Error Resilience, in: *IEEE International Conference*
51 *on Image Processing*, Cairo, Egypt, 2009, pp. 3109–3112.
52
53
54
55 720 [38] S. Chaabouni, J. Benois-Pineau, O. Hadar, Prediction of visual saliency in
56 video with deep cnns, Vol. 9971, 2016, pp. 99711Q–99711Q–14.
57
58
59
60
61
62
63
64
65

- 1
2
3
4
5
6
7
8
9 [39] L. Mai, H. Le, Y. Niu, F. Liu, Rule of thirds detection from photograph,
10 in: Multimedia (ISM), 2011 IEEE International Symposium on, 2011, pp.
11 91–96.
12
13
14
15 725 [40] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Er-
16 han, V. Vanhoucke, A. Rabinovich, Going Deeper with Convolutions, in:
17 Computer Vision and Pattern Recognition (CVPR), 2015.
18
19
20 [41] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-
21 scale image recognition, CoRR abs/1409.1556. [arXiv:1409.1556](https://arxiv.org/abs/1409.1556).
22 URL <http://arxiv.org/abs/1409.1556>
23 730
24
25 [42] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recogni-
26 tion, 2016 IEEE Conference on Computer Vision and Pattern Recognition
27 (CVPR) (2016) 770–778.
28
29
30 [43] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick,
31 S. Guadarrama, T. Darrell, Caffe: Convolutional Architecture for Fast Fea-
32 735 ture Embedding, in: Proceedings of the ACM International Conference on
33 Multimedia, MM '14, Orlando, FL, USA, November 03 - 07, 2014, 2014,
34 pp. 675–678.
35
36
37
38
39 [44] A. Borji, L. Itti, State-of-the-art in Visual Attention Modeling, IEEE
40 740 Transactions on Pattern Analysis and Machine Intelligence 35 (1) (2013)
41 185–207.
42
43
44
45 [45] S. Marat, T. Phuoc, L. Granjon, N. Guyader, D. Pellerin, A. Guerin-
46 Dugue, Modelling spatiotemporal saliency to predict gaze direction for
47 short videos., International Journal of Computer Vision (82) (2009) 231–
48 243.
49 745
50
51
52 [46] D. Purves, G. J. Augustine, D. Fitzpatrick, L. C. Katz, A. S. LaMantia,
53 J. O. McNamara, S. M. Williams, Neuroscience, 2nd edition, Sinauer As-
54 sociates, 2001.
55
56
57
58
59
60
61
62
63
64
65

- 1
2
3
4
5
6
7
8
9 [47] M. A. Fischler, R. C. Bolles, Random sample consensus: a paradigm for
10 750 model fitting with applications to image analysis and automated cartogra-
11 phy, *Magazine Communications of the ACM* 24 (6) (1981) 381–395.
12
13
14 [48] A. Krizhevsky, Learning multiple layers of features from tiny images, Ph.D.
15 thesis, University of Toronto (2009).
16
17
18 [49] S. Chaabouni, J. Benois-Pineau, F. Tison, C. Ben Amar, Prediction of vi-
19 755 sual attention with deep CNN for studies of neurodegenerative diseases, in:
20 2016 14th International Workshop on Content-Based Multimedia Indexing
21 (CBMI), 2016, pp. 1–6.
22
23
24 [50] J. Yosinski, J. Clune, Y. Bengio, H. Lipson, How transferable are fea-
25 tures in deep neural networks?, in: Z. Ghahramani, M. Welling, C. Cortes,
26 760 N. Lawrence, K. Weinberger (Eds.), *Advances in Neural Information Pro-
27 cessing Systems 27*, Curran Associates, Inc., 2014, pp. 3320–3328.
28
29
30
31
32 [51] M. D. Zeiler, R. Fergus, Visualizing and Understanding Convolutional Net-
33 works, in: D. Fleet, T. Pajdla, B. Schiele, T. Tuytelaars (Eds.), *Com-
34 puter Vision ECCV 2014: 13th European Conference, Zurich, Switzer-
35 land, September 6-12, 2014, Proceedings, Part I*, Springer International
36 765 Publishing, Cham, 2014, pp. 818–833.
37
38
39
40 [52] G. Mesnil, Y. Dauphin, X. Glorot, S. Rifai, Y. Bengio, I. J. Goodfel-
41 low, E. Lavoie, X. Muller, G. Desjardins, D. Warde-Farley, P. Vincent,
42 A. Courville, J. Bergstra, Unsupervised and Transfer Learning Challenge:
43 a Deep Learning approach, in: *JMLR W& CP: Proceedings of the Unsu-
44 770 pervised and Transfer Learning challenge and workshop, Vol. 27*, 2012, pp.
45 97–110.
46
47
48
49
50 [53] M. Marszałek, I. Laptev, C. Schmid, Actions in Context, in: *2009 IEEE
51 Conference on Computer Vision and Pattern Recognition*, pp. 2929–2936.
52
53
54 775 [54] S. Mathe, C. Sminchisescu, Actions in the Eye: Dynamic Gaze Datasets
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9 and Learnt Saliency Models for Visual Recognition, *IEEE Transactions on*
10 *Pattern Analysis and Machine Intelligence* 37 (7) (2015) 1408–1424.

11
12 [55] L. Itti, CRCNS Data Sharing: Eye movements during free-viewing of natu-
13 ral videos, in: *Collaborative Research in Computational Neuroscience An-*
14 *780* *annual Meeting*, Los Angeles, California, 2008.

15
16
17
18 [56] S. Ioffe, C. Szegedy, Batch normalization: Accelerating deep network train-
19 ing by reducing internal covariate shift, 2015, pp. 448–456.

20
21
22 [57] L. Itti, C. Koch, E. Niebur, A Model of Saliency-Based Visual Attention for
23 Rapid Scene Analysis, *IEEE Transactions on Pattern Analysis and Machine*
24 *785* *Intelligence* 20 (11) (1998) 1254–1259.

25
26
27 [58] X. Hou, J. Harel, C. Koch, Image Signature: Highlighting Sparse Salient
28 Regions., *IEEE Trans. Pattern Anal. Mach. Intell.* 34 (1) (2012) 194–201.

29
30
31 [59] J. Harel, C. Koch, P. Perona, Graph-based visual saliency, in: *Advances in*
32 *neural information processing systems* 19, MIT Press, 2007, pp. 545–552.

33
34
35 [60] H. J. Seo, P. Milanfar, Static and space-time visual saliency detection by
36 *790* self-resemblance, *Journal of Vision* (9(12):15) (2009) 1–27.

37
38
39 [61] O. Le Meur, T. Baccino, Methods for comparing scanpaths and saliency
40 maps: strengths and weaknesses, *Behavior Research Methods* 45 (1) 251–
41 266.
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58