

An Integration of Bottom-up and Top-Down Salient Cues on RGB-D Data:

Saliency from Objectness vs. Non-Objectness

Nevrez Imamoglu · Wataru Shimoda · Chi Zhang · Yuming Fang ·
Asako Kanezaki · Keiji Yanai · Yoshifumi Nishida

Preprint. This work includes the accepted version content of the paper published in
Signal Image and Video Processing (SIViP), Springer, Vol. 12, Issue 2, pp 307-314, Feb 2018.
DOI <https://doi.org/10.1007/s11760-017-1159-7>

Abstract Bottom-up and top-down visual cues are two types of information that helps the visual saliency models. These salient cues can be from spatial distributions of the features (space-based saliency) or contextual / task-dependent features (object based saliency). Saliency models generally incorporate salient cues either in bottom-up or top-down norm separately. In this work, we combine bottom-up and top-down cues from both space and object based salient features on RGB-D data. In addition, we also investigated the ability of various pre-trained convolutional neural networks for extracting top-down saliency on color images based on the object dependent feature activation. We demonstrate that combining salient features from color and dept through bottom-up and top-down methods gives significant improvement on the salient object detection with space based and object based salient cues. RGB-D saliency integration framework yields promising results compared with the several state-of-the-art-models.

Keywords Salient object detection · Multi-model saliency · Saliency from objectness

This paper is based on the results obtained from a project commissioned by the New Energy and Industrial Technology Development Organization (NEDO), Japan.

N. Imamoglu, A. Kanezaki, and Y. Nishida
National Institute of advanced Industrial Science and Technology, Tokyo, Japan
E-mail: nevrez.imamoglu@aist.go.jp

C. Zhang and Y. Fang
Jiangxi University of Finance and Economics, Nanchang, China

W. Shimoda and K. Yanai
The University of Electro-Communications, Tokyo, Japan

1 Introduction

Visual attention is an important mechanism of the human visual system that assists visual tasks by leading our attention or finding relevant features from significant visual cues [1,2,3,4]. Perceptual information can be classified as bottom-up (unsupervised) and top-down (supervised or prior knowledge) visual cues. These salient cues can be from spatial distributions of the features (space-based saliency) or contextual / task dependent features (object based saliency) [1,2,3,4].

Many researches have been done on computing saliency maps for image or video analysis [5,6,7,8,9,10]. Saliency detection models in the literature generally demonstrate computational approaches either in bottom-up or top-down separately without integration of spatial and object based saliency information from image and 3D. Therefore, in this work, we introduce a multi-modal salient object detection framework that combines bottom-up and top-down information from both space and object based salient features on RGB-D data (Fig.1).

In addition, regarding top-down saliency computation on color-images, we investigate salient object detection capability of various Convolutional Neural Networks (CNNs) trained for object classification or semantic object segmentation on large-scale data. Unlike the state-of-the-art deep-learning approaches to achieve salient object detection, we simply take advantage of the pre-trained CNN without the need of additional supervision to regress CNN features to ground truth saliency maps. The assumption is that prior-knowledge of CNNs on known objects can help us to detect salient objects, which are not included as trained object classes of the networks. We demonstrate that this can be done

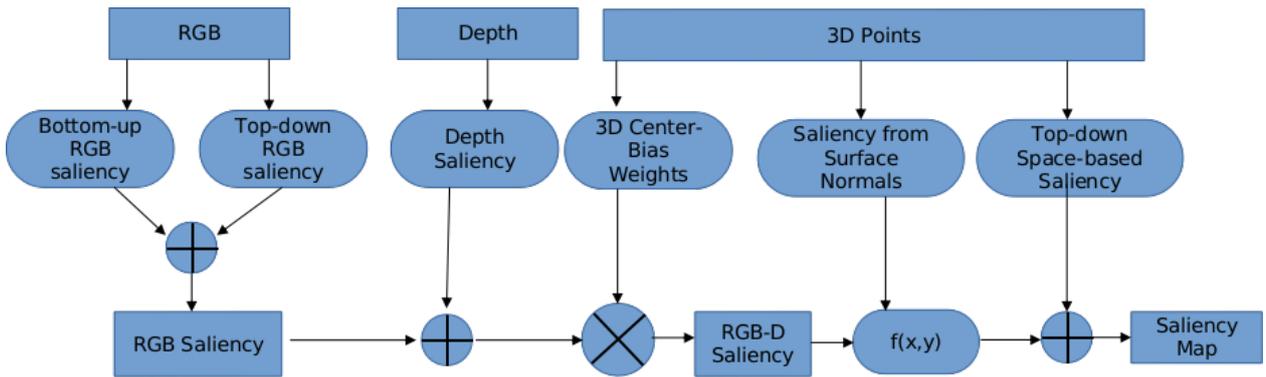


Fig. 1 Flowchart of the proposed saliency computation for integrating various perceptual salient cues

by object dependent features from both objectness and non-objectness scores of the CNNs.

Finally, we examined integration of salient cues from depth and 3D point cloud information based on the spatial distribution of the observed scene in space. First, we use patch similarity based saliency computation on depth image for combining salient cues from color and depth images. Then, we apply weighting operation on color-depth salient cues based on two information: i) the distributions of normal vectors, and ii) center-bias weighting from joint 3D distribution of the observed scene as an improvement to the separate calculation on 2D image and depth demonstrated in [11]. Finally, as top-down space-based saliency [12], spatial working memory of the robot is also included to support the attention cues based on the changes detected in the environment.

In summary, combining color and depth information through bottom-up and top-down processes yielded significant improvement on salient object detection task from both space based and object based approaches. Our experimental evaluations on two different data-set from [12] and [13] demonstrate promising results compared with the several state-of-the-art-models.

2 Related Works

Inspired by the studies on attention mechanism such as [3], first bottom-up (unsupervised) computational model of saliency was developed by Itti et. al [6,4], which was computing center-surround differences in multi-scale pyramid structure of images processed with various filters. Many researches followed the path of [3] for computing saliency maps with bottom-up and top-down approaches on color and depth images [4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,19,20,21,22,23,24]. Among these models, bottom-up approaches rely on finding center-surround differences or contrast based on pixel

or patch similarities such as WT [14], CA [15], SF [17], PCA [18], Y2D [19,11]. Some of the works applied prior knowledge to obtain improved saliency results such as center-bias weighting Y3D [19,11,12] or foreground/background location prior on super-pixels such as MR [16], RBD [20], RGBD [13].

Top-down saliency on color images have been explored in many works by taking advantage of bottom-up low-level features and/or high-level top-down features as training data. Supervised approaches on training data are applied to regress these extracted features to the ground truth salient object regions [7,8,9,10]. After the recent developments and success of deep networks on image classification or segmentation tasks [25,26,27,28], similar methods were also implemented on features of CNNs trained for object detection or classification tasks as in [21,22,23,24]. For example, MDF [24] extracts multi-scale features from a deep CNN trained for large-scale object classification, then these features are used in an end-to-end (pixel-level) training with ground-truth saliency maps to yield saliency predictions for each pixel. However, these processes on trained CNNs require fine-tuning or another training to regress CNN features to infer saliency maps [21,22,23,24].

Object or class dependent top-down saliency maps from color images can be obtained from objectness by using backward propagation on weakly supervised CNN models [29,30], which are generally trained for recognizing objects on the scene. In addition to the weakly supervised approaches, fully supervised models [27,28] (i.e. fully convolutional neural networks) for semantic image segmentation result in several score maps to show the likelihood of representing each object class in different score or activation maps. However, all these models [27,28,29,30] give likelihood maps or score maps for each class separately as CNN outputs. Therefore, in this work, we extend these works simply by combining these objectness based maps to one saliency map

to investigate how these simple pooling handles general salient object detection tasks. In addition to objectness based saliency, we also show that score values for the background scene can aid salient object detection by simply using the negative likelihood of the background data as non-objectness map of CNNs. Our experiments demonstrate that these feed-forward computations from fully supervised models [27,28] or backward process (gradients from back-propagation for CNN layers) from weakly supervised models [29,30] can be representative of salient features. The intuition is that knowledge on the learned objects or learned features from CNNs can help to detect salient objects in the scene despite being unknown by the CNN as an object class. In contrast to the works in [21,22,23,24], the advantage of proposed simple usage of likelihood maps to obtain saliency is not having any additional training to regress CNN features to ground truth saliency maps.

Saliency features from depth [11,13,19] or 3D [11,12] information can improve the detection of salient cues on the perceived scene in addition to the color saliency. Saliency of depth images can be obtained similar to the color image saliency computation [13,19]. For example, center bias-weighting can be applied to improve RGB-D saliency computation by applying center-bias on 2D image and depth [11] separately as in [19]. However, this can also be done in 3D jointly so we implemented center-bias weighting from joint 3D distribution of the observed scene to enhance saliency map computation. If there is a prior information on the spatial distribution of the scene as mean of the memory of the observer, detected changes in the environment can also affect the attention cues such as new objects in the scene, positional changes of the objects, and etc. [35,36,37]. However, due to difficulty of obtaining data that enable prior information on the scene, it is not very common to take advantage of top-down space-based saliency based on finding the spatial changes in the environment as in [12]. It is possible if only the observer has the memory of the observed scene so salient cues can be obtained through visual working memory by finding the changes in the spatial arrangements. If there is an object as a change in the environment, space-based saliency will be enough to detect salient object. However, if there are multiple new objects, new spatial arrangement of the environment, sensory error during the computation of changes, or etc., other salient cues will also be crucial to obtain reliable salient object detection results.

3 Multi-model salient object detection

Multi-modal saliency integration (Fig.1) consists of fusion of attention cues in RGB, Depth, and 3D data. Proposed saliency integration can be given as:

$$\mathbf{S} = F(\mathbf{S}_{RGBD} + \mathbf{S}_{SbS}) \quad (1)$$

$$\mathbf{S}_{RGBD} = F((\alpha \times \mathbf{S}_{RGB} + (1 - \alpha) \times \mathbf{S}_D) \times \mathbf{W}_{cb})^{\mathbf{S}_N} \quad (2)$$

$$\mathbf{S}_{RGB} = \alpha \times \mathbf{S}_{RGB}^{TD} + (1 - \alpha) \times \mathbf{S}_{RGB}^{BU} \quad (3)$$

\mathbf{S} is the proposed saliency map, \mathbf{S}_{RGBD} is the color saliency (\mathbf{S}_{RGB}) and bottom-up depth saliency (\mathbf{S}_D) integration obtained as in Eq.2, \mathbf{w}_{cb} is the 3D center-bias weighting, \mathbf{S}_N is salient cues obtained from the statistical distribution of the surface normal calculated for each image pixel, \mathbf{S}_{SbS} is the top-down space-based saliency that requires visual memory of the environment, $F(\cdot)$ is the normalization function to scale the saliency values within 0-1 range, $\alpha = 0.7$ is empirically selected for weighted averaging process to give top-down salient cues more impact on the integration through averaging with a bottom-up saliency model. Saliency of color images (\mathbf{S}_{RGB} in Eq.2 and Eq.3) is obtained by two types of information; i) top-down RGB saliency (\mathbf{S}_{RGB}^{TD}), and ii) bottom-up RGB saliency (\mathbf{S}_{RGB}^{BU}).

3.1 Saliency from color images

Proposed computation of top-down color saliency maps and explanation of the bottom-up saliency model used in this paper (see Eq.3) will be explained in the following sections.

3.1.1 Top-down image saliency from objectness and non-objectness in CNN models

Weakly or fully supervised convolutional neural networks (CNNs) can handle the process of salient object detection on a given scene even though given network may not be trained to classify the objects in the scene. However, feature representations of learned objects from the training data-set can help to represent similar objects on saliency detection task whether the CNN model knows the object/s of the scene or not. Then, we investigated efficiency of weakly and fully supervised models by introducing simple top-down saliency map computations using objectness or non-objectness values from pre-trained networks in the literature.

Saliency computation with fully supervised CNNs through objectness and non-objectness: Two fully supervised CNNs (DeconvNet [27] and SegNet [28]) are tried to obtain top-down color image saliency, in which both CNNs are trained with PASCAL VOC data-set [31] for semantic segmentation of 20 objects. Both models give object likelihood values at pixel level for each of 20 classes (i.e. objectness maps) and 1 likelihood map for the background class (i.e. non-objectness). In this part, we will demonstrate that top-down saliency from color images can be generated by both of these objectness or non-objectness likelihood maps. By using the objectness map from the DeconvNet [27] and SegNet [28], salient features can be obtained as below:

$$\mathbf{S}_{RGB}^{TD} = F \left(\frac{1}{C} \times \sum_{c=1}^C \mathbf{O}_c \right)^\lambda \quad (4)$$

$$\lambda = F \left(\arg \max_c \left(\left(\frac{\mathbf{O}_c - \mu_c}{\sigma_c} \right)^2 \right) \right)^{\frac{1}{2}} \quad (5)$$

In Eq.4, c is the object or class index, C is the number of objects/classes that can be recognized by the CNN model, \mathbf{O}_c is the score map for each class c , λ (given in Eq.5) is the parameter for amplifying by taking the maximum of the normalized objectness value (see Eq.5) over all classes, μ_c is the mean objectness of the class c , σ_c is the standard deviation of \mathbf{O}_c for each object score maps.

In addition to the saliency through objectness, we also create saliency by using the negative non-objectness likelihood map in a very simple way. Again, we use both DeconvNet [27] and SegNet [28] to provide non-objectness likelihood map. Top-down saliency from non-objectness also gives quite good results even though it is not the best (see experimental results section). So, non-objectness based top-down salient features can be given as below, where \mathbf{NO}_c is the non-objectness likelihood map scaled to 0-1 range to obtain, \mathbf{S}_{RGB}^{TD} , top-down object saliency through non-objectness.

$$\mathbf{v}_i^c = F_{\text{scale}}(-\mathbf{NO}_c) \quad (6)$$

Saliency computation with weakly CNNs through objectness: One of the main problems of fully supervised CNN models as stated in [30] is to require a large set of training data since ground truth labels should be assigned for each pixel of each image in the training data-set. Therefore, saliency computations using weakly supervised models are demanded, and have big advantage on creating large training data. Unlike DeconvNet or SegNet with only 20 classes of object recognition capability, it is easier to create large training data for weakly

supervised models with many object category such as ImageNet with 1000 object class.

Saliency through back-propagation process in [29, 30] demonstrated that regions which have high score derivatives respond to object location. In other words, these regions are related to object based attention cues. Therefore, we use VGG16 CNN model trained with ImageNet by Simonyan et al. [25, 29]. Simonyan et al. [29] regarded the derivatives of the class score with respect to the input image as class saliency maps. However, we use the derivatives of relatively upper intermediate layers which are expected to retain more high-level semantic information by extending our previous work from Shimoda et al. [30]. The difference is to obtain saliency for all objects in one saliency map rather than computing class level saliency map for each object class category. Finally, we average them to obtain one saliency map. The procedure can be given as follows [30, 29]:

$$\mathbf{v}_i^c = \frac{\partial \mathbf{S}_c}{\partial \mathbf{L}_i} \quad (7)$$

$$m_{i,x,y}^c = \max(\mathbf{w}_{i,h_i(x,y,k)}^c) \quad (8)$$

$$g_{x,y}^c = \frac{1}{L} \sum_{i=1}^L \tanh(a \cdot m_{i,x,y}^c) \quad (9)$$

$$\mathbf{S}_{TD}^{RGB} = \frac{1}{K} \sum_{c=1}^K g_{x,y}^c \quad (10)$$

The class score derivative \mathbf{v}_i^c (Eq.7) of a feature map of the i^{th} layer ($i=\{3,4,5\}$) is the derivative of class score with respect to the layer. \mathbf{v}_i^c can be computed by guided back-propagation (GBP) [32, 30] instead of back propagation (BP). In the GBP, only positive loss values are propagated back to the previous layers through ReLUs. GBP can emphasize salient cues on the objects such as edges or boundaries of the objects along with the textures [30]. After obtaining \mathbf{v}_i^c , we up-sample it to \mathbf{w}_i^c with bi-linear interpolation so that the size of a 2-D map of becomes the same as an input image as in Eq.8, where $h_i(x, y, k)$ is the index of the element of \mathbf{w}_i^c , k represents kernel. $m_{i,x,y}^c$ is aggregated for each target layer to obtain saliency feature maps $g_{x,y}^c$ as in Eq.9. In Eq.9, a is a scaling factor defined as 3 as in [30], L is the number of layer to aggregate. Then, in contrast to distinct class saliency maps [30], here, we introduce to combine all class dependent salient cues to obtain an objectness based final saliency map. For this reason, top $k = 3$ class of recognition result as score maps are combined to obtain over all top-down saliency (Eq.10).

3.1.2 Bottom-up image saliency from low-level features

Saliency from low-level features of color images, we use our work in [14] since this part is not the main contribution of this study. However, we demonstrate that including bottom-up salient cues can improve the saliency obtained from color images while combining with the top-down salient cues. This model uses Wavelet Transform (WT) [14] to obtain local and global salient features. WT model [14] only relies on bottom-up low-level features representing attention areas such as edges and color contrasts without any prior information or knowledge. But, any bottom-up model would be replaced in this framework to investigate which bottom-up and top-down models are more complementary to each other.

3.2 Salient cues from depth image and 3D data for RGB-D saliency

Depth Saliency: Depth saliency, \mathbf{S}_D , is computed with same approach as described by Yuming et. al.[11]. This approach takes advantage of comparing patch features with the surrounding patches by distance weighting, where model utilizes DCT based saliency computation. Similar to the bottom-up color saliency calculation of the same work.

Other than the color and depth saliency, we can take advantage of 3D spatial distribution of the scene for improving the saliency computation with additional information such as 3D center-bias weighting and salient cues from surface normal.

Center Bias Weighting: Yuming et. al. [11] showed that integration of center bias theorem to scene content can improve the performance of saliency maps. They used two types of center bias: i) on image space that applies center bias on pixel indexes regarding X and Y of the scene, and ii) on the range of depth values. However, if 3D data are available for each pixel, center bias weights can be obtained in 3D space jointly in a more uniform space rather than having center bias weighting on image and depth separately. Therefore, in this work, we utilized adaptive 3D center bias weighting depending on the visible depth information. The center bias weighting can be calculated as below:

$$\mathbf{W}_{CB} = e^{\left(c_h \times \frac{\mathbf{h}}{2\sigma^2} - c_v \times \frac{\mathbf{v}}{2\sigma^2} - \frac{(\mathbf{d} - \min(\mathbf{d} | \mathbf{d} > 0))}{2\sigma^2}\right)} \quad (11)$$

\mathbf{W}_{CB} (Eq.11) is the center-bias weight matrix. In Eq.11, \mathbf{h} , \mathbf{v} , and \mathbf{d} are 3D data matrix, \mathbf{d} is depth, \mathbf{h} is height or vertical axis, and \mathbf{v} is for the width or horizontal axis. All zero values of \mathbf{d} are assigned to minimum non-zero value in \mathbf{d} . σ is the deviation around

the focus area regarding the center-bias, which is defined as $\sigma = \eta \times \max(\mathbf{d})$, and η is empirically assigned to 0.25. c_v and c_h are the scaling constant on horizontal and vertical points to give more priority to depth on 3D center-bias weighting, and they are also used during normal calculations for smoothing the the data to avoid non-number values or noise since Kinect sensor can not measure close distance and it does not give accurate values at more than 5 meters depth measurement.

Weighting from the distribution of Surface Normal Vectors: After we obtain center-bias weighting to improve attention cues, we also find salient cues for 3D points by checking the distribution of the surface normal of each point. First, we calculate normal of each 3D point by creating a surface with the points within a defined radius [33]. Then, inspired by [34], for each point represented with the surface normal, we calculate the Mahalanobis distance (Eq.12) to the cluster of all 3D points (i.e. distribution of the normal vector of all 3D points).

$$S_N = \mathbf{F}_{\text{scale}} (\mathbf{n} \times \mathbf{\Sigma}^{-1} \times \mathbf{n}^T) \quad (12)$$

\mathbf{n} is the 3D normal vectors of each point corresponding to each pixel at color image, $\mathbf{\Sigma}$ is co-variance matrix of \mathbf{n} . To remove the noise on the \mathbf{S}_N , we apply 2D median filter with $\{5,5\}$ window after reshaping the \mathbf{S}_N to 2D saliency map. Also, we do enhancement by increasing the saliency values of points around the peak salient cues with values higher than 0.8 [21, 24].

3.3 Top-down space-based saliency from 3D points

We apply our previous work [12] for top-down space-based saliency model relying on detecting changes in the environment from 3D Kinect observations. Data-set to create space-based saliency is obtained from a mobile robot monitoring system, which is established with Pioneer 3-DX mobile robot with a Laser Range Finder (LRF) for localization and mapping and a Kinect sensor on a rotating platform for observing the scene [12]. Local Kinect point cloud scene is projected on the 2D global map (i.e. obstacles and free regions in the room). Global map as the spatial memory of the robot is created by using SLAM [38,39] with LRF sensor [12]. Using the changes in the environment based on the projected data and depth information, space-based attention map is obtained to be the top-down salient cues from 3D points (see [12] for details). The performance improvement through the top-down space-based saliency are given in experimental results.

On the other hand, change based saliency is only reliable when we have prior information of the envi-

Table 1 Evaluation of the objectness/non-objectness color saliency maps using Area Under Curve (AUC) metric

DO	DNO	SO	SNO	GBP	DOC	SNOC	GBPC
0.8962	0.8826	0.7424	0.9044	0.9256	0.9281	0.9219	0.9368

ronment as stated before. For example, if a mobile observer (e.g. robot) enters a new environment with no prior knowledge or if there is a sensory error in localization, change detection based saliency will not be reliable and change may not be the only needed attention cues. Hence, we need other salient cues to analyze the scene. In sum, to define focus of attention in the scene, we need to take advantage of a multi-model saliency computation framework as proposed in this work. These bottom-up and top-down salient cues includes features obtained from color image, depth image, and 3D data. In the following section, we demonstrate that our multi-modal salient feature fusion can give very reliable results for salient object detection.

4 Experimental Results

Proposed multi-modal saliency framework is tested in two different data-sets. First one is RGB-D saliency data-set (RGB-D-ECCV2014) used for evaluation of salient object detection work in [13]. RGB-D-ECCV2014 data-set, 1000 color images and depth data, is used to compare our proposed saliency framework with the several state of the art models. However, in this comparison, we exclude change detection related top-down space-based saliency since RGB-D-ECCV2014 data-set does not have any information for any reference environment memory to aid comparison of current spatial observation with past state or spatial placement of the scene. For the second data-set (ROBOT-TCVA2015), we used the data from [12] with RGB color images and 3D Kinect data, which is recorded from a mobile robot while a person was doing various indoor activities in a room. This data also includes 2D global map of the environment with robot pose and Kinect pose in the room for each frame. So, this information can help to project local Kinect data on global 2D map to detect possible attention changes in the room.

Evaluation of top-down saliency on color images: We introduce top-down image saliency computation from the fully-supervised (DeconvNet [27] and SegNet [28]) or weakly-supervised (VGG-16 [32,30]) pre-trained models (see Section 3.1). In this section, we compare proposed DO and DNO (objectness and non-objectness based saliency using DeconvNet [27]), SO and SNO (objectness and non-objectness based saliency using Seg-

Net [18]), and GBP (objectness based saliency using Guided Back-propagation [32,30]). Regarding the non-objectness based saliency using GBP, we tried to train the weakly supervised CNN model by including class output for the background or non-objectness class as in fully supervised models; however, we could not get reliable results from these training attempts. Therefore, we did not include trials for saliency from non-objectness and GBP combinations in weakly supervised models.

Comparison of AUC results based on the ECCV2014 RGBD data-set are given in Table.1. An interesting observation is seen on the SegNet [28], in which non-objectness based saliency (SNO) outperformed objectness based saliency (SO) on SegNet [28] model. Moreover, SO results in considerably poor results as a top-down saliency approach through objectness. This observation on SO performance shows that SegNet [28] still open to improvement for object representation and classification, which can be achieved by introducing multi-task learning as a future work. Among these five proposed top-down color image saliency trials, GBP, SNO and DO had the top three AUC values. So, we will use these top three performing model for further evaluation while we apply our multi-model saliency framework. Then, for comparison, we combined each of the top-down saliency approaches GBP, SNO, and DO with bottom-up model to obtain color image saliency map \mathbf{S}_{RGB} (Eq.3). We refer these color saliency map variations using GBP, SNO, and DO as GBPC, SNOC, and DOC, respectively.

Comparison using RGB-D-ECCV2014 data-set: In this section, we will demonstrate the effectiveness of the proposed saliency model as an RGB-D saliency framework without the use of top-down space-based saliency (see Fig.1), when prior environment data (memory) is not available. We use Area Under Curve (AUC) measure obtained from the Receiver Operating Characteristic (ROC) curve [40,14,11] as evaluation metric to compare the performances of the state-of-the-art models and proposed framework. We compared our results with various RGB and RGB-D based saliency models. The saliency works, WT [14], CA [15], MR [16], SF [17], PCA [18], Y2D [19], Y3D [11], RBD [20], RGBD [13], MDF [24] are included in our comparisons.

On RGB-D-ECCV2014 data-set, it is clear that proposed RGB-D saliency performances improved while

Table 2 Area Under Curve (AUC) based evaluation of the selected models in the literature and our proposed saliency computations (SNOP, DOP, GBPP) within our multi-modal framework

SF	WT	CA	Y2D	RGBD	PCA	Y3D
0.7637	0.8453	0.8488	0.8859	0.9033	0.9089	0.9094
RBD	MR	MDF	SNOP	DOP	GBPP	
0.9170	0.9283	0.9328	0.9339	0.9398	0.9491	

**Fig. 2** (a) sample color images with their (b) ground truths, and saliency results of (c) our framework with GBPP, and other selected models (d) MDF [24] (e) MR [16] (f) RBD [20] (g) RGBD [13] (h) Y3D [11] (i) PCA [18] (j) SF [17] (k) CA [15]

our GBPP, SNOP, and DOP (see Table.2) having AUC values 0.9491, 0.9339, and 0.9398, which are higher compared to their color only respective proposed variants (see Table.1). Among the state-of-the-art models, MDF [24] has the best AUC performance, followed by MR [16]. In summary, proposed models (SNOP, DOP, GBPP), outperformed the state-of-the-art saliency models compared with. And in overall evaluation, our GBPP using weakly supervised CNN trained for 1000 object class has the best AUC performance on this data-set. In Fig.2, some color images and their corresponding saliency maps from the proposed GBPP and state-of-the-art models are given.

Comparison using ROBOT-TCVA2015 data-set: We validated that proposed framework gives promising results on a RGB-D public data-set (RGB-D-ECCV2014 [22]) in previous experimental results. However, we were not able to use top-down space based saliency integration previously. Therefore, we will express the improvement of selective attention cues depending on the spatial changes in the environment.

We will use ROBOT-TCVA2015 data-set [12] for this purpose, which consists of frames recorded in a room with a subject doing daily activities. The activities include tasks such as standing, sitting, walk-

ing, bending, using cycling machine, walking on treadmill, and lying-down, and etc. Since the environment is similar within an activity, we selected 10 frames randomly for representing all activities in these test samples. Then, we manually created ground truth binary images, where the subject is the focus of attention. Then, we tested proposed framework fully with the salient cues obtained from the changes in the environment. ROBOT-TCVA2015 includes Kinect data and global map with robot pose recorded for all frames. So, this allows us to create top-down space-based saliency by projecting the local Kinect data to global map to find changes and attention values to these changes [12].

Proposed framework (Fig.1) can be implemented fully by combining salient cues from color images (bottom-up and top-down saliency of color images), depth (bottom-up depth saliency), center-bias weighting, normal vector (saliency weighting from the distribution normal vectors), and spatial changes in the environment (top-down space based saliency) as resulting GBPP-SbS in this experiment. For comparison on ROBOT-TCVA2015 data, four best performing state-of-the-art models are selected from the previous analysis, which are MR [16], PCA [18], RBD [20], and MDF [24] models. From our proposed variants, the best performing case, GBPP, is used to compare with other models and also to check



Fig. 3 Saliency results of (a) sample images with (b) ground truth using models: (c) our GBPP-SbS (d) MDF [24] (e) MR [16] (f) RBD [20] (g) PCA [18]

Table 3 Evaluation of the selected models and our GBPP-SbS using Area Under Curve (AUC) metric

MR	PCA	RBD	MDF	GBPP	GBPP-SbS
0.8060	0.8659	0.7518	0.8657	0.9468	0.9592

the improvement when we combine top-down space-based saliency with GBPP which is labelled as GBPP-SbS. In Fig.3, some sample images, ground truth (GT) salient object (person in the data), and their corresponding saliency examples for some of the state-of-the-art models and our saliency framework (GBPP and GBPP-SbS) are given.

Our proposed GBPP-SbS shows the best AUC performance (0.9592) for the ROBOT-TCVA2015 test data among the all compared models (see Table.3). In this test, AUC performances of the selected state-of-the-art models decrease drastically compared to the test results on RGB-D-ECCV2014 data-set in previous section. Perhaps, real-time data from an uncontrolled environment effected their accuracy on detecting salient cues due to noise and high illumination change conditions in ROBOT-TCVA2015 data. MR [16], PCA [18], RBD [20], and MDF [24] have AUC performances as 0.8657, 0.8659, 0.7518, and 0.8060 respectively. On the other hand, top-down space-based saliency from detected changes improves the AUC performance of the GBPP from 0.9468 to 0.9592 for the proposed GBPP-SbS saliency maps.

5 Conclusion

Proposed work demonstrates a saliency framework that takes advantage of various attention cues from RGB-

D data. The model demonstrated its reliability from two different data-sets compared to the state-of-the-art models. However, even though saliency results on mobile robot data having promising performance, the current framework is not suitable for real-time computation. Therefore, we would like to extend and improve the model for real-time mobile robot surveillance. In summary, we applied proposed saliency integration framework step-by-step to obtain saliency on color images, then RGB-D, and finally RGB-D with top-down space based saliency. Evaluation from AUC metric shows importance of multi-model saliency from both spatial and object based salient cues. Especially, saliency analysis on CNNs from objectness and non-objectness shows interesting findings for these networks trained for object classification or segmentation.

Acknowledgements Dr. Nevrez Imamoglu thanks Dr. Boxin Shi at AIST (Tokyo, Japan) for discussions on 3D data processing.

References

1. G. D. Logan, The CODE theory of visual attention: An integration of space-based and object-based attention, *Psychological Review*, vol.103, no.4, pp.603-649, 1996.
2. R. Desimone and J. Duncan Neural mechanisms of selective visual attention, *Annual Review of Neuroscience*, vol.18, pp.193-222,1995.

3. J. Wolfe, Guided search 2.0: A revised model of guided search, *Psychonomic Bull. Rev.*, vol.1, no.2, pp. 202-238, 1994.
4. L. Itti, Models of bottom-up and top-down visual attention, Ph.D. Dissertation, Dept. of Computat. Neur. Syst., California Inst. of Technol., Pasadena, 2000.
5. L. Zhang and W. Lin, *Selective Visual Attention: Computational Models and Applications*, Wiley-IEEE Press, 2013.
6. L. Itti, C. Koch, and E. Niebur, A model of saliency-based visual attention for rapid scene analysis, *IEEE Transactions on PAMI*, vol.20, no.11, 1998.
7. J. Yang and M.-H. Yang, Top-down visual saliency via joint CRF and dictionary learning, in *Proc. of 2013 IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*, pp.2296-2303.
8. S. Frintrop, *VOCUS: A Visual Attention System for Object Detection and Goal-Directed Search*, Springer, 2006.
9. V. Navalpakkam and L. Itti, An integrated model of top-down and bottom-up attention for optimizing detection speed, in *Proc. of 2006 IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*.
10. A. Borji, Boosting bottom-up and top-down visual features for saliency estimation, in *Proc. of 2012 IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*.
11. Y. Fang, J. Wang, Y. Yuan, J. Lei, W. Lin, P. Le Callet, Saliency-based stereoscopic image retargeting, *Information Sciences*, vol.372, pp.347-358, 2016.
12. N. Imamoglu, E. Dorrnzoro, M. Sekine, K. Kita, W. Yu, Spatial visual attention for novelty detection: Space-based saliency model in 3d using spatial memory, *IPSJ Transactions on Computer Vision and Applications*, vol.7, pp.35-40, 2015.
13. H. Peng, B. Li, W. Xiong, W. Hu and R. Ji, RGBD salient object detection: a benchmark and algorithms, in *Proc. of 2014 European Conference on Computer Vision (ECCV)*, pp.92-109.
14. N. Imamoglu, W. Lin, Y. Fang, "A Saliency Detection Model Using Low-Level Features Based on Wavelet Transform," *IEEE Transactions on Multimedia*, vol.15, issue.1, pp.96-105, 2013.
15. S. Goferman, L. Z. Manor, A. Tal, Context-aware saliency detection, in *Proc. of 2010 IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*.
16. C. Yang, L. Zhang, H. Lu, X. Ruan, M. H. Yang, Saliency detection via graph-based manifold ranking, in *Proc. of 2013 IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*.
17. F. Perazzi, P. Krahenbuhl, Y. Pritch, A. Hornung, Saliency filters: Contrast based filtering for salient region detection, in *Proc. of 2012 IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*.
18. R. Margolin, A. Tal, L. Zelnik-Manor, What makes a patch distinct?, in *Proc. of 2013 IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*.
19. Y. Fang, Z. Chen, W. Lin, C.-W. Lin, Saliency detection in the compressed domain for adaptive image retargeting, *IEEE Transactions on Image Processing*, vol.21, no.9, pp.3888-3901, 2012.
20. Wangjiang Zhu, Shuang Liang, Yichen Wei, Jian Sun, Saliency optimization from robust background detection, *2014 IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*.
21. E. Vig, M. Dorr, and D. Cox. Large-scale optimization of hierarchical features for saliency prediction in natural images, *2014 IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*.
22. R. Zhao, W. Ouyang, H. Li, and X. Wang, Saliency detection by multi-context deep learning, in *Proc. of 2015 IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*.
23. G. Li and Y. Yu, Deep contrast learning for salient object detection, *2016 IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*.
24. G. Li and Y. Yu, Visual saliency based on multi-scale deep features, *2015 IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*.
25. K. Simonyan, A. Vedaldi, A. Zisserman, Very deep convolutional networks for large-scale image recognition, in: *Proc. of 2015 International Conference on Learning Representations*.
26. Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, Caffe: Convolutional architecture for fast feature embedding, *arXiv preprint arXiv: 1408.5093*, 2014.
27. H. Noh, S. Hong, B. Han, Learning deconvolution network for semantic segmentation, in *Proc. of 2015 International Conference on Computer Vision*.
28. V. Badrinarayanan, A. Kendall and R. Cipolla SegNet: A deep convolutional encoder-decoder architecture for image segmentation, *arXiv preprint arXiv:1511.00561*, 2015.
29. K. Simonyan, A. Vedaldi, A. Zisserman, Deep inside convolutional networks: Visualizing image classification models and saliency maps, in *Proc. of 2014 Int. Conference on Learning Representations*.
30. W. Shimoda and K. Yanai, Distinct class saliency maps for weakly supervised semantic segmentation, accepted to appear in *Proc. of 2016 European Conference on Computer Vision (ECCV)*.
31. M. Everingham, L. Van Gool, C. K. Williams, J. Winn, A. Zisserman, The pascal visual object classes (VOC) challenge, *International Journal of Computer Vision*, vol.88, no.2, pp.303-338, 2010.
32. J. T. Springenberg, A. Dosovitskiy, T. Brox, M. Riedmiller, Striving for simplicity: The all convolutional net, in: *Proc. of 2015 International Conference on Learning Representations (ICLR)*.
33. P. S. Rungta, Kinect cloud normals: Towards surface orientation estimation, M.S. Thesis in Computer Science in the Graduate College of the University of Illinois at Urbana-Champaign, 2011.
34. C. H. Lee, A. Varshney, and D. Jacobs, Mesh saliency, in *Proc. of ACM SIGGRAPH*.
35. S. H. Oh and M. S. Kim The role of spatial working memory in visual search efficiency, *Psychonomic Bulletin and Review*, vol.11, no.2, pp.275-281, 2004.
36. M. M. Chun and Y. Jiang, Contextual cueing: Implicit learning and memory of visual context guides spatial attention, *Cognitive Psychology*, vol.36, pp.287-311, 1998.
37. M. M. Chun and J. M. Wolfe, *Visual Attention, Handbook of Sensation and Perception (Chapter 9)*, Goldstein, E.B. (Ed.), pp.273-310, Blackwell Publishing, 2005.
38. Robot Operating System (ROS): (online), available from <http://wiki.ros.org>
39. S. Thrun, and W. Burgard, *Probabilistic Robotics*, The MIT Press Cambridge, 2005.
40. T. Liu, J. Sun, N.-N. Zheng, X. Tang, H.-Y. Shum, Learning to detect a salient object, *2007 IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*.