

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/51925347>

# Context-Aware Saliency Detection

Article in IEEE Transactions on Software Engineering · December 2011

DOI: 10.1109/TPAMI.2011.272 · Source: PubMed

---

CITATIONS

1,331

---

READS

970

3 authors, including:



[Lihi Zelnik-Manor](#)

Technion - Israel Institute of Technology

84 PUBLICATIONS 6,293 CITATIONS

SEE PROFILE

Some of the authors of this publication are also working on these related projects:



Sparse modelling [View project](#)



Collages [View project](#)

# Context-Aware Saliency Detection

Stas Goferman  
Technion

stasix@gmail.com

Lihi Zelnik-Manor  
Technion

lihi@ee.technion.ac.il

Ayellet Tal  
Technion

ayellet@ee.technion.ac.il

## Abstract

We propose a new type of saliency – context-aware saliency – which aims at detecting the image regions that represent the scene. This definition differs from previous definitions whose goal is to either identify fixation points or detect the dominant object. In accordance with our saliency definition, we present a detection algorithm which is based on four principles observed in the psychological literature. The benefits of the proposed approach are evaluated in two applications where the context of the dominant objects is just as essential as the objects themselves. In image retargeting we demonstrate that using our saliency prevents distortions in the important regions. In summarization we show that our saliency helps to produce compact, appealing, and informative summaries.

## 1. Introduction

Please take a look at the images on the top row of Figure 1. How would you describe them? Probably you’d say “a smiling girl”, “a figure in a yellow flower field”, and “a weight lifter in the Olympic games” (or something similar)<sup>1</sup>. Each title describes the essence of the corresponding image – what most people think is important or *saliency*.

A profound challenge in computer vision is the detection of the salient regions of an image. The numerous applications (e.g., [1, 21, 17, 20]) that make use of these regions have led to different definitions and interesting detection algorithms. Classically, algorithms for saliency detection focused on identifying the fixation points that a human viewer would focus on at the first glance [9, 8, 24, 3, 6, 12]. This type of saliency is important for understanding human attention as well as for specific applications such as auto focusing. Others have concentrated on detecting a single dominant object of an image [13, 7, 5]. For instance, in Figure 1, such

<sup>1</sup>These descriptions were obtained by collecting titles given by 12 different people. See samples in the second row of Figure 1

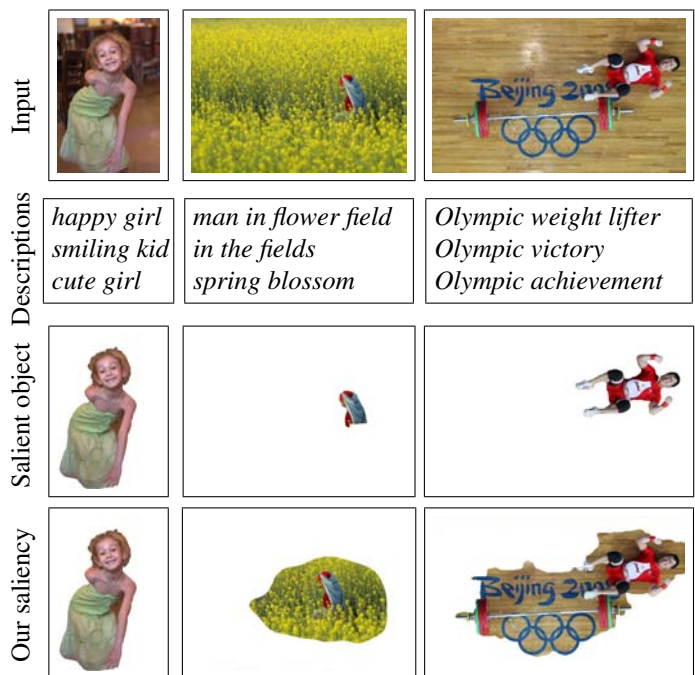


Figure 1. Our context-aware saliency results (bottom) comply with the descriptions that people provided (samples in the second row) for the input images (top). People tend to describe the *scene* rather than the *dominant object*. Classical saliency extraction algorithms aim at the third row, which might miss the essence of the scene. Conversely, we maintain all the essential regions of the image.

methods aim to extract the “girl”, the “figure”, and the “athlete” (third row). This type of saliency is useful for several high-level tasks, such as object recognition [20] or segmentation [18].

There are, however, applications where the context of the dominant objects is just as essential as the objects themselves. Examples include image classification [14], summarization of a photo collection [17], thumbnailing [21], and retargeting [19]. For these applications, the detected regions in Figure 1 should correspond to the titles you gave above. The regions on the bottom row of Figure 1 match these titles better than the regions on the third row.



(a) Input (b) Local [24] (c) Global [7] (d) Local-global [13] (e) Our context-aware  
Figure 2. Comparing different approaches to saliency

This calls for introducing a new type of saliency – context-aware saliency. Here, the goal is to identify the pixels that correspond to the bottom row (and to the titles). According to this concept, the salient regions should contain not only the prominent objects but also the parts of the background that convey the context.

We differentiate between three types of images, as illustrated in Figure 1. In the girl’s case, the background is not interesting, hence, we expect the extracted salient region to coincide with the salient object. In the flower-field’s case, the texture of the flowers is essential for understanding the content. However, only a small portion of it – the portion surrounding the figure – suffices. In the weight lifter’s case, some of the contextual background is vital for conveying the scene. This is not necessarily the portion surrounding the athlete, but rather a unique part of the background (the weights and the olympic logo). Therefore, detecting the prominent object together with naive addition of its immediate surrounding will not suffice.

This paper proposes a novel algorithm for context-aware saliency detection. The underlying idea is that salient regions are distinctive with respect to both their local and global surroundings. Hence, the unique parts of the background, and not only the dominant objects, would be marked salient by our algorithm (e.g., the Olympics logo in Figure 1). Moreover, to comply with the Gestalt laws, we prioritize regions close to the foci of attention. This maintains the background texture, when it is interesting, such as in the case of the flower field in Figure 1.

We demonstrate the utility of our context-aware saliency in two applications. The first is retargeting[1, 19, 16], where we show that our saliency can successfully mark the regions that should be kept untouched. The second is summarization [17, 25, 2, 15, 4], where we demonstrate that saliency-based collages are informative, compact, and eye-pleasing.

The contribution of this paper is hence threefold. First, we introduce principles for context-aware saliency (Section 2). Second, we propose an algorithm that detects this saliency (Section 3) and present results on images of various types (Section 4). Last but not least, we demonstrate the applicability of our saliency (Section 5).

## 2. Principles of context-aware saliency

Our context-aware saliency follows four basic principles of human visual attention, which are supported by psychological evidence [22, 26, 10, 11]:

1. Local low-level considerations, including factors such as contrast and color.
2. Global considerations, which suppress frequently-occurring features, while maintaining features that deviate from the norm.
3. Visual organization rules, which state that visual forms may possess one or several centers of gravity about which the form is organized.
4. High-level factors, such as human faces.

Related work typically follows only some of these principles and hence might not provide the results we desire. The biologically-motivated algorithms for saliency estimation [9, 8, 24, 3, 6, 12] are based on principle (1). Therefore, in Figure 2(b), they detect mostly the intersections on the fence. The approaches of [7, 5] focus on principle (2). Therefore, in Figure 2(c), they detect mostly the drops on the leaf. In [13] an algorithm was proposed for extracting rectangular bounding boxes of a single object of interest. This was achieved by combining local saliency with global image segmentation, thus can be viewed as incorporating principles (1) and (2). In Figure 2(d) they detect as salient both the fence and the leaf, with higher importance assigned to the leaf.

We wish to extract the salient objects together with the parts of the discourse that surrounds them and can throw light on the meaning of the image. To achieve this we propose a novel method for realizing the four principles. This method defines a novel measure of distinctiveness that combines principles (1),(2),(3). As illustrated in Figure 2(e) our algorithm detects as salient the leaf, the water-drops and just enough of the fence to convey the context. Principle (4) is added as post-processing.

### 3. Detection of context-aware saliency

In this section we propose an algorithm for realizing principles (1)–(4). In accordance with principle (1), areas that have distinctive colors or patterns should obtain high saliency. Conversely, homogeneous or blurred areas should obtain low saliency values. In agreement with principle (2), frequently-occurring features should be suppressed. According to principle (3), the salient pixels should be grouped together, and not spread all over the image.

This section is structured as follows (Figure 3). We first define single-scale local-global saliency based on principles (1)–(3). Then, we further enhance the saliency by using multiple scales. Next, we modify the saliency to further accommodate principle (3). Finally, principle (4) is implemented as post-processing.

**3.1 Local-global single-scale saliency:** There are two challenges in defining our saliency. The first is how to define distinctiveness both locally and globally. The second is how to incorporate positional information.

According to principles (1)–(2), a pixel is salient if its appearance is unique. We should not, however, look at an isolated pixel, but rather at its surrounding patch, which gives an immediate context. For now we consider a single patch of scale  $r$  at each pixel. Thus, a pixel  $i$  is considered salient if the appearance of the patch  $p_i$  centered at pixel  $i$  is distinctive with respect to all other image patches.

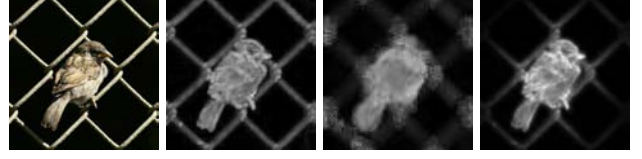
Specifically, let  $d_{color}(p_i, p_j)$  be the Euclidean distance between the vectorized patches  $p_i$  and  $p_j$  in CIE  $L^*a^*b$  color space, normalized to the range  $[0, 1]$ . Pixel  $i$  is considered salient when  $d_{color}(p_i, p_j)$  is high  $\forall j$ .

According to principle (3) the positional distance between patches is also an important factor. Background patches are likely to have many similar patches both near and far-away in the image. This is in contrast to salient patches which tend to be grouped together. This implies that a patch  $p_i$  is salient when the patches similar to it are nearby, and it is less salient when the resembling patches are far away.

Let  $d_{position}(p_i, p_j)$  be the Euclidean distance between the positions of patches  $p_i$  and  $p_j$ , normalized by the larger image dimension. Based on the observations above we define a dissimilarity measure between a pair of patches as:

$$d(p_i, p_j) = \frac{d_{color}(p_i, p_j)}{1 + c \cdot d_{position}(p_i, p_j)}, \quad (1)$$

where  $c = 3$  in our implementation. This dissimilarity measure is proportional to the difference in appearance and inverse proportional to the positional distance. Pixel  $i$  is considered salient when it is highly dissimilar to all other image patches, i.e., when  $d(p_i, p_j)$  is high  $\forall j$ .



(a) Input (b) Scale 1 (c) Scale 4 (d) Final  
Figure 3. The steps of our saliency estimation algorithm

In practice, to evaluate a patch’s uniqueness, there is no need to incorporate its dissimilarity to all other image patches. It suffices to consider the  $K$  most similar patches (if the most similar patches are highly different from  $p_i$ , then clearly all image patches are highly different from  $p_i$ ). Hence, for every patch  $p_i$ , we search for the  $K$  most similar patches  $\{q_k\}_{k=1}^K$  in the image, according to Equation (1). A pixel  $i$  is salient when  $d(p_i, q_k)$  is high  $\forall k \in [1, K]$ . The single-scale saliency value of pixel  $i$  at scale  $r$  is defined as ( $K = 64$  in our experiments):

$$S_i^r = 1 - \exp\left\{-\frac{1}{K} \sum_{k=1}^K d(p_i^r, q_k^r)\right\}, \quad (2)$$

**3.2 Multi-scale saliency enhancement:** Background pixels (patches) are likely to have similar patches at multiple scales, e.g., in large homogeneous or blurred regions. This is in contrast to more salient pixels that could have similar patches at a few scales but not at all of them. Therefore, we incorporate multiple scales to further decrease the saliency of background pixels, improving the contrast between salient and non-salient regions.

For a patch  $p_i$  of scale  $r$ , we consider as candidate neighbors all the patches in the image whose scales are  $R_q = \{r, \frac{1}{2}r, \frac{1}{4}r\}$ . Among all these patches, the  $K$  most similar patches according to Equation (1) are found and used for computing the saliency. Hence, Equation (2) can be rewritten as (where  $r_k \in R_q$ ):

$$S_i^r = 1 - \exp\left\{-\frac{1}{K} \sum_{k=1}^K d(p_i^r, q_k^{r_k})\right\}, \quad (3)$$

Furthermore, we represent each pixel by the set of multi-scale image patches centered at it. Let  $R = \{r_1, \dots, r_M\}$  denote the set of patch sizes to be considered for pixel  $i$ . The saliency at pixel  $i$  is taken as the mean of its saliency at different scales:

$$\bar{S}_i = \frac{1}{M} \sum_{r \in R} S_i^r, \quad (4)$$

where  $S_i^r$  is defined in Equation (3). The larger  $\bar{S}_i$  is, the more salient pixel  $i$  is and the larger is its dissimilarity (in various levels) to the other patches.

In our implementation, we scale all the images to the same size of 250 pixels (largest dimension) and take patches of

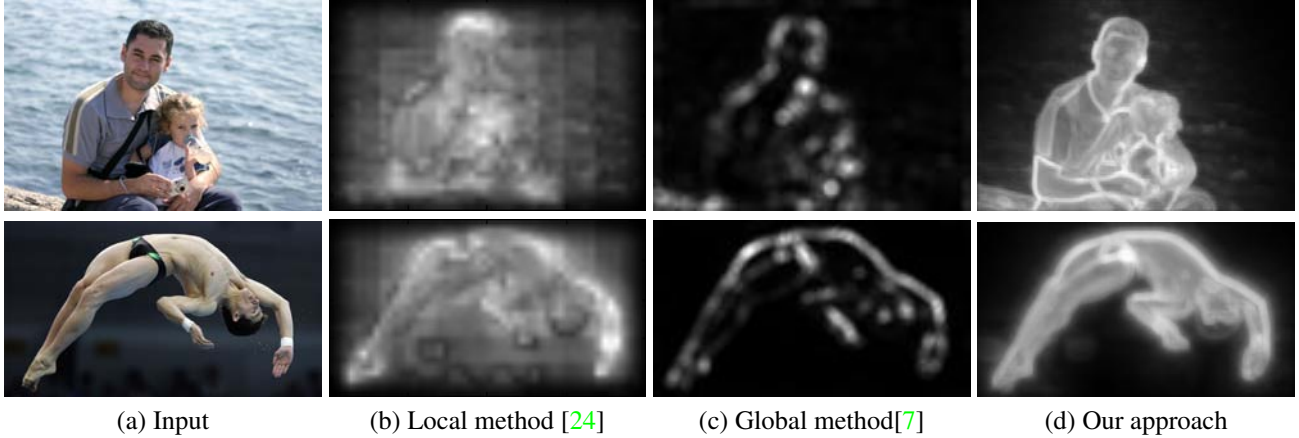


Figure 4. Comparing saliency results on images of a single object over an uninteresting background

size  $7 \times 7$  with 50% overlap. We use four scales:  $R = \{100\%, 80\%, 50\%, 30\%\}$ . The smallest scale allowed in  $R_q$  is 20% of the original image scale.

Figures 3(b)–(c) demonstrate the difference between the saliency maps obtained at different scales. While the fine-scale result detects all the details, including those of the background, the coarse scale result detects mostly the bird.

**3.3 Including the immediate context:** According to Gestalt laws, visual forms may possess one or several centers of gravity about which the form is organized [11] (principle (3)). This suggests that areas that are close to the foci of attention should be explored significantly more than far-away regions. When the regions surrounding the foci convey the context, they draw our attention and thus are salient.

We simulate this visual contextual effect in two steps. First, the most attended localized areas are extracted from the saliency map produced by Equation (4). A pixel is considered attended if its saliency value exceeds a certain threshold ( $\bar{S}_i > 0.8$  in the examples shown in this paper).

Then, each pixel outside the attended areas is weighted according to its Euclidean distance to the closest attended pixel. Let  $d_{foci}(i)$  be the Euclidean positional distance between pixel  $i$  and the closest focus of attention pixel, normalized to the range  $[0, 1]$ . The saliency of a pixel is redefined as:

$$\hat{S}_i = \bar{S}_i(1 - d_{foci}(i)). \quad (5)$$

Note, that the saliency of non-interesting regions, such as blurred or homogeneous regions, remains low, since  $\bar{S}$  of Equation (4) will dominate. However, the saliency of interesting background in the neighborhood of the salient objects will be increased by Equation (5). This explains why in Figure 1 parts of the flower field were detected as salient in the center example, whereas the girl on the left was segmented

accurately. In Figure 3(d) this final step enhances the bird and attenuates the far parts of the background wire.

**3.4 High-level factors:** Finally, the saliency map should be further enhanced using some high-level factors, such as recognized objects or face detection. In our implementation, we incorporated the face detection algorithm of [23], which generates 1 for face pixels and 0 otherwise. The saliency map of Equation (5) is modified by taking the maximum value of the saliency map and the face map.

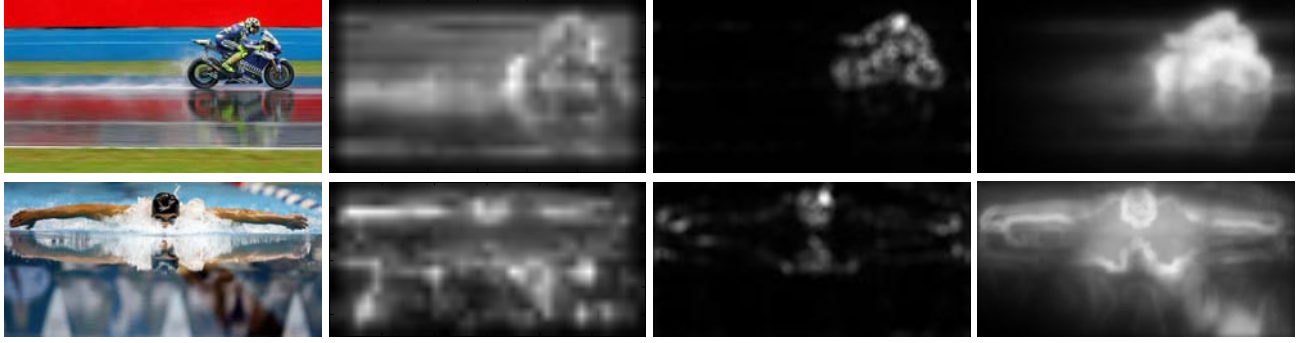
## 4. Results

This section evaluates the results of our approach. Figures 4–6 compare our results with the biologically-inspired local-contrast approach of [24] and the spectral residual global approach of [7]. Later on in Figure 8 we compare our results with the single-object detection of [13].

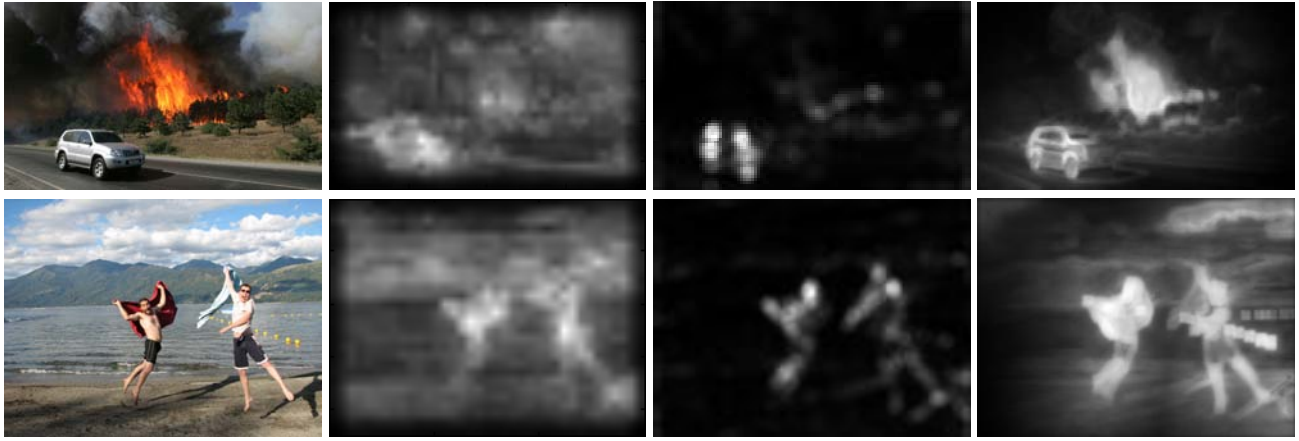
As will be shown next, the method of [24] detects as salient many non-interesting background pixels since it does not consider any global features. The approach of [7] fails to detect many pixels on the prominent objects since it does not incorporate local saliency. Our approach consistently detects with higher accuracy the pixels on the dominant objects and their contextual surroundings. In all the results presented here, our saliency maps were computed using Equation (5) without face detection for a fair comparison.

We distinguish between three cases. The first case (Figure 4) includes images that show a single salient object over an uninteresting background. For such images, we expect that only the object’s pixels will be identified as salient. In [24], some pixels on the objects are very salient, while other pixels – both on the object and on the background – are partially salient as well. In [7] the background is nicely excluded, however, many pixels on the salient objects are





(a) Input (b) Local method [24] (c) Global method[7] (d) Our approach  
Figure 5. Comparing saliency results on images in which the immediate surroundings of the salient object is also salient



(a) Input (b) Local method [24] (c) Global method[7] (d) Our approach  
Figure 6. Comparing saliency results on images of complex scenes

not detected as salient. Our algorithm manages to detect the pixels on the salient objects and only them.

The second case (Figure 5) includes images where the immediate surroundings of the salient object shed light on the story the image tells. In other words, the surroundings are also salient. Unlike the other approaches, our results capture the salient parts of the background, which convey the context. For example the motor-cyclist is detected together with his reflection and part of the race track, and the swimmer is detected together with the foam he generates.

The third case includes images of complex scenes. For instance, Figure 6 shows an image of a car in a fire scene and an image of two cheering guys by the lake and mountains. It can be observed that our approach detects as salient both the vehicle and the fire in the first scene and the guys with part of the scenery in the other one.

To obtain a quantitative evaluation we compare ROC curves on the database presented in [7]. This database includes 62 images of different scenes where ground-truth was obtained by asking people to “select regions where objects are presented”. In part of the images only the dominant object was

marked while in others also parts of the essential context was selected. Even-though this database is not perfectly suited for our task Figure 7 shows that our algorithm outperforms both [7] and [24].

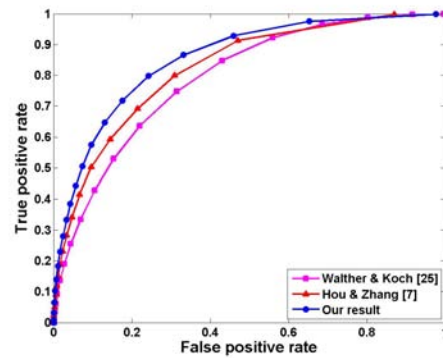


Figure 7. ROC curves for the database of [7].

Methods like [13] are not designed for such complex scenes, but rather for single dominant-object images. We do not have access to their code, hence we cannot show their re-



Figure 8. Comparing our saliency results with [13]. Top: Input images. Middle: The bounding boxes obtained by [13] capture a single main object. Bottom: Our saliency maps convey the story .

sults on Figures 5-6. Instead, comparisons are shown on images from their paper (Figure 8). In [13], a large database of single-object images is presented with impressive extraction results. In the left two images of Figure 8, they successfully extract the "man" and the "bird". Conversely, our saliency maps indicate that the images show "two men talking" (as both are marked salient) and a "bird on a branch feeding its fledglings", hence providing the context. The image of the woman demonstrates another feature of our algorithm. While [13] detect the upper body of the woman

(the black dress is captured due to its salient color), our algorithm marks as salient the entire woman as well as some of the stone wall, thus capturing her posing to the camera.

## 5. Applications

Many applications require saliency maps as input. In this section we show via two applications that our proposed context-aware saliency is beneficial.

### 5.1. Image retargeting

Image retargeting aims at resizing an image by expanding or shrinking the non-informative regions [1, 19, 16]. Therefore, retargeting algorithms rely on the availability of saliency maps which accurately detect all the salient image details.

Using context-aware saliency for retargeting could assure that the dominant objects, as well as their meaningful neighborhoods, will remain untouched in the resized image. Distortions, if and when introduced, will exist only in regions of lower significance.

Seam carving is a popular retargeting technique that repeatedly carves out seams in a certain direction [19]. To get pleasing results, removal/addition of seams should not introduce salient features. The selection and order of seams attempt to protect the content of the image, according to the saliency map. We ran the original code of [19] and compared their results with those produced after replacing their saliency map with ours.

Figure 9 presents a couple of results. Differently from [19], our saliency guarantees that the salient objects (the fish and

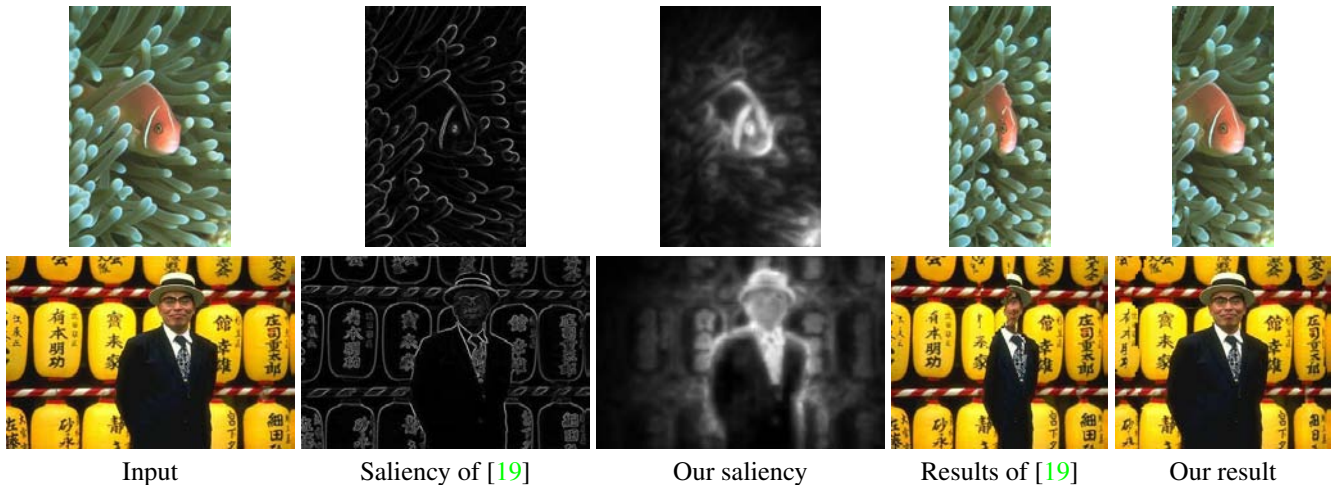


Figure 9. Seam carving of 100 "vertical" lines. The salient objects are distorted by [19] in contrast to our results.

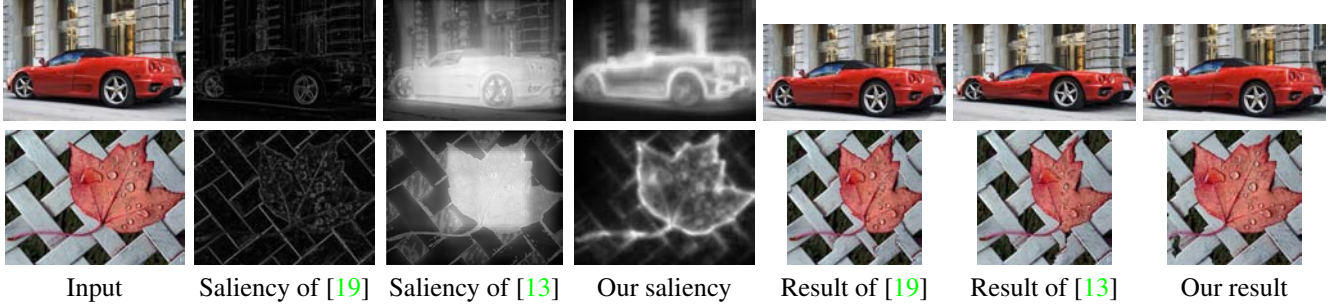


Figure 10. In our saliency maps the details of the car and the leaf are detected more accurately, hence they are not distorted by retargeting.

the man) are not distorted. The improved results can be explained by comparing the saliency maps. In the saliency maps of [19], the background appears important due to the edges it contains. Consequently, the seam carving algorithm prefers to carve out parts of the salient object. On the other hand, our saliency maps differentiate between the non-salient background and the salient object and its close salient background. Both are maintained after resizing, resulting in eye-pleasing images.

Further comparisons are provided in Figure 10, where the saliency map of seam-carving is replaced by that of [13]. In [13] the dominant object is detected, but the object details are not outlined accurately. Therefore, carving out seams through the car and the leaf (the dominant objects), does not generate new salient features and hence these seams are selected. In contrast, when our saliency is used, seams through the car and leaf would introduce salient features. In this case, these seams are avoided, leaving the objects untouched and resulting in less distortions.

## 5.2. Summarization through collage creation

Collages have been a common form of artistic expression since their first appearance in China around 200 BC. Manually creating a collage is a difficult and time consuming task, since the pieces should be nicely cut and matched. Therefore, automation is a welcomed tool. Today, with the advent of large image collections, collages have value also as a summarization tool. In summaries, the salient objects as well as informative pieces of the background should be maintained, whereas the non-meaningful background should be excluded.

Earlier work on automating collage creation extracts rectangular salient regions and assembles them in various fashions [17, 25, 2, 15]. This produces beautiful collages, however, since the extracted regions are rectangular, the variety of possible compositions is limited and uninteresting regions are included. In [17] the rectangular shapes are modified by graph cuts and alpha blending. This creates nicer transitions between images, however, the uninteresting regions

(typically from the background) cannot be eliminated. This approach to assemblage, while informative, is not compact and does not match in spirit the way in which many artists construct collages.

In [4] a method was proposed for automating collage creation, which is inspired by artistic collage work and glues together the salient cutouts. Given a collection of images, this technique consists of three stages. First, the saliency maps of the images are computed. Then, regions-of-interest (ROI) are extracted by considering both saliency and image-edge information. Finally, these non-rectangular ROIs are assembled, allowing slight overlaps.

Our context-aware saliency is very beneficial for Stage 1 of this algorithm. If the background is not interesting, it is not included in the collage. But, if the background is essential for understanding the context, it is included.

Figure 11 illustrates an example of an automatic summarization result, in which our saliency was used [4]. Note that in some images (e.g., the castle), the objects are accurately extracted from the background, whereas in other images (e.g., the boys running towards the ocean), the informative views are included jointly with the people. The resulting collage is compact, pleasing and informative. Making a good use of the space was made possible by using our saliency maps to extract the ROIs.

## 6. Conclusion

This paper proposes a new type of saliency – context-aware saliency – which detects the important parts of the scene. This saliency is based on four principles observed in the psychological literature: local low-level considerations, global considerations, visual organizational rules, and high-level factors. The paper further presents an algorithm for computing this saliency.

There exists a variety of applications where the context of the dominant objects is just as essential as the objects themselves. This paper evaluated the contribution of context-



