MULTI-SCALE REGION-BASED SALIENCY DETECTION USING W2 DISTANCE ON N-DIMENSIONAL NORMAL DISTRIBUTIONS

Lei Zhu^{†‡}, Dominik A. Klein[†], Simone Frintrop[†], Zhiguo Cao[‡], and Armin B. Cremers[†],

[†]Inst. of Comput. Sci. III, Rheinische Friedrich-Wilhelms Universität Bonn, Germany [‡]Inst. for Pattern Recognit. and AI, Huazhong University of Sci. and Tech., Wuhan, China

ABSTRACT

We present a new segment-based method for saliency detection based on multi-size superpixels that combines local and global saliency cues. We extract superpixels at several scales and represent each superpixel with a normal distribution in CIE-Lab space estimated from its associated pixels. Global saliency is computed by grouping similar superpixels to estimate the spatial distribution of colors, while local saliency detection is achieved by determining the center-surround contrast of neighboring superpixels. Both methods rely on the Wasserstein distance on L_2 norm (W_2) to measure perceptual (dis-)similarity between superpixels. Additionally, we propose a *Saliency Flow* technique to refine the local saliency map. Our approach uses very few empirical parameters and outperforms 6 recent state-of-the-art saliency detection methods in terms of several evaluations on a widely used benchmark.

Index Terms— Saliency detection, Wasserstein distance, Color distribution, Center-surround contrast, Clustering.

1. INTRODUCTION

Saliency detection has become a major research area in recent years and has been used for several high level vision applications [1, 2]. Many computational attention methods have built their saliency models based on local center-surround differences [3, 4, 5]. On the other hand, the global compactness of color is another efficient evidence in saliency detection [6, 7, 8]. From the perspective of computation, the classical approaches often calculate a *pixel-based* saliency by use of sliding integration windows to estimate center and surrounding appearances [3, 4, 5]. These windows are rigid and ignore boundaries present in the image content. Another type we call *segment-based* approaches, which apply content-aware segmentation as a preliminary step before calculating saliency for each extracted segment instead of each pixel [6, 8, 9, 10].

In this work, we propose a new segment-based saliency detection approach based on superpixels that are computed on multiple scales. The color occurrence of pixels in each superpixel is approximated by a multivariate normal distribution. This assumption is appropriate for superpixels since the clustered pixels have similar properties in the selected feature space. The difference between superpixels is measured with the Wasserstein distance on L_2 norm (W_2 -distance) [11]. Based on this metric, the global and the local saliencies are computed in a coherent manner. As shown in Fig. 1b, global saliency is determined with a parameter-free approach. This is in contrast to others which employ the global compactness evidence that usually needs several parameters to control the color sensitivity of element distributions [6], or determine the dominant colors [8]. In local saliency estimation, a locally constrained and finite random walk based procedure named *Saliency Flow* is proposed to refine



Fig. 1. From left to right: (a) the original image. (b) global saliency map at single scale ((c) and (d) correspond to the same scale). (c) local saliency map. (d) local saliency map after the *Saliency Flow* method. (e) final saliency map of combining all scales.

our local saliency map via balancing the inside of probable objects as well as smoothing their outside borders with similar surrounding background regions. Fig. 1d shows two examples of using *Saliency Flow* to refine the local saliency maps in Fig. 1c. Furthermore, we introduce a method for measuring the effectiveness and determining the convergence state of this iterating procedure.

We show that our approach outperforms 6 recent methods for saliency detection on a widely used benchmark.

2. THE COMPUTATIONAL SALIENCY FRAMEWORK

This section describes the architecture of our approach in detail. As shown in Fig. 2, firstly, the multi-size superpixels are extracted from images according to predefined scale levels. Secondly, the visual appearance of each superpixel in every scale is approximated by a three-dimensional normal distribution in CIE-Lab space. Differences between superpixel-appearances are measured by the W_2 -distance. Both global and local saliency maps are obtained based on this metric and further fused into a single saliency map per scale. Thirdly, we combine these single scale saliency maps over all scales to form a final saliency map.

2.1. Multi-size Superpixel Extraction and Measurement

We use the SLIC algorithm introduced in [12] to extract superpixels of a given number. Since too small regions make the appearance distribution estimation less meaningful, and are not able to reliably capture texture information, at least 64 pixels are ensured in each superpixel at the finest scale. Then, we repeatedly conduct the SLIC algorithm for several scales starting from the finest scale by decreasing the number of demanded superpixels by a factor of two. Here, three scales are analyzed in our experiments. At scale *t*, the *i*th superpixel is represented as a tuple $S_i^t = (N_S(\mu, \Sigma), \vec{c}_S = (x, y))_i^t$, where μ_i^t and Σ_i^t are the mean value and the covariance matrix of



Fig. 2. Diagram of our proposed algorithm

the normal distribution \mathcal{N}_{Si}^{t} in CIE-Lab space estimated by a MLfitting of the feature values associated with the pixels of the *i*th superpixel, respectively. \vec{c}_{Si} is the spatial center of those pixels. Recently, Klein and Frintrop [11] introduced the \mathcal{W}_2 -distance in their pixel-based saliency detection system. We follow this idea and apply it for our segment-based measure to compute the distance between superpixels. According to [13], the \mathcal{W}_2 -distance between two normal distributions as used here for appearance models of superpixels is computed as

$$\mathcal{W}_2(\mathcal{N}_{\mathcal{S}_i}^t, \mathcal{N}_{\mathcal{S}_j}^t) = \sqrt{\|\mu_i - \mu_j\|^2 + \operatorname{tr}(\Sigma_i + \Sigma_j) - 2\operatorname{tr}\left(\sqrt{\Sigma_i \Sigma_j}\right)}, \quad (1)$$

where $tr(\cdot)$ refers to the trace of a matrix.

2.2. Global Saliency: The Spatial Distribution of Colors

For all color components in an image, the rare, spatially centralized distributed ones always appear more salient. The spatial distribution of colors could be directly assessed based on the dissimilarities between superpixels in the whole image [6]. Our method is more sophisticated in the way that we add a higher level clustering step based on superpixels and rate the spatial intra-cluster distances. For this purpose, we employ the APC (Affinity Propagation Clustering [14]) to identify clusters and estimate the intra-cluster probabilities by the *responsibilities* between each superpixel and all cluster centers [15]. Hereby, the responsibility $r(s_i^t, c_k^t)$ denotes how well-suited superpixel c_k^t is to serve as *exemplar* for superpixel s_i^t .

With the *affinity matrix* composed of negative squared W_2 distances based on Eq. (1), APC selects representative *exemplar* superpixels as the cluster centers. Similar to the definition of superpixels, the k^{th} cluster at scale t forms a tuple $C_k^t = (\mathcal{N}_c(\mu, \Sigma), \vec{c}_c)_k^t$, where \vec{c}_{ck}^t is the center of gravity with respect to the membership probabilities of all superpixels in image coordinates, which is computed as $\vec{c}_{ck}^t = \sum_{i=1}^{M(t)} P_g(\mathcal{C}_k^t | \mathcal{S}_i^i) \cdot \vec{c}_{\mathcal{S}_i}^t$, where M(t) is the number of superpixels and $P_{g}(\mathcal{C}_{k}^{t}|\mathcal{S}_{i}^{t})$ refers to the membership probability of the i^{th} superpixel to the k^{th} cluster at scale t.

Let $C^t = \{c_1^t, ...\}$ denotes the set of superpixels that has been chosen as exemplars at scale t, while $S^t = \{s_1^t, ...\}$ is the disjoint set of non-exemplars, respectively. The membership probability of the j^{th} non-exemplar superpixel to the k^{th} cluster is computed according to how well the exemplar c_k^t matches with s_j^t .

We first linearly normalize the range of responsibilities of all superpixels in S^t to [-1,0]

$$\hat{S}(s_j^t, c_k^t) = \frac{r(s_j^t, c_k^t) - \max_{s_i^t \in S^t} \left\{ r(s_i^t, c_k^t) \right\}}{\max_{s_i^t \in S^t} \left\{ r(s_i^t, c_k^t) \right\} - \min_{s_i^t \in S^t} \left\{ r(s_i^t, c_k^t) \right\}} \quad (2)$$

í

and, since r corresponds to log-probabilities, further exponentially rescale them as

$$\hat{r}_e(s_j^t, c_k^t) = \exp\left(\hat{r}(s_j^t, c_k^t) \middle/ \operatorname{Var}_{s_j^t \in S^t}\left(\hat{r}(s_j^t, c_k^t)\right)\right), \quad (3)$$

where $\operatorname{Var}_{s_j^t \in S^t} \left(\hat{r}(s_j^t, c_k^t) \right)$ is the variance of the normalized responsibilities send to c_k^t . The rescaled responsibilities between two exemplars are defined as

$$\hat{r}_e(c_j^t, c_k^t) = \begin{cases} 1, & \text{if } j = k\\ 0, & \text{otherwise} \end{cases}.$$
(4)

Eq. (3) and Eq. (4) define responsibilities from all superpixels to each exemplar. With these, the intra-cluster probabilities of each superpixel can be obtained as

$$P_{g}(\mathcal{C}_{k}^{t}|\mathcal{S}_{i}^{t}) = \hat{r}_{e}(\mathcal{S}_{i}^{t},\mathcal{C}_{k}^{t}) / \sum_{k=1}^{\mathrm{K}(t)} \hat{r}_{e}(\mathcal{S}_{i}^{t},\mathcal{C}_{k}^{t}),$$
(5)

where K(t) is the number of clusters at scale t. By weighting all superpixels with their intra-cluster probabilities, the probability of being salient for cluster C_k^t is obtained by scoring the relative spatial spreading between all superpixels and the k^{th} cluster:

$$P_{g}(\mathcal{C}_{k}^{t}) = 1 / \sum_{j=1}^{K(t)} \sum_{i=1}^{M(t)} P_{g}(\mathcal{C}_{k}^{t} | \mathcal{S}_{i}^{t}) \cdot ||\vec{c}_{\mathcal{S}_{i}}^{t} - \vec{c}_{\mathcal{C}_{j}}^{t}||^{2}.$$
(6)

Finally, the global, superpixel level saliency can be represented as the joint probability of Eq. (5) and Eq. (6) as

$$P_{g}(\mathcal{S}_{i}^{t}) = \sum_{k=1}^{K(t)} P_{g}(\mathcal{C}_{k}^{t}) \cdot P_{g}(\mathcal{C}_{k}^{t}|\mathcal{S}_{i}^{t}).$$
(7)

2.3. Local Saliency: The Local Contrast of Superpixels

Sometimes, the global method cannot distinguish the salient object from the background at every scale (cf. Fig. 3b). An alternative way to solve this problem is introducing local information. We follow the classical center-surround theory to compute the local saliency by regarding each superpixel as *center* and other superpixels as *surround* which are weighted by their Euclidean distances to the *center*.

At scale t, the local contrast of superpixel S_i^t is obtained by accumulating all the appearance distances from other superpixels

$$P_{\mathrm{l}}(\mathcal{S}_{i}^{t}) = \frac{1}{g} \sum_{j=1}^{\mathrm{M}(t)} \mathcal{W}_{2}(\mathcal{N}_{\mathcal{S}i}^{t}, \mathcal{N}_{\mathcal{S}j}^{t}) \cdot \exp\left(-\frac{||\vec{c}_{\mathcal{S}i}^{t} - \vec{c}_{\mathcal{S}j}^{t}||^{2}}{\sigma_{c}^{2}(t)}\right), \quad (8)$$



(a) image/result (b) global saliency (c) local saliency (d) saliency flow

Fig. 3. From top to bottom: (a) the original image and the final saliency map. (b) the global saliency maps in lowest and finest scales ((c) and (d) correspond to the same scale). (c) the local saliency maps. (d) the local saliency maps after the *Saliency Flow*.

where $g = \sum_{j=1}^{M(t)} \exp\left(-||\vec{c}_{s\,i}^{t} - \vec{c}_{s\,j}^{t}||^{2} / \sigma_{c}^{2}(t)\right)$. To ensure the spatial radius of neighbor range to be scaled proportionally, we associate $\sigma_{c}(t)$ to the average distance between adjacent superpixels.

Inspired by the Random Walk related theories, we propose a procedure called *Saliency Flow* to refine the local saliency map by balancing the saliency values inside probable proto-objects as well as smoothing the locally standing out background regions.

At each scale, we construct an undirected graph model $G_t = \{V_t, E_t\}$ comprising a set of vertices with all superpixels: $V_t = \{S_1^t, \ldots, S_{M(t)}^t\}$ and edges $E_t = \{e_{i,j;t}\}$ which represent the pairwise relations between superpixels in both color and spatial space:

$$e_{i,j;t} = \begin{cases} \exp\left(-\mathcal{W}_2(\mathcal{N}_{\mathcal{S}_i^t}, \mathcal{N}_{\mathcal{S}_j^t})^2 / \sigma_d^2(t)\right), \text{ if } \mathcal{S}_i^t \text{ connects to } \mathcal{S}_j^t \\ 0, & \text{otherwise.} \end{cases}$$
(9)

Here, we employ the *Three-Sigma Rule* and adaptively choose σ_d following $3 \cdot \sigma_d(t) = \max \{ W_2(S_i^t, S_j^t) \}$. Each column of E_t is normalized by its L_1 norm for a probability description of saliency transferring from one superpixel to another. The balanced local saliency after *n* steps of *Saliency Flow* can be expressed as:

$$f(0)^{t} = \vec{p}^{t} = \left(P_{l}(V_{t;1}), \dots, P_{l}(V_{t;M(t)})\right)^{\mathsf{T}}$$

$$f(n)^{t} = f(n-1)^{t} \cdot E_{t} = \vec{p}^{t} \cdot (E_{t})^{n}.$$
(10)

In each step, the saliency of superpixels gravitates towards their similar looking neighbors. Taking superpixels located at both sides of the object's boundaries for example, the superpixels at the background side have a high probability to flow to the adjacent regions of background while the superpixels at the object side will more likely flow towards the inside of object.

We observed that the saliency differences between similar neighbors are quickly smoothed after some iterations. However, the flow continues to equalize the differences also between non-similar neighbors which finally results in a uniform saliency distribution of all superpixels when $n \rightarrow +\infty$. Therefore, an appropriate *n* should not only ensure enough iterations to smooth the inner regions of an object but also terminate the iteration in time to prevent much "saliency mass" flowing between object and background. We introduce a measure of locally absolute exchanged saliency to estimate the effect of saliency flow and determine *n* as follows:

$$\tau_f(n) = \operatorname{mean}\left\{\sum_{i,j} \left| f(n)_i^t - f(n)_j^t \right| e_{i,j;t} \right\},$$

where $\mathcal{S}_i^t, \mathcal{S}_j^t \in \mathcal{C}_k^t$, and (11)
 $\operatorname{argmin}_n \left\{ \tau_f(n-1) - \tau_f(n) \right\} < c_{\min},$



Fig. 4. Precision-Recall curves to measure the effectiveness of saliency segmentation based on fixed thresholds. Please see Fig. 5a-5i for visual comparisons.

where $c_{\min} = 10^{-5}$ was chosen in our experiments. τ_f counts the absolute exchanged "saliency mass" between each pair of superpixels which not only spatially connect to each other but also are assigned to the same cluster. That is, it considers the locality both in spatial and color space.

2.4. Feature Fusion and Multi-size Combination

At each scale t, both global map $P_g(S_i^t)$ and local saliency map $f(n)_i^t$ are normalized to [0, 1] and combined into the single scale saliency map by superpixel-wise multiplication. Furthermore, the arithmetic mean of the normalized pixel level saliency maps over all scales is used to generate the final saliency map.

3. EXPERIMENTS

In this section, we evaluate our salient region detection method with 6 recently proposed saliency detection approaches (the numbers following the names of these methods indicate their publishing years): CoDi12 [11], SF12 [6], XL12 [10], Rare12 [16], RC11 [9], IG09 [17]. All methods are evaluated on 1000 images from the subset of the MSRA Salient Object Database [17, 18]. All result images of other methods except Rare12 [16] directly came from their project home pages, and we produced the result of Rare12 [16] using their published code.

3.1. Saliency Region Segmentation by Fixed Thresholds

We normalize all the saliency maps to the range of [0, 255] and obtain 256 binary results by varying the segmentation threshold. The precision and recall rates are firstly computed for each image at every threshold then averaged over the overall benchmarks. Fig. 4 compares the precision versus recall curves of our method with other algorithms. It is clear that our method achieves a higher precision value at every given recall than other methods. Notice that, compared to others, XL12 [10] has a much larger minimum recall value, which means that it produces a quite uniform saliency distribution inside its detected foreground and Fig. 5e shows an example with this situation. However, the higher false positives make it still achieving a much lower precision than our method achieved at all recalls.



Fig. 5. Visual comparison of all evaluated methods ("GT" in column *i* refers to the ground truth).



Fig. 6. Precision, Recall, *F*-measure and MCC to measure the effectiveness of saliency segmentation based on adaptive threshold. From left to right, all methods are sorted by descending MCC measures.

3.2. Saliency Region Segmentation by Adaptive Thresholds

Instead of using all possible thresholds for testing every image, Achanta et al. [17] proposed a simple and adaptive way of selecting a threshold only depending on the saliency values, which defines the threshold as 2 times of average gray value in the saliency map. Furthermore, based on the obtained precision and recall values, the *F*-measure is computed for each image and averaged over the whole database, which is defined as $(1 + \beta^2) \times P \times R/(\beta^2 \times P + R)$, where *P* and *R* refer to the precision and recall, respectively. Similar to [17] and [6], β^2 is fixed at 0.3 in our evaluation. Fig. 6 compares the precision, recall and *F*-measure (the green bars) values of all evaluated approaches. As one can see, our method is also ahead of the others in this evaluation.

From the perspective of precision, recall and *F*-measure computations, neither of them considers the impact of true negative saliency assignments. These evaluations bias to benefit the approaches which can correctly find the true foreground rather than others which fail to do it but can avoid pushing the not salient regions forward. A fair comparison may be achieved by computing the *Matthews correla*- tion coefficient (MCC) [19], which is defined as:

$$MCC = \frac{tp \times tn - fp \times fn}{\sqrt{(tp + fp)(tp + fn)(tn + fp)(tn + fn)}},$$
 (12)

where tp, tn, fp, and fn refer to true positives, true negatives, false positives and false negatives, respectively. MCC is regarded as a fair measurement in binary classifications especially for two unbalanced categories. Eq. (12) returns a value ranging in [-1, 1] which stands for a continuous measure from totally irrelevant (-1) to precisely overlapping (+1) in our evaluation. The violet bars in Fig. 6 show the MCC values of each evaluated method. Again, our method outperforms others in this evaluation. It is also interesting that the evaluation result of the MCC measure differs from the one according to the *F*-measure for several approaches.

4. CONCLUSION

We have presented a bottom-up region based saliency detection by evaluating two attributes of pre-segmented regions: global rareness and local prominence with an effective W_2 -distance measure. In the global part, by taking advantage of APC, several practical problems caused by non-Euclidean metric are solved properly. In the local part, we also proposed the "Saliency Flow" procedure to address the problem of discontinuous saliency assignment inside proto-objects. The used two evidences work in different ways and the experiments also show that they usually give prominence to same salient objects while suppress different background regions in an image. That is, the global and the local considerations independently contribute to the final result with equivalent importance.

On average, our algorithm takes 0.7s for processing an image. The whole project is implemented in C++ and timings have been tested on an Intel Core i5-2410M(obile) at 2.3 GHz. Most processing time is spent on the distance computation and the clustering. In the future, we plan to use a sparse affinity matrix based on clipped distances to achieve real-time performance.

5. REFERENCES

 Simone Frintrop, Erich Rome, and Henrik I. Christensen, "Computational visual attention systems and their cognitive foundations, A survey," *ACM Trans. Appl. Percept.*, vol. 7, no. 1, pp. 1–39, Jan. 2010.

- [2] Ali Borji and Laurent Itti, "Salient object detection: A benchmark," in *Proc. European Conf. Comput. Vis.*, 2012, pp. 414– 429.
- [3] L. Itti, C. Koch, and E. Niebur, "A model of saliency-based visual attention for rapid scene analysis," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 20, no. 11, pp. 1254–1259, Nov. 1998.
- [4] Simone Frintrop, VOCUS: A Visual Attention System for Object Detection and Goal-directed Search, vol. 3899 of Lecture Notes in Artificial Intelligence (LNAI), Springer, Berlin/Heidelberg, 2006.
- [5] Dominik A. Klein and Simone Frintrop, "Center-surround divergence of feature statistics for salient object detection," in *Proc. IEEE Int. Conf. Comput. Vis.*, Nov. 2011, pp. 2214–2219.
- [6] Federico Perazzi, Philipp Krahenbuhl, Yael Pritch, and Alexander Hornung, "Saliency filters: Contrast based filtering for salient region detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, Jun. 2012, pp. 733–740.
- [7] Hsin-Ho Yeh and Chu-Song Chen, "From rareness to compactness: Contrast-aware image saliency detection," in *Proc. IEEE Int. Conf. Image Process.*, Sep. 2012, pp. 1077–1080.
- [8] Zhixiang Ren, Yiqun Hu, Liang-Tien Chia, and Deepu Rajan, "Improved saliency detection based on superpixel clustering and saliency propagation," in *Proc. ACM Int. Conf. Multimedia*, New York, USA, 2010, number 2, pp. 1099–1102.
- [9] Ming-Ming Cheng, Guo-Xin Zhang, Niloy J. Mitra, Xiaolei Huang, and Shi-Min Hu, "Global contrast based salient region detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, Jun. 2011, pp. 409–416.
- [10] Yu-lin Xie, Hu-chuan Lu, and Ming-Hsuan Yang, "Bayesian Saliency via Low and Mid Level Cues," *IEEE Trans. Image Process.*, pp. 1–10, Sep. 2012.
- [11] Dominik A. Klein and Simone Frintrop, "Salient Pattern Detection using W2 on Multivariate Normal Distributions," in *Proc. DAGM-OAGM Conf.*, Aug. 2012, vol. 7476, pp. 246– 255.
- [12] Radhakrishna Achanta, Appu Shaji, Kevin Smith, Aurelien Lucchi, Pascal Fua, and Sabine Süsstrunk, "SLIC superpixels compared to state-of-the-art superpixel methods.," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 34, no. 11, pp. 2274– 82, Nov. 2012.
- [13] D. C. Dowson and B. V. Landau, "The Fréchet distance between multivariate normal distributions," *J. Multivar. Anal.*, vol. 12, no. 3, pp. 450–455, Sep. 1982.
- [14] Brendan J. Frey and Delbert Dueck, "Clustering by passing messages between data points," *Science*, vol. 315, pp. 972– 976, Feb. 2007.
- [15] T. Geweniger, D. Zühlke, B. Hammer, and Thomas Villmann, "Fuzzy variant of affinity propagation in comparison to median fuzzy c-means," in *Proc. Int. Workshop on Adv. Self-Organizing Maps.* 2009, pp. 72–79, Springer-Verlag.
- [16] Nicolas Riche, Matei Mancas, Bernard Gosselin, and Thierry Dutoit, "Rare: A new bottom-up saliency model," in *Proc. IEEE Int. Conf. Image Process.*, Sep. 2012, pp. 641–644.
- [17] Radhakrishna Achanta, Sheila Hemami, F. Estrada, and S. Süsstrunk, "Frequency-tuned salient region detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, Jun. 2009, pp. 1597–1604.

- [18] Tie Liu, Zejian Yuan, Jian Sun, Jingdong Wang, Nanning Zheng, Xiaoou Tang, and Heung-Yeung Shum, "Learning to detect a salient object.," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 33, no. 2, pp. 353–367, Feb. 2011.
- [19] BW Matthews, "Comparison of the predicted and observed secondary structure of T4 phage lysozyme," *Biochim. Biophys. Acta*, vol. 405, no. 2, pp. 442–451, 1975.