

# Learning the Relation between Interested Objects and Aesthetic Region for Image Cropping

Peng Lu, Hao Zhang, Xujun Peng, and Xiaofu Jin

**Abstract**—As one of the fundamental techniques for image editing, image cropping discards irrelevant contents and remains the pleasing portions of the image to enhance the overall composition and achieve better visual/aesthetic perception. In this paper, we primarily focus on improving the efficiency of automatic image cropping, and on further exploring its potential in public datasets with high accuracy. From this perspective, we propose a deep learning based framework to learn the objects composition from photos with high aesthetic qualities, where an interested object region is detected through a convolutional neural network (CNN) based on the saliency map. The features of the detected interested objects are then fed into a regression network to obtain the final cropping result. Unlike the conventional methods that multiple candidates are proposed and evaluated iteratively, only a single interested object region is produced in our model, which is mapped to the final output directly. Thus, low computational resources are required for the proposed approach. The experimental results on the public datasets show that as a weakly supervised method, the proposed network outperforms the other weakly supervised methods on FLMS and FCD datasets and achieves comparable results to the existing methods on CUHK dataset. Furthermore, the proposed method is more efficient than these methods, where the processing speed is as fast as 20ms per image.

**Index Terms**—Deep learning, aesthetics, image composition, convolutional network.

## I. INTRODUCTION

Image cropping, which aims at removing unexpected regions and non-informative noises from a photo/image, by modifying its aspect ratio or through improving the composition, is one of the basic image manipulation processes for graphic design, photography and image editing. Nowadays, with the proliferation of hand-held smart devices, users are more eager to capture photos obtaining not only the theme that the image needs to express but also the good composition and appealing colors, to facilitate semantic searching and to make audiences enjoyable. This trend attracts increasing interests of image cropping from both research community and industries.

However, cropping an image to obtain appropriate composition for achieving better visual quality is notoriously difficult, primarily driven by three facts: (1) to determine the main object/theme of a given image is a nontrivial task, which needs deep domain knowledge and sophisticate skills; (2) assessment of aesthetic of the cropped image is highly subjective such that different viewers might have various opinions for the same

cropped photo, or even the same viewer might have opposite feelings for the same image at different time; (3) vast amount of cropping candidate areas can be extracted from the image which causes the solution space is exponentially increased.

To tackle these problems, many researchers seek to propose novel approaches to automatically crop images with high aesthetic score. These existing researches can be roughly grouped into four main categories: *sliding-judging* based, *determining-adjusting* based [1], *sequential decision-making* based [2] and *detecting-determining* based methods.

The sliding-judging based approaches normally exhaustively scan the entire image using windows with different size and aspect ratio to produce abundant candidate regions [3], [4]. For each candidate, a classifier or ranker is applied to evaluate its visual/aesthetic quality and the one with the highest score is considered as the optimal cropping result. However, the low computational efficiency of these approaches limits their success. In order to avoid greedily searching against all possible sub-windows, determining-adjusting based approaches attempt to propose a small set of candidate windows with high probabilities to narrow down the searching space for the optimal cropping rectangle. Normally, a seed candidate is initially determined by face, salient object, or attention detection algorithms [5], [6], [7], from which the surrounding areas are scanned and evaluated by classifiers or rankers to select the region with the highest visual/aesthetic quality. Although determining-adjusting based approaches have higher efficiency than sliding-judging based methods, they still encounter the same problems of multiple candidates generation and selection. To avoid evaluating a large amount of proposals, sequential decision-making based approaches use aesthetics aware reward function to guide the searching for cropping windows and decision-making iterations are reduced to as low as dozens for crops prediction [2]. Unlike all existing cropping approaches, by discovering the relation between interested objects and the aesthetic quality of cropped image, detecting-determining based approaches find the optimal cropping rectangle based on the detected interested object region directly without any multiple proposals and evaluations [8].

As can be seen from [8] that by employing interested objects which represent those areas attracting most attentions from the viewers, and its relation to aesthetic areas, the detecting-determining based image cropping approaches demonstrate the promising results for both accuracy and efficiency. However, the brute force searching technique [6] used for the interested object localization (IOL) is still a bottle-neck for the efficiency of this type of methods. And the multiple stages training and inference scheme also limits its applicability. Notably, unlike

P. Lu, H. Zhang and X. Jin are with the School of Computer Science, Beijing University of Posts and Telecommunications, Beijing, China. E-mail: plu@bupt.edu.cn

X. Peng is with the Information Sciences Institute, University of Southern California, Marina del Rey, CA, USA. E-mail: xpeng@isi.edu

94 the general objects detection approaches, which specially focus  
95 on one or multiple objects in its entirety, IOL focuses more on  
96 the psychological feelings from the perspective of viewers and  
97 is more suitable for the image cropping task. Normally, IOL  
98 is computed based on the saliency map detection. Similar to  
99 but differ from the salient objects detection [9], [10], [11],  
100 saliency map detection attempts to calculate the “saliency  
101 map” that simulates the eye movement behaviors of human,  
102 but salient objects detection tends to bias to particular objects  
103 in the image, which results in different assessment criteria  
104 and ground-truth used for these two different tasks. Thus,  
105 the saliency map detection generally has higher generalization  
106 capability because it is not relied on particular objects. The  
107 difference between them is detailed addressed in [1] and [12].

108 Thus, in this paper we propose a weakly supervised end-  
109 to-end image cropping framework to address the problems  
110 of detecting-determining based approaches, where the ground  
111 truths of cropping bounding boxes are not required in our  
112 system. Particularly, the proposed image cropping system uses  
113 a deep neural network to extract the saliency map of the image,  
114 which is followed by the proposed IOL layer to determine the  
115 region containing the interested objects in the image. Then a  
116 regression network is employed to map the interested object  
117 region to the final cropping rectangle based on its feature.  
118 The proposed cropping method only has one pass to achieve  
119 the optimal cropping result, without iterative searching or  
120 scanning on multiple proposals of different windows, which  
121 highly improves the computational efficiency and obtains  
122 comparative accuracy performance.

123 In summary, we make three contributions to the literature:

- 124 • We propose a probabilistic framework to model the  
125 relationship between the interested objects in the image  
126 and the area with high aesthetic quality. Based on this  
127 relation, the task of image cropping can be considered as  
128 an optimization problem to maximize the joint probability  
129 of interested objects and the cropped area with high  
130 aesthetic quality for the given image;
- 131 • The searching technique which is normally used for the  
132 IOL is the main obstacle for the end-to-end training and  
133 inference because it is not differentiable. An IOL layer  
134 is proposed based on the saliency map to avoid this type  
135 of heuristic searching scheme. The proposed layer can  
136 effectively find the location of interested objects and is  
137 differentiable;
- 138 • Based on the proposed probabilistic framework and IOL  
139 layer, an end-to-end cropping system is designed, which  
140 is not relied on the cropping annotation datasets for  
141 training but achieves the state-of-the-art accuracy and  
142 high efficiency with 50 frame per second (FPS) on public  
143 datasets.

144 The remainder of this paper is organized as follows. Section  
145 II briefly covers the related works to image cropping. Section  
146 III is an in-depth introduction of the proposed methodology.  
147 The experimental setup, results analysis and discussion are  
148 presented in section IV. Finally, we conclude our work in  
149 section V.

## II. PREVIOUS WORK

### A. Saliency Map Detection

150 Most vertebrates have the ability to move their eyes and  
151 predict fixation with limited time and resources, which enables  
152 them to focus on the most informative region and extract the  
153 most relevant features for the particular scene [13], [14]. This  
154 phenomenon inspired researchers to obtain the cropping areas  
155 of image relied on the saliency map prediction.  
156

157 Generally, saliency map is produced prior to the salient  
158 object detection, as demonstrated in [12], where each pixel in  
159 the map indicates the confidence of the fixation. In [15], Harel  
160 *et al.* proposed a graph-based visual saliency model depended  
161 on Markovian chain assumption, which calculated and normal-  
162 ized the activation map by measuring the dissimilarity between  
163 neighboring pixels. The reported ROC curve for this method  
164 beat the classical attention map detection approach proposed  
165 by Itti *et al.*, where multiple empirical features were fed into a  
166 neural network to select the proper attended locations [16]. In  
167 the same manner, Judd *et al.* defined a set of hand-crafted  
168 features to represent low-, mid- and high-level perception  
169 of human visual system and fed them into a support vector  
170 machine (SVM) to distinguish positive and negative salient  
171 pixels [17].  
172

173 However, the drawback of these mentioned approaches is  
174 that strong domain knowledge and experiences are required  
175 for design of those hand-tuned features, which is a obstacle  
176 to extend their applicability. Therefore, in [18], Vig *et al.*  
177 proposed to utilize CNN to learn the representations for salient  
178 and non-salient regions. With labeled feature vectors, an L2-  
179 regularized, linear, L2-loss SVM was trained to predict the  
180 probability of fixation of images in their work.

### B. Aesthetic Assessment

181 Besides salient objects that affect the performance of image  
182 cropping, aesthetic, which represents the degree of beauty,  
183 is another key factor to determine the quality of cropped  
184 images. Early work for aesthetic assessment can be dated  
185 back to the researches of color harmony theories [19] and  
186 photographic composition [20]. In recent years, many auto-  
187 matic image aesthetic assessment algorithms were proposed,  
188 where hand-crafted global features, such as spatial distribution  
189 of edges, color distributions, hue count etc. [21] and local  
190 features, e.g. wavelet-based texture and shape convexity [22]  
191 were employed. To take the advantage of both global and  
192 local features, Zhang *et al.* combined structural cues of these  
193 two levels for photo aesthetic evaluation [23]. Particularly,  
194 graphlet-based local structure descriptors were constructed and  
195 projected onto a manifold to preserve the global layout of the  
196 image, which was embedded into a probabilistic framework  
197 to assess image aesthetic. However, these representations con-  
198 sider whole image indiscriminately ignoring the importance  
199 of main subjects in the image. To remedy this problem, Luo  
200 *et al.* suggested extracting different subject areas prior to the  
201 aesthetic evaluation and treating them using different aesthetic  
202 features [24]. Furthermore, genetic image descriptors were  
203 also applied to facilitate the aesthetic assessment task. For  
204 instance, Marchesotti *et al.* developed two types of local image  
205

signatures originated from Bag-of-visual-words and Fisher vectors by incorporating SIFT and color information into them [25].

Under the assumption that semantic recognition task can help the aesthetic assessment, Kao *et al.* proposed a multi-task framework where two tasks were trained simultaneously while the representations were shared by two networks [26]. This idea was also applied by Lu *et al.* in their work of color harmony modeling, which used both bag-of-visual-words features and semantic tag information to boost the aesthetic assessment performance through colors [27]. Moreover, in order to overcome the problem of contaminated tags, a semi-supervised deep active learning algorithm was proposed in [28], where a large set of object patches were extracted and linked to the semantical tags to benefit image aesthetic assessment.

### C. Regression Networks

Although photos can be cropped depended on the obtained salient objects only, they are not necessary to be with high aesthetic quality because the aspect of aesthetic is normally ignored for saliency detection. To tackle this problem, one feasible solution is to determine a seed cropping window according to the saliency map and propose a set of candidates around this seed window. Thereafter, every candidate is evaluated by its aesthetic quality and the one with the highest aesthetic score is considered as the final cropping result, as the methods described in [4], [29]. However, iterative assessment for each candidate window's aesthetic score increases the computational complexity. Thus, a more practical and efficient approach is to map the seed salient region to the final cropping window directly using regression models, where the aesthetic information is integrated into the system.

In the field of computer vision [30], [31], regression method is widely used for object detection. Girshick *et al.* combined regions with CNN features (R-CNN) to find different objects in the image [32], where bounding-box regression technique proposed in [33] was applied for a selective search region proposal to refine the detection window. To speed up R-CNN, Girshick improved their work by using a fully connected network to predict the bounding-box regression offsets and confidence of each proposal [34]. Instead of performing classification for detection problem, Redmon *et al.* framed object detection as a regression problem, where the input image was divided into small patches initially and the bounding boxes offsets as well as their probabilities for each class were predicted directly in one neural network, which was called YOLO [35]. The final detections were obtained by merging bounding boxes for the same class. To make YOLO better and faster, Redmon and Farhadi shrink the CNN and used region proposal network (RPN) to generate more anchor boxes for boosting recall and localization accuracy, where regression network was remained for location/confidence prediction [36]. Unlike YOLO, Liu *et al.* detected different objects by evaluating a small set of boxes that were produced through multi-resolution CNN feature maps. The final bounding boxes of objects were also obtained by regressing to offsets for the centers of the default

boxes [37]. The similar ideas of using regression networks for objects detection can be found in [38].

### D. Image Cropping & Recomposition

As an important procedure to enhance the visual quality of photos, image cropping and recombination benefit from the development of salient object detection, aesthetic assessment and other computer vision techniques. To combine visual composition, boundary simplicity and content preservation into a photo cropping system, saliency map and salient object were used to encode the spatial configuration and content information, and gradient values were applied to measure the simplicity of image by Fang *et al.* [3]. In this method, image was densely cropped, evaluated and merged by the mentioned features to obtain the optimal cropping results. By segmenting the entire image into small regions, a region adjacency graph (graphlets) was constructed by Zhang *et al.* to represent the aesthetic features of the image, from which the image was cropped through a probabilistic model [39], [40]. Zhang *et al.* also extended the idea of graphlets in the semantic space for image cropping, which was created based on the category information of the images [41]. In the semantic space, semantically representative graphlets were selected sequentially and evaluated by a pre-trained aesthetic prior model to guide the cropping process. Unlike the other algorithms that evaluated multiple candidate cropping areas, Samii *et al.* searched against a high quality image database to find exemplar photos with similar spatial layouts as the query image, and matched the composition of the query image to each of exemplars by minimizing composition distance in a high-level context feature space to calculate the optimal crop areas [42]. In [43], Wang *et al.* applied similar concept for photo cropping that exploited sparse auto-encoder to discover the composition basis from a database containing well-composed images. Differ from [42], Wang's method organized the learning and inference in a cascade manner to achieve higher efficiency. By considering that perspective effect is one of the most commonly used techniques for photography, Zhou *et al.* developed a hierarchical segmentation method integrating photometric cues with perspective geometric cue to detect the dominant vanishing point in the image, which was employed for image re-composition or cropping [44].

Recently, thanks to the development of DNN, the research of image cropping tends to utilize deep learning approaches. By imitating the process of professional photographic, Chen *et al.* proposed a ranking CNN to harvest unambiguous pairwise aesthetic ranking examples on the web and applied this network to find the optimal cropping result from many candidate regions [45]. Instead of generating the attention map for cropping, Kao *et al.* proposed to use aesthetic map, which was extracted via a CNN, and gradient energy map to accomplish the image cropping task, by learning the composition rules through a SVM classifier [46]. In [47], Guo *et al.* designed a cascaded cropping regression (CCR) approach to crop the image, where a deep CNN was applied to extract features from images and the cropping areas were predicted by the proposed CCR algorithm. Inspired by human's decision making, Li *et al.*

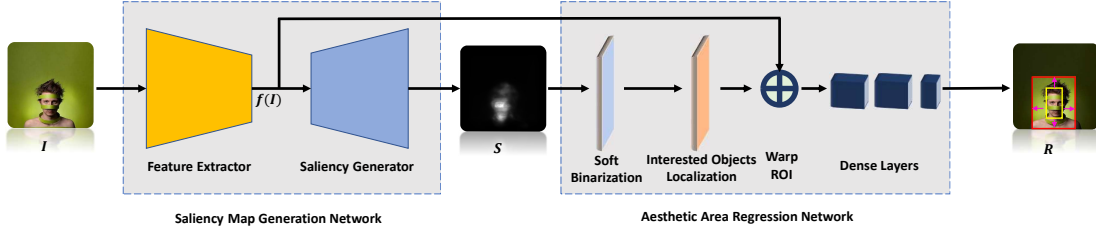


Fig. 1. Architecture of the proposed saliency map detection and aesthetic area regression network.

318 designed a weakly supervised aesthetic aware reinforcement  
 319 learning framework to address the problem of image cropping,  
 320 where the photo was initially cropped and repeatedly updated  
 321 based on the current observation and the historical experience  
 322 [2]. In [8], Lu *et al.* proposed a regression network based  
 323 cropping method, which mapped initial detected saliency rect-  
 324 angle to a cropping area with high aesthetics quality. Unlike  
 325 the conventional photo cropping method that only produced  
 326 a single output, in [48], Wei *et al.* proposed a system that  
 327 returned multiple cropping outputs based on a teacher-student  
 328 framework. In this framework, the teacher was trained to  
 329 evaluate candidate anchor boxes, and the scores from the  
 330 teacher were used to supervise the training of student, a  
 331 view proposal net. The combination of these two networks  
 332 effectively improve the cropping performance. The interest  
 333 readers can refer [49] for more comprehensive surveys.

### 334 III. PROPOSED APPROACH

#### 335 A. Motivation & System Overview

336 By studying the procedure of professional photography, we  
 337 can see that the theme is firstly determined by the photographer  
 338 prior to other actions. To express this theme, the objects along  
 339 with the compatible backgrounds are selected subsequently.  
 340 Once the main objects contained in the photo are given, the  
 341 other parameters for the photography, such as exposure time,  
 342 composition, colors, etc., will be set for the final shooting.

343 Based on this observation, the process of image cropping to  
 344 obtain the high aesthetic quality can be decomposed into two  
 345 steps: detection of the interested objects  $\mathcal{S}$  in the image  $\mathcal{I}$  and  
 346 prediction of the aesthetic areas  $\mathcal{R}$  of the image based on the  
 347 objects of interest  $\mathcal{S}$ . This process can be formally expressed  
 348 as:

$$P(\mathcal{R}, \mathcal{S}|\mathcal{I}) = P(\mathcal{S}|\mathcal{I}) \cdot P(\mathcal{R}|\mathcal{S}, \mathcal{I}), \quad (1)$$

349 where  $\mathcal{S} = \{S_{i,j} | i \times j \in |\mathcal{I}|\}$  denotes the interested objects  
 350 of the image,  $|\mathcal{I}|$  means the number of pixels in the image,  
 351 and  $S_{i,j} \in \{0, 1\}$  represents whether a given pixel belongs to  
 352 the objects of interest.  $P(\mathcal{S}|\mathcal{I})$  is the probability of interested  
 353 objects  $\mathcal{S}$  given an image  $\mathcal{I}$ , and  $P(\mathcal{R}|\mathcal{S}, \mathcal{I})$ , which reveals the  
 354 hidden relationship between the interested objects and the final  
 355 cropping region, denotes the probability of  $\mathcal{R}$  with respect to  
 356 the image  $\mathcal{I}$  and the detected objects of interest  $\mathcal{S}$ .

357 Thus, the aesthetic region of an image can be obtained if  
 358  $P(\mathcal{R}, \mathcal{S}|\mathcal{I})$  is calculated. Hence, a probabilistic model based  
 359 cropping system can be designed whose parameters can be  
 360 expressed as  $\Theta$ , and the image cropping task can be considered  
 361 as the maximum likelihood (ML) estimation problem for a

given training image  $\mathcal{I}_i$ , along with its ground truths  $\mathcal{S}_i$  and  
 $\mathcal{R}_i$ :

$$\begin{aligned} \Theta &= \arg \max_{\Theta} \sum_{i=1}^N P(\mathcal{R}^{(i)}, \mathcal{S}^{(i)} | \mathcal{I}^{(i)}; \Theta) \\ &= \arg \max_{\Theta} \sum_{k=1}^N P(\mathcal{S}^{(k)} | \mathcal{I}^{(k)}; \Theta_s) \cdot P(\mathcal{R}^{(k)} | \mathcal{S}^{(k)}, \mathcal{I}^{(k)}; \Theta_r) \\ &= \arg \max_{\Theta} \sum_{k=1}^N (\log P(\mathcal{S}^{(k)} | \mathcal{I}^{(k)}; \Theta_s) \\ &\quad + \log P(\mathcal{R}^{(k)} | \mathcal{S}^{(k)}, \mathcal{I}^{(k)}; \Theta_r)), \end{aligned} \quad (2)$$

364 where superscript  $k$  indicates the index of training sample and  
 365 ground truth,  $N$  is the total number of training samples, and  
 366  $\Theta = [\Theta_s, \Theta_r]^T$  are the parameters of the model.

367 Based on this analysis, we design an end-to-end DNN based  
 368 image cropping system that follows the probability framework  
 369 as described in Eq. 2. In the proposed cropping system,  
 370 two main components are constructed, where the saliency  
 371 map generation network  $H(\mathcal{I}; \Theta_s)$  in the Fig. 1 is served  
 372 to predict  $\mathcal{S}$  given image  $\mathcal{I}$ . And aesthetic area regression  
 373 network  $G(\mathcal{I}, \mathcal{S}; \Theta_r)$  containing the proposed IOL layer, ROI  
 374 warping pooling layer and fully connected layers is used  
 375 as a regressor to produce final cropping outputs  $\mathcal{R}$  based  
 376 on  $\mathcal{I}$  and  $\mathcal{S}$ . These two components are corresponding to  
 377 the photographer's actions of the objects decision and final  
 378 cropping areas selection.

Thus, maximizing the Eq. 2 is equivalent to minimizing the  
 loss  $L_{total}$  of the neural network:

$$\begin{aligned} \Theta &= \arg \min_{\Theta} \mathcal{L}_{total} \\ &= \frac{1}{N} \sum_{k=1}^N \arg \min_{\Theta} (\mathcal{L}_s(\hat{\mathcal{S}}^{(k)}, \mathcal{S}^{(k)}) + \lambda \mathcal{L}_r(\hat{\mathcal{R}}^{(k)}, \mathcal{R}^{(k)})), \end{aligned} \quad (3)$$

379 where  $\mathcal{L}_s(\cdot)$  represents the loss from the inconsistency between  
 380 predicted  $\hat{\mathcal{S}}^{(k)}$  by the saliency map detection network and  
 381 ground truth  $\mathcal{S}^{(k)}$  of the images  $\mathcal{I}^{(k)}$ ,  $\mathcal{L}_r(\cdot)$  is the loss  
 382 caused by the difference between predicted aesthetic region  
 383  $\hat{\mathcal{R}}^{(k)}$ , and the ground truth region  $\mathcal{R}^{(k)}$ , and  $\lambda$  is the weight  
 384 controlling the influence from these two networks and we use  
 385  $\lambda = 1$  in this work.

386 As can be seen from the Fig. 1, unlike the conventional  
 387 image cropping methods that explicitly or implicitly generate  
 388 and evaluate multiple candidate cropping regions, the proposed  
 389 system takes the input image to extract the interested object

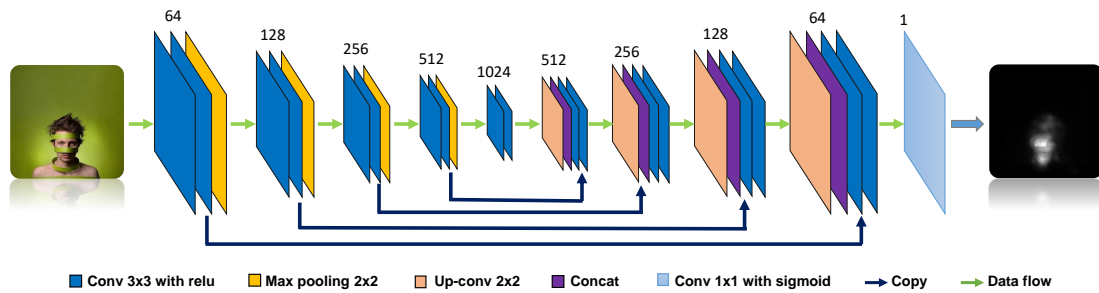


Fig. 2. The U-Shaped network implemented in this work for saliency map detection.

390 region and maps this single area to the final output region  
 391 by regression network directly. Thus, in this framework the  
 392 data flows through the network only once without extensively  
 393 assessing multiple candidates, which highly improves the  
 394 efficiency and maintains the accuracy in the meantime.

### 395 B. Saliency Map Generation Network

396 Saliency map detection aims at predicting visually interested  
 397 objects in an image that attract human attention. In the  
 398 proposed system, we adopt a modified U-shaped network to  
 399 produce the saliency map. As a variant of widely used fully  
 400 convolutional encoder-decoder, U-shaped network is originally  
 401 designed for semantic segmentation on biomedical images  
 402 [50]. It merges feature maps from convolutional layers to  
 403 deconvolutional layers gradually during the upsampling pro-  
 404 cedure. Thus, different types of features are preserved for the  
 405 semantic labeling task.

406 Particularly, in our implementation, the encoder for the U-  
 407 shaped network is composed by four fundamental blocks,  
 408 where every two convolutional layers followed by a max  
 409 pooling layer are stacked to form the basic block. Similarly, a  
 410 decoder is constructed by four basic blocks where every two  
 411 deconvolutional layers and a upsampling layer are used. For  
 412 each fundamental block in the encoder, its feature maps are  
 413 copied and concatenated directly to the corresponding block  
 414 in the decoder with the same size of feature dimensions. Thus,  
 415 the encoder is employed to extract features for the image  
 416 and the saliency map is produced based on the decoder. The  
 417 detailed structure of the U-shaped network we implemented is  
 418 illustrated in Fig. 2.

### 419 C. Aesthetic Area Regression Network

420 Based on the image feature obtained through feature ex-  
 421 tractor and saliency map detected by the saliency generator,  
 422 the relation between the interested objects and the area with  
 423 high aesthetic quality can be learned through the proposed  
 424 IOL layer and regression layers, which are described in the  
 425 following subsections with details.

426 1) *Soft Binarization Layer (SBL)*: To make our cropping  
 427 system less sensitive to the presence of outliers in the saliency  
 428 map, we introduce a function  $\rho(x; \sigma)$  to enhance the quality  
 429 of interested objects in saliency map, which is defined by:

$$429 \rho(x; \sigma) = \frac{x^2}{x^2 + \sigma^2}. \quad (4)$$

By selecting proper scale parameter  $\sigma$ ,  $\rho(x; \sigma)$  function  
 maps small value of pixel in saliency map to 0, and saliency  
 map will be saturated to 1 with larger pixel value. In Fig. 3,  
 we demonstrate a sample saliency map and its corresponding  
 enhanced version, from which we can observe that the dif-  
 ference between the interested objects in the image and the  
 background is enlarged and they can be distinguished with  
 minimum efforts.

As the derivative of  $\rho(x; \sigma)$  function is calculated by:

$$\frac{\partial \rho(x; \sigma)}{\partial x} = \frac{2x\sigma^2}{(x^2 + \sigma^2)^2}, \quad (5)$$

this operation can be easily integrated into the proposed neural  
 network pipeline without blocking the backpropagation. And  
 in our work,  $\sigma$  is empirically set as 0.01.

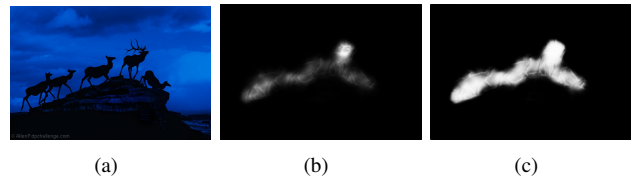


Fig. 3. Sample saliency map and enhanced interested object image for color image with high aesthetic scores. (a) The original high quality image. (b) The corresponding saliency map. (c) Enhanced saliency map by using function  $\rho(x; \sigma)$ , where the interested objects are easily distinguished from backgrounds.

2) *Interested Object Localization (IOL) Layer*: Based on  
 the obtained soft binarization saliency map  $\mathcal{S}$  that shows each  
 pixel's probability to be the interested objects, it is necessary  
 to model the  $P(\mathcal{R}|\mathcal{S}, \mathcal{I})$  to reveal the relation between the  
 interested objects and the final cropping window. In order to  
 represent this relation, it needs to extract the features of the  
 interested objects first. To achieve this goal, the location of  
 those interested objects needs to be determined. Generally, to  
 locate the interested objects in the image, researchers extract  
 the saliency map first and then search and locate the salient  
 region based on it. However, most salient region localization  
 methods use heuristic searching technique to scan all possible  
 candidate regions, which are prohibitively slow. Even by using  
 many speed up algorithms to reduce the searching space [32],  
 [34], [51], [6], [52], those approaches are not differentiable  
 which are infeasible in an end-to-end image cropping pipeline  
 to allow the backpropagation. Thus, in this work, we propose  
 an IOL layer that can effectively detect areas with interested

460 objects in the image and is differentiable for end-to-end  
461 training.

462 Inspired by the mean shift algorithm that is used to locate  
463 and track the face regions in videos [53], a region generation  
464 algorithm is proposed in this work to perform the interested  
465 object region creation.

466 Given a saliency map  $\mathcal{S}$  extracted by U-shaped network, the  
467 center of mass  $(c_x, c_y)$  for this map can be calculated by:

$$c_x = \frac{M_{10}}{M_{00}}, \quad c_y = \frac{M_{01}}{M_{00}},$$

and the standard deviation for the center of mass are obtained  
according to:

$$\sigma_x = \sqrt{\frac{M_{20}}{M_{00}} - c_x^2}, \quad \sigma_y = \sqrt{\frac{M_{02}}{M_{00}} - c_y^2},$$

where moments  $M_{00}$ ,  $M_{01}$ ,  $M_{10}$ ,  $M_{20}$  and  $M_{02}$  are calculated  
based on:

$$M_{00} = \sum_{i,j} S_{i,j} \quad (6)$$

$$M_{10} = \sum_{i,j} i \cdot S_{i,j}, \quad M_{01} = \sum_{i,j} j \cdot S_{i,j} \quad (7)$$

$$M_{20} = \sum_{i,j} i^2 \cdot S_{i,j}, \quad M_{02} = \sum_{i,j} j^2 \cdot S_{i,j}. \quad (8)$$

468 Therefore, a region that includes the energy of the saliency  
469 map can be defined through its top-left corner  $(x_{tl}^s, y_{tl}^s)$  and  
470 bottom-right corner  $(x_{br}^s, y_{br}^s)$  by using a Gaussian-like win-  
471 dow:

$$(x_{br}^s, y_{br}^s) = (c_x + \gamma\sigma_x, c_y + \gamma\sigma_y) \quad (9)$$

472 and

$$(x_{tl}^s, y_{tl}^s) = (c_x - \gamma\sigma_x, c_y - \gamma\sigma_y), \quad (10)$$

473 where  $\gamma$  is a hyper-parameter controlling the amount of energy  
474 contained in the window and maintaining the integrity of  
475 interested objects in the image. In this work,  $\gamma = 3.0$  is applied  
476 to include over 99% energy from the interested objects in the  
477 image.

478 In Fig. 4, the examples for areas of interested object  
479 obtained by the IOL layer with different  $\gamma$  are illustrated. From  
480 these figures we can see that the IOL layers with  $\gamma = 1.5$   
481 can only cover partial of interested objects in the image. And  
482 when  $\gamma = 3.0$ , most areas of interested objects can be included  
483 whereas the background of the image is still excluded.

484 To allow backpropagation of the loss pass through this  
485 region generation layer, the gradient of the coordinates for  
486 interested object region w.r.t  $\mathcal{S}$  can be defined. For coordinate  
487  $x$  of bottom-right corner, the partial derivative is given by:

$$\frac{\partial x_{br}^s}{\partial S_{i,j}} = \frac{\partial c_x}{\partial S_{i,j}} + \gamma \frac{\partial \sigma_x}{\partial S_{i,j}}, \quad (11)$$

where  $\frac{\partial c_x}{\partial S_{i,j}}$  and  $\frac{\partial \sigma_x}{\partial S_{i,j}}$  are further calculated based on follow-  
ing equations:

$$\begin{aligned} \frac{\partial c_x}{\partial S_{i,j}} &= \frac{1}{M_{00}} \frac{\partial M_{10}}{\partial S_{i,j}} - \frac{M_{10}}{M_{00}^2} \frac{\partial M_{00}}{\partial S_{i,j}} \\ &= \frac{i}{M_{00}} - \frac{M_{10}}{M_{00}^2} \end{aligned} \quad (12)$$

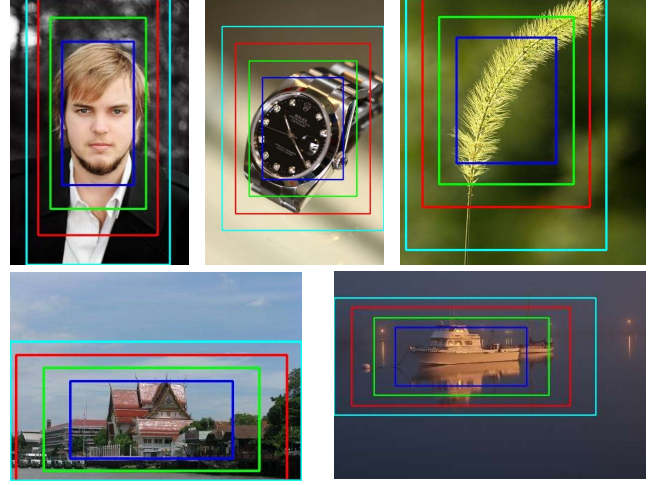


Fig. 4. The regions extracted by IOL layer with different  $\gamma$ . Blue box corresponds to  $\gamma = 1.5$ , green box is the area extracted by using  $\gamma = 2.0$ , red box is for  $\gamma = 2.5$  and cyan box means the region obtained by  $\gamma = 3.0$ .

and

$$\begin{aligned} \frac{\partial \sigma_x}{\partial S_{i,j}} &= \frac{1}{2\sqrt{\frac{M_{20}}{M_{00}} - c_x^2}} \\ &\times \left( \frac{1}{M_{00}} \frac{\partial M_{20}}{\partial S_{i,j}} - \frac{M_{20}}{M_{00}^2} \frac{\partial M_{00}}{\partial S_{i,j}} - 2c_x \frac{\partial c_x}{\partial S_{i,j}} \right) \\ &= \frac{1}{2\sqrt{\frac{M_{20}}{M_{00}} - c_x^2}} \\ &\times \left\{ \frac{i^2}{M_{00}} - \frac{M_{20}}{M_{00}^2} - 2c_x \left( \frac{i}{M_{00}} - \frac{M_{10}}{M_{00}^2} \right) \right\} \end{aligned} \quad (13)$$

and similar partial derivatives can be applied for  $\frac{\partial y_{br}^s}{\partial S_{i,j}}$ ,  $\frac{\partial x_{tl}^s}{\partial S_{i,j}}$   
and  $\frac{\partial y_{tl}^s}{\partial S_{i,j}}$ .

This provides the proposed IOL layer a mechanism that  
allows loss gradients to flow back to the input of network to  
update the model's parameters  $\Theta$ .

3) *Dense Layers*: Although the interested object region  
produced by the IOL layer contains most significant objects  
in the image, it is mostly far from having high aesthetic  
quality. So, based on the observation that the professional  
photographer tends to adjust the scene area for final shooting  
according to the interested objects, and the discovery that "*one  
may still roughly infer the extent of an object if only the middle  
of the object is visible*" [51], three fully connected layers  
are implemented to map the interested object region to the  
eventual cropping window with high visual quality based on  
its feature.

In our implementation, the region of interest (RoI) warping  
pooling layer followed by fully connected layers is used to  
estimate final cropping areas. The RoI warping pooling layer  
is proposed in [54], which takes two inputs: the coordinates  
of predicted interested object region and the feature maps  
generated from the bottle layer of U-shaped network. Prior  
to feeding into the RoI warping pooling layer, the coordinates  
of predicted interested object region are reduced 16 times to  
match the size of feature maps from bottle layer in U-shaped

513 network. In the RoI warping pooling layer, only features from  
 514 interested object region are extracted, which are consequently  
 515 passed to two fully connected layers with ReLU activation,  
 516 whose sizes are 2048 and 1024, respectively. The last layer of  
 517 this regression network is a fully connected layer who has 4  
 518 units with linear activation function, which predicts the four  
 519 coefficients defined by Eq. 16 and Eq. 17.

4) *Aesthetic Area Representation*: To represent the relation between the detected interested object region and areas with high aesthetic qualities, we use the approach described in [8]. Given a detected interested object region, whose size is  $w^s \times h^s$ , if its corresponding high aesthetic quality image's size is  $w^a \times h^a$ , and their top-left and bottom-right corners are  $(x_{tl}^s, y_{tl}^s)$ ,  $(x_{br}^s, y_{br}^s)$ ,  $(x_{tl}^a, y_{tl}^a)$  and  $(x_{br}^a, y_{br}^a)$ , respectively, the offsets between the corners of these two rectangles  $R((x_{tl}^s, y_{tl}^s), (x_{br}^s, y_{br}^s))$  and  $R((x_{tl}^a, y_{tl}^a), (x_{br}^a, y_{br}^a))$  can be represented as:

$$(\Delta x_t, \Delta y_t) = (x_{tl}^s, y_{tl}^s) - (x_{tl}^a, y_{tl}^a) \quad (14)$$

$$(\Delta x_b, \Delta y_b) = (x_{br}^a, y_{br}^a) - (x_{br}^s, y_{br}^s). \quad (15)$$

520 Hence, the height and width of these two rectangles can be  
 521 expressed as:

$$h^a = h^s + \Delta y_t + \Delta y_b = h^s + \alpha_t \cdot h^a + \alpha_b \cdot h^a \quad (16)$$

522 and

$$w^a = w^s + \Delta x_t + \Delta x_b = w^s + \beta_t \cdot w^a + \beta_b \cdot w^a, \quad (17)$$

523 where  $\mathcal{O} = [\alpha_t, \alpha_b, \beta_t, \beta_b]$  are four coefficients.

524 In our implementation, the above coefficients  $\mathcal{O} =$   
 525  $[\alpha_t, \alpha_b, \beta_t, \beta_b]$  are used to represent the final aesthetic area  
 526 and can be learned through a neural network.

During the testing, the corner coordinates  $(x_{tl}^s, y_{tl}^s)$ ,  $(x_{br}^s, y_{br}^s)$  of the interested object region of input image and four coefficients  $\mathcal{O} = [\alpha_t, \alpha_b, \beta_t, \beta_b]$  are predicted by the proposed end-to-end network. Then, the width and height of final aesthetic region can be expressed as follows:

$$h^a = \frac{y_{br}^s - y_{tl}^s}{1 - \alpha_t - \alpha_b} \quad (18)$$

$$w^a = \frac{x_{br}^s - x_{tl}^s}{1 - \beta_t - \beta_b} \quad (19)$$

Thus, the coordinates of top-left and bottom-right corners of aesthetic area can be calculated by:

$$\begin{aligned} x_{tl}^a &= x_{tl}^s - \beta_t \cdot w^a & y_{tl}^a &= y_{tl}^s - \alpha_t \cdot h^a \\ x_{br}^a &= x_{br}^s + \beta_b \cdot w^a & y_{br}^a &= x_{br}^s + \alpha_b \cdot h^a \end{aligned}$$

#### 527 D. Loss Functions for the Cropping System

528 As introduced in subsection III-A, the total loss of the  
 529 proposed neural network based cropping system is given by:

$$\mathcal{L}_{total} = \frac{1}{N} \sum_{k=1}^N (\mathcal{L}_s(\cdot) + \lambda \mathcal{L}_r(\cdot)) \quad (20)$$

530 where  $\mathcal{L}_s(\cdot)$  is the loss from saliency map detection network  
 531 and  $\mathcal{L}_r(\cdot)$  is the loss from aesthetic regression network,  $N$   
 532 means the total training number, and  $\lambda$  is the weight control-  
 533 ling the influence from these two networks.

To train the U-shaped based network  $H(\mathcal{I}, \Theta_s)$ , the binary cross-entropy of each pixel is calculated:

$$\begin{aligned} H(\mathcal{I}_{i,j}; \Theta_s) &= (\mathbf{W}_s, \mathbf{b}_s) \\ &= -S_{i,j} \log p(\mathcal{I}_{i,j}; (\mathbf{W}_s, \mathbf{b}_s)) \\ &\quad - (1 - S_{i,j}) \log (1 - p(\mathcal{I}_{i,j}; (\mathbf{W}_s, \mathbf{b}_s))), \end{aligned} \quad (21)$$

534 where  $[\mathbf{W}_s, \mathbf{b}_s]$  are weights of U-shaped saliency map de-  
 535 tection network,  $p(\mathcal{I}_{i,j}; (\mathbf{W}_s, \mathbf{b}_s))$  stands for the predicted  
 536 confidence for the interested objects of each pixel, and  $\hat{\mathcal{S}}_{i,j} =$   
 537  $p(\mathcal{I}_{i,j}; (\mathbf{W}_s, \mathbf{b}_s))$  holds for the detected saliency map  $\hat{\mathcal{S}}$ .

Thereafter, the loss for a given image  $\mathcal{I}^{(k)}$  can be expressed as:

$$\begin{aligned} \mathcal{L}_s(\hat{\mathcal{S}}^{(k)}, \mathcal{S}^{(k)}) &= \mathcal{L}_s(\mathbf{W}_s, \mathbf{b}_s) \\ &= \sum_{\mathcal{I}_{i,j}^{(k)}} H(\mathcal{I}_{i,j}^{(k)}; (\mathbf{W}_s, \mathbf{b}_s)), \end{aligned} \quad (22)$$

538 where superscript  $k$  is the index of the training sample.

And as described in section III-C4, unlike other image cropping methods that train a ranker or classifier to evaluate the cropping areas' aesthetic quality by using training samples with high/low qualities, the proposed aesthetic region regression network uses a regressor to predict the cropping window, where only features from high aesthetic images are required and learned. Thus, in our training phase for the regression network, the interested object region  $\mathcal{R}((x_{tr}^s, y_{tr}^s), (x_{bl}^s, y_{bl}^s))$  for high aesthetic quality image is firstly detected by IOL layer. Then, the region of original high quality image  $\mathcal{R}((x_{tr}^a, y_{tr}^a), (x_{bl}^a, y_{bl}^a))$  is used to calculate the offsets coefficients  $\mathcal{O} = [\alpha_t, \alpha_b, \beta_t, \beta_b]$ , where  $(x_{tr}^a, y_{tr}^a) = (0, 0)$ ,  $(x_{bl}^a, y_{bl}^a) = (w^a, h^a)$  and  $w^a \times h^a$  is high quality image's size. And these offsets coefficients are used to supervise the training of the proposed regression network, where L2 loss is applied according to:

$$\begin{aligned} \mathcal{L}_r(\hat{\mathcal{R}}^{(k)}, \mathcal{R}^{(k)}) &= \mathcal{L}_r(\mathbf{W}_r, \mathbf{b}_r) \\ &= \|\hat{\mathcal{O}}^{(k)} - \mathcal{O}^{(k)}\|^2 \end{aligned} \quad (23)$$

539 where  $[\mathbf{W}_r, \mathbf{b}_r]$  are system weights for aesthetic regression  
 540 network,  $\mathcal{O}^{(k)}$  is the ground truth of offsets coefficients of an  
 541 image  $\mathcal{I}^{(k)}$  and  $\hat{\mathcal{O}}^{(k)}$  is the corresponding predicted offsets  
 542 coefficients.

## 543 IV. EXPERIMENTS

### 544 A. Databases and Evaluation Protocol

545 We conducted our experiments on the following four  
 546 databases.

1) *Training database*: In this experiment, the AVA database [55] was used for training. The AVA database, which was originally designed for aesthetic visual analysis, gathered more than 250,000 images from www.dpchallenge.com. Each image in AVA set contains plenty of meta-data, including multiple aesthetic scores from reviewers, semantic labels for over 60 categories, etc. In this work, we utilized AVA database to train the proposed end-to-end image cropping network, where only images whose average aesthetic scores were greater than or equal to 6 were selected for training, which resulted in a

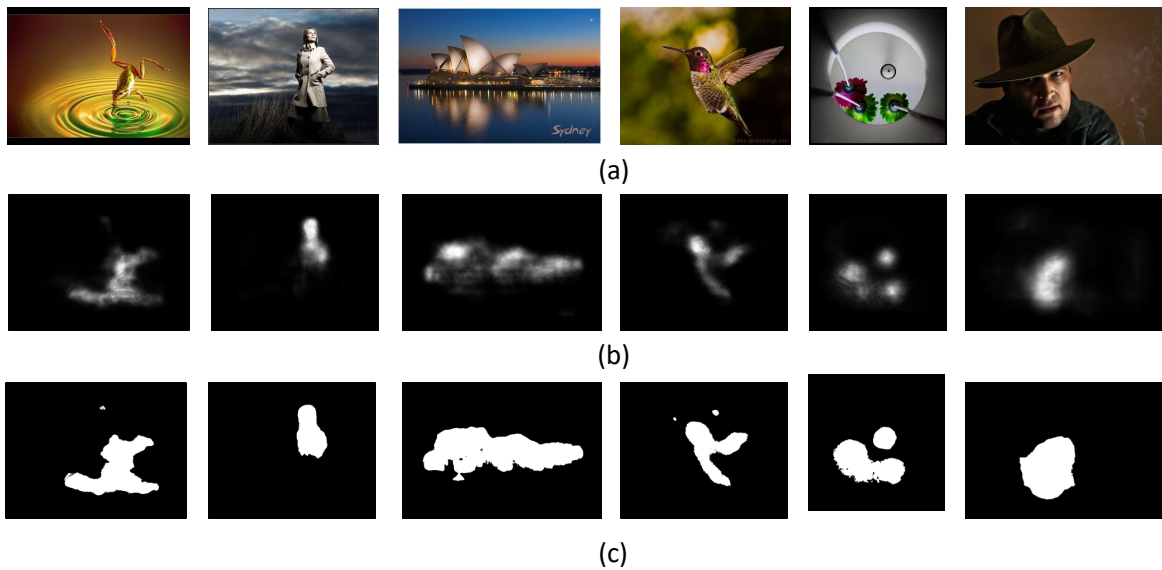


Fig. 5. Sample images along with their ground truths from AVA database. (a) Sample images with high aesthetic scores from AVA database. (b) Corresponding saliency maps for sample images. (c) Binarized interested object image for AVA samples where threshold is 0.12.

557 training set with 50,189 high qualities images. Sample images  
558 from AVA database can be found in Figure 5 (a).

559 However, the AVA database was originally designed for  
560 aesthetic evaluation and only aesthetic scores were provided  
561 as the ground truths for each image. So in order to train the  
562 proposed neural networks with this database, the synthetic  
563 ground truths of interested object image and offsets of final  
564 crop window w.r.t. the interested object region for each image  
565 were produced initially. The preprocessing details for the  
566 training database is described in section IV-B accordingly.

567 2) *Test databases:* In our experiments, three public  
568 databases were applied for evaluation purpose.

569 The FCD database [56] was constructed to facilitate the  
570 aesthetic cropping task, where thousands images were col-  
571 lected from Flickr and cleaned by annotators. For each cleaned  
572 image, the cropping area was labeled by professional pho-  
573 tographers and validated by multiple professional annotators  
574 who had passed Human Intelligence Tasks qualification test.  
575 Only those images that were ranked as preferable by at least  
576 4 professional annotators were selected in the final cropping  
577 database. In our experiments, 334 samples were applied for  
578 evaluation purpose among this database.

579 FLMS database [3] collected 500 images from Flickr and  
580 the best cropping areas of each image were manually annotated  
581 by 10 experienced editors. In this work, we used FLMS  
582 database to evaluate the cropping performance.

583 Furthermore, to measure the proposed image cropping  
584 method, CUHK-ICD database [29] was employed. In this  
585 database, 950 images were captured by amateur photographers  
586 but cropped by 3 professional editors. All images in this  
587 database were used for evaluation in our work.

588 To quantitatively evaluate the cropping performance, the  
589 intersection over union (IoU) and boundary displacement error  
590 (BDE) were employed, where IoU is defined as:

$$IoU = \frac{A' \cap \hat{A}}{A' \cup \hat{A}} \quad (24)$$

and BDE is defined by:

$$BDE = \sum_{k=1}^4 \left\| B'_k - \hat{B}_k \right\| / 4. \quad (25)$$

592 Here,  $A'$  means the ground truth of the cropping area,  $\hat{A}$   
593 represents the predicted cropping region, and  $B'_k$  and  $\hat{B}_k$  are  
594 the normalized boundary coordinates for ground truth and  
595 predicted crop windows, respectively.

### B. Neural Networks Training

596 Because the proposed image cropping system contained two  
597 main components conceptually, the parameters' search space  
598 is large with the joint training of entire network, which causes  
599 the low efficiency and unstable training. So, in this work, a  
600 corresponding three-stage training scheme was applied, where  
601 the U-shaped saliency map detection network  $H(\mathcal{I}, \Theta_s)$  and  
602 the regression network  $G(\mathcal{I}, \Theta_r)$  were trained sequentially and  
603 the entire network was fine-tuned afterwards.

604 To train the U-shaped network, the images from AVA  
605 database were employed, where 50,189 images with their  
606 synthetic saliency maps were fed into the network for training.  
607 The synthetic saliency maps of AVA database were obtained  
608 by using method in [57], where an existing single branched  
609 DNN model was applied to detect the saliency maps of  
610 each image for AVA. Based on the obtained saliency map,  
611 the binarized image can be calculated by using a simple  
612 thresholding approach with empirical threshold, which was  
613 employed to guide the training of the proposed U-shaped  
614 network. In Figure 5 (b) and 5 (c), the corresponding saliency  
615 maps and binarized interested objects for the sample images  
616 from AVA database are shown.

617 In this experiment, SGD optimization scheme was applied  
618 and the training rate was fixed to  $1 \times 10^{-4}$  for 4 epochs.

619 Once the U-shape network was learned, the obtained  
620 weights were locked for the second stage training of regression  
621



622 network. To train this regression network, the same training  
623 images with high qualities from AVA database were fed  
624 into U-shaped network to create saliency maps, from which  
625 an interested object region can be estimated based on the  
626 proposed IOL layer subsequently. Then, the coordinates of  
627 this interested object region were passed to RoI warping layer,  
628 where the corresponding features from U-shaped network were  
629 extracted and sent to the following fully connected layers, as  
630 illustrated in Figure 1.

631 The pre-calculated ground truths for offsets, which were  
632 obtained based on the method described in subsection III-D,  
633 were used to guide the training of regression network to  
634 predict the offset between interested object region and the final  
635 cropping rectangle. We used SGD optimizer with learning rate  
636 of  $1 \times 10^{-4}$  for 6 epochs in this training stage.

637 Finally, we used training images from AVA database along  
638 with the synthetic binarized saliency maps and pre-calculated  
639 offsets ground truths to fine-tune the entire network from end  
640 to end. The SGD optimizer was used in this stage with learning  
641 rate of  $1 \times 10^{-5}$  for 2 epochs and the U-shaped saliency map  
642 detection network and aesthetic area regression network have  
643 the same loss weight.

644 In this three-stage training phase, the input images were  
645 resized so that the shorter side of the image was 224 but the  
646 original aspect ratio was maintained.

### 647 C. Results Evaluation & Analysis

648 1) *Comparison with the state-of-the-art approaches:* To  
649 analyze the performance of the proposed end-to-end image  
650 cropping model, we compared the proposed method with other  
651 state-of-the-art cropping approaches, which were used as our  
652 baselines.

653 In table I, the IoUs and BDEs of the proposed cropping sys-  
654 tem and other state-of-the-arts cropping approaches on three  
655 public datasets are demonstrated, where \* denotes weakly  
656 supervised cropping methods that do not use bounding boxes  
657 from annotated cropping datasets for training. As can be  
658 seen from this table, our cropping method obtained better  
659 cropping performances than any other approach on FLMS  
660 set. On FCD dataset, the proposed method achieved best  
661 performance among weakly supervised cropping approaches.  
662 And for CUHK-ICD database, the proposed method had  
663 competitive IoU and BDE performance on this evaluation  
664 set, which shows the effectiveness of the proposed cropping  
665 method.

666 In Fig. 6, multiple cropping results along with the cor-  
667 responding detected saliency maps from the evaluation sets  
668 are demonstrated, where red boxes represent the optimized  
669 cropping window predicted by the proposed system and the  
670 green boxes show the detected IOLs based on Eq. 9 and Eq.  
671 10 for saliency maps. From these images, we can see that  
672 the cropped images obtain better composition and aspect ratio  
673 than the original images, especially for those amateur captured  
674 low quality images.

675 2) *Ablation test:* To investigate the effectiveness of the  
676 proposed soft binarization layer (SBL), one cropping system  
677 was training where the SBL was removed. We illustrated and

678 compared the cropping systems with/without SBL in the last  
679 two row of table I. From these numbers, we can observe  
680 that the IoU on CUHK-ICD dataset for the cropping system  
681 with SBL is higher than its counterpart without SBL for more  
682 than 5.0 on average. And the cropping results by the system  
683 with SBL on other test sets are also superior than the system  
684 without SBL. From these results, we can see that SBL can  
685 effectively help the cropping system to filter the noises and  
686 find the interested objects more accurately.

687 In the proposed image cropping framework, the U-shape  
688 based saliency map generation network can be easily replaced  
689 by other state-of-the-art saliency detection modules. Thus,  
690 in our experiments, we re-trained the SALICON saliency  
691 detection network, which was introduced in [57] and applied  
692 to generate the synthetic ground truth for AVA database in  
693 section IV-A, to detect the saliency maps for test images and  
694 consequently feed them into the IOL layer and aesthetic area  
695 regression network to produce the final cropping window. In  
696 table II, the overall cropping performances by combining the  
697 SALICON saliency map detection network and the proposed  
698 aesthetic area regression network are listed, where we can  
699 see it provides similar cropping results compared with the  
700 U-shape based saliency detection network, which shows the  
701 generalization capability of the proposed framework.

702 Need to note that in the ablation test, in order to avoid  
703 the size of feature maps extracted by SALICON saliency map  
704 detection network being too small, the input images of the  
705 neural networks were resized to ensure the shorter side of the  
706 image was 512 with the original aspect ratio.

707 By analyzing three tables and the structure of SALICON  
708 network and U-shaped network, it can be concluded that the  
709 cropping performance differences between these two saliency  
710 detection modules rely on the resolution of extracted features  
711 from these two networks. For the SALICON network, it  
712 applies the VGG-16 to extract down-sampled feature maps  
713 for images, which provides coarse details of the interested  
714 objects. But U-shaped saliency detection network extracts the  
715 feature map whose size is the same as the input image, that  
716 maintains more details of the interested objects. Therefore, for  
717 the cleaned high resolution images, such as photos from AVA  
718 database, U-shaped saliency detection network tends to extract  
719 more pleasant features of the interested objects in the image to  
720 help the cropping task. And for the noisy low quality images,  
721 SALICON network acts more like a noise suppressor to extract  
722 smoothed features to boost the cropping performance, as we  
723 observed from FCD database.

724 3) *Investigation of image's size and aspect ratio:* In many  
725 other research articles, it is claimed that the aesthetic quality  
726 of images is highly relied on the size or aspect ratio of the  
727 images [59], [60]. Thus, we carried out several experiments  
728 to investigate cropping performance with different image size  
729 and aspect ratio. In these experiments, we trained three models  
730 by keeping the original aspect ratio of the images but resizing  
731 the image till the shorter side of the image is 224, 384 or  
732 512. Other three models were trained by resizing the image  
733 to square, whose size is  $224 \times 224$ ,  $384 \times 384$  or  $512 \times 512$ ,  
734 respectively.

735 In table III, we list the IoUs and BDEs of different models

TABLE I  
THE CROPPING PERFORMANCE FOR DIFFERENT APPROACHES ON CUHK-ICD, FLMS AND FCD DATABASES.

Approach	CUHK-ICD						FLMS		FCD	
	Photographer1		Photographer2		Photographer3		IoU	BDE	IoU	BDE
	IoU	BDE	IoU	BDE	IoU	BDE				
*ATC [5]	0.605	0.108	0.628	0.100	0.641	0.095	0.720	0.063	0.58	0.10
*AIC [6]	0.469	0.142	0.494	0.131	0.512	0.123	0.640	0.075	0.47	0.13
*MPC [58]	0.603	0.106	0.582	0.112	0.608	0.110	0.410	N/A	N/A	N/A
*A2-RL [2]	0.802	0.052	0.796	0.054	0.790	0.054	0.820	N/A	N/A	N/A
*ABP-AA [1]	0.815	0.031	0.810	0.030	0.830	0.029	0.810	0.057	0.65	0.08
*VFN-SW [45] [2]	0.740	0.069	0.719	0.076	0.713	0.077	N/A	N/A	0.633	0.098
*Lu <i>et al.</i> [8]	0.827	0.032	0.816	0.035	0.805	0.036	0.843	0.029	0.659	0.062
VEN [48]	N/A	N/A	N/A	N/A	N/A	N/A	0.837	0.041	0.735	0.072
LCC [29]	0.748	0.066	0.728	0.072	0.732	0.071	0.630	N/A	N/A	N/A
*proposed w/o SBL	0.777	0.039	0.766	0.043	0.759	0.043	0.820	0.031	0.655	0.060
*Proposed w/ SBL	0.822	0.031	0.815	0.034	0.802	0.035	0.846	0.026	0.673	0.058

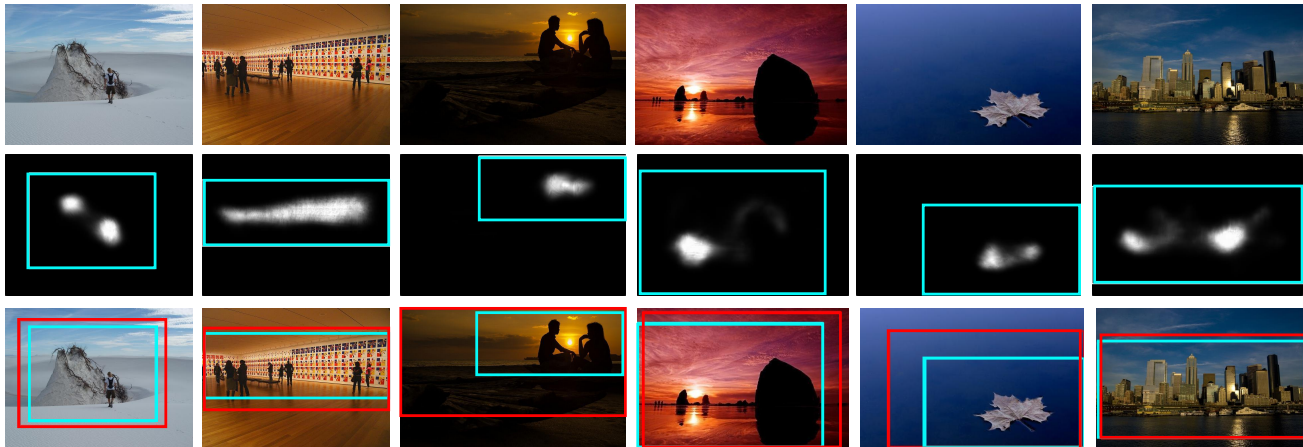


Fig. 6. Cropping rectangle produced by the proposed system.

TABLE II  
THE CROPPING PERFORMANCE USING SALICON BASED SALIENCY DETECTION NETWORK ON CUHK-ICD, FLMS AND FCD DATABASES.

Method	CUHK-ICD						FLMS		FCD	
	Photographer1		Photographer2		Photographer3		IoU	BDE	IoU	BDE
	IoU	BDE	IoU	BDE	IoU	BDE				
Salicon + Regression	0.819	0.032	0.808	0.036	0.799	0.037	0.838	0.028	0.666	0.060
U-shaped + Regression	0.825	0.032	0.820	0.034	0.806	0.036	0.845	0.028	0.664	0.060

736 with various input image size on the three public test sets.  
 737 From this table, it can be seen that both IoUs and BDEs  
 738 for different input size of images have similar performances and  
 739 no significant difference can be found between these models,  
 740 which means the proposed image cropping model is insensitive  
 741 to the size and aspect ratio of the input image.

742 By digging into table III, we observed that the overall IoU  
 743 and BDE scores on CUHK-ICD database were getting better  
 744 with larger input image size of neural networks, whilst the  
 745 cropping performance was degraded on FCD database with  
 746 larger input image size. The main reason of this phenomenon  
 747 is that the images in the FCD database were collected from  
 748 Flickr’s website, containing more irrelevant background noises  
 749 than the training database and other two evaluation databases.  
 750 With a larger input size, more detailed features of images from  
 751 FCD database, including non-interested background noises,  
 752 can be discovered by the neural networks. But these features  
 753 of noises cannot be effectively represented by the neural  
 754 network which was trained based on the clean images from  
 755 AVA database, and can be easily mis-represented as objects’

756 features. This causes the proposed cropping method tending  
 757 to generate larger crop windows to include more details when  
 758 the input image size is big, which degrades the performance of  
 759 the system on FCD database. But for the FLMS database, each  
 760 test image had multiple annotated ground truths and the best  
 761 cropping result was calculated using the ground truth which  
 762 provided best performance. So, FLMS set is less sensitive  
 763 to the neural network’s input image size. With regard to  
 764 CUHK-ICD database, it was constructed by high aesthetic  
 765 quality images similar to AVA database, which in turn can be  
 766 sufficiently embedded by the proposed cropping networks with  
 767 larger input image size to attain better cropping performance.

768 Compromised by the cropping performance across three  
 769 evaluation sets and the computation efficiency, it is preferable  
 770 to resize the input image such that the shorter side is 224 for  
 771 the proposed cropping approach.

772 4) Investigation of parameters  $\sigma$  and  $\gamma$ : From subsection  
 773 III-C1, we can see that the quality of interested objects  
 774 within the obtained saliency map can be enhanced by function  
 775  $\rho(x; \sigma)$ . To investigate the impact from scale parameter  $\sigma$

TABLE III  
THE CROPPING PERFORMANCE FOR DIFFERENT ASPECT RATIO AND IMAGE SIZE ON CUHK-ICD, FLMS AND FCD DATABASES.

Input size	CUHK-ICD						FLMS		FCD	
	Photographer1		Photographer2		Photographer3		IoU	BDE	IoU	BDE
	IoU	BDE	IoU	BDE	IoU	BDE				
224 × 224	0.825	0.031	0.818	0.034	0.805	0.036	0.840	0.028	0.672	0.059
384 × 384	0.827	0.031	0.817	0.034	0.804	0.036	0.843	0.028	0.670	0.059
512 × 512	0.828	0.031	0.822	0.034	0.806	0.036	0.842	0.028	0.665	0.061
min(w, h) = 224	0.822	0.031	0.815	0.034	0.802	0.035	0.846	0.026	0.673	0.058
min(w, h) = 384	0.823	0.032	0.818	0.034	0.804	0.036	0.844	0.027	0.670	0.059
min(w, h) = 512	0.825	0.032	0.820	0.034	0.806	0.036	0.845	0.028	0.664	0.060

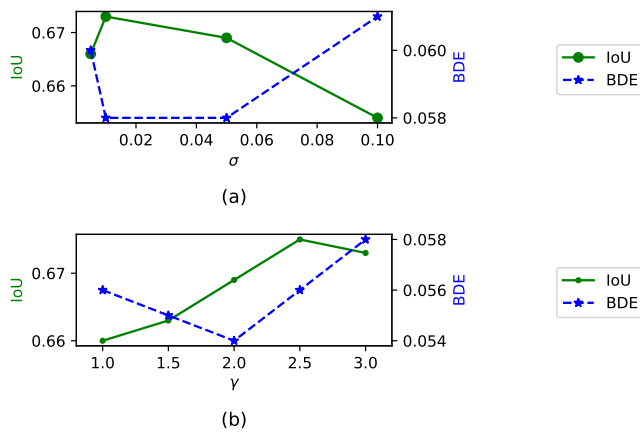


Fig. 7. The cropping performance for different  $\sigma$  and  $\gamma$  on FCD database.

of function  $\rho(x; \sigma)$  for the cropping system, cropping performances (IoUs and BDEs) with different  $\sigma$ s on test set FCD are shown in Fig. 7 (a), where function  $\rho(x; \sigma)$  with  $\sigma = \{0.005, 0.01, 0.05, 0.1\}$  are used. As can be seen from this figure, too small or large  $\sigma$  can cause lower cropping performance because more noises are introduced into IOLs or more interested objects are filtered out. Thus, a proper  $\sigma$  close to 0.01 provides us a better cropping result.

Similar to  $\sigma$ , the  $\gamma$  in Eq. 9 and Eq. 10 controls the amount of energy contained in the IOLs where larger  $\gamma$  causes more saliency areas but includes more background area also, and smaller  $\gamma$  results in the loss of integrity for salient object in IOLs. This effect can be seen in Fig. 7 (b), from which we can observe the cropping performances are getting degraded when  $\gamma$  is larger than 3.0 and smaller than 2.0 on the FCD test set.

5) *Efficiency analysis:* As one of the main contributions of this work is to use an end-to-end neural network to accomplish the image cropping task, without iteratively evaluating multiple candidates' aesthetic qualities, which has lower computational cost. So we measured time efficiency of the proposed system with different input size on the FLMS set, where the experiments were implemented with Keras on a server with Intel(R) Xeon(R) E5-2620 CPU @ 2.10GHz, 64Gb Memory and Nvidia 2080Ti GPU. We also compared our cropping method with other approaches w.r.t. speed, where the FPSs are shown in table IV. From this table, we notice that when the input image is resized to  $224 \times 224$ , the overall time for image cropping of our system is less than 20ms. Thus, the proposed system can reach over 50fps on average for real-time

processing, which is much faster than other state-of-the-arts approaches and shows its high efficiency.

Furthermore, by comparing the time efficiency with the cropping method presented in [8] which is relied on a brute force search algorithm [6], the proposed cropping system is five times faster.

TABLE IV  
THE TIME EFFICIENCY COMPARISON OF THE PROPOSED CROPPING SYSTEM WITH DIFFERENT SETTINGS AND OTHER METHODS.

Method	FPS
A2RL [2]	4
ABP-AA [1]	5
VFN [45]	0.5
Lu <i>et al.</i> [8]	10
proposed [ $224 \times 224$ ]	52
proposed [ $384 \times 384$ ]	29
proposed [ $512 \times 512$ ]	18
proposed [ $\min(w, h) = 224$ ]	40
proposed [ $\min(w, h) = 384$ ]	22
proposed [ $\min(w, h) = 512$ ]	14

#### D. Subjective analysis

Because image's aesthetics is difficult to represent from the subjective perspective, such that different person might have different views for the same cropping results based on their tastes, education backgrounds, etc. So, in our work, a subjective comparative experiment was carried out.

In this experiment, 200 images were randomly collected from three (CUHK-ICD, FLMS and FCD) test sets. For each image, the proposed cropping method, along with the algorithms AIC [6], A2-RL [2] and VEN [48] was employed to obtain four cropping results. Then, 10 users were recruited, including 5 males and 5 females. All users had no prior knowledge of the experiment content and the databases. For each participant, the four cropping results of each test image were presented, where the order of the cropping images was randomized and the users were asked to vote the most pleasing one in terms of their aesthetics. Finally, 2000 votes from 10 participants were received and shown in figure 8.

As can be seen from this figure, the proposed method had gained most votes (792/2000) among four state-of-the-art cropping approaches, which shows the proposed system provided more pleasing cropping results than the other methods in respect to the aesthetics.

#### E. Case Study

To analyze the effects of different contents for the cropping performance, both success and failure cases in the evaluation are demonstrated.

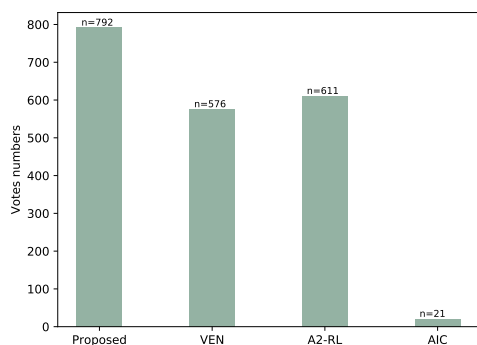


Fig. 8. Votes received from users for different state-of-the-art cropping methods.

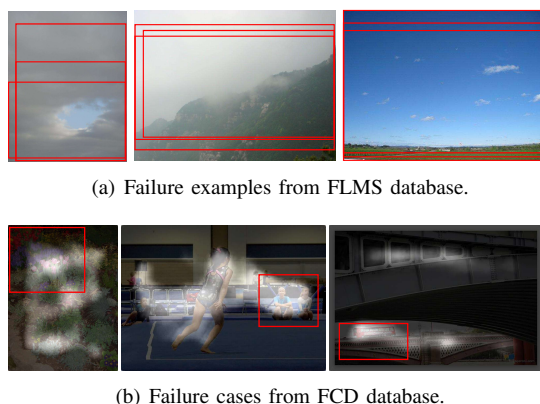


Fig. 9. Failure examples from the evaluation sets. (a) Failure samples from FLMS database, where red boxes are cropping windows by annotators. (b) Failure images from FCD database, where red boxes are ground truth and light areas are the detected saliency maps.

As shown in Fig. 6, when the interested object region are obtained from the saliency map successfully, the relations between the interested objects and the final cropping window can be learned by the proposed cropping system with sufficient training samples, where the area with the high aesthetic scores can also be inferred consequently.

Although the proposed image cropping approach works well on the majority of testing images, several failure cases can be found in the evaluation, which can be categorized into two types of errors broadly.

The first type of failures are mostly from the FLMS database, as shown in the Fig. 9(a), where only spurious texture regions exist in the image and it is hard to find enough salient pixels to determine the interested objects in the image. In our implementation, if no visual fixation is found, we use center areas that cover the 70% of entire image as our interested object region to feed into regression network to obtain the final cropping rectangle. The other type of failure cases can be seen in FCD database, where multiple interested objects are located by the saliency map detection network, as shown in Fig. 9(b), but only partial of these saliency area is included into the ground truth and most parts are missing, which causes the low IoUs and high BDEs.

## V. CONCLUSION

In this paper, an end-to-end automatic image cropping system is proposed to learn the relationship between the interested objects and the areas with high aesthetic scores in an image through a DNN. Conceptually, the saliency map is initially detected by using a U-shaped neural network, which is then passed into a soft binarization layer to separate objects from the background. Based on this enhanced saliency map, an interested object region is determined by the proposed IOL layer, which is fed into a ROI warping pooling layer and following dense layers along with the features of the interested objects, to predict the optimal cropping region with high aesthetic scores.

As a weakly supervised cropping method, the proposed algorithm outperforms other weakly supervised state-of-the-art cropping methods w.r.t IoU and BDE metrics. Moreover, because the proposed approach finds the final cropping areas based on the hidden relationship between interested objects and areas with high aesthetics quality through neural networks, which avoids to iteratively evaluate multiple cropping candidates, high processing efficiency is achieved with 50 FPS.

Our future research will be exploring other cropping metric instead of IoU and BDE to measure the performance of aesthetics based cropping system.

## REFERENCES

- [1] W. Wang, J. Shen, and H. Ling, "A deep network solution for attention and aesthetics aware photo cropping," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 41, no. 7, pp. 1531–1544, July 2019.
- [2] D. Li, H. Wu, J. Zhang, and K. Huang, "A2-RL: Aesthetics aware reinforcement learning for image cropping," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8193–8201.
- [3] C. Fang, Z. Lin, R. Mech, and X. Shen, "Automatic image cropping using visual composition, boundary simplicity and content preservation models," in *Proceedings of the 22Nd ACM International Conference on Multimedia*, 2014, pp. 1105–1108.
- [4] M. Nishiyama, T. Okabe, Y. Sato, and I. Sato, "Sensation-based photo cropping," in *Proceedings of the 17th ACM International Conference on Multimedia*, ser. MM '09, 2009, pp. 669–672.
- [5] B. Suh, H. Ling, B. B. Bederson, and D. W. Jacobs, "Automatic thumbnail cropping and its effectiveness," in *Proceedings of the 16th Annual ACM Symposium on User Interface Software and Technology*, 2003, pp. 95–104.
- [6] J. Chen, G. Bai, S. Liang, and Z. Li, "Automatic image cropping: A computational complexity study," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 507–515.
- [7] W. Wang and J. Shen, "Deep cropping via attention box prediction and aesthetics assessment," in *2017 IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 2205–2213.
- [8] P. Lu, H. Zhang, X. Peng, and X. Peng, "Aesthetic guided deep regression network for image cropping," *Signal Processing: Image Communication*, vol. 77, pp. 1 – 10, 2019.
- [9] X. Wang, H. Ma, and X. Chen, "Salient object detection via fast r-cnn and low-level cues," in *2016 IEEE International Conference on Image Processing (ICIP)*, 2016, pp. 1042–1046.
- [10] N. Liu and J. Han, "Dhsnet: Deep hierarchical saliency network for salient object detection," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 678–686.
- [11] Y. Ji, H. Zhang, and Q. J. Wu, "Salient object detection via multi-scale attention cnn," *Neurocomputing*, vol. 322, pp. 130 – 140, 2018.
- [12] Y. Li, X. Hou, C. Koch, J. M. Rehg, and A. L. Yuille, "The secrets of salient object segmentation," in *2014 IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 280–287.
- [13] L. Itti and C. Koch, "Comparison of feature combination strategies for saliency-based visual attention systems," in *Human Vision and Electronic Imaging IV*, vol. 3644, 1999, pp. 473 – 482.

- [14] —, “A saliency-based search mechanism for overt and covert shifts of visual attention,” *Vision Research*, vol. 40, no. 10, pp. 1489 – 1506, 2000.
- [15] J. Harel, C. Koch, and P. Perona, “Graph-based visual saliency,” in *Proceedings of the 19th International Conference on Neural Information Processing Systems*, ser. NIPS’06, 2006, pp. 545–552.
- [16] L. Itti, C. Koch, and E. Niebur, “A model of saliency-based visual attention for rapid scene analysis,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 20, no. 11, pp. 1254–1259, 1998.
- [17] T. Judd, K. Ehinger, F. Durand, and A. Torralba, “Learning to predict where humans look,” in *2009 IEEE 12th International Conference on Computer Vision*, 2009, pp. 2106–2113.
- [18] E. Vig, M. Dorr, and D. Cox, “Large-scale optimization of hierarchical features for saliency prediction in natural images,” in *2014 IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 2798–2805.
- [19] P. Moon and D. E. Spencer, “Geometric formulation of classical color harmony,” *J. Opt. Soc. Am.*, vol. 34, no. 1, pp. 46–59, Jan 1944.
- [20] T. Grill and M. Scanlon, *Photographic composition*. New York, N.Y.: Amphot, 1983.
- [21] Y. Ke, X. Tang, and F. Jing, “The design of high-level features for photo quality assessment,” in *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR’06)*, 2006, pp. 419–426.
- [22] R. Datta, D. Joshi, J. Li, and J. Z. Wang, “Studying aesthetics in photographic images using a computational approach,” in *Computer Vision – ECCV 2006*, A. Leonardis, H. Bischof, and A. Pinz, Eds., 2006, pp. 288–301.
- [23] L. Zhang, Y. Gao, R. Zimmermann, Q. Tian, and X. Li, “Fusion of multichannel local and global structural cues for photo aesthetics evaluation,” *IEEE Transactions on Image Processing*, vol. 23, no. 3, pp. 1419–1429, March 2014.
- [24] X. Tang, W. Luo, and X. Wang, “Content-based photo quality assessment,” *IEEE Transactions on Multimedia*, vol. 15, no. 8, pp. 1930–1943, 2013.
- [25] L. Marchesotti, F. Perronnin, D. Larlus, and G. Csurka, “Assessing the aesthetic quality of photographs using generic image descriptors,” in *2011 International Conference on Computer Vision*, 2011, pp. 1784–1791.
- [26] Y. Kao, R. He, and K. Huang, “Deep aesthetic quality assessment with semantic information,” *IEEE Transactions on Image Processing*, vol. 26, no. 3, pp. 1482–1495, 2017.
- [27] P. Lu, X. Peng, X. Zhu, and R. Li, “An EL-LDA based general color harmony model for photo aesthetics assessment,” *Signal Processing*, vol. 120, pp. 731 – 745, 2016.
- [28] Z. Liu, Z. Wang, Y. Yao, L. Zhang, and L. Shao, “Deep active learning with contaminated tags for image aesthetics assessment,” *IEEE Transactions on Image Processing*, in press.
- [29] J. Yan, S. Lin, S. B. Kang, and X. Tang, “Learning the change for automatic image cropping,” in *2013 IEEE Conference on Computer Vision and Pattern Recognition*, 2013, pp. 971–978.
- [30] A. Voukoudimos, N. Doulamis, A. Doulamis, and E. Protopapadakis, “Deep learning for computer vision: A brief review,” *Computational intelligence and neuroscience*, vol. 2018, p. 7068349, 2018.
- [31] Y. Guo, Y. Liu, A. Oerlemans, S. Lao, S. Wu, and M. S. Lew, “Deep learning for visual understanding: A review,” *Neurocomputing*, vol. 187, pp. 27 – 48, 2016.
- [32] R. Girshick, J. Donahue, T. Darrell, and J. Malik, “Rich feature hierarchies for accurate object detection and semantic segmentation,” in *2014 IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 580–587.
- [33] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan, “Object detection with discriminatively trained part-based models,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 32, no. 9, pp. 1627–1645, 2010.
- [34] R. Girshick, “Fast R-CNN,” in *2015 IEEE International Conference on Computer Vision (ICCV)*, 2015, pp. 1440–1448.
- [35] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, “You only look once: Unified, real-time object detection,” in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 779–788.
- [36] J. Redmon and A. Farhadi, “Yolo9000: Better, faster, stronger,” in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 6517–6525.
- [37] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, “Ssd: Single shot multibox detector,” in *Computer Vision – ECCV 2016*, B. Leibe, J. Matas, N. Sebe, and M. Welling, Eds., 2016, pp. 21–37.
- [38] T. Lin, P. Goyal, R. Girshick, K. He, and P. Dollr, “Focal loss for dense object detection,” in *2017 IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 2999–3007.
- [39] L. Zhang, M. Song, Q. Zhao, X. Liu, J. Bu, and C. Chen, “Probabilistic graphlet transfer for photo cropping,” *IEEE Transactions on Image Processing*, vol. 22, no. 2, pp. 802–815, 2013.
- [40] L. Zhang, M. Song, Y. Yang, Q. Zhao, C. Zhao, and N. Sebe, “Weakly supervised photo cropping,” *IEEE Transactions on Multimedia*, vol. 16, no. 1, pp. 94–107, 2014.
- [41] L. Zhang, Y. Gao, R. Ji, Y. Xia, Q. Dai, and X. Li, “Actively learning human gaze shifting paths for semantics-aware photo cropping,” *IEEE Transactions on Image Processing*, vol. 23, no. 5, pp. 2235–2245, 2014.
- [42] A. Samii, R. Měch, and Z. Lin, “Data-driven automatic cropping using semantic composition search,” *Comput. Graph. Forum*, vol. 34, no. 1, pp. 141–151, 2015.
- [43] P. Wang, Z. Lin, and R. Mech, “Learning an aesthetic photo cropping cascade,” in *2015 IEEE Winter Conference on Applications of Computer Vision*, 2015, pp. 448–455.
- [44] Z. Zhou, S. He, J. Li, and J. Z. Wang, “Modeling perspective effects in photographic composition,” in *Proceedings of the 23rd ACM International Conference on Multimedia*, 2015, pp. 301–310.
- [45] Y.-L. Chen, J. Klopp, M. Sun, S.-Y. Chien, and K.-L. Ma, “Learning to compose with professional photographs on the web,” in *Proceedings of the 2017 ACM on Multimedia Conference*, 2017, pp. 37–45.
- [46] Y. Kao, R. He, and K. Huang, “Automatic image cropping with aesthetic map and gradient energy map,” in *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2017, pp. 1982–1986.
- [47] G. Guo, H. Wang, C. Shen, Y. Yan, and H. M. Liao, “Automatic image cropping for visual aesthetic enhancement using deep neural networks and cascaded regression,” *IEEE Transactions on Multimedia*, vol. 20, no. 8, pp. 2073–2085, 2018.
- [48] Z. Wei, J. Zhang, X. Shen, Z. Lin, R. Mech, M. Hoai, and D. Samarasinghe, “Good view hunting: Learning photo composition from dense view pairs,” in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
- [49] M. B. Islam, W. Lai-Kuan, and W. Chee-Onn, “A survey of aesthetics-driven image repositioning,” *Multimedia Tools Appl.*, vol. 76, no. 7, pp. 9517–9542, 2017.
- [50] O. Ronneberger, P. Fischer, and T. Brox, “U-net: Convolutional networks for biomedical image segmentation,” in *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015*, N. Navab, J. Hornegger, W. M. Wells, and A. F. Frangi, Eds., 2015, pp. 234–241.
- [51] S. Ren, K. He, R. Girshick, and J. Sun, “Faster R-CNN: Towards real-time object detection with region proposal networks,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 39, no. 6, pp. 1137–1149, 2017.
- [52] N. Doulamis, A. Doulamis, D. Kalogeras, and S. Kollias, “Low bit-rate coding of image sequences using adaptive regions of interest,” *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 8, no. 8, pp. 928–934, 1998.
- [53] G. R. Bradski, “Computer vision face tracking for use in a perceptual user interface,” 1998.
- [54] J. Dai, K. He, and J. Sun, “Instance-aware semantic segmentation via multi-task network cascades,” in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 3150–3158.
- [55] N. Murray, L. Marchesotti, and F. Perronnin, “Ava: A large-scale database for aesthetic visual analysis,” in *2012 IEEE Conference on Computer Vision and Pattern Recognition*, 2012, pp. 2408–2415.
- [56] Y. Chen, T. Huang, K. Chang, Y. Tsai, H. Chen, and B. Chen, “Quantitative analysis of automatic image cropping algorithms: A dataset and comparative study,” in *2017 IEEE Winter Conference on Applications of Computer Vision (WACV)*, March 2017, pp. 226–234.
- [57] X. Huang, C. Shen, X. Boix, and Q. Zhao, “Salicon: Reducing the semantic gap in saliency prediction by adapting deep neural networks,” in *2015 IEEE International Conference on Computer Vision (ICCV)*, 2015, pp. 262–270.
- [58] J. Park, J. Y. Lee, Y. W. Tai, and I. S. Kweon, “Modeling photo composition and its application to photo re-arrangement,” in *2012 19th IEEE International Conference on Image Processing*, 2012, pp. 2741–2744.
- [59] S. Ma, J. Liu, and C. W. Chen, “A-lamp: Adaptive layout-aware multi-patch deep convolutional neural network for photo aesthetic assessment,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 722–731.
- [60] L. Mai, H. Jin, and F. Liu, “Composition-preserving deep photo aesthetics assessment,” in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 497–506.