

Learning to Compose with Professional Photographs on the Web

Yi-Ling Chen¹ Jan Klopp² Min Sun³ Shao-Yi Chien² Kwan-Liu Ma¹

¹University of California, Davis ²National Taiwan University ³National Tsing Hua University

<https://github.com/yiling-chen/view-finding-network>

Abstract

Photo composition is an important factor affecting the aesthetics in photography. However, it is a highly challenging task to model the aesthetic properties of good compositions due to the lack of globally applicable rules to the wide variety of photographic styles. Inspired by the thinking process of photo taking, we formulate the photo composition problem as a view finding process which successively examines pairs of views and determines their aesthetic preferences. We further exploit the rich professional photographs on the web to mine unlimited high-quality ranking samples and demonstrate that an aesthetics-aware deep ranking network can be trained without explicitly modeling any photographic rules. The resulting model is simple and effective in terms of its architectural design and data sampling method. It is also generic since it naturally learns any photographic rules implicitly encoded in professional photographs. The experiments show that the proposed view finding network achieves state-of-the-art performance with sliding window search strategy on two image cropping datasets.

1. Introduction

¹“Aesthetics is a beauty that is found by a relationship between things, people and environment.”

Naoto Fukasawa

In the past decade, a considerable amount of research efforts have been devoted to computationally model aesthetics in photography. Most of these methods aim to either *assess* photo quality by resorting to well-established photographic rules [7, 16, 8, 23] or even to *manipulate* the image content to improve visual quality [1, 20, 37, 11]. However, to model photographic aesthetics remains a very challenging task due to the lack of a complete set of programmable rules to assess photo quality. In recent years, large-scale datasets with

¹© 2017. This is the authors’ version of this work. It is posted here for your personal use. Not for redistribution. The definitive version was published in *ACM MM '17*, <https://doi.org/10.1145/3123266.3123274>

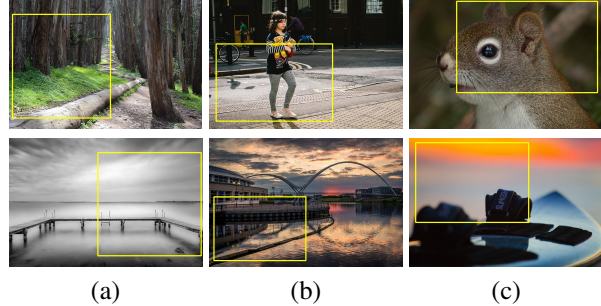


Figure 1. Professional photographs on the web are typically compliant with certain photographic rules. On the other hand, a crop of the image is highly likely to ruin the original composition, e.g., (a) symmetry, (b) rule of thirds, (c) object emphasis. By pairing a professional picture with a corresponding crop, it enables us to leverage human knowledge of photo composition under a learning-to-rank framework.

peer-rated aesthetic scores [26, 17] enable aesthetics modeling with learning based approaches [26, 21, 22, 24, 17]. However, the peer-rated aesthetic scores are subject to the bias between subjects, since comparing the aesthetic of arbitrary pairs of image is inevitably ambiguous sometimes. To mitigate the bias, one way is to get more signal than noise by enlarging the dataset. However, it is a daunting task to collect significantly more images with peer-rated aesthetic scores.

Rethink about the most basic behavior of photo taking: a photographer repeatedly *moves* the camera² and *judges* if the current view is more visually pleasing than the previous one until the desired view is obtained. The above observation reveals the essential property of photo composition – to *successively rank a pair of views with gradually altered contents*. Unlike most existing methods, which typically try to differentiate the aesthetics of *distinct* images, comparing the aesthetics relationship of visually similar views is relatively easy and less ambiguous. However, to collect a large amount of ranking samples by human raters for training effective models will inevitably face the aforementioned

²More specifically, the camera movement may include shift and zoom in/out to properly frame the desired view.

challenges – subjectiveness and scalability.

One key observation is that professional photographs are typically compliant with certain photographic rules (see some examples in Figure 1), which are inherently positive examples of good composition. On the other hand, a crop of professional photographs is highly likely to ruin the original composition. In other words, a pair of a professional photograph and its corresponding crop is highly likely to possess definite visual preference in terms of aesthetics. Thanks to the abundant professional photographs on the web, it is thus possible to harvest many unambiguous pairwise aesthetic ranking examples *for free*. Additionally, the model can naturally learn more photographic rules encoded in more training data without the necessity of explicitly modeling any new hand-crafted features.

Based on the above observations, we formulate the *learning-to-compose* problem as a pairwise view ranking process. We show that it can be effectively solved by a simple and powerful *view finding network* (VFN), which is trained to honor images of good composition and avoid those of bad composition. VFN is composed of a widely used object classification network [18] optionally augmented with a spatial pyramid pooling (SPP) layer [19, 12]. A costless data augmentation method is proposed to collect large-scale ranking samples from the unlimited high-quality images on the web. Without using any complex hand-crafted features, VFN learns the best photographic practices from examples by *relating* different views in terms of aesthetic ordering. To evaluate the capability of VFN for view finding, we evaluate its performance on two image cropping databases [5, 36]. We demonstrate that with simple sliding window search, VFN achieves state-of-the-art performance in cropping accuracy.

To summarize, our main contributions are as follows: We revisit the extensively studied problem of modeling photo aesthetics and composition and provide new key insights. The resulting technical solution is surprisingly simple yet effective. We show that a large number of automatically generated pairwise ranking constraints can be utilized to effectively train an aesthetics-aware deep ranking network. The proposed method significantly outperforms state-of-the-art methods as demonstrated by a quantitative evaluation on two public image cropping datasets.

2. Previous Work

Photo composition is an essential factor influencing the aesthetics in photography. A considerable amount of methods have been developed to assess photo quality [7, 16, 8, 23, 27]. Early works typically exploit “hand-crafted” features that mimic certain well-known photographic practices (*e.g.*, rule of thirds, visual balance etc.) and combine them with low-level statistics (*e.g.*, color histogram and wavelet analysis) to accomplish content-based aesthetic analysis.

More recently, generic image descriptors [25] and deep activation features [9] originally targeted at recognition are shown to be generic and outperform rule-based features in aesthetics prediction and style recognition [15]. With the advance of deep learning, recent works [14, 21, 22, 24, 17] train end-to-end models without explicitly modeling composition and achieve state-of-the-art performance in the recently released large scale Aesthetics Visual Analysis dataset (AVA) [26].

Compared to traditional photo quality assessment methods, which typically exploit photo composition as a high-level cue, some photo *recomposition* techniques attempt to actively enhance image composition by rearranging the visual elements [1], applying crop-and-retarget operations [20] or providing on-site aesthetic feedback [37] to improve the aesthetics score of the manipulated image.

Photo composition has also been extensively studied in *photo cropping* [38, 10, 4] and *view recommendation* [2, 6, 33] methods. Generally speaking, these methods aim at the same problem of finding the best view among a number of candidate views within a larger scene and mainly differ in how they differentiate a good view from the bad ones. Traditionally, *attention-based* approaches exploit visual saliency detection to identify a crop window covering the most visually significant objects [34, 32]. Some hybrid approaches employ a face detector [39] to locate the region-of-interest or fitting saliency maps to professional photographs [29]. On the other hand, *aesthetics-based* approaches aim to determine the most visually pleasing candidate window by resorting to photo quality classifiers [28], optimizing composition quality [10], or learning contextual composition rules [6]. In [36], a change-based method is proposed to model the variations before and after cropping so as to discard distracting content and improve the overall composition. In [5], the authors first investigate the use of learning-to-rank methods for image cropping. Unlike our method, they intentionally avoid professional pictures and relied on human raters to rank crops without obvious visual preference, resulting in a moderate-sized database.

To summarize, the main challenges faced by previous methods include 1) the limited applicability of rule-based features, and 2) the difficulty of obtaining composition information for training. The existing methods or databases build their training data by relying on a few experts [36, 10] or crowd-sourcing [26, 17, 5] to annotate and validate the training data, which makes it difficult to scale. In this work, we tackle these problems with a generic model powered by large-scale training data that is easy to obtain.

3. Approach

We model the photo composition or *view finding* process with View Finding Network (VFN). VFN, which is composed of a CNN augmented with a ranking layer, takes two

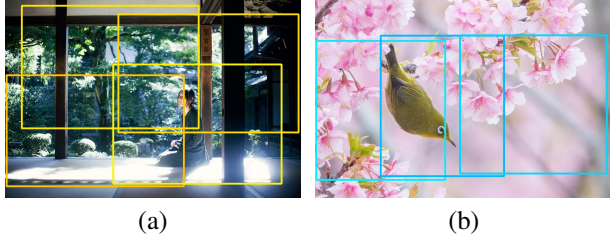


Figure 2. Examples of crop generation: (a) border crops, (b) square crops. Best viewed in color. Note that the rectangles indicate the crops corresponding to a single scale of crops.

views as input and predicts the more visually pleasing one in terms of composition. VFN learns its visual representations (*i.e.*, optimizes the weights of the CNN) by minimizing the disorder of image pairs with known aesthetic preference. Ideally, by examining extensive examples, VFN learns to compose as human professionals learned their skills.

3.1. Mining Pairwise Ranking Units

Every beginner to photography learns by seeing good examples, *i.e.*, professional photographs with perfect composition. One key observation is that the visual appearance of such golden examples typically achieves a state of *dangerous visual balance*. It implies that any deviations away from the current view will highly likely degrade the aesthetics – an inverse process of how the photographer obtained the optimal (current) view. It is thus possible to costlessly mine numerous image pairs with known relative aesthetic ranking. Figure 1 demonstrates several exemplary crops that possess less aesthetics due to violating the photographic heuristics encoded in the original image.

Based on the above observation, we empirically devise the following crop sampling strategies when given a source image I : 1) We always form pairs of the original image and a crop because the aesthetic relationship between two random crops is hard to define and thus requires human validation [5]. 2) To enrich the example set required when choosing the best view among different views, we include crops of varying scales and aspect ratios. 3) To best utilize the information in I , we aim to maximize the coverage of crops over I while minimizing the overlap between crops.

The resulting crop sampling procedure can be illustrated by Figure 2. Denote a crop of I as C and (x, y, w, h) indicates its *origin*, *width*, and *height*, respectively. For each image I , we generate a set of *border crops* and *square crops*. A border crop is created by first placing a uniformly resized window of I at the four corners. On the other hand, several square crops (we set the number to 3 in our experiments) are created along the long axis of I and evenly spaced. The parameters of C are then added with a small amount of random perturbation. Note that the above procedure is by no

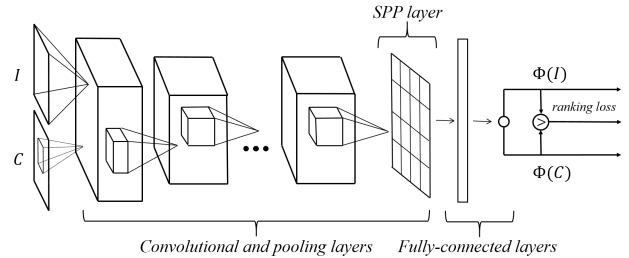


Figure 3. Architecture of View Finding Network.

means the optimal way to generate crops since it is impossible to test all possible configurations. Nevertheless, different sampling configurations consistently achieve better results than existing methods in our experiments. Please also refer to the supplementary material for more details of a series of experiments conducted to obtain the crop sampling configurations.

3.2. View Finding Network

Given an image I_j and its corresponding crops C_j^n , the objective of VFN is to learn a mapping function $\Phi(\cdot)$ that relates I_j and C_j^n according to their aesthetic relationship,

$$\Phi(I_j) > \Phi(C_j^n). \quad (1)$$

Notice that here we assume that I_j is always higher ranked than C_j^n in terms of aesthetics. We can thus define the following *hinge loss* for an image pair (I_j, C_j^n) :

$$l(I_j, C_j^n) = \max \{0, g + \Phi(C_j^n) - \Phi(I_j)\}, \quad (2)$$

where g is a gap parameter that regularizes the minimal margin between the ranking scores of I_j and C_j^n . We set $g = 1$ throughout all experiments. To learn $\Phi(\cdot)$, we minimize the total loss which sums up l over all training pairs.

Compared to many existing CNN models for aesthetics assessment [21, 22, 24, 17], the architecture of VFN is extremely simple, as illustrated in Figure 3. The convolutional layers of VFN are adopted from the popular AlexNet [18]. The output of the convolutional layers is then fed into two fully-connected layers followed by a *ranking layer*. The ranking layer is parameter-free and merely used to evaluate the hinge loss of an image pair. During training, the model updates its parameters such that $\Phi(\cdot)$ minimizes the total ranking loss in Equation (2). Once the network is trained, we discard the ranking layer and simply use $\Phi(\cdot)$ to map a given image I to an aesthetic score that differentiates I with other visually similar views.

On top of the last convolutional layer, we *optionally* append a *spatial-pyramid pooling* (SPP) layer [12]. SPP (also known as *spatial pyramid matching* or SPM) [19] is a widely used method to learn discriminative features by

dividing the image with a coarse-to-fine pyramid and aggregating the local features. It enhances the discrimination power of features by considering the global spatial relations. Notably, unlike [12, 24], we still use fix-sized input image in VFN (*i.e.*, the input image/patch are first resized to 227×227). We simply apply the SPP technique to accomplish data aggregation on the convolutional activation features.

Since photo composition is a property affected by both small (*e.g.*, a small object like a flagpole that may destroy the composition) and large structures (*e.g.*, visually significant objects in the scene) in the images, we thus choose the pooling regions of sizes 3×3 , 5×5 and 7×7 with the stride set to one pixel smaller than the pooling size (*e.g.*, 2 for 3×3 pooling regions). The multi-resolution pooling filters retain composition information at different scales. In addition, we empirically found that without SPP, the larger feature space causes the model more prone to overfitting. We apply both *max-pooling* and *average-pooling* in our experiments.

The pooled features are 12,544-dimensional and then fed into the first fully-connected layer. `fc1` is followed by a ReLU and has an output dimension of 1,000. We choose a relatively small feature dimension since the ranking problem is not as complex as object classification. Besides, as shown in [3], convolutional activation features can be compressed without considerable information loss while wide fully connected layers tend to overfit. `fc2` has only a single neuron and simply outputs the final ranking scores.

3.3. Training

To train our network, stochastic gradient descent algorithm with momentum is employed. We start from AlexNet [18] pre-trained on the ImageNet ILSVRC2012 dataset [31] and the fully connected layers are initialized randomly according to [13]. Momentum is set to 0.9 and the learning rate starts at 0.01 and is reduced to 0.002 after 10,000 iterations, with each mini-batch comprised of 100 image pairs. A total of 15,000 iterations is run for training and the validation set is evaluated every 1,000 iterations. The model with the smallest validation error is selected for testing. To combat overfitting, the training data is augmented by random horizontal flips as well as slight random perturbations on brightness and contrast. We implement and train our model with the TensorFlow³ framework.

4. Experimental Results

4.1. Training Data

To build the training data, we opt to download pictures shared by professional photographers on the Flickr website⁴. We exploited the Flickr API that returns the “interest-

ing photos of the day” and crawled 31,860 images⁵ during a period of 2,230 consecutive days. The initial data set is then manually curated to remove non-photographic images (*e.g.*, cartoons, paintings etc.) or images with post-processing affecting the composition (*e.g.*, collage, wide outer frame). The resulting image pool consists of 21,045 high-quality images and covers the most common categories in photography. We randomly selected 17,000 images for training and the rest images are used for validation. As described in Section 3.1, we generate 8 border crops and 6 square crops for each image corresponding to two scales $s = \{0.5, 0.6\}$. Each crop is paired with the corresponding original image and thus there are 294,630 image pairs in total. The image pair collection is then used to train the VFN. Note that the above procedure is inexpensive and it is very easy to expand the dataset.

4.2. Performance Evaluations

To validate the effectiveness of our model for view finding, we evaluate its *cropping accuracy* on two public image cropping databases, including Flickr Cropping Database (FCDB) [5] and Image Cropping Database (ICDB) [36], and compare against several baselines.

4.2.1 Evaluation Metrics

We adopt the same evaluation metrics as [36, 5], *i.e.*, *average intersection-over-union (IoU)* and *average boundary displacement* to measure the *cropping accuracy* of image croppers. IoU is computed by $area(\hat{C}_i \cap C_i) / area(\hat{C}_i \cup C_i)$, where \hat{C}_i and C_i denote the ground-truth crop window and the crop window determined by the baseline algorithms for the i -th test image, respectively. Boundary displacement is given by $\sum_{j=1}^4 \|\hat{B}_i^j - B_i^j\| / 4$, where \hat{B}_i^j and B_i^j denote the four corresponding edges between \hat{C}_i and C_i . Additionally, we report α -*recall*, which is the fraction of best crops that have an overlapping ratio greater than α with the ground truth. In all of our experiments, we set α to 0.75.

For the simplicity and fairness of comparison, we follow the sliding window strategy of [5] to evaluate the baselines and VFN. Similarly, we set the size of search window to each scale among $[0.5, 0.6, \dots, 0.9]$ of the test images and slide the search window over a 5×5 uniform grid. The ground truth is also included as a candidate. The optimal crops determined by individual methods are compared to the ground truth to evaluate their performance.

4.2.2 Baseline Algorithms

Following [36], we compare with two main categories of traditional image cropping methods, *i.e.*, attention-based

³<https://www.tensorflow.org/>

⁴<https://www.flickr.com/>

⁵We only kept those images with Creative Common license and more than 100 “favorite” counts.

Method	IoU	Disp.	α -recall
eDN [35]	0.4929	0.1356	12.68
AlexNet_finetune	0.5543	0.1209	16.092
MNA-CNN [24]	0.5042	0.1361	0.0747
RankSVM+AVA [5]	0.5270	0.1277	12.6437
RankSVM+FCDB [5]	0.602	0.1057	18.1034
AesRankNet [17]	0.4843	0.1401	0.0804
VFN	0.6842	0.0843	35.0575
VFN+AVA (SPP-Max)	0.544	0.124	12.93
VFN (SPP-Avg)	0.6783	0.0859	35.0575
VFN (SPP-Max)	0.6744	0.0872	33.9080

Table 1. Performance comparison on FCDB [5]. The best results are highlighted in bold.

and aesthetics-based approaches. Additionally, we compare with several ranking-based image croppers [5].

- Attention-based: For attention-based methods, we choose the best performing method (eDN) reported in [5], which adopts the saliency detection method described in [35] and searches for the best crop window that maximizes the difference of saliency score between the crop and the outer region of the image.
- Aesthetics-based: We choose to fine-tune AlexNet [18] for binary aesthetics classification with the AVA dataset [26] as the baseline of this category and follows the configuration suggested by [22, 17]. We simply utilize the softmax confidence score to choose the best view. The methods of [36, 24] also fall into this category. We compare with these methods by the accuracy reported in the original paper [36] or use the pre-trained model to evaluate its performance in both datasets [24].
- Ranking-based: We adopt two variants of RankSVM-based image croppers using deep activation features [9] and trained on the AVA and FCDB datasets [5], which differ in their data characteristics. AVA characterizes the aesthetics preference between *distinct* images while FCDB provides the ranking order between crop pairs in the *same* images. Additionally, we compare with the recent work of aesthetics ranking network [17]. We use the pre-trained model released by the authors and utilize the ranking scores of the sliding windows to determine the best crop.

4.2.3 Performance Evaluation

We evaluate cropping accuracy of VFN and several baselines on FCDB and ICDB, which differ in data characteristics and annotation procedure. The test set of FCDB contains 348 images. Each image was labeled by a photogra-

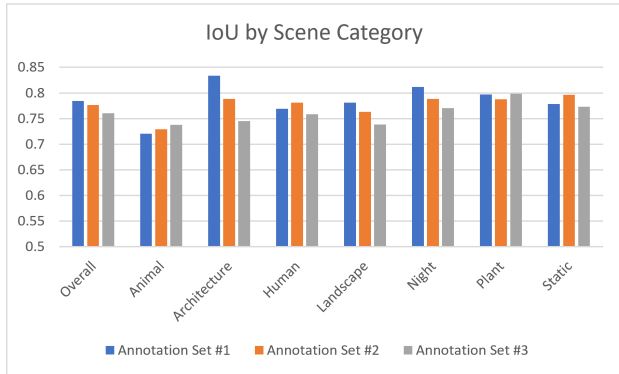


Figure 4. Performance of VFN (SPP-Max) on ICDB by category.

phy hobbyist and then validated by 7 workers on Amazon Mechanical Turk. On the other hand, ICDB includes 950 images, each annotated by 3 experts. The images of ICDB are typically of *iconic* views and thus more object-centric. Compared to ICDB, FCDB is considered to be more challenging for image cropping methods because the annotations reflect the tastes of various photographers and the images contain more contextual information.

Table 1 and 2 summarize the benchmark results. Generally, the performance of each category is consistent with [5]. The attention-based method (eDN) performs poorly due to the lack of aesthetic consideration and aesthetics-based methods based on a photo quality classifier (AlexNet_finetune) achieves only moderate performance. Surprisingly, the aesthetics ranking network [17] and MNA-CNN [24] methods also do not perform well in the benchmark. This is most possibly because these networks are trained to predict the aesthetic rating of distinct images, which does not reflect the relations between different views with large overlaps. We validate this by training VFN with the traditional dataset [26] and will discuss the results soon. Additionally, some image attributes (*e.g.*, color) assessed by the model may not be very discriminating for similar views.

All variants of VFN trained by our data sampling technique significantly outperform the other baselines. The best performing baseline method is the change-based algorithm [36], which achieves very good results in their dataset (ICDB). Notably, the model of [36] is trained on the first annotation set and evaluated on the same images of all annotation sets. On the other hand, the images of ICDB are totally unseen to our models. In addition, a crop selection procedure which selects an initial set of good candidate windows is incorporated in [36], while VFN is evaluated by a fixed set of sliding windows.

We conduct a more fine-grained performance analysis of VFN by the scene category annotations provided by ICDB, which further divide the dataset into seven categories: *animal*, *architecture*, *human*, *landscape*, *night*, *plant* and *static*.

Method	Annotation Set #1			Annotation Set #2			Annotation Set #3		
	IoU	Disp.	α -recall	IoU	Disp.	α -recall	IoU	Disp.	α -recall
eDN [35]	0.5535	0.1273	27.3684	0.5128	0.1419	20.1053	0.5257	0.1358	22.4211
AlexNet_finetune	0.5687	0.1246	23.0526	0.5536	0.1296	22.7368	0.5544	0.1288	20.6316
MNA-CNN [24]	0.4693	0.1555	0.0716	0.4553	0.1615	0.0642	0.4610	0.1590	0.0684
RankSVM+AVA [5]	0.5801	0.1174	18.7368	0.5678	0.1225	18.6316	0.5665	0.1226	18.9474
RankSVM+FCDB [5]	0.6683	0.0907	33.4737	0.6618	0.0932	32.1053	0.6483	0.0973	31.2632
AesRankNet [17]	0.4484	0.1631	0.0863	0.4372	0.168	0.0747	0.4408	0.1655	0.0863
LearnChange [36]	0.7487	0.0667	–	0.7288	0.072	–	0.7322	0.0719	–
VFN	0.7720	0.0623	58.8421	0.7638	0.0654	56.4211	0.7487	0.0692	53.7895
VFN+AVA (SPP-Max)	0.5273	0.1387	18.21	0.5268	0.14	19.0526	0.5261	0.1389	18
VFN (SPP-Avg)	0.7837	0.0588	61.5789	0.7729	0.0627	58.1053	0.7514	0.0681	54.1053
VFN (SPP-Max)	0.7847	0.0581	59.7895	0.7763	0.0614	58.1053	0.7602	0.0653	54.8421

Table 2. Performance evaluation on ICDB [36]. The best results are highlighted in bold.

Figure 4 illustrates the average IoU scores of VFN (SPP-Max) over various annotation sets and scene categories. We consider VFN as a generic aesthetics model since it shows no *bias* to specific annotation sets or categories. It performs generally well across subjects despite image cropping’s subjective nature. Since no category-specific features are exploited, VFN generalizes across categories as well. Nevertheless, we still can see some insufficiency of VFN, *e.g.*, consistently lower accuracy in *animal* and higher performance variance in *architecture*. Considering the data-driven nature of VFN, this phenomenon can be most possibly accounted for insufficient or unbalanced distribution of training images among various categories.

Some more interesting observations can be made from the benchmarking results, as discussed below.

View finding is intrinsically a problem of ranking pairwise views in the same context. Intuitively, image rankers trained on aesthetics relations derived from distinct images, such as RankSVM+AVA and [17], do not necessarily perform well in ranking visually similar views. To further validate such an assumption, we additionally trained a VFN with ranking units purely sampled from AVA. To mitigate the ambiguity of ranking relationship between images, we choose the 30,000 highest and lowest ranked images from AVA and randomly select a pair of images from each pool to train VFN. The resulting model (VFN+AVA) can be regarded as the counterpart of RankSVM+AVA. As shown in Table 1 and 2, the performance of VFN+AVA drastically degrades compared to other variants of VFN. It also confirms that our data sampling technique contributes to the most significant leap in performance.

Performance gain due to top level pooling is dependent on the characteristics of test data. VFN achieves the best results in ICDB and FCDB with and without SPP, respectively. Recall that the images in ICDB are largely object-centric with iconic views. Since pooled features are typically equipped with certain *invariance* (*e.g.*, translation),

it is thus beneficial to discriminate scenes with significant objects. On the other hand, since the images in FCDB possess richer contextual information, the greater feature space without pooling is thus more capable of capturing more subtle variations in photo composition, resulting in higher performance.

4.3. Applications

Automatic image cropping The ability of VFN makes it very suitable to facilitate the process of identifying unattractive regions in an image to be cut away so as to improve its visual quality. As demonstrated by the quantitative evaluation in Section 4.2.3, VFN achieves state-of-the-art performance in two image cropping datasets. Figure 7 illustrates several examples of applying VFN to crop images from FCDB and compares the results with several baselines. One can see that VFN successfully selects more visually pleasing crop windows compared to other baseline algorithms. Some of the results by VFN are arguably no worse than the ground truth (*e.g.*, the 2nd and 3rd row in Figure 7). Currently, only sliding windows with the same aspect ratio as the original image are used for evaluation, which limits VFN’s ability to identify other possible good compositions, as the ground truth shown in the 1st row of Figure 7. Nevertheless, VFN selects a preferable view with rule-of-thirds composition in this example when compared with other baselines. However, a crop selection procedure that adaptively determines the parameters of crop windows is still desirable for VFN to maximize its performance.

View recommendation VFN is aesthetics-aware and very sensitive to the variation of image composition. Figure 5 demonstrates an example of applying VFN to an image and its artificially “corrupted” version. We generate a heatmap by evaluating sliding windows and smoothing the ranking scores corresponding to the raw pixels. As one can see, the altered image composition causes VFN to shift its

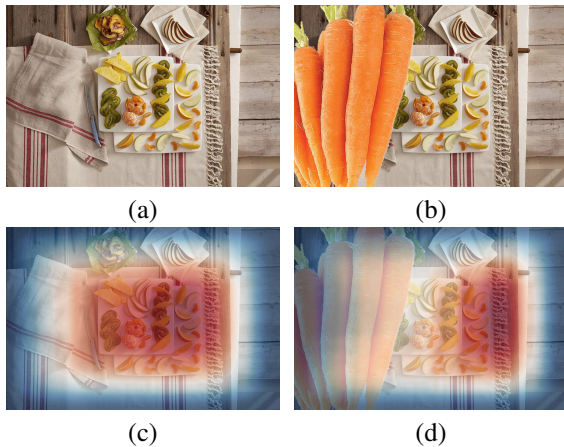


Figure 5. VFN is aesthetics-aware and capable of differentiating good/bad views in terms of photo composition. Given an source image (a) and a corrupted image (b), VFN produces higher response to the visually pleasing regions, as demonstrated in the corresponding heatmaps (c)(d).

attention to the untouched region. Due to its aesthetics-awareness, VFN is very suitable to be applied for view suggestion in panoramic scenes or even 360 video, as demonstrated in Figure 6. In this example, VFN identifies a visually attractive view while ignoring large unimportant areas in the scene. Unlike [2], which requires a *query* or *template* image to locate similar views in the panoramic image, our model is able to suggest a good view based on a much larger *database* (*i.e.*, the training images).

4.4. Discussion

Unlike traditional approaches, VFN learns to compose without explicitly modeling photo composition. In a sense, it is accomplished by *avoiding* the views violating photographic rules encoded in professional photographs. Take the fifth row of Figure 7 as an example, the baseline methods inappropriately cut through visually significant subjects. Previous methods explicitly deal with such situation by modeling *cut-through* feature [36] or *border simplicity* [21]. However, VFN naturally ignores these views because the proposed crop sampling method covers such cases and they are always penalized in our ranking model. Due to the principle of pairing a good source image and a bad crop, there is thus the concern that the learned model is biased to favor larger views with more image content. However, according to the benchmark, such tendency is not observed and the ranking model works well regardless of the scales.

Currently, VFN does not take full advantage of the SPP technique. Its performance can potentially be further improved if the constraint of fixed-size input can be removed and the input images do not need to undergo undesired transformations (*e.g.*, cropping or scaling) that typically

cause damages to image composition [24].

We have shown that VFN is generic across categories in Section 4.2.3. It is considered that the generalization capability of VFN partially benefits from the object classification capability of the pre-trained AlexNet [18], which provides rich information to learn category-specific features that discriminate aesthetic relationships.

Limitations and Future Works The main limitation of VFN comes from its data sampling methodology, which only samples a sparse set of possible pairs of views. The success of VFN can be accounted for that the aesthetic relations between the sampled pairs (*i.e.*, a source image and a random crop) are definite. However, it remains a challenging task for VFN to rank similar views whose aesthetic relation is ambiguous (*e.g.*, two random crops or two nearly identical views). Empirically, we found that evaluating a finer set of sliding windows with VFN causes the performance to degrade instead, which is possibly caused by the confusion between very similar views. To maximize the performance of VFN, it is considered to incorporate a view selection procedure, which can effectively eliminate most unnecessary candidates and produces a sparse set of good candidates. VFN currently needs to evaluate a number of proposal windows to accomplish view finding. For future work, we plan to incorporate techniques like Faster R-CNN [30] to improve its time efficiency.

5. Conclusion

In this work, we considered one of the most important problems in computational photography – automatically finding a good photo composition. Inspired by the thinking process of photo taking, a deep ranking network is proposed to learn the best photographic practices by leveraging human knowledge from the abundant professional photographs on the web. We develop a costless and effective method to sample high-quality ranking samples in an unsupervised manner. Without any hand-crafted features, the proposed method is simple and generic. The resulting aesthetics-aware model is evaluated on two image cropping datasets and achieves state-of-the-art performance in terms of cropping accuracy.

References

- [1] S. Bhattacharya, R. Sukthankar, and M. Shah. A framework for photo-quality assessment and enhancement based on visual aesthetics. In *ACM Multimedia*, pages 271–280, 2010. 1, 2
- [2] Y.-Y. Chang and H.-T. Chen. Finding good composition in panoramic scenes. In *ICCV*, pages 2225–2231, 2009. 2, 7
- [3] K. Chatfield, K. Simonyan, A. Vedaldi, and A. Zisserman. Return of the Devil in the Details: Delving Deep into Convolutional Nets. In *BMVC*, pages 1–11, 2014. 4



Figure 6. An example of applying VFN to a panorama image. The yellow rectangles in the left column images indicate the crop with the maximum score among 2112 uniformly sampled candidate crops of different sizes and aspect ratios. The resulting crop is shown in right column. Best viewed in color.

- [4] J. Chen, G. Bai, S. Liang, and Z. Li. Automatic image cropping: A computational complexity study. In *CVPR*, 2016. [2](#)
- [5] Y.-L. Chen, T.-W. Huang, K.-H. Chang, Y.-C. Tsai, H.-T. Chen, and B.-Y. Chen. Quantitative analysis of automatic image cropping algorithms: A dataset and comparative study. In *WACV*, 2017. [2](#), [3](#), [4](#), [5](#), [6](#), [9](#)
- [6] B. Cheng, B. Ni, S. Yan, and Q. Tian. Learning to photograph. In *ACM Multimedia*, pages 291–300, 2010. [2](#)
- [7] R. Datta, D. Joshi, J. Li, and J. Z. Wang. Studying aesthetics in photographic images using a computational approach. In *ECCV*, pages 288–301, 2006. [1](#), [2](#)
- [8] S. Dhar, V. Ordonez, and T. L. Berg. High level describable attributes for predicting aesthetics and interestingness. In *CVPR*, pages 1657–1664, 2011. [1](#), [2](#)
- [9] J. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng, and T. Darrell. DeCAF: A deep convolutional activation feature for generic visual recognition. In *ICML*, pages 647–655, 2013. [2](#), [5](#)
- [10] C. Fang, Z. Lin, R. Mech, and X. Shen. Automatic image cropping using visual composition, boundary simplicity and content preservation models. In *ACM Multimedia*, pages 1105–1108, 2014. [2](#)
- [11] Y. W. Guo, M. Liu, T. T. Gu, and W. P. Wang. Improving photo composition elegantly: Considering image similarity during composition optimization. *Comput. Graph. Forum*,

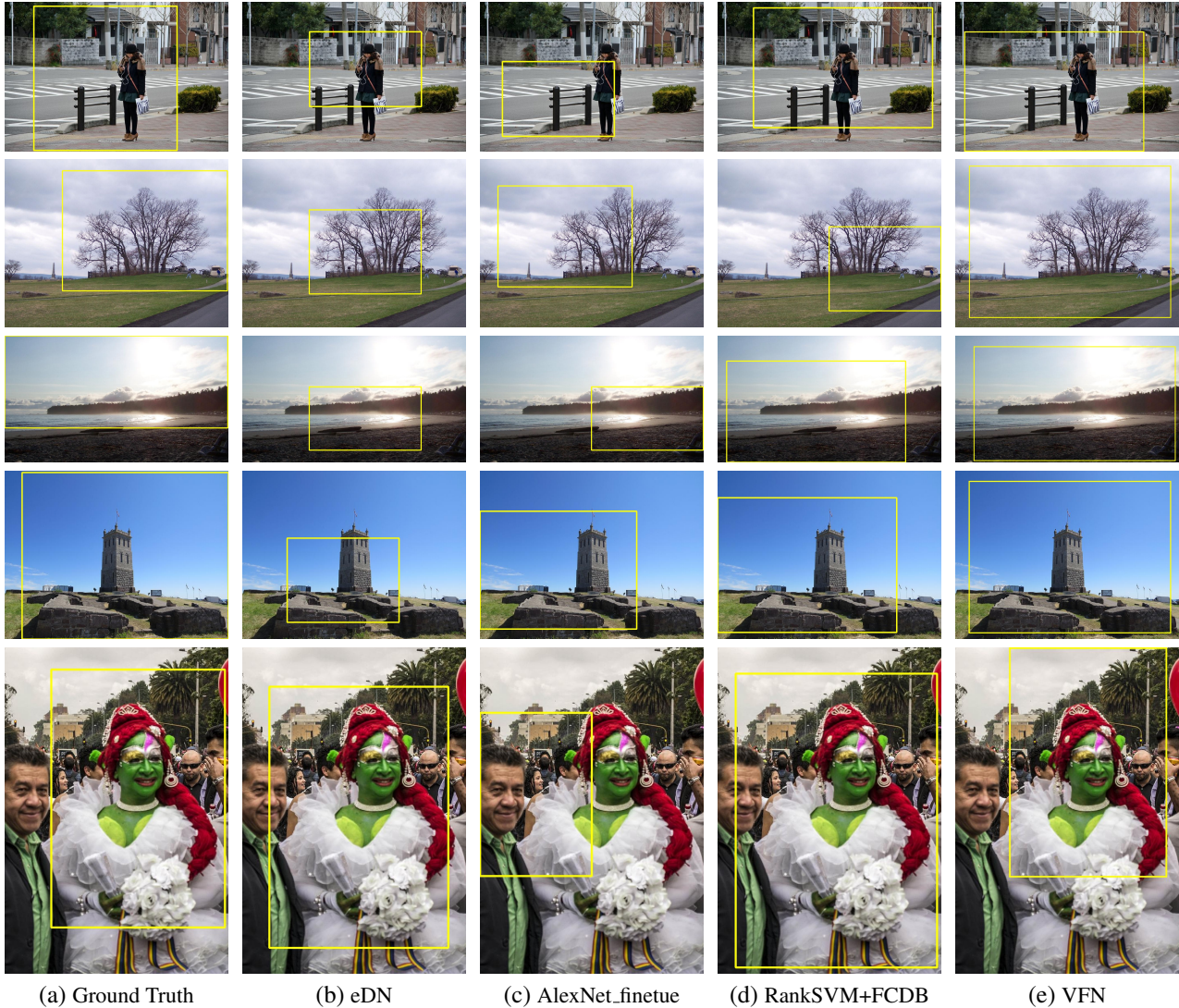


Figure 7. Image cropping examples from FCDB [5]. The best crops determined by various methods are drawn as yellow rectangles. Best viewed in color.

31(7):2193–2202, Sept. 2012. 1

[12] K. He, X. Zhang, S. Ren, and J. Sun. Spatial pyramid pooling in deep convolutional networks for visual recognition. In *ECCV*, pages 346–361, 2014. 2, 3, 4

[13] K. He, X. Zhang, S. Ren, and J. Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *ICCV*, pages 1026–1034, 2015. 4

[14] L. Kang, P. Ye, Y. Li, and D. Doermann. Convolutional neural networks for no-reference image quality assessment. In *CVPR*, pages 1733–1740, 2014. 2

[15] S. Karayev, M. Trentacoste, H. Han, A. Agarwala, T. Darrell, A. Hertzmann, and H. Winnemoeller. Recognizing image style. In *BMVC*, pages 1–20, 2014. 2

[16] Y. Ke, X. Tang, and F. Jing. The design of high-level features for photo quality assessment. In *CVPR*, pages 419–426, 2006. 1, 2

[17] S. Kong, X. Shen, Z. Lin, R. Mech, and C. Fowlkes. Photo aesthetics ranking network with attributes and content adaptation. In *ECCV*, 2016. 1, 2, 3, 5, 6

[18] A. Krizhevsky, I. Sutskever, and G. E. Hinton. ImageNet classification with deep convolutional neural networks. In *NIPS*, pages 1106–1114, 2012. 2, 3, 4, 5, 7

[19] S. Lazebnik, C. Schmid, and J. Ponce. Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories. In *CVPR*, pages 2169–2178, 2006. 2, 3

[20] L. Liu, R. Chen, L. Wolf, and D. Cohen-Or. Optimizing photo composition. *Computer Graphics Forum (Proc. of Eurographics '10)*, 29(2):469–478, 2010. 1, 2

[21] X. Lu, Z. Lin, H. Jin, J. Yang, and J. Z. Wang. RAPID: Rating pictorial aesthetics using deep learning. In *ACM Multimedia*, pages 457–466, 2014. 1, 2, 3, 7

[22] X. Lu, Z. Lin, X. Shen, R. Mech, and J. Z. Wang. Deep

- multi-patch aggregation network for image style, aesthetics, and quality estimation. In *ICCV*, pages 990–998, 2015. 1, 2, 3, 5
- [23] W. Luo, X. Wang, and X. Tang. Content-based photo quality assessment. In *ICCV*, pages 2206–2213, 2011. 1, 2
- [24] L. Mai, H. Jin, and F. Liu. Composition-preserving deep photo aesthetics assessment. In *CVPR*, 2016. 1, 2, 3, 4, 5, 6, 7
- [25] L. Marchesotti, F. Perronnin, D. Larlus, and G. Csurka. Assessing the aesthetic quality of photographs using generic image descriptors. In *ICCV*, pages 1784–1791, 2011. 2
- [26] N. Murray, L. Marchesotti, and F. Perronnin. AVA: A large-scale database for aesthetic visual analysis. In *CVPR*, pages 2408–2415, 2012. 1, 2, 5
- [27] M. Nishiyama, T. Okabe, I. Sato, and Y. Sato. Aesthetic quality classification of photographs based on color harmony. In *CVPR*, pages 33–40, 2011. 2
- [28] M. Nishiyama, T. Okabe, Y. Sato, and I. Sato. Sensation-based photo cropping. In *ACM Multimedia*, pages 669–672, 2009. 2
- [29] J. Park, J.-Y. Lee, Y.-W. Tai, and I. S. Kweon. Modeling photo composition and its application to photo rearrangement. In *ICIP*, 2012. 2
- [30] S. Ren, K. He, R. Girshick, and J. Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In *NIPS*, 2015. 7
- [31] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. 4
- [32] F. Stentiford. Attention based auto image cropping. In *ICVS Workshop on Computation Attention and Applications*, 2007. 2
- [33] H.-H. Su, T.-W. Chen, C.-C. Kao, W. H. Hsu, and S.-Y. Chien. Preference-aware view recommendation system for scenic photos based on bag-of-aesthetics-preserving features. *IEEE Transactions on Multimedia*, 14(3-2):833–843, 2012. 2
- [34] B. Suh, H. Ling, B. B. Bederson, and D. W. Jacobs. Automatic thumbnail cropping and its effectiveness. In *ACM UIST*, pages 95–104, 2003. 2
- [35] E. Vig, M. Dorr, and D. Cox. Large-scale optimization of hierarchical features for saliency prediction in natural images. In *CVPR*, pages 2798–2805, 2014. 5, 6
- [36] J. Yan, S. Lin, S. B. Kang, and X. Tang. Learning the change for automatic image cropping. In *CVPR*, pages 971–978, 2013. 2, 4, 5, 6, 7
- [37] L. Yao, P. Suryanarayan, M. Qiao, J. Z. Wang, and J. Li. OSCAR: On-site composition and aesthetics feedback through exemplars for photographers. *International Journal of Computer Vision*, 96(3):353–383, 2012. 1, 2
- [38] L. Zhang, M. Song, Q. Zhao, X. Liu, J. Bu, and C. Chen. Probabilistic graphlet transfer for photo cropping. *IEEE Transactions on Image Processing*, 22(2):802–815, 2013. 2
- [39] M. Zhang, L. Zhang, Y. Sun, L. Feng, and W. Ma. Auto cropping for digital photographs. In *ICME*, pages 2218–2221, 2005. 2