

Intelligent Portrait Composition Assistance

— Integrating Deep-learned Models and Photography Idea Retrieval

Farshid Farhat, Mohammad Mahdi Kamani, Sahil Mishra, James Z. Wang*

The Pennsylvania State University, University Park, Pennsylvania, USA

ABSTRACT

Retrieving photography ideas corresponding to a given location facilitates the usage of smart cameras, where there is a high interest among amateurs and enthusiasts to take astonishing photos at anytime and in any location. Existing research captures some aesthetic techniques and retrieves useful feedbacks based on one technique. However, they are restricted to a particular technique and the retrieved results have room to improve as they are confined to the quality of the query. There is a lack of a holistic framework to capture important aspects of a given scene and help a novice photographer by informative feedback to take a better shot in his/her photography adventure. This work proposes an intelligent framework of portrait composition using our deep-learned models and image retrieval methods. A highly-rated web-crawled portrait dataset is exploited for retrieval purposes. Our framework detects and extracts ingredients of a given scene representing as a correlated semantic model. It then matches extracted semantics with the dataset of aesthetically composed photos to investigate a ranked list of *photography ideas*, and gradually optimizes the human pose and other artistic aspects of the composed scene supposed to be captured. The conducted user study demonstrates that our approach is more helpful than other feedback retrieval systems.

KEYWORDS

Photographic Composition; Image Aesthetics; Smart Camera; Portrait Photography; Deep Learning; Image Retrieval.

1 INTRODUCTION

Art still has many ambiguous aspects out of the known sciences, and the beauty of the art is coming from the virgin novelty by artists. It is still daunting for a machine to compose an impressive original song, painting or script. However, high-resolution photography has been made ubiquitous by recent technologies, such as high-quality smart camera phones. Also, the aesthetics of photography is known as a collection of rules in artistic literature [1–3] such as balance, geometry, symmetry, the rule of thirds, framing, and etc. Digital photography is of great interest among most people using

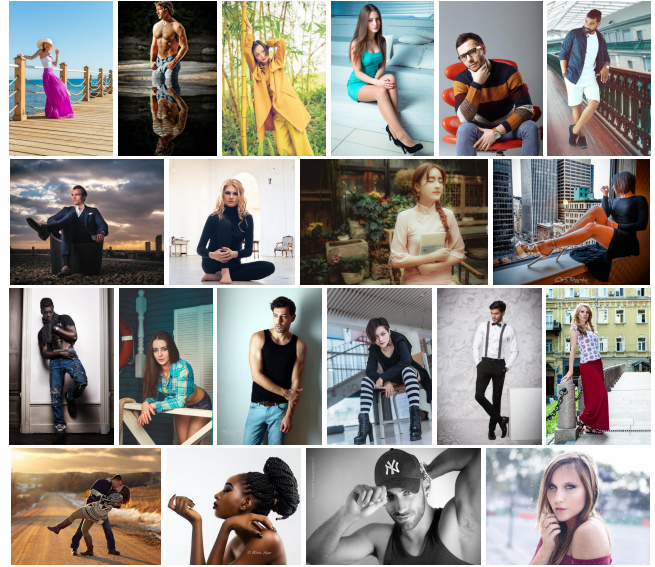


Figure 1: Portrait images given various scenes with several pose ideas for a better composed photo. Images from the 500px website are selected by our framework.

social networking and photo sharing websites such as Facebook, Google Photos, Twitter, Instagram, etc., but getting a striking photo involves experience and skills, and is often not easy.

While there are many styles for photography [2, 4] around the world, selecting proper *photography ideas* for a given scene remains a challenging problem, and yet to be adequately investigated. The major problem with taking a good portrait photo in a given location is the lack of a local photographer guide conveying us to capture a good portrait pose. In fact, professional photographers usually have expertise and creativity in making good positions intuitively [5–7]. Through reading books about photography, one can get familiar with some common composition rules such as balancing, framing, the rule of thirds, etc., but it can still be difficult to select and apply techniques for making genuine photos, in a way similar to the gap between reading and writing a novel.

Some basic rules of photography composition inspired from art books [1, 2, 6] have been used by multimedia and vision researchers as aesthetics features for assessment and evaluation of photos [8–11]. Other approaches manipulate the taken photo in an online system [12, 13] for auto-composition or re-composition. The techniques include smart cropping [14–20], warping [21, 22], patch re-arrangement [23–25], cutting and pasting [13, 17], and seam carving [26, 27], but they can barely help an amateur photographer capture a brilliant photo. More specifically in portrait photography, there are rule-based assessment models [28, 29] using known photography basics to evaluate portraits, and facial assessment

*F. Farhat and S. Mishra are with the School of Electrical Engineering and Computer Science. M. M. Kamani and J. Z. Wang are with the College of Information Sciences and Technology. Emails: {fuf111, mqk5591, szm5707, jwang}@psu.edu .

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

ThematicWorkshops'17, October 23–27, 2017, Mountain View, CA, USA

© 2017 ACM. ISBN 978-1-4503-5416-5/17/10...\$15.00

DOI: <https://doi.org/10.1145/3126686.3126710>

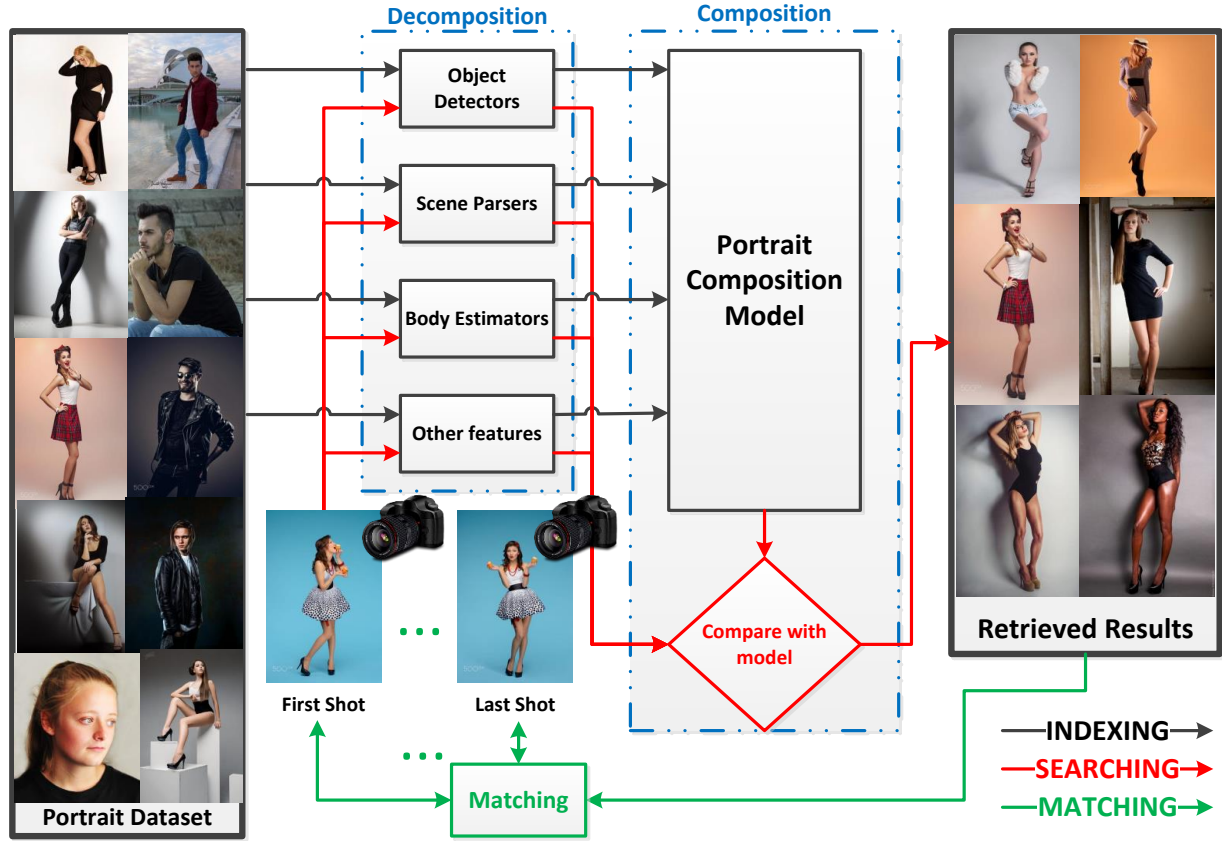


Figure 2: The flowchart of our portrait composition assistance: Black flows show the indexing process, red flows show the searching process, and green flows show the matching process. Decomposition step extracts the image semantics and features, and composition step searches for well-posed images in the dataset based on the semantics and other features.

models [30–34] exploiting facial features including smile, age, gender, and etc. Perhaps on-site photographic feedbacks [35–37] can help amateur photographers better by retrieving similar-aesthetic images as a qualitative composition feedback, but their retrieval system is limited to perspectiveness [38] or triangle technique [39].

In this paper, we focus on a framework that helps people make a better shot for their portrait photo with regard to their current location. Given a prior shot from the photographer or the camera viewfinder, our portrait composition assistance outputs some highly-rated prior-composed photos as an assessed feedback. Figure 1 shows some highly-rated portrait poses, many taken by professionals, collected from the 500px website and selected by our framework. These 20 portraits are captured in various locations and scenes. Each of them has its own photography idea(s) such as a woman with hat (1st image) has made a apropos pose at the heart of the leading lines (fence), or a woman sitting with crossed ankles bended legs (4th image) where this pose creates a nice S-shape. These techniques are believed to make portrait photography more appealing. Specifically, we address aesthetic retrieval and evaluation of the human poses in portrait photography, and try to improve the quality of the next shot by providing meaningful and constructive feedback to an amateur photographer. Figure 2 shows

the flowchart of our approach to assist an amateur photographer in getting a better shot. Based on the first shot as a query, some highly-rated well-posed results are retrieved from our designated dataset using portrait composition model containing the dataset features, and the results are illustrated to the photographer to help compose a better shot, and the last shot is captured when the current pose is matched with one of the results closely. The details of the flowchart has been explained in later sections. The main **contributions** are as follows:

- A holistic framework to intelligently assist amateur photographers to compose a better portrait using our proposed deep-learned models and image retrieval system.
- Various improved deep-learned detectors including object detector, scene parser and pose estimator to extract semantics are integrated.
- Scene construction by composing the semantics and retrieving the desired images from the dataset. We match the scene ingredients with semantic retrievals to optimize the final pose.
- Creating a large dataset containing over 320,000 highly-rated aesthetically composed portraits, with several categories and various scenes.

2 RELATED WORK

General Visual Aesthetics Assessment: While there are many books in art or photography to guide people mastering the challenges of taking professional photographs, the conducted research in technical fields mostly focus on the evaluation and manipulation of the images, after the photo is taken. Basic image aesthetics and composition rules in art [1–3] as visual semantic features have first been studied computationally by Datta *et al.* [8]. Luo *et al.* [9] and Wong and Low [10] attempt to leverage a saliency map method to focus on the features of the salient part as the more appealing part of the image. Marchesotti *et al.* [11] show that generic image descriptors can be very useful to assess image aesthetics, and build a generic dataset for aesthetics assessment called as Aesthetic Visual Analysis (AVA) [40].

Image Re-Composition: Auto-composition or re-composition systems [12, 13] can passively change the taken photo for a better composition. Cropping techniques [14–16] separate the region of interest (ROI) by the help of saliency map or eye fixation, basic aesthetic rules [17], or visual aesthetics features in the salient region [18–20]. As another type of re-composition, Liu *et al.* [21] use warping, *i.e.*, representing the image as a triangular or quad mesh, to map the image into another mesh while keeping the semantics and perspectiveness. Also, R2P [22] detects the foreground part in reference and input image, and tries to re-target the salient part of the image to the best fitted position using a graph-based algorithm. Furthermore, patch re-arrangement techniques patch two ROIs in an image together. Pure patch rearrangement [23–25] detects the group of pixels on the borders of the patch and matches this group to the other vertical or horizontal group of pixels near the patched area. Cut and paste methods [13, 17] remove the salient part, and re-paint the foreground with respect to salient part and borders, and then paste it to the desired position in the image. Seam carving [26, 27] replaces useless seams.

Portrait Aesthetics Assessment: While there have been works in image aesthetics assessment, few considered portrait photography in depth. Even in this domain, they haven’t tried to explore a novel method to solve the problem in photographic portraiture, rather they just combine and use well-known features or modified trivial ones to apply in the facial domain. We can categorize them into two main groups: rule-based evaluation models [28, 29] exploit known photography rules to assess portraits, and facial evaluation models [30–34] use visual facial features like smiling, age, gender, etc. Khan and Vogel [28] show a small set of proper spatial features can perform better than a large set of aesthetics features. Also, their feature importance analysis interestingly shows their spatial features mostly affect the system accuracy. Males *et al.* [29] explore aesthetic quality of headshots by means of some famous photography rules and low-level facial features. Xue *et al.* [30] study the design inferring portrait aesthetics with more appealing facial features like smiling, orientation, to name but a few. Similarly, Harel *et al.* [41] exploit traditional features like hue, saturation, brightness, contrast, simplicity, sharpness, and the rule of thirds. They also extract saliency map by graph-based visual saliency, and calculate standard deviation and main subject coincidence of the saliency map. The other facial evaluation models [31–33] use well-known low-level aesthetics features such as colorfulness, sharpness and

contrast, as well as high-level facial features such as gender, age, and smile. Their idea is based on exploiting these features for all segmented parts of the face including hair, face, eyes, and mouth. Redi *et al.* [34] interestingly show that the beauty of the portrait is related to the amount of art used in it not the subject beauty, age, race, or gender. While using a large dataset from AVA [40], they exploit a high-dimensional feature vector including aesthetics rules, biometrics and demographics features, image quality features, and fuzzy properties. Based on lasso regression output, eyes sharpness and uniqueness have the highest rank to be a good portrait.

Feedback on Photographic System: An aesthetics assessor system may find a metric value to evaluate an input image, but the way it conveys this information to photographer is more crucial, since the photographer probably has no idea about how to improve the aesthetics features of the image. That is why providing meaningful feedback to enhance the future shots and not just aesthetics assessment is our final goal in this work. Giving feedback on a photographic system firstly is mentioned by Joshi *et al.* [35], as they suggest a real-time filter to trace and aesthetically rate the camera shots, and then the photographer retake a better shot. On-site composition and aesthetics feedback system (OSCAR) [36, 37] helps smartphone users improve the quality of their taken photos by retrieving similar-aesthetic images as a qualitative composition feedback. Also it gives color combination feedback for having good colorfulness in the next taken photo, and outputs the overall aesthetics rating of the input photo as well. OSCAR is assumed to fulfill future needs of an amateur photographer, but giving such feedbacks might be unrelated or unrealistic to the user, and also it is restricted to a pretty small database in terms of coverage, diversity and copyright. Xu *et al.* [42] suggest to use three-camera array to enhance the quality of the taken photos by the rule of thirds. In fact, the smartphone interface using the camera array information shows some real-time guideline to the user for taking photo from another position. More recently general aesthetic techniques including perspectiveness [38] and triangle [39] methods are exploited to retrieve proper images as an on-site guidance to amateur photographers, but they are restricted to basic ideas in photography while the pose and scene content are ignored.

3 THE METHOD

In this section, we describe the way that we come up with our proposed framework to assist an amateur photographer intelligently capture beautiful photos from the scenes around. More specifically our proposed approach focuses on scene semantics and aesthetic features in portrait images, while in our opinion our ideas can be extended to genres. The flow of proposed framework (Figure 2) includes indexing, searching, and matching.

3.1 Our Approach

A portrait image not only contains a face but also may contain human body including head, trunk, arms, hands, and feet. Beauty of a portrait depends on the foreground positions of the human parts as well as the constellation of the background objects. The goal of portrait composition assistance is to aid a photographer to capture a better portrait given his or her current photography location or setup. The system input is an amateurishly taken photo by the

photographer or an automatically captured photo from camera viewfinder. The system output is a feedback (e.g. image, animation, comment, etc.) to guide the photographer to get better shots in next shots. A useful feedback as a side information can be any professional photo taken in a similar scene having a good pose with respect to the location/setup. We name such feedback as a *photography idea* because master photographers usually have their own ideas and each taken photo can be categorized as a new idea.

While most of image aesthetics studies focus on image assessment and manipulation of captured photos as mentioned in Section 2, there is a lack of innovative active helping with an amateur photographer to take a better shot. Also, available photographic feedback systems [35–39] have limitations in filtering unrelated photography categories or covering a broad range of photography ideas. For instance, a general retrieval system [35–37] consists of mixed photography categories including portrait, landscape, closeup, to name but a handful. Hence, this leads to an unrelated output from the input, or a feedback which is limited to a narrow range topic such as perspective photos [38] or photos having triangles [39]. The current available frameworks could not remedy the thirst of the beginners for getting professional-looking snapshots. Also, the more challenging part of the problem is that this treatment is not only a single point but also an ambiguous region because of the subjectivity in art. Expressly, there is no unique solution for an unique input, and based on various unseen tastes and manners of the subject, there may be a range of related feedbacks.

Our approach is inspired from the strategy of professional photography [5–7] — artists gradually make a subject perfect for the last shot while they usually have a “to-do” list and a “not-to-do” list in their mind. But the difference is that we do not assume the photographer has access to a studio to compose a new environment for the subject, and some background objects are static or naturally composed before. For example, when we are in the woods, the trees and sky are invariant with respect to our abilities. However, human bodies, animals, or some objects are posable, dynamic, or movable. Furthermore, the number of photography ideas for any given location is not limited to any boundary. Even if we assume that the number of photography ideas in the world is limited, this number would be very high (e.g. over 10K). To our knowledge, the performance of the deep learning models to classify an idea among high number of correlated ideas degrades substantially. Similarly, there is no accurate food category detector from dish images, because the number of food categories is high (e.g. over 65K) and the recipe retrieval is done after detection of ingredients [43].

Our method includes *decomposition* and then *composition* of semantic objects, and *matching* the proper pose. Like a chess puzzle, we should understand the constellation of the scene, and then move toward the best position. Similarly, we decompose the input shot from the amateur photographer or camera viewfinder into as many as observable objects. Indeed, we extract high-level semantics of the scene in the shot, and then realize these semantics as a whole with available photography ideas in the dataset. After finding the proper photography idea based on the current scene ingredients, in the next step of our method, known as *matching*, we follow the subject via viewfinder to match his/her pose with the available ideas, and automatically shot the scene, similar to the “smile shot” mode in smart cameras.

3.2 The Dataset

The most valuable resource of this work is the collected dataset, because it contains a large number of innovative photography ideas from around the world. We tried many photo-sharing websites for photography purposes including Flickr, Photo.net, DPChallenge, Instagram, Google Photos, Pinterest, and Unsplash. However, none of them could properly cover several categories of portrait photography comprising full body, upper body, facial, group, couple or any two, side-view, hand-only, leg-only, to mention but a few. The process of searching, collecting, and updating the dataset is time-consuming and taxing, hence, automating this process is quite helpful.

Our dataset is constructed by crawling the 500px website which contains photos from millions of photographers around the world expanding their social networks of colleagues while exploiting technical and aesthetic skills to make money. To get the JSON file list of the images sorted by rating, we wrote a distributed multi-IP, block-free Python script mostly using keywords including portrait, pose, human, person, woman, man, studio, model, fashion, male, female, and so on. We end up with over 320,000 images dataset where the number of images is still growing.

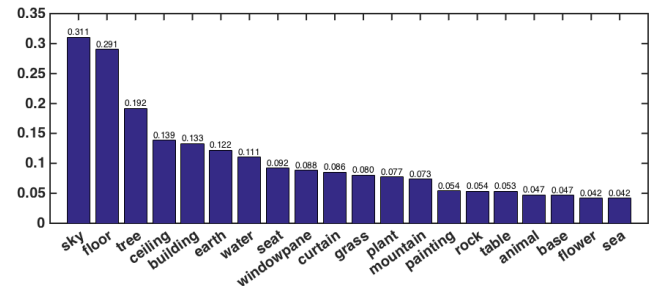


Figure 3: The distribution of the high-repetitive semantics in our dataset.

Finally, we construct a large dataset for photography ideas specially for the aforementioned portrait categories (full body, upper body, facial, group, couple or any two, side-view, hand-only, and leg-only) from highly-rated images taken by professional photographers. If we consider the semantics with area greater than the 1% of the image area, we could calculate the probability of the highly-repeated semantics in our dataset (i.e. the frequency of the semantic divided by the total number of images). These probabilities are shown in Figure 3, while we have removed “person” (probability=0.9) and “wall” (probability=0.78) from the figure, because they are dominant semantics in most of the images. Definitely having diverse semantics with high frequency in our dataset makes the proposed recommendations with respect to the query shot more helpful. After collecting the dataset and filtering the portraits, we describe the way to retrieve the corresponding results for the query image taken by the camera viewfinder in the following subsections.

3.3 Decomposition: Semantics Extraction

Extracting the detected objects in the scene as semantics of the scene is our goal in this subsection. Then, we construct available scenes in our dataset from the detected semantics and match these semantics with a sub-collection of retrieved photos in our dataset.

To achieve this goal, we explain our decomposition strategy which takes the query image from the user and gives a sorted weight list of detected semantics.

While deep learning based models can help computer vision researchers map from nearly unbounded random data domain to a nice classified range, there are still many restrictions to exploit them for applied problems. As mentioned in Section 3, there is *no limit* for innovation in art domains such as portrait photography. Hence, it is very difficult if not impossible for available deep learning architectures to learn all of these correlated ideas and classify based on the input query with high accuracy. While the number of ideas increases, mean average precision (MAP) falls abruptly with the rate of $O(\frac{1}{n})$. Also, manual idea labeling of a large dataset is costly in terms of time or money.

To tackle the problem of classifying to a large number of ideas, we detect as many objects as possible in the scene instead of photography ideas. In fact, we believe the scene captured in viewfinder consists of various static and dynamic objects. State-of-the-art deep-learned detectors (YOLO [44], PSPNet [45] and RTMPPE [46]) are customized for our purpose. YOLO [44] neural network trained on MSCOCO dataset [47] partitions the query photo into several bounding boxes predicting their probabilities. Pyramid scene parsing network (PSPNet) [45] as the winner of scene parsing challenge on ADE20K dataset [48] uses global context information through a pyramid pooling module. PSPNet predicts the scene objects in the pixel level. Real-time multi-person 2D pose estimation (RTMPPE) predicts vector fields to represent the associative locations of the anatomical parts via two sequential prediction process exposing the part confidence maps and the vector fields on MSCOCO [47] and MPII [49] datasets. To improve the accuracy, we have re-trained YOLO, PSPNet, and RTMPPE models on extended MSCOCO and ADE20K datasets by adding some of failed cases from our 500px dataset as an augmented training dataset. The illustration of some sample results are in Figure 4, where YOLO object names are shown in a red rectangle with a probability, RTMPPE pose is shown as a colorful connection of skeleton joints, and PSPNet scenes are colorized pixel-wisely based on the pixel codename.

We unify the outputs of the detectors in terms of pixel-level tensors, i.e., our modified YOLO outputs MSCOCO object IDs among 80 categories (from 1 to 80) and their scores as the minus logarithm of their NOT probability ($-\log(1-p)$) for each pixel of the image is representing as a tensor. Also our version of PSPNet outputs ADE20K object IDs among 150 categories (from 1 to 150) and the score for each pixel of the image is represented as a tensor. Similarly, our version of RTMPPE gives 18 anatomical part IDs with their scores as a tensor. So, for any image ($I_{m \times n}$) we have:

$$\begin{aligned} T_{m \times n \times 2}^{I,od} &= \begin{bmatrix} t_{i,j,k}^{I,od} \end{bmatrix}, t_{i,j,1}^{I,od} = C_{i,j}^{I,id}, t_{i,j,2}^{I,od} = -\log_2(1 - p_{i,j}^{I,od}), \\ T_{m \times n \times 2}^{I,sp} &= \begin{bmatrix} t_{i,j,k}^{I,sp} \end{bmatrix}, t_{i,j,1}^{I,sp} = A_{i,j}^{I,id}, t_{i,j,2}^{I,sp} = -\log_2(1 - p_{i,j}^{I,sp}), \\ T_{m \times n \times 2}^{I,pe} &= \begin{bmatrix} t_{i,j,k}^{I,pe} \end{bmatrix}, t_{i,j,1}^{I,pe} = J_{i,j}^{I,id}, t_{i,j,2}^{I,pe} = -\log_2(1 - p_{i,j}^{I,pe}), \end{aligned}$$

where I is an input image, m is the number of rows, n is the number of columns in the image, $T_{m \times n \times 2}^{I,od}$ is corresponding tensor of object detector (e.g. YOLO), $C_{i,j}^{I,id} \in \{1..80\}$ is MSCOCO ID of the pixel at (i, j) , $p_{i,j}^{I,od}$ is the MSCOCO ID probability of the pixel at (i, j) ;

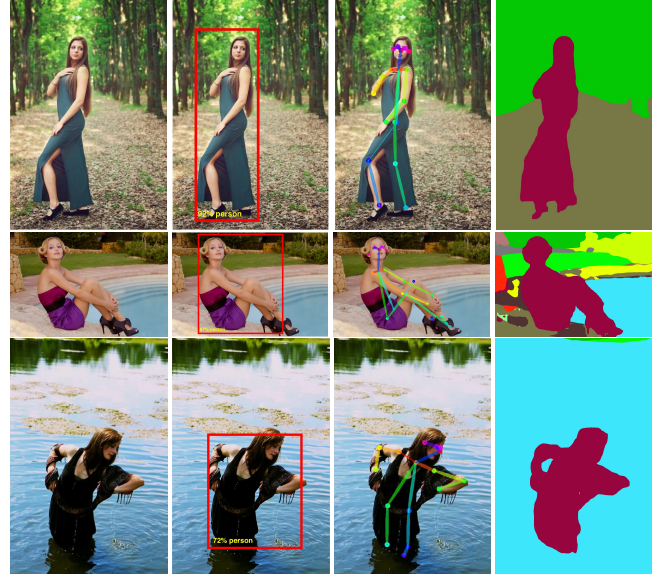


Figure 4: Sample results of our dataset where respectively from left to right are: original thumbnail, YOLO, RTMPPE, and PSPNet illustrations.

$T_{m \times n \times 2}^{I,sp}$ is tensor of scene parser (e.g. PSPNet), $A_{i,j}^{I,id} \in \{1..150\}$ is ADE20K ID of the pixel at (i, j) , $p_{i,j}^{I,sp}$ is the ADE20K ID probability of the pixel at (i, j) ; $T_{m \times n \times 2}^{I,pe}$ is tensor of pose estimator (e.g. RTMPPE), $J_{i,j}^{I,id} \in \{1..18\}$ is the joint ID of the pixel at (i, j) , and $p_{i,j}^{I,pe}$ is the joint ID probability of the pixel at (i, j) .

To auto-tag or auto-label the image, we integrate these unified results in terms of the objects, their coordinates, and their scores (or probabilities). The number of the detectable objects is 210 objects by merging MSCOCO (80 categories) and ADE20K (150 categories) objects and de-duplicating 20 objects. Also we have 18 joints from RTMPPE including nose, neck, right shoulder, right elbow, right wrist, left shoulder, left elbow, left wrist, right hip, right knee, right ankle, left hip, left knee, left ankle, left eye, right eye, left ear, and right ear. YOLO detection for full-body small persons in the image is poor, but it can detect big limbs of the body as a person well. RTMPPE detection for occluded bodies is poor but the detection for full-body persons is acceptable. Also PSPNet detection for objects, not persons, is relatively good compared to others.

First, we detect any person in the image. Our detector's integration scheme has LOW and HIGH thresholds for each detector. These thresholds are trained by a random set of highly-rated ground-truth portraits. If the average score of the pixel with person/limb ID in the image is higher than its HIGH threshold, there is a person in the image, otherwise if the average score of the pixels with person/limb ID in the image is lower than the corresponding LOW threshold, the image will be ignored from indexing or searching process, and we wait for another image for indexing or another shot for searching. We call this detector's integration scheme as *hysteresis* detection. Actually it is assured that the confidence ratio of the person in the image with his/her limbs is in a good condition using hysteresis detection. Using this filtering on our dataset, about 90% (280K+) of the images are passed. The 3D histogram of our portrait dataset in

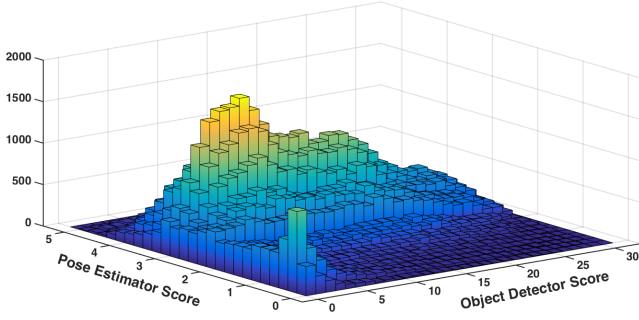


Figure 5: The 3D histogram of the portrait images binned by the object detector and pose estimator scores.

Figure 5 illustrates the frequency of the images smart-binned by the normalized object detector and pose estimator scores. In fact, it shows the effectiveness of integrating the detectors to capture the usefulness of the dataset images more precisely, because we are unifying the detection results for more broad range of ideas rather than intersecting them to have more confident narrow range of ideas.

Second, we estimate the types of the portrait in the image respectively as two (couple or any two persons), group (more than two persons), full-body, upper-body, facial, side-view, faceless, headless, hand-only, and leg-only. The number of persons is estimated by the max number of person IDs higher than their corresponding HIGH thresholds via YOLO and RTMPPE models. Otherwise, if the image contains a nose, two eyes, a hand and a leg OR a nose, an eye, two hands and two legs, it will be categorized as full body. Such combinations are learned after trying some random images as ground truth, because RTMPPE model is not perfect and, in some cases, the limbs are occluded. After examining full-body, if the image contains a nose and two eyes and one hand, it will be categorized as upper-body. After category indexing of our dataset, the distribution of the portrait category with respect to the number of corresponding images in each category by total number of images from previous step is shown in Figure 6. Consequently, the number of images for some categories like full-body, upper-body, facial, group, two, and side-view are adequate.

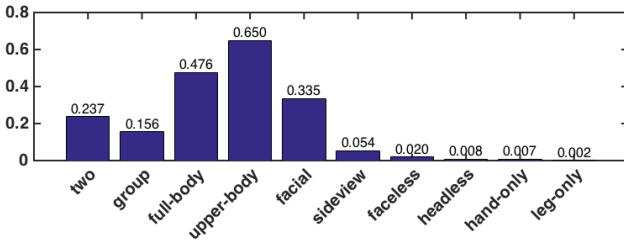


Figure 6: The distribution of the portrait categories with respect to the number of corresponding images.

Third, we seek to use other semantics. Some of them are coming from YOLO model (80 categories) and others from PSPNet (150 categories). At most there are 210 semantics (including person). To rank the order of the semantics, we exploit the max score ID multiply by the saliency map (S) features [50] with our centric

distance (D) feature to get our weighted saliency map (W).

$$W(i, j) = \max(T_{*,*,2}^{I,od}, T_{*,*,2}^{I,sp}) \times S(i, j) \times D(i, j), \quad (1)$$

$$D(i, j) = 1/K \times e^{-\|i, j\| - c\|_k}, \quad (2)$$

$$c = \frac{\sum_{i,j} S(i, j) \cdot [i, j]}{\sum_{i,j} S(i, j)}, \quad (3)$$

where $W(i, j)$ is our weighted saliency ID map, \max operation is on the 2nd matrix (score matrix) of the tensors, $S(i, j)$ is the saliency map in [50], and $D(i, j)$ is our centric distance feature, K is a tunable constant, c is the center of mass coordinate, and $\|\cdot\|_k$ is the k -th norm operator where $k = 1$ in our experiments. Based on various sorted semantics called as portrait scenes, we have indexed our portrait-categorized dataset from the previous step, and Figure 7 depicts the number of portrait scenes for some of highly-frequent portrait categories.

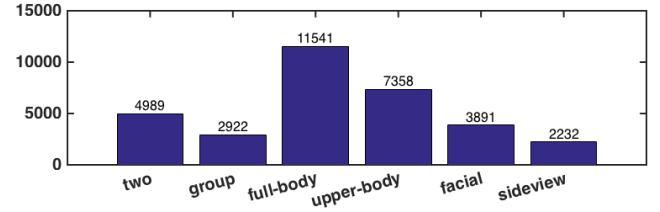


Figure 7: The frequency of the portrait scenes with respect to the highly-requested portrait categories.

Our weighted saliency map can make the detected objects in-order, because we can sum up the scores from the semantics and sort them based on their accumulated weights. The output of this step is an ordered list of detected semantics in the query image. In the next step, we will find the closest image to this auto-tagged ordered object list of the query image.

3.4 Composition: Scene Construction

The goal of composition step is to build up a new scene consisting of the desired objects in it. The input of this stage is the ordered weight list of the semantics as well as the image query. The output of this stage is a collection of well-posed images corresponding to the query image. As we focus on portraits, we desire the targeted image to contain a well-posed portrait with similar semantics. The “interaction” between the subject(s) and the objects around is important, because the proposed pose by the system is dependent on them.

As we have collected a dataset containing generally well-posed portrait, we should dig into the dataset and look for an image with similar object constellation, and the existence of this professional-quality dataset makes it possible that the retrieved photos contain mostly good aesthetic photography ideas. Our image retrieval system is not supposed to find images with similar colors, patterns, or poses, but it tries to find images with better poses with similar semantics. So the location of the movable objects doesn’t matter, but the detected objects are important.

To look for a semantically composed version with respect to the query, we exploit the ordered weight list of the detected objects in Eq. 1 as well as the other common feature vectors [51] of the

Algorithm 1 Semantic Retrieval

```
1: procedure SemanticRetrieval( $Q \in \text{ViewFinder}$ )
2:    $T_{m \times n \times 2}^{Q, od} \leftarrow \text{Object\_Detect}(Q)$ 
3:    $T_{m \times n \times 2}^{Q, sp} \leftarrow \text{Scene\_Parse}(Q)$ 
4:    $T_{m \times n \times 2}^{Q, pe} \leftarrow \text{Pose\_Estimate}(Q)$ 
5:    $T_{m \times n \times 2}^{Q, pf} \leftarrow \text{Common\_Features\_Extract}(Q)$ 
6:    $c^Q \leftarrow \frac{\sum_{i,j} \|[i,j]-[0,0]\|_k \times S(i,j)}{\sum_{i,j} \|[i,j]-[0,0]\|_k}$ 
7:    $D^Q(i,j) \leftarrow 1/K \times e^{-\|[i,j]-c^Q\|_k}$ 
8:    $W^Q(i,j) \leftarrow \max \left( T_{*,*,2}^{I, od}, T_{*,*,2}^{I, sp} \right) \times S^Q(i,j) \times D^Q(i,j)$ 
9:    $V^{Q, pref} \leftarrow M_{PCM} \cdot W^Q$ 
10:   $\text{Retrieved\_Indexes} \leftarrow \text{Index\_Sort}(V^{Q, pref})$ 
11:   $\text{Show\_Top4}(\text{Retrieved\_Indexes})$ 
12: end procedure
```

query image, because we do not assume the query image taken by an amateur photographer to be well-posed enough to be the query base for our image retrieval system. In fact, we just want to understand the location/setup around the subject, and then based on the scene ingredients, a well-posed image taken by a professional is proposed to the photographer.

Based on the ordered weight list of the detected objects in the image, we can apply the same operations on all of our images in the dataset, and then retrieve the highest-ranked candidates as the results. The operations are the same as mentioned for decomposition, including features from YOLO, PSPNet, and RTMPPE to detect the possible persons, their joints, and other objects with their scores. Also, we should compute our weighted saliency map to rank the detected objects. The ordered weight list of the semantics with other common features of our portrait dataset, known as *portrait composition model* (PCM) in Figure 1, is represented as a large “number of images by number of features” matrix (M_{PCM}). Similarly, the ordered weight list and other common features of the image query is represented as a long feature vector (W^I). The distance of the W^I from each row of M_{PCM} is defined as inner vector product. Consequently, we have:

$$V^{\text{pref}} = M_{PCM} \cdot W^I, \quad (4)$$

where V^{pref} is our preference vector, and if we sort it based on the vector values, the indexes of the rows represent the highest-ranked candidates as our feedbacks to the image query (I). The whole process of semantic retrieval for an input image query (Q) has been shown in Algorithm 1. Our experimental results show the quality of the results.

Notes on Indexing and Searching: Our semantic retrieval system is equivalent to the flows of indexing and searching in Figure 2. Practically, there are many challenges in image retrieval systems [52–54] as well as in our case. To improve the speed of our image retrieval system, we compute the decomposition step for all images in our dataset. Indexing procedure is lengthy for the first time, but at the time of update it is fast enough because the semantic and feature extraction for an image is real-time using GPU. Furthermore, indexing procedure for our retrieval system organizes the dataset images into categorized folders labeled by the sorted semantic list.

Consequently, the composition step is also fast, as it just extracts the query image ordered semantic list, and then finds the target folder containing similar images, and finally retrieves first five best results from the folder with respect to the top-4 indexes in sorted preference vector (Eq. 4) explained in Algorithm 1. As we just include semantics with normalized score higher than 10 percent of the total semantics score, the number of ordered semantics are limited (typically less than 5 names). Also, naturally it is impossible to cover all semantic combinations, so the scene semantics are limited in number.

3.5 Matching

Professional photographer starts to pose the subject from head to toe step-by-step, while there are many to-do list and not-to-do list for portrait photography in his/her mind. We want to create the same environment in a smart camera to accompany an amateur photographer gradually to his/her perfect shot. From semantic retrieval subsection, we retrieve the proper photography ideas given a query image from the camera viewfinder, and we assume that the photographer has chosen some of the retrieved images as a desired set, and disregarded the others as an ignored set. Now we explain how to capture the proper pose of the subject in the scene, and trigger the “pose shot” for the camera.

The variant component in our framework is the human pose. The relative positions of the human body components (including nose, eyes, ears, neck, shoulders, elbows, wrists, hips, knees, and ankles) with respect to the nose position as portrait origin are consisting our pose model. Preferably, we would like to start from the position of the nose ($J_0 = (0, 0)$) that is connected to neck (J_1), right eye (J_2), and left eye (J_3) are connected to right ear (J_4), left ear (J_5) as they are on a plane of the head. Also, shoulders (J_6 and J_7) can be recognized by a length and an angle from neck, and similarly elbows (J_8 and J_9) from shoulders, wrists (J_{10} and J_{11}) from elbows, hips (J_{12} and J_{13}) from neck, knees (J_{14} and J_{15}) from hips, and ankles (J_{16} and J_{17}) from knees, i.e. they are connected as follows:

$$\text{Pre}(J_i) = J_0, \quad i \in \{0, 1, 2, 3\}, \quad (5)$$

$$\text{Pre}(J_i) = J_1, \quad i \in \{6, 7, 12, 13\}, \quad (6)$$

$$\text{Pre}(J_i) = J_{i-2}, \quad i \in \{4, 5, 8, 9, 10, 11, 14, 15, 16, 17\}. \quad (7)$$

So, we can always calculate the absolute position using 2D polar coordinates as follows:

$$J_i = J_j + r_{i,j} \cdot e^{i\theta_{i,j}}, \quad i \in \{0..17\}, \quad (8)$$

where $j = \text{Pre}(i)$ i.e. part j is the previous part connected to part i , $r_{i,j}$ is the length from joint J_i to joint J_j , $\theta_{i,j}$ is the angle between the line from joint J_i to joint J_j and the line from joint J_j to joint $\text{Pre}(J_j)$, and the line crossing J_0 is the image horizon. i is the unit imaginary number. Note that for a 2D human body $r_{i,j}; \forall i, j$ are fixed, but $\theta_{i,j}; \forall i, j$ can be changed to some fixed not arbitrary extents. Also having 3D pose-annotated/estimated single depth images, similarly we can calculate the relative 3D position of the joints using spherical coordinates. So, we have such action boundaries for joints as follows:

$$\theta_{i,j}^{\min} \leq \theta_{i,j} \leq \theta_{i,j}^{\max}, \quad j = \text{Pre}(i), \quad (9)$$

$$\phi_{i,j}^{\min} \leq \phi_{i,j} \leq \phi_{i,j}^{\max}, \quad j = \text{Pre}(i). \quad (10)$$

As a result, a human body pose (J) is represented by:

$$J^k = (J_1^k, J_2^k, \dots, J_{17}^k), \quad (11)$$

where J^k is the pose for k -th person (or k -th image with one person), and $\forall i \in \{1..17\} : J_i^k$ is the i -th coordinate of the k -th person. Also we need a distance metric to calculate the difference between two pose features. So we define the distance metric as follows:

$$D(J^k, J^l) = \sum_{i=1}^{17} \|J_i^k - J_i^l\|_q, \quad (12)$$

where $D(\cdot)$ is the distance operator, where J^k is the pose feature for k -th person (or k -th image with one person), $\forall i \in \{1..17\} : J_i^k$ is the i -th coordinate of the k -th person, and $\|\cdot\|_q$ (usually L1-norm or L2-norm) is the L_q -norm function of two equal-length tuples.

Now, the camera viewfinder may take and hold several photos gradually from the scene, and finally choose the best among them to save onto the camera storage. Actually our matching algorithm searches among the taken photos to get the nearest pose to one of the collected ideas. It is an integer programming problem to find the best seed among all photography ideas. Given the distance operator of two pose features explored in 12, we can construct our optimization problem as the maximum over all taken photos of the difference of the minimum distance of the ignored set and minimum distance of the desired set. Mathematically, we compute the following optimization problem subject to 9 and 10:

$$I_w = \arg \max_{\forall I_i \in I^t} \left(\min_{\forall Q_j^g \in Q^g} D(J^{Q_j^g}, J^I_i) - \min_{\forall Q_k^d \in Q^d} D(J^{Q_k^d}, J^I_i) \right),$$

where I_w is the wish image, I^t is the set of taken photos, Q^g is the set of ignored retrieved ideas, Q^d is the set of desired retrieved ideas, $D(\cdot)$ is the distance operator in 12, and J^x is the pose for x -th image with one person in 11. The optimization problem in continuous mode (not over taken images set) may have (a) solution(s) in feasible region, and in L1-norm case, it is equivalent to multiple linear programming problems but the complexity of the problem will be exponential, and also the solution is not always a desired pose.

3.6 User Study

Currently, there is no other similar or comparable system in the literature to contrast with our proposed framework. To evaluate the functionality and the performance of our method, and measure how much the recommended photos make sense and are helpful to the photographer, we conduct a quantitative user study based on the human subject results to compare our method with state-of-the-art semantic and scene retrieval based on CNN [51] and KNN-SVM [55]. We select a variety of image queries based on many types of portrait categories such as background scene and semantics, single versus group, full-body, upper-body, facial, standing versus sitting, and male versus female. All 4096 generic descriptors via public CNN model [56] trained on ImageNet [57] are extracted for our large dataset images as well as the features of KNN-SVM-based method [55]. Using a PHP-based website with a usage guidance, the hypothesis tests are asked, and the outputs of the methods are randomly shown in each row to be chosen by more than thirty participants who are graduate students. Our framework received

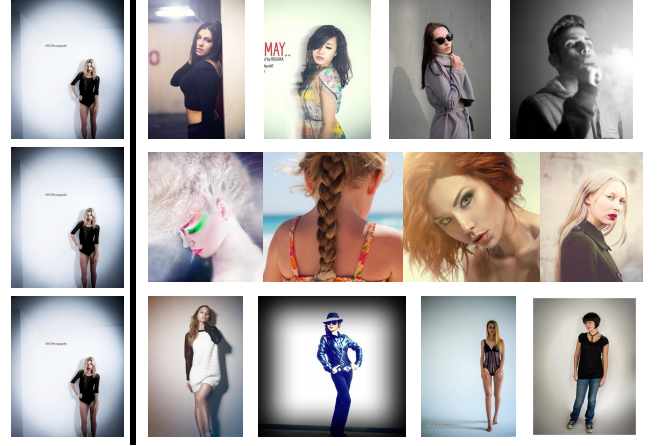


Figure 8: The results of CNN (1st row), KNN-SVM (2nd row), and our method (3rd row) for a sample shot at the left side.

65.3% of the 1st ranks among the tests compared to 27.1% CNN as 2nd rank and 7.6% KNN-SVM as 3rd rank. Figure 8 illustrates the results of all methods (CNN: 1st, KNN-SVM: 2nd, ours: 3rd row) with respect to a similar shot at the left side. As it is realized from Figure 8 and our study, the other methods cannot capture portrait categories, scene structure, and corresponding poses of the query shot very well, because the badly-posed full-body query shot is suggested as upper-body, facial, and back poses by other category-agnostic methods. As we hierarchically index our dataset by portrait images, portrait categories, and then semantic categories, our semantic-aware framework accessing to our indexed dataset can retrieve related photography ideas.

4 CONCLUSIONS AND FUTURE WORK

We have collected a large dataset for portrait photography ideas and introduced a new framework for portrait composition assistance which aids amateur photographer to capture a better shot. As the number of photography ideas are increasingly high, directly retrieving and matching the viewfinder photo with an image in our dataset is not straightforward. Furthermore, the retrieving system not only finds similar images but also searches for images with similar semantics through decomposition and composition stages. After providing feedbacks for the photographer, the camera tries to match the final pose with one of the retrieved feedbacks, and make a pose-shot. The performance of our framework has been evaluated by a user study. Another merit of this work is the integration of the deep-based detectors which can make the whole process automatic.

This work can be extended to other photography genres such as fashion, close-up, and landscape photography using other appropriate detectors. The criteria for one genre are generally different from those for another. For instance, the pose is crucial in portrait photography, while leading lines and vanishing points can be important in landscape or cityscape photography. It can also be interesting to investigate a metric to estimate the potential novelty of the current shot based on computing similarity to other shots.

Acknowledgments: The authors would like to thank the participants of the human subject study.

REFERENCES

- [1] David A Lauer and Stephen Pentak. *Design Basics*. Wadsworth Publishing, 2011.
- [2] Roberto Valenzuela. *Picture Perfect Practice: A Self-Training Guide to Mastering the Challenges of Taking World-Class Photographs*. New Riders, 2012.
- [3] Bert Krages. *Photography: The Art of Composition*. Skyhorse Publishing, Inc., 2012.
- [4] Patrick Rice. *Master Guide for Professional Photographers*. Amherst Media, 2006.
- [5] Jeff Smith. *Posing for Portrait Photography: A Head-To-Toe Guide for Digital Photographers*. Amherst Media, 2012.
- [6] Roberto Valenzuela. *Picture Perfect Posing: Practicing the Art of Posing for Photographers and Models*. New Riders, 2014.
- [7] Christopher Grey. *Master Lighting Guide for Portrait Photographers*. Amherst Media, 2014.
- [8] Ritendra Datta, Dhiraj Joshi, Jia Li, and James Z Wang. Studying aesthetics in photographic images using a computational approach. *European Conference on Computer Vision*, pages 288–301, 2006.
- [9] Yiwen Luo and Xiaoou Tang. Photo and video quality evaluation: Focusing on the subject. In *European Conference on Computer Vision*, pages 386–399, 2008.
- [10] Lai-Kuan Wong and Kok-Lim Low. Saliency-enhanced image aesthetics class prediction. In *IEEE Conference on Image Processing*, pages 997–1000, 2009.
- [11] Luca Marchesotti, Florent Perronnin, Diane Larlus, and Gabriela Csuska. Assessing the aesthetic quality of photographs using generic image descriptors. In *IEEE Conference on Computer Vision*, pages 1784–1791, 2011.
- [12] Subhabrata Bhattacharya, Rahul Sukthankar, and Mubarak Shah. A framework for photo-quality assessment and enhancement based on visual aesthetics. In *ACM Conference on Multimedia*, pages 271–280, 2010.
- [13] Subhabrata Bhattacharya, Rahul Sukthankar, and Mubarak Shah. A holistic approach to aesthetic enhancement of photographs. *ACM Transactions on Multimedia Computing, Communications, and Applications*, 7(1):21, 2011.
- [14] Bongwon Suh, Haibin Ling, Benjamin B Bederson, and David W Jacobs. Automatic thumbnail cropping and its effectiveness. In *ACM symposium on User Interface Software and Technology*, pages 95–104, 2003.
- [15] Anthony Santella, Maneesh Agrawala, Doug DeCarlo, David Salesin, and Michael Cohen. Gaze-based interaction for semi-automatic photo cropping. In *ACM SIGCHI Conference on Human Factors in Computing Systems*, pages 771–780, 2006.
- [16] Fred Stentiford. Attention based auto image cropping. In *ICVS Workshop on Computational Attention & Applications*, volume 1. Citeseer, 2007.
- [17] Mingju Zhang, Lei Zhang, Yanfeng Sun, Lin Feng, and Wei-ying Ma. Auto cropping for digital photographs. In *IEEE Conference on Multimedia and Expo*, 2005.
- [18] Jaesik Park, Joon-Young Lee, Yu-Wing Tai, and In So Kweon. Modeling photo composition and its application to photo re-arrangement. In *IEEE Conference on Image Processing*, pages 2741–2744, 2012.
- [19] A Samii, R M  ch, and Z Lin. Data-driven automatic cropping using semantic composition search. In *Computer Graphics Forum*, volume 34, pages 141–151. Wiley Online Library, 2015.
- [20] Jianzhou Yan, Stephen Lin, Sing Kang, and Xiaoou Tang. Learning the change for automatic image cropping. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 971–978, 2013.
- [21] Ligang Liu, Yong Jin, and Qingbiao Wu. Realtime aesthetic image retargeting. In *Computational Aesthetics*, pages 1–8, 2010.
- [22] Hui-Tang Chang, Po-Cheng Pan, Yu-Chiang Frank Wang, and Ming-Syan Chen. R2p: Recomposition and retargeting of photographic images. In *ACM Conference on Multimedia Conference*, pages 927–930, 2015.
- [23] Connelly Barnes, Eli Shechtman, Adam Finkelstein, and Dan Goldman. Patchmatch: A randomized correspondence algorithm for structural image editing. *ACM Transactions on Graphics*, 28(3):1–24, 2009.
- [24] Taeg Sang Cho, Moshe Butman, Shai Avidan, and William T Freeman. The patch transform and its applications to image editing. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8, 2008.
- [25] Yael Pritch, Eitam Kav-Venaki, and Shmuel Peleg. Shift-map image editing. In *International Conference on Computer Vision*, volume 9, pages 151–158, 2009.
- [26] YW Guo, M Liu, TT Gu, and WP Wang. Improving photo composition elegantly: Considering image similarity during composition optimization. In *Computer Graphics Forum*, volume 31, pages 2193–2202. Wiley Online Library, 2012.
- [27] Ke Li, Bo Yan, Jun Li, and Aditi Majumder. Seam carving based aesthetics enhancement for photos. *Signal Processing: Image Communication*, 39:509–516, 2015.
- [28] Shehroz S Khan and Daniel Vogel. Evaluating visual aesthetics in photographic portrait. In *Symposium on Computational Aesthetics in Graphics, Visualization, and Imaging*, pages 55–62. Eurographics Association, 2012.
- [29] Matija Males, Adam Hedi, and Mislav Grgic. Aesthetic quality assessment of headshots. In *International Symposium ELMAR*, pages 89–92, 2013.
- [30] Shao-Fu Xue, Hongying Tang, Dan Tretter, Qian Lin, and Jan Allebach. Feature design for aesthetic inference on photos with faces. In *IEEE Conference on Image Processing*, pages 2689–2693, 2013.
- [31] Arnaud Lienhard, Marion Reinhard, Alice Caplier, and Patricia Ladret. Photo rating of facial pictures based on image segmentation. In *IEEE Conference on Computer Vision Theory and Applications*, volume 2, pages 329–336, 2014.
- [32] Arnaud Lienhard, Patricia Ladret, and Alice Caplier. Low level features for quality assessment of facial images. In *Conference on Computer Vision Theory and Applications*, pages 545–552, 2015.
- [33] Arnaud Lienhard, Patricia Ladret, and Alice Caplier. How to predict the global instantaneous feeling induced by a facial picture? *Signal Processing: Image Communication*, 39:473–486, 2015.
- [34] Miriam Redi, Nikhil Rasiwasia, Gaurav Aggarwal, and Alejandro Jaimes. The beauty of capturing faces: Rating the quality of digital portraits. In *IEEE Conference on Automatic Face and Gesture Recognition*, volume 1, pages 1–8, 2015.
- [35] Dhiraj Joshi, Ritendra Datta, Elena Fedorovskaya, Quang-Tuan Luong, James Z Wang, Jia Li, and Jiebo Luo. Aesthetics and emotions in images. *IEEE Signal Processing Magazine*, 28(5):94–115, 2011.
- [36] Lei Yao, Poonam Suryanarayan, Mu Qiao, James Z Wang, and Jia Li. Oscar: On-site composition and aesthetics feedback through exemplars for photographers. *International Journal of Computer Vision*, 96(3):353–383, 2012.
- [37] Jia Li, Lei Yao, and James Z Wang. Photo composition feedback and enhancement. In *Mobile Cloud Visual Media Computing*, pages 113–144. Springer, 2015.
- [38] Zihan Zhou, Farshid Farhat, and James Z Wang. Detecting dominant vanishing points in natural scenes with application to composition-sensitive image retrieval. *IEEE Transactions on Multimedia*, 2017.
- [39] Siqiong He, Zihan Zhou, Farshid Farhat, and James Z Wang. Discovering triangles in portraits for supporting photographic creation. *IEEE Transactions on Multimedia*, 2017.
- [40] Naila Murray, Luca Marchesotti, and Florent Perronnin. AVA: A large-scale database for aesthetic visual analysis. In *Conference on Computer Vision and Pattern Recognition*, pages 2408–2415, 2012.
- [41] Jonathan Harel, Christof Koch, and Pietro Perona. Graph-based visual saliency. In *Advances in Neural Information Processing Systems*, pages 545–552, 2006.
- [42] Yan Xu, Joshua Ratcliff, James Scovell, Gheric Speigener, and Ronald Azuma. Real-time guidance camera interface to enhance photo aesthetic quality. In *ACM Conference on Human Factors in Computing Systems*, pages 1183–1186, 2015.
- [43] Jingjing Chen and Chong-Wah Ngo. Deep-based ingredient recognition for cooking recipe retrieval. In *ACM on Multimedia Conference*, pages 32–41, 2016.
- [44] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 779–788, 2016.
- [45] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing network. *arXiv preprint arXiv:1612.01105*, 2016.
- [46] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. *arXiv preprint arXiv:1611.08050*, 2016.
- [47] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Doll  r, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision*, pages 740–755, 2014.
- [48] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [49] Mykhaylo Andriluka, Leonid Pishchulin, Peter Gehler, and Bernt Schiele. 2d human pose estimation: New benchmark and state of the art analysis. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2014.
- [50] Laurent Itti, Christof Koch, and Ernst Niebur. A model of saliency-based visual attention for rapid scene analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20(11):1254–1259, 1998.
- [51] Ali Sharif Razavian, Hossein Azizpour, Josephine Sullivan, and Stefan Carlsson. Cnn features off-the-shelf: an astounding baseline for recognition. In *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 806–813, 2014.
- [52] Arnold WM Smeulders, Marcel Worring, Simone Santini, Amarnath Gupta, and Ramesh Jain. Content-based image retrieval at the end of the early years. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(12):1349–1380, 2000.
- [53] Michael S Lew, Nicu Sebe, Chabane Djeraba, and Ramesh Jain. Content-based multimedia information retrieval: State of the art and challenges. *ACM Transactions on Multimedia Computing, Communications, and Applications*, 2(1):1–19, 2006.
- [54] Ritendra Datta, Dhiraj Joshi, Jia Li, and James Z Wang. Image retrieval: Ideas, influences, and trends of the new age. *ACM Computing Surveys*, 40(2):5, 2008.
- [55] Joani Mitro. Content-based image retrieval tutorial. *arXiv preprint arXiv:1608.03811*, 2016.
- [56] Ken Chatfield, Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Return of the devil in the details: Delving deep into convolutional nets. *arXiv preprint arXiv:1405.3531*, 2014.
- [57] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009.