

# CAPTAIN: Comprehensive Composition Assistance for Photo Taking

Farshid Farhat · Mohammad Mahdi Kamani · James Z. Wang

Received: date / Accepted: date

**Abstract** Many people are interested in taking astonishing photos and sharing with others. Emerging high-tech hardware and software facilitate ubiquitousness and functionality of digital photography. Because composition matters in photography, researchers have leveraged some common composition techniques to assess the aesthetic quality of photos computationally. However, composition techniques developed by professionals are far more diverse than well-documented techniques can cover. We leverage the vast underexplored innovations in photography for computational composition assistance. We propose a comprehensive framework, named CAPTAIN (Composition Assistance for Photo Taking), containing integrated deep-learned semantic detectors, sub-genre categorization, artistic pose clustering, personalized aesthetics-based image retrieval, and style set matching. The framework is backed by a large dataset crawled from a photo-sharing Website with mostly photography enthusiasts and professionals. The work proposes a sequence of steps that have not been explored in the past by researchers. The work addresses personal preferences for composition through presenting a ranked-list of photographs to the user based on user-specified weights in the similarity measure. The matching algorithm recognizes the best shot among a sequence of shots with respect to the user's preferred

style set. We have conducted a number of experiments on the newly proposed components and reported findings. A user study demonstrates that the work is useful to those taking photos.

**Keywords** Computational Composition · Image Aesthetics · Photography · Deep Learning · Image Retrieval

## 1 Introduction

Digital photography is of great interest to many people, regardless of whether they are professionals or amateurs. For example, people on social networks often share their photos with their friends and family. It has been estimated that over a billion photos are taken every year, and many people take photos with smartphones primarily. Smartphones' increasing computing power and ability to connect to more powerful computing platforms via the network make them potentially useful as a composition assistant to amateur photographers (Yao et al., 2012).

Besides, emerging technologies, including artificial intelligence (AI)-chips and AI-aware mobile applications, provide more opportunities for composition assistance. Taking stunning photos often needs expertise and experience at a level that professional photographers have. Like in other visual arts, a lack of common alphabet similar to music notes or mathematical equations makes transferring of knowledge in photography difficult. To many amateurs, as a result, photography is mysterious and gaining skills is not easy and cannot be done quickly. Nonetheless, many people are fascinated about professional-quality photos and desire to have the ability to create similar-quality photos themselves.

---

F. Farhat (✉)  
School of Electrical Engineering and Computer Science  
The Pennsylvania State University, University Park, PA, USA  
E-mail: fuf111@psu.edu

M.M. Kamani · J.Z. Wang  
College of Information Sciences and Technology  
The Pennsylvania State University, University Park, PA, USA  
E-mail: mqk5591@psu.edu

J.Z. Wang  
E-mail: jwang@ist.psu.edu

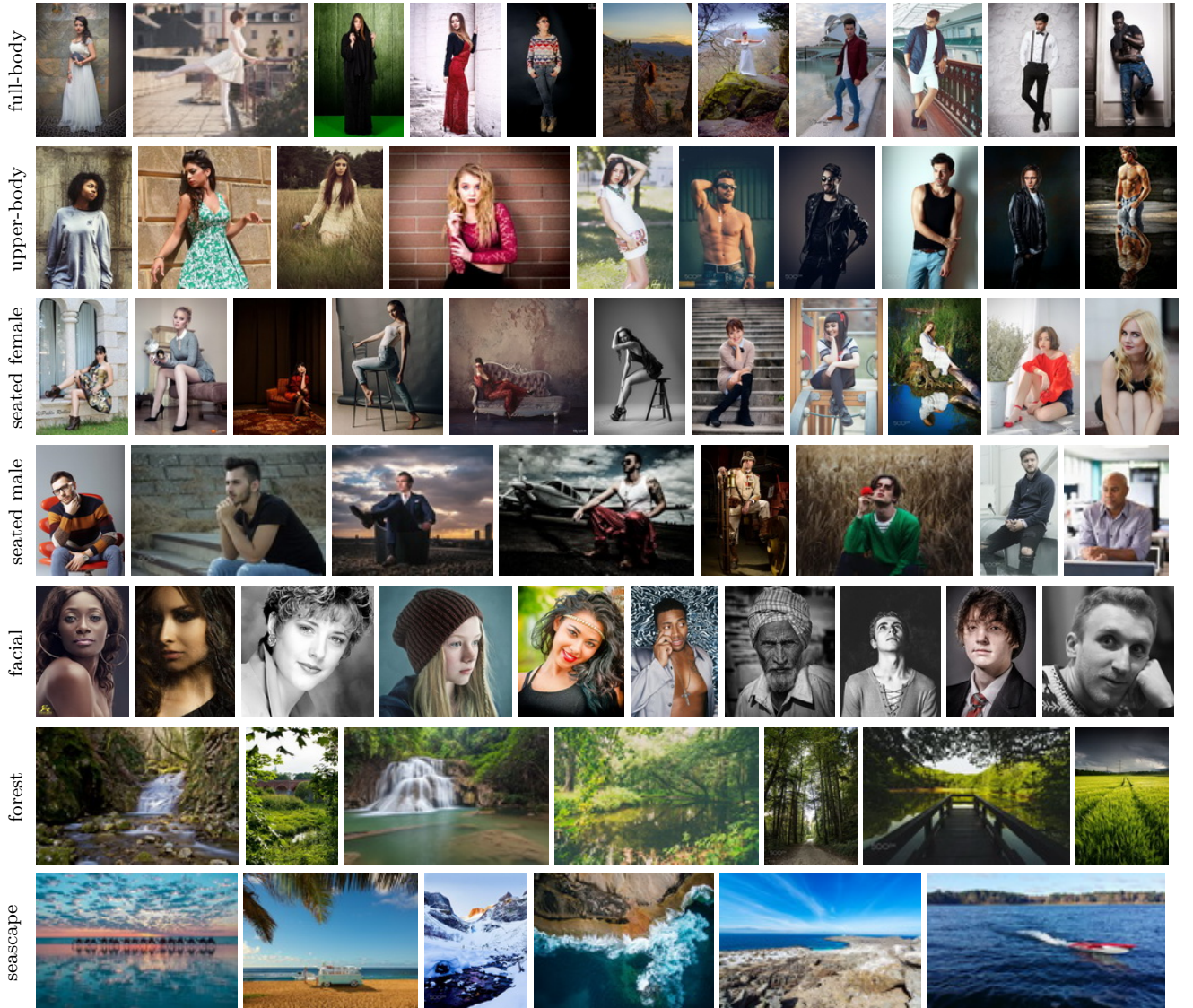


Fig. 1: Photos retrieved from the dataset based on the photo category and/or subject gender. Each retrieved result shows a collection of photography ideas that can be used by an amateur to compose photos for a given situation. Photos of the dataset were crawled from the 500px Website.

Because aesthetics in photography is strongly linked to human creativity, it is daunting for an artificial intelligence (AI) to compose photographs at a given scene or a given studio setup that can impress people in a way professional photographers do. In our work, we attempt to connect human creativity as demonstrated through their creative works with AI.

Aesthetics and composition in photography have generally been heuristically explored and known as a collection of rules or principles such as balance, geometry, symmetry, the rule of thirds, and framing (Lauer and Pentak, 2011; Valenzuela, 2012; Krages, 2012). It is well known that professional photographers take a lot of pic-

tures, and through their practice they gain experience and knowledge which in turn enable them to be creative. They have written about their knowledge in photography books (Smith, 2012; Valenzuela, 2014; Grey, 2014). Some composition rules or principles have been well articulated and many amateurs make use of these principles in their photo taking. However, we argue that the set of known rules or principles can hardly cover the creativity and experience of thousands of photographers around the world (Valenzuela, 2012; Matthew, 2010; Rice, 2006).

In order to capture an aesthetically appealing photo, photographers often integrate different visual elements.

For instance, the beauty of a full-body portrait depends on the foreground positions of the human limbs as well as the constellation of the background objects (if any). A good portrait is often a product of an appropriate color palette, an appealing composition of shapes, and an interesting human pose. There is no unique photography idea for a given situation, and people have different opinions on those ideas depending on their cultural background, gender, age, experience, and emotional state. As a result, if the aesthetic quality of photos is quantified by one number, one is making an unrealistic assumption that different people share the same opinions on the same photo.

For an amateur, it would be helpful if an AI can help select *photography ideas*, from thousands of such ideas available through online photos taken with professional quality, for a given scene or a given studio setup. The key technical difficulties for accomplishing this goal are (1) finding a suitable mapping between a professional-quality photo of a scene and the underlying photography ideas, (2) there are virtually unlimited number of photography ideas, and (3) to provide meaningful and intuitive in-situ assistance to the photographer based on personal preference. Our work tackles these challenges using a data-driven approach based on retrieval from a large dataset of professional-quality photos.

The multimedia and computer vision research communities have been leveraging some of the photography composition rules or principles for aesthetics and composition assessment (Datta et al., 2006; Ke et al., 2006; Luo and Tang, 2008; Wong and Low, 2009; Marchesotti et al., 2011). Other approaches manipulated (modified) the photo to comply with artistic rules (Bhattacharya et al., 2010, 2011), and such systems are referred to as auto-composition or re-composition. The techniques include smart cropping (Suh et al., 2003; Santella et al., 2006; Stentiford, 2007; Zhang et al., 2005; Park et al., 2012; Samii et al., 2015; Yan et al., 2013), warping (Liu et al., 2010; Chang et al., 2015), patch re-arrangement (Barnes et al., 2009; Cho et al., 2008; Pritch et al., 2009), cutting and pasting (Bhattacharya et al., 2011; Zhang et al., 2005), and seam carving (Guo et al., 2012; Li et al., 2015b). However, they do not help an amateur photographer capture a more impressive photo to begin with. More specifically in portrait photography, there have been rule-based assessment models (Khan and Vogel, 2012; Males et al., 2013) using known photography basics to evaluate portraits, and facial assessment models (Xue et al., 2013; Lienhard et al., 2014, 2015b,a; Redi et al., 2015) exploiting features including smile, age, and gender from face. On-site feedback systems (Yao et al., 2012; Li et al., 2015a) have been developed to help amateur photographers by retrieving im-

ages with similar composition, but the system is limited to basic composition categories (e.g. horizontal, vertical, diagonal, textured, and centered). More recently, perspective-related techniques (Zhou et al., 2017b), the triangle technique (He et al., 2018) and some portrait composition techniques (Farhat et al., 2017) have also been exploited.

We investigate a holistic framework for helping people take a better shot with regard to their current photography location and need. The framework addresses the differences in preferences of the users through adjusting the ranking process used to retrieve exemplars. After getting a first shot from the camera, our framework provides some highly-scored related photos as pre-composed “recipes” (i.e. photography ideas) for the user to consider. As an example, regarding some personalized criteria (such as photo category and subject gender), Figure 1 shows sample results retrieved from the photo dataset. These photos illustrate various locations, scenes, and categories. One can argue that while photos in the same row have the same category or gender, each individual photo has a photography idea(s) that is different from those used in other photos of the same row. For example, in the 2nd and 3rd photos from the right in the 1st row, the subjects cross their legs and bend one of the knees to form a triangle in the resulting photo. As mentioned before, the triangle technique is a popular technique used by professionals. While both uses the technique, the way they use them is different, forming different photography ideas. Similarly, the first subject in the 3rd row sits beside the arch with an apropos pose, creating a triangle, but is different from the triangles formed in the earlier examples.

We address the complexity of transferring photography idea(s) to a user through providing useful exemplar-based feedbacks. Specifically, we break down the scene that the user wants to take a photo from into composition primitives, and then build them up for a better-composed shot using highly-rated similar photos from the dataset. To accommodate the user’s individual preferences, we perform personalized aesthetics-related image retrieval (PAIR). Figure 2 shows the flowchart of our approach for assisting photographers in taking an improved photo. Based on the first shot as a query, some highly-rated photos are retrieved from the collected dataset using the user-specified preferences (USP) and our composition model (CM). Then, the results are shown to the photographer to select some of them as a user-preferred style set. The camera then takes a sequence of shots, from which the one that is the closest match to the style set is chosen. The details of the procedure will be explained later. The **main contributions** of our work are as follows:



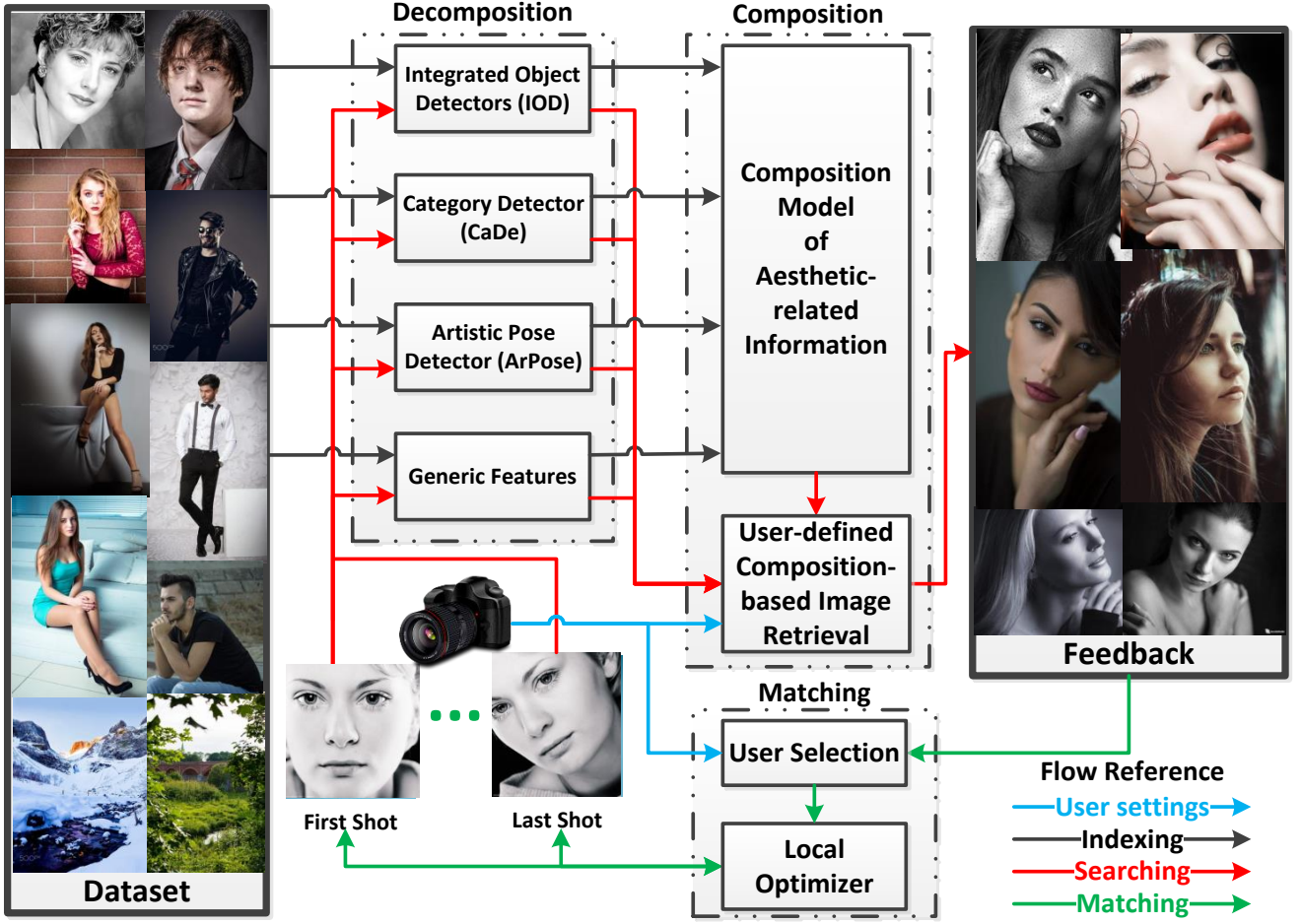


Fig. 2: The flowchart of our composition assistance framework: The blue, black, red, and green flows show the user settings, the indexing process, the searching process, and the matching process, respectively. The decomposition step (the box with dashed line) extracts the aesthetics-related information of the images and computes our composition model. The composition step (another box with dashed line) searches for well-composed images in our dataset based on the aesthetics-related information and user-specified preferences. The matching step (the other small box at the bottom) considers the next shots, and finds the shot that is closest to the user-preferred style set.

- We propose a new framework that finds the mapping between a photo and its potential underlying photography idea(s) through decomposing the photo into composition ingredients. Through such a framework, it is possible to leverage the virtually unlimited number of photography ideas available on the Internet. We design the *decomposition* step to extract composition primitives of a query shot using various detectors including newly developed integrated object detector (IOD), category detector (CaDe), and artistic pose detector (ArPose). The IOD consists of a collection of performance-enhanced detectors, the integration of which substantially boosts the detection accuracy by hysteresis detection and makes them unified and compatible with any detector. The CaDe has top-down hi-

erarchical clustering and multi-class categorization to leverage sub-genre information. The ArPose performs pose clustering to extract pose information using skeleton context features.

- We address the complexity of transferring photography knowledge, caused by the existence of abundant, diverse, and correlated photography ideas for any given scene, by providing meaningful and intuitive feedback to amateur photographers. We design the *composition* step to perform personalized aesthetics-related image retrieval from a large managed dataset containing over 200,000 highly-rated aesthetically composed photos covering a large number of photography ideas and different categories.
- We accommodate the user’s personal preferences for composition, through showing a ranked-list of pho-

tos to the user based on user-specified preferences (USP). The framework further helps the user to select a shot among a sequence of shots that is optimal with respect to the user’s preferred style set.

Our proposed framework is not limited to any specific photography genre. It can be generalized to other genres such as architectural or closeup photography.

In the remainder of the paper, we explain the realization of our framework in detail. After discussing related work, we describe the dataset that we have collected and indexed, which serves as the main asset for the computer program to analyze the available photography ideas on the Internet. Next, we explain our *decomposition* and *composition* strategies, which are the most significant parts of the design. After that, our *matching* stage will be presented as the concluding stage in the photo taking process. We provide qualitative results as we describe the method, and experimental results afterward to show how each part of our framework compares with the state of the art.

## 2 Related Work

Existing work closely related to ours is categorized into four groups, covering aesthetic quality assessment, composition, portrait, and on-site feedback, respectively.

### 2.1 Aesthetic Quality Assessment

Books on professional photography (Lauer and Pentak, 2011; Valenzuela, 2012; Krages, 2012; Matthew, 2010; Rice, 2006; Smith, 2012; Valenzuela, 2014; Grey, 2014) guide people to master skills of taking striking photos practically in various situations. However, learning through them takes a lot of time and practice. Existing technical approaches attempt to automatize this process, but they are limited and mostly focused on offline evaluation or active manipulation of the photos after they are taken. Basic image aesthetics and composition rules in visual art (Lauer and Pentak, 2011; Valenzuela, 2012; Krages, 2012), including geometry, color palette, and the rule of thirds, have first been studied computationally by Datta et al. (2006) and Ke et al. (2006) as visual aesthetic features.

Luo and Tang (2008), and Wong and Low (2009) attempted to leverage a saliency map method, and considered the features of the salient parts, because more appealing parts of an image often reside in the prominent region. Marchesotti et al. (2011) showed that generic image descriptors was useful to assess image aesthetics, and built a generic dataset for composition assessment

- the Aesthetic Visual Analysis (AVA) dataset (Murray et al., 2012). Deep learning based approaches (Lu et al., 2015; Mai et al., 2016; Talebi and Milanfar, 2018; Liu et al., 2018) exploit customized architectures to train image aesthetic-quality models with annotated datasets, and the outcome is an estimation for actual (average) aesthetic rating of an image.

### 2.2 Image Auto-Composition and Re-Composition

Auto-composition systems (Bhattacharya et al., 2010, 2011) actively manipulate and then re-compose the taken photo for a better view. Cropping techniques (Suh et al., 2003; Santella et al., 2006; Stentiford, 2007) separate the region of interest (ROI) by the help of a saliency map, an eye fixation, basic aesthetic rules (Zhang et al., 2005), or visual aesthetics features in the salient region (Park et al., 2012; Samii et al., 2015; Yan et al., 2013). Warping (Liu et al., 2010) is another type of re-composition that represents image as a triangular or quad mesh, to map the image into another mesh while keeping the semantics and perspective unchanged. Also, R2P (Chang et al., 2015) detects the foreground part in the reference image. Then, it re-targets the salient part of the image to the best-fitted position using a graph-based algorithm.

Furthermore, patch re-arrangement techniques mend two ROIs in an image together. Pure patch rearrangements (Barnes et al., 2009; Cho et al., 2008; Pritch et al., 2009) detect a group of pixels on the border of the patch and match this group to the other vertical or horizontal group of pixels near the patched area. Also, cut-and-paste methods (Bhattacharya et al., 2011; Zhang et al., 2005) remove the salient part, and re-paint the foreground with respect to the salient part and the borders, and then paste it to the desired position in the image. Another auto-composition system, seam carving (Guo et al., 2012; Li et al., 2015b), replaces useless seams.

### 2.3 Assessment of Portrait Aesthetics

While there exist prior studies on image aesthetics assessment, few considered portrait photography in depth, despite the fact that the portion of portrait genre is very high in the photography domain. Even in this domain, prior works have not explored a novel method to solve the problem in photographic portraiture, rather than combining and using well-known features or modifying trivial ones to apply in the facial domain. We categorize prior works into two main groups: rule-based evaluation models (Khan and Vogel, 2012; Males et al.,

2013) exploit known photography rules to assess portraits, and facial evaluation models (Xue et al., 2013; Lienhard et al., 2014, 2015b,a; Redi et al., 2015) use visual features on face like smiling, age, gender, etc.

In a rule-based evaluation model, Khan and Vogel (2012) show that a small set of face-centered spatial features extend the rule of thirds and perform better than a large set of aesthetics-related features. Actually, their dataset containing 500 images from Flickr scored by 40 people is limited for a general conclusion. Their aesthetic features, especially spatial features, are close to well-known photography rules which were widely investigated before.

Males et al. (2013) explore the aesthetic quality of head-shots by means of some famous photography rules and low-level facial features. More specifically, sharpness and depth of field, the rule of thirds as a composition rule, contrast, lightness, hue counts and face size are exploited as their fundamental features. Unfortunately, the experimental results of the paper are limited, and it is hard to conclude for general cases. Xue et al. (2013) study the design inferring portrait aesthetics with appealing facial features like smiling, orientation, to name but a few. Similarly, Harel et al. (2006) 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. Then, they calculate the standard deviation and the main subject coincidence of the saliency map.

The other facial evaluation models (Lienhard et al., 2014, 2015b,a) use well-known low-level aesthetic features such as colorfulness, sharpness, and contrast, as well as high-level face-related 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. (2015) 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. Using a dataset derived from AVA (Murray et al., 2012), they exploit a high-dimensional feature vector including aesthetic rules, biometrics and demographic features, image quality features, and fuzzy properties. Based on lasso regression output, eyes sharpness and uniqueness features have the highest rank for a good portrait.

#### 2.4 On-site Feedback on Photographic System

An aesthetic assessor may find a metric to evaluate aesthetic quality of an image, but the way it conveys this information to photographer is also crucial. Because, an amateur photographer probably has no idea about how

to improve the image composition. That is why providing meaningful feedback to enhance the next shots and not just image aesthetic assessment is one of our main intentions.

Giving feedback on a photographic system firstly has been introduced by Joshi et al. (2011), as they suggest a real-time filter to trace and aesthetically rate the camera shots, and then the photographer retakes a better shot. On-site composition and aesthetics feedback system (Yao et al., 2012; Li et al., 2015a) helps smartphone users improve the quality of their taken photos by retrieving similarly composed images as a qualitative composition feedback. Also, it gives color combination feedback for having colorfulness in the next photo, and outputs the overall aesthetic rating of the input photo as well. OSCAR (Yao et al., 2012) is assumed to fulfill future needs of an amateur photographer, but giving such a feedback may be unrelated or unrealistic to the user, and also it is restricted to a small database in terms of coverage, diversity, and copyright.

Xu et al. (2015) suggest using a 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 a photo from another position. More recently general aesthetic techniques including perspective-related techniques (Zhou et al., 2016) and triangle technique (He et al., 2018) are exploited to retrieve proper images as an on-site guidance to amateur photographers, but they are limited to basic ideas in photography while many aspects such as human pose or scene content are ignored, and these methods just try to retrieve similar photos to query photo having perspective or triangles, but the retrieved results may not be necessarily useful to amateur photographer.

### 3 The Dataset

The most valuable resource used by the computer program developed in this work is the collected dataset because it contains a large number of innovative photography ideas from around the world. We have attempted many photo-sharing websites for photography purposes including Flickr, Photo.net, DPChallenge, Instagram, Pinterest, and Unsplash. However, none of them properly cover several categories such as full-body and upper-body in portrait photography as well as urban in landscape photography. The process of searching, collecting, and updating the dataset is time consuming and taxing, hence, automating this process is quite helpful.

### 3.1 Portrait and Landscape Dataset

The dataset is gradually collected 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 by marketing their photographs. To get the file list and then the images sorted by rating, we have implemented a distributed multi-IP address, block-free Python script that collects the photos having tags such as portrait, pose, human, person, woman, man, studio, model, fashion, male, female, landscape, nature and so on.

Nearly half a million images for the current dataset have been collected. The dataset has diverse photography ideas specially for the aforementioned portrait categories (full body, upper body, facial, group, couple or any two-body, side-view, hand-only, and leg-only) and landscape categories (nature, urban, etc.) from highly-rated images taken by mostly photography enthusiasts and professionals. Figure 3 illustrates some sample images from the dataset including portrait categories such as facial, full-body, seated, upper-body, no-face, side-view, group, hand-only and leg-only as well as landscape categories such as nature, plain, water, sky, and trees.

Figure 12 in Section 7.1 shows the logarithmic distribution of the view counts, the distribution of the ratings, the logarithmic distribution of the vote counts, and the logarithmic distribution of the favorite counts of the dataset respectively. The probability of a bin represents the frequency of the images which reside in the interval starting from the current bin threshold to the next bin threshold divided by the total number of images. As a result, more than 90% of the images were viewed more than 100, and nearly half of the images in the dataset had a rating between 40 to 50, which is a high rating.

### 3.2 Automating Dataset Annotation

The number of images in the dataset (about half a million by the end of 2017) is large. While we have manually annotated around 10% of the dataset for training, verification, and testing purposes, to annotate the rest, we leverage multiple highly-accurate detectors to automate and accelerate the process. However, the accuracy of the annotation is not perfect, but it is high enough for getting feedback by our aesthetics-based image retrieval from the dataset. Also, the redundancy across our designed detectors makes the annotation process more accurate that we will discuss later in Section 4.

To automate image categorization, we formulate the problem as a multi-class model for support vector machines (SVM). Therefore, we train an SVM model using radial basis function (RBF or Gaussian) kernel to predict image category (e.g. facial, upper-body, full-body, urban, nature, etc.) from multiple classes. Pose features (later explained in detail) are extracted to train our SVM model on a random subset of images from the dataset including 20K diverse images which are manually labeled. Figure 18 depicts the distribution of the categories with respect to the number of corresponding images in each category divided by the total number of images. Consequently, the number of images for some categories like full-body, upper-body, facial, group, two, and side-view is higher than the others to cover the diversity of our image retrieval.

Instead of directly labeling photography ideas, we detect all detectable semantic classes in the scene using the object detectors and the scene parsers. In fact, we believe the scene snapshot captured by camera consists of various static and dynamic objects that constructs the constellation of the scene. The object detector partitions each shot into several boundaries (not necessarily segments) with a detection probability for each, while these boundaries can also have overlapping region. The scene parser predicts the potential objects available in each shot in pixel level, *i.e.* each pixel has an object label detectable by scene parser.

We enhanced the deep-learned model of object detector YOLO (Redmon et al., 2016) and scene parser PSPNet (Zhao et al., 2017) to annotate the dataset. To improve the accuracy, we have trained our purpose-driven architecture of the object detector on an extended dataset including a subset of common failure cases (CFC) from the dataset with MSCOCO dataset (Lin et al., 2014). Similarly, we have trained a customized architecture of the scene parser on CFC with ADE20K dataset Zhou et al. (2017a) as an augmented training set.

After getting all automatized annotations of the images in the dataset, we just keep those detected objects having an area greater than the 1.15% of the image area. The probabilities of the highly-repeated semantic classes in the dataset (*i.e.* the frequency of the semantic class divided by the total number of images) are shown in Figure 16, while we have removed “person” (probability=0.9) and “wall” (probability=0.78) from the figure because they are dominant semantic classes in most of the images. Definitely having diverse semantic classes with high frequency in the dataset makes the proposed recommendations with respect to the query shot more helpful. After collecting the dataset, filtering unrelated images including low-quality or nudity, and



Fig. 3: Sample images from the collected dataset. First row: facial, full-body, seated, and upper-body. Second row: no-face, side-view, group, hand-only, and leg-only. Third row: nature, plain, water, sky, and trees.

auto-annotating them, we start indexing to extract the aesthetics-related information from them to accelerate the retrieval process.

#### 4 Photo Decomposition

To suggest a better composed photo to the user, we decompose query image from camera (*i.e.* *shot*) into composition ingredients called aesthetics-related information. This information includes high-level features (such as semantic classes, photography categories, human poses, subject gender, photo tags, and photo rating) as well as low-level features (such as color, texture, and etc). To accelerate the retrieval process from the dataset based on query image, we perform the decomposition procedure on all images in the dataset as an offline process, called *indexing*, shown as black arrows in Figure 2. We construct the composition model (CM) after indexing the whole dataset. If new images join the dataset, we index them and update our CM. In the *searching* step shown as red arrows in Figure 2, we decompose query image, and compare with our CM. Then, we retrieve the highly-ranked photos from the dataset based on the decomposed values of the shot and user-specified preferences (USP).

Through this section, we describe our integrated object detector (IOD) to determine semantic classes in query image more comprehensively and more accurately than a single object detector. Also, our category detector (CaDe) specifies the photography genre and style. Furthermore, our artistic pose clustering (Ar-Pose) extracts pose information specially for portrait photography. The other properties such as rating, tags,

and gender in the shot are extracted from the image descriptor as a JSON file. For the low-level features, we collect all 4096 generic descriptors via public pre-trained CNN model (Chatfield et al., 2014) on ImageNet (Deng et al., 2009) and the conventional features of Mitro (2016)’s method as shown in the following equation. Note that there is no limit to collect any other aesthetics-related information from query image to extend our work depending on image style or functionality.

$$F_{I,vgg} = [f_{I,1}^{vgg} \ f_{I,2}^{vgg} \ \dots \ f_{I,4096}^{vgg}]^T, \quad (1)$$

where  $F_{I,vgg}$  is a vector containing generic features of image  $I$ , and  $f_{I,i}^{vgg} \ \forall i$  is  $i$ -th generic feature. The superscript “ $T$ ” represents the transpose of the vector/matrix. Also, we extract available statistical data via the image properties including rating, view counts, and gender. Then, we similarly have them as follows:

$$F_{I,stat} = [f_{I,1}^{rating} \ f_{I,2}^{views}]^T, \quad (2)$$

$$F_{I,gender} = [f_{I,1}^{male} \ f_{I,2}^{female} \ f_{I,3}^{unknown}]^T, \quad (3)$$

where  $F_{I,stat}$  is a vector containing the statistical data of image  $I$  including its rating  $f_{I,1}^{rating}$  and its view counts  $f_{I,2}^{views}$ . Furthermore,  $F_{I,gender}$  is a vector containing the gender specification of image  $I$  represented by  $[1 \ 0 \ 0]$  as male,  $[0 \ 1 \ 0]$  as female, or  $[0 \ 0 \ 1]$  as unknown.

##### 4.1 Integrated Object Detectors (IOD)

Deep learning based models help computer vision researchers map from an unbounded correlated data (*e.g.*



an image) to a bounded classified range (*object labels*), but there are many restrictions to exploit them for applied problems. As mentioned before, there is *no limit* to innovation in visual arts. 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. As the number of ideas increases, mean average precision (MAP) falls abruptly at the rate of  $O(\frac{1}{n})$ . Also, manual idea labeling of a large dataset is costly in terms of computational time and available budget (Farhat and Tootaghaj, 2017; Farhat et al., 2016a; Tootaghaj and Farhat, 2017; Farhat et al., 2016b).

To tackle the problem of classifying to a large number of ideas, we detect as many components 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 as its high-level features. We improve the detection accuracy by training our customized object detector on an augmented dataset including a subset of common failure cases (CFC) from the 500px dataset with MSCOCO dataset (Lin et al., 2014). We train our scene parser on our CFC with ADE20K dataset (Zhou et al., 2017a) as an extended training dataset, and also our human pose estimator is trained on our CFC with MSCOCO dataset (Lin et al., 2014) and MPII dataset (Andriluka et al., 2014).

We start from state-of-the-art deep-learned detectors, YOLO (Redmon et al., 2016), PSPNet (Zhao et al., 2017) and RTMPPE (Cao et al., 2017), and we extend, improve and integrate them for our purpose. YOLO network partitions the query photo into several bounding boxes predicting their probabilities. Pyramid scene parsing network (PSPNet) uses global context information through a pyramid pooling module, and 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 by means of two sequential prediction process exposing the part confidence maps and the vector fields.

Figure 4 illustrates a small subset of common failure cases (CFC) across object detector (YOLO), human pose estimator (RTMPPE), and scene parser (PSPNet). Occasionally RTMPPE misses at facial photos to detect human parts like neck in close-up photos, and it is not very accurate at “two” or “group” categories to associate parts overlapping. The non-person detection of YOLO under 30% probability is sometimes not reliable. PSPNet detection is partially not accurate enough at photos with many objects, as it partitions the photo into small chunks and it never considers overlapped area. Generally, to improve the accuracy of the detectors, we have

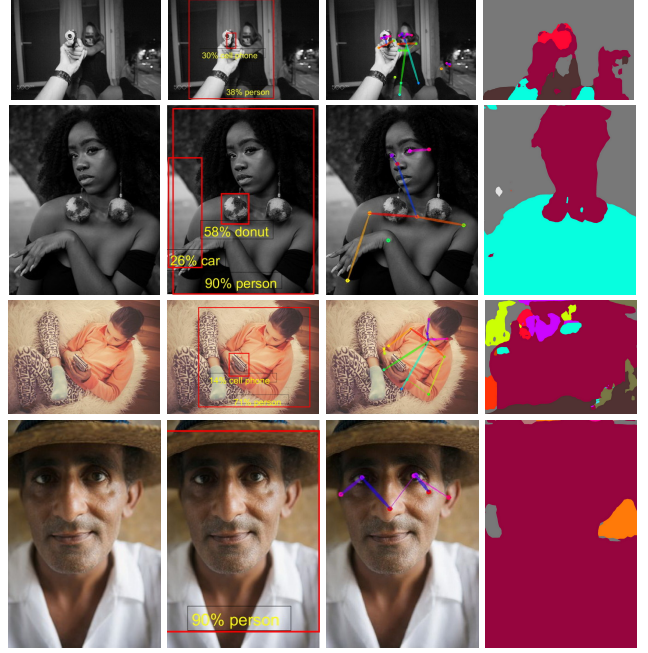


Fig. 4: Sample failed cases generated by detectors developed in the computer vision field. In each row, the images are the original, the YOLO result, the RTMPPE result, and the PSPNet result, respectively.

changed the transformation parameters of the architecture such as maximum rotation, crop size, scale min and max. Because higher rotation and bigger portrait are more important in our work.

For comparison with YOLO, we use the regular MAP on all intended objects. Table 1 in Section 7.2.1 shows the MAP and the average accuracies of some objects by our trained model versus pre-trained YOLO model. For comparison with RTMPPE, we measure MAP of all body parts (left and right are combined) as mentioned in DeeperCut (Insafutdinov et al., 2016). Table 2 in Section 7.2.2 compares MAP performance between ours and RTMPPE on a subset of testing images randomly selected from the 500px dataset. For scene parsing evaluation, we measure pixel-wise accuracy (PixAcc) and mean of class-wise intersection over union (CIoU), where the performance values of our trained scene parser are 78.6% PixAcc and 42.5% CIoU better than PSPNet with 101-depth ResNet (74.9% PixAcc and 40.8% CIoU) listed in Table 3 of Section 7.2.3.

Figure 5 illustrates four different qualitative results, where YOLO object names are shown in a red rectangle with a probability, RTMPPE poses are shown as a colorful connection of skeleton joints, and PSPNet scenes are colored pixel-wisely based on the pixel codename.

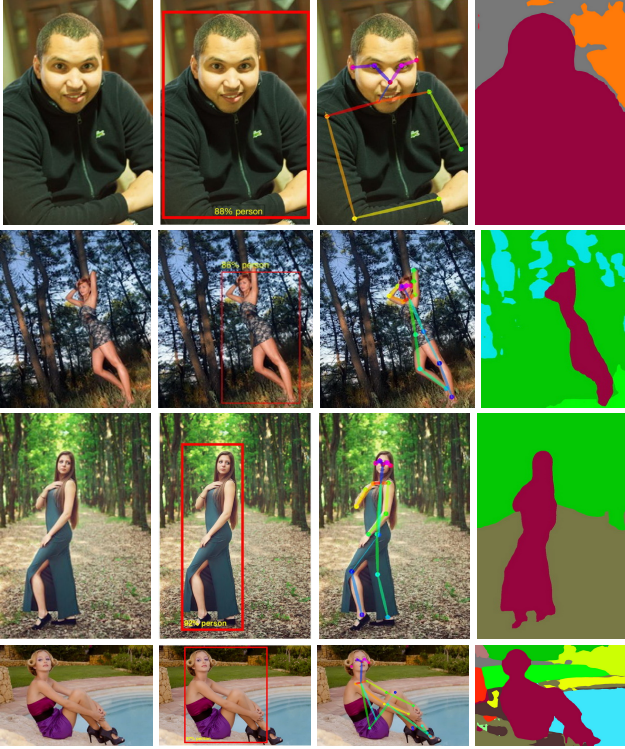


Fig. 5: Qualitative results generated by the enhanced integrated object detector (IOD) show refined samples after training. In each row, the images are the original, the object detector result, the pose estimator result, and the scene parser result, respectively.

#### 4.1.1 Value Unification

To work more conveniently on the outputs of our customized detectors in next steps, we need to unify the outputs in terms of pixel-level tensors. Our object detector outputs MSCOCO object-IDs among 80 categories (from 1 to 80). We define their scores as the minus logarithm of their NOT probability ( $-\log(1-p)$ ) for each pixel of the image. The object-ID and its score for each pixel is represented as a  $m \times n \times 2$  tensor. Also, our scene parser outputs ADE20K object-IDs among 150 categories (from 1 to 150), and the object-ID with its score for each pixel of the image is represented as a tensor. Similarly, our human pose estimator gives 18 anatomical part IDs with their scores as a tensor. Thus, for any image ( $I_{m \times n}$ ) we have:

$$T_{m \times n \times 2}^{I,od} = [t_{i,j,k}^{I,od}], \quad (4)$$

$$t_{i,j,1}^{I,od} = C_{i,j}^{I,id}, t_{i,j,2}^{I,od} = -\log_2(1 - p_{i,j}^{I,od}),$$

$$T_{m \times n \times 2}^{I,sp} = [t_{i,j,k}^{I,sp}], \quad (5)$$

$$t_{i,j,1}^{I,sp} = A_{i,j}^{I,id}, t_{i,j,2}^{I,sp} = -\log_2(1 - p_{i,j}^{I,sp}),$$

$$T_{m \times n \times 2}^{I,pe} = [t_{i,j,k}^{I,pe}], \quad (6)$$

$$t_{i,j,1}^{I,pe} = J_{i,j}^{I,id}, t_{i,j,2}^{I,pe} = -\log_2(1 - p_{i,j}^{I,pe}),$$

where  $I$  is an input image,  $m$  is the number of rows,  $n$  is the number of columns in the image,  $T^{I,od}$  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  $(i, j)$ ,  $T^{I,sp}$  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  $(i, j)$ ,  $T^{I,pe}$  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 500px dataset in indexing step, we combine these unified results in terms of the semantic classes, their coordinates, and their scores (or probabilities). The number of the detectable classes is 210 semantic objects by merging MSCOCO (80 categories) and ADE20K (150 categories) objects and deduplicating 20 semantic objects (such as person and sky). 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 a small full-body person in the image is poor, but it detects big limbs of the body (as a person label) well. RTMPPE detection for occluded bodies is poor but the detection for a full-body person is acceptable. Also, PSPNet detection for objects, not a person, is relatively good compared to others.

#### 4.1.2 Hysteresis Detection

To expand our framework coverage, using our available detectors, we detect the potential objects in the image. Our detector's integration scheme has LOW (usually with the probability less than 0.09) and HIGH (usually with the probability higher than 0.44) thresholds for each binary (object,detector). These thresholds are tuned by a random set of highly-rated ground-truth images. If the average probability (score) of the pixels with object ID X in the image is higher than its HIGH threshold, there is an object ID X in the image, otherwise if the average probability (score) of the pixels with object ID X in the image is lower than the corresponding LOW threshold, there is no object ID X in the image, and we examine other objects for indexing purpose or another image for searching purpose. We call our detector's integration scheme as *hysteresis* detection.

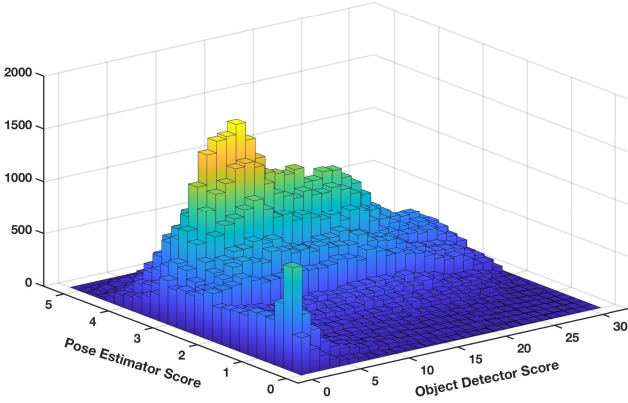


Fig. 6: The 2D histogram of the portrait images binned by the object detector and pose estimator scores.

Hysteresis detection guarantees that confidence ratio of existence of an object in an image is at least higher than a tuned detection threshold for one of the detectors. A person may be detected by one detector but not by another one. Therefore, if we do not want to miss any shared detectable object, we consider the union of the detectors, *i.e.*, if any detector outputs the object with HIGH probability, the object is considered inside the image. For example, if we use the hysteresis detection of the 500px dataset for “person” object, 90.5% (400K+) of the images pass. The 2D histogram of the portrait dataset in Figure 6 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 these detectors to use the coverage of the dataset more precisely, because we are unifying the detection results for a more broad range of ideas rather than intersecting them to have a more confident narrow range of ideas.

To tune the boundaries, we conduct the experiments and consider detection probability and normalized area as two features of a dominant object from the object detector, and detection score and normalized area as two features of a dominant object from the pose estimator. From the 2D ROC results in Figure 13, we infer that *probability* = 45% for the object detector and *area* = 10% for the pose estimator are the optimized cut-off thresholds to decide about the existence of a person in the image.

#### 4.1.3 Object Importance

To prioritize the prominence of the objects in the image, we seek to use the importance map of the objects, because the subject of the image should be more important even if its detection probability is lower. To rank

the order of the objects, we exploit the max score multiply by a saliency map ( $S$ ) features with our centric distance ( $D$ ) feature to get our weighted saliency map ( $W$ ).

$$W^I(i, j) = \max \left( T_{*,*,2}^{I,od}, T_{*,*,2}^{I,sp} \right)_H S^I(i, j) D^I(i, j), \quad (7)$$

$$D^I(i, j) = \frac{1}{K} e^{-\| [i,j] - c^I \|_k}, \quad (8)$$

$$c^I = \frac{\sum_{i,j} S^I(i, j) \cdot [i, j]}{\sum_{i,j} S^I(i, j)}, \quad (9)$$

where  $W^I(i, j)$  is our weighted saliency map pointwisely for image  $I$ ,  $\max(\cdot)_H$  operation is a hysteresis max on the 2nd plane of the tensors (score matrix),  $S^I(i, j)$  is the saliency map of image  $I$  inspired by Itti et al. (1998), and  $D^I(i, j)$  is our centric distance feature of image  $I$ ,  $K$  is a tunable constant equal to  $\sum_{i,j} e^{-\| [i,j] - c^I \|_k}$  for image  $I$ , the binary value  $c^I$  is the center of the mass coordinate, and  $\|\cdot\|_k$  is the  $k$ -th norm operator where  $k = 1$  in our experiments.

Our weighted saliency map makes the detected objects prioritized, because we sum up the scores from the semantic classes, and we end up with a total score for each semantic class. The output of this step is a weighted vector of detected semantic classes (undetected object has zero weight) in the query image. We show it as the following vector where the elements represent the importance (normalized as a probability) of the corresponding object in the image:

$$F_{I,iod} = [f_1^{imp} \ f_2^{imp} \ \dots \ f_{210}^{imp}], \quad (10)$$

$$f_k^{imp} = \frac{\sum_{\forall objID(i,j)=k} W(i, j)}{\sum_{\forall i,j} W(i, j)}, \quad (11)$$

$$\forall k \in \{1, 2, \dots, 210\},$$

where  $f_k^{imp}$  is the importance (*imp*) value of  $k$ -th object which is the summation of the weighted saliency of the pixel  $(i, j)$  with ID  $k$ , *i.e.*,  $objID(i, j) = k$ , and  $F_{I,iod}$  is the importance vector of all objects. As we index the 500px dataset by integrated object detectors, the union of the objects by the detectors is recognized. The distribution of the highly-repeated semantic classes in the 500px dataset is shown in Figure 16 in Section 7.2.5, where the “person” and “wall” were removed from the figure as mentioned before.

#### 4.2 Category Detector (CaDe)

The photo categories in portrait include two (couple or two people), group (more than two people), full-body, upper-body, facial, side-view, faceless, headless, hand-only, and leg-only, which are ten classes. In landscape

photography, there are sea, mountain, forest, cloud, and urban, which are five classes. While we focus on portrait and landscape photography genres, we believe that this work can be extended to other genres as well.

Knowing the photo genres and categories help our framework guide the photographer more adequately because it retrieves better-related results based on the photographer preferences. The downside can be the low coverage or a limited number of contents on the leaves of this hierarchical tree of the photo styles, but the comprehensive dataset addresses this potential issue.

#### 4.2.1 Top-down Hierarchical Clustering

To distinguish a portrait from a landscape photo, the number of people in the image is estimated by the max (union) number of person-IDs higher than their corresponding HIGH thresholds across the detectors in integrated object detector (IOD). If the score for detecting a person is lower than a LOW threshold for all detectors in IOD, there is no person in the image. Then, if there is a water-like, mountain-like, plant-like, cloud-like, or building-like object in the image with total area higher than 26.5% (optimized for landscape), the landscape category will be recognized as well. Otherwise, if there is a person in the image, detecting the right category by a decision tree is one heuristic approach, *i.e.*, 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 determined after trying tens of random images as ground truth, because pose estimator 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. But, we do not follow this approach, because this hierarchical approach for portrait images is not very accurate, as some leaves of its decision tree have some correlations like full-body and group categories. Also, the coverage is not fair, since upper levels like full-body attract most of the photos, and the rest will remain for the lower levels.

#### 4.2.2 Portrait Multi-class Categorization

Our more efficient and accurate approach to automate portrait categorization formulates the problem as a multi-class model for support vector machines (MCMSVM). The inputs are our feature vectors and the corresponding class labels, and the trained MCMSVM is a fully trained multi-class error-correcting output codes (ECOC) model, while we are using 10 portrait categories or

unique class labels, it needs 45 ( $= 10(10 - 1)/2$ ) binary SVM models with radial basis function (RBF or Gaussian) as its kernel and a one-vs-one coding design. We have annotated 5% (about 25K+) of portrait photos uniformly selected at random from the dataset as the ground truth of the portrait categories. Then, we train an MCMSVM with the feature vectors and the corresponding labels of 80% (about 20K) of our ground truth and leave the rest for testing our MCMSVM. Our feature vector for each photo includes 40 different features as follows:

- General MAX: (1,2) max scores for detected people from IOD, (3,4) max areas for the detected people from IOD.
- Intersected Area: (5) intersected area between highly probable people from IOD, (6,7) scores of the highly probable people for each detector in IOD, (8,9) areas of the highly probable people for each detector in IOD.
- Number of people: (10,11) number of people higher than HIGH threshold for each detector in IOD, (12,13) number of people with area higher than 5% for each detector in IOD, (14) max number of people by score from IOD, (15) max number of people by area from IOD, (16) max of (14) and (15).
- Limb Features: (from 17 to 40) the limbs respectively including nose, neck, right shoulder, right elbow, right wrist, right hand, left shoulder, left elbow, left wrist, left hand, right hip, right knee, right ankle, right leg, left hip, left knee, left ankle, left leg, right eye, left eye, eyes, right ear, left ear, ears which add up to 40 features.

The output of this step for an image query is the following unitary vector that shows its category (facial, full-body, upper-body, two, group, side-view, leg, no-face, hand, and no-head) as:

$$F_{I, \text{cade}} = [f_1^{\text{facial}} \ f_2^{\text{fullbody}} \ f_3^{\text{upperbody}} \ f_4^{\text{two}} \ f_5^{\text{group}} \ f_6^{\text{sideview}} \ f_7^{\text{leg}} \ f_8^{\text{noface}} \ f_9^{\text{hand}} \ f_{10}^{\text{nohead}}], \quad (12)$$

where  $F_{I, \text{cade}}$  shows the unitary category vector of the image  $I$  by CaDe detector, and only one of the vector element is one and the rest are zero.

The mean average accuracy of our category detection is shown in Table 4 in Section 7.2.6 for the dataset images divided by various styles. The CaDe indexing of the dataset results the distribution of the portrait categories shown in Figure 18 in Section 7.2.6. Consequently, the number of photos in the categories containing full-body, upper-body, facial, group, two, and side-view is adequate.



### 4.3 Artistic Pose Clustering (ArPose)

*Posing*, one of the essential ingredients of the portrait photography, could substantially differentiate between amateur and professional shots. Having little experience in portrait photography, finding correct postures or coming up with novel poses is hard for amateur photographers. Hence, it is vital for our system to have an understanding of different poses and how to categorize them. Recently, there have been numerous efforts in the computer vision field for human pose estimation of images and videos. With the rise of deep learning models, these approaches are getting more accurate and more robust. One of the state-of-the-art algorithms for pose estimation is RTMPPE (Cao et al., 2017), where they use VGG features as an input and then exploit a two-stage CNN to find the probability of joints and their connections together. Their architecture 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 (Lin et al., 2014) and MPII (Andriluka et al., 2014) datasets. Although RTMPPE extract body joints in images, these joints are merely considered as our features for pose detection. Hence, we use two sets of features on top of RTMPPE in order to define the distance between different poses. These set of features are scale invariant, thus regardless of the scale of the human body in images, we measure the similarity of two poses. These features are defined as follows:

- **Joint to Line Distance (J2L):** Li et al. (2017) apply this distance in their action recognition system from body joints skeleton. They capture the distance of each joint from any line that connects two other joints. To have the scale invariant distance, we normalize these distances with the maximum J2L distance in each body in the picture. Having the joint  $j_l$  and the line crossing two joints,  $j_m$  and  $j_n$ , Joint to Line Distance is calculated as follows:

$$J2L(l, m, n) = 2S_{\Delta_{lmn}} / \|j_m - j_n\|_2, \quad (13)$$

where  $S_{\Delta_{lmn}}$  is the area under the triangle formed by three joints. Based on the total number of joints in each body, which is 18, and the total number of different distances is  $18 \times \binom{17}{2} = 2448$ .

- **Skeleton Context (SC):** Kamani et al. (2016, 2017) introduce a scale invariant feature applied to a skeleton matching task. Skeleton context is a polar histogram of each point in the skeleton indicating the angular and distance distribution of other points in the skeleton around that point. We benefit from the angular distribution of each point and create an  $18 \times 18$  angular matrix for each body in the image.

These features are designed to capture the relative position of each joint with respect to other points, hence, they are used as a measure of distance between different poses. Next, we use these features to cluster images based on various poses.

#### 4.3.1 Pose Clustering

To rank each body posture in images, and find the nearest professional poses to the amateur one in the query image, we use a clustering method. The clustering method should be able to distinguish between different poses and group similar ones using the features explained in Section 4.3. In order to do so, we use two clustering algorithms, Kmeans and Deep Embedding (Xie et al., 2016). We compare the result of these two clustering on this task. In order to do the clustering, we first need to determine the number of clusters. There are several heuristic methods to estimate the optimal number of clusters for each dataset, including but not limited to *elbow* and *silhouette* methods. In this clustering task, having too many clusters would diminish the novelty and diversity of the results, in a sense that it tries to have samples as close as possible to one cluster. On the other hand, keeping the number of clusters low would affect the quality of clustering, such that irrelevant poses might appear in the same cluster. The result of our experiment using elbow method shows that the optimal number of cluster heads is around 10-15 as depicted in Figure 19 in Section 7.2.7.

After finding the number of clusters, we set up two clustering algorithms, namely, Kmeans and Deep Embedding Clustering (DEC). As for the Kmeans, the only parameter that we should set is the number of clusters, but in DEC we should setup the auto-encoder network in addition to the number of clusters. As suggested by Xie et al. (2016) and tested by ourselves, the network with 4 layers of encoder consisting of 500, 500, 2000, and 10 neurons in each unit performs astonishingly well on the clustering task of different supervised datasets including but not limited to MNIST (LeCun et al., 1998), STL (Coates et al., 2011), and REUTERS (Lewis et al., 2004). Although DEC works great on these supervised datasets, it has not been tested on an actual unsupervised dataset, simply because there is not a gold standard to evaluate the performance on those datasets. However, visual data like the unsupervised portrait dataset reveals how these algorithms perform, based on human vision evaluation of the output. Hence, we compare the results of this deep model for clustering with the base clustering algorithm, Kmeans.

In Kmeans, to define the probability that each sample is in the cluster or the degree to which each sample

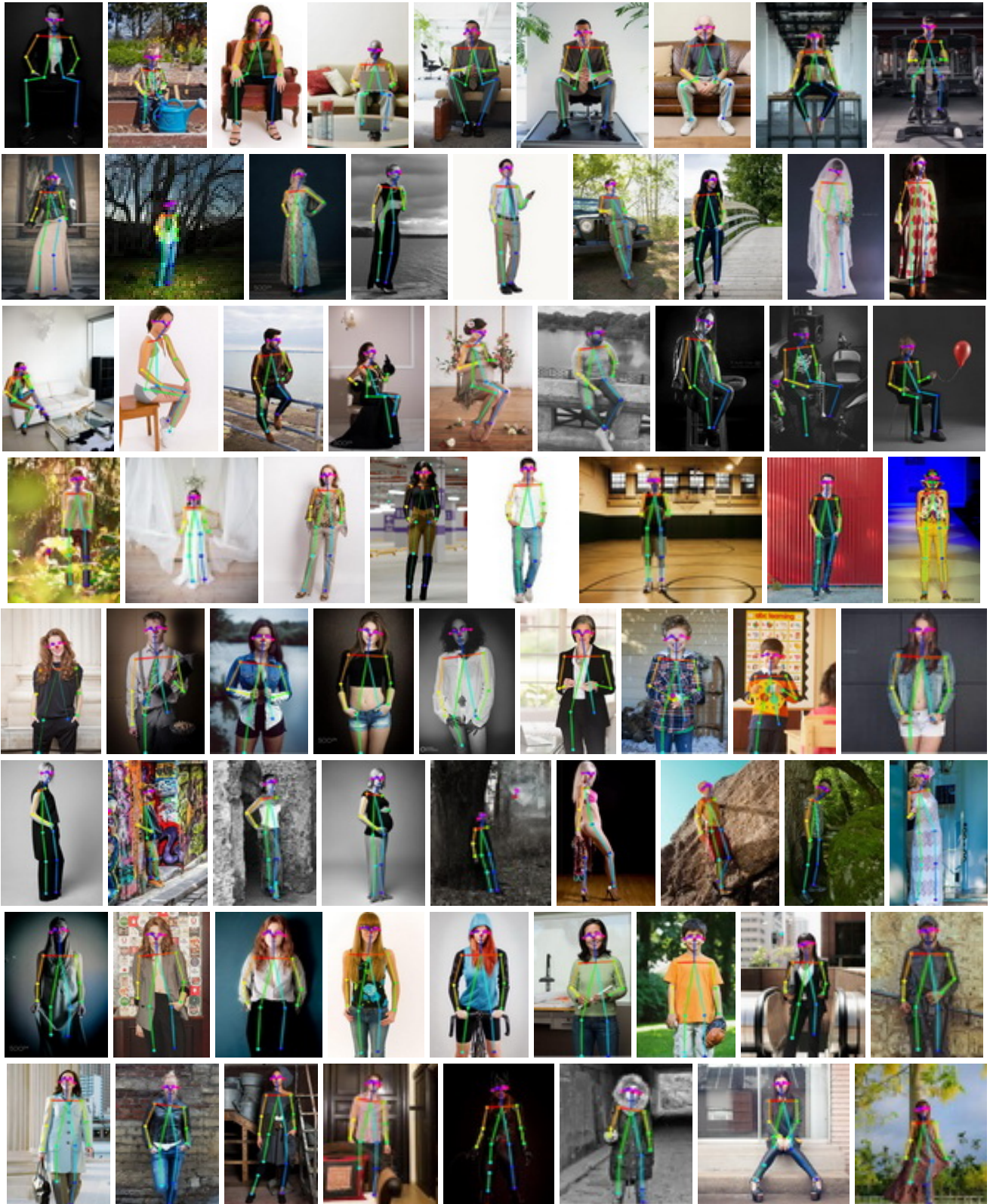


Fig. 7: First 8 clusters derived from our algorithm on the portrait dataset. Each row represents the top poses of each cluster fitted in a line.



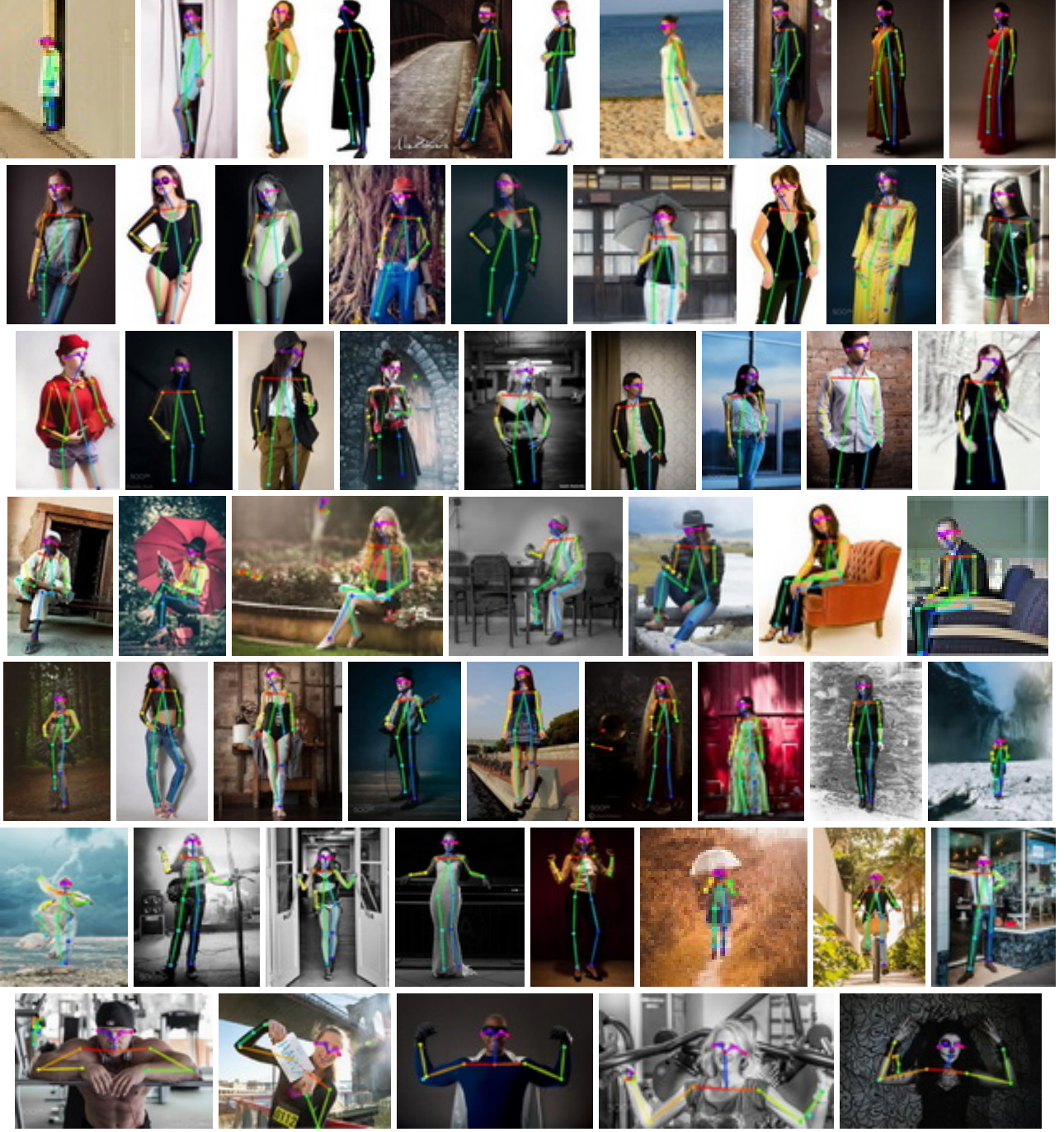


Fig. 8: The rest of 7 clusters derived from our algorithm on the portrait dataset. Each row represents the top poses of each cluster fitted in a line.

belongs to a cluster, we use the same quantity in fuzzy C-means clustering (Dunn, 1973):

$$q_{ij} = \frac{1}{\sum_{k=1}^K \left( \frac{\|x_i - c_j\|_2}{\|x_i - c_k\|_2} \right)^{\frac{2}{m-1}}}, \quad (14)$$

where  $x_i$  is the sample,  $c_j$  is the center of the cluster  $j$ , and  $m$  is positive real number greater than 1 which

defines the smoothness of the function.  $q_{ij}$  represents the probability or degree to which each sample belongs to a cluster. Also, DEC has a similar quantity defined by Xie et al. (2016), using Student's t-distribution:

$$q_{ij} = \frac{(1 + \|z_i - c_j\|_2^2/\alpha)^{-\frac{\alpha+1}{2}}}{\sum_{j'} (1 + \|z_i - c_{j'}\|_2^2/\alpha)^{-\frac{\alpha+1}{2}}}, \quad (15)$$

in which  $z_i$  is the embedded version of  $x_i$ , and  $\alpha$  is the degree of freedom in Student's t-distribution. Using these metrics we estimate the probability that each sample belongs to a cluster.

The qualitative results of Kmeans-based clustering algorithm are depicted in Figure 7 and Figure 8, showing the top ranked poses in all fifteen major clusters, and those of DEC algorithm are depicted in Figure 9 showing the best ranked poses of first 5 clusters. The top images are ranked based on their probability, calculated as above. As shown in the figures, Kmeans clusters surprisingly better represent poses in images, that is, different pose clusters distinguish between different poses and each cluster represents visually the same pose. However, DEC fails to accomplish the goal of clustering task based on poses of human subjects. Since the input features are intelligently chosen to be related to the goal, the input space is linearly separable, however, the result of the DEC shows information loss in the autoencoder. We tried the Kmeans algorithm with PCA to reduce the dimension of the input space to 10 (as it is in the output of the autoencoder in DEC), and still the results of the Kmeans surpasses DEC's. Through that, we successfully categorize the input portrait image and retrieve similar poses close to the query or novel ideas in that pose category based on the probability of the poses.

#### 4.4 Construction of Composition Model

In order to index all photos in our dataset, we decompose their values and construct our composition model (CM). If more photos are added to the dataset, we execute the decomposition for them, and update our composition model. In fact,  $F_{I_i,vgg}$ ,  $F_{I_i,iod}$ ,  $F_{I_i,cade}$ ,  $F_{I_i,arpose}$  and other aesthetic information vectors for all ( $\forall i$ ) images are calculated and appended to corresponding matrices respectively including generic feature matrix  $M_{vgg}$ , integrated object detector matrix  $M_{iod}$ , category detector matrix  $M_{cade}$ , artistic pose detector matrix  $M_{ap}$ , statistics matrix  $M_{stat}$  and gender matrix  $M_{gnd}$ . Algorithm 1 describes different steps of our decomposition method precisely to make each row of our composition model. We have:

$$\begin{aligned} &\forall feat \in \{vgg, iod, cade, ap, stat, gnd\}, \\ &\forall i \in \{1, \dots, N\}, \\ &F_{I_i} = \begin{bmatrix} F_{I_i,vgg} \\ F_{I_i,iod} \\ F_{I_i,cade} \\ F_{I_i,ap} \\ F_{I_i,stat} \\ F_{I_i,gnd} \end{bmatrix}, \end{aligned} \quad (16)$$

where “*feat*” is feature type from the set  $\{vgg, iod, cade, ap, stat, gnd\}$ ,  $F_{I_i}$  is the feature vector of the image  $I_i$ . Then, we compute the corresponding feature matrix.

$$M_{feat} = \begin{bmatrix} F_{I_1,feat}^T \\ F_{I_2,feat}^T \\ \dots \\ F_{I_N,feat}^T \end{bmatrix}, \quad (17)$$

$$\begin{aligned} M &= [M_{vgg} \ M_{iod} \ M_{cade} \ M_{ap} \ M_{stat} \ M_{gnd}], \quad (18) \\ &= \begin{bmatrix} F_{I_1}^T \\ F_{I_2}^T \\ \dots \\ F_{I_N}^T \end{bmatrix}, \end{aligned}$$

where matrix  $M_{feat}$  is the corresponding feature matrix containing feature vector of each image in each row. The final feature matrix  $M$  is the composition model matrix which is the concatenation of all feature matrices or equivalently all feature vectors.

### 5 Composition of Visual Elements

The goal of composition step is to retrieve related photography ideas as highly-rated photos from our collected 500px dataset satisfying the proximity to the decomposed aesthetics-related information of the query image. In fact, the input to this step is the decomposed values of the image query and user-specified preferences (USP) with our composition model. The output of this stage is a collection of well-composed images from the dataset. If we focus on portraits, we desire a feedback that contains well-posed portraits with similar semantics but better aesthetic quality w.r.t USP. The “interaction” between the subject(s) and the objects in the image is important because system's proposed composition depends on them.

As we have collected the 500px dataset containing generally well-composed images, we should dig into the dataset and look for images with “pretty” similar color, pattern, category, pose, or object constellation where the term “pretty” is framed by USP to address the user's needs and subjectivity. The existence of this professional-quality dataset makes it possible that the retrieved photos have highly-accepted photography ideas by the people. Our image retrieval system is not supposed to find images with exactly similar colors, patterns, or poses, but it finds images with better composition having similar semantic classes. Thus, the location of the movable objects does not matter, but the detected objects are important.



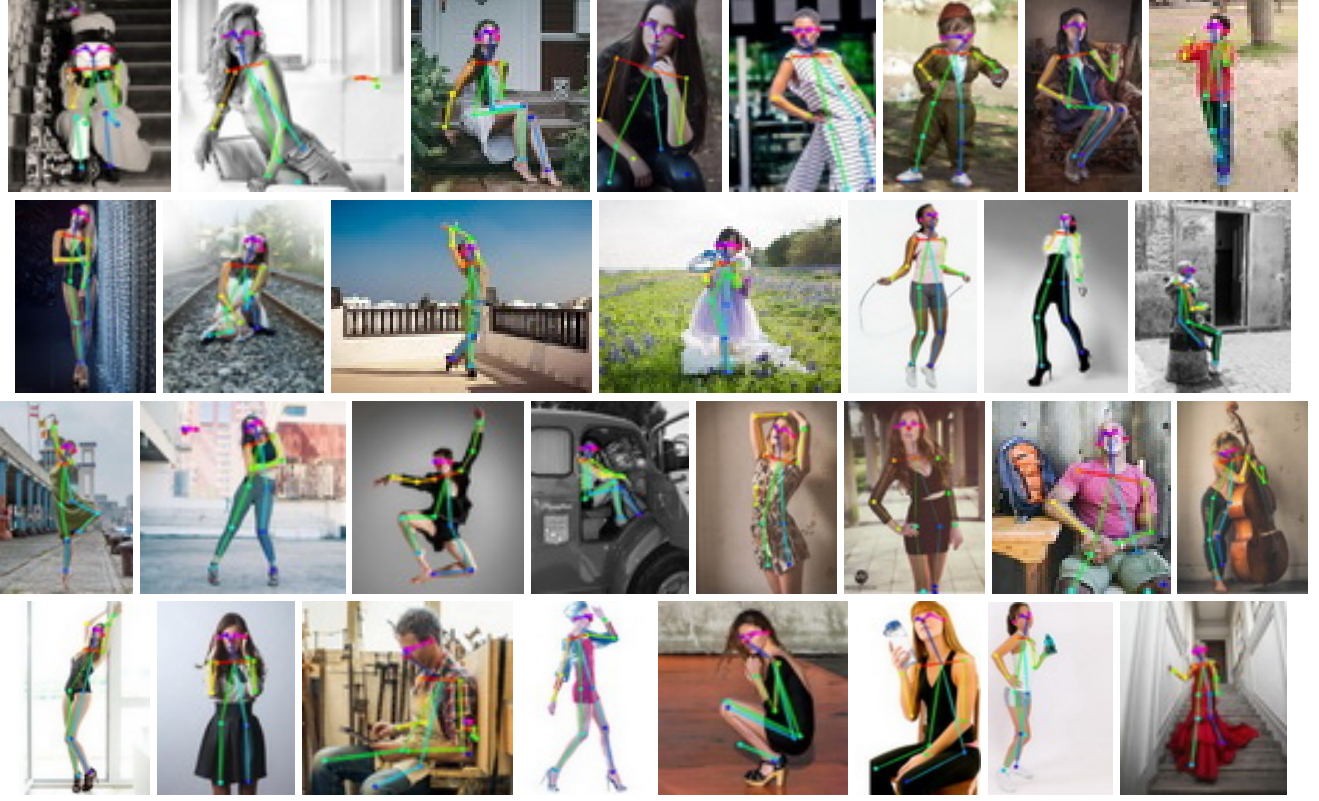


Fig. 9: Qualitative results of four clusters derived from the portrait dataset by DEC algorithm. Each row represents the top poses in each cluster. The DEC algorithm fails in clustering poses because the poses in the same cluster are not all consistent.

To look for a semantically composed version with respect to the image query, we exploit all of the decomposed values of the image query. Because we do not assume the query image taken by an amateur photographer to be well-composed enough to be the basic query for our personalized aesthetics-based image retrieval (PAIR). In fact, we just want to understand the location/setup around the subject, and then based on the scene ingredients, a well-composed image taken by a professional is proposed to the photographer.

### 5.1 Similarity Scores and Normalization

Having our composition model for all images in the 500px dataset and the query image, we first calculate the similarity score between the query image and any image in the dataset. The similarity metric is different for each detector. For generic CNN descriptors is just matrix  $M_{vgg}$  by the query vector  $F_{vgg}$  multiplication. Similarly, category detector has a matrix by vector multiplication. For integrated object detectors, we use Gaussian function after masking unrelated objects. For statistics and gender information, it is trivial as in the

following equations.

$$S_{vgg}(I, Q) = F_{I,vgg}^T F_{Q,vgg}, \quad (19)$$

$$S_{cade}(I, Q) = F_{I,cade}^T F_{Q,cade}, \quad (20)$$

$$S_{iod}(I, Q) = e^{-(\sum (F_{I,iod} \circ \text{sgn}(F_{Q,iod}) - F_{I,iod}))^2}, \quad (21)$$

$$S_{stat}(I, Q) = \text{FirstColumn}(I_{stat}) = I_{stat}^{rating}, \quad (22)$$

$$S_{gender}(I, Q) = \begin{cases} 1, & \text{if } F_{I,gender} = F_{Q,gender} \\ -1, & \text{otherwise} \end{cases} \quad (23)$$

where  $F^T$  means the transpose of  $F$ ,  $e$  is a mathematical constant about 2.72, the  $\circ$  operation is the element-wise multiplication,  $\text{sgn}(\cdot)$  is the sign function operating on each element separately. Also,  $S_{vgg}(I, Q)$ ,  $S_{cade}(I, Q)$ ,  $S_{iod}(I, Q)$ , and  $S_{stat}(I, Q)$  are similarity score values between image  $I$  and image  $Q$  respectively for generic CNN descriptors, category detection, integrated object detectors, and statistics and gender information.

The similarity score function is easily generalized to a function between two different set of images, *i.e.*,  $I_{m \times 1}$  and  $Q_{n \times 1}$  can be a set of images not only one image, and the output will be a  $m \times n$  matrix. Since we want to score the similarity between the images in

**Algorithm 1** Decomposition

---

**Input:** image query  $Q$ .  
**Output:** the feature vector of the query  $F_Q^T$ .

- 1: **procedure** *Decomposition*( $Q$ )
- 2:   Get the generic features of the query:  
 $F_{Q,vgg} \leftarrow [f_{Q,1}^{vgg} \ f_{Q,2}^{vgg} \ \dots \ f_{Q,4096}^{vgg}]^T$
- 3:   Get stat and gender features:  
 $F_{Q,stat} \leftarrow [f_{Q,1}^{rating} \ f_{Q,2}^{views}]^T$   
 $F_{Q,gnd} \leftarrow [f_{Q,1}^{male} \ f_{Q,2}^{female} \ f_{Q,3}^{unknown}]^T$
- 4:   Get the tensor of the object detector ( $od$ ):  
 $T_{m \times n \times 2}^{Q,od} \leftarrow \text{Object\_Detect}(Q)$
- 5:   Get the tensor of the scene parser ( $sp$ ):  
 $T_{m \times n \times 2}^{Q,sp} \leftarrow \text{Scene\_Parse}(Q)$
- 6:   Get the tensor of the pose estimator ( $pe$ ):  
 $T_{m \times n \times 2}^{Q,pe} \leftarrow \text{Pose\_Estimate}(Q)$
- 7:   Get the center of the mass coordinate of the query:  
 $c^Q \leftarrow \frac{\sum_{i,j} \|[i,j] - [0,0]\|_k \times S(i,j)}{\sum_{i,j} \|[i,j] - [0,0]\|_k}$
- 8:   Get the centric distance feature of the query:  
 $D^Q(i,j) \leftarrow \frac{1}{K} e^{-\|[i,j] - c^Q\|_k}$
- 9:   Get the weighted saliency ID map for the query:  
 $C^Q(i,j) \leftarrow S^Q(i,j) D^Q(i,j)$   
 $W^Q(i,j) \leftarrow \max(T_{*,*,2}^{Q,od}, T_{*,*,2}^{Q,sp}) C^Q(i,j)$
- 10:   Get the IOD feature vector as Eq. 12:  
 $f_k^{imp} \leftarrow \frac{\sum_{\forall objID(i,j)=k} W(i,j)}{\sum_{\forall i,j} W(i,j)}, \forall k \in \{1, 2, \dots, 210\}$   
 $F_{Q,iod} \leftarrow [f_1^{imp} \ f_2^{imp} \ \dots \ f_{210}^{imp}]$
- 11:   Get the category feature  $F_{Q,cade}$  as Eq. 12.
- 12:   Get the artistic pose feature  $F_{Q,ap}$  as Eq. 13.
- 13:   Get the whole feature vector:  
 $F_Q^T = [F_{Q,vgg}^T \ F_{Q,iod}^T \ F_{Q,cade}^T \ F_{Q,ap}^T \ F_{Q,stat}^T \ F_{Q,gnd}^T]$
- 14:   If  $Q$  is a dataset image, add  $F_Q^t$  to the last row of composition model matrix  $M$  in Equation 19, otherwise the output is used in Algorithm 2 to retrieve better composed photos and match with the final shot.
- 15: **end procedure**

---

the 500px dataset (say  $\mathbb{I}$ ) and an image query ( $Q$ ), in the above equations, vector  $F_{I,det}^t$  will be substituted by matrix  $M_{det}$ , and the output will be a similarity vector, while  $det$  can be any detector as follows:

$$det \in \{vgg, iod, cade, arpose, stat, gender\}$$

To make the scores uniform across various detectors, we normalize each detector score vector dividing by the summation of the whole output. Thus, each detector's similarity score is like a probability distribution over all images. We have:

$$S_{feat}^N(\mathbb{I}, Q) = \frac{S_{feat}(\mathbb{I}, Q)}{\sum_{i \in \mathbb{I}, q \in Q} S_{feat}(i, q)}, \quad (24)$$

$$feat \in \{vgg, iod, cade, arpose, stat, gender\}, \quad (25)$$

where  $S_{feat}^N(\mathbb{I}, Q)$  is a normalized similarity score matrix between each image in  $\mathbb{I}$  and each image in  $Q$  for detector  $feat$  which can be any of the mentioned detectors. Also, we combine the similarity scores across

various detectors to create a tensor of similarity scores for each pair of images from  $(\mathbb{I}, Q)$ . We have:

$$S_{d \times m \times n}^N(\mathbb{I}, Q) = [S_{vgg}^N \ S_{iod}^N \ S_{cade}^N \ S_{arpose}^N \ S_{stat}^N \ S_{gender}^N], \quad (26)$$

where  $S_{d \times m \times n}^N(\mathbb{I}, Q)$  is a tensor of size  $d \times m \times n$  where  $d$  is the number of the  $feat$  detectors ( $\|feat\|$  here is 6),  $m$  is the number of the images in  $\mathbb{I}$ , and  $n$  is the number of the images in  $Q$ .

## 5.2 User Preferences and Ranking

The user-specified preferences are a probability vector containing the weights of the decomposed vectors of the image query. We have:

$$W_{USP} = [W_{vgg} \ W_{iod} \ W_{cade} \ W_{arpose} \ W_{stat} \ W_{gender}]^T, \quad (27)$$

where  $W_{USP}$  is a  $d \times 1$  vector showing the weights of the user for each  $feat$  detector, and “ $t$ ” shows the transpose operation. Then, to retrieve the highest-ranked candidates as the results, the normalized similarity score matrix is multiplied by the USP vector. Consequently, we have:

$$V_{pref}(\mathbb{I}, Q) = W_{USP}^T S^N(\mathbb{I}, Q), \quad (28)$$

where  $V_{pref}(\mathbb{I}, Q)$  is the user's preferred image vector, and if we sort it in a descending manner with respect to the vector values, the indexes of the rows represent the highest-ranked candidates with the nearest feedbacks to the image query ( $Q$ ). The whole process of the photography idea retrieval for an input image query ( $Q$ ) is shown in Algorithm 2, and our experimental results show the quality of the results. Also, sample results for some queries with USP are shown in Figure 10.

## 5.3 Offline Indexing and Real-time Searching

Our framework consists of the flows of indexing, searching, and matching shown in Figure 2. Practically, there are many challenges in image retrieval systems (Smeulders et al., 2000; Lew et al., 2006; Datta et al., 2008) as well as in our case. To improve the performance of our image retrieval system, we compute the decomposition step for all images (*i.e.*, indexing as an offline process). Indexing procedure is lengthy for the first time, but at the time of update, it is faster because the detections for an image is real-time using GPU. Furthermore, indexing procedure for our retrieval system organizes the decomposed values of the images into categorized matrices. Consequently, the composition step is real-time using GPU, as it just extracts the decomposed values



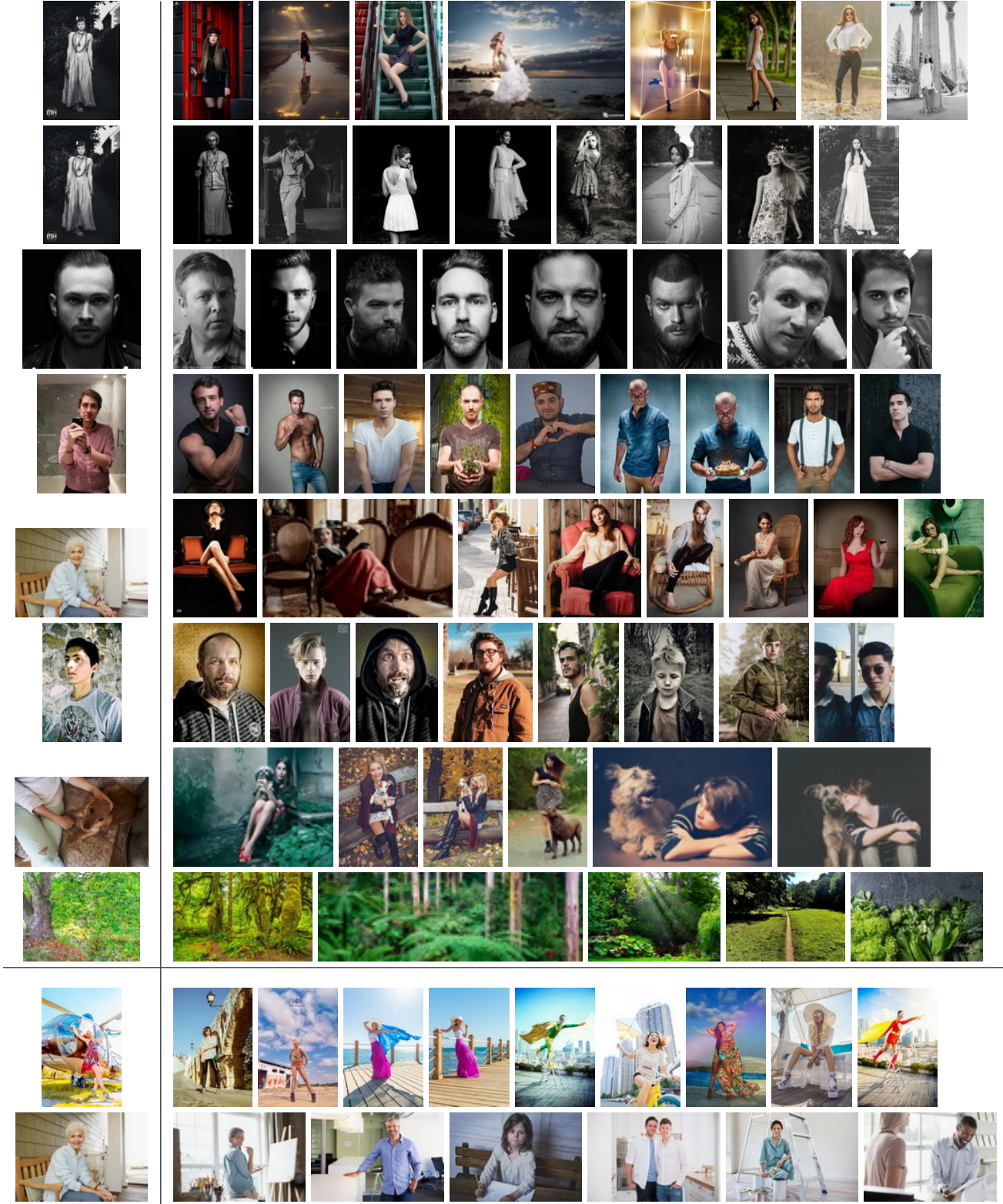


Fig. 10: Qualitative results of the composition step. Each row shows a query image and the images retrieved w.r.t. the query and user-specified preferences (USP). Last two rows show failed cases because of an undetected object (e.g. helicopter) or incompletely entered USP. The USP for each row is: 1)  $W_{cade} = W_{stat} = 0.5$ , 2)  $W_{cade} = W_{vgg} = 0.5$ , 3) like (2), 4) like (2), 5)  $W_{cade} = W_{iod} = 0.5$ , 6)  $W_{gender} = W_{vgg} = 0.5$ , 7)  $W_{gender} = W_{iod} = 0.5$ , 8)  $W_{stat} = W_{vgg} = 0.5$ , 9)  $W_{cade} = W_{vgg} = W_{iod} = 0.33$ , and 10)  $W_{vgg} = W_{stat} = 0.5$ .

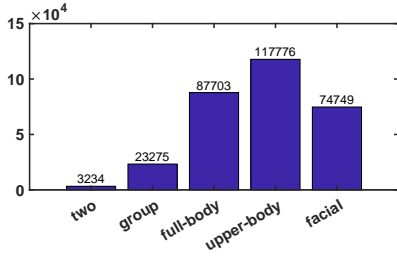


Fig. 11: The frequency of the portrait shots with respect to the highly-requested portrait categories.

of the query image similar to updating process, and then finds the target similar images, and finally retrieves best results from the dataset with respect to USP and normalized similarity score vector (Eq. 28). Having various semantic classes for portrait genre in the dataset, we have indexed the 500px dataset from previous step. Figure 11 depicts the results as the number of portrait shots for some of the highly-frequent portrait categories.

## 6 Matching

Professional photographers arrange the photograph from top to down (general to detail) step-by-step, while there are many to-do lists and not-to-do lists for 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 composition step, we retrieve the proper photography ideas given a query shot from the camera. We assume that the photographer has chosen some of the retrieved images as a *preferred style set*, and disregarded the others as an ignored set. Now we explain how to capture the proper moment of the subject in the scene, and trigger the “moment shot” for the camera.

### 6.1 Pose Shot

The best-fitting genre for matching step is portrait photography that we start with, and then we extend for general genre. The variant component in our framework is the human pose. In this scenario, it is assumed that the user has no personal preference on human pose, *i.e.*, the user has given zero weight to the ArPose detector to see various proper poses via PAIR, and then the user has selected a preferred style set from the available choices. In this case, we need a mechanism to continue guiding the user to a desired shot by tracking his/her pose via camera viewfinder.

The relative positions of the human body components (including nose, eyes, ears, neck, shoulders, el-

bows, 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 ( $J_4$ ) and left ear ( $J_5$ ) as they are on a plane of the head. Also, shoulders ( $J_6$  and  $J_7$ ) are 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.* these joints are connected as follows:

$$Pre(J_i) = J_0, i \in \{0, 1, 2, 3\}, \quad (29)$$

$$Pre(J_j) = J_1, j \in \{6, 7, 12, 13\}, \quad (30)$$

$$Pre(J_k) = J_{k-2}, \quad (31)$$

$$k \in \{4, 5, 8, 9, 10, 11, 14, 15, 16, 17\}.$$

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

$$J_i = J_j + r_{i,j} \cdot e^{i\theta_{i,j}}, i \in \{0..17\}, \quad (32)$$

where  $j = Pre(i)$  *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  $Pre(J_j)$ , 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. Similarly, having 3D pose-annotated/estimated single depth images, we can calculate the relative 3D position of the joints using spherical coordinates. Thus, we have such action boundaries for joints as follows:

$$\theta_{i,j}^{min} \leq \theta_{i,j} \leq \theta_{i,j}^{max}, j = Pre(i), \quad (33)$$

$$\phi_{i,j}^{min} \leq \phi_{i,j} \leq \phi_{i,j}^{max}, j = Pre(i). \quad (34)$$

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

$$J^k = (J_1^k, J_2^k, \dots, J_{17}^k), \quad (35)$$

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. Thus, we define the distance metric as follows:

$$D(J^k, J^l) \doteq \sum_{i=1}^{17} \|Phase(J_i^k) - Phase(J_i^l)\|_q, \quad (36)$$

where  $D(\cdot)$  is the distance operator,  $J^k$  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. Our chosen function to track phase,  $Phase(J_i)$ , is  $\sin(\theta_{i,i-1})$  which  $\theta_{i,i-1}$  (the angle between current joint and the previous one) is from  $-\Pi/2$  to  $+\Pi/2$ . Because the length of the limb may change by going far or near, but the angles between consecutive limbs matter for posing purposes.

Now, the camera 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. The problem is formulated as an integer programming problem to find the best seed among all photography ideas. Given the distance operator of two pose features explored in 36, we construct our optimization problem by maximizing the difference of the minimum distance of the ignored set and the minimum distance of the preferred style set of taken photos. Mathematically, we compute the following optimization problem subject to 33 and 34:

$$I_w = \arg \max_{\forall I_i \in I^t} \left( \min_{\forall Q_j^g \in Q^g} D(J^{Q_j^g}, J^{I_i}) - \min_{\forall Q_k^d \in Q^d} D(J^{Q_k^d}, J^{I_i}) \right), \quad (37)$$

where  $I_w$  is the wish-image,  $I^t$  is the set of taken photos,  $Q^g$  is the set of ignored ideas,  $Q^d$  is the set of preferred ideas,  $D(\cdot)$  is the distance operator in 36, and  $J^x$  is the pose for  $x$ -th image with one person in 35. The optimization problem in continuous mode (not over all taken image 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 is exponential. Further, the solution does not always give the desired shot.

## 6.2 User Favorite Shot

Given a query shot from the camera, related photography ideas have been already retrieved. Suppose that the photographer selects a *preferred style set*, denoted as  $\mathbb{C} = \{C_1, C_2, \dots, C_m\}$ , and we also have a set of shots from camera called next query shots, denoted as  $\mathbb{Q} = \{Q_1, Q_2, \dots, Q_n\}$ . The problem of finding the user favorite shot among query shots while satisfying the closest similarity score to the *preferred style set* is an integer programming. We have:

$$Q_{fav} = \arg \max_{q \in \mathbb{Q}} \sum_{j \in \{1, \dots, m\}} W_{USP}^T S^N(C_j, q), \quad (38)$$

where  $Q_{fav}$  is the favorite shot,  $\mathbb{Q}$  is the set of the query shots by camera,  $W_{USP}^t$  is the transpose of the

---

## Algorithm 2 Composition and Matching

---

**Input:** query  $Q$ , user pref.  $W_{USP}$ , and the set of the images in the 500px dataset  $\mathbb{I}$ .

**Output:** user favorite shot  $Q_{fav}$ .

- 1: **procedure** *IdeaRetrieval*( $Q, W_{USP}, \mathbb{I}$ )
  - 2:   Get  $F_Q^t$  from Algorithm 1.
  - 3:   Get the similarity score through Eq. 23, 25, and 27:  
 $S^N(\mathbb{I}, Q) = [S_{vgg}^N, S_{iod}^N, S_{cade}^N, S_{arpose}^N, S_{stat}^N, S_{gender}^N]$
  - 4:   Get the preferred image vector through Eq. 28:  
 $V_{pref}(\mathbb{I}, Q) = W_{USP}^t S^N(\mathbb{I}, Q)$
  - 5:    $Retrieved\_Indexes \leftarrow Index\_Sort(V_{pref}(\mathbb{I}, Q))$
  - 6:    $Show\_Top(Retrieved\_Indexes)$
  - 7:   Now, the user selects some of the retrieved results, and the camera takes multiple shot as  $\mathbb{Q}$ .
  - 8:   Find the favorite shot through Eq. 38:  
 $Q_{fav} = \arg \max_{q \in \mathbb{Q}} (\sum_{j \in \{1, \dots, m\}} (W_{USP}^t S^N(\mathbb{C}, q)))$
  - 9:   Take  $Q_{fav}$  as user favorite shot.
  - 10: **end procedure**
- 

user-specified preference vector, and  $S^N(C_j, q)$  is the similarity vector between query  $q$  and each photo  $C_j$  in the preferred style set of the user  $\mathbb{C}$ . The computational detail of the composition and matching steps has been explained in Algorithm 2 and inspired from (Diyanat et al., 2011; Farhat et al., 2011, 2012; Farhat and Ghaemmaghami, 2014).

Mathematically  $\mathbb{Q}$  set is not finite, or its size  $n$  is not bounded. Also, there are many constraints such as color value ranges, human pose angles, category limits, and etc. In the reality, the number of the query shots is limited, and the matching solver gradually determines and updates the user favorite shot. But the solution is not necessarily optimal, because finding the optimal shot needs the whole shot space which is impractical. The good news is that the user can follow the retrieved professional shots to optimize his/her photography adventure, and the last shot would be close enough to the optimal shot.

Our approach to giving hints to the user includes two steps: 1) defining the query shot space with dynamic parameters in the scene like movable objects or human pose, 2) finding the max over the defined space. This second step is similar to pose shot approach, and some extra parameters such as photographic lighting may be adjustable as well. The solution of the problem can give a hint to the user to make a change in his/her lighting condition, pose, or any other dynamic parameter.

## 7 Experiments

In the following sub-sections, we describe our experimental results which are categorized into different components of our method including (i) the dataset, (ii)

the decomposition step, and (iii) composition step. Furthermore, the decomposition step includes object detector, pose estimator, scene parser has multiple parts to demonstrate the effectiveness of our method compared to an state-of-the-art available approach.

### 7.1 Dataset Properties

We have collected the images in the portrait and landscape categories from 500px Website and saved them as smaller images where their highest dimension has been resized to 500 pixels. Then, we have collected available metadata for each image including the number of views, the average ratings, the number of vote clicks, and the number of favorite clicks. Both jobs are very time consuming, and we are still crawling and managing to update the dataset for better coverage with up-to-date data.

We conduct statistical experiments to get the properties of the collected dataset. The best way to show the distribution of the number of views, the number of votes, and the number of favorites is using the logarithmic bin interval versus the frequency or probability of each bin interval. Because, these properties change dramatically in linear scale, and their trends can be captured intuitively in logarithmic x-axis. But, we show the average ratings using a normal bin interval. Figure 12 illustrates the distributions of the view counts, ratings, vote counts, and favorite counts of the dataset. Each bar represents a bin where its interval is from the corresponding number written under the bin to right before the number written under the next bin. The last bin interval is from its corresponding number to the next predictable number in the sequence of the axis.

Figure 12 shows that most of the images have been seen more than 100 times, i.e., 500px Website has a live community, while many images have at least 1-10 votes or favorite clicks. Having a rating higher than 10 is considered high by us, because the rating trend changes its slope direction from bin 0-9 to bin 10-19 negatively, and after that, the slope will positively grow until bin 40-49. Most of the images in the dataset have a rating more than 40 which is a very high rating, and it indicates that the 500px community of the photographer has many highly-rated photos.

### 7.2 Decomposition Analysis

To show the improvement and the effectiveness of our decomposition step, we conduct some experiments on various detectors used in our framework including object detector, human pose estimator, and scene parser. Also, we examine our hysteresis detection, category detector, and pose clustering.

Method	MAP	person	seat	plant	animal	car
YOLO	52.0	73.5	40.1	33.2	71.2	54.8
Ours	60.2	78.1	53.2	46.8	69.8	58.2

Table 1: The accuracy comparison between our object detector model versus the YOLO model on the 500px dataset to detect some of its known objects.

#### 7.2.1 Object Detection

Our network for object detection is inspired by YOLO as it is fast compared to the others (Redmon et al., 2016). The output of the network is some bounding boxes where detected objects are respectively with their detection probabilities. As we have tested, non-person object detection of YOLO under 30% probability is not accurate enough on the 500px dataset, and any wrong detection affects all pixels in the bounding box based on our method. As a result, we divide the input image into bigger chunks of  $5 \times 5$  grid for a higher accuracy, and small objects are less important for detection as a secondary subject of photography. We implement our model as 24 convolutional layers with two fully connected layers. We train the whole network on ImageNet (Deng et al., 2009) for about a week, and three times on a labeled subset of 768 common failure cases (CFC) from the 500px dataset.

We evaluate our model compared to YOLO on a test subset from the dataset. We use the regular MAP on all intended objects. Table 1 shows the MAP and the average accuracies of some objects (person, seat, plant, animal, and car) for our trained model versus YOLO model. The “seat” average accuracy is the average for “seat, bench, and chair”, “plant” average accuracy is the average for “plant, tree, and grass”, and “animal” average accuracy is the average for “bird, cat, dog, cow, and sheep”.

#### 7.2.2 Pose Estimator

Our pose estimator architecture has two parallel lines predicting limb confidence map and encoding limb-to-limb association, which is inspired by RTMPPE architecture (Cao et al., 2017). We adjust the transformation parameters of the architecture including maximum rotation degree to 60, crop size to 500, scale min to 0.6 and scale max to 1.0. Since higher rotation degrees and bigger persons are used frequently in our work. Then, we train our model on MSCOCO (Lin et al., 2014), MPII (Andriluka et al., 2014) and three times on our 317 CFC from the 500px dataset.

To evaluate the performance of our pose estimator model, we leverage MAP of all limbs like Deeper-Cut (Insafutdinov et al., 2016). The comparison results

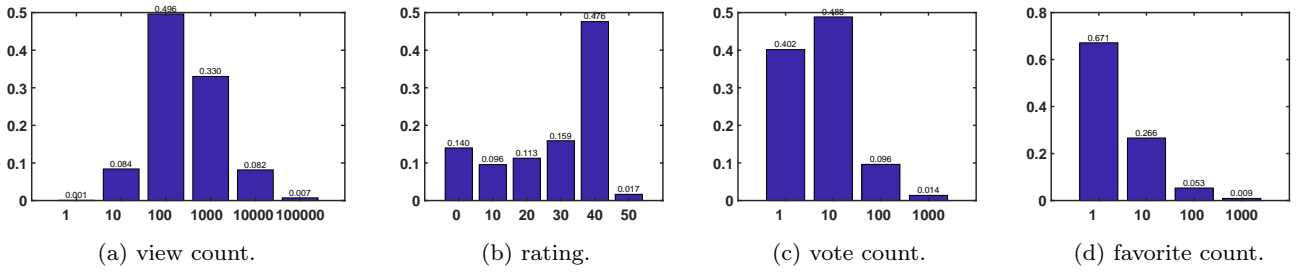


Fig. 12: The properties of the managed 500px dataset respectively from left to right including (a) the logarithmic distribution of the view counts, (b) the distribution of the ratings, (c) the logarithmic distribution of the vote counts and (d) the logarithmic distribution of the favorite counts.

Method	MAP	hea	sho	elb	wri	hip	kne	ank
RTMPPE	73.3	90.4	83.3	74.6	65.1	70.9	66.3	62.6
Ours	74.3	92.2	86.7	75.3	64.8	70.4	67.3	63.4

Table 2: The comparison between our pose estimator model versus the RTMPPE model on the 500px dataset to detect body parts.

of the MAP performance between RTMPPE and our approach on a subset of 507 testing images from our dataset are shown in Table 2, where the left limb and the right limb are merged.

### 7.2.3 Scene Parser

To parse any scene, we ignore confused categories like building and skyscraper. We place the related objects in the same object category. Also, our scene parser architecture exploits a 4-level pyramid pooling module (Zhao et al., 2017) with sizes of 11, 22, 33 and 44 respectively. We do not consider detecting small objects in the scene since they are mostly not the main subject of the photographer.

To train our model, we use ADE20K dataset (Zhou et al., 2017a) with 576 common failure cases annotated by LabelMe (Russell et al., 2008). To evaluate scene parsing performance, pixel-wise accuracy (PixAcc) and mean of class-wise intersection over union (CIoU) are measured. The performance values of our scene parser model versus PSPNet (Zhao et al., 2017) with 101-depth ResNet is shown in Table 3 where indicates better PixAcc and CIoU than PSPNet has been achieved on the 500px dataset.

### 7.2.4 Hysteresis Detection

Hysteresis detection covers more photos by allowing the union of all images above HIGH thresholds across

Method	Pixel Accuracy (%)	Mean IoU (%)
PSPNet	74.9	40.8
Ours	78.6	42.5

Table 3: The accuracy comparison between our scene parser model versus the PSPNet model with 101-depth ResNet on the 500px dataset.

the detectors. We show how we configure these tunable thresholds, while we trade-off between the total coverage and the partial accuracies by the detectors. When we have more than one detector with common detectable objects, we consider multiple features from all detectors to decide about the detection of the common objects. For example, “person” is a common object between object detector and pose estimation. We perform our pose estimator on our ground-truth images with a person or without any person from the dataset, and calculate (a) the detection score (as mentioned in Eq. 7) and (b) the normalized area (*i.e.* the detected object area divided by the image area) of the dominant person (*i.e.* the person with the highest score) detected in each image as our pose estimator features. Also, we perform our object detector on those images, and compute (c) the detection probability and (d) the normalized areas of the dominant person (*i.e.* the person with the highest probability) detected in each image as our object detector features.

Figure 13 shows the distributions of the features obtained from our pose estimator and our object detector for “person” as a common object between the pose estimator and the object detector. In some images, no person is detected by the pose estimator and the object detector, because the pose estimator or the object detector have detection error or actually there is no person in the image. We consider such detections as non-person object detection. Figure 14 shows the distributions of those features obtained from our pose estimator and our object detector, when there is no person in our

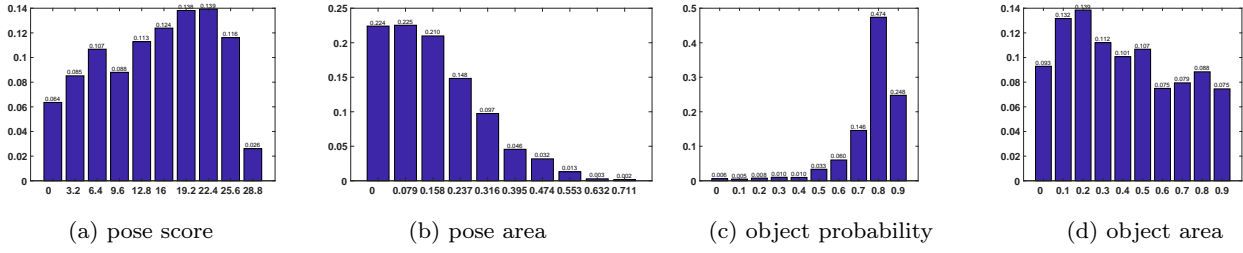


Fig. 13: The distributions of (a) the score obtained from the pose estimator, (b) the normalized area obtained from the pose estimator, (c) the detection probability obtained from the object detector, and (d) the normalized area obtained from the object detector for our ground-truth images with a “person” as a common detectable object by the pose estimator and the object detector.

ground-truth images. We have removed the frequency of the first component, *i.e.* score or area == 0, from all of the curved in Figure 14, because probability of zero score/area is very high and we want to bold the probabilities of the other score/area values.

Figure 13a shows pose estimator’s score does not have enough sensitivity to detect a person, because the distribution is similar to a uniform probability mass function (PMF). Similarly, Figure 13d shows object detector’s normalized area does not have enough sensitivity to detect a person, because the distribution is pretty uniform. But, object detector’s probability in Figure 13c and pose estimator’s normalized area in Figure 14b are not similar to uniform distribution, and we can infer some cut-off thresholds from them.

First, we derive the 2D probability density function (PDF) of the these mutual features including the normalized area by the pose estimator and the detection probability by the object detector. Second, we determine the 2D receiver operating characteristic (ROC) curve of these two parameters as a heat map. Finally, we search on the heat map and find the optimal point for these two mutual features. As shown in Figure 15, it can be inferred from 2D ROC of these two features that the optimal cut-off thresholds for “person” object detection in an image are object detector’s probability 40% and pose estimator’s normalized area 10% that leads to a detection accuracy (*i.e.* 92.2%) higher than other detectors’ accuracy solely.

### 7.2.5 Indexing with Integrated Object Detection

As shown in Figure 2, we do indexing of the 500px dataset by performing various detections. It is observed that once we integrate the object detectors in the dataset, all of the potentially detectable objects appear in the output. In addition, by collecting and distributing the results, the distribution (*i.e.*, the frequency of the total) of the semantic classes detected in our dataset except

for the highly-repetitive ones (“person” and “wall”) is illustrated in Figure 16.

The detectable objects by the IOD are not complete list of all available objects in any photo, but they can cover mostly-used objects in the photos. To investigate the case, we manually extract the available objects in 1600 random images from the dataset. Figure 17 shows the distribution of the objects highly available in the photos of the dataset.

### 7.2.6 Portrait Category Detection

As mentioned in Section 4.2, we start with top-down hierarchical clustering to specify the genre of the input image, and then we do multi-class categorization for portrait images. We train our model having 40 suggested features on a set of 6407 annotated portraits from the 500px dataset, and we test the model on another set of 1508 annotated portraits. The mean average accuracy of the model is listed in Table 4 categorized by various styles.

Also, we just consider the first 16 features for object detectors including general max and number of detected people in the image as mentioned in Section 4.2, and train a model using the same ground truth as before. The current model is our baseline model, because it can be used for any other object detector, as the features can be defined in other object detector domains as well. To compare rationally with this baseline, we test the same set of images from our ground truth. The second line in Table 4 listed the baseline results. Because we remove limb features, the baseline has no ability to detect sub-genres such as hand-only, leg-only, no-face, and sideview.

After portrait category indexing of the 500px dataset, the distribution of the portrait categories with respect to the number of corresponding images in each category divided by the total number of images is shown in Figure 18 stating the number of photos in full-body,



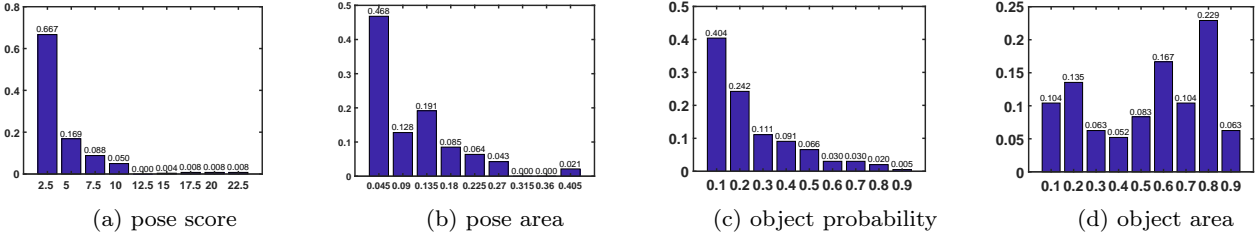


Fig. 14: For the ground-truth images without any “person”, the distributions of (a) the highest score (if any) obtained from the pose estimator, (b) the highest normalized area of the highest score object (if any) obtained from the pose estimator, (c) the detection probability for the dominant object (if any) obtained from the object detector, and (d) the normalized area of the dominant object (if any) obtained from the object detector.

	facial	full-body	group	hand	leg	no-face	sideview	two	upper-body
Our CaDe	94.35	92.40	85.90	44.0	74.57	79.35	67.50	74.74	90.81
16-feat Baseline	61.81	77.50	62.81	N/A	N/A	N/A	N/A	51.92	47.95

Table 4: The accuracy results of our category detector (CaDe) for ground-truth images from the 500px dataset compared to a baseline.

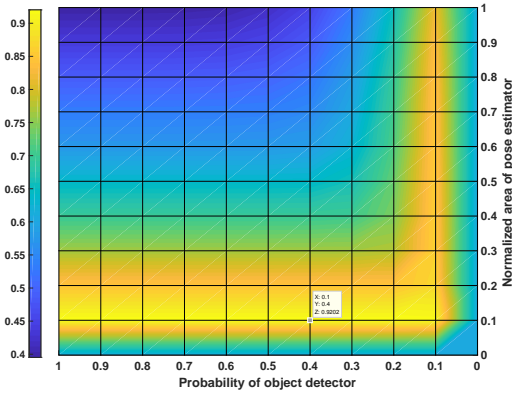


Fig. 15: The 2D ROC curve for the normalized area of the pose estimator and the detection probability of the object detector as a heat map.

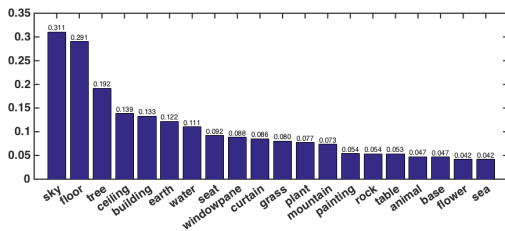


Fig. 16: The distribution of the highly-repetitive semantic classes detected in the dataset.

upper-body, facial, two, group, and side-view categories are high, but there are not enough samples for face-less, head-less, leg-only, and hand-only categories which is not a big deal, because these categories are not very popular.

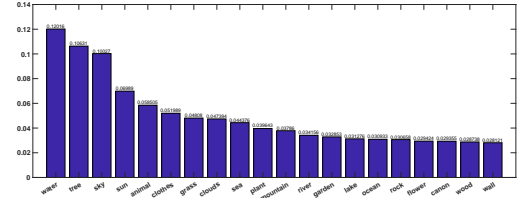


Fig. 17: The distribution of the mostly-available objects in the photos of the dataset which are manually extracted.

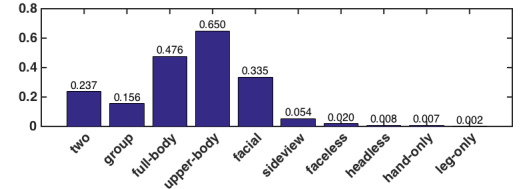


Fig. 18: The distribution of the portrait categories with respect to the number of corresponding images.

### 7.2.7 Artistic Pose Clustering

Regarding artistic pose clustering, we conduct an experiment to cluster similar professional poses using our features explained in Section 4.3. We do the clustering with a various number of cluster heads, and we find the optimal number of cluster heads for the 500px dataset using elbow method. That being said, we use the elbow method and do the clustering 40 times with the different number of clusters ranged from 1 to 40. This method calculates the sum of squared errors (the distance of each point to the center of its cluster) and it is expected to see an elbow pattern in the plot of this

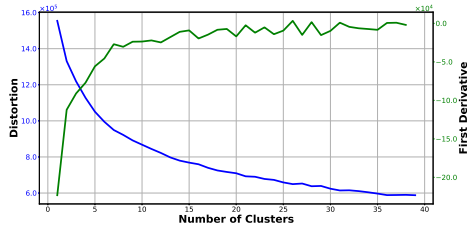


Fig. 19: The result of elbow method on the dataset. We could spot the elbow around 13-17 clusters. We are also showing the first derivative of the distortion, to show where it is going to flatten out.

error when the number of clusters is increasing. The result of this method on the 500px dataset is depicted in Figure 19 in Section 7.2.7, which indicates that the best choice for the number of clusters in this dataset is between 13 and 17.

### 7.3 User Study for Composition

Many computer vision papers compare with their peers to show the effectiveness of their approach, but currently there exists no other similar or comparable system in the literature to compare with our proposed framework. To evaluate the functionality and the performance of our method, and measure how much the recommended photos are relevant to the query and helpful to the photographer, we conduct a quantitative user study based on the human subject results to compare our method with two other reasonable baselines. The first baseline is a semantic and scene retrieval method based on state-of-the-art CNN model (Sharif Razavian et al., 2014) and the other baseline is a non-CNN retrieval method based on the color, shape, and texture features (Mitro, 2016). To create the baselines, all 4096 generic descriptors via public CNN model (Chatfield et al., 2014) trained on ImageNet (Deng et al., 2009) are extracted for the 500px dataset images as well as the features of non-CNN method (Mitro, 2016).

We select a variety of image queries (38 queries) based on many types of categories such as background scene and semantics, single versus group, full-body, upper-body, facial, standing versus sitting, and male versus female. To be fair, we do not use customized queries as shown in Figure 10, and we just focus on a single highlighted feature in each image query with the same question throughout the study. 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 87 participants. Our framework received 74.20% of the 1st ranks among the tests

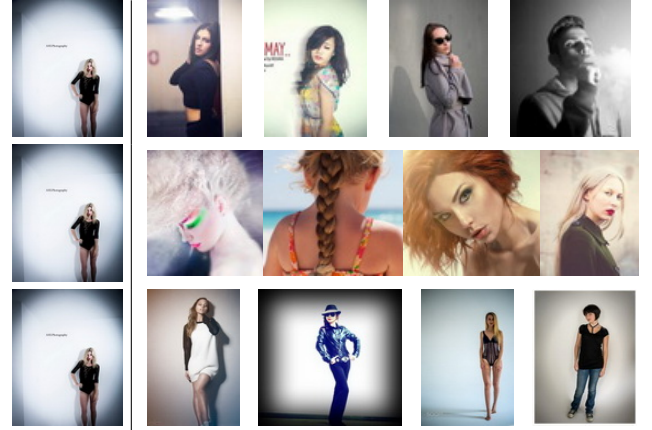


Fig. 20: The results of CNN (1st row), non-CNN (2nd row), and our method (3rd row) for a sample shot at the left side.

compared to 20.76% CNN as 2nd rank and 5.04% non-CNN as 3rd rank. Figure 20 illustrates the results of all methods (CNN: 1st, non-CNN: 2nd, ours: 3rd row) with respect to a similar shot at the left side. As it is realized from Figure 20 and our study, because the other methods cannot capture genre categories, scene structure, and corresponding poses of the query shot, it is common as shown in the figure that mixed categories are suggested by other semantic/category/pose-agnostic methods. As we hierarchically index the 500px dataset and recognize the right genre, category, pose and semantic classes, our semantic/category/pose-aware framework accessing to the indexed dataset can retrieve better related photography ideas.

As mentioned, the expected value of the accepted recommended photos by the participants with respect to the total number of recommendations including the baselines is 74.20%. More accurately, the histogram of the acceptability rate for the queries of the user study is shown in Figure 21. The x-axis shows the acceptability rate ranged from 0 to 1 with 0.1-width bins, *i.e.*, what percentage of the participants has accepted our recommended photos for some queries. The y-axis shows the frequency of our accepted recommendations by the total number of the examined corresponding queries (*i.e.* probabilities) which fall into each bin. The histogram has indicated that 23.2% of our recommended photos were accepted by over 90% of the participants, 44.6% of them with over 80%, 58.9% of them with over 70%, and 92.8% of them with over 50%. Consequently, the majority of the recommended photos are accepted with a mean of 74.20%.

### 7.4 Runtime Analysis

For training purposes, we mostly used an NVIDIA Tesla K40 GPU which took couple of days for the intensive

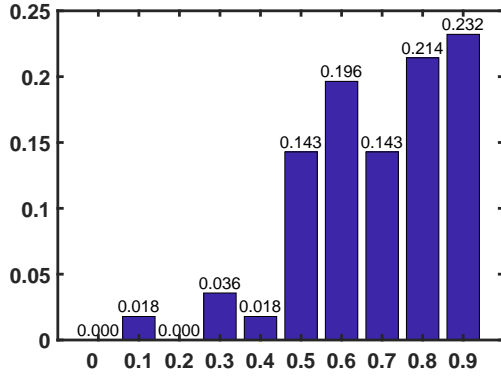


Fig. 21: The histogram of the acceptability rate for our recommended photos versus the total number of recommendations including other baselines.

computations on the dataset. To analyze the runtime performance of our method, we collect images with various styles and categories. Offline processing (indexing) of the 500px dataset is not included in the runtime results. We next perform IOD (including object detection, pose estimator, and scene parser), CaDe, ArPose as three parallel processes on those images to measure the average duration of the decomposition step. We then perform PAIR in composition step to get a ranked list of retrieved results for each image. Also, we randomly select the preferred style set, and perform the matching algorithm for each image as a single available shot. We sum up the average time for these three sequential steps to get the average runtime. The runtime is roughly proportional to the number of decomposed images because the IOD step which uses deep-learned models for detection is time-consuming. Among the detectors in IOD, pose estimator takes longer than the others and pretty independent from the number of people in the image with an average duration around 94.7 ms. Using a parallel structure shown in Figure 2, the end-to-end runtime is around 103.5 ms.

## 8 Conclusions and Discussions

We have collected a large dataset for portrait and landscape photography ideas and introduced a new framework for composition assistance which guides amateur photographers to capture better shots. As the number of photography ideas increases, retrieving and matching the camera photo directly with an image in the dataset becomes more challenging. Furthermore, the retrieving system not only finds similar images but also searches for images with similar semantic constellations with better composition through decomposition and composition steps. After providing feedback for the photographer, the camera tries to match the final pose

with one of the retrieved feedbacks, and make an astonishing shot. The performance of our framework has been evaluated by various experiments. Another merit of this work is the integration of the deep-based detectors which can make the whole process automatic.

### 8.1 Genre Extendibility

The general idea behind this work can be extended to other photography genres such as candid, fashion, close-up, and architectural 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 architectural photography.

### 8.2 Enhancement in User Interaction

The user-specified preferences (USP) should be quantified by the individual, but it may be difficult for them to accurately adjust the detectors' importance for their personal preference. They may want to do hierarchical preference, as some results are eliminated in each branch when going down the user-specified decision tree. Qualitatively they check the results and come up with a better decision, but it may be time-consuming for them. One can optimize the decision weights for a specific user after learning his/her behavior, and then they can just request for the ordering (not the weight values). We believe that there is still room to improve the interactions with the individual.

### 8.3 Clustering of Photography Ideas

Some future directions include working on an unsupervised learning approach that can cluster all the images based on various photography ideas. Recognizing the ideas is not easy for amateurs, and one shot can have multiple ideas. After clustering, we might detect new ideas. Therefore, it would be interesting to explore a metric to estimate the potential novelty of the current shot based on computing similarity to other shots.

### 8.4 Innovative Shot

Another promising idea is to design a system where the camera automatically detects an innovative situation and takes a shot. Conventional methods in machine learning just use the history of the field to help amateurs take professional photos, and of course, these

approaches can not go beyond it. After recognizing new photography ideas, the system can think of it as a compact space, not a finite discrete space, and it attempts to find a solution in this compact space. Fortunately, the complexity of the problem can change from an integer programming to a linear programming, but the way we define these compact spaces is hard based on the complexity of finding new photography ideas.

**Acknowledgements** This research has been supported in part by Penn State University. We gratefully acknowledge the support of NVIDIA Corporation with the donation of the Tesla K40 GPUs used for this research. The work would not have been possible without the large quantity of creative photography ideas of the thousands of photographers who share their visual creations on the Internet. The authors would like to thank the participants of the user study.

## References

- Andriluka M, Pishchulin L, Gehler P, Schiele B (2014) 2d human pose estimation: New benchmark and state of the art analysis. In: IEEE conference on computer vision and pattern recognition
- Barnes C, Shechtman E, Finkelstein A, Goldman D (2009) Patchmatch: A randomized correspondence algorithm for structural image editing. *ACM Transactions on Graphics* 28(3):1–24
- Bhattacharya S, Sukthankar R, Shah M (2010) A framework for photo-quality assessment and enhancement based on visual aesthetics. In: ACM conference on multimedia, pp 271–280
- Bhattacharya S, Sukthankar R, Shah M (2011) A holistic approach to aesthetic enhancement of photographs. *ACM Transactions on Multimedia Computing, Communications, and Applications* 7(1):21
- Cao Z, Simon T, Wei SE, Sheikh Y (2017) Realtime multi-person 2d pose estimation using part affinity fields. In: IEEE conference on computer vision and pattern recognition, pp 7291–7299
- Chang HT, Pan PC, Wang YCF, Chen MS (2015) R2p: recomposition and retargeting of photographic images. In: ACM conference on multimedia, pp 927–930
- Chatfield K, Simonyan K, Vedaldi A, Zisserman A (2014) Return of the devil in the details: Delving deep into convolutional nets. *arXiv preprint arXiv:14053531*
- Cho TS, Butman M, Avidan S, Freeman WT (2008) The patch transform and its applications to image editing. In: IEEE conference on computer vision and pattern recognition, pp 1–8
- Coates A, Ng A, Lee H (2011) An analysis of single-layer networks in unsupervised feature learning. In: International conference on artificial intelligence and statistics, pp 215–223
- Datta R, Joshi D, Li J, Wang JZ (2006) Studying aesthetics in photographic images using a computational approach. *European conference on computer vision* pp 288–301
- Datta R, Joshi D, Li J, Wang JZ (2008) Image retrieval: Ideas, influences, and trends of the new age. *ACM Computing Surveys* 40(2):5
- Deng J, Dong W, Socher R, Li LJ, Li K, Fei-Fei L (2009) Imagenet: A large-scale hierarchical image database. In: IEEE conference on computer vision and pattern recognition, pp 248–255
- Diyanat A, Farhat F, Ghaemmaghani S (2011) Image steganalysis based on svd and noise estimation: Improve sensitivity to spatial lsb embedding families. In: TENCON 2011-2011 IEEE Region 10 Conference, IEEE, pp 1266–1270
- Dunn JC (1973) A fuzzy relative of the isodata process and its use in detecting compact well-separated clusters. *Journal of Cybernetics* 3(3):32–57
- Farhat F, Ghaemmaghani S (2014) Towards blind detection of low-rate spatial embedding in image steganalysis. *IET Image Processing* 9(1):31–42
- Farhat F, Tootaghaj DZ (2017) Carma: Contention-aware auction-based resource management in architecture. *arXiv preprint arXiv:171000073*
- Farhat F, Diyanat A, Ghaemmaghani S, Aref MR (2011) Multi-dimensional correlation steganalysis. In: Multimedia Signal Processing (MMSP), 2011 IEEE 13th International Workshop on, IEEE, pp 1–6
- Farhat F, Diyanat A, Ghaemmaghani S, Aref MR (2012) Eigenvalues-based lsb steganalysis. *The ISC International Journal of Information Security* 4(2):97–106
- Farhat F, Tootaghaj D, He Y, Sivasubramaniam A, Kandemir M, Das C (2016a) Stochastic modeling and optimization of stragglers. *IEEE Transactions on Cloud Computing*
- Farhat F, Tootaghaj DZ, Arjomand M (2016b) Towards stochastically optimizing data computing flows. *arXiv preprint arXiv:160704334*
- Farhat F, Kamani MM, Mishra S, Wang JZ (2017) Intelligent portrait composition assistance: Integrating deep-learned models and photography idea retrieval. In: Proceedings of the on Thematic Workshops of ACM Multimedia 2017, ACM, pp 17–25
- Grey C (2014) Master lighting guide for portrait photographers. Amherst Media
- Guo Y, Liu M, Gu T, Wang W (2012) Improving photo composition elegantly: Considering image similarity during composition optimization. *Computer Graphics Forum* 31(7):2193–2202



- Harel J, Koch C, Perona P (2006) Graph-based visual saliency. In: *Advances in neural information processing systems*, pp 545–552
- He S, Zhou Z, Farhat F, Wang JZ (2018) Discovering triangles in portraits for supporting photographic creation. *IEEE Transactions on Multimedia* 20(2):496–508
- Insafutdinov E, Pishchulin L, Andres B, Andriluka M, Schiele B (2016) Deepercut: A deeper, stronger, and faster multi-person pose estimation model. In: *European conference on computer vision*, Springer, pp 34–50
- Itti L, Koch C, Niebur E (1998) A model of saliency-based visual attention for rapid scene analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 20(11):1254–1259
- Joshi D, Datta R, Fedorovskaya E, Luong QT, Wang JZ, Li J, Luo J (2011) Aesthetics and emotions in images. *IEEE Signal Processing Magazine* 28(5):94–115
- Kamani MM, Farhat F, Wistar S, Wang JZ (2016) Shape matching using skeleton context for automated bow echo detection. In: *IEEE International Conference on Big Data*, IEEE, pp 901–908
- Kamani MM, Farhat F, Wistar S, Wang JZ (2017) Skeleton matching with applications in severe weather detection. *Applied Soft Computing*
- Ke Y, Tang X, Jing F (2006) The design of high-level features for photo quality assessment. In: *IEEE conference on computer vision and pattern recognition*, vol 1, pp 419–426
- Khan SS, Vogel D (2012) Evaluating visual aesthetics in photographic portraiture. In: *Symposium on computational aesthetics in graphics, visualization, and imaging*, Eurographics Association, pp 55–62
- Krages B (2012) *Photography: the art of composition*. Skyhorse Publishing, Inc.
- Lauer DA, Pentak S (2011) *Design basics*. Wadsworth Publishing
- LeCun Y, Bottou L, Bengio Y, Haffner P (1998) Gradient-based learning applied to document recognition. *Proceedings of the IEEE* 86(11):2278–2324
- Lew MS, Sebe N, Djeraba C, Jain R (2006) Content-based multimedia information retrieval: State of the art and challenges. *ACM Transactions on Multimedia Computing, Communications, and Applications* 2(1):1–19
- Lewis DD, Yang Y, Rose TG, Li F (2004) Rcv1: A new benchmark collection for text categorization research. *Journal of Machine Learning Research* 5(Apr):361–397
- Li C, Wang P, Wang S, Hou Y, Li W (2017) Skeleton-based action recognition using lstm and cnn. In: *IEEE international conference on multimedia & expo workshops*, IEEE, pp 585–590
- Li J, Yao L, Wang JZ (2015a) Photo composition feedback and enhancement. In: *Mobile cloud visual media computing*, Springer, pp 113–144
- Li K, Yan B, Li J, Majumder A (2015b) Seam carving based aesthetics enhancement for photos. *Signal Processing: Image Communication* 39:509–516
- Lienhard A, Reinhard M, Caplier A, Ladret P (2014) Photo rating of facial pictures based on image segmentation. In: *IEEE conference on computer vision theory and applications*, vol 2, pp 329–336
- Lienhard A, Ladret P, Caplier A (2015a) How to predict the global instantaneous feeling induced by a facial picture? *Signal Processing: Image Communication* 39:473–486
- Lienhard A, Ladret P, Caplier A (2015b) Low level features for quality assessment of facial images. In: *International conference on computer vision theory and applications*, pp 545–552
- Lin TY, Maire M, Belongie S, Hays J, Perona P, Ramanan D, Dollár P, Zitnick CL (2014) Microsoft coco: Common objects in context. In: *European conference on computer vision*, pp 740–755
- Liu L, Jin Y, Wu Q (2010) Realtime aesthetic image retargeting. In: *Computational aesthetics*, pp 1–8
- Liu Z, Wang Z, Yao Y, Zhang L, Shao L (2018) Deep active learning with contaminated tags for image aesthetics assessment. *IEEE Transactions on Image Processing*
- Lu X, Lin Z, Jin H, Yang J, Wang JZ (2015) Rating image aesthetics using deep learning. *IEEE Transactions on Multimedia* 17(11):2021–2034
- Luo Y, Tang X (2008) Photo and video quality evaluation: Focusing on the subject. In: *European conference on computer vision*, pp 386–399
- Mai L, Jin H, Liu F (2016) Composition-preserving deep photo aesthetics assessment. In: *IEEE conference on computer vision and pattern recognition*, pp 497–506
- Males M, Hedi A, Grgic M (2013) Aesthetic quality assessment of headshots. In: *International symposium ELMAR*, pp 89–92
- Marchesotti L, Perronnin F, Larlus D, Csurka G (2011) Assessing the aesthetic quality of photographs using generic image descriptors. In: *IEEE conference on computer vision*, pp 1784–1791
- Matthew B (2010) 101 quick and easy ideas taken from the master photographers of the twentieth century. Course Technology CENGAGE Learning
- Mitro J (2016) Content-based image retrieval tutorial. arXiv preprint arXiv:160803811

- Murray N, Marchesotti L, Perronnin F (2012) AVA: a large-scale database for aesthetic visual analysis. In: Conference on computer vision and pattern recognition, pp 2408–2415
- Park J, Lee JY, Tai YW, Kweon IS (2012) Modeling photo composition and its application to photo rearrangement. In: IEEE conference on image processing, pp 2741–2744
- Pritch Y, Kav-Venaki E, Peleg S (2009) Shift-map image editing. In: International conference on computer vision, vol 9, pp 151–158
- Redi M, Rasiwasia N, Aggarwal G, Jaimes A (2015) The beauty of capturing faces: Rating the quality of digital portraits. In: IEEE conference on automatic face and gesture recognition, vol 1, pp 1–8
- Redmon J, Divvala S, Girshick R, Farhadi A (2016) You only look once: Unified, real-time object detection. In: IEEE conference on computer vision and pattern recognition, pp 779–788
- Rice P (2006) Master guide for professional photographers. Amherst Media
- Russell BC, Torralba A, Murphy KP, Freeman WT (2008) Labelme: A database and web-based tool for image annotation. *International Journal of Computer Vision* 77(1-3):157–173
- Samii A, Mëch R, Lin Z (2015) Data-driven automatic cropping using semantic composition search. *Computer Graphics Forum* 34(1):141–151
- Santella A, Agrawala M, DeCarlo D, Salesin D, Cohen M (2006) Gaze-based interaction for semi-automatic photo cropping. In: ACM conference on human factors in computing systems, pp 771–780
- Sharif Razavian A, Azizpour H, Sullivan J, Carlsson S (2014) Cnn features off-the-shelf: An astounding baseline for recognition. In: IEEE conference on computer vision and pattern recognition workshops, pp 806–813
- Smeulders AW, Worring M, Santini S, Gupta A, Jain R (2000) Content-based image retrieval at the end of the early years. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 22(12):1349–1380
- Smith J (2012) Posing for portrait photography: a head-to-toe guide for digital photographers. Amherst Media
- Stentiford F (2007) Attention based auto image cropping. In: ICVS workshop on computational attention & applications, Citeseer, vol 1
- Suh B, Ling H, Bederson BB, Jacobs DW (2003) Automatic thumbnail cropping and its effectiveness. In: ACM symposium on user interface software and technology, pp 95–104
- Talebi H, Milanfar P (2018) Nima: Neural image assessment. *IEEE Transactions on Image Processing* 27(8):3998–4011
- Tootaghaj DZ, Farhat F (2017) Cage: A contention-aware game-theoretic model for heterogeneous resource assignment. In: Computer Design (ICCD), 2017 IEEE International Conference on, IEEE, pp 161–164
- Valenzuela R (2012) Picture perfect practice: a self-training guide to mastering the challenges of taking world-class photographs. New Riders
- Valenzuela R (2014) Picture perfect posing: practicing the art of posing for photographers and models. New Riders
- Wong LK, Low KL (2009) Saliency-enhanced image aesthetics class prediction. In: IEEE conference on image processing, pp 997–1000
- Xie J, Girshick R, Farhadi A (2016) Unsupervised deep embedding for clustering analysis. In: International conference on machine learning, pp 478–487
- Xu Y, Ratcliff J, Scovell J, Speiginer G, Azuma R (2015) Real-time guidance camera interface to enhance photo aesthetic quality. In: ACM conference on human factors in computing systems, pp 1183–1186
- Xue SF, Tang H, Tretter D, Lin Q, Allebach J (2013) Feature design for aesthetic inference on photos with faces. In: IEEE conference on image processing, pp 2689–2693
- Yan J, Lin S, Kang S, Tang X (2013) Learning the change for automatic image cropping. In: IEEE conference on computer vision and pattern recognition, pp 971–978
- Yao L, Suryanarayan P, Qiao M, Wang JZ, Li J (2012) Oscar: On-site composition and aesthetics feedback through exemplars for photographers. *International Journal of Computer Vision* 96(3):353–383
- Zhang M, Zhang L, Sun Y, Feng L, Ma Wy (2005) Auto cropping for digital photographs. In: IEEE conference on multimedia and expo
- Zhao H, Shi J, Qi X, Wang X, Jia J (2017) Pyramid scene parsing network. In: IEEE conference on computer vision and pattern recognition, pp 2881–2890
- Zhou B, Zhao H, Puig X, Fidler S, Barriuso A, Torralba A (2017a) Scene parsing through ade20k dataset. In: IEEE conference on computer vision and pattern recognition
- Zhou Z, Farhat F, Wang JZ (2016) Detecting dominant vanishing points in natural scenes with application to composition-sensitive image retrieval. *arXiv preprint arXiv:160804267*
- Zhou Z, Farhat F, Wang JZ (2017b) Detecting dominant vanishing points in natural scenes with application to composition-sensitive image retrieval. *IEEE Transactions on Multimedia* 19(12):2651–2665