

Two-Class Weather Classification

Cewu Lu, *Member, IEEE*, Di Lin, *Student Member, IEEE*, Jiaya Jia, *Senior Member, IEEE*,
Chi-Keung Tang, *Senior Member, IEEE*

Abstract—Given a single outdoor image, we propose a collaborative learning approach using novel weather features to label the image as either sunny or cloudy. Though limited, this two-class classification problem is by no means trivial given the great variety of outdoor images captured by different cameras where the images may have been edited after capture. Our overall weather feature combines the data-driven convolutional neural network (CNN) feature and well-chosen weather-specific features. They work collaboratively within a unified optimization framework that is aware of the presence (or absence) of a given weather cue during learning and classification. In this paper we propose a new data augmentation scheme to substantially enrich the training data, which is used to train a latent SVM framework to make our solution insensitive to global intensity transfer. Extensive experiments are performed to verify our method. Compared with our previous work and the sole use of a CNN classifier, this paper improves the accuracy up to 7 – 8%. Our weather image dataset is available together with the executable of our classifier.

Index Terms—weather understanding, image classification, structure SVM.

1 INTRODUCTION

We address the problem of two-class weather classification from a single outdoor image. This seemingly easy task for humans – to tell whether a given image is sunny or cloudy – turns out to be challenging. This paper attempts to provide technical insight and solutions to address the above issues, while acknowledging that our work is a first but significant step for weather understanding from single images. Note that naive schemes based on image brightness or color/intensity statistics (Figures 1 and 2) are doomed to fail in this two-class classification problem. While hardware solutions relying on expensive sensors are employed, for centuries human vision is still the most powerful tool for weather observation. If we can exploit existing surveillance and smartphone cameras, which are found almost everywhere, it may be possible to turn human weather observation into a powerful and cost-effective computer vision application.

Previously, in [37], we demonstrated that careful engineering of well-chosen weather-specific features employed in supervised learning can adequately address this two-class weather classification problem. In this paper, we further investigate the efficacy of the state-of-the-art convolutional neural network (CNN) in solving the problem. The CNN approach as well as the CNN feature is data-driven, which is in contrast to hand-picked weather features in [37]. As will be shown in the experimental section, we found that the concatenation of the CNN feature and weather-specific features reports the best performance, namely, a 7–8% improvement over the sole use of a CNN classifier that is



Fig. 1. (a) A sunny image with mean lightness 32.41. (b) A cloudy image with mean lightness 58.25.

trained end-to-end using the given training data. The data-driven CNN feature and the weather-specific features work together to exploit the synergies between the two. Based on this new overall weather feature, our approach consists of the following three technical contributions:

First, we describe the design and implementation of various weather cues, which are used to form the weather feature [37]. These everyday weather cues (such as sky, shadow, reflection, contrast and haze) are what humans are still using for weather observing – a hazy or grayish sky characterizes a cloudy day while hard shadow cast on the ground indicates a sunny day, as illustrated in Figure 3(a). Conversely, in the absence of any weather cues, even we humans may not be confident in labeling the weather type, as illustrated in Figure 3(b). In this paper, we concatenate the CNN feature with the above weather-specific features to form the overall weather feature in training and testing.

Given the overall weather feature, the next question is how to properly learn the classifier. The main issue is that the weather cues used in this paper may not be all available in an image – e.g., not every outdoor image has a sky region – which is problematic to a discriminative training process adopted by traditional classifiers such as SVM. To address this issue, our second technical contribution consists of a *collaborative learning* framework using homogeneous voters – the outdoor images are clustered where images in the same

- Cewu Lu is with the Department of Computer Science, Stanford University, Stanford, CA 94305, USA.
- Jiaya Jia and Di Lin are with the Department of Computer Science and Engineering, The Chinese University of Hong Kong.
- Chi-Keung Tang is with the Department of Computer Science and Engineering, The Hong Kong University of Science and Technology.

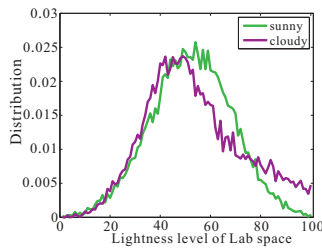


Fig. 2. Pixel intensity distributions in the lightness L channel in the LAB color space of 5K cloudy images and 5K sunny images. It is almost impossible to draw a decision boundary between the two types of weather.

cluster are similar in terms of weather cues. This allows us to build a classifier in a conventional way thanks to the homogeneity in each cluster. The final labeling is the weighted voting result of the cluster classifier outputs. The cluster closer to the testing image is given a higher weight. As will be explained in the following, homogeneous voters are learned under a unified optimization framework.

To make our system more robust to training images harvested from the web, we propose a novel strategy to enrich our training set by synthesizing for each training image its *weather counterpart images*, which belong to a subclass of images of the same scene taken under different camera settings and/or after photo editing that can be characterized by a *global color/intensity transfer*. The training image and its synthesized weather counterparts together are then used in *latent SVM learning* which encourages each training sample and its counterparts to have the *same* weather label. Our synthesis strategy is scalable. That is, the production of the weather counterparts given a training image is fully automatic, which requires no further data collection or annotation by humans.

Finally, we perform quantitative comparison with a number of typical baselines including SVM, Adaboost [58], [62], and prior weather-related methods [26], [61], [46]. Our final contribution consists of a 10,000-image weather dataset in which the images are properly selected and annotated. This is used to evaluate our learning and labeling strategy.

This manuscript extends its conference version [37] along the following dimensions:

- The overall weather feature combines the data-driven CNN feature and handcrafted weather features.
- A data-driven approach for synthesizing weather counterparts to make scalable data collection and training; the new system is insensitive to global intensity transfer and achieves improvement over [37].
- A latent SVM framework is proposed to capture a wide variety of global intensity transfer.
- More experiments are conducted to evaluate the proposed method.

The paper is organized as follows. In Section 2 we review the related work. In Section 3, we introduce our weather-specific and data driven CNN features. In Section 4, we describe our weather dataset and weather counterparts generation. Section 5 presents the collaborative and latent SVM learning

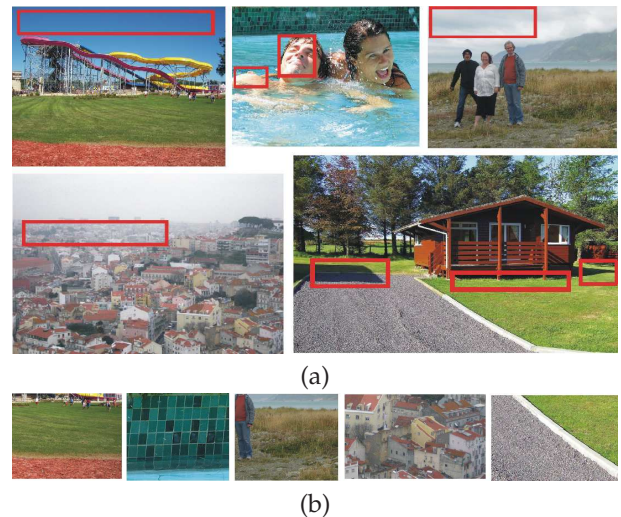


Fig. 3. Weather cues. (a) Common weather cues in red rectangles. (b) Regions in (a) lacking any weather cues.

of our weather classifier. Section 6 discusses our results on weather classification. We conclude this paper in Section 7.

2 RELATED WORK

This section gives an overview of the related work on weather understanding, which can be regarded as a category in scene recognition. The background of the convolutional neural network (CNN) is also investigated since the CNN feature is an important component in our computational framework.

2.1 Weather Understanding

2.1.1 Weather Understanding with Hand-crafted Feature

Weather understanding plays a vital role in many real-world applications such as navigation control in self-driving cars. Automatic understanding of weather conditions enhances road safety by, for instance, controlling the vehicle speed in response to real-time weather situation [61], [46]. In [20], a built-in weather understanding component was found in an accurate navigation system that involves sky detection. Raindrops have been a frequently used cue for weather recognition, and in [23], discriminative raindrop templates were learned to infer weather situation. In [50], a photometric stereo-based method was proposed to estimate weather situation. Multiple images were required to estimate the illumination situation of a given site. Therefore, only a few sites (e.g. popular tourist sites) can meet this requirement. In [40] Narasimhan *et al.* proposed a physics-based model to capture multiple scattering of light rays from a source to the camera. This model works well for cases where scattering effect is strong, such as fog, haze, mist and rain at night. In [24], 40 transient attributes were studied and a model was learned to predict these attributes given a single image. The approach is standard “features + SVM” scheme where the features used are standard (e.g. HOG, SIFT) which may not capture well weather characteristics. In [65] a multi-class

weather classification method was described using multiple weather features and multiple kernel learning.

In [7], [51], weather types were recognized by the motion of cloud, snow flakes, etc. These weather-related phenomena depict periodic movement which provides specific weather patterns for detection.

Though the above methods have shown good performance in their respective applications, custom devices or conditions were often required. Exploiting existing smartphones which are cheaper, and surveillance cameras which can be found or installed almost everywhere, can make it possible to turn general weather observation into a powerful and cost-effective computer vision application.

2.1.2 Weather Understanding with CNN

Our problem is related to classification which can also be solved by Convolutional Neural Networks (CNN). Deep CNN has led a series of breakthroughs in image classification [22]. It is a feed-forward, end-to-end multilayered neural network inspired by the organization of the animal visual cortex. The convolutional neural network has found applications in image and video recognition tasks, such as video classification [19], object detection [11], and action recognition [1]. Excellent models, such as AlexNet [22], Network-in-Network [35], VGG [52] and ResNet [14] have been developed. Unlike general image classification (e.g. object classification), weather classification relies on weather-sensitive cues (as we shall demonstrate in the experiment section) which is somewhat similar to fine-grained recognition in the level of details required, while the deep CNN is excellent in capturing global scene semantics which needs to be integrated to make such recognition succeed [64], [34]. There is only a handful of research attempts in applying CNN to weather understanding. In [9], convolutional neural networks (VGG) is directly used in classifying weather, and in [56] to predict outdoor ambient temperature and the time of the year. The CNN filter may however miss subtle weather cues inherent in the input image.

2.2 Scene Understanding

Weather understanding is a specific case of scene recognition. General scene recognition focuses on discovering discriminative scene structure. As explained in [45], scene structure can be regarded as a combination of parts which are called regions of interest. These discriminative parts provide a powerful representation of the scene. Thus exploiting them in relevant tasks has recently become a popular trend.

The work [48], [30], [31], [53], [32], [17] discovered parts with specific visual concepts. Each learned part is expected to represent a single or a cluster of visual objects, which is beneficial to alleviate visual ambiguity. Meanwhile, unsupervised discovery of discriminative parts has received much attention. Though handcrafted part filters are easy to comprehend, they strongly rely on human labeling and are not scalable. Unsupervised frameworks [55], [18], [21], [28], [42], [43], [49], [66], [32], [33] can be more practical and efficient especially for large data sets.

Recently, various mid-level representations have been employed to enhance the discriminative power in classification [3], [30], [31], [66]. State-of-the-art methods [53], [32], [17] proposed discriminative parts and used them to construct mid-level representation, e.g., response maps obtained from convolution with part filters. These mid-level representations are fed into discriminative classifiers and evaluated on different scene classification datasets. Mid-level representation can be a better alternative or complementary to traditional low-level representations [6], [36], [41], [58], [27], [44], [62], because mid-level representation is capable of differentiating among a large variety of inter and intra categories as described in [63].

We note the methods in scene classification cannot be employed to solve our problem. Though weather is part of the scene, it is not as concrete (e.g., no closed boundary) as objects such as trees, buildings, and mountains, thus risking information loss when we apply these methods.

2.3 Weather Applications

Weather cues have been used to enable various applications. In [2], deep convolutional neural networks was used to estimate transient attributes including weather, time of the day, season and subjective properties of a given scene. In [12], the interaction between the appearance of an outdoor scene and the ambient temperature was studied, where the statistical correlations between image sequences from outdoor cameras and temperature measurements were derived. In [15], weather conditions with the scene structure and position of the sun were used to estimate the time and location the image was captured. Snow recognition was studied in [57], which enables the production of satellite maps of snowfall using geo-tagged, time stamped images from Flickr. Cloud cues were explored in [59], where a method was presented for estimating the geometry of an outdoor scene. Another work that used cloud cues is [16], where cloud motion enables geometric calibration of static outdoor cameras.

3 THE OVERALL WEATHER FEATURE

We compute for each image the *overall weather feature*, a 4717-D feature consisting of two parts, namely the CNN feature and five weather features. The feature vector is formed by concatenating the six components

$$[\mathbf{f}_{sk}; \mathbf{f}_{sh}; \mathbf{f}_{re}; \mathbf{f}_{co}; \mathbf{f}_{ha}; \mathbf{f}_{cn}] \quad (1)$$

where the first five features, namely, sky, shadow, reflection, contrast and haze, correspond to a key weather cue to be defined shortly. We incorporate the CNN feature [22] \mathbf{f}_{cn} to describe the image in general, which is extracted from a learned two-class weather CNN model. Since not all of these cues are necessarily present in a given outdoor image, we also compute the *existence vector*

$$[v_{sk}; v_{sh}; v_{re}; v_{ha}; v_{cn}], \quad (2)$$

where each scalar score in $[0, 1]$ indicates the confidence that the corresponding weather cue is present in the given image and in particular, v_{cn} is the confidence score of the CNN classifier. Since image contrast difference exists in both sunny and cloudy photos, v_{co} is always 1 and excluded.

3.1 Weather Feature

3.1.1 Sky

If present, the sky is the most important cue for weather labeling. A clear, cloudless sky is blue as air molecules scatter blue light more than red light. Cloud is made of tiny water droplets which make the sky look grayish white.

To define v_{sk} , the sky region is detected in a pixel-wise manner in the following steps. We respectively collect 20,000 sky and non-sky patches, each of size 15×15 , and extract a 131 dimensional feature, which contains the SIFT descriptor (128D) and mean HSV color (3D). This feature was suggested in [54]. Then a random forest classifier is learned on the two patch classes. Now, given an image, we uniformly sample 15×15 patches and test their labels (sky or non-sky) as seeds. Sky region can be segmented by implementing graph cuts on those seeds (see Figure 4(a)–(b)). Let A be the sky to image area ratio. We set $v_{sk} \in [0, 1]$ as

$$v_{sk} = \begin{cases} 1 & \text{if } A > 0.5 \\ \min\{2A, 1\} & \text{otherwise} \end{cases} \quad (3)$$

To define the \mathbf{f}_{sk} vector we have considered various alternatives. Straightforward color histogram feature in the sky region suffers from two defects. First, possible sky colors (both cloudy and sunny) are sparse, thus yielding most color bins with the zero value (Figure 4(c)). Second, no adequate consideration is given to color contrast. In this paper, we define \mathbf{f}_{sk} using color-pair dictionary coding as follows.

We collect 2,000 images with detected sky regions. Neighborhood pixels in pairs are extracted from the sky region to form a large number of 6D vectors, each of them consisting of a total of 6 RGB values. This process results in about 100,000 pixel pairs. We then learn a sky color-pair dictionary $\mathbf{D} \in \mathbb{R}^{6 \times 256}$ on the vectors using the method described in [38], thus producing a set of neighborhood-pixel vectors sparsely coded over the learned dictionary, expressed as

$$\min_{\beta_i} \|p_i - \mathbf{D}\beta_i\|_2^2 + \lambda \|\beta_i\|_1, \quad (4)$$

where $p_i \in \mathbb{R}^{6 \times 1}$ is the i^{th} vector, $\beta_i \in \mathbb{R}^{256 \times 1}$ is the sparse code over \mathbf{D} . We solve Eq. (4) using [58]. Our final \mathbf{f}_{sk} is filtered by max pooling of all β_i . That is, the j^{th} bin of our feature is set to $\max_i \{\beta_{i,j}\}$ where $\beta_{i,j}$ is the j^{th} bin of β_i .

Max pooling can preserve subtle sun-to-cloud contrast in the feature representation. Figure 4(d) shows a typical \mathbf{f}_{sk} plot. In comparison to color histogram, our 256-D \mathbf{f}_{sk} covers the full range of the histogram and encodes color contrast information as well. The advantage over color histogram was demonstrated in [37].

3.1.2 Shadow

Hard shadow boundaries form another useful cue because they are often found in outdoor photos shot in sunny days. To compute v_{sh} and \mathbf{f}_{sh} , we resort to shadow detection tools. Unlike sky detection, shadow detection in an image is still a challenging problem. Our extensive evaluation indicates while working well in sunny images, state-of-the-art shadow detection often fails for cloudy images, where dark regions are often misclassified as shadow as shown in Figure 5.

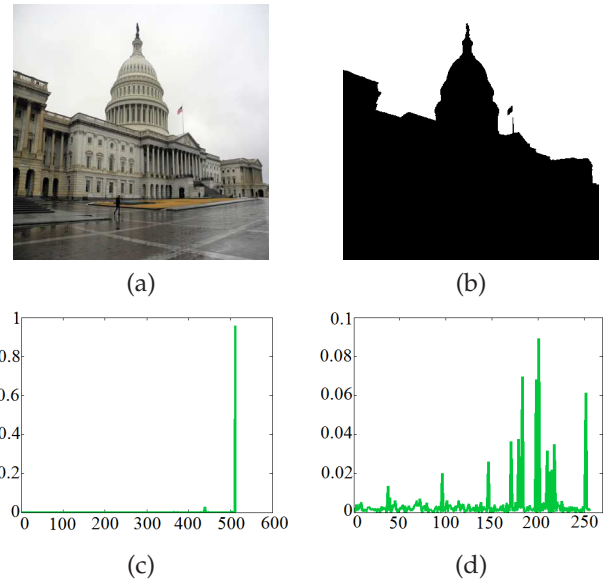


Fig. 4. Sky. (a) input image, (b) detected sky region, (c) color histogram of the sky, (d) plot of \mathbf{f}_{sk} .



Fig. 5. Shadow detection results of [25] for (a) a cloudy image and (b) a sunny image. Shadow detection in cloudy images is vulnerable to false detection.

Notwithstanding, we apply [25], rank the resulting shadow boundary confidence scores and take the 10^{th} highest score to set v_{sh} . This serves as a rough relative indicator in our method. A larger v_{sh} represents possibly stronger shadow presence. High precision is not needed in the estimation.

Using a data-driven approach we design our \mathbf{f}_{sh} by relying on the shadows detected in the training images restricted to sunny outdoor photos. If a given boundary is similar to those training shadow boundaries, we regard this as a shadow boundary typical of a sunny image.

In detail, initially, for all of the sunny images in the training set, we apply [25] to detect shadow boundaries and generate their corresponding confidence scores and boundary descriptors. For each image, we keep only the top 10 most confident shadow boundaries, and save them to the pool \mathcal{P} which has $10V$ samples, where V is the number of sunny images in the training set.

Given a boundary, we measure its likelihood to be a shadow boundary typical of a sunny photo by the mean distance to its K -nearest ($K = 5$) neighbors in \mathcal{P} . Two examples of K -nearest neighbor matching are shown in Figure 6. The Euclidean distance between the two boundaries descriptor vectors was used [25]. Given an image, we obtain its top 10 most confident shadow boundaries and compute their likelihood as described above to form the 10-D \mathbf{f}_{sh} vector.

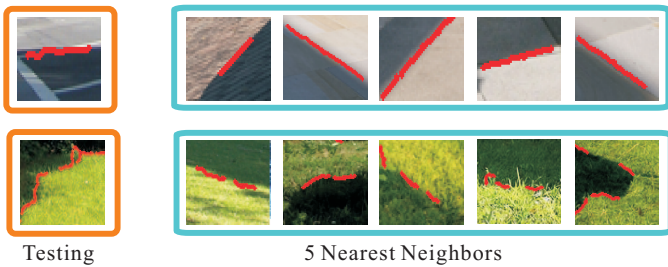


Fig. 6. K -nearest neighbor matching in \mathcal{P} . Shown in the blue rectangles are the five nearest neighbors.

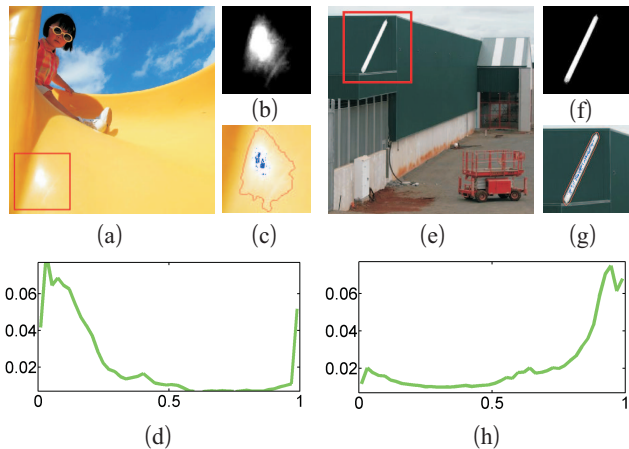


Fig. 7. Reflection cue. A sunny image with strong sunlight reflection in (a) versus a cloudy image with inherently white regions in (e). (b) and (f) are the corresponding alpha mattes. In (c) and (g), red and blue points indicate background and foreground seeds used in alpha matting. (d) and (h) are distributions of the alpha maps, taken as the \mathbf{f}_{re} cue.

3.1.3 Reflection

Strong sunlight reflected from shiny objects is another powerful cue. Except for a perfect mirror reflector, sunlight reflection is usually characterized by a brightly lit region in the image where pixels in the region center are brightest and saturated in nearly all color channels. The reflection intensity decays from the center toward the boundary of the reflection region. An example was shown in Figure 7, which compares strong sunlight reflection with the reflection from a white matte/dull object.

We set v_{re} to 1 if white pixels are present in the image and 0 otherwise. To construct \mathbf{f}_{re} , we apply image matting [29] at the detected white pixels. The definite foreground region consists of white pixels, and definite background region consists of a closed curve enclosing the foreground seeds. We then estimate the closed curve under the constraint that the distance between pixels along the curve and enclosed foreground seeds should be larger than a threshold (0.5 in our experiments). This closed curve can be computed by simple dynamic programming. An example was shown in Figure 7(b)–(c).

Given the matting result (e.g., Figure 7(b) and (f)) we plot the alpha matte distributions as shown in (d) and (h), and then assign the 100-bin alpha matte histogram as our 100-D \mathbf{f}_{re} vector.

3.1.4 Contrast

Outdoor images captured in sunny and cloudy days exhibit different global and local saturation contrast. To compute \mathbf{f}_{co} , we utilize contrast information encoded as the percentile in image saturation. For example, a value at the 20th saturation percentile means that 20% of the image pixels are grayer. Clearly, if all saturation percentiles are the same for a given image, the saturation contrast is low. If on the other hand the 50th percentile is at 100 (saturation level) while the 49th percentile is 0, this image is very likely to have a high saturation contrast. In our paper, we use the C channel of LCH color space as our saturation map.

We collect all saturation percentile ratios to build \mathbf{f}_{co} and leave the selection process to the final classifier. Specifically, we denote p_i as the i^{th} percentile in the saturation map. The set of all saturation percentile ratios is given by $\{r | r = p_i/p_j, \forall i > j\}$, where i and j are multiples of 5. We thus obtain 171 percentile ratios in total, which are used to form our 171-D \mathbf{f}_{co} vector. An example is shown in [37].

3.1.5 Haze

Cloudy weather may come with haze. Haze priors have been well studied in computer vision: the dark channel prior presented in [13] is effective. Similarly, we compute the dark channel as

$$\mathcal{J}^k(x) = \min_{r,g,b} \{ \min_{y \in \Omega(x)} \{ \mathcal{J}^c(y) \} \}, \quad (5)$$

where \mathcal{J}^c is a color channel and $\Omega(x)$ is a local patch (with 8×8) centered at x . Most haze-free regions have a low intensity in the dark channel. We measure the haze level and set v_{ha} of a given image as the median value of its dark channel.

We define the \mathbf{f}_{ha} component with the consideration that haze becomes thicker when a region is distant from the camera. These regions commonly exist at the top of an outdoor image. We consider haze location by using a spatial pyramid scheme. The input image is resized into 512×512 . The dark channel in each image is uniformly partitioned into 2^2 , 4^2 , and 8^2 non-overlapping regions to obtain 84 sub-regions. We use the median value of the dark channel intensity in these regions to form the 84-D \mathbf{f}_{ha} vector. An example of haze feature is shown in [37].

3.2 CNN Feature

Our new method in this paper includes the CNN feature which incorporates global discriminative information of the image. We train a CNN model on the two-class image set; the model we used is the AlexNet [22]. As in standard setting [11], the 7th layer neurons of the AlexNet model is extracted to form our 4096D feature \mathbf{f}_{cn} . This feature is effective when the CNN classifier is confident with its prediction. With the confidence scores of sunny and cloudy image s_{cnn} and c_{cnn} ($s_{cnn} + c_{cnn} = 1$), we measure v_{cn} as

$$v_{cn} = 2 \max\{s_{cnn}, c_{cnn}\} - 1, \quad (6)$$

where v_{cn} is in the range of $[0, 1]$. Thus, only discriminable feature of the CNN model contributes to the system.

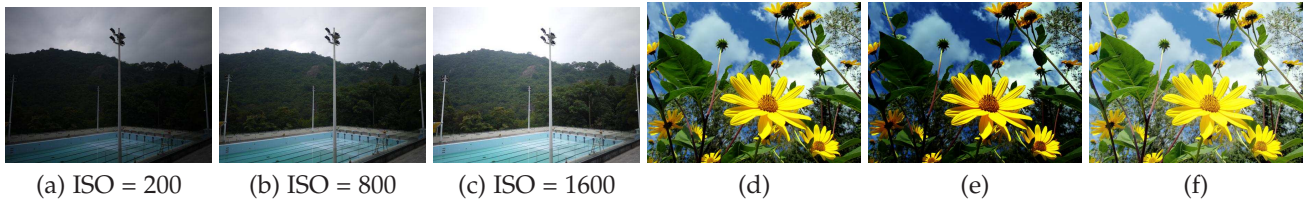


Fig. 8. Camera parameters and photo editing. (b) and (c) are weather counterparts of (a). These weather counterparts are obtained by adjusting camera parameters while capturing the same scene. Here, a high ISO value makes the picture sensitive to light. (e) and (f) are weather counterparts of (d). These weather counterparts are obtained via editing the original photos.

4 WEATHER IMAGES AND COUNTERPARTS DATASET

We created a new weather dataset that contains 10,000 images for training and testing. The dataset and the classifier executable are publicly available. The training images collected from the web were taken by different cameras, under different settings, and might have been edited as well. As most of our proposed weather features are designed based on scene illumination, we describe a new strategy to make our system insensitive to camera settings and photo editing that can be defined by a global transfer function.

A robust weather recognition system should output a uniform weather label for an image as well as its *weather counterpart images*, which are images of the identical scene transformed by a global intensity mapping, and hence they should have the same weather label. Figure 8 shows examples of weather counterpart images.

To make our data collection scalable, we propose to automatically generate weather counterparts for the training data. In the following, we first describe the construction of the weather dataset, and then present a perceptual study to validate the dataset, followed by detailing our learning-based weather counterpart generation.

4.1 Weather Image Dataset

Our weather dataset contains sunny and cloudy images obtained from three sources: Sun Dataset [60], Labelme Dataset [47] and Flickr. The minimum and maximum dimensions of the images are respectively 600 and 1500.

To avoid bias, the helpers recruited to collect and label images were unaware of the purposes or methods used in our experiments. They worked with their own understanding, and collected each 14,000 outdoor images, in which sunny and cloudy images are in equal proportion.

We discarded very similar images by first computing the color histogram distance for all of the image pairs, and then rejected those identical or highly similar. As a result, 1,121 sunny images and 812 cloudy images were rejected. Next, we asked two helpers to independently check the remaining images (5,879 sunny and 6,188 cloudy). Images labeled as ambiguous weather condition by either or both of the helpers were discarded. A total of 5,467 sunny images and 5,612 cloudy images remained after this round. Finally, we asked the third helper to pick 5,000 sunny and 5,000 cloudy images in the final dataset.

	mean (Average)	variance (Average)
Sunny	0.88	0.03
Cloudy	0.85	0.02

TABLE 1
The mean and variance of the mean user-assigned probabilities not equal to 0 or 1.

	Average regression error
LLC [58]	0.442 ± 0.027
ScSPM [62]	0.437 ± 0.031
CNN classifier [22]	0.032 ± 0.002
Ours	0.026 ± 0.003

TABLE 2
Regression error (mean ± variance) of different methods. The momentum, weight decay and learning rate of CNN are 0.0001, 0.9 and 0.005 respectively. The SVM step in LLC and ScSPM is replaced by SVR.

4.2 Perceptual Validation

We first validate our weather image dataset using a user study. A total of 11 participants were recruited for the validation experiment. The participants consisted of 5 males and 6 females whose ages ranged from 26 to 41. All of the subjects reported normal or corrected-to-normal vision with no color-blindness, and reported that they were familiar with the outdoor scenes to be tested in the study. The participants were volunteers who were unaware of the purpose of the experiment.

We asked the participants to assign a sunny/cloudy score to each image where the two assigned scores should sum up to 1. We found that 97.6% of the images were assigned a score of 1 for the designated class by all of the subjects, which indicates that the weather type of most of the images in our dataset are unambiguous.

For the remaining 3.4% images, we compute the mean and variance of the user-assigned “probability scores” across different subjects, and report the average scores in Table 1. If we take 0.5 as the threshold for the “mean probability,” the user-assigned weather type has 100% accuracy. We also regress these user-assigned probability scores, and the regression errors are shown Table 2. Note on the one hand in the evaluation in the following sections, we do not use these user-assigned probabilities as the measurement metric due to the fact that the assigned scores are subjective, although the average score across different subjects is used here. On the other hand, we believe this may lead to interesting and worthwhile future work on using user-assigned probability scores (i.e., user’s observation) in performance evaluation.

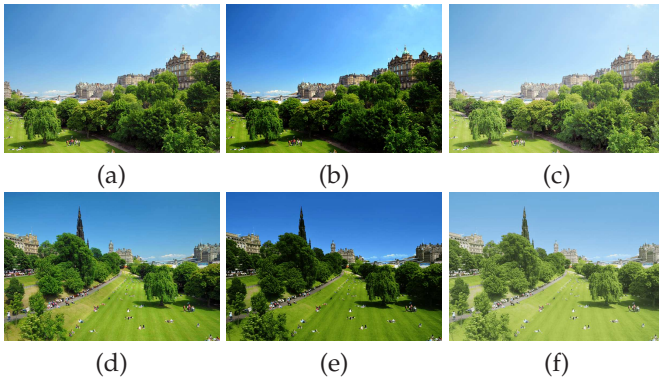


Fig. 9. Examples of weather counterpart mapping. (b) and (c) are weather counterparts of (a). We build the mapping function from (a) to (b)-(c) respectively denoted as \mathcal{H}_1 and \mathcal{H}_2 . Given the testing image (d), (a) is among those visually similar to (d) in the weather counterpart dataset. (e) and (f) are respectively the mapped results of \mathcal{H}_1 and \mathcal{H}_2 with input (a).

4.3 Weather Counterpart Image Generation

Now, we build a dataset for learning weather counterpart images. We capture 1,000 outdoor images with different camera parameters. To guarantee pixel-wise alignment, a tripod was used during capture. For each image, 4 weather counterpart images were captured using different camera parameters. The ISO setting is the most important camera parameter and we found it sufficient to capture a large set of tone variances in different cameras.

In addition, we recruited seven helpers to edit the 1,000 outdoor images using Photoshop. Five of them were without computer vision background. To avoid editing bias, these subjects were given no specific instruction except that their editing should not change the weather label of the image. The editing operation mainly includes gamut mapping, tone mapping, sharpening, blurring, etc. For each image, 4 edited versions are produced.

Therefore, for each of the two cases we have 5 (1 original + 4 additional) images, which are ordered differently to form 20 image pairs after permutation. Finally we obtain a total of $(1000 + 1000) \times 20$ image pairs, which are named *weather counterpart image pairs*.

Weather Counterpart Mapping Functions We now learn the mapping relationship from a given image to its weather counterpart. We assume the RGB color of an outdoor image (denoted as I) and its weather counterpart image (denoted as I^o) has the mapping relationship expressed as

$$\{r^o, g^o, b^o\} = \mathcal{H}(r, g, b), \quad (7)$$

where $\{r, g, b\}$ and $\{r^o, g^o, b^o\}$ are the input and mapped RGB color vectors. Eq. (7) has only 3 variables, so we can build a 3D array to turn this mapping operator into a table look-up operation.

The table construction is as follows. The intensity in each variable (channel) is quantized into 256 bins, which produces the mapping table with dimension $256 \times 256 \times 256$. We define \mathbf{v}_p and \mathbf{v}_p^o as the quantized RGB vector of pixel

color at p in I and I^o respectively. For an input (r, g, b) , by defining $\mathbf{c} = [r, g, b]^T$, the output of the mapping table is

$$\mathcal{H}(r, g, b) = \frac{1}{K} \sum_{p \in \mathcal{D}(\mathbf{c})} \{\exp[-\frac{1}{\sigma^2} \|\mathbf{v}_p - \mathbf{c}\|_2^2] \mathbf{v}_p^o + \gamma \mathbf{c}\}, \quad (8)$$

where

$$\mathcal{D}(\mathbf{c}) = \{q \mid \|\mathbf{v}_q - \mathbf{c}\|_2^2 \leq \delta, \forall q \in \Theta\}. \quad (9)$$

Here Θ is the pixel set of the image. The K in Eq. (8) is a normalization factor, which is given by

$$K = \sum_{p \in \mathcal{D}(\mathbf{c})} \exp[-\frac{1}{\sigma^2} \|\mathbf{v}_p - \mathbf{c}\|_2^2] + \gamma. \quad (10)$$

In our experiments, we set $\delta = 45$, $\gamma = 0.1$ and $\sigma = 5$. The output of $\mathcal{H}(r, g, b)$ is the weighted combination of pixel-pairs mapping whose input RGB vector is close to $[r, g, b]^T$. Therefore, the mapping relationship of I and I^o can be smoothly transferred to table \mathcal{H} . If we cannot find sufficient number of RGB vectors close to $[r, g, b]^T$ in the original image, $\sum_{p \in \mathcal{D}(\mathbf{c})} \exp[-\frac{1}{\sigma^2} \|\mathbf{v}_p - \mathbf{c}\|_2^2]$ is small, so that we trust more the original input \mathbf{c} in this case.

Weather Counterparts Generation We produce 40,000 mapping functions for all 40,000 weather counterpart pairs in our dataset. Each function captures a color transform from the first image to the second in the corresponding weather counterpart pair. Given a training image, we use the mapping where the input images are similar to produce the weather counterpart images for training.

We use the color GIST descriptors [8] to extract the color and context information. We pick $d = 50$ mapping functions whose input image is closest to the given image according to color GIST features, and then use the d mapping functions to produce d weather counterpart images. Figure 9 shows examples of weather counterpart images.

5 COLLABORATIVE LEARNING WITH HOMOGENEOUS VOTERS

Traditional classifiers such as SVM cannot achieve good performance on our overall weather feature because they assume all of the components are present simultaneously in every image, which may unfortunately not be the case. For example, outdoor images do not always contain the sky region. Images lacking one or more weather cues would significantly affect SVM's classification performance.

Our learning strategy is to partition the training images into disjoint clusters of homogeneous voters, so that voters closer to a given testing image have more weights when the weather label is considered.

5.1 Voting Scheme

Our training outdoor images are first partitioned into homogeneous clusters according to the existence vector of each image as defined in Eq. (2). The partitioned sets thus correspond to different weather cue patterns, such as "reflection + shadow", "sky + haze", and "sky + reflection + shadow". Images in the same cluster/pattern are the homogeneous.

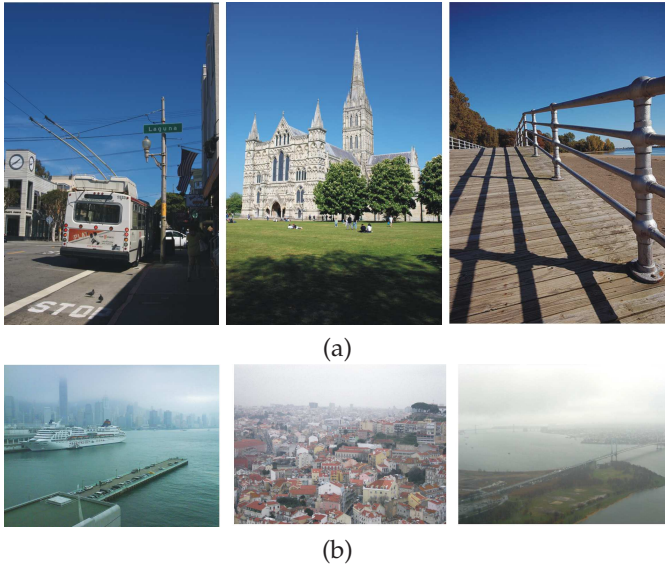


Fig. 10. Sample images found in two clusters. (a) “sky + shadow” cluster with center $\{0.90, 0.87, 0.26, 0.11\}$. (b) “sky + haze” cluster with center $\{0.94, 0.24, 0.27, 0.84\}$, where $\{v_{sk}, v_{sh}, v_{re}, v_{ha}\}$ is composed of the respective existence scores.

In implementation, we partition the set of training images into M subsets $\{\Omega_1, \dots, \Omega_M\}$ based on the existence vectors using hierarchical clustering [10]. We set the cluster error threshold to 0.5 in terms of Euclidean distance. M can be found automatically. We denote the set of cluster center vectors as $\{\hat{e}_1, \dots, \hat{e}_M\}$. Figure 10 shows sample images of two converged clusters and their cluster centers.

In the testing phase, given an overall weather feature \mathbf{x} with existence vector e , the training data whose existence vectors are similar to e should be used. So our classifier is implemented using a weighted voting scheme, expressed as

$$h(\mathbf{x}, e) = \text{sign}[\sum_{i=1}^M s(\hat{e}_i, e) \hat{h}_i(\mathbf{x})], \quad (11)$$

where $\text{sign}[\cdot]$ is the function outputting 1 (resp. -1) for non-negative (resp. negative) input, $s(\hat{e}_i, e)$ is a similarity function under parameter σ :

$$s(\hat{e}_i, e) = \frac{\exp(-\frac{\|\hat{e}_i - e\|_2^2}{2\sigma^2})}{\sum_i^M \exp(-\frac{\|\hat{e}_i - e\|_2^2}{2\sigma^2})}, \quad (12)$$

and $\hat{h}_i(\cdot)$ (defined shortly) is the homogeneous voter trained using the data in Ω_i . Our classifier Eq. (11) gives a larger weight to the homogeneous voter whose existence vector pattern is similar to that of the testing data.

5.2 Collaborative Learning

For training image i , we denote the overall weather feature as \mathbf{x}_i , and the weather label as $y_i \in \{-1, +1\}$, where -1 and +1 correspond respectively to “cloudy” and “sunny”. For each homogeneous voter, we model $\hat{h}_i(\cdot)$ as

$$\hat{h}_i(\mathbf{x}) = \text{sign}(\sum_{j=1}^p \omega_{j,i} \mathbf{x}(j) + b_i), \quad (13)$$

where $\mathbf{x}(j)$ is the j^{th} element of vector \mathbf{x} . If each homogeneous voter works independently without information sharing, the classifier in Eq. (13) can be modeled as a standard SVM [5], expressed as

$$\begin{aligned} \min_{\omega_{j,i}, b_i, \zeta_{i,k}} \quad & \sum_{j=1}^p \omega_{j,i}^2 + C \sum_{k \in \Omega_i} \zeta_{i,k} \\ \text{s.t.} \quad & y_k (\sum_{j=1}^p \omega_{j,i} \mathbf{x}_k(j) + b_i) \geq 1 - \zeta_{i,k}, \zeta_{i,k} \geq 0, \forall k \in \Omega_i, \end{aligned} \quad (14)$$

where $p = 4717$ is the dimension of the overall weather feature and C is a constant.

In our framework, we do not train each $\hat{h}_i(\mathbf{x})$ independently because this will lead to a large bias. Our voters work collaboratively to determine the classification result and we optimize them together in a unified framework.

By removing sign from $\hat{h}_i(\mathbf{x})$, we make the system linear, which updates Eq. (11) into

$$h(\mathbf{x}, e) = \text{sign}[\sum_{i=1}^M s(\hat{e}_i, e) (\sum_{j=1}^p \omega_{j,i} \mathbf{x}_k(j) + b_i)]. \quad (15)$$

We make this change because a voter should not be restricted to output binary values. This also helps to indicate ambiguous situation where sunny and cloudy features are present at the same time, see Figure 18.

5.3 Latent SVM Learning

Now, for each training sample, we produce d weather counterpart images and extract the 4717-D weather feature on them. Given training image t , we denote the weather feature of its l^{th} weather counterpart image as \mathbf{x}_t^l ($l = \{1, \dots, d\}$). We also define $\mathbf{x}_t^0 = \mathbf{x}_t$.

For each training sample, we require that its weather counterparts to have the same weather label in the training phase. This requires us to prevent all the weather counterpart features from falling into the margin during the training stage. To this end, we define a latent variable to indicate which weather counterpart image can produce the minimum classification margin, and encourage the minimum margin to be large during optimization. According to the max-margin strategy, we can write the constraints as

$$\begin{aligned} \min_{c \in \{0, \dots, d\}} \quad & \{y_k (\sum_{j=1}^p \omega_{j,i} \mathbf{x}_k^c(j) + b_i)\} \geq 1 - \zeta_{i,k}, \\ & \zeta_{i,k} \geq 0, \forall k \in \Omega_i, \forall i = 1, \dots, M, \end{aligned} \quad (16)$$

and

$$\begin{aligned} \min_{m \in \{0, \dots, d\}} \quad & \{y_t [\sum_{i=1}^M s(\hat{e}_i, e_t) (\sum_{j=1}^p \omega_{j,i} \mathbf{x}_t^m(j) + b_i)]\} \geq 1 - \xi_t, \\ & \xi_t \geq 0, \forall i = 1, \dots, M, \forall t = 1, \dots, N \end{aligned} \quad (17)$$

where N is the number of training images, c and m are two latent variables to indicate which weather counterpart image produces the minimum classification margin.

Denoting the latent variables as $\mathcal{C} = \{c(1), \dots, c(N)\}$ and $\mathcal{M} = \{m(1), \dots, m(N)\}$ for Eqs. (16) and (17) respectively, the final objective function for $h(\mathbf{x}, e)$ is written as

$$\min_{\omega_{j,i}, b_i, \xi_t, \zeta_{i,k}, \mathcal{C}, \mathcal{M}} \sum_{i=1}^M \sum_{j=1}^p \omega_{j,i}^2 + C_1 \sum_{i=1}^M \sum_{k \in \Omega_i} \zeta_{i,k} + C_2 \sum_{t=1}^N \xi_t \quad (18)$$

s.t.

$$\min_{c(k) \in \{0, \dots, d\}} \{y_k (\sum_{j=1}^p \omega_{j,i} x_k^{c(k)}(j) + b_i)\} \geq 1 - \zeta_{i,k},$$

$$\zeta_{i,k} \geq 0, \forall k \in \Omega_i, \forall i = 1, \dots, M \quad (19)$$

$$\min_{m(t) \in \{0, \dots, d\}} \{y_t [s(\hat{e}_i, e_t) (\sum_{j=1}^p \omega_{j,i} x_t^{m(t)}(j) + b_i)]\} \geq 1 - \xi_t$$

$$\xi_t \geq 0, \forall i = 1, \dots, M, \forall t = 1, \dots, N \quad (20)$$

where C_1 and C_2 are constants.

Eq. (18) can be regarded as latent SVM whose latent variables are \mathcal{C} and \mathcal{M} . We solve Eq. (18) by iteratively optimizing $\{\omega_{j,i}, b_i, \xi_t, \zeta_{i,k}\}$ as a standard SVM problem, and optimize the latent values $\{\mathcal{C}, \mathcal{M}\}$ in the constraining conditions. The following steps are adopted:

- 1) Keep $\{\omega_{j,i}, b_i, \xi_t, \zeta_{i,k}\}$ fixed, optimize the latent $\{\mathcal{C}, \mathcal{M}\}$ subject to the constraints (19) and (20). This is a simple minimization operation.
- 2) Keep $\{\mathcal{C}, \mathcal{M}\}$ fixed, optimize $\{\omega_{j,i}, b_i, \xi_t, \zeta_{i,k}\}$ by solving a standard SVM problem which can be solved using Lagrange multipliers [5].

For the second step, voter collaboration is characterized by Eq. (20), which forces all of the voters to work together in the classification. The effectiveness of each voter is governed by Eq. (19). It guarantees that each voter is learned from its corresponding homogeneous data. Eqs. (19) and (20) can accomplish good classification performance. We solve Eq. (18) using different σ s for $s(\cdot)$ in Eq. (12). In the final stage, we pick the σ with the minimum energy for the objective function (18).

Similar to other latent SVM solvers, the two-step iteration converges to a satisfactory $\{\omega_{j,i}, b_i\}$ in our experiments. This is because the difference among feature weather counterpart images is much smaller than the difference among all of the training samples. In the two-step iteration, compared to the first step (latent variables optimization), the second step plays a more important role in model update with SVM weights. The latent variable optimization can be regarded as a fine-tuning step in the feature space, in order to capture the variance introduced by different camera parameters and image editing. Empirically, the update stops after 8 – 10 iterations.

6 EXPERIMENTS

We report the classification results under different evaluation settings, and further validate our dataset using a perceptual study. In the following, the *CNN feature* is the feature extracted from the CNN model used as input to our system. The *CNN classifier* refers to the CNN model trained end-to-end on our dataset. The visualization of the learned convolutional filters can be found in the supplementary material.

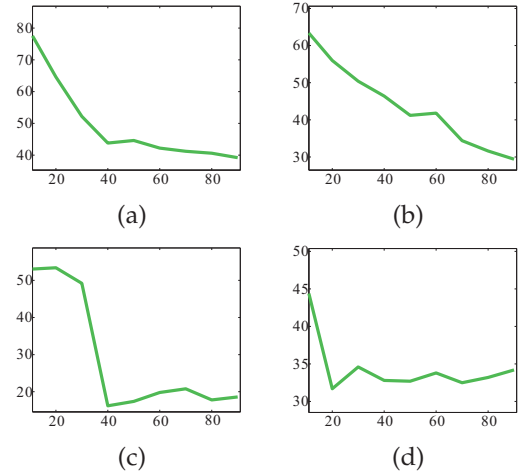


Fig. 11. (a)–(d) are respectively the performance curves of sky, shadow, reflection, and haze classifiers. The x -axis values are the respective percentages of selected images (with the highest existence score) in the dataset. The y -axis is the classification accuracy (in %).

6.1 Classification Results

The training and classification were done using the weather dataset constructed. We adopted the cross validation scheme where in each round, 80% of the data were selected randomly as the training set, with the remaining 20% as the testing set. We ran 5 rounds of experiment and recorded the mean and variance of the classification accuracy.

On two-class labeling, even random guess can reach 50% accuracy. We use the normalized accuracy given by $\max\{(a - 0.5)/(1 - 0.5), 0\}$, where a is the raw accuracy obtained. Thus, the normalized accuracy is within $[0, 1]$ and random guess is expected to get zero.

6.1.1 Individual Features and Scores

We use SVM to evaluate individual weather features. Note that it may not be fair for the CNN feature whose performance is sensitive to choice of parameters. Nevertheless we include the CNN feature in this section, and defer the discussion of its performance under different parameter settings in a later section. The momentum, weight decay and learning rate of the CNN are 0.0001, 0.9 and 0.005 respectively to produce the best performance.

Table 3 tabulates the classification results. Although the overall weather feature vector (4717D) is lopsided to the CNN feature (4096D), and that using the CNN feature alone reports better performance than any of the individual non-CNN features, we shall show that the combination of CNN and non-CNN features reports the best performance by taking the advantages of both features. Not surprisingly, the sky is the most important weather cue among the five non-CNN features. We believe that this is due to the fact that sky detection is relatively easier and more stable. The majority of failure cases are images without a prominent sky region. In addition, the reflection and shadow classifiers also work well. The performance of the contrast classifier on the other hand depends on the complexity of the scene.

We note that the haze cue is weaker than the sky and contrast cues mainly due to the fact that many images in our

Feature	Normalized accuracy
Sky	44.3 \pm 1.9
Shadow	41.1 \pm 2.2
Reflection	27.0 \pm 2.1
Contrast	40.5 \pm 2.0
Haze	34.1 \pm 1.9
CNN feature	83.6 \pm 2.0

TABLE 3

Classification results (mean \pm variance) using individual features.

Feature	Normalized accuracy
Sky	86.3 \pm 1.9
Shadow	86.9 \pm 2.1
Reflection	89.0 \pm 2.0
Contrast	86.2 \pm 1.9
Haze	87.7 \pm 1.5
CNN feature	61.4 \pm 2.0
All	91.4 \pm 1.6

TABLE 4

Classification results (mean \pm variance) with individual weather cues being left out. "All" means individual weather cues are included.

Classifier	Normalized accuracy
CNN classifier (without DA)	77.8 \pm 2.0
CNN classifier (with DA)	83.3 \pm 2.1
Ours (without DA)	84.0 \pm 1.8
Ours (with DA)	91.4 \pm 1.6

TABLE 5

Classification results (mean \pm variance) with and without data augmentation; "with DA" and "without DA" respectively stand for training with and without data augmentation. The learning rate, momentum and weight decay of the CNN are respectively 0.0001, 0.9 and 0.005.

	Normalized accuracy
SVM	41.2 \pm 2.2
Adaboost	36.4 \pm 2.3
LLC [58]	0.3 \pm 0.1
ScSPM [62]	0.2 \pm 0.1
CNN classifier [22]	83.3 \pm 1.8
Ours	91.4 \pm 1.6

TABLE 6

Classification results (mean \pm variance) of different methods. The momentum, weight decay and learning rate of the CNN are 0.0001, 0.9 and 0.005 respectively.

dataset simply do not exhibit detectable haze. To confirm this, we select 415 images with haze v_{ha} score larger than 0.7 and 415 sunny images. The haze classifier performance is improved up to 84.2% in normalized accuracy when it is applied to these 830 images. We also found that the haze cue can help identify sunny images as well in classification, since many sunny images have vivid color which exhibits low dark-channel intensities.

Next, we evaluate individual existence scores, which are used to form Eq. (2). For each individual feature, we select s percent of the images with the highest existence score in the dataset, and apply SVM classification on this image subset. Figure 11 shows the performance with varying s of each individual classifier. The plot indicates that our existence score design is effective – each individual feature is more useful when it has a higher existence score.

6.1.2 Ablation Study

We also conducted an ablation study, that is, to study the performance of the system when a given weather cue is left out. Table 4 verifies that all of the proposed cues are useful in accounting for the overall weather classification performance, since removing any one of them results in a performance drop. In particular, the momentum, weight decay and learning rate of our CNN are 0.0001, 0.9 and 0.005 respectively. We note that the CNN feature plays an important role among all of the features due to the larger drop in comparison to other features.

6.1.3 Data Augmentation

We enrich the data set by introducing weather counterpart images in the training phase. In this section we compare the classification results on training using the data with and without data augmentation. Table 5 shows that our data augmentation leads to about 7% improvement in normalized accuracy. Note that we use the best parameters for the CNN classifier and our classifier.

Therefore, for the rest of our evaluation in this section we will use the augmented training dataset.

6.1.4 Comparison

We report our overall classification performance compared with typical baseline systems, weather related systems, and the CNN classifier. For the latter, we will show and explain that relying on the CNN without the proposed weather-specific features will result in a significant performance drop, despite that CNN is powerful in encoding global scene structure and characteristics.

Comparison with Baseline Systems The first baseline is to implement SVM directly on the 4717-D weather feature. We test both the linear and non-linear versions with different kernels and report the results with the best performance. The second baseline is the traditional Adaboost, which combines several classifiers to build a stronger one. We take each feature bin as a weak classifier. Another two baseline methods based on dictionary learning [39] are typical image classification methods, namely LLC [58] and ScSPM [62].

For the SVM baseline we tried different parameters: $C \in [0.0001, 100]$ and we report the best one ($C = 0.05$). We tried different non-linear kernels, including the Gaussian kernel, RBF kernel and Polynomial kernel. The best kernel is the polynomial kernel which produces a performance similar to the linear solver. For the Adaboost baseline, we took each single dimension in the 1645-D feature as a weak classifier, and tried 50 different permutations of the weak classifiers. We found that the variance in accuracy is very small – 0.16 only. For LLC and ScSPM, we use the codes provided in [58] and [62] with default parameters. We also tried three dictionary sizes, 512, 1024 and 2048. We found all of the parameters cannot yield a result significantly greater than 0 normalized accuracy.

Table 6 lists the classification results. Figures 12 and 13 show a few examples, where we test 5 different σ values in Eq. (12), that is, $\{0.5, 0.1, 0.01, 0.05, 0.001\}$, and select the best result with the lowest energy in Eq. (18).

For traditional image classification methods LLC [58] and ScSPM [62], the normalized accuracies are close to 0. This is because these methods rely on scene structure and do not consider illumination information. SVM and Adaboost



Fig. 12. Detection results: cloudy images.



Fig. 13. Detection results: sunny images.

	Normalized accuracy
Lalonde <i>et al.</i> [26]	46.5 ± 1.7
Yan <i>et al.</i> [61]	24.6 ± 2.6
Roser and Moosmann [46]	26.2 ± 2.3
Laffont <i>et al.</i> [24]	21.4 ± 1.9
Ours	91.4 ± 1.6

TABLE 7
Classification statistics of different methods.

do not yield significant improvement over single weather classifiers, such as those of sky or shadow (c.f. Tables 3 and 6). We also find that the use of kernel SVM yields similar performance.

Comparison with Related Methods We also compare our classifier with weather-related methods. The first related work is Lalonde *et al.* [26]. Note that their system is not designed for weather classification; one component, namely, the sun visibility prediction, can be regarded as a coarse weather estimator. We implemented this component and tested it on our dataset. Another two vehicle-based weather classifiers [61], [46] were also compared. The comparison with [24] is also provided. It can output 40 attributes of image feature including “sunny” and “cloudy”. We train the regressor on our data.

Table 7 tabulates the classification statistics. For the method of [26], the assumption that an outdoor scene is composed of ground, sky, and vertical surfaces may not be satisfied (see a few exceptions in Figures 12 and 13). For the work of [61], [46], the weather estimators are specially designed for driver assistance. They rely on vehicle-mounted image priors, which cannot properly deal with general natural images.

The framework proposed in [24], where off-the-shelf but not weather-specific features are used, produces less effective results. But it still outperforms [61], [46] since it does not rely on on vehicle-mounted priors. We also observe that the precisions for “cloudy” and “sunny” attributes reported on the dataset of [24] is high (> 0.95), because the testing images in [24] are dominated by a prominent sky region which is relatively easy to recognize.

Comparison with CNN In this section, we train a CNN classifier end-to-end to perform the two-class classification.

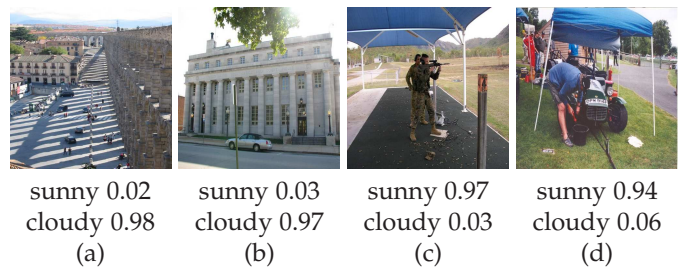


Fig. 14. Four CNN failure cases. In (a) and (b), Sunny image mis-detected as cloudy. In (c) and (d) cloudy image mis-detected as sunny. “sunny score” and “cloudy score” refer to output CNN confidence score in the sunny and cloudy class.

The standard structure provided by [22] is adopted, and we fine tune on the AlexNet. We report the result with the best parameter setting.

For the sake of fairness, the CNN feature in our classifier adopts the same set of parameters with the one being compared. The results are reported with various hyperparameters.

First, in Table 8 we show the performance under different learning rates; the most important parameter for CNN training. The other two important parameters are namely the momentum and weight decay. Table 9 summarizes their effects in the comparison experiments. We found that the proposed non-CNN weather features are complementary with the CNN feature in the overall weather feature vector. That is, while the CNN feature is capable of capturing global image characteristics, it may not encode well weather characteristics which are better represented by non-CNN weather cues. The result shows that our system which combines CNN feature and non-CNN features leads to about 8% improvement.

To further verify this point, we manually label non-CNN weather cues to explore their full potentials. For the shadow feature, we manually label shadow regions in all of the images using bounding boxes, followed by shadow extraction [25] within the box region. The existence score is 1 (with shadow) or 0 (without shadow). For the reflection feature, we label reflection regions manually with a 1-dimensional feature. If reflection occurs the feature bin is 1; otherwise 0.

learning rate	CNN classifier	Ours
0.001	81.5 ± 1.9	88.7 ± 2.0
0.0001	83.3 ± 1.8	91.4 ± 1.6
0.00001	82.0 ± 2.0	89.5 ± 1.9

TABLE 8

Comparison with CNN classifier under different learning rates, given momentum and weight decay are respectively 0.9 and 0.005.

parameters	CNN classifier	Ours
$\nu = 0.9, \xi = 0.005$	83.3 ± 1.8	91.4 ± 1.6
$\nu = 0.9, \xi = 0.0025$	81.4 ± 1.8	89.6 ± 1.7
$\nu = 0.45, \xi = 0.005$	77.1 ± 2.1	85.2 ± 1.8
$\nu = 0.45, \xi = 0.0025$	78.6 ± 1.9	86.7 ± 2.0

TABLE 9

Comparison with CNN classifier, given learning rate is 0.0001, ν and ξ are respectively momentum and weight decay.

	Normalized accuracy
CNN classifier	83.3 ± 1.8
Ours (Sky)	93.5 ± 1.8
Ours (Shadow)	96.3 ± 1.7
Ours (Reflection)	95.2 ± 1.6
Ours (Sky + Shadow + Reflection)	97.4 ± 1.2

TABLE 10

Comparison with manual feature localization. In the table, the feature inside the parentheses are manually labeled as described in the text. The learning rate, momentum and weight decay for the CNN classifier are respectively 0.0001, 0.9 and 0.005.

The existence score is the same as the feature bin value. For the sky feature, with a few exceptions most of the sky regions can be correctly segmented, and we manually label the missing sky regions. Table 10 tabulates the results, which demonstrates that the proposed features are indeed effective in weather recognition and perform better than the CNN classifier. Thus, we believe that with the continuing improvement of low-level vision techniques, the proposed non-CNN weather cues will improve the overall performance significantly by working in synergy with the CNN feature. Figure 14 shows some fail cases of the CNN classifier which can otherwise be correctly labeled by our method. For example, Figure 14(a) and (b) are recognized as cloudy by the CNN classifier due to its globally gray color tone. But our system can look into more details such as shadows to produce the correct weather label. In Figure 14(c) and (d), although shadows are found, they are not strong cast shadow caused by the sun. On the other hand, the small sky region can be correctly detected by our method.

6.2 Comparison on Two-Class Laffont Dataset

In [24], a dataset including 40 transient attributes was proposed. In this dataset, 8,571 images in total from 101 webcams are annotated by crowd-sourcing. “Sunny” and “cloudy” are two of the attributes in their setting. That is, for each image, confidence score of “sunny” and “cloudy” are available. As suggested in [24], attributes with confidence score larger than 0.8 is considered strong positive attributes. The images with strong positive sunny/cloudy attributes are selected to form the two-class Laffont dataset. We found in this dataset none of the images is labeled both “sunny” and “cloudy”, so the images are unambiguous. The resulting dataset contains 1,729 cloudy images and 1,085 sunny images. Figure 15 shows example images of the *two-class*



Fig. 15. Examples of the two-class Laffont dataset.

	Normalized Accuracy
[24]	92.4 ± 1.3
CNN classifier	96.4 ± 0.5
Ours without CNN feature	98.2 ± 0.1
Ours	98.6 ± 0.2

TABLE 11

Comparison on the two-class Laffont Dataset

Laffont dataset. We apply different methods on the dataset by using 80% of it for training and 20% for testing, and report the cross-validation results. Table 11 shows that our method has a better performance than [24] on the two-class dataset. In particular, note that the performance of our method and the CNN classifier are close to 100%. This is because the images in the dataset typically have a prominent sky region which makes the classification easier.

6.3 Application in Weather Monitoring using Surveillance Cameras

We verify our technique in a real-world application: real-time weather monitoring. Surveillance cameras can be found almost everywhere, so we believe running our fast and cost-effective method on these cameras can effectively monitor real-time weather in urban areas. This is particularly useful where solar panels are extensively installed on the rooftops of many buildings; sunny and cloudy weather provides important guidance for optimizing power transfer in the main grids. That is, when cloudy weather is detected for an extended period of time, the main power grid may start to take over early to maintain stable power supply and avoid outage. We apply our weather predictor in surveillance videos footage, and collected 2000 surveillance images, 1000 of them are sunny and the other 1000 are cloudy. The normalized classification accuracy is 93.2%. Sample images are shown in Figure 16.

7 CONCLUSION AND FUTURE WORK

We have presented a learning-based approach for classifying two types of weather. This apparently simple two-class weather labeling problem is not trivial given the great variety of outdoor images. The feature cues we used resonate well with our own common sense in judging weather conditions. Because some of the feature cues may be unavailable in images, the key to our computational framework is a collaborative learning strategy where voters closer to the testing image in terms of weather information/structure are given more weight in classification. We have also incorporated the powerful CNN feature into our overall weather feature. To resist variations caused by different camera parameters and photo editing, a latent SVM framework is proposed to learn from various synthesized weather counterpart images. Our experimental results showed that this is an effective strategy, which we believe has good potential beyond weather classification.

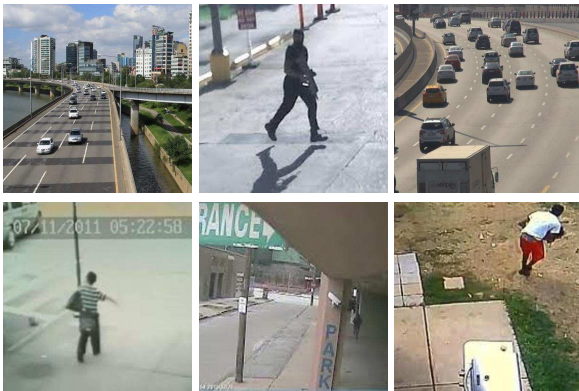


Fig. 16. Some detection results for surveillance images. The first and second row are sunny and cloudy images respectively



Fig. 17. (a) Sunny image mis-detected as cloudy, and (b) cloudy image mis-detected as sunny.

Our current approach is limited to label two weather types. More research needs to be engaged in generalizing the approach to labeling more conditions on larger dataset [4]. For example, Figure 18 shows two images where sunny and cloudy features are present at the same time. They may be labeled as “partly sunny” or “partly cloudy” and in fact, our system labels (a) as sunny, with the rescaled SVM sunny score 0.641 (and cloudy score 0.359), while labeling (b) as cloudy with the rescaled SVM cloudy score 0.716 (and sunny score 0.284), which we believe are reasonable for two-class weather classification.

We hope this paper will spark interest and subsequent work along this line of research. Executable and the weather dataset are available at the project website.

8 ACKNOWLEDGMENTS

The research was supported by the Research Grants Council of the Hong Kong Special Administrative Region (Project nos 619313 and 412911).

REFERENCES

- [1] M. Baccouche, F. Mamalet, C. Wolf, C. Garcia, and A. Baskurt. Sequential deep learning for human action recognition. In *Human Behavior Understanding*, 2011.
- [2] R. Baltenberger, M. Zhai, C. Greenwell, S. Workman, and N. Jacobs. A fast method for estimating transient scene attributes. In *2016 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1–8, 2016.
- [3] Y.-L. Boureau, F. Bach, Y. LeCun, and J. Ponce. Learning mid-level features for recognition. In *CVPR*, 2010.



Fig. 18. Sunny or cloudy?

- [4] W. T. Chu, X. Y. Zheng, and D. S. Ding. Image2weather: A large-scale image dataset for weather property estimation. In *2016 IEEE Second International Conference on Multimedia Big Data (BigMM)*, pages 137–144, April 2016.
- [5] C. Cortes and V. Vapnik. Support-vector networks. *Machine learning*, 1995.
- [6] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *CVPR*, 2005.
- [7] K. Derpanis, M. Lecce, K. Daniilidis, and R. Wildes. Dynamic scene understanding: The role of orientation features in space and time in scene classification. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012.
- [8] M. Douze, H. Jégou, H. Sandhawalia, L. Amsaleg, and C. Schmid. Evaluation of gist descriptors for web-scale image search. In *ACM International Conference on Image and Video Retrieval*, 2009.
- [9] M. Elhoseiny, S. Huang, and A. Elgammal. Weather classification with deep convolutional neural networks. In *IEEE International Conference on Image Processing (ICIP)*, pages 3349–3353, 2015.
- [10] M. Fionn. A survey of recent advances in hierarchical clustering algorithms. *The Computer Journal*, 1983.
- [11] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014.
- [12] D. Glasner, P. Fua, T. Zickler, and L. Zelnik-Manor. Hot or not: Exploring correlations between appearance and temperature. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [13] K. He, J. Sun, and X. Tang. Single image haze removal using dark channel prior. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009.
- [14] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *arXiv preprint arXiv:1512.03385*, 2015.
- [15] M. Islam, N. Jacobs, H. Wu, and R. Souvenir. Images+ weather: Collection, validation, and refinement. In *IEEE Conference on Computer Vision and Pattern Recognition Workshop on Ground Truth*, volume 6, page 2, 2013.
- [16] N. Jacobs, M. T. Islam, and S. Workman. Cloud motion as a calibration cue. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1344–1351, 2013.
- [17] M. Juneja, A. Vedaldi, C. Jawahar, and A. Zisserman. Blocks that shout: Distinctive parts for scene classification. In *CVPR*, 2013.
- [18] H. Kang, M. Hebert, and T. Kanade. Discovering object instances from scenes of daily living. In *ICCV*, 2011.
- [19] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and L. Fei-Fei. Large-scale video classification with convolutional neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014.
- [20] H. Katsura, J. Miura, M. Hild, and Y. Shirai. A view-based outdoor navigation using object recognition robust to changes of weather and seasons. In *IROS*, 2003.
- [21] G. Kim and A. Torralba. Unsupervised detection of regions of interest using iterative link analysis. 2009.
- [22] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems (NIPS)*, 2012.
- [23] H. Kurihata, T. Takahashi, I. Ide, Y. Mekada, H. Murase, Y. Tamatsu, and T. Miyahara. Rainy weather recognition from in-vehicle camera images for driver assistance. In *Intelligent Vehicles Symposium, 2005. Proceedings. IEEE*, 2005.
- [24] P.-Y. Laffont, Z. Ren, X. Tao, C. Qian, and J. Hays. Transient attributes for high-level understanding and editing of outdoor scenes. *ACM Transactions on Graphics (TOG)*, 33(4):149, 2014.
- [25] J.-F. Lalonde, A. Efros, and S. Narasimhan. Detecting ground shadows in outdoor consumer photographs. In *European Conference on Computer Vision (ECCV)*, 2010.
- [26] J.-F. Lalonde, A. Efros, and S. Narasimhan. Estimating the natural illumination conditions from a single outdoor image. *IJCV*, 2012.

- [27] S. Lazebnik, C. Schmid, and J. Ponce. Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories. In *CVPR*, 2006.
- [28] Y. J. Lee and K. Grauman. Object-graphs for context-aware category discovery. In *CVPR*, 2010.
- [29] A. Levin, D. Lischinski, and Y. Weiss. A closed form solution to natural image matting. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2006.
- [30] L.-J. Li, H. Su, L. Fei-Fei, and E. P. Xing. Object bank: A high-level image representation for scene classification & semantic feature sparsification. In *NIPS*, 2010.
- [31] L.-J. Li, H. Su, Y. Lim, and L. Fei-Fei. Objects as attributes for scene classification. In *ECCV*, 2012.
- [32] Q. Li, J. Wu, and Z. Tu. Harvesting mid-level visual concepts from large-scale internet images. In *CVPR*, 2013.
- [33] D. Lin, C. Lu, R. Liao, and J. Jia. Learning important spatial pooling regions for scene classification. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014.
- [34] D. Lin, X. Shen, C. Lu, and J. Jia. Deep lac: Deep localization, alignment and classification for fine-grained recognition. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [35] M. Lin, Q. Chen, and S. Yan. Network in network. *arXiv preprint arXiv:1312.4400*, 2013.
- [36] D. G. Lowe. Distinctive image features from scale-invariant keypoints. *IJCV*, 2004.
- [37] C. Lu, D. Lin, J. Jia, and C.-K. Tang. Two-class weather classification. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014.
- [38] C. Lu, J. Shi, and J. Jia. Online robust dictionary learning. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2013.
- [39] C. Lu, J. Shi, and J. Jia. Scale adaptive dictionary learning. *IEEE Transactions on Image Processing (TIP)*, 2013.
- [40] S. G. Narasimhan and S. K. Nayar. Shedding light on the weather. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2003.
- [41] A. Oliva and A. Torralba. Modeling the shape of the scene: A holistic representation of the spatial envelope. *IJCV*, 2001.
- [42] M. Pandey and S. Lazebnik. Scene recognition and weakly supervised object localization with deformable part-based models. In *ICCV*, 2011.
- [43] S. Parizi, J. Oberlin, and P. Felzenszwalb. Reconfigurable models for scene recognition. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012.
- [44] F. Perronnin, Y. Liu, J. Sánchez, and H. Poirier. Large-scale image retrieval with compressed fisher vectors. In *CVPR*, 2010.
- [45] A. Quattoni and A. Torralba. Recognizing indoor scenes. In *CVPR*, 2009.
- [46] M. Roser and F. Moosmann. Classification of weather situations on single color images. In *IEEE Intelligent Vehicles Symposium*, 2008.
- [47] B. Russell, A. Torralba, K. Murphy, and W. T. Freeman. Labelme: a database and web-based tool for image annotation. *IJCV*, 2008.
- [48] B. C. Russell, W. T. Freeman, A. A. Efros, J. Sivic, and A. Zisserman. Using multiple segmentations to discover objects and their extent in image collections. In *CVPR*, 2006.
- [49] F. Sadeghi and M. F. Tappen. Latent pyramidal regions for recognizing scenes. In *ECCV*, 2012.
- [50] L. Shen and P. Tan. Photometric stereo and weather estimation using internet images. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009.
- [51] N. Shroff, P. Turaga, and R. Chellappa. Moving vistas: Exploiting motion for describing scenes. In *Computer Vision and Pattern Recognition (CVPR)*, 2010 *IEEE Conference on*, 2010.
- [52] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [53] S. Singh, A. Gupta, and A. A. Efros. Unsupervised discovery of mid-level discriminative patches. In *ECCV*, 2012.
- [54] L. Tao, L. Yuan, and J. Sun. Skyfinder: Attribute-based sky image search. In *ACM SIGGRAPH*, 2009.
- [55] S. Todorovic and N. Ahuja. Unsupervised category modeling, recognition, and segmentation in images. *PAMI*, 2008.
- [56] A. Volokitin, R. Timofte, L. Van Gool, and D. CVL. Deep features or not: Temperature and time prediction in outdoor scenes. 2016.
- [57] J. Wang, M. Korayem, and D. Crandall. Observing the natural world with flickr. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 452–459, 2013.
- [58] J. Wang, J. Yang, K. Yu, F. Lv, T. Huang, and Y. Gong. Locality-constrained linear coding for image classification. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2010.
- [59] S. Workman, R. Souvenir, and N. Jacobs. Scene shape estimation from multiple partly cloudy days. *Computer Vision and Image Understanding*, 134:116–129, 2015.
- [60] J. Xiao, J. Hays, K. A. Ehinger, A. Oliva, and A. Torralba. Sun database: Large-scale scene recognition from abbey to zoo. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2010.
- [61] X. Yan, Y. Luo, and X. Zheng. Weather recognition based on images captured by vision system in vehicle. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009.
- [62] J. Yang, K. Yu, Y. Gong, and T. Huang. Linear spatial pyramid matching using sparse coding for image classification. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009.
- [63] B. Yao, G. Bradschi, and L. Fei-Fei. A codebook-free and annotation-free approach for fine-grained image categorization. In *CVPR*, 2012.
- [64] N. Zhang, E. Shelhamer, Y. Gao, and T. Darrell. Fine-grained pose prediction, normalization, and recognition. *arXiv preprint arXiv:1511.07063*, 2015.
- [65] Z. Zhang and H. Ma. Multi-class weather classification on single images. In *IEEE International Conference on Image Processing (ICIP)*, pages 4396–4400, 2015.
- [66] Y. Zheng, Y.-G. Jiang, and X. Xue. Learning hybrid part filters for scene recognition. In *ECCV*, 2012.

Cewu Lu received the BS and MS degrees from Chongqing University of Posts and Telecommunications and Graduate University of Chinese Academy of Sciences in 2006 and 2009 respectively, and the PhD degree in 2013 in computer science and engineering from the Chinese University of Hong Kong. His team received the fourth place in ILSVRC 2014 among the 38 participating teams. He received the best paper award of NPAR 2012 and served as an associate editor for journal *gtCVPR* and reviewers for several major computer vision and graphics conferences and journals such as *TPAMI* and *TOG*. His research interests include activity recognition, object detection and image/video processing. He is currently a research fellow in the Department of Computer Science at Stanford University.

Di Lin received the bachelor's degree in software engineering from Sun Yat-sen University in 2012. He is currently pursuing the PhD degree in the Chinese University of Hong Kong. His research interests are computer vision and machine learning.

Jiaya Jia received the PhD degree in Computer Science from the Hong Kong University of Science and Technology in 2004 and is currently a full professor in Department of Computer Science and Engineering at the Chinese University of Hong Kong (CUHK). He was a visiting scholar at Microsoft Research Asia from March 2004 to August 2005 and conducted collaborative research at Adobe Systems in 2007. He heads the research group in CUHK, focusing specifically on computational photography, 3D reconstruction, practical optimization, and motion estimation. He currently serves as an associate editor for the *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)* and served as an area chair for *ICCV* 2011, *ICCV* 2013, and *CVPR* 2016.

Chi-Keung Tang received the MSc and PhD degrees in computer science from the University of Southern California (USC), Los Angeles, in 1999 and 2000, respectively. Since 2000, he has been with the Department of Computer Science at the Hong Kong University of Science and Technology (HKUST) where he is currently a full professor. He was an adjunct researcher at the Visual Computing Group of Microsoft Research Asia. His research areas are computer vision, computer graphics, and human-computer interaction. He was an associate editor of *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, and was on the editorial board of *International Journal of Computer Vision (IJCV)*. He served as an area chair *ICCV* 2007, *ICCV* 2009, *ICCV* 2011, *ICCV* 2015, and as a technical papers committee member for the inaugural *SIGGRAPH Asia* 2008, *SIGGRAPH* 2011, *SIGGRAPH Asia* 2011, *SIGGRAPH* 2012, *SIGGRAPH Asia* 2014, and *SIGGRAPH Asia* 2015.