

Incorporating multiple SVMs for automatic image annotation

Xiaojun Qi*, Yutao Han

Computer Science Department, Utah State University, Logan, UT 84322-4205, USA

Received 24 September 2005; received in revised form 11 March 2006; accepted 28 April 2006

Abstract

In this paper, a novel automatic image annotation system is proposed, which integrates two sets of support vector machines (SVMs), namely the multiple instance learning (MIL)-based and global-feature-based SVMs, for annotation. The MIL-based bag features are obtained by applying MIL on the image blocks, where the enhanced diversity density (DD) algorithm and a faster searching algorithm are applied to improve the efficiency and accuracy. They are further input to a set of SVMs for finding the optimum hyperplanes to annotate training images. Similarly, global color and texture features, including color histogram and modified edge histogram, are fed into another set of SVMs for categorizing training images. Consequently, two sets of image features are constructed for each test image and are, respectively, sent to the two sets of SVMs, whose outputs are incorporated by an automatic weight estimation method to obtain the final annotation results. Our proposed annotation approach demonstrates a promising performance for an image database of 12 000 general-purpose images from COREL, as compared with some current peer systems in the literature.

© 2006 Pattern Recognition Society. Published by Elsevier Ltd. All rights reserved.

Keywords: Image annotation; Image sub-blocking; Support vector machines; Multiple instance learning; Edge histogram descriptors; Color histogram

1. Introduction

Image annotation refers to the labeling of images with a set of predefined keywords. It is mainly used for visual information management and can be applied in a variety of domains such as entertainment, commerce, education, biomedicine, military, web image classification and search, etc. In particular, image annotation can aid in image retrieval since annotated keywords greatly narrow the semantic gap between low-level features and high-level semantics.

Automatic image annotation is a challenging task due to various imaging conditions, complex and hard-to-describe objects, a highly textured background, and occlusions. In general, most approaches use learning-based techniques to train manually categorized images and test the uncategorized images based on the training results. Because the selected training images are usually very limited and cannot represent all the aspects of real life, automatic annotation may not achieve high accuracy using the current computer vision and

image processing technologies [1]. The relevance feedback approach [2] is believed to be more effective by refining the initially assigned labels via several rounds of feedback provided by the user through an interactive interface. However, relevance-feedback-based annotation systems may be a burden to users, especially when more information is required than just Boolean feedback (relevant or not-relevant). In this paper, we exclusively focus on the techniques for automatic image annotation without using the relevance feedback to refine the labeled images.

Our proposed system utilizes an automatic weight estimation method to fuse the results from two complementary sets of support vector machines (SVMs) for automatic annotation. One set of SVMs is obtained by training the bag features of the image blocks using the multiple instance learning (MIL) technique. In MIL, each image is a bag and its segmented regions or sub-images are instances. As a result, the image features obtained from MIL are referred to as bag features. In our proposed system, we divide the images into non-overlapping blocks and apply the MIL technique upon the block-based color and texture features to extract the bag features. An enhanced diversity density (DD) method

* Corresponding author. Tel.: +1 435 7978155; fax: +1 435 7973265.
E-mail address: xqi@cc.usu.edu (X. Qi).

and a faster searching algorithm are further employed in the MIL technique to improve the efficiency and accuracy for such extraction. These extracted bag features are then fed into a set of SVMs for finding the optimum hyperplanes to annotate each training image. More details about MIL and DD are explained in Section 3.1. To address any inaccuracy issues related to the image sub-blocking, another set of SVMs is obtained by training the global color and texture features using the color and edge histogram-based technique. That is, we compute the MPEG-7 scalable color descriptor (SCD) and the modified MPEG-7 edge histogram descriptor (EHD) as the global features for the same training images and send these global features to another set of SVMs to find the optimum hyperplanes for annotating each training image. Once these two complementary sets of SVMs are configured, we can annotate each uncategorized test image by constructing both bag features of the image blocks and the global color and texture features and sending these two features to their respective sets of SVMs. The decision values yielded from these two sets of SVMs are then mapped to the likelihood values, which are further incorporated by a novel automatic weight estimation method to obtain the final annotation results.

The remainder of the paper is organized as follows. Section 2 briefly reviews the related work. Section 3 describes the general framework of our proposed annotation system. Section 4 illustrates the experimental results. Section 5 concludes with a brief discussion of our approach and some proposed directions for future work.

2. Related work

Many automatic image annotation systems have been developed since the early 1990s. In these systems, the images are represented by either global features, or block-based local and spatial properties, or region-based local features. Several relevant image annotation systems are briefly reviewed here.

2.1. Global-feature-based image annotation

Global image features have been widely used in image annotation. Huang et al. [3] categorize images by using a classification tree, which captures the spatial correlation of colors in an image. In [4], the k -nearest neighbor classifier on the color histogram is used to discriminate between indoor and outdoor images. Chapelle et al. [5] apply SVMs on the global $16 \times 16 \times 16$ -bin HSV color histograms to categorize images. Gdalyahu and Weinshall [6] use the local curve matching method for the shape silhouette classification, where objects in an image are represented by their contours. Vailaya et al. [7] use Bayesian classifiers on the color and edge direction histograms to, respectively, classify sunset/forest/mountain images and city/landscape images. Chang et al. [8] extract color features (i.e., dominant colors,

color histograms, color means, and color variances), shape features (i.e., elongation and spreadness), and wavelet-based texture features for each image and use SVMs and Bayes point machines (BPM) for automatic image annotation. These global features are easy to compute and effective for certain tasks. However, they have several major drawbacks as follows:

1. They lack information about the spatial feature distribution. Color correlogram [9] adds certain spatial information to address this issue. However, they are very limited in terms of the spatial layout of the objects.
2. They are sensitive to intensity variations and distortion.
3. They fail to narrow down the semantic gap, i.e., the difference between users' high-level semantic concepts and the low-level visual features, due to their limited descriptive power based on objects.

2.2. Block/region-based image annotation

To address these issues associated with the global features, a number of block-based and region-based methods have been exploited. In block-based systems, the image is divided into sub-blocks and features for each block are extracted. Gorkani and Picard [10] first divide the image into 16 non-overlapping equal-sized blocks and compute dominant orientations for each block. The image is then classified as city or suburb according to the majority orientations of the blocks. In [11], a one-dimensional hidden Markov model (HMM) trained on vector-quantized color histograms of image blocks is used for the indoor/outdoor classification. Li and Wang [12] propose an automatic linguistic indexing of pictures (ALIP) system which uses the two-dimensional multi-resolution HMM on features of image blocks for classification. Murphy et al. [13] build four graphical models to relate block-based features to objects and perform joint scene and object recognition. Cusano et al. [14] divide the image into a fixed number of partially overlapping subdivisions (tiles) and apply a multi-class SVM to classify the image into one of the seven predefined categories.

In region-based systems, an image is first segmented into homogeneous regions and features for each region are then extracted. Modestino and Zhang [15] use a Markov random field model to capture the spatial relationships between regions and apply a maximum a posteriori rule to interpret images. Minka and Picard [16] develop a system to generate several groups of possible regions based on different feature combinations. The best feature combinations for representing the sample semantic categories are discovered via the supervised learning of various parts of the images. Smith and Li [17] annotate images by applying a composite region template descriptor matrix on the spatial orderings of regions, whose attributes are represented by symbols in a finite pattern library. Barnard and Forsyth [1] apply a hierarchical statistic model on semantically meaningful regions to generate keywords for annotation. Brank [18] proposes

an annotation system by applying the SVM to the sparse vectors, whose elements indicate the possibilities of all the clustered regions in each training image. In addition, Brank presents another annotation system by applying the SVM to a generalized kernel converted from a similarity measure between the segmented images. Jeon et al. [19] use the cross media relevance model to predict the probability of generating a word given the regions in an image. Barnard et al. [20] apply the probabilistic modeling techniques on multi-modal data sets to automatically categorize images or image regions with relevant keywords by learning a joint distribution over each region of the training images and its corresponding text annotations. Their proposed modeling techniques combine a number of related statistical approaches with a variety of latent variable models. The image regions are segmented by the normalized cut method.

2.3. MIL-based image annotation

Recently, MIL has been applied for automatic image annotation. MIL was originally studied by Dietterich et al. [21] in the area of drug activity prediction and has recently received much attention in machine learning. In MIL, each image is a bag and its segmented regions or sub-images are instances. A bag is labeled positive if at least one instance is positive and is labeled negative if all instances are negative. The goal of the MIL is to generate a hypothesis to accurately predict the labels of unknown bags. Maron and Ratan [22] apply the MIL framework for the natural image classification. They use the DD algorithm to learn the user's concept based on a collection of positive and negative examples. The learning results are then used for image classification. Zhang and Goldman [23] develop the expectation maximization-diverse density (EM-DD) algorithm to improve the annotation speed and facilitate the scale-up to large data sets. However, the MIL is degraded to a "single instance learning" since EM only estimates one instance that is responsible for the label of the bag. Andrews et al. [24] propose an MI-SVM approach for annotation. The instance-based (i.e., region-based) image features are iteratively fed into SVMs until no updates for all the positive training images. The converged instance-based features are then used to annotate an image. In both DD and EM-DD algorithms, an optimum point in the feature space with a global maximum DD value is used to represent the object of interest. However, this optimum point may not correctly represent the object of interest due to the imperfect or incorrect image segmentation. To address this issue, Chen and Wang [25] propose the DD-SVM method for image annotation. This method constructs the SVMs by using bag-based image features as inputs. These bag-based features consist of multiple local maxima obtained by employing the EM-DD method on image regions. Experimental results [24,25] demonstrate that the DD-SVM method achieves the best annotation accuracy and the MI-SVM method achieves better annotation accuracy than both EM-DD and DD methods.

2.4. Summary

In spite of their successes, all these annotation systems have their shortcomings. The global-feature-based systems cannot precisely represent the semantics of an image. The region-based systems often break an object into several regions or put different objects into a single region due to inaccurate image segmentation. The block-based and MIL-based systems have the similar problems as the region-based systems. To our knowledge, the DD-SVM method [25] is the only one that makes an effort to partially address the inaccurate segmentation issues. However, this approach is computationally expensive in terms of both image segmentation and system training. In this paper, we propose an efficient and effective automatic annotation system which addresses both inaccurate segmentation and expensive training issues. Specifically, several complementary non-overlapping sub-block schemes are exploited to represent possible objects of interest. An improved MIL technique, which combines an enhanced DD method and a faster searching algorithm, is then applied to obtain the bag features for the SVMs to annotate an image. The global-feature-based SVMs are further incorporated to compensate any inaccuracy issues associated with the sub-blocking schemes and the MIL-based SVMs. As a result, the proposed annotation system provides more robustness against any issues associated with the shortcomings of the peer systems.

3. Proposed approach

3.1. MIL-based SVMs

3.1.1. Image sub-blocking and block feature extraction

It is well known that automatic image segmentation is an open problem in computer vision [26,27] and no system achieves perfect segmentation results. In addition, it is computationally expensive. Therefore, segmentation is not performed on our proposed system. Instead, we divide the image into blocks. Several block representation schemes were thoroughly studied and experiments have been performed on each scheme. Fig. 1 shows the final layout of the blocks chosen in our system, which has been proven to be efficient and effective. As shown in Fig. 1, each image is divided into five non-overlapping blocks. In general, this number of

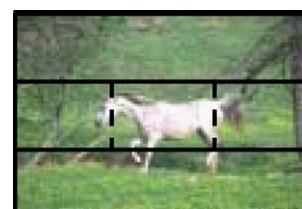


Fig. 1. Image sub-blocking.

blocks (i.e., 5) is close to the number of segmented regions for a majority of images. During image sub-blocking, each image is first horizontally divided into three blocks since the orientations of the objects in most images are horizontal. The middle block is further evenly divided into three sub-blocks to magnify the details of the main objects since they are usually located in the center. Even though the proposed sub-blocking system may divide an object into different blocks or put multiple objects into one block, which is likely for images whose main objects are not located in the center or whose orientation is vertical, our proposed MIL-based SVMs and global-feature-based SVMs will partially resolve this issue. That is, multiple local maxima are found by our enhanced MIL technique to represent the possible objects of interest. Moreover, global features are extracted independent of the sub-blocking scheme to represent the overall color and texture of the image. This integration makes the proposed system more robust against the inaccurate image sub-blocking issue. As a result, the annotation accuracy is not affected by the nature of the images in terms of the location of the main objects and the orientation of the images.

For each block, the mean and standard deviation in the LUV color space are computed as the color feature. The texture feature is calculated by the average energy in each high-frequency subband after applying two level wavelet decompositions on the luminance component of the image. Therefore, the length of the feature vector for each block is 12 (i.e., 6 color features and 6 texture features).

3.1.2. Multiple instance learning (MIL) and bag features

In MIL, the user labels the bag (i.e., image), which usually contains many instances (i.e., regions), as positive or negative. The goal of the MIL is to find what is common in all positive images, but not in any negative images. Maron and Ratan [22] develop the DD method to solve the MIL problem by converting this goal to a maximization problem. That is, with the assumption of n labeled bags and the hypothesis t , the DD value is calculated as

$$DD(t) = \prod_{i=1}^n \Pr(B_i, l_i | t) = \prod_{i=1}^n (1 - |l_i - \text{Label}(B_i | t)|),$$

$$\text{Label}(B_i | t) = \max_j \left\{ \exp \left[- \sum_{d=1}^m (s_d(B_{ijd} - t_d))^2 \right] \right\}, \quad (1)$$

where B_i denotes the i th bag, l_i denotes the actual label (0 or 1) of the i th bag, B_{ij} denotes the j th instance of bag i , B_{ijd} denotes the feature value of instance B_{ij} on dimension d , S_d denotes the value of feature weight vector S on dimension d , t_d denotes the value of t on dimension d , n denotes the number of instances, and m denotes the number of features. The maximization of Eq. (1) is to find the optimum t that leads to the maximum DD value for representing the user's interest in the feature space.

Two observations can be easily made by analyzing Eq. (1). That is

- If any instance in one negative bag B_i is close to t , $\text{Label}(B_i | t)$ will be close to 1 and therefore $\Pr(B_i, l_i | t)$ (i.e., the probability of the image being negative under the hypothesis t) will be close to 0.
- If all instances in one positive bag B_i are away from t , $\text{Label}(B_i | t)$ will be close to 0 and therefore $\Pr(B_i, l_i | t)$ (i.e., the probability of the image being positive under the hypothesis t) will be close to 0.

These two observations adversely drop the DD value close to 0 even though all the other instances in negative bags are far away from t and all the other positive bags fit t very well. As a result, we modify Eq. (1) so the DD value will not be drastically affected by several aforementioned abnormal instances. The enhanced definition of DD is

$$DD(t) = \sum_{i=1}^n \Pr(B_i, l_i | t). \quad (2)$$

By substituting multiplication in Eq. (1) with addition in Eq. (2), a robust DD method, which is more resistant to the presence of outliers, is obtained.

This enhanced DD method is further combined with the EM method [23] to speed up the searching process for finding the maximum DD value. That is, given the initial hypothesis t , the EM method selects one representative instance B_{ij}^* from each bag B_i by

$$B_{ij}^* = \arg \max_j \left(\exp \left[- \sum_{d=1}^m (s_d(B_{ijd} - t_d))^2 \right] \right). \quad (3)$$

This selected instance most likely corresponds to the annotation label of the image. Consequently, one chosen representative instance instead of all instances from each bag is used for calculating the DD value as follows:

$$DD(t) = \sum_{i=1}^n \Pr(B_i, l_i | t) = \sum_{i=1}^n (1 - |l_i - \text{Label}(B_i | t)|),$$

$$\text{Label}(B_i | t) = \exp \left[- \sum_{d=1}^m (s_d(B_{ijd}^* - t_d))^2 \right]. \quad (4)$$

In our proposed system, a simplex search method [28] is applied on Eq. (4) to locate the optimum point t which yields the local maximum DD value. This method is faster than the gradient-based method as used in other peer systems [22,23,25] since it is a direct search without using any numerical or analytical gradients. As a result, the time for searching the optimum point t is greatly shortened and the DD value is more resistant to the outliers, which ensure that our proposed enhanced EM-DD algorithm is more scalable and robust.

The bag feature of an image is further constructed by using the instance prototypes (IPs) which consist of all the appropriate local maxima and their weights. In order to find

each local maximum and its corresponding weight, we start the search from every instance of all positive bags with the same initial weights. The local maximum and its weight are automatically updated upon each search. Once all local maxima are found, the IPs are obtained by replacing clumped local maxima with their average and removing local maxima whose DD values are too small. These IPs $\{(x_k^*, w_k^*) : k = 1, \dots, n\}$ approximately represent all the possible objects of interest and are used to construct the bag feature $\phi(B_i)$ of each image $B_i = \{x_{ij} : j = 1, \dots, N_i\}$:

$$\phi(B_i) = \begin{bmatrix} \exp\left(\min_{j=1, \dots, N_i} \|x_{ij} - x_1^*\|_{w_1^*} / 30\right) \\ \exp\left(\min_{j=1, \dots, N_i} \|x_{ij} - x_2^*\|_{w_2^*} / 30\right) \\ \vdots \\ \exp\left(\min_{j=1, \dots, N_i} \|x_{ij} - x_n^*\|_{w_n^*} / 30\right) \end{bmatrix}, \quad (5)$$

where x_k^* 's are the feature values of the k th IP, w_k^* 's are the weights of the k th IP, n is the number of IPs, x_{ij} is the j th block features of image i , N_i is the block numbers (i.e., 5) in image i , and $\|\cdot\|_{w^*}$ is the weighted Euclidean distance. The exponential function is used here to properly scale the values of the bag features to the range between 0 and 1.

The detailed construction of the image bag feature is summarized in Fig. 2 in an algorithmic view. The input to the algorithm is a set of labeled training images whose block-based features are extracted according to Section 3.1.1. The output is a set of feature vectors, where each vector corresponds to the bag feature of a training image.

3.1.3. Support vector machines (SVMs)

Once the bag features for all training images are obtained using Eq. (5), they are further fed into SVMs, which will find a hyperplane that separates the training data by a maximal margin. That is, given m training data $\{x_i, y_i\}$'s, where $x_i \in R^n$ and $y_i \in \{-1, 1\}$, SVMs need to solve the following optimization problem:

$$\min_{\omega, b, \xi} \left(\frac{1}{2} \omega^T \omega + C \sum_{i=1}^l \xi_i \right) \\ \text{s.t. } y_i (\omega^T \phi(x_i) + b) > 1 - \xi_i, \quad \xi_i > 0, \quad (6)$$

where C is the penalty parameter of the error term, ω is the coefficient vector, b is a constant, ξ_i is a parameter for handling non-separable data, and $K(x_i, y_i) = \phi(x_i)^T \phi(x_j)$ is the kernel function. The non-linear SVMs with the Gaussian radial basis function (RBF) kernel are used in our system since they yield excellent results compared with linear and polynomial kernels [29]. This RBF kernel is defined as

$$K(x_i, y_i) = \exp(\gamma \|x_i - x_j\|^2), \quad \gamma > 0. \quad (7)$$

As a result, two SVMs related parameters C and γ need to be predetermined. We combine the 3-fold cross-validation and

1. Initialize IP to be an empty set
2. For each instance p in all positive training images
 - 2.1 Set the hypothesis vector x to be p
 - 2.2 Set its associated weight vector w to be all 1's
 - 2.3 Repeat
 - 2.3.1 Select the representative instance for each training image using Eq. (3)
 - 2.3.2 Find the optimal hypothesis and weight pair (x, w) that maximizes the DD value in Eq. (4) by using the simplex search method
 Until the hypotheses from the current and previous iterations converge or the maximum iteration is achieved
 - 2.4 Add (x, w) in to IP
 Endfor
3. For each (x, w) in IP
 - 3.1 If (its DD value < a predefined threshold)
 - 3.1.1 Remove it from IP
 - 3.2 Else
 - 3.2.1 Set $simIP$ to be an empty set
 - 3.2.2 Find all (x', w') 's in IP that are close to (x, w)
 - 3.2.3 Put (x, w) and all (x', w') 's into $simIP$ and remove all (x', w') 's from IP
 - 3.2.4 Replace (x, w) with the average of $simIP$
 Endif
 Endfor
4. For each training image
 - 4.1 Calculate its bag feature using Eq. (5)
 Endfor

Fig. 2. The algorithmic view of the steps to construct the image bag features.

grid-search algorithms [30] to find the best C and γ for the image annotation task by testing exponentially growing sequences of $C = 2^{-5}, 2^{-3}, \dots, 2^{15}$ and $\gamma = 2^{-15}, 2^{-13}, \dots, 2^3$ on several sets of pre-labeled training images. The pair that gives the minimum 3-fold cross-validation error is selected as the optimal parameters and is used in our proposed image annotation system.

Since the SVMs are designed for the binary classification, an appropriate multi-class method is needed to handle several classes as in image annotation. Two common approaches are "one against one" (i.e., apply pairwise comparisons between classes) and "one against the others" (i.e., compare a given class with all the others put together). We use "one against the others" in the proposed system as it achieves comparable performance with a faster speed than "one against one". That is, n SVMs will be generated for n categories to accomplish the annotation task.

In order to obtain the likelihood for an unknown image to be in each predefined category, the outputs (i.e., decision values) from each corresponding SVM, which are real values in the range of $[-\infty, +\infty]$, will be mapped into probabilities by training the parameters of a sigmoid function [31].

This sigmoid function is defined as

$$P(y = 1|f) = \frac{1}{1 + \exp(Af + B)}, \quad (8)$$

where P is the mapped probability, y is the actual categorical label, f is the output from SVMs, and A and B are two parameters estimated from a set of (f_i, y_i) s. For simplicity, these (f_i, y_i) s are generated as follows:

1. Generate 10 SVMs by applying the “one against the others” scheme on the global color and texture features of 1000 images from 10 predefined categories, where each category contains 100 images.
2. Feed the same 1000 training images as inputs to the 10 SVMs generated in step 1 to produce 10 sets of (f_i, y_i) s.

The maximum likelihood estimation method [31] is then applied on these 10 sets of (f_i, y_i) s to estimate 10 pairs of (A, B) s. Consequently, 10 sigmoid functions can be generated based on these 10 estimated pairs of (A, B) s. Fig. 3 illustrates the relationship between the decision values obtained from 10 SVMs and the mapping probability values computed by using 10 sigmoid functions along with their associated (A, B) s.

Since the 10 sigmoid functions are clumped together as shown in Fig. 3, we decided to use one representative sigmoid function to map the output from each SVM to the likelihood of the input image to be in each associated predefined category. To this end, we use the average of 10 probability values as the mapping probability value for any corresponding given decision value. A grid-searching algorithm [30] is then applied to find the best-fit (A, B) pair by testing on linearly growing sequences of $A = \{-4, -3.95, -3.9, \dots, -1.05, -1\}$ and $B = \{0.2, 0.22, \dots, 0.62, 0.64\}$ where the ranges of A and B are obtained from the 10 sigmoid functions. The pair that gives the minimum validation error is selected as the optimal parameters and is used consistently to map the output of each SVM to the probability. The final representative sigmoid function with the optimal parameters, i.e., $(A, B) = (-3, 0.28)$, is also shown in Fig. 3 as the thick black line.

3.2. Global-feature-based SVMs

Inaccurate image segmentation/blocking may make the IPs-based bag feature representation imprecise and therefore decrease the annotation accuracy of the MIL-based SVMs approach. Chen and Wang [25] address this issue by negating the labels of all bags and starting the search from every instance in all negative bags. This additional reverse procedure improves the annotation accuracy by around 2.2% for the 10-category database of 1000 images. However, it takes at least 9 times longer for the 10-category training compared with the scheme without negation. In our proposed system, we add global-feature-based SVMs, which require almost

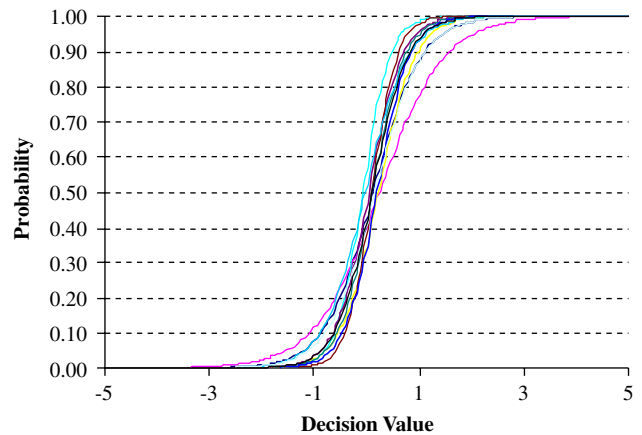


Fig. 3. Sigmoid-based probability functions by using different pairs of (A, B) s.

no additional computational cost, to address the inaccurate image sub-blocking issues.

In order to compensate the limitations associated with the specific color and texture representations, we construct the global features in a different manner as used in creating the regional features. To this end, two MPEG-7 descriptors are adopted in our system.

The SCD is one of the four MPEG-7 normative color descriptors [32]. It uses the HSV color histograms to represent an image since the HSV color space provides an intuitive representation of color and approximates human’s perception. We directly adopt the 128-bin SCD in our system.

The EHD is one of the three normative texture descriptors used in MPEG-7 [32]. It captures the spatial distribution of edges in an image. Five types of edges, namely, vertical, horizontal, 45° diagonal, 135° diagonal, and non-directional, have been used to represent the edge orientation in 16 non-overlapping sub-images. The normative EHD is therefore a total of 5×16 histogram bins, which represent the local edge distribution for each sub-image. Based on the EHD, we construct a gEHD (global EHD) to partially address the rotation, scaling, and translation related issues. This gEHD represents the edge distribution of the entire image and therefore has five bins. So the total length of our global features is 133, i.e., 128-bin SCD and 5-bin gEHD. This length is similar to the length of our MIL-based bag features.

After the global features of all the training images are obtained, they are fed into another set of SVMs to find optimum hyperplanes to distinguish one category from the others. This set of SVMs is designed by the same approaches used in the MIL-based SVMs.

3.3. Fusion approach

The fusion approach combines the outputs from the MIL-based and global-feature-based SVMs to obtain the final annotation results. For each test image, two sets of image

features are constructed and sent to the corresponding SVMs. The probability values computed from the MIL-based SVMs and the global-feature-based SVMs are further combined to obtain the final probability vector \vec{p} ,

$$\vec{p} = \vec{w} * \vec{p}_m + (1 - \vec{w}) * \vec{p}_g, \quad (9)$$

where \vec{p}_m is the probability vector obtained from the MIL-based SVMs, \vec{p}_g is the probability vector obtained from the global-feature-based SVMs, $*$ denotes the inner product operation, and \vec{w} determines the contribution from the MIL-based SVMs and is automatically estimated by applying the likelihood normalization method [33], which is adapted to fit our annotation system. The length of all these vectors equals the number of predefined categories. Specifically, each element in \vec{p}_m and \vec{p}_g indicates the probability of a test image to be classified as each corresponding category by using MIL-based and global-feature-based SVMs, respectively. Each element in \vec{w} indicates the category-based contribution from the MIL-based SVMs. As a result, each element in \vec{p} indicates the final probability of a test image to be classified as each corresponding category.

In our proposed system, we independently estimate the weights contributed from the MIL-based and global-feature-based SVMs. Two matrices, namely the MIL-based likelihood matrix L_m and the global-feature-based likelihood matrix L_g , are first created by testing the same training images on the corresponding sets of SVMs. The size of both matrices is $N \times K$, where N is the number of testing images and K is the number of predefined categories. Each value $L_m(n, c)$ indicates the probability of image n to be classified as category c by using the MIL-based SVMs. Similarly, each value $L_g(n, c)$ indicates the probability of image n to be classified as category c by using the global-feature-based SVMs. The MIL-based weight vector \vec{w}_m and the global-feature-based weight vector \vec{w}_g are then computed as

$$\begin{aligned} \vec{w}_m &= [w_{m,1}, w_{m,2}, w_{m,3}, \dots, w_{m,K-1}, w_{m,K}], \\ \vec{w}_g &= [w_{g,1}, w_{g,2}, w_{g,3}, \dots, w_{g,K-1}, w_{g,K}], \end{aligned} \quad (10)$$

where $w_{m,k}$ and $w_{g,k}$, $k = 1, 2, \dots, K$, are the normalized output likelihood for category k by using MIL-based and global-feature-based SVMs, respectively. They are computed as

$$\begin{aligned} w_{m,k} &= \frac{(1/NK) \sum_{n=1}^N \sum_{c=1}^K L_m(n, c)}{(1/N) \sum_{n=1}^N L_m(n, k)}, \\ w_{g,k} &= \frac{(1/NK) \sum_{n=1}^N \sum_{c=1}^K L_g(n, c)}{(1/N) \sum_{n=1}^N L_g(n, k)}, \end{aligned} \quad (11)$$

where each denominator is the average of likelihood obtained from the respective SVMs for category k , and each numerator is the average likelihood over all K categories for the respective SVMs. That is, the output likelihood for every category is normalized in accordance with the average of the probability values calculated over all testing images.

Once both \vec{w}_m and \vec{w}_g are estimated, the final weight vector \vec{w} in Eq. (9) can be obtained by

$$\vec{w} = \vec{w}_m / (\vec{w}_m + \vec{w}_g), \quad (12)$$

where $/$ denotes the element-wise division operation.

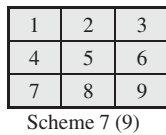
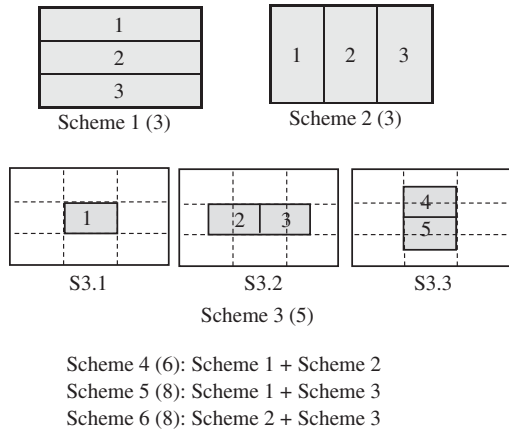
4. Experimental results

We have tested our annotation algorithm on 120 categories from the COREL database. Each category has 100 images. The 60 categories, which have distinct semantic topics such as Horse, Beach, and the like, are further selected. The images from the remaining categories form a more difficult set as the content of each category is not distinct, such as England, Japan, and so on. To correctly evaluate the annotation performance on these difficult categories, manual adjustment is necessary. Therefore, we only report the performance on the 60 categories that have distinct semantic topics.

4.1. Evaluation of image sub-blocking scheme

To evaluate the effect of image sub-blocking, a total of 16 alternative sub-blocking schemes are studied. The block representations of these schemes are shown in Fig. 4. Basically, we evenly divide an image into nine non-overlapping blocks as shown in Scheme 7 of Fig. 4. As the main objects are usually located around the middle of the image, the center block 5 is shifted to the left, right, up, or down by a half block to accommodate the possible shifts of the main objects (refer to S3.2 and S3.3 of Fig. 4). In addition, three horizontal groupings ($\{1, 2, 3\}$, $\{4, 5, 6\}$, $\{7, 8, 9\}$) and three vertical groupings ($\{1, 4, 7\}$, $\{2, 5, 8\}$, $\{3, 6, 9\}$) are considered for the possible wide or tall objects and their shifts (refer to scheme 1 and scheme 2 of Fig. 4). Finally, 16 different sub-blocking schemes are obtained by the different combinations of the above blocks.

These 16 alternative sub-blocking schemes together with our proposed sub-blocking scheme (shown in Fig. 1) are tested on the first 10 categories including African People and Villages, Beach, Historical Buildings, Buses, Dinosaurs, Elephants, Flowers, Horses, Mountains, and Food. For each category, 50 images are randomly selected as training images for the MIL-based SVMs and the remaining 50 images are used as test images. We repeat the above procedure 5 times and calculate the average annotation accuracy, which is shown in Fig. 5 by the diamonds connected via a blue line. It is clearly observed that the annotation accuracy increases in general as the number of blocks increases. The highest annotation accuracy is yielded when the number of blocks is the maximum (i.e., scheme 16 with 19 blocks). Specifically, the annotation accuracy obtained from the MIL-based SVMs increases from 75.1 to 84.3% when the number of blocks increases from 3 to 19. Our proposed sub-blocking scheme achieves an annotation accuracy of 82.1%, which differs from the highest accuracy by 2.2%.



- Scheme 8 (11): Scheme 3 + Scheme 4
- Scheme 9 (11): Scheme 7 + S3.2
- Scheme 10 (11): Scheme 7 + S3.3
- Scheme 11 (12): Scheme 7 + Scheme 1
- Scheme 12 (12): Scheme 7 + Scheme 2
- Scheme 13 (13): Scheme 7 + S3.2 + S3.3
- Scheme 14 (16): Scheme 13 + Scheme 1
- Scheme 15 (16): Scheme 13 + Scheme 2
- Scheme 16 (19): Scheme 13 + Scheme 4

Fig. 4. A total of 16 alternative sub-blocking schemes.

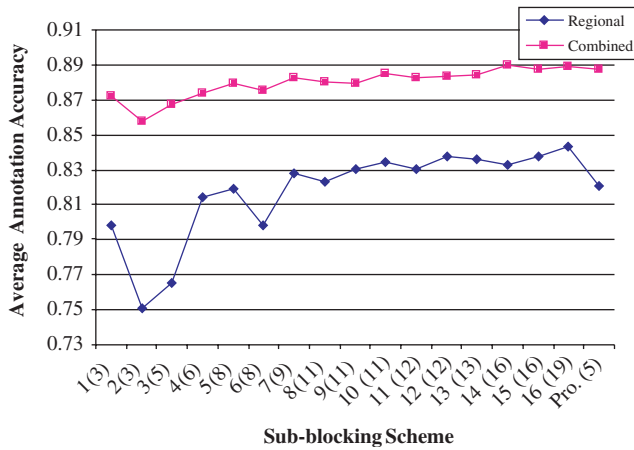


Fig. 5. Average annotation accuracy for different sub-blocking schemes.

The average annotation accuracy of our proposed system, which combines both MIL-based and global-feature-based SVMs, is shown in Fig. 5 as the squares connected by a purple line. Two observations can be easily made:

- The addition of the global-feature-based SVMs consistently improves the performance of the MIL-based SVMs. In specific, it, respectively, improves the

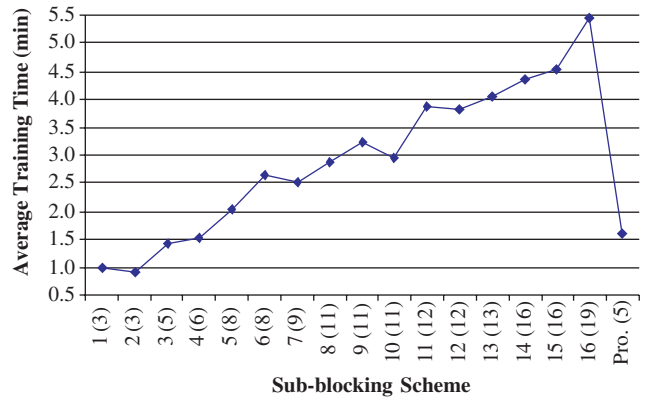


Fig. 6. Average training time for different sub-blocking schemes.

annotation accuracy from 82.1 to 88.8% and from 84.3 to 88.9% for our proposed sub-blocking scheme and scheme 16.

- The addition of the global-feature-based SVMs stabilizes the annotation fluctuations associated with MIL-based SVMs, which are mainly caused by the inaccurate image sub-blocking. That is, the difference between the highest and lowest annotation accuracy is 4.1% and 9.2% for our proposed fusion approach and the MIL-based approach, respectively.

Consequently, we experimentally prove that our fusion approach not only improves the overall annotation accuracy but also partially addresses the inaccurate sub-blocking issues.

Fig. 6 shows the average time for training one binary MIL-based SVM using different sub-blocking schemes, where the time is measured in minutes on a Pentium IV 3.06 GHz PC running the Windows XP professional edition and Matlab 7.0. It clearly shows that the training time increases in general as the number of blocks increases. In specific, the longest versus the shortest training time is 5.44 against 0.92. However, the corresponding annotation accuracy is 88.9% and 85.8%. Compared with these two extreme sub-blocking schemes, our proposed scheme takes about 1.62 min in training and yields 88.8% accuracy in annotation. The relatively short training time and comparable accuracy validates our chosen scheme as the best compromise in terms of efficiency and accuracy for a large-size image database.

4.2. Annotation results

To measure the effectiveness of our annotation system, the training procedure mentioned in Section 4.1 is repeated on 6000 images from the COREL database. Specifically, Table 1 shows the average annotation results of images from the first 10 categories, which have distinct semantics and have been widely used in the peer retrieval or annotation systems. This experiment demonstrates the following: (1) The overall annotation accuracy is 88.8%. (2) The accuracy of

Table 1

The confusion matrix of the proposed annotation system, where each row lists the average percentage of the images in one category classified into each of the 10 categories

	Africa	Beach	Building	Buses	Dinosaur	Elephant	Flower	Horse	Mountain	Food
Africa	0.808	0.012	0.04	0.012	0.012	0.068	0	0.008	0.012	0.028
Beach	0.024	0.772	0.028	0.024	0.004	0.012	0.004	0.008	0.112	0.012
Building	0.048	0.048	0.804	0.012	0.004	0.024	0.012	0	0.02	0.028
Buses	0	0	0	0.996	0	0	0	0	0	0.004
Dinosaur	0	0	0	0	1	0	0	0	0	0
Elephant	0.012	0	0.008	0	0	0.916	0	0.024	0.02	0.02
Flower	0.004	0	0	0	0	0	0.98	0.012	0	0.004
Horse	0	0.016	0	0	0	0	0	0.984	0	0
Mountain	0.004	0.156	0.044	0.012	0	0.036	0.008	0.004	0.732	0.004
Food	0.028	0.008	0.004	0.004	0.04	0	0.016	0.004	0.012	0.884

Numbers on the diagonal show the categorization accuracy for each category. The off-diagonal numbers indicate the classification errors.

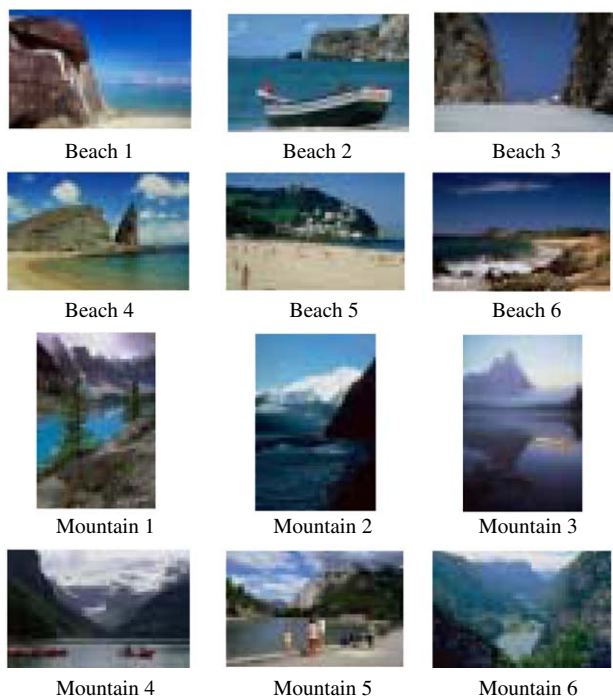


Fig. 7. Sample images from Beach and Mountain categories, where all listed beach images are misclassified as mountain and vice versa.

4 categories, namely Buses, Dinosaur, Flower, and Horse, is equal or close to 100%. (3) Eight categories achieve an average annotation accuracy of above 80.4% with the exception of Beach and Mountain categories, whose average annotation accuracy is about 75%. In these two categories, 11.2% of Beach is misclassified as Mountain and 15.6% of Mountain is misclassified as Beach. This misclassification is mainly caused by the fact that Beach and Mountain categories are semantically similar, i.e., most mountain and beach images have blue sky and similar textures. Moreover, some beach images may contain mountain and some mountain images may contain beach as well. Fig. 7 further illustrates the above observations using the misclassified beach and mountain images.

Table 2

Comparisons of three systems

	Accuracy (%)	Feature length	Training time
HistSVM	79.8	4096	~1
DD-SVM	81.5	~150	~15
Proposed	88.8	~200	~1.5

The proposed system is also compared with the DD-SVM system [25] and our implemented HistSVM system [5] using images from the first 10 categories. Table 2 summarizes the performance of these three systems in terms of the overall average annotation accuracy, the estimated feature length, and the approximate average training time in minutes for one binary SVM. It clearly shows that our proposed system performs the best. It outperforms the HistSVM system by 11.28% in the overall accuracy. In addition, its feature length is about 20 times shorter than that of HistSVM. Our system also improves the accuracy by 8.96% over the DD-SVM system, which is about 10 times slower than our system for the case of 10 categories as discussed in Section 3.2.

The average annotation accuracy for each of the first 10 predefined categories by using the above three systems is further plotted in Fig. 8. It clearly illustrates that the proposed system achieves the best average accuracy in all categories except category 1, where HistSVM performs the best. This is mainly due to the distinct colors of the most images (e.g., African People and Villages) in category 1, which make the color-based systems, such as the HistSVM system, more effective.

To validate the proposed automatic weight estimation method and the proposed fusion approach, the overall average annotation accuracy for the first 10 categories obtained by assigning different weights to the global-feature-based and MIL-based SVMs is plotted in Fig. 9, where G and M , respectively, represent global and MIL weights. For simplicity, it is assumed that each category contributes the same to each element in any given global-feature-based weight vector \vec{w}_g and MIL-based weight vector \vec{w}_m .

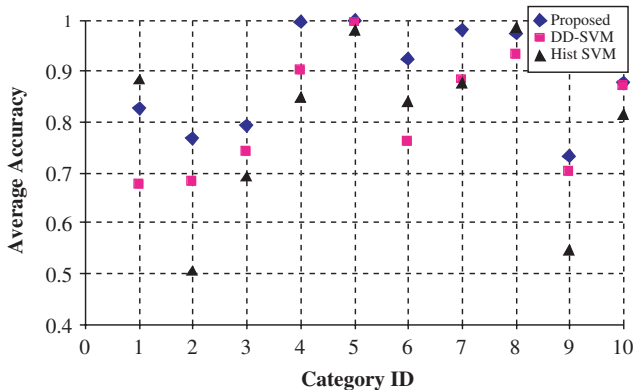


Fig. 8. Comparisons of the average annotation accuracy for each category by using three systems.

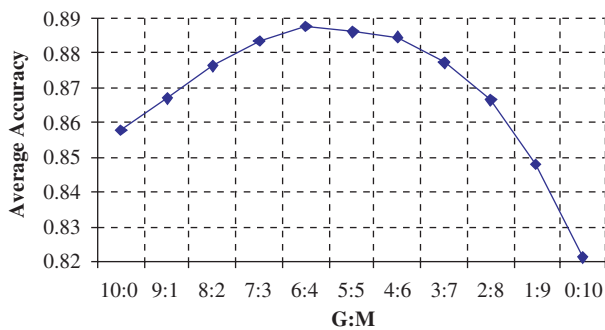


Fig. 9. Average annotation accuracy by assigning different weights to global and regional features.

That is, the two weight vectors will be $\vec{w}_g = [0.2, 0.2, \dots, 0.2]$ and $\vec{w}_m = [0.8, 0.8, \dots, 0.8]$ if $G : M = 2 : 8$. Under this simplification, the system achieves the best annotation accuracy of 88.7% when $G : M = 6 : 4$ as shown in Fig. 9. This accuracy is 0.1% less than the accuracy yielded from the proposed automatic weight estimation method, which approximates the actual contribution from each category based on the normalized likelihood. As a result, the effectiveness of the proposed automatic weight estimation method is experimentally confirmed. It is also observed from Fig. 9 that either our global-feature-based SVMs system (i.e., $G : M = 10 : 0$) or our MIL-based SVMs system (i.e., $G : M = 0 : 10$) alone achieves the respective average accuracy of 85.8% or 82.1%, which is better than both DD-SVM and HistSVM systems. Moreover, it clearly shows the efficacy of the fusion approach as it improves the global-feature-based and MIL-based SVMs systems by 3.5% and 8.2%, respectively.

4.3. Sensitivity to the number of categories

The scalability of the method is tested by performing image annotation experiments over data sets with different number of categories. Two experiments are performed.

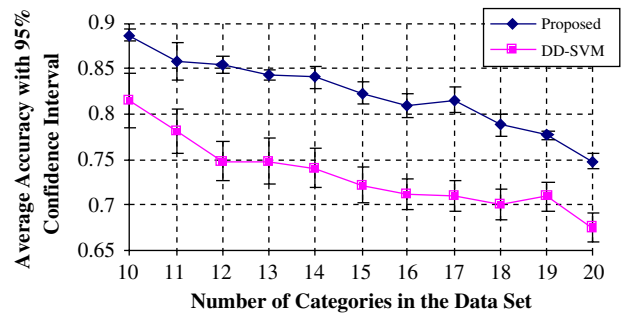


Fig. 10. Comparison of the two methods on the robustness to the number of categories.

In the first experiment, a total of 11 data sets are created. The number of categories in a data set varies from 10 to 20. These data sets are arranged in the same manner as in Ref. [25] for fair comparisons. That is, the first 10 categories used in our previous experiments form the first data set; the first 11 categories, including the first 10 categories plus the Dog category, form the second data set; etc. The overall average annotation accuracy and the 95% confidence intervals of our system and the DD-SVM system are shown in Fig. 10 by running both systems 5 times on each data set. We observe a decrease in average annotation accuracy in both systems as the number of categories increases. In specific, when the number of categories doubles (i.e., increases from 10 to 20 categories), the decrease in average annotation accuracy of our proposed system versus the DD-SVM system is from 88.8 to 74.8% against from 81.5 to 67.5%, respectively. However, our method consistently outperforms the DD-SVM method in all the 11 data sets by 8.9%, 10.0%, 14.1%, 12.6%, 13.4%, 14.1%, 13.7%, 15.1%, 12.5%, 9.4%, and 10.8%, respectively. Since the accuracy improvement from the second to the eleventh data set is always higher than the improvement from the first data set containing the least number of categories, the proposed system is certainly less sensitive to the increasing category numbers.

In the second experiment, the proposed block-based fusion approach is compared with the segmentation-based fusion approach. For the latter one, an image is first segmented into homogeneous regions by applying the unsupervised k -means algorithm on colors. The image bag features are extracted in the same manner as our block-based MIL method. The same fusion approach is also applied to ensure fair comparisons. Fig. 11 shows the overall average annotation accuracy of both block-based and segmentation-based fusion approaches upon the 11 data sets used in the first experiment. It clearly shows that our block-based approach achieves a little bit better accuracy than the segmentation-based approach in all 11 data sets. Furthermore, it takes almost no time for our block-based approach to divide an image into 5 blocks when comparing with the time spent in segmentation. The training time for both block-based and segmentation-based systems is comparable due to the fact

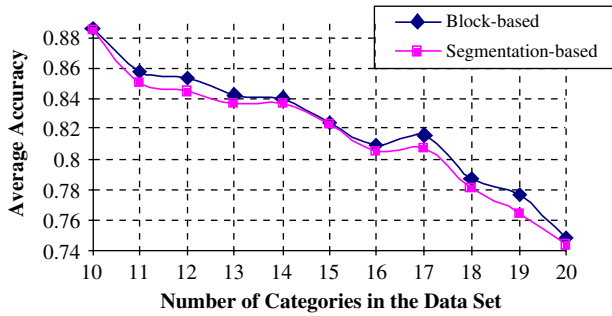


Fig. 11. Comparison of block-based and segmentation-based fusion approaches.

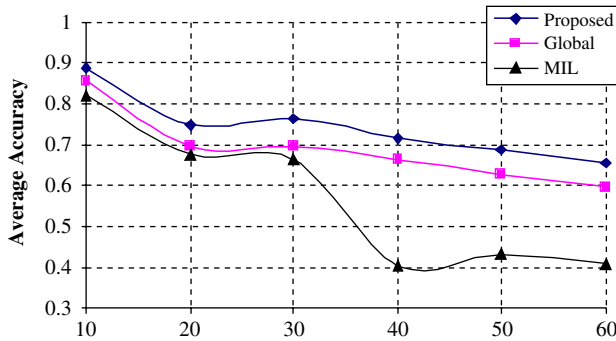


Fig. 12. Comparisons of the average annotation accuracy of three methods on different number of categories.

that the average number of regions obtained from image segmentation (i.e., 4.4) is similar to the number of blocks obtained from the proposed sub-blocking scheme (i.e., 5). Our proposed block-based fusion approach is therefore empirically validated in terms of both accuracy and efficiency.

To further verify the robustness of the proposed system, six additional data sets are created where the number of categories in a data set varies from 10 to 60 with a step size of 10 (i.e., the number of categories equals to 10, 20, 30, 40, 50, 60). Fig. 12 plots the overall average annotation accuracy of the proposed system, global-feature-based SVMs system (Global), and MIL-based SVMs system (MIL) for each of the six data sets. The expected decrease in average accuracy is clearly observed. In particular, average accuracy drops from 88.8 to 65.7% when the number of categories increases from 10 to 60. This reasonable accuracy decrease indicates the scalability and robustness of the proposed system. Fig. 12 also demonstrates the validity of the fusion approach as it consistently improves the MIL-based and global-based systems. For instance, the fusion approach, respectively, improves the global-based and MIL-based systems by 10.3% and 61.4% for 60 categories. Three interesting observations are: (1) The average accuracy of both fusion and MIL-based approaches increases a bit when the number of categories increases from 20 to 30. (2) The average accuracy of the MIL-based approach drastically decreases when the number of categories increases from 30 to

40. (3) The average accuracy of the global-based approach stays very stable as the number of categories increases. The first two “abnormal” observations are mainly caused by the fact that images in categories 21–30 are more distinct and the possibility for image blocks from two categories to be similar is greatly increased for a large-size image database. However, the last observation demonstrates the effectiveness of the extracted global features.

4.4. Sensitivity to the diversity of training images

To test the performance of the proposed system as the diversity of training images varies, we created four data sets, each containing 1000 images from 10 categories. Among these four data sets, one consists of distinct categories, including Antique, Bus, Dinosaur, Fashion, Flower, Food, Horse, Owls, Snow-mountain, and Sunset; one is composed of similar semantics, including Beach, Building, Snow-mountain, Ski, Waterfall, Subsea, Space, Minerals, Firework, and Sunset; the other two data sets are the combinations of the first two data sets, i.e., 5 out of 10 categories have distinct semantics and the remaining 5 categories have similar semantics. Specifically, two groups of similar categories are separately generated for the third and fourth data set. The first group consists of images from Beach, Building, Snow-mountain, Ski, and Waterfall, where both beach and building images contain blue sky; both mountain and ski images contain snow; all mountain and ski images and some beach and waterfall images contain mountain; and both beach and waterfall images contain water. The second group consists of images from Subsea, Space, Minerals, Firework, and Sunset, where subsea, space, and minerals images have very dark background with some objects of different shape or color in the center; firework and sunset images have some bright color scattered in the dark background. Some sample images from these 18 categories are shown in Fig. 13.

The overall annotation accuracy by applying the proposed system on the above four data sets is summarized in Table 3. As expected, the proposed system achieves a high annotation accuracy of 93.5% for the first data set with distinct categories. A relatively low annotation accuracy (i.e., 74.4%) is yielded for the second data set with similar categories. Comparable annotation accuracy (i.e., 81.6% vs. 85.8%) is obtained for the two data sets with the combination of similar and dissimilar categories. This experimental result clearly demonstrates the robustness of the proposed system to the diversity of training images.

4.5. Speed

The proposed annotation system has been implemented using Matlab 7.0 on a Pentium IV 3.06GHz PC running Windows XP operating system. Training one binary MIL-based SVM using 500 images takes about 1.5 min as shown

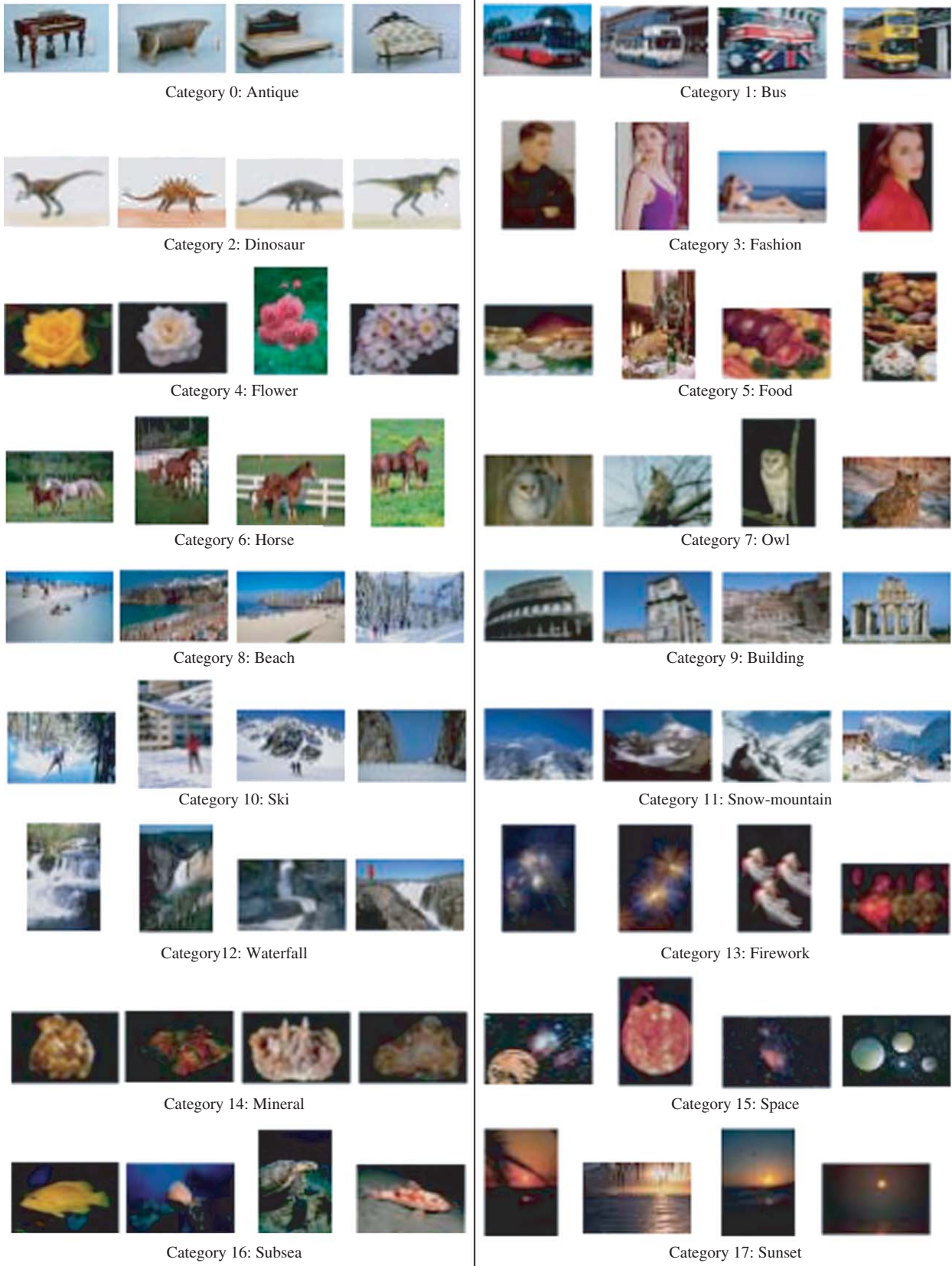


Fig. 13. Sample images from 18 categories with similar and dissimilar semantics.

Table 3
Average accuracy of the proposed system on diverse data sets

Data set	Accuracy (%)
Distinct categories	93.5
Similar categories	74.7
Combination1	81.6
Combination2	85.8

in Table 2. However, most time is spent on searching the hypotheses with the maximum DD values, which can be easily reduced to at least one-twentieth if the searching algorithm is implemented by C language. In general, the training time linearly increases with the number of training images since each image contains 5 blocks. Training the global-feature-based SVMs takes almost no extra time when comparing with the MIL-based SVM training.

The time complexities of other non-training related operations are summarized below:

- Sub-blocking images takes $O(C * N)$ with N being the number of images in the database and C being the number of blocks for each image.
- Calculating the block-based color features takes $O(C * N * d1)$ with $d1$ being the dimensionality of the block-based color feature.
- Calculating the block-based texture features takes $O(N * Row * Col)$ with Row and Col being the dimensionality of the image itself.
- Calculating the global color features takes $O(N * d2)$ with $d2$ being the dimensionality of the global color feature.
- Calculating the global edge features takes $O(N * e)$ with e being the total number of edge types.
- Sorting the annotation results for one testing image takes $O(Cat * \log Cat)$ with Cat being the number of categories to be annotated.

For the COREL database, we have $N = 6000$, $C = 5$, $d1 = 6$, $d2 = 128$, $e = 5$, and $Cat = 60$.

5. Conclusions

In this paper, we present an efficient and effective automatic image annotation system, which combines MIL-based SVMs and global-feature-based SVMs. The main contributions are:

- Using novel block-based features, instead of the expensive segmentation-based features, for MIL.
- Employing more robust DD definition in MIL.
- Applying a faster search algorithm (i.e., a simplex search method) to speed up the process of finding the maximum DD values.
- Combining the IPs-based bag features with SVMs to approximately represent all possible objects of interest.

- Integrating the global-feature-based SVMs with the MIL-based SVMs to address the inaccurate image sub-blocking related issues.
- Proposing a novel and fast automatic weight estimation method, which applies the likelihood normalization for weight optimization.
- Constructing the global and block features in a different manner to compensate the limitations associated with the specific color and texture representations.
- Using multi-category SVMs to classify images by a set of probability values for each category.

The proposed system has been validated by testing 6000 general-purpose images with 60 semantically distinct categories. The experimental results indicate that our system outperforms peer systems (e.g., the DD-SVM system and the HistSVM system) in the literature in terms of accuracy, efficiency, robustness, and scalability.

The proposed system can be easily integrated into the image retrieval system, where both annotated keywords and the query image(s) can be combined as the query.

References

- [1] K. Barnard, D. Forsyth, Learning the semantics of words and pictures, Proceedings of the International Conference on Computer Vision, 2001, pp. 408–415.
- [2] G. Sychay, E. Chang, K. Goh, Effective image annotation via active learning, Proceedings of the IEEE International Conference on Multimedia and Expo, Lausanne, Switzerland, 2002, pp. 209–212.
- [3] J. Huang, S.R. Kumar, R. Zabih, An automatic hierarchical image classification scheme, Proceedings of the Sixth ACM International Conference on Multimedia, Bristol, England, 1998, pp. 219–228.
- [4] M. Szummer, R.W. Picard, Indoor–outdoor image classification, Proceedings of the IEEE International Workshop on Content-based Access of Image and Video Database, 1998, pp. 42–51.
- [5] O. Chapelle, P. Haffner, V.N. Vapnik, Support vector machines for histogram-based image classification, IEEE Trans. Neural Networks 10 (1999) 1055–1064.
- [6] Y. Gdalyahu, D. Weinshall, Flexible syntactic matching of curves and its application to automatic hierarchical classification of silhouette, IEEE Trans. PAMI 21 (12) (1999) 1312–1328.
- [7] A. Vailaya, M.A.T. Figueiredo, A.K. Jain, H.J. Zhang, Image classification for content-based indexing, IEEE Trans. Image Process. 10 (1) (2001) 117–130.
- [8] E. Chang, K. Goh, G. Sychay, G. Wu, CBSA: content-based soft annotation for multimodal image retrieval using bayes point machines, IEEE Trans. Circuits Syst. Video Technol. 13 (1) (2003) 26–38.
- [9] J. Huang, S. Ravi Kumar, M. Mitra, W. Zhu, R. Zabih, Image indexing using color correlogram, IEEE International Conference on Computer Vision and Pattern Recognition, 1997, pp. 762–768.
- [10] M.M. Gorkani, R.W. Picard, Texture orientation for sorting photos ‘at a glance’, Proceedings of the 12th International Conference on Pattern Recognition, 1994, pp. 459–464.
- [11] H. Yu, W. Wolf, Scenic classification methods for image and video database, Proceedings of the SPIE International Conference on Digital Image Storage and Archiving Systems, 1995, pp. 363–371.
- [12] J. Li, J.Z. Wang, Automatic linguistic indexing of pictures by a statistical modeling approach, IEEE Trans. PAMI 25 (2003) 1075–1088.

- [13] K. Murphy, A. Torralba, W. Freeman, Using the forest to see the trees: a graphical model relating features, objects, and scenes, *Advances in Neural Information Processing Systems*, vol. 16, Cambridge, MA, MIT Press, Vancouver, BC, 2003.
- [14] C. Cusano, G. Ciocca, R. Schettini, Image annotation using SVM, *Proceedings of Internet Imaging IV, SPIE 5304*, 2004.
- [15] J.W. Modestino, J. Zhang, A Markov random field model-based approach to image interpretation, *IEEE Trans. PAMI 14 (6) (1992)* 606–615.
- [16] T.P. Minka, R.W. Picard, Interactive learning using a society of models, *Pattern Recognition 30 (3) (1991)* 565–581.
- [17] J.R. Smith, C.S. Li, Image classification and querying using composite region templates, *Int. J. Comput. Vision Image Understanding 75 (1–2) (1999)* 165–174.
- [18] J. Brank, Image categorization based on segmentation and region clustering, *Proceedings of the First Starting AI Researchers Symposium (STAIRS)*, Lyon, France, 2002, pp. 145–154.
- [19] J. Jeon, V. Lavrenko, R. Manmatha, Automatic image annotation and retrieval using cross-media relevance models, *Proceedings of the 26th International ACM SIGIR Conference*, Toronto, Canada, 2003, pp. 119–126.
- [20] K. Barnard, P. Duygulu, D. Forsyth, N. Freitas, D.M. Blei, M.I. Jordan, Matching words and pictures, *J. Mach. Learn. Res. 3 (2003)* 1107–1135.
- [21] T.G. Dietterich, R.H. Lathrop, T. Lozano-Perez, Solving the multiple instance problem with axis-parallel rectangles, *Artif. Intell. 89 (1–2) (1997)* 31–71.
- [22] O. Maron, A.L. Ratan, Multiple-instance learning for natural scene classification, *Proceedings of the 15th International Conference on Machine Learning*, San Francisco, CA, 1998, pp. 341–249.
- [23] Q. Zhang, S.A. Goldman, W. Yu, J. Fritts, Content-based image retrieval using multiple instance learning, *Proceedings of the 19th International Conference on Machine Learning*, Sydney, Australia, 2002, pp. 682–689.
- [24] S. Andrews, I. Tsochantaris, T. Hofmann, Support vector machines for multiple-instance learning, in: *Advances in Neural Information Processing Systems*, MIT Press, Cambridge, MA, 2002, pp. 561–568.
- [25] Y. Chen, J.Z. Wang, Image categorization by learning and reasoning with regions, *J. Mach. Learn. Res. 5 (2004)* 913–939.
- [26] J. Shi, J. Malik, Normalized cut and image segmentation, *IEEE Trans. PAMI 22 (8) (2000)* 888–905.
- [27] S. Zhu, A.L. Yuille, Region competition: unifying snakes, region growing, and Bayes/MDL for multiband image segmentation, *IEEE Trans. PAMI 18 (9) (1996)* 884–900.
- [28] J.C. Lagarias, J.A. Reeds, M.H. Wright, P.E. Wright, Convergence properties of the Nelder–Mead simplex method in low dimensions, *SIAM J. Optim. 9 (1998)* 112–147.
- [29] B. Scholkopf, K. Sung, C. Burges, F. Girosi, P. Niyogi, T. Poggio, V. Vapnik, Comparing Support Vector Machines with Gaussian Kernels to Radial Basis Function Classifiers, *A.I. Memo 1599*, MIT, Cambridge, MA, 1996.
- [30] C. Hsu, C. Chang, C. Lin, A practical guide to support vector classification, Technical Report, Department of Computer Science and Information Engineering, National Taiwan University, Taiwan, 2003.
- [31] J.C. Platt, Probabilistic output for support vector machines and comparisons to regularized likelihood methods, in: A. Bartlett, P. Schölkopf, B.E. Schuurmans (Eds.), *Advances in Large Margin Classifiers*, MIT Press, Cambridge, MA, 2000, pp. 61–74.
- [32] B.S. Manjunath, P. Salembier, T. Sikora, *Introduction to MPEG-7 Multimedia Content Description Interface*, Wiley, New York, 2002.
- [33] S. Tamura, K. Iwano, S. Furui, A stream-weight optimization method for multi-stream HMMs based on likelihood value normalization, *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing*, 2005, pp. 469–472.

About the Author—XIAOJUN QI received her Ph.D. degree in Computer Science in 2001 from Louisiana State University. She is an assistant professor in Computer Science Department at Utah State University since 2002. Her research interests include image processing, pattern recognition, computer vision, and machine learning.

About the Author—YUTAO HAN received his M.S. degree in Electrical Engineering from Utah State University in 2003. He currently works on his M.S. degree in Computer Science at Utah State University. His research interests include signal and image processing, pattern recognition, machine learning, and networking.