



the database is almost fully annotated. Note that the classifier has to be incremental and should operate in real time basis. As a classifier, we will use one of the most successful classifiers, Support vector machine (SVM) along with a K nearest neighbor (KNN).

The paper is organized as follows: Section 2 presents a detailed description of our proposed architecture for annotation. Section 3 presents discussions and comments on future work.

## 2. Proposed Architecture for Annotation

Our system is composed of five components (See Figure 1). Each component performs a well defined task. Grid segmentation segments the images into n number of grids. In this system we use 16 grids. Next we apply feature extraction module. Then we provide ground truth for the classifier. Finally, the system will support keyword based query or query by example. For our experiments, we use Corel dataset [1].

### 2.1. Feature Extraction

First, feature extraction is employed for the entire image without doing any segmentation/partition. Next, we extract features for each grid. For grid case, we focus on color and texture features. We ignore shape because our focus is not on object recognition/boundary detection. We extract texture information of each grid using three representative energies in high frequency bands of the wavelet transforms. To obtain these moments, a Daubechies 4-wavelet transform is applied to the L component of the image. After one-level transformation, a grid is decomposed into four frequency bands: the LL (low low), LH (low high), HL, and HH bands. Here, we resample all images to the same resolution before computing wavelet coefficients.

Here, we exploit three sets of features. Feature set 1 is used for the whole image. We use average RGB, LUV color features (dimension=6), texture features (=3), and 10 moments including spatial, central, moment Hu invariant. As a whole, we have a total of 19 features. Feature set 2 and 3 will be used for each grid. For feature set 2, we convert the basic color components (R, G, B) to HSV. With HSV, we calculate the color-histogram bins (7\*7) of H and S, respectively. We construct 7X7 bins (=49). Along with 6 color features (RGB, LUV) and texture features (=3), we will end with 58 features/grid for feature set 2. Finally, for feature set 3, we increase histogram bin sizes. Here, we use 10X10 bins. Hence, the total number of features

will be 109. In the feature extraction phase, we extract the features for each grid using the INTEL OPENCV software.

### 2.2. Ground Truth Annotation

In ground truth phase, we try to choose quality images that are related to a specific concept and annotate them manually. The goal here is to choose the grids that are mostly related to specific concept. For example, in Figure 2, we choose ground truth for the concept "Bear". As we can see, not all the 16 grids are related to "Bear"; hence, we only choose grid 1 and 5 for the image in part A, and grids 5, 6, and 10 for image in part B. On the other hand, for the whole image case, we choose the whole image as a ground truth.

### 2.3. SVM Training

After collecting reasonable number of ground truth for each concept, we start the process of training. We use Support Vector Machines (SVM) for training. The training process is very fast because only the basic ground truth grids are involved here. For the purpose of feature selection and testing, our system provides several options. For example, the user has three different feature sets (see Section 2.1: Feature Extraction). Also the user has the choice to choose between SVM and K nearest neighbor (KNN) in training and prediction. In addition, the user can choose between the whole and 16 grid schemes.

### 2.4. K Nearest Neighbor (KNN) Training

K Nearest Neighbor is a supervised learning algorithm where the result of new instance query is classified based on majority of K-nearest neighbor category. The purpose of this algorithm is to classify a new object based on its attribute values and training samples. K Nearest neighbor algorithm uses neighborhood classification as the prediction value of the new query instance. KNN works based on the minimum distance from the query instance to the training samples to determine the K-nearest neighbors. In image annotation, we compute the K nearest images for a given image. We use Euclidean distance as a similarity measure. Depending on the features set of the images, the computation of the distance differs. For Feature set 1 where we extract features for the whole image, we compute the Euclidean distance directly as follows:

$$d(x_i, x_j) = \sqrt{\sum (x_i - x_j)^2} \quad (1)$$

where  $x_i$  and  $x_j$  are the feature set of image  $i$  and image  $j$ . For feature set 2 where we partition the image into grids, we compute the distance by summing up the Euclidean distance between each two corresponding grids in the images as follows:

$$gd(n, x, y) = \sum_{k=1}^n d(x_k, y_k) \quad (2)$$

where  $n$  is the number of grids and  $d(x, y)$  is the Euclidean distance between feature sets  $x$  and  $y$ , and  $x_k$  is the grid  $k$  of image  $x$ . After gathering  $K$  nearest neighbors, we take simple majority voting of these  $K$ -nearest neighbors to be the prediction of the query instance.

## 2.5. Keyword/Example based Search

Our system facilitates a keyword search and query by example. The user clicks on a concept keyword and the related/closed images to that concept are displayed. Figure 3 presents a keyword search results. The user clicks on a concept keyword and the related/closed images to this concept are displayed. Prediction here is based also on either SVM or KNN in which the user can choose. Notice that the grids on the right side of the images in Figure 3 represents the feedback box in our architecture. In this portion of our interface, the user can express his satisfaction of the results and can add new ground truth to the database. Here positive feedback is taken into consideration. For example, checking a grid means that grid is relevant to the concept that is being queried. After a reasonable number of ground truth is added, new training round is performed. The idea here is to increase the training set incrementally by utilizing users' feedback. User can also retrieve relevant images based on query by example (QBE). In QBE testing phase, we test our system by searching based on Query-By-Example. We have a fixed list of example images. Once we click on one of these images, the closest images are displayed. The similar image will be determined based on the Euclidean distance between QBE image and the test image. As mentioned in Section 2, the distance is computed differently depending on the feature set chosen. For feature set 1 where we extracted features for the whole image, we compute the distance using Eq. 1. For feature set 2 and 3 where we partition the images into grids, we use Eq. 2 to compute the distance. The lower the distance, the more similar will be the image with the query image. For example, Figure 4 presents the results of searching the closest images to a given image. The image on the top is the

example, and the two rows of images are the closest images to it.

We consider the color features in feature set 2 and 3 more than the shape and texture. Because of that different concepts with similar colors might be shown in the results. For example, our system might mistakenly show building results for the concept people. That is because the brown color is dominant in many building and people images. To mitigate that, we capture the user feedback at the end of each query results. The user can evaluate the results by choosing those grids/images that are relevant to the example image (or keyword) the most. As a result of that, the training set increases and a new training round is applied. The system performance enhances using the feedback, training set is incrementally increased, and the search results improves.

## 3. Discussion and Future Work

The result we obtained is promising, and depends on the feature we use. In the case of whole image extraction (for feature set 2 & 3), we notice that the results are concept dependent. For example, images of concepts with similar colors might be predicted incorrectly. For example, the rock may be confused by the building because both the rock and the building have similar colors. On the other hand, concepts that has distinct color features, such as snow and sky, get perfect results, 100% accuracy (see Figure 3, which shows the prediction results for the concept "snow"). On the other hand, for grid based feature extraction case, results will be satisfactory if images have multiple concepts.

In future, we would like to extend our work along the following directions. First, we would like to extract more low level features using MPEG-7 and test the annotation accuracy. Next, we would like to modify SVM/KNN classifier so that it can be used in incremental basis along with real time usage.

## Acknowledgement

This research was supported in part by gift from NOKIA Research Center, Irving, Texas, USA.

## 4. References

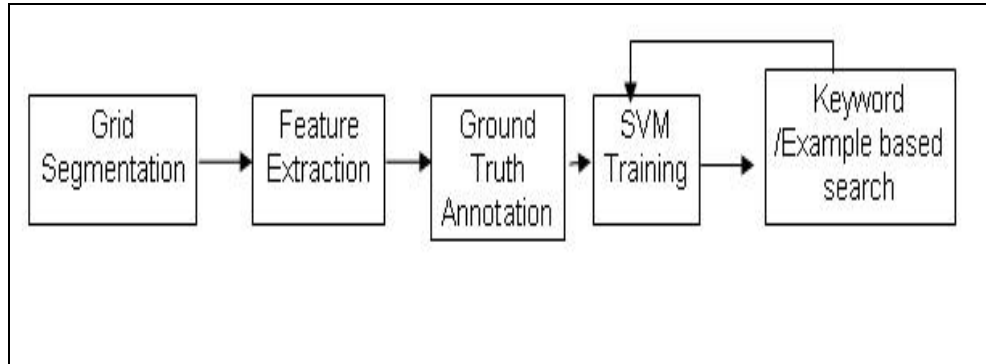
- [1] <http://kdd.ics.uci.edu/databases/CorelFeatures/CorelFeatures.data.html>

[2] A. Natsev, M. Naphade, J. Tesic, "Learning the Semantics of Multimedia Queries and Concepts from a Small Number of Examples," in Proc. of The 13th Annual ACM International Conference on Multimedia (MM 2005), Singapore, November 2005, pp. 598-607.

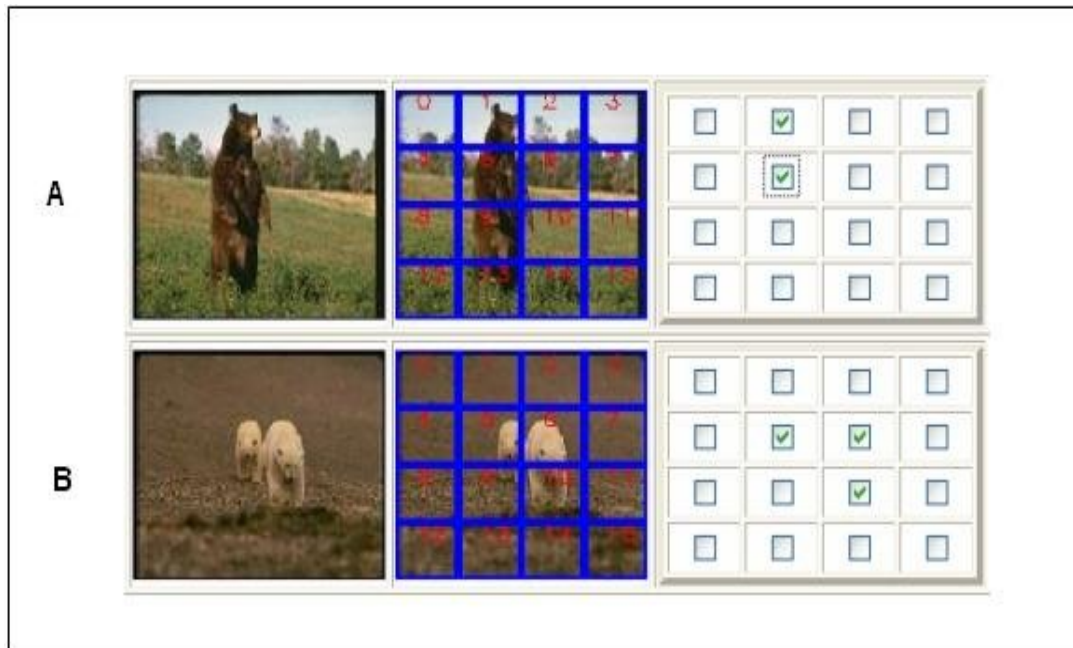
[3] J. R. Smith and S. Chang, "VisualSEEK: A Fully Automated Content-Based Image Query System," *ACM Multimedia 1996*, pp. 87-98.

[4] L. Zhu, A. Rao, and A. Zhang, "Theory of Keyblock-based Image Retrieval," *ACM Transactions on Information Systems*, vol. 20, No. 2, April 2002, pp. 224-257.

[5] Y. Chin, L. Khan, L. Wang, and M. Awad, "Image Annotations By Combining Multiple Evidence & WordNet" in Proc. of The 13th Annual ACM International Conference on Multimedia (MM 2005), Singapore, November 2005, pp. 706-715.



**Figure 1 Architecture of Our Image Annotation Mechanism**



**Figure 2 Ground Truth Collection**

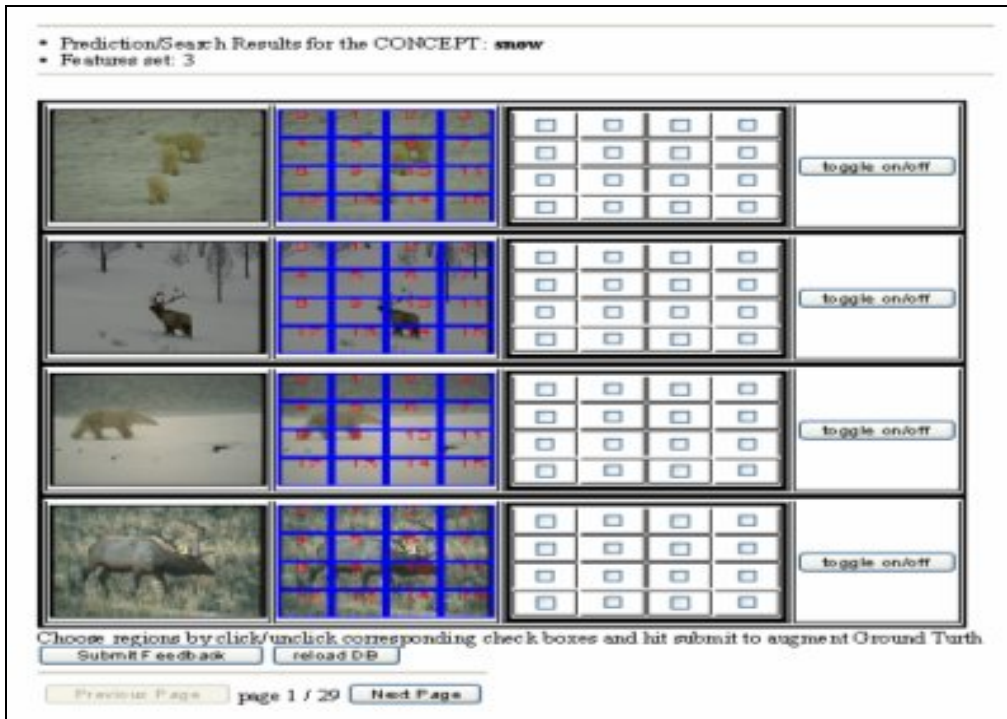


Figure 3 Keyword query based search

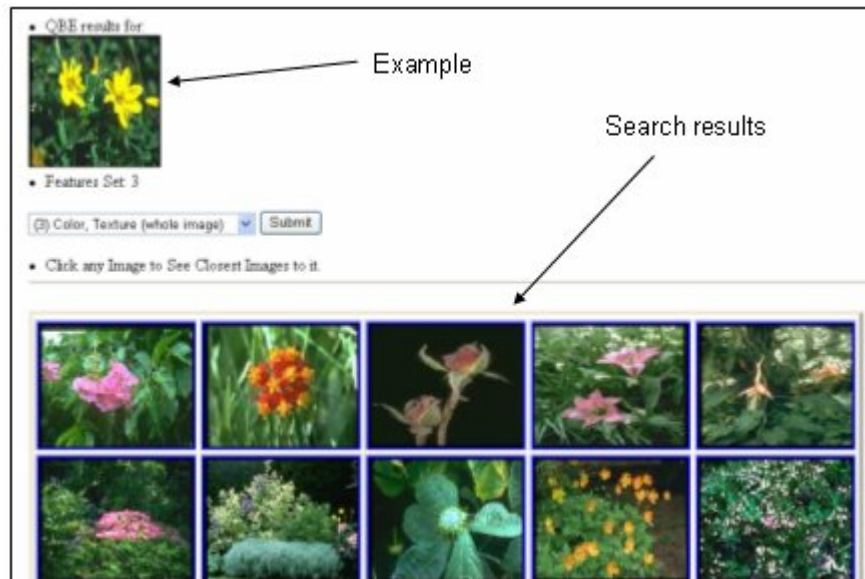


Figure 4 Query by Example Search Results