# Image Similarity Measure using Color Histogram, Color Coherence Vector, and Sobel Method

**Kalyan Roy[1], Joydeep Mukherjee[2]**

[1]Jadavpur University, School of Education Technology,
M Tech Scholar, Kolkata, India
*kalyanroy.08.13@gmail.com*

[2]Jadavpur University, School of Education Technology,
Asst. Professor, Kolkata, India
*joym761@email.com*

**Abstract:** *Image Retrieval means searching, browsing, and retrieving the images from an image database. Two different methods are used for image retrieval, namely text based image retrieval and content based image retrieval techniques. But now Text based search technique is old one. In Content Based Image Retrieval many visual feature like color, shape, and texture are extracted, next when we query an image its feature are compared with the stored feature and we get most similar kind of image. In our proposed method we firstly extract low level image feature like- color histogram, color coherence vector. Then we add edge detection technique sobel edge detection method to get better output. Finally we use Manhattan distance to find the similar images from our database.*

**Keywords:** Histogram, Color Coherence Vector, Sobel Edge detection method, Manhattan distance.

## 1. Introduction

Content Based Image Retrieval is a popular research area after increasing large amount of multimedia information. In different area of Medical Diagnosis, Military, Retail Catalogs we can see the application of image. The importance of Content-Based Image Retrieval is motivated by the increasing desire for retrieving images from growing digital image databases over the Internet.

Content Based Image Retrieval, means extracting a range of images which is relevant with the given image from a large database of images. In Content Based Image Retrieval we first need to extract the feature of image data and store it in a table in row-column format. Then perform any similarity measurement algorithm to find similar kind of image data. Here firstly one question comes in mind "what is feature?" Feature means some visual similarity which can uniquely identify an image from others. There exist many image features like color, shape, texture etc. We need to extract these features from image. Here we get another term "feature extraction" which means mapping the image pixels into the feature space. Using this extracted feature we can measure similarity between indexed image and query image.

Usually, human feel more sensitive to the color feature than to texture and shape. Also, computer used to describe images by RGB form, color feature extraction can save much time because of computing more easily. So, retrieval systems using color is most popular. There are many methods to retrieve image using color. This color feature can be extracted in many ways like Histogram [1], Color moments [2], Color Correlogram [3], color coherence vector [4] etc. There also many more method to extract Texture feature, which can be extracted using Gray Level Co-occurrence matrix (GLCM) [5], Gabor Filter Response [5] etc.

But only color feature does not consider the spatial information of pixel, which may produce a problem in similar color distribution in different image. So, it is better to consider any other feature to retrieve similar image. Edge detection technique can be added with color feature to get more accuracy. There are many edge detection techniques like Robert, Sobel, Prewitt, Kirsch, Robinson, Marr-Hildreth, LoG and Canny Edge Detection [6]. Many researchers tells that Marr-Hildreth, LoG and Canny produce the same result where Kirsch, Robinson produce the same also. But in our paper we consider sobel edge detection method as it is simple to calculate and give prominent result.

In the paper [7] retrieving similar image is done based on color and shape based method. But in color histogram does not consider the spatial information of a pixel which may result similar color distribution for different image. So in this paper we propose an image retrieving method based on color histogram analysis, color coherence vector and sobel edge detection technique.

**Our first approach** on feature extraction is to extract the color histogram value from an image. The color histogram for an image is constructed by quantizing the colors within the image and counting the number of pixels of each color. Then we take a summation of it and find the mean and standard deviation from the color histogram. Finally it stored in a 1D array. This value is calculated for every image in the database.

**Our second approach** on feature extraction is color coherence vector value analysis. In CCV each histogram bin is partitioned into two types: coherent and incoherent. All pixels which informally fall into a similar colored region are called coherent type. Otherwise it calls incoherent in type.

Here in our method we measure Red, Green, Blue, Red+ Green, Red+ Blue, Blue+ Green, White and Black Coherent-

Incoherent type and store it in a 1D array. This value is calculated for every image in the database.

**Our third approach** in feature extraction is edge detection of an image. For this purpose we used Sobel's Edge detection technique which gives most effective and prominent results. And it gives low computation complexity for our proposed system. Here we get another 1D array.

Finally we combine each all the three extracted feature in one matrix and repeated the process for every image. Then we apply Manhattan distance to measure similarity between the images.

The rest of the paper organized as follows, section 2- related work surveyed by us , section 3- System overview , section 4- Proposed Work, next section is about experimental result that contains snapshot of results, and this is followed by future work.

## 2. Related Work

In past years, many papers have been presented to measure similarity between query image and image database. Most of the Content Based Image Retrieval system stores the image features in the feature vector. Once the features are represented as a vector it can be used to determine the similarity between the images. Content Based Image Retrieval systems used different technique to measure similarity between the images. In retrieval stage, we also find the feature vector for query image and find the similarity between the query vector and early store feature vector. The similarity measure is used between the query image and all stored images in a database. After this process the images which are minimal distance according to query image are retrieved and ranked according to distance. Image feature is categorizes in two group. Pixel domain and compress domain feature extraction technique. In pixel domain technique we consider visual feature such as color, texture and shape.

Color is one of the most outstanding features of the image, it is the most important human visual content and it is very easy to calculate. It is widely concerned by many researchers because it does not effected by natural rotation, scaling and translation of a image. Many methods can be used to describe color feature such as color histogram [8], color correlogram [9], color moment [10], color structure descriptor (CSD), and scalable color descriptor (SCD) [11]. The color distributions in the images can be represented using the Color Histograms. The histogram of the query image and the database images are compared for the retrieval. But this method fails if two different images have the same color distributions. To avoid this color coherence vectors can be used which store spatial information of an image in vector format.

Color Histogram only describes global color distribution of a image but cannot reflect the spatial information of an image. To overcome this problem flaw, Stricker and Orengo[12] combine color moments and cumulative color histogram method, both of which effectively describe color spatial distributions. To find similar images from large scale image

database Smith and Chang [13] came up with color set as an alternative of cumulative color histogram. But color histogram and color moments cannot express color spatial location of an image, so Pass [14] gives another method named color coherence vector. Each single feature only represents partial attribution of an image. So in our method we combine color histogram and color coherence vector as a combined feature extraction method to get better accuracy and accurate result.

Shape detection of an image is an important feature for object recognition. Shape description or representation of edge is an important issue to measure similarity between images. There exists different edge detection technique to detect edges of objects in the image. Like Robert, Sobel, Prewitt, Kirsch, Robinson, Marr-Hildreth, LoG and Canny Edge Detection [15]. In shape detection technique region based method are most commonly used. Region based method are based on continuity of image pixel. These techniques divide the entire image into sub regions depending on some rule like the entire pixel in one region must have the same gray level. From the experiment it is observed that Kirsch and Robinson and Prewitt gives almost same result where Marr-Hildreth, LoG, Canny falls in a group and Robert and Sobel falls in another group. It is seen that canny gives best result but very high computation complexity. Compare with that Sobel gives almost same result with low computation complexity. So we consider Sobel edge detection method as a feature extraction method with color histogram and color coherence vector.

## 3. System Overview

In our system we propose a multi feature extraction method e.g. color histogram, color coherence vector and canny edge detection technique. After that for similarity measure we take Manhattan distance. The overall system overview is shown in the figure 1.
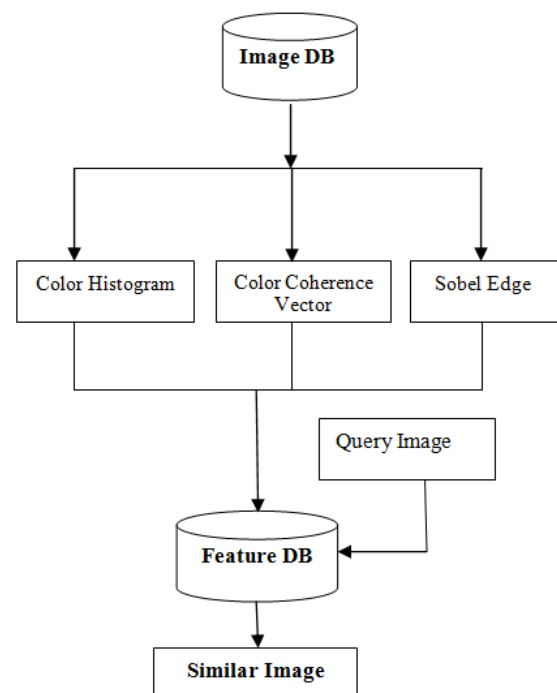


**Figure 1:** System Overview

## 3.1 Color Histogram Analysis Method

Color Histogram is a commonly used feature in image retrieval. It is very popular because Color histograms are computationally trivial to compute. Small changes in camera view point does not does effect in the histogram. Many researchers in computer vision and image processing investigate histogram.

In our method we consider all images are scaled to 256 X 256 in height and width. We discretize the colors pace of the image such that there are n distinct colors. In our method we consider 8-bit image. So, there will be 256 different bins of colors. We take a summation of these color bins and then find mean and standard deviation from this value.

A color histogram H is a vector $<h1; h2; \ldots; h_i>$ in which each bucket $h_j$ contains the number of pixels of color j in the image. In our proposed approach image histogram feature is extracted for three color space Red, green and Blue. And the histogram value matrix holds no of pixels of 256 different bins in Figure 2.
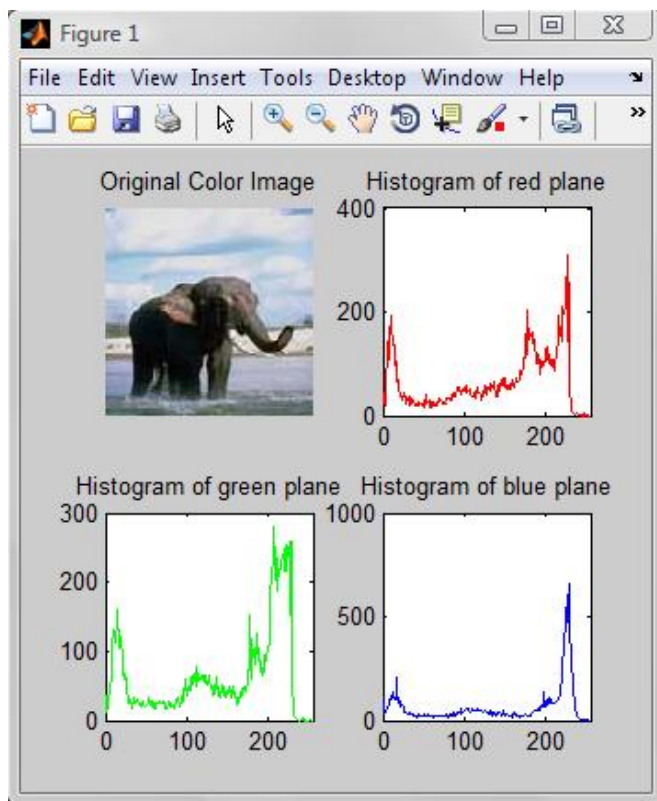


**Figure 2:** Color Histogram Analysis

### 3.2 Color Coherence Vector Analysis

There exist many CBIR systems, which retrieve images by color histogram. But only histogram based method cannot give best results if two pictures have same color distribution. For this reason we Consider Color Coherence Vector also as a color feature extraction method. We define a color's coherence as the degree to which pixels of that color are members of large similarly-colored regions. We refer to these significant regions as coherent regions, and observe that they are of significant importance in characterizing images. Our coherence vector classifies image pixels as either coherent or incoherent in type. Coherent pixels are a part of some sizable contiguous region, while incoherent pixels are not. A coherent pixel is part of a large group of pixels of the same color, while an incoherent pixel is not. We determine the pixel groups by computing connected components. In this method some region is created based on color information. Details Method to extract color coherence vector is described in the section-Proposed method.

### 3.3 Edge Detection Method

Image segmentation is the process of partitioning an image into multiple region or sets of similar pixel. Essentially, in image segmentation all different objects which have same color, texture, or intensity features are merged together. There are different approaches to find the similarity between images. These are (i) finding region based on discontinuity in intensity level (ii) finding threshold value to distinguish different segment and (iii) based on the finding region directly. Which method we used it depends on the problem being considered. Segmentation method based on finding the region directly finds the abrupt changes in the intensity value. So these methods are called as edge or boundary based method. Edge detection techniques are generally used to find the discontinuous area in gray level image. There are various techniques of edge detection available i.e. canny method, LoG method, Sobel method, Prewitt method, Robert method, Krish method, Robinson Method etc. Here we seen that Marr-Hildreth, LoG, Canny falls in a group and Robert and Sobel falls in another group. Canny gives good output with high computational complexity, where sobel give almost same output with low complexity. For this reason we consider sobel edge detection method introduced by sobel in 1970.
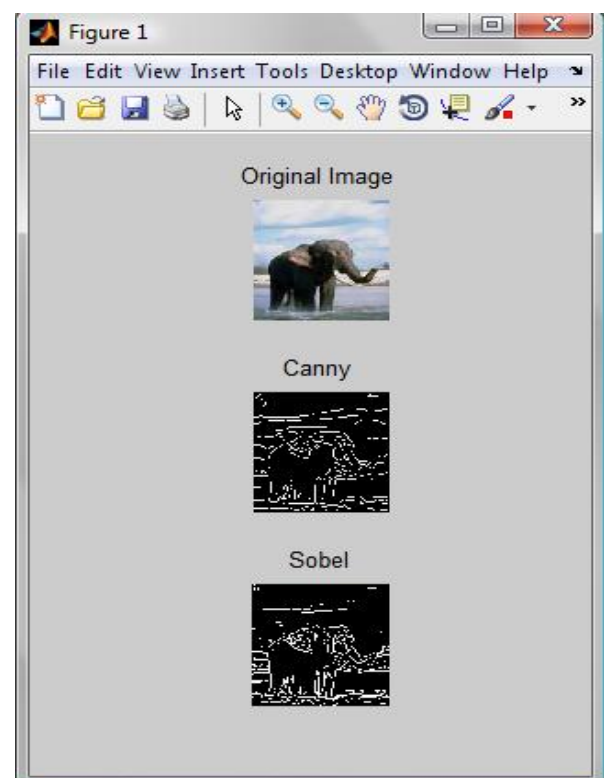


**Figure 3:** Edge Detection technique

## 4. Proposed Method

As mentioned earlier, our proposed approach uses three feature extraction methods i.e. Image color histogram extraction, image color coherence vector values extraction and Image edge detection using sobel edge detection technique. In our proposed work we consider 8 bit color image.

### 4.1 Image color Histogram Analysis

Step 1: Read the image.

Step 2: Store the value of each red, green, blue component in three different arrays.

Step 3: Find the image histogram of red, green, blue component by Matlab's own histogram computational method.

Step 4: Store the histogram value of red, green, blue component in three arrays.

Step 5: Calculate the sum of 256 different bin of red component. If there are N no of color bin of red component then we get M (summation of each color bin) by the below said formula:

$$M = \sum_{i=1}^{N} hi$$

Apply this method for green and blue component also.

Step 6: Then find the mean of color histogram using

$$\bar{x} = \frac{\sum_{i=1}^{N} (i*hi)}{M}$$

Step 7: Then find the standard deviation of color histogram using

$$\sigma = \sqrt{\frac{1}{M} \left( \sum_{i=1}^{N} hi * (i - \bar{x})2 \right)}$$

Step 8: Finally store this three value- summation, mean, standard deviation in a 1D array.

Step 9: Repeat Step 1-8 for every images in the database.

### 4.2 Image color Coherence Vector Analysis

Step 1: Read the image.

Step 2: Calculate the size of the image

Step 3: Compute 4 different array from original array which is created after reading the image. In the first array, 1st row value of the original image will be stored as last row. In second array, last row value of the original image will be stored as first row. In the next two array column value of the original image will be replaced in this manner.

Step 4: Convert the entire 5 image array into double data class in MATLAB.

Step 5: Take the average of 5 image array.

Step 6: Then make a matrix which is combination of Red, Green, Blue, (Red +Green), (Red +Blue), (Blue +Green), White and Black value of an image and they are initializes as default value.

Step 7: Now we measure the magnitude of difference from the average matrix which is get in Step5 and the value said in step6 for each 8 distinct element.

Step 8: Then we calculate the threshold value which will give best result.

Step 9: Using the threshold value calculate how many Red channel connected pixel we get and we called it as a Red Coherent(Re), and find also how many pixel do not contain any red information and we called it as a Red incoherent(Ri). We continue this process for all 8 different color information said in step6.

Step 10: Finally we get 16 different Coherent and Incoherent value from the image and store it in a previously created array.

Step 11: Repeat step 1 to 10 for every image in the database.

### 4.3 Image Sobel Edge Detection Method

Step 1: Read the image.

Step 2: Convert the RGB image into gray color space.

Step 3: Then convert the value into double data space in MATLAB.

Step 4: Using sobel mask we find the x-direction and y-direction derivative $G_x$ and $G_y$.

$G_x = (Z_7 + 2Z_8 + Z_9) - (Z_1 + 2Z_2 + Z_3)$

$G_y = (Z_3 + 2Z_6 + Z_9) - (Z_1 + 2Z_4 + Z_7)$

Step 5: Then compute the gradient value for each image pixels in x-direction and y-direction.

Step 6: Compute the summation of x-direction and y-direction gradient value and add this 2 value in the previously created 1D array.

Step 7: Repeat Step 1-6 for every images in the database.

## 5. Results and Discussion

All the features are extracted from different image and stored into a single matrix. In the training phase we need to train our system to recognize an image properly by its extracted feature values. Here in our system we extract feature of 50 image of size 256 X 256.Next when we query an image, its

feature is extracted and stored in a matrix. Then we measure the similarity with the stored feature using Manhattan distance. Our experiment gives most similar kind of image with a prominent result. In the below figure (Figure 4 and Figure 5) we summarize our results.
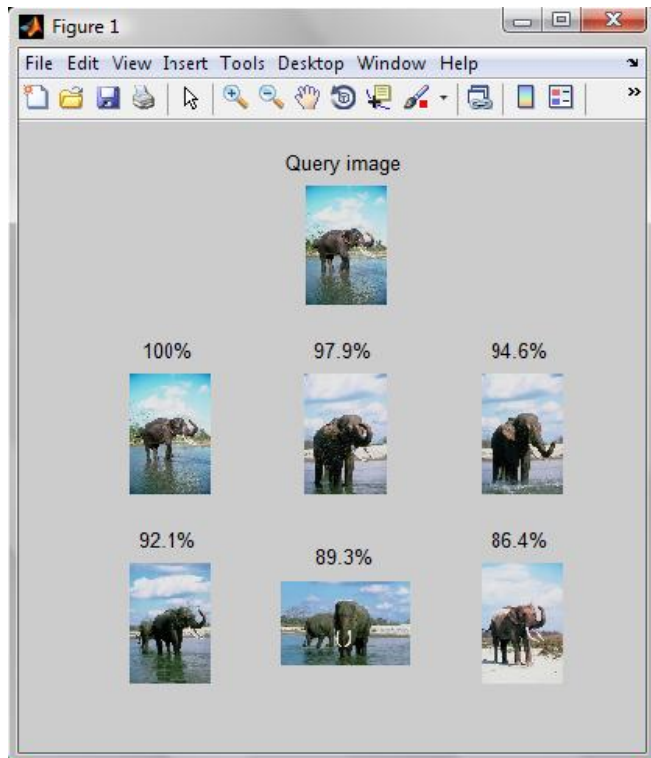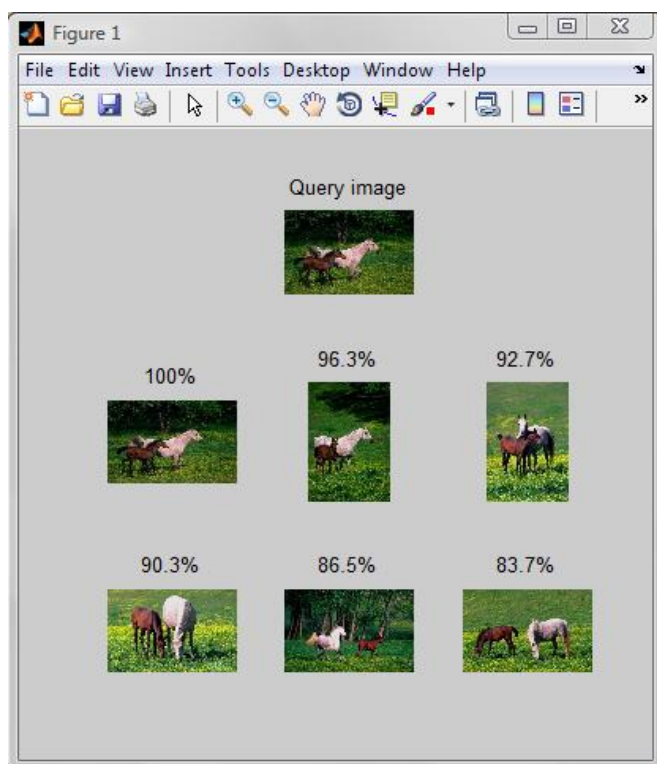


**Figure 4:** Find most similar images



**Figure 5:** Find most similar images

# 6. Future Work

Here in this experiment we use Manhattan distance for similarity measure. Some other approaches i.e. neural network method, Chebyshev distance may also give better accuracy. Here we use stored database for image feature extraction and calculate this feature using MATLAB. Properly maintained database may reduce the computational speed. This feature concentrates only retrieving the similar images, but in future this work can be enhanced to clustering the similar kind of images.

## References

[1] M L Swain, H H Ballard, "color indexing", Internationaln Journal of Computer Vision, 7(1):pp.11-32, 1991.

[2] M Stricker, M Orengo, "Similarity of color images", In: W Niblack, R C Jain eds. Pro of SPIE Storage and Retrieval for Image and Video Databases, Vol 2420. San Jose, CA, USA: SPIE Press, p381-392, 1995.

[3] J Huang, S R Kumar, M Mitra et al, "Image indexing using color correlograms", Proc of IEEE Conf on Computer Vision and Pattern Recognition, San Jose, Puerto Rico, USA: IEEE CS Press, p762-768,

[4] 1997

[5] Reza Ravani. Mohamad Reza Mirali &, Maryam Baniasadi, *Parallel CBIR System Based on Color Coherence Vector.* International Conference on Systems, Signals and Image Processing, 2010,

[6] D.S. Guru, Y.H. Sharath & S. Manjunath, Texture Features and KNN in Classification of Flower Images, IJCA Special Issue, 2010

[7] Ehsan Nadernejad. Sara Sharifzadeh &, Hamid Hassanpour, *Edge Detection Techniques: Evaluations and Comparisons.* Applied Mathematical Sciences, 2008

[8] You Fu-cheng. Zhang Cong &,, *The Technique of Color and Shape-based Multi- feature Combination of TradeMark Image Retrieval,*, 2010

[9] R. Brunelli and O. Mich, "Histograms Analysis for Image Retrieval," Pattern Recognition, Vol.34, No.8, Pp1625-1637, 2001

[10] J. Huang, S. R. Kumar, M. Mitra, W. J. Zhu, and R. Zabih, "Image indexing using color correlograms," IEEE Int. Conf. Computer Vision and Pattern Recognition, San Juan, Puerto Rico, Jun. 1997, Pp.762–768.

[11] A. Del Bimbo, M. Mugnaini, P. Pala, and F. Turco, Picasso, "Visual querying by color perceptive Regions", In Proceedings of the 2nd International Conference on Visual Information Systems, San Diego, December'97, pages 125-131, 1997

[12] ISO/1EC JTC 1/SC29/WG 11/N4062: "MPEG-7, Requirements document" 2001

[13] Fazal Malik &, Baharum Bin Baharudin, *Feature Analysis Of Quantized Histogram Color Feature For Content-Based Image Retrieval Based on Laplacian Filter.*, International Conference on System Engineering and Modeling, 2012,

[14] SMITH J R, CHANG S F. Tools and techniques for color image retrieval [A]. In SPIE of the Storage & Retrieval for Image and Video Databases IV, February 1996, 2670: 426-437

[15] PASS G, ZABIHR. Histogram refinement for content based image retrieval [C]. In Proc IEEE Workshop on Applications of Computer Vision, 1996, (3): 96-102.

[16] O.R. Vincent & O. Folorunso, *A Descriptive Algorithm for Sobel Image Edge Detection,* Informing Science & IT Education Conference (InSITE), 2009,

## Author Profile

**Kalyan Roy** has received his B.E degree in Computer Science and Engineering from Burdwan University. Presently he is perusing his M. Tech degree from School of Education Technology, Jadavpur University, Kolkata. His research interests comprise of Image processing, Audio Signal Processing, Genetics Algorithm.

**Joydeep Mukherjee** obtained M. Tech degree from Jadavpur University. Currently he works as an Asst. Prof. in School of Education Technology, Jadavpur University, Kolkata. His research interests comprise of Digital image processing, Character Recognition, Pattern Recognition, E-Learning, M-Learning.