# A Hybrid Image Mining Technique using LIM-based Data Mining Algorithm

C. Lakshmi Devasena
Department of Software Systems
Karpagam University
Coimbatore-21

M. Hemalatha
Department of Software Systems
Karpagam University
Coimbatore-21

## ABSTRACT
The field of image retrieval and mining has become a vibrant research area due to speedy enhancement in the volume of digital image databases. Nowadays, a large portion of information is in visual form; it is essential and certainly pleasing to search for images by content. Image mining has a variety of applications in various sectors like medical diagnosis, biology, remote sensing, space research, etc. This research is to determine the exact images while mining an image (multimedia) database and proposes a novel approach for mining images using LIM based image matching technique with neural networks. This process is independent of too many parameter setting to generate a robust solution. It is designed and implemented on MATLAB and is tested with the images of various databases. Appropriate measures were devised to evaluate the performance of the system. The performances of the LIM based image matching technique results were noteworthy and comparable. While comparing with the number of false retrievals with the correct retrievals, the anticipated system performance level will be suited for several simple day to day multimedia database applications and image mining systems. The image mining system derived from the LIM based image matching technique provided promising results.

## General Terms
Multimedia Data Mining – Image Mining.

## Keywords
Discrete Cosine Transform, Image Mining, Image Signature, Lorenz Information Measure.

## 1. INTRODUCTION
The growth of Internet not only causes an explosive growing volume of digital multimedia data, but also provides people more ways to get those images. The significance of an effective technique in searching and retrieving images from the huge collection cannot be exaggerated. One approach for indexing and retrieving image data is done using manual text annotations. The annotation can be used to seek out images indirectly. But there are numerous problems with this approach. First, it is very difficult to depict the contents of an image or a video scene using only a few keywords. Second, the manual annotation process is very slanted, confusing, and deficient. Those problems have created great demands for automatic and effective techniques for Content-Based Image Retrieval (CBIR) System. Most of the system use low-level image features such as color, texture, shape, edge, etc., to index and retrieving images.

It's because the low-level features can be computed in an efficient manner.

The unique objective of computer vision is to recognize a single image in a scene by identifying the objects and their structure and spatial arrangements in that scene. This is referred as image understanding.

The goal of this project is to derive a new method for detecting images. Here we used the method Lorenz Information Measure (LIM) for representing features extracting from the images for retrieval.

This research proposes a novel approach for recognizing an image. This approach attempts to make the process mostly independent of any parameter setting to generate a robust solution. The proposed model of the system is designed and implemented on MATLAB. The performance of the proposed system is tested and produces promising results.

## 2. LITERATURE REVIEW
In [1] the basic and advanced concepts of data mining are described which includes the various types of basic algorithms used for mining. In [2] the basic thoughts of digital image processing applied in the field of Image retrieval and Image mining as well as the feature extraction techniques like gray level and histogram equalization are given. The earlier research work in image mining and retrieval are summarized in the survey papers [3], [4] [5], [6] and [7]. The concepts regarding content based image retrieval with high level semantic features are explained in paper [8]. The trends and developments in content based image retrieval systems are given in [9]. The theoretical aspects of content based image retrieval are derived from the paper [10]. The idea about Lorenz Information Measure is taken from the research papers [11] and [12]. The theoretical concepts and applications of Discrete Cosine Transform are derived from [13], [14] and [15]. Our proposed method combines Lorenz Information Measure with Discrete Cosine Transform and produces the result.

## 3. METHODOLOGY AND DESIGN
There is lots of image processing and statistical and machine learning techniques involved in the design of the proposed system.

### 3.1 Techniques Used in Proposed Model
#### 3.1.1 Histogram
Gray level / Color histogram shows the frequency of occurrence of each gray level / color in the image versus the gray level itself

and provides a global description of the appearance of the image. The histogram with gray levels in the range [0, L-1] of a digital image is a discrete function.

$$P\ (r_k) = n_k / n$$

where,
$r_k$ - Gray level K

$n_k$ - Number of pixels in the image with the gray level $r_k$.

n - Total number of pixels contained in the image.

K = 0, 1, 2… L-1.

L = 256.

$P\ (r_k)$ gives an estimate of the probability of occurrence of gray level $r_k$.

### 3.1.2  Color Histograms

The focal method of symbolizing color information of images in CBIR systems is done through color histograms. A color histogram is a type of bar graph, where each bar represents a particular color of the color space being used. The bars in a color histogram are referred to as bins and they represent the x-axis. The number of bins depends on the number of colors there are in an image. The y-axis indicates the number of pixels there are in each bin. That is, the number of pixels in an image with a particular color.  An example of a color histogram in the HSV color space is shown in Fig 1.
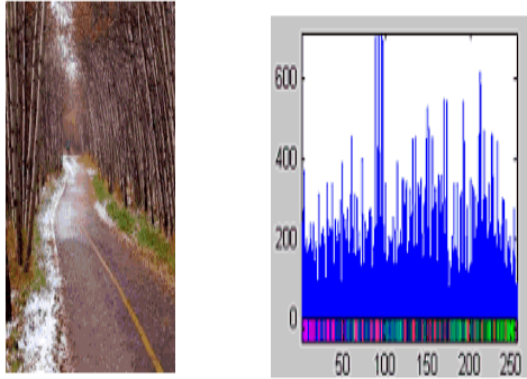


**Fig 1 : Sample Image and its Corresponding Histogram**

### 3.1.3  The Lorenz Information Measure (LIM)

Generally, Lorenz Information Measure (LIM) widely used in economics. Rorvig was the first to suggest use of general features extracted from the images for retrieval and represented as LIMs.  The Lorenz Information Measure (LIM) $(P_1,…,P_n)$ is defined to be the area under the Lorenz information curve. The area of LIM $C_a$ is greater than the area of LIM $C_b$. Clearly, 0 < = LIM $(P_1,……,P_n)$ < = 0.5. If the probability vector is $(P_1,……,P_n)$, then LIM $(P_1,……,P_n)$ can be measured by the first ordering $P_i$'s, and then calculating the area under the piecewise linear curve. Because LIM $(P_1,……,P_n)$ (which can be expressed as the sum of $f(P_i)$, and $f(P_i)$) is a continuous convex function, LIM $(P_1,……,P_n)$ is considered as an information measure.

Spontaneously, the LIM can be considered as a universal content-based information measure. To calculate the area of

histograms, the histogram intervals are arranged from low to high, and the resulting off-diagonal shape measured through differentiation. It is based on the concept of using minimum number of gray level changes to convert a picture into a desired histogram i.e. with a single value.

### 3.1.4  The Discrete Cosine Transform

The DCT can be used to Create feature based Image Profile. The Discrete Cosine Transform is a real domain transform which symbolizes the entire image as the coefficients of distinct frequencies of cosines (which are the source vectors for this transform). The DCT of the image is calculated by taking 8x8 blocks of the image, which are then transformed individually. The two dimensional DC Transform of an image gives the result matrix such that top left corner signifies lowest frequency coefficient whereas the bottom right corner is the highest frequency.

The 1-D *discrete cosine transform* (DCT) is defined as

$$C(u) = \alpha(u) \sum_{x=0}^{N-1} f(x) \cdot \cos\left[\frac{(2x+1)u\pi}{2N}\right]$$

In the same way, the inverse of the DCT is defined as

$$f(x) = \sum_{u=0}^{N-1} \alpha(u) C(u) \cdot \cos\left[\frac{(2x+1)u\pi}{2N}\right]$$

where
$$\alpha(u) = \begin{cases} \sqrt{1/N} & \text{for } u = 0 \\ \sqrt{2/N} & \text{for } u = 1,2,...,N-1 \end{cases}$$

The equivalent 2-D DCT, and its inverse are defined as

$$C(u,v) = \alpha(u)\alpha(v) \sum_{x=0}^{N-1}\sum_{y=0}^{N-1} f(x,y) \cdot \cos\left[\frac{(2x+1)u\pi}{2N}\right]\cos\left[\frac{(2x+1)v\pi}{2N}\right]$$

and

$$f(x,y) = \sum_{u=0}^{N-1}\sum_{v=0}^{N-1} \alpha(u)\alpha(v) C(u,v) \cdot \cos\left[\frac{(2x+1)u\pi}{2N}\right]\cos\left[\frac{(2x+1)v\pi}{2N}\right]$$

The benefit of DCT is that it can be expressed without complex numbers. 2-D DCT is also separable (like 2-D Fourier transform), i.e. it can be obtained by two subsequent 1-D DCT in the same way than Fourier transform..

## 3.2  The LIM based Image Matching Technique

In the proposed method, twelve content-based image features are derived from LIM of histogram of the following the image.1) Red,    2) Green, 3) Blue, 4) Hue, 5) Saturation, 6)

Luminance, 7) DCT of Red, 8) DCT of Green, 9) DCT of Blue, 10) DCT of Hue, 11) DCT of Saturation, 12) DCT of Luminance. In this model, there will be 4 major steps to perform image retrieval based on the similarity:

1. Load Query Image (Image we want to search for or find images similar to this)

2. Generate Signatures of Key Image using Lorenz Information Measure(LIM) and other Suitable Techniques

3. For every images in the database Load and generate the signatures

4. Calculate Distance for Key Image Signature and Database Image Signature using a Suitable Distance Measure and find the possible matches.

### 3.2.1  The Image Signature Creation

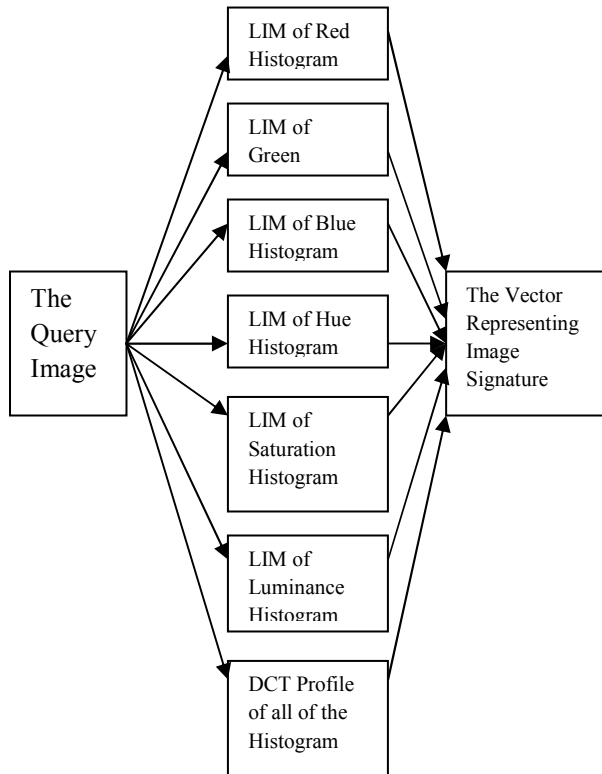The diagram explaining Image Signature Creation of the proposed system is shown in Fig 2.



**Fig 2: Image Signature Creation**

### 3.2.2  Proposed LIM based Image Matching Model.

The diagram explaining the proposed LIM based Image Matching Model is shown in Fig 3 and the main interface of the proposed system is shown in Fig 4.
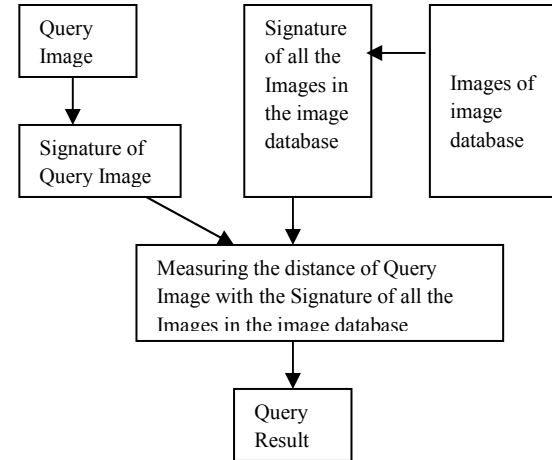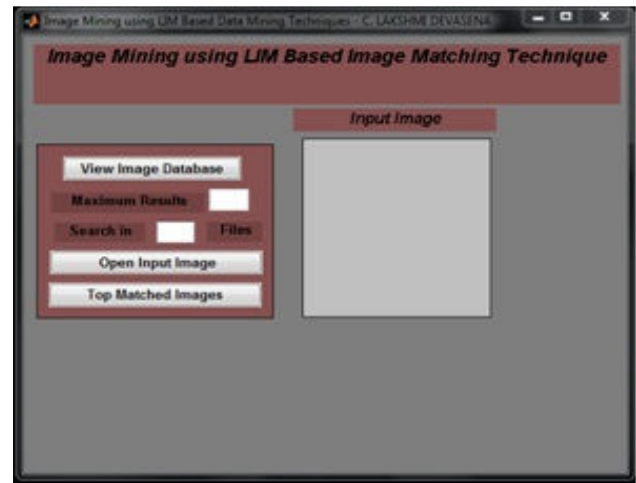


**Fig 3:  Image Matching Model**



**Fig 4: The Main Interface**.

## 4.  RESULTS & DISCUSSION

In this chapter, the results of the LIM based Image/Frame Matching technique has been presented.

## 4.1.  Evaluating the Performance of LIM based Matching Technique

### 4.1.1  Image Database Used.

The Image Database used in this work contains the Images arbitrarily selected from different categories.  The Images were collected from internet sources.  The test images were prepared from some of the images in the image database. An image in test image collection is nothing but a portion of a particular original image in the dataset.  The proposed system is tested with the images in the original data base as well as the newly created ones.

## 4.1.2 The Sample Test

**Step 1: Selection of the Input Image.**

The Input image was selected from the image database or from the newly created test image set (Fig 5). In this particular case, the query image was taken from the newly created test data set.
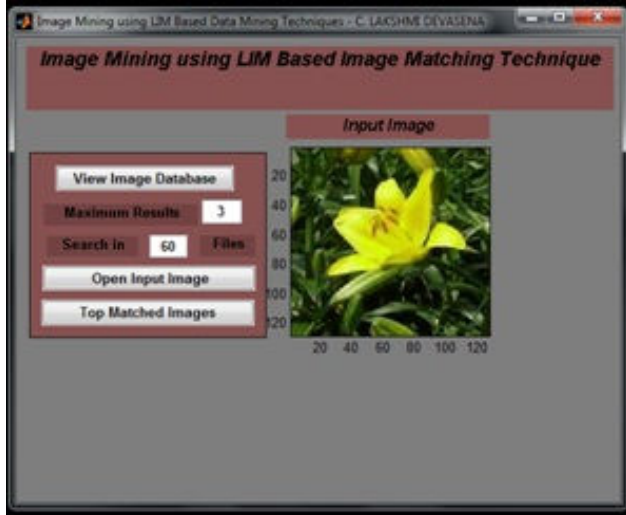


**Fig 5: Input Image Selected**

**Step 2: The Different Histograms used for the Creation of LIM Profile of the Image**

The Fig 6 shows the different feature sets which are used to create the Image Signature of that particular query image. In this method, 12 features were used for the creation of image signature.
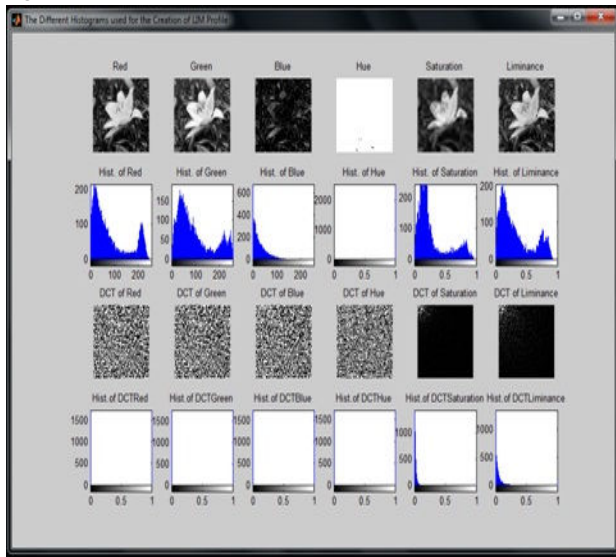


**Fig 6: Image Profiles for Image Signature Creation**

**Step 3: The Search Results**

The Fig 7 shows final results after this particular image query using LIM based image matching method. During the search process, the image signature of each and every image in the target database were derived and compared with the image signature of the query image. For finding the suitable matches, the simple distance measure (sum of square of distances) was used. The result shows the first three nearby solutions of the proposed method. The number of nearby solutions can be increased or decreased by specifying it in the textbox.



**Fig 7: Search Results – LIM Based Image Matching Technique**

## 4.1.3 The Time Study

The Size of the Image Files in the Image Database is 128 x 128. The time taken to search images was studied with different sizes of image database and it is shown in Table1.

**Table 1: The Time study**

| S.No. | Number of images in Image Database | Time Taken to Search (in sec) |
|---|---|---|
| 1 | 10 | 0.57 |
| 2 | 20 | 1.13 |
| 3 | 30 | 1.74 |
| 4 | 40 | 2.21 |
| 5 | 50 | 2.73 |
| 6 | 60 | 3.27 |
| 7 | 100 | 4.45 |
| 8 | 200 | 6.25 |

The following bar chart (Fig 8) shows the Performance of the LIM based Image Matching algorithm with respect to the increase in the size of the image database.
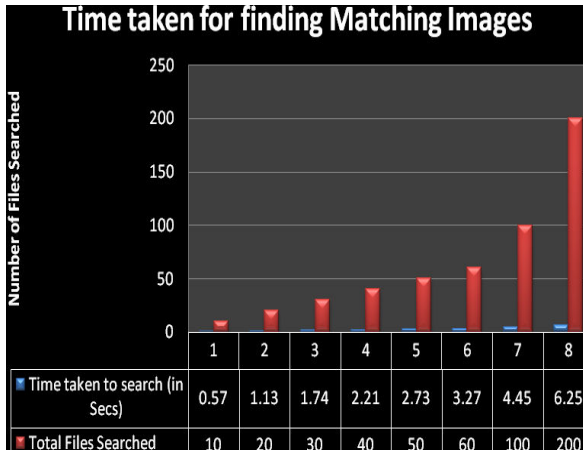
**Fig 8: The Bar Chart Showing Performance**.

## 5. CONCLUSIONS

The area of image and video storage and retrieval within the multimedia domain is an area that is growing rapidly. The vastness of images and videos contained within these purpose built image and video retrieval systems (and other systems) is growing by the day. Content based retrieval and novelty detection at present are still very much a research topic. The technology is exciting but immature, and few operational image and video archives have yet shown any serious interest in adoption. The crucial question that this report attempts to answer is whether content based retrieval will turn out to be a flash in the pan, or the wave of the future. This is why there is a need for developers to urgently address these issue's and provides the global community with full proof contents based retrieval systems. The systems available at the moment are still in the prototype stages and have quite a long way to go before they are considered reliable.

The proposed LIM based Image matching technique has been successfully implemented and evaluated on Matlab. The Images which were used to test the performance of the system were collected from internet sources. The Image Data base used in this work contains different kinds of images. A test image in this collection is nothing but a portion of a particular original image. The proposed system is tested with the images in the original data base as well as the newly created ones. The performance of the LIM based image matching technique has been studied with different images. Suitable measures were formulated to evaluate the performance of the system. The arrived results were significant and comparable. While comparing with the number of false retrievals with the correct retrievals the proposed system seems to be achieved a performance level which will be suited for several simple day to day multimedia database applications and systems. The issues related with improvement of speed by applying other algorithms can be addressed in future works.

## 6. REFERENCES

[1] Tiawei Han and Micheline Kamber. 2001 Data Mining Concepts & techniques.

[2] Fafael C.Gonzalez and Richard E. Woods. 1993 Digital Image Processing 2nd edition. Addision Wesley.

[3] C. Lakshmi Devasena, T.Sumathi, Dr. M. Hemalatha 2010 Image/Video Retrieval Technique: Grand Challenges and Trends. Proceedings of National Conference on Applications of Data Mining in National Security, NCOADMINS 2010, ISBN No: 987-93-80697-51-2, 96-105.

[4] C. Lakshmi Devasena, T.Sumathi, Dr. M. Hemalatha 2011 An Experiential Survey on Image Mining Tools, Techniques and Applications. International Journal on Computer Science and Engineering (IJCSE), ISSN : 0975-3397, Vol. 3 No. 3 Mar 2011, 1155 – 1167.

[5] Wynne Hsu, MongLiLee, Ji Zhang, 2002 Image Mining : trends and developments. Journal of Intelligent Information systems, 2002, 7-23.

[6] Ji Zhang, Wynne Hsu, Mong Li Lee. 2001 Image Mining: Issues, Frameworks And Techniques.

[7] Hu Min  Yang Shuangyuan  2010 Overview of image mining research. 978-1-4244-6002-1 IEEE Explore, 24-27 Aug. 2010, 1868 – 1870.

[8] Prof. Sharvari Tamane 2008 Content Based Image Retrival using High Level Semantic Features. Proceedings of the 2nd National Conference –INDIA Com 2008.

[9] P. Geetha and Vasumathi Narayanan. 2008 A Survey of Content-Based Video Retrieval. Journal of Computer Science 4 (6), ISSN 1549-3636, 474-486.

[10] Stéphane Marchand-Maillet. 2000 Content-based Video Retrieval: An overview.

[11] McMurray, T.  Pearce, J.A. 2002 Theoretical and experimental comparison of the Lorenz information measure, entropy, and the mean absolute error. IEEE Explore ISBN: 0-8186-6250-6, 2002 24-29.

[12] Ki Tai Jeong, Mark Rorvig, Jeho Jeon and Neena Weng 2001 Image Retrieval by Content Measure Metadata Coding.

[13] Syed Ali Khayam. 2003. The Discrete Cosine Transform (DCT): Theory and Application. Department of Electrical & Computer Engineering, Michigan State University.

[14] Andrew B. Watson. 1994 Image Compression Using the Discrete Cosine Transform. Mathematica Journal, 4(1), 81-88.

[15] Wen-Hsiung Chen, Smith. C, Fralick S. 2003 A Fast Computational Algorithm for the Discrete Cosine Transform. 10.1109/TCOM.1977.1093941 IEEE Explore, ISSN: 0090-6778, 25-9 Jan 2003, 1004-1009.