# Content Base Image Retrieval by Clustering

**Junaid Khan**
PRMIT  & R Badnera
Amravati - India

*Abstract: The purpose of this report is to describe our research and solution to the problem of designing a Content Based Image Retrieval, CBIR system. It outlines the problem, the proposed solution, the final solution and the accomplishments achieved. Due to the enormous increase in image database sizes, as well as its vast deployment in various applications, the need for CBIR development arose. Firstly, this report outlines a description of the primitive features of an image; texture, colour, and shape. These features are extracted and used as the basis for a similarity check between images. The algorithms used to calculate the similarity between extracted features, are then explained. Our final result was a built software application, with an image database, that utilized texture and colour features of the images in the database as the basis of comparison and retrieval. The structure of the final software application is illustrated. Furthermore, the results of its performance are illustrated by a detailed example.*

*Keywords: Algorithms, Design, Experimentation, Human Factors, Model, Method.*

## I. INTRODUCTION

The aim of this project is to review the current state of the art in content-based image retrieval (CBIR), a technique for retrieving images on the basis of automatically-derived features such as color, texture and shape. Our findings are based both on a review of the relevant literature and on discussions with researchers in the field.

The need to find image from a collection a desired is shared by many professional groups, including journalists, design engineers and art historians. While the requirements of image users can vary considerably, it can be useful to characterize image queries into three levels of abstraction: primitive features such as color or shape, logical features such as the identity of objects shown and abstract attributes such as the significance of the scenes depicted. While CBIR systems currently operate effectively only at the lowest of these levels, most users demand higher levels of retrieval.

Content-based image retrieval (CBIR), also known as query by image content (QBIC) and content-based visual information retrieval (CBVIR) is the application of computer vision to the image retrieval problem, that is, the problem of searching for digital images in large databases. "Content-based" means that the search will analyze the actual contents of the image. The term 'content' in this context might refer colors, shapes, textures, or any other information that can be derived form the image itself. Without the ability to examine image content, searches must rely on metadata such as captions or keywords. Such metadata must be generated by a human and stored alongside each image in the database Problems with traditional methods of image indexing  have led to the rise of interest in techniques for retrieving images on the basis of automatically-derived features such as color, texture and shape – a technology now generally referred to as Content-Based Image Retrieval (CBIR). However, the technology still lacks maturity, and is not yet being used on a significant scale. In the concepts which are presently used for CBIR system are all under research.

Researcher at content base image retrieval has gain tremendous momentum   over last decades. A lot of worked carried out on image retrieval by researcher expending on depth or breadth [1]-[5].  The term Content Based Image Retrieval (CBIR) seems to have originated with the work of Kato [6] for the automatic retrieval of the images from a database, based on the color and

shape present. Since then, the term has widely been used to describe the process of retrieving desired images from a large collection of database, on the basis of syntactical image features (color, texture and shape).   The techniques, tools and algorithms that are used, originate from the fields, such as statistics, pattern recognition, signal processing, data mining and computer vision. CBIR is the most important and effective image retrieval method and widely studied in both academia and industry arena. In this we propose an image retrieval system, called Wavelet-Based Color Histogram Image Retrieval (WBCHIR), based on the combination of color and texture features.  The color histogram for color feature and wavelet representation for texture and location information of an image. This reduces the processing time for retrieval of an image with more promising representatives. The extraction of color features from digital images depends on an understanding of the theory of color and the representation of color in digital images. Color spaces are an important component for relating color to its representation in digital form. Absence of hard evidence on the effectiveness of CBIR techniques in practice, opinion is still sharply divided about their usefulness in handling real-life queries in large and diverse image collections.
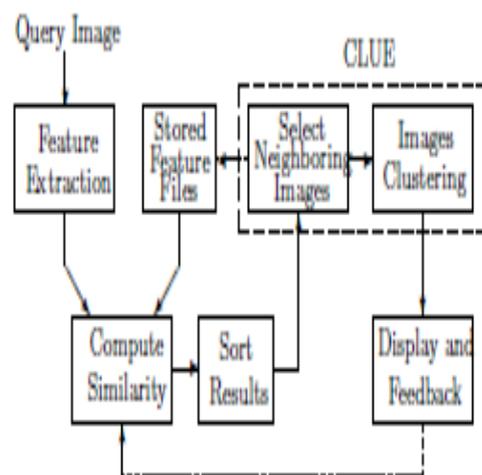
## II. PREVIOUS WORK

In the past decade, many general-purpose image retrieval systems have been developed. Examples include QBIC System [6], Photobook System [16], Blobworld System [3], Virage System [9], VisualSeek and WebSeek Systems [20], the PicHunter System [5], Netra System [14], MARS System [15], and simplicity Systems [22]. A typical CBIR system views the query image and images in the database (target images) as a collection of features, and ranks the relevance between the query image and any target images in proportion to feature similarities. Nonetheless, the meaning of an image is rarely self-evident. Images with high feature similarities to the query image may be very different from the query in terms of the interpretation made by a user (*user semantics* or, in short, *semantics*). This is referred to as the *semantic gap*, which reflects the discrepancy between the relatively limited descriptive power of low level imagery features and the richness of user semantics. Depending on the degree of user involvement in the retrieval process, generally, two classes of approaches have been proposed to reduce the semantic gap: relevance feedback and image database preprocessing using statistical classification.

A relevance-feedback-based approach allows a user to interact with the retrieval algorithm by providing the information of which images he or she thinks are relevant to the query [5, 17]. Based on the user feedbacks, the model of similarity measure is dynamically updated to give a better approximation of the perception subjectivity. Empirical results demonstrate the effectiveness of relevance feedback for certain applications. Nonetheless such a system may add burden to a user especially when more information is required than just Boolean feedback (relevant or non-relevant). Statistical classification methods group images into semantically meaningful categories using low level visual features so that semantically-adaptive searching methods applicable to each category can be applied [18, 21, 22, 12]. For example, Sem Query system [18] categorizes images into different set of clusters based on their heterogeneous features. Vailaya et al. [21] organize vacation images into a hierarchical structure. At the top level, images are classified as indoor or outdoor. Outdoor images are then classified as city or landscape that is further divided into sunset, forest, and mountain classes. The simplicity system [22] classifies images into graph, textured photograph, or non-textured photograph, and thus narrows down the searching space in a database. ALIP system [12] uses categorized images to train hundreds of two-dimensional multiresolution hidden Markov models each corresponding to a semantic category. Although these classification methods are successful in their specific domains of application, the simple ontology built upon them could not incorporate the rich semantics of a sizable image database. There has been work on attaching words to images by associating the regions of an image with object names based on region-term co-occurrence [2]. But as noted by the authors in [2], the algorithm relies on semantically meaningful segmentation. And semantically precise image segmentation by an algorithm is still an open problem in computer vision [19, 23].

## III. MOTIVATION

Apart from presenting information independently, each medium can also express a different characteristic of the same event, so that the media taken together describe the existence, development and result of an event in its entirety. So, there have to be features of information, attributes and the relationships in multimedia data sets that are not within our intuitive grasp. Multimedia data mining involves intelligent data analysis, aimed at finding these features, attributes and relationships in order to construct models for making decisions, taking countermeasures and achieving fusion analysis. "Based on the data stored in them, multimedia databases are used in content-based image retrieval, sound delivery system, video on demand system, World Wide Web and identifying the password command voice based user interface, etc. Multimedia Mining focuses on the following five fronts: Image Mining, Video Mining, Audio Mining, Web Mining and Multi-Media Integrated Mining. Image mining involves the introduction of data mining technology into the image field of study, to discover the information and knowledge hidden in a large quantity of image data. It is the process of identifying hidden, valid, novel, potentially useful, and ultimately understandable semantics of information and knowledge from extensive image data. Content Base image retrieval system is a technique in which instead of firing a query user has to fire an image and get the respective information which are the nearest match of the image which is being fired. As we retrieve an image from the database which depends on some featute like colour histogram, texture and image density. When the user enter feature define above then image search engine check nearest match and show the result.



## IV. LITERATURE REVIEW

CBIR systems are based on a free hand sketch (Sketch based image retrieval – SBIR). With the help of the existing methods, describe possible solution how to design and implement a task spastic descriptor, which can handle the informational gap between a sketch and a colored image, making an opportunity for the efficient search hereby. The used descriptor is constructed after such special sequence of preprocessing steps that the transformed full color image and the sketch can be compared. We have studied EHD, HOG and SIFT. Experimental results on two sample databases showed good results. Overall, the results show that the sketch based system allows users an intuitive access to search-tools [13-16]. The retrieve using sketches in frequently our purpose is to develop a content based image retrieval system, which can use databases. The user has a drawing area where he can draw those sketches, which are the base of the retrieval method. Using a sketch based system can be very important and efficient in many areas of the life. In some cases we can recall our minds with the help of gores or drawing. In these systems the user draws color sketches and blobs on the drawing area. The images were divided into grids, and the color and texture features were determined in these grids. The applications of grids were also used in other algorithms, for example in the edge histogram descriptor (EHD) method [13]. The disadvantage of these methods is that they are not invariant opposite

rotation, scaling and translation. Lately the development of difficult and robust descriptors was emphasized. Another research approach is the application of fuzzy logic or neural networks. In these cases the purpose of the investment is the determination of suitable weights of image features. The CBIR technology can be used in several applications such as digital libraries, crime prevention, photo sharing sites. Such a system has great value in apprehending suspects and indentifying victims in forensics and law enforcement. A possible application is matching a forensic sketch to a gallery of mug shot images. The area of retrieve images based on the visual content of the query picture intensive recently, which demands on the quite wide methodology spectrum on the area of the image processing.

### 4.1 INDEXING:

Indexing the whole set of images using K-means Clustering algorithm. Indexing is done using an implementation of the Document Builder Interface. A simple approach is to use the Document Builder Factory, which creates Document Builder instances for all available features as well as popular combinations of features (e.g. all JPEG features or all avail-able features).In a content based image retrieval system, target images are sorted by feature similarities with respect to the query (CBIR).In this indexing, we propose to use K-means clustering for the classification of feature set obtained from the histogram. Histogram provides a set of features for proposed for Content Based Image Retrieval (CBIR). Hence histogram method further refines the histogram by splitting the pixels in a given bucket into several classes. Here we compute the similarity for 8 bins and similarity for 16 bins. Standard histograms, because of their efficiency and insensitivity to small changes, are widely used for content based image retrieval. But the main disadvantage of histograms is that many images of different appearances can have similar histograms because histograms provide coarse characterization of an image [17]. The k-means algorithm takes the input parameter, k, and partitions a set of n objects into k clusters so that the resulting intra cluster similarity is high but the inter cluster similarity is low. Cluster similarity is measured in regard to the mean value of the objects in a cluster, which can be viewed as the cluster's center of gravity. "How does the k-means algorithm work". The k-means algorithm proceeds as follows. First, it randomly selects k of the objects, each of which initially represents a cluster mean or center. For each of the remaining objects, an object is assigned to the cluster to which it is the most similar, based on the distance between the object and the cluster mean. It then computes the new mean for each cluster [15].

### 4.1.1 THE K-MEANS ALGORITHM:

Algorithm: k-means. The k-means algorithm for partitioning based on the mean value of the objects in the cluster.

Input: The number of clusters k and a database containing n objects.

Output: A set of k clusters that minimizes the squared-error criterion.

Method:

(1) Arbitrarily choose k objects as the initial cluster centers:

(2) Repeat

(3) (Re) assign each object to the cluster to which the object is the most similar, based on the mean value of the objects in the cluster;

(4) Update the cluster means, i.e., calculate the mean value of the objects for each cluster;

(5) Until no change.

## 4.2 Searching:

This module performs intensive searching of images from the database. User gives his query image and based on various algorithms a set of images are generated. The query is taken up and according to the algorithm given a set of related images is generated. This module searches for the images in the database according to the specified algorithm. To generate more relevant images, the number of search images can be decreased. Now if the user is not satisfied by the images generated, he/she can perform the search test again and again, until it generates

### 4.2.1 Types of   Technique Access Image:

For the purpose of simplifying the explanations, we call a CBIR system using CLUE a Content-Based Image Clusters Retrieval (CBICR) system. From a data-flow view point, a general CBICR system can be characterized by a diagram in Figure. The retrieval process starts with feature extraction. The features for target images (images in the database) are usually computed are forehand and stored as feature files. Using these features together with a Neighboring images are selected by nearest neighbors method (NNM). It first chooses k nearest neighbors of the query image I as seeds. The nearest neighbors for each seed are then found. Finally, the neighboring images are selected to be all the distinct images among seeds and their r nearest neighbors, i.e., distinct images in k(r + 1) target images. Data representation is typically the first step to solve any clustering problem. In the field of computer vision, two types of representations are widely used. One is called the *geometric representation*, in which data items are mapped to some real normed vector space. The other is the *graph representation*. It emphasizes the pair wise relationship, but is usually short of geometric interpretation. When working with images, the geometric representation has a major limitation: it requires that the images be mapped to points in some real normed vector space.

### First method:

 It can retrieve an image not already existing in the database it support online image addition.

### Main and Similarity methods:

 Main and Similarity method only permit query image already in the database, not a new image. If using clustering approach, it first compute distance to the initial set, for the relevant ones, then go further into clusters to compute distances, then rank the whole results. Now it only implements to first level. It doesn't work from Biren's retrieval interface. Biren wrote some test programs to perform the experiments. Those test programs are not connected to the actual system... files with names CbirTest1, 2, 3 are the test programs, and the results in them are not direct as far as the main and similarity method retrieval is concerned. The results of those test programs are formatted and listed in a manner to report them in his thesis, and those results can be used to derive the main and similarity method retrieval results. In his "ImageRetrievalProcess.java", there is some code for main method also along with first method, but it does not work. He tested but may not have removed that part of the code from the program. For the first level retrieval, I rewrite the codein"ImageRetrievalProcessForMainAndNew.java". The real distance is the same for all methods. The estimated distance computation is same for all methods too. Only the high level feature vectors for the three methods are different.

### 4.2.2 Image Database systems:

Set of images are collected, analyzed and stored in multimedia information systems, office systems, Geographical information systems(GIS), robotics systems , CAD/CAM systems, earth resources systems,  medical databases, virtual reality systems, information retrieval systems, art gallery and museum catalogues, animal and plant atlases, sky star maps, meteorological maps, catalogues in shops and many other places.  There are sets of international organizations dealing with

different aspects of image storage, analysis and retrieval. Some of them are: AIA (Automated Imaging/Machine vision), AIIM (Document imaging), ASPRES (Remote Sensing/ Protogram) etc.

There are also many international centers storing images such as : Advanced imaging, Scientific/Industrial Imaging, Microscopy imaging, Industrial Imaging etc. There are also different international work groups working in the field of image compression, TV images, office documents, medical images, industrial images, multimedia images, graphical images, etc.

### 4.2.3 Logical Image Representation in Database Systems:

The logical image representation in image databases systems is based on different image data models. An image object is either an entire image or some other meaningful portion (consisting of a union of one or more disjoint regions) of an image. The logical image description includes: meta, semantic, color, texture, shape, and spatial attributes. color attributes could be represented as a histogram of intensity of the pixel colors. A histogram refinement technique is also used by partitioning histogram bins based on the spatial coherence of pixels. Statistical methods are also proposed to index an image by color correlograms, which is actually a table containing color pair's archival material, tape recordings and information files. A number of less widely-known schemes have been devised to classify images and drawings for specialist purposes. Examples include the Vienna classification for trademark images [World Intellectual Property Organization, 1998], used by registries Worldwide to identify potentially conflicting trademark applications, and the Opitz coding system for machined parts [Opitz et al, 1969], used to identify families of similar parts which can be manufactured together.

### References

1.  R. Datta, D. Joshi, J. Li and J. Z. Wang, "Image retrieval: Ideas, influences, and trends of the new age", ACM computing Survey, vol.40, no.2, pp.1-60, 2008.

2.  J. Eakins and M. Graham, "Content-Based Image Retrieval", Technical report, JISC Technology Applications Programme, 1999.

3.  Y. Rui, T. S. Huang and S.F. Chang, "Image Retrieval: Current Techniques, Promising Directions and Open Issues. Journal of Visual Communication and Image Representation. 10(4): pp. 39-62. 1999.

4.  A. M. Smeulders, M. Worring and S. Santini, A. Gupta and R. Jain, "Content Based Image Retrieval at the End of the Early Years", IEEE Transactions on Pattern Analysis and Machine Intelligence, 22(12): pp. 1349-1380, 2000.

5.  Y. Liu, D. Zang, G. Lu and W. Y. Ma, "A survey of content-based image retrieval with high-level semantics", Pattern Recognition, Vol-40, pp-262-282, 2007.

6.  A. Natsev, R. Rastogi and K. Shim, "WALRUS: A Similarity Retrieval Algorithm for Image Databases", In Proceeding. ACM SIGMOD Int. Conf. Management of Data, pp-395–406, 1999.

7.  S. Ardizzoni, I. Bartolini, and M. Patella, "Windsurf: Region based Image Retrieval using Wavelets", In IWOSS'99, pp. 167–173, 1999.

8.  G. V. D. Wouwer, P. Scheunders and D. V. Dyck, "Statistical texture characterization from discrete wavelet representation", IEEE Transactions on Image Processing, Vol.8, pp-592–598, 1999.

9.  S. Livens, P. Scheunders, G. V. D. Wouwer and D. V. Dyck, "Wavelets for texture analysis, an overview", Proceedings of Sixth International Conference on Image Processing and Its Applications, Vol. 2, pp-581–585, 1997.

10. R. C. Gonzalez and E.W. Richard, Digital Image Processing, Prentice Hall. 2001.

11. Yang Changchun, Yi Li, "A Data Mining Model and Methods Based on Multimedia Database," Internet Technology and Applications, 2010 IEEE International Conference on, vol., no., pp.1-4, 20-22 Aug. 2010.

12. S. Andrews, T. Hofmann, and I. Tsochantaridis. Multipleinstance learning with generalized support vector machines.Artificial Intelligence, pages 943–944, 2002

13. M. Eitz, K. Hildebrand, T. Boubekeur, and M. Alexa, "An evaluation of descriptors for large-scale image retrieval from sketched feature lines," Computers and Graphics, vol. 34, pp. 482– 498, October 2010

14. R.Hu, M. Barnard, and J. Collomosse, "Gradient _eld descriptor for sketch based image retrieval and localization," International Conference on Image Processing, pp. 1–4, 2010.

15. A.K. Jain, J.E. Lee, and R. Jin, "Sketch to photo matching: a feature-based approach," Proc. SPIE, Biometric Technology for Human Identi_cation VII, vol. 7667, pp. 766702–766702, 2010.

16. T. Hashimoto, A. R¨ovid, G. Ohashi, Y. Ogura, H. Nakahara, and A.R. V´arkonyi-K´oczy, "Edge detection based image retrieval method by sketches," Proc. of the Internatonal Symposium on Flexible Automation, pp. 1–4, 2006.

*Junaid  et al.,*

*International Journal of Advance Research in Computer Science and Management Studies*
*Volume 2, Issue 1, January 2014  pg. 708-714*

17. R. Fabbri, L.D.F. Costa, J.C. Torelli, and O.M. Bruno, "2D Euclidean distance transform algorithms: a comparative survey," ACM Computing Surveys, vol. 44, pp. 1–44, February 2008.

18. G. Sheikholeslami, W. Chang, and A. Zhang, "SemQuery: Semantic Clustering and Querying on Heterogeneous Features for Visual Data," IEEE Trans. Knowledge and Data Engineering, vol. 14, no. 5, pp. 988–1002, 2002.

19. J. Shi and J. Malik, "Normalized Cuts and Image Segmentation," IEEE Trans. Pattern Anal. Machine Intell., vol. 22, no. 8, pp. 888–905, 2000.

20. J. R. Smith and S.-F. Chang, "VisualSEEK: A Fully Automated Content-Based Query System," Proc. 4th ACM Int'l Conf. on Multimedia, pp. 87–98, 1996.

21. A. Vailaya, M. A. T. Figueiredo, A. K. Jain, and H.-J. Zhang, "Image Classification for Content-Based Indexing," IEEE Trans. Image Processing, vol. 10, no. 1, pp. 117–130, 2001.

22. J. Z. Wang, J. Li, and G. Wiederhold, "SIMPLIcity: Semantics-Sensitive Integrated Matching for Picture LIbraries," IEEE Trans. Pattern Anal. Machine Intell., vol. 23, no. 9, pp. 947–963, 2001.