# Parallel Image Searching and Retrieval

## Putul Siddharth[1], Rahul Mahajan[2]

*[1] VIT University, Chennai, India*

*E-mail: putul.siddharth2017@vitstudent.ac.in,*

*rahul.mahajan2017@vitstudent.ac.in*

## Abstract –

In this project we are going to take a set of images as in input form and extracting the features vectors in a parallel fashion using multiple threads.

Then we export these resultant feature vectors of all images to an excel sheet.
Now a query image is taken and in order to obtain the similar image, its feature vector is matched with the feature vector present in excel sheet which in return show the required result.

This is similar to the one we do in the web browser where we input our elements and the browser searches for the similar one in his record and present the same.

We are going to implement this project using MATLAB. As it has emerged as one of the powerful language used for technical computing. Parallel computing with MATLAB has an interactive environment which provide high performance computing procedures and helps to present the output in an efficient and designed way.

## Literature Review –

**1. A parallel content-based image retrieval system using spark and tachyon frameworks
Journal of King Saud University - Computer and Information Sciences**

**Explanation –** The first proposed approaches are centralized CBIR systems, such as QBIC (Flickner et al., 1995), Virage (Bach et al., 2670), MARS (Rui et al., 1997), TinEye[1],etc. However, with the increasing generation of image data being produced every day, the conventional systems of CBIR become so greedy in time consumption. This drawback has motivated the researchers to increase the performance of CBIR systems using a set of solutions such as: parallel interactive retrieval system based on semantic and visual descriptor (Hong et al., 2017), image retrieval optimization using meta-heuristic algorithms (Abdel-Basset et al.,

2018, Abdel-Basset et al., 2018) or neutrosophic sets (Abdel-Basset et al., 2018, Abdel-Basset et al., 2018, Abdel-Basset et al., 2018) and parallel retrieval system for large-scale images using cloud computing (Yang et al., 2018). We present here some of the proposed parallel and distributed systems. Lu et al. (2007) used parallel computing techniques to execute similarity comparison and feature extraction of visual features rooted in a cluster architecture. The experiments indicate that the use of parallel computing techniques can enhance considerably the speed of CBIR systems. In Zhang et al. (2010), Zhang et al. present an implementation for a Distributed Image Retrieval System (DIRS) deployed on Hadoop distributed computing environment. In the indexing step of DIRS system, the authors adopt distributed database HBase as a storage layer. In addition, they used MapReduce computing model of Hadoop to improve the performance of searching step in large image data bases. Gu and Gao (2012) proposed a CBIR system upon Hadoop, HBase and Lucene frameworks. Raju et al. (2015) developed a CBIR system based on Hadoop MapReduce. They used the local-tetra patterns (LTPr) (Murala et al., 2012) as a feature vector technique to represent the images in the CBIR system. Costantini and Nicolussi (2015) proposed a novel Hadoop and spark based system of image retrieval for large scale image. In their work, the authors use Hadoop in the indexing phase and Spark in the searching step. Sakr et al. (2016) used an efficient Hadoop MapReduce technique to search the closest images for the query image. Duan et al. (2016) proposed an improved CBIR system based on Apache Spark. They used Avro framework to combine the image files, where an Avro file can persist in memory in order to allow future actions to be much faster. Lately, Lagiewka et al. (2017) introduced a distributed CBIR system in relational databases. Their system is composed of a set of machines, each of them uses the Apache Hadoop software framework with HDFS. Beside the parallel and distributed systems cited above, many other methods have been attempted by researches to optimize the k-NN computation method, especially in big data. In Dong et al. (2011), Dong et al. present NN-Descent method, which is a simple and efficient K-Nearest Neighbor Graph (K-NNG) construction with arbitrary similarity measures for large-scale applications, where data structures need to be distributed over the network. Song et al. (2015) introduced a distributed k-NN method for massive data using MapReduce programming model. Maillo et al. (2015) proposed a parallel k-NN algorithm based on MapReduce programming model for large scale data classification denoted as MR-kNN. The authors in Ding and Boykin (2017) have developed a custom K-Nearest Neighbor algorithm using Hadoop for deploying applications in a parallel way on a cluster.

## 2. Content Based Image Retrieval: A Review by Piyush A. Dahake1, S. S. Thakare2

**1(Electronics and Telecommunication Department, GCOE, Amravati, India)**

**2(Assistant Professor, Electronics and Telecommunication Department, GCOE, Amravati, India)**

**Explanation –**An approach is proposed for retrieval based on combination of color, texture and edge features of image. Performance evaluation of studied image retrieval techniques and proposed technique is done using parameters like sensitivity, Specificity, Retrieval score, Error rate and accuracy[1]. In this image retrieval system extraction is based on the averaging method clustering image, revised averaging algorithm to reduce the complexity of extraction and efficiency[2]. Gabor wavelet transform is mostly combining of features of image and the Gabor Wavelet Transform is degrade into distinct scaling and orientation with various of filters to minimize the unwanted information of the images[3]. In this methodology only the color feature get extracted from image and at first image is divided into 16 equal sized segment after that the average value of each color component is considered into account[4]. Rather than the transform and averaging techniques an unsupervised learning technique is also used i.e. First a Self Organizing Map (SOM) and then Latent Dirichlet Allocation (LDA)[5].

## 3. Image Retrieval Through Sketches Based on Descriptor by Dipika Birari1, Kushal Mandge2 (Assistant Professor, Computer Engineering, PGMCOE, Pune, Maharashtra)1 (Assistant Professor, E&TC, DYPSOE, Pune, Maharashtra)2

**Explanation –** Eitz et al.[4] performed random sampling on images and then proposed the SHoG descriptor to describe each sampling point. Only the gradient value near the most dominant edge line is retained in SHoG. Hu and Collomosse[5] introduced dense gradient field on which they computed a multi-scale HOG feature (GF-HOG). GF-HOG is also utilized to describe regions which are generated by hierarchical image segmentation. Bozas [2] divided an image into overlapping patches and computed a HOG feature for each patch. In addition to HOG-based descriptors, Eitz utilized shape context to perform retrieval. Contour consistency filtering based on shape context descriptor was performed by Chen et al. Chalechale et al. performed angular partitioning on the edge image. Fourier transform was applied to achieve rotational invariance. Eitz [10] proposed a descriptor known as structure tensor, which was designed to find a single vector that is closest to the parallel direction of the majority of the edges in a local region. The MinHash method was used to build an index structure. In addition to HOG-based descriptors, Visual Saliency Weighting (VSW) was employed by Furuya and Ohbuchi [9] to emphasize the object of interest. Saavedra and Bustos represented sketches by six types of key shape. Zhou et al. extracted multi-scale features on candidate regions and built a hierarchical index structure to achieve coarse-to-fine retrieval.

## 4. CLASSIFICATION OF BIOMEDICAL IMAGES USING CONTENT BASED IMAGE RETRIEVAL SYSTEMS by Yinghui Zhang 1, Fengyuan Zhang 2, Yantong Cui 3, Ruoci Ning 4

1 University of California, Berkeley, USA 2 Northeastern University, China
3 Anshan No.1 High School, China 4 Marian High School, USA

**Explanation –**According to Rani et al. [6], The main objective of their work is to use Support Vector Machine in a more efficient manner so that results can be produced in lesser execution time. In their proposed work, images are pre-processed before they are stored in the database. This process helps to enhance the quality of the images by removal of the noise. Then RGB model is used for clustering the images. According to Chaudhary et al. [7] , In their proposed work, integrated approach is used to extract color and texture feature from images. The help of the multi featured extraction is taken so that efficient image retrieval process can be performed. A higher order of color moment is used for the extraction of color features. Texture extraction and face recognition are done using local binary pattern (LBP). According to Trojacanec et al. [8], In their proposed work, an improvement of the two-level CBIR architecture is performed. In the first level, the clinician's voting is included, and in the second level, inclusion to the previously computed agents' voting is performed. The main motive of these two levels is to increase the efficiency of the proposed work and accuracy of the retrieval process. According to Antaniet al. [9], Their research work is focused on developing techniques for hybrid text/image retrieval from the survey of text and image data. Their research also shows various challenges that come into existence, when a particular CBIR system is developed especially for the biomedical images. According to Ramos et al. [12], In their proposed work in radiology, reports are done which supervise the CBIR. Inferring the relationship between the patients and subsequently applying them to supervise the metric learning algorithms are the main phases of the proposed work. Their research work calculates the text distances between exam reports in the exam image space to supervise a metric learning algorithm. Using this method, there is a consistent increase in the CBIR performance.

## 5. An Empirical Study on Video Retrieval and its Impacts in Contemporary World of Image Processing

**Explanation –** Yuxin Peng et. al [1] developed a new method for detection of sizzling incident and overview of information videos. This method is mostly created by considering two graph algorithms as its base: normalized cut (NC) and optimal matching (OM). First, OM used to calculate the visual relationship stuck among all pair of actions below the one-to-one identify restriction between the video shots. Then, the actions of the news are displayed as a full slanted graph and normalize Cut is approved out to divider the chart into incident clusters globally. The planned come up to has been tested on information videos of 10 hours and has been establish to be successful. 2. Alan F. Smeaton [2] stated that a short analysis of the nature of the video analysis, retrieval and indexing. It includes the fact to research directions, to consider which the capable structure for the process of searching is and browsing of video records based on the content of the video, so easy as surfing the (text) web pages. 3. Yuk Ying Chung et.al [3] developed and executed a method called content based video retrieval system by considering the following methods, D4 Daubechies wavelet transform, Haar wavelet, and five various kinds of clustering techniques. The experimental output reveals that the Haar wavelet with 3- Level

transform ensures the enhanced result, which has the accuracy rate of retrieval (89%). 4. Ritendra Dutta et. al [4] proposed that nearly 300 key hypothetical and experimental contributions in the present decade corresponding to the image retrieval and automatic image annotation, and in the process includes the study of the content of the related subfields. They also conferred the important issues involved in the adaptation of past techniques of image retrieval to create the systems that can be suitable in the real world.

## Proposed Work –

In this project we are going to retrieve the images from our database on the basis of given query images parallel. For this, we divided our work into three modules:-
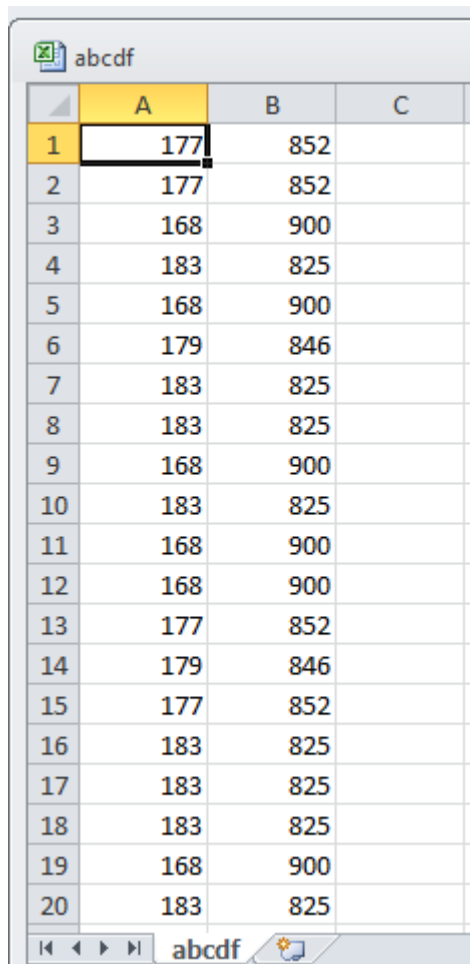
**Module1:- Create the database of images**

At first create the database of images by downloading from our respective browser. We use two different database to check time complexity of image retrieval.

## Module 2:- Read the images from our database

After creating the database, it is time to read images one -by-one from their respective directory. Additionally, we are storing the feature vector namely rows and columns of each images and exporting it to the excel sheet. We take any image from our database to be chosen for our query image.

| abcdf | | |
|---|---|---|
| | A | B | C |
| 1 | 177 | 852 | |
| 2 | 177 | 852 | |
| 3 | 168 | 900 | |
| 4 | 183 | 825 | |
| 5 | 168 | 900 | |
| 6 | 179 | 846 | |
| 7 | 183 | 825 | |
| 8 | 183 | 825 | |
| 9 | 168 | 900 | |
| 10 | 183 | 825 | |
| 11 | 168 | 900 | |
| 12 | 168 | 900 | |
| 13 | 177 | 852 | |
| 14 | 179 | 846 | |
| 15 | 177 | 852 | |
| 16 | 183 | 825 | |
| 17 | 183 | 825 | |
| 18 | 183 | 825 | |
| 19 | 168 | 900 | |
| 20 | 183 | 825 | |

abcdf

## Module 3:- Get the distance of each of the images from their respective query image

Now we need to find the distance of each image present in our database with the query image. As in image processing field, similarity is defined on the basis of metrics- shorter the distance, more similar is the image. We use Euclidean distance as our distance metric tool as it is most widely used. Despite of that we also develop our new distance metric which is used for calculating the distances.

New Distance Metric:

$dist1\ (i) = sinh\ (sum\ (abs\ (row\_q - xx(i)) + abs\ (col\_q - yy\ (i))))$

**Module 4:- Retrieve the top 6 similar image**

After finding the distances, we need to sort them. As in MATLAB many inbuilt library are there for sorting purpose but if we sort the distances directly, we won't be able to know that these sorted distances belong to which images. That's why we use bubble sort to sort distance along with their feature vector so that it became easy for us know which images have the sorted distance compared to query image.

**Module 5:- Find the Accuracy, Precision, recall, f1-measure for the obtained result**

When we get the retrieved image, we find the Accuracy, Precision, recall, f1-measure for the obtained result using their corresponding formula.

|  |  | Predicted | |
|---|---|---|---|
|  |  | **Negative** | **Positive** |
| **Actual** | **Negative** | True Negative | False Positive |
|  | **Positive** | False Negative | True Positive |

$$\text{Precision} = \frac{True\ Positive}{True\ Positive + False\ Positive}$$

$$\text{Recall} = \frac{True\ Positive}{True\ Positive + False\ Negative}$$

$$\text{F1} = 2 \times \frac{Precision * Recall}{Precision + Recall}$$

# Result and Discussion –

We have taken set of images as in input form (**query image**) and extracting the features vectors which output (**retrieved/similar images**) in a parallel fashion using multiple threads.

## Query Image –



## Retrieved Images –

Also we have calculated and displayed dimension of images retrieved and its distance from query image. Further we have calculated various performance metrics such as **Recall**, **Accuracy** and **F1 Score** and **Time Taken**.



# References –

**1.** https://www.researchgate.net/publication/258431660_Parallel_Content_Based_Sub-Image_Retrieval_Using_Hierarchical_Searching (Parallel Content Based Sub-Image Retrieval Using Hierarchical Searching research paper).

**2.** YouTube

**3.** Wikipedia