



**IJIRCCCE**

e-ISSN: 2320-9801 | p-ISSN: 2320-9798



# INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN COMPUTER & COMMUNICATION ENGINEERING

**Volume 9, Issue 11, November 2021**

**ISSN** INTERNATIONAL  
STANDARD  
SERIAL  
NUMBER  
INDIA

**Impact Factor: 7.542**

 9940 572 462

 6381 907 438

 [ijircce@gmail.com](mailto:ijircce@gmail.com)

 [www.ijircce.com](http://www.ijircce.com)

# Performance Evaluation of Vehicle Classification using Colour and Deep Features

Mutyala Anusha<sup>1</sup>, Sarakanam S Manikanta<sup>2</sup>

M.Tech Student, Department of ECE, Chaitanya Institute of Science and Technology Kakinada, India<sup>1</sup>

Assistant Professor, Department of ECE, Chaitanya Institute of Science and Technology Kakinada, India<sup>2</sup>

**ABSTRACT:** Vehicle similarity identification from a large image database is a critical task. The solution to this problem is the use of a vehicle Image Retrieval (VIR) System. The images are described through their content, there is three predominant content existing in an image like color, shape, and texture. In this paper, we are evaluating the performance of the VIR system using two methods. The first method consists of colour and texture features. The second method consists of Use of CNN. The feature extraction technique is achieved based on an input query image from the database and features are saved in a feature dataset. A proposed strategy retrieves similar images from a database that fulfills the user's desire. The similarity measurement can be done using the Euclidean distance and hashing technique. The overall performance of the retrieval system has been analyzed through the parameters accuracy, Precision and Mean Average Precision. The experimental result shows encouraging results using CNN which leads to improving accuracy

**KEYWORDS:** vehicle recognition, classification, CNN, color moments, texture, feature extraction.

## 1. INTRODUCTION

The explosive increase and ubiquitous accessibility of visual data on the Web have led to the prosperity of research activity in image search or retrieval. With the ignorance of visual content as a ranking clue, methods with text search techniques for visual retrieval may suffer inconsistency between the text words and visual content. Content-based image retrieval (VIR), which makes use of the representation of visual content to identify relevant images, has attracted sustained attention in recent two decades. Such a problem is challenging due to the intention gap and the semantic gap problems. Numerous techniques have been developed for content-based image retrieval in the last decade. With the universal popularity of digital devices embedded with cameras and the fast development of Internet technology, billions of people are projected to the Web sharing and browsing photos. The ubiquitous access to both digital photos and the Internet sheds bright light on many emerging applications based on image search. Image search aims to retrieve relevant visual documents to a textual or visual query efficiently from a large-scale visual corpus.

Although image search has been extensively explored since the early 1990s [1], it still attracts lots of attention from the multimedia and computer vision communities in the past decade, thanks to the attention on scalability challenge and emergence of new techniques. Traditional image search engines usually index multimedia visual data based on the surrounding meta data information around images on the Web, such as titles and tags. Since textual information may be inconsistent with the visual content, content-based image retrieval (VIR) is preferred and has been witnessed to make great advance in recent years. In content-based visual retrieval, there are two fundamental challenges, i.e., intention gap and semantic gap. The intention gap refers to the difficulty that a user suffers to precisely express the expected visual content by a query at hand, such as an example image or a sketch map. The semantic gap originates from the difficulty in describing high-level semantic concept with low-level visual feature [2] [3] [4].

To narrow those gaps, extensive efforts have been made from both the academia and industry. From the early 1990s to the early 2000s, there have been extensive study on content-based image search. The progress in those years has been comprehensively discussed in existing survey papers [5] [6] [7]. Around the early 2000s, the introduction of some new insights and methods triggers another research trend in VIR. Specially, two pioneering works have paved the way to the significant advance in content-based visual retrieval on large-scale multimedia database. The first one is the introduction of invariant local visual feature SIFT [8]. SIFT is demonstrated with excellent descriptive and discriminative power to capture visual content in a variety of literature. It can well capture the invariance to rotation and scaling transformation and is robust to illumination change. The second work is the introduction of the Bag-of-Visual-Words (BoW) model [9]. Leveraged from information retrieval, the BoW model makes a

compact representation of images based on the quantization of the contained local features and is readily adapted to the classic inverted file indexing structure for scalable image retrieval.

Based on the above pioneering works, the last decade has witnessed the emergence of numerous works on multimedia content-based image retrieval [10] [11] [12] [13] [9] [14] [15] [16] [17] [18] [19] [20] [21] [22] [23] [24] [25] [26] [27] [28] [29]. Meanwhile, in industry, some commercial engines on content-based image search have been launched with different focuses, such as TinEye1, Ditto2, Snap Fashion3, ViSenze4, Cortica5, etc. TinEye is launched as a billion-scale reverse image search engine in May, 2008. Until January of 2017, the indexed image database size in TinEye has reached up to 17 billion. Different from TinEye, Ditto is specially focused on brand images in the wild. It provides an access to uncover the brands inside the shared photos on the public social media web sites.

**Color Features:** Images are largely categorized into grayscale images and color images. In a grayscale image color pixel having a solely grayscale area while in a color image three color intensity ranges are used. In the color image red, green and blue intensities are used. Color histogram, color coherence, and color moments are important methods used for image retrieval

**Texture Features:** It measures the homogeneity of a pixel over repeated patterns in the image. We can form a retrieval system the use of two tactics particularly structural and frequency-based approaches. Shape Features: It gives edges or outlines of an object existing in an image. Region and boundary-based techniques are used in the retrieval systems based totally on shape features.

**Neural Network:** A neural network consists of the input layer, hidden layer, and output layer. Convolution Neural Network is used for feature extraction from images

## II. LITERATURE SURVEY

Even though Multimedia databases (MMD) is among the fastest growing emerging technologies in the field of database systems. New technologies pose numerous challenges, and MMD has its share of challenges. Most of MMD challenges are around Content-based Image Retrieval (VIR) systems. VIR is a technique for retrieving images on the basis of automatically-derived features such as color, texture and shape. Moreover, multimedia objects contain encoding of raw sensorial data, which compromise the efficient indexing and retrieval. As a result of which, Query by Image Content (QBIC) technique using image descriptors for indexing and retrieval of multimedia objects were proposed by various studies to address this problem. However, an effective and precise performance evaluation benchmarking for this technique remains exclusive.

Since the invention of the Internet, and the availability of image capturing devices such as smartphones, digital cameras, image scanners and geospatial satellite devices, the size of digital image storage is increasing rapidly. Efficient image searching, browsing and retrieval tools are required by end users from various domains, including remote sensing, fashion design, criminology, publishing, medicine, architecture, etc. It is for these reasons that, many general purpose image retrieval systems have been developed. Therefore, for the same reasons we explore the in-depth survey of content-based image retrieval technology, descriptor technology and performance measure framework technology in order to gain an insight of this domain field.

The main object of a Content-Based Image Retrieval (VIR) system, also known as Query by Image Content (QBIC), is to help users to retrieve relevant images based on their contents. VIR technologies provide a method to find images in large databases by using unique descriptors from a trained image. The image descriptors include texture, color, intensity and shape of the object inside an image. The urgency of efficient image searching, browsing and retrieval techniques by users from large repositories such as the internet, meteorological images and geospatial images is real.

It is reported by [5] that, there are two retrieval frameworks: text-based and content-based. In the text-based approach, the images are manually annotated by text descriptors, which are then used by a database management system to perform image retrieval. There are two disadvantages with this approach. The first is that a human labor at a considerable level is required for manual annotation. The second is the inaccuracy in annotation due to the subjectivity of human perception. To overcome these disadvantages in text-based retrieval system, content-based image retrieval (VIR) was introduced.

It is asserted by [24], that content-based imageretrieval (VIR), also known as query by image content(QBIC) and content-based visual information retrieval(CBVIR), is the application of computer visiontechniques to the image retrieval problem. It is atechnique which uses visual features of image such ascolor, shape, texture, etc. to search user required imagefrom large image database according to user's requestsin the form of a query image. Images are retrieved onthe basis of similarity in features where features of thequery specification are compared with features from theimage database to determine which images matchsimilarly with given features.

It is defined by [18] that, in computer vision, visual descriptors or image descriptors are defined as the descriptions of the visual features of the contents in images, videos, or algorithms or applications that produce such descriptions. They describe elementary characteristics such as the shape, the color, the texture or the motion, among others. It is describe by [28], that visual descriptors are divided in two main groups: General information descriptors, which they contain low level descriptors which give a description about color, shape, regions, textures and motion, and specific domain information descriptors which they give information about objects and events in the scene.

In their book [6] describe the general information descriptors as consisting of a set of descriptors that covers different basic and elementary features like: color, texture, shape, motion, location and others. The color descriptor is the most basic quality of visual content. Five tools are defined to describe color; Dominant Color Descriptor (DCD), Scalable Color Descriptor (SCD), Color Structure Descriptor (CSD), Color Layout Descriptor (CLD), and Group of frame (GoF) or Group-of-pictures (GoP). The Texture descriptors are used to characterize image, textures, or regions. They observe the region homogeneity and the histograms of these region borders. The set of descriptors is formed by: Homogeneous Texture Descriptor (HTD), Texture Browsing Descriptor (TBD), and Edge Histogram Descriptor (EHD). The Shape descriptor contains important semantic information due to human's ability to recognize objects through their shape.

However, this information can only be extracted by means of a segmentation similar to the one that the human visual system implements. These descriptors describe regions, contours and shapes for 2D images and for 3D volumes. The shape descriptors are formed by; Region based Shape Descriptor (RSD), Contour-based Shape Descriptor (CSD) and 3-D Shape Descriptor (3-D SD). While, the Motion descriptors are defined by four different descriptors which describe motion in video sequence. The descriptor set is formed by; Motion Activity Descriptor (MAD), Camera Motion Descriptor (CMD), Motion Trajectory Descriptor (MTD), and Warping and Parametric Motion Descriptor (WMD and PMD). Finally, the Location descriptor element's location in the image is used to describe elements in the spatial domain.

### III. PROPOSED METHOD

The process of Vehicle image retrieval (VIR) system consists of the following six main stages of: image acquisition, image pre-processing, feature extraction, similarity matching, resultant retrieval image and user interface and feedback.

#### 3.1 Image acquisition

It is the process of acquiring a digital image from the image database. The image database consists of the collection of n number of images depends on the user range and choice.

#### 3.2 Image pre-processing

It is the process of improving the image in way that increases the chances for success of the other processes. The image is first processed in order to extract the features, which describe its contents. The processing involves filtering, normalization, segmentation, and object identification. Image segmentation is the process of dividing an image into multiple parts. The output of this stage is a set of significant regions and objects.

#### 3.3 Feature Extraction

It is the process where features such as shape, texture, colour, etc. are used to describe the content of the image. The features further can be classified as low-level and high-level features. In this stage visual information is extracted from the image and saves them as features vectors in a features database. For each pixel, the image description is found in the form of feature value (or a set of values called a feature vector) by using the feature extraction. These feature vectors are used to compare the query with the other images and retrieval.

### 3.4 Similarity Matching

It is a process that entails the information about each image is stored in its feature vectors for computation process and these feature vectors are matched with the feature vectors of query image (the image to be searched in the image database whether the same image is present or not or how many are similar kind images exist or not) which helps in measuring the similarity. This step involves the matching of the above stated features to yield a result that is visually similar with the use of similarity measurement method called as Distance method. There are various distance methods available such as Euclidean distance, City Block Distance, and Canberra Distance.

### 3.5. Resultant Retrieved images

It is the process that searches the previously maintained information to find the matched images from database. The output will be the similar images having same or very closest features as that of the query image.

### 3.6 User interface and feedback

It is the process which governs the display of the outcomes, their ranking, the type of user interaction with possibility of refining the search through some automatic or manual preferences scheme etc. The Figure 1 below demonstrates the VIR System and its various components.

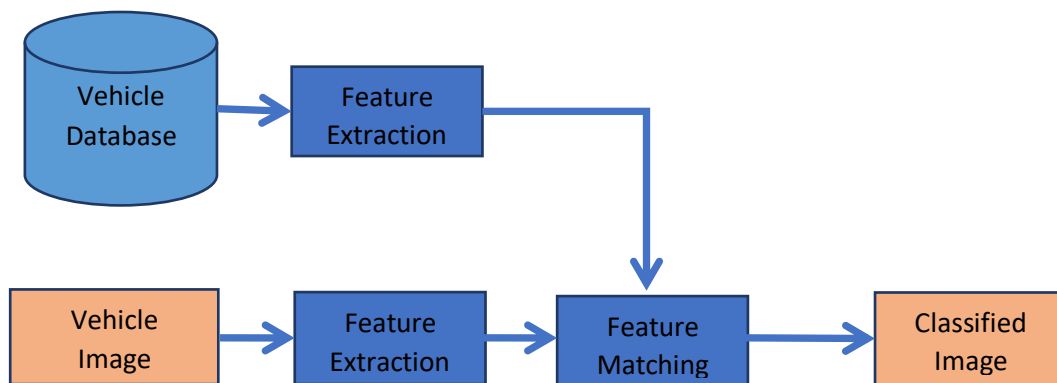


Figure.1 Vehicle Image Retrieval System

## IV. FEATURE EXTRACTION

There has been tremendous work on various ways to deal with the recognition of different kinds of features in images. These features can be classified as follows:

### 4.1 Low Level Features

Features in this category are all application independent, e.g. color, texture, and shape. According to concept level, they can be further divided into:

- Pixel-level features

features be determined at every pixel, for example color, area, and the first and second derivatives of gray-scale values at every pixel. Pixels are extracted and stored in an array. The array contains the RGB components of each pixel. Each pixel in the image is then processed to identify the feature vectors of the image. Edges were used as the only feature vector.

- Local features

The local image description is established on the reason that images can be described by characteristics registered on regions of the image. Can be determined over the consequences of image division and edge detection algorithms. Object shape is an example of such feature [14].

- Global features

The global image descriptor is composed by color and texture features being computed.

- Texture Feature Extraction

The second element of the new system is the texture feature. For this purpose, EHD algorithm is used. Texture is an important feature of expected images. A variety of techniques have been proposed for estimating texture comparability. These strategies ascertain proportions of image texture, for example, the level of differentiation, coarseness, directionality and consistency [8]; or periodicity, directionality and randomness.

4.2. CNN Based Features

This is our proposed frame work for utilizing features from a pretrained deep CNN. We extract features from Pre-trained VGG16 deep CNN model for image retrieval task. A deep CNN model usually consists of many layers that incrementally calculate features. As outlined in Fig. 2. deep CNN model incrementally learns the features through layers of convolutions and subsampling.

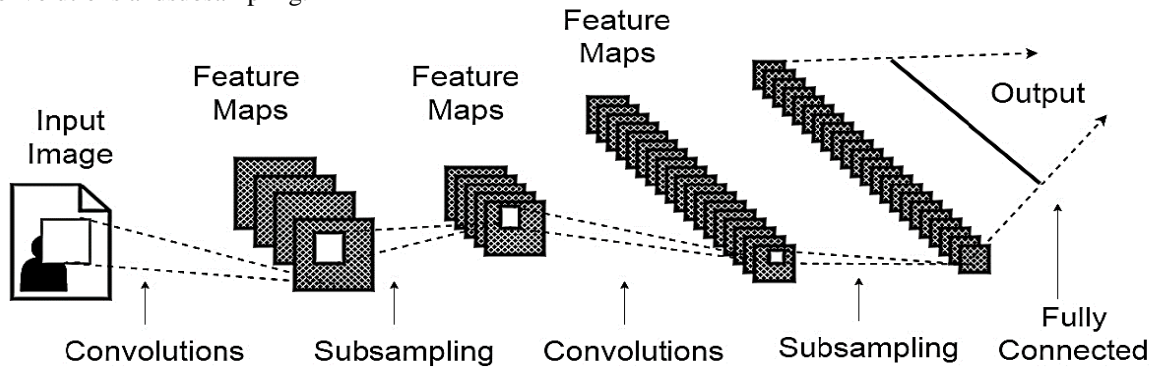


Figure.2 Architecture of Convolution neural networks

In this work, VGG16 deep CNN model is implemented from Python Keras package. It is a 16-layer deep CNN created by the Visual Geometry Group from University of Oxford [29]. VGG16 model is trained on ImageNet, which is a very large-scale dataset containing 3+ million digital images distributed across 5000+ categories. VGG16 model consists of 5 convolution blocks and each convolution block contains two convolution layers (size 3X3) and one maxpooling layer (size 2X2). The final classification step of the model consists of fully connected (FC) layers. Our algorithm extracts 4096 features from fully connected FC2. This is output of second and penultimate fully connected layer of the pre-trained VGG16 CNN model. The feature extraction is done for each image in the dataset and query images.

V. EXPERIMENTAL RESULTS

Our experiment on baseline VIR with handcrafted features (Colour, texture and shape) yields us an average precision of 73.25% across all classes of the weather images dataset. The proposed VIR method which uses pre-trained VGG16 deep CNN features achieves an average precision of 86.73% across all classes of the dataset. The improvement in precision rate is observed across all image classes. Fig. 3. depicts the improvement in precision as recorded across different retrieval sizes. The improvement in precision rate in Clear image class is lower as compared to other three classes (Cloudy, Rain and Sunrise), where precision improvement is profound. Experimental results show that our proposed VIR frame work using features from pre-trained VGG16 CNN model performs better than traditional VIR using handcrafted features (Color, texture and shape). The improvement in performance is seen across the fetch sizes and image classes.

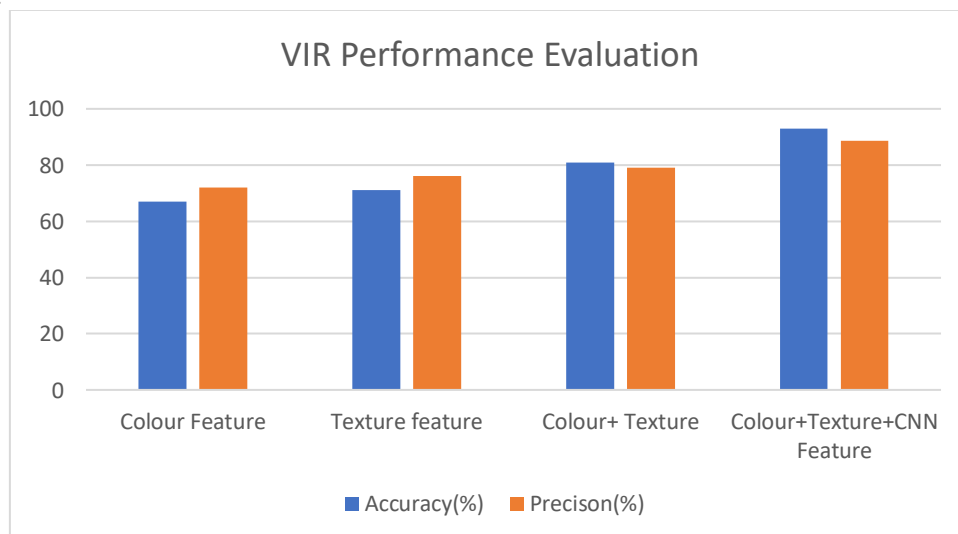


Figure.3 Accuracy and precision plots of VIR system

**Figure.4.** Proposed method Experimental Results

## V. CONCLUSION

This system works for searching and retrieving images. Regarding the “huge” size of the database, our system provided good results. Using more performance measures fine adjustments can be made with more features and possibly provide the users with the best options of retrieval as default parameters, the system attempts to present a hybrid technique for VIR, which uses the combination of Feature extraction with better image retrieval accuracy. The proposed system matches the images if the dominant color is similar. This limitation can be resolved by using more than one feature options to represent the image. In the next section, some ideas to enhance the system have been stated. The present work can be extended by improving the recognition rate by increasing the feature vectors and using a combined approach to retrieve similar images. The present implementation has an application in lot of fields such as military, medicine, crime detection, etc. An embedding of number plate recognition program with this method will help to identify vehicles automatically, which will help in finding stolen vehicles. Given an image database of vehicles, the program can retrieve similar images of cars from database in accordance to input image. Furthermore using number plate recognition program the user can search for number plates by giving the number as query, and retrieve information about the vehicle. Further studies regarding measuring the performance of more options and 3D visualization of these search results are currently being investigated.

## REFERENCES

- [1] Y. Rui, T. S. Huang, M. Ortega, and S. Mehrotra, “Relevance feedback: a power tool for interactive content-based image retrieval,” *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 8, no. 5, pp. 644–655, 1998.
- [2] A. Alzubi, A. Amira, and N. Ramzan, “Semantic content-based image retrieval: A comprehensive study,” *Journal of Visual Communication and Image Representation*, vol. 32, pp. 20–54, 2015. <http://acmmm13.org/submissions/call-for-multimedia-grandchallenge-solutions/msr-bing-grand-challenge-on-image-retrieval/scientific-track22.tp://tianchi.liyun.com/competition/introduction.htm?spm=5176.100069.5678.1.SmufkG&traceId=231510&lang=en US18>
- [3] X. Li, T. Uricchio, L. Ballan, M. Bertini, C. G. Snoek, and A. D. Bimbo, “Socializing the semantic gap: A comparative survey on image tag assignment, refinement, and retrieval,” *ACM Computing Surveys (CSUR)*, vol. 49, no. 1, p. 14, 2016.
- [4] Z. Lin, G. Ding, M. Hu, and J. Wang, “Semantics-preserving hashing for cross-view retrieval,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 3864–3872.
- [5] A. W. Smeulders, M. Worring, 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*, vol. 22, no. 12, pp. 1349–1380, 2000.
- [6] M. S. Lew, N. Sebe, C. Djeraba, and R. Jain, “Content-based multimedia information retrieval: State of the art and challenges,” *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, vol. 2, no. 1, pp. 1–19, 2006.
- [7] Y. Liu, D. Zhang, G. Lu, and W.-Y. Ma, “A survey of content based image retrieval with high-level semantics,” *Pattern Recognition*, vol. 40, no. 1, pp. 262–282, 2007.
- [8] D. G. Lowe, “Distinctive image features from scale invariant keypoints,” *International Journal of Computer Vision*, vol. 60, no. 2, pp. 91–110, 2004.
- [9] J. Sivic and A. Zisserman, “Video Google: A text retrieval approach to object matching in videos,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2003, pp. 1470–1477.

- [10] D. Nister and H. Stewenius, "Scalable recognition with a vocabulary tree," in IEEE Conference on Computer Vision and Pattern Recognition, vol. 2, 2006, pp. 2161–2168.
- [11] J. Philbin, O. Chum, M. Isard, J. Sivic, and A. Zisserman, "Object retrieval with large vocabularies and fast spatial matching," in IEEE Conference on Computer Vision and Pattern Recognition, 2007, pp. 1–8.
- [12] H. Jegou, M. Douze, and C. Schmid, "Hamming embedding and weak geometric consistency for large scale image search," in European Conference on Computer Vision, 2008, pp. 304–317.
- [13] W. Zhou, H. Li, Y. Lu, and Q. Tian, "Large scale image search with geometric coding," in ACM International Conference on Multimedia, 2011, pp. 1349–1352.
- [14] O. Chum, J. Philbin, J. Sivic, M. Isard, and A. Zisserman, "Total recall: Automatic query expansion with a generative feature model for object retrieval," in International Conference on Computer Vision, 2007, pp. 1–8.
- [15] J. Philbin, O. Chum, M. Isard, J. Sivic, and A. Zisserman, "Lost in quantization: Improving particular object retrieval in large scale image databases," in IEEE Conference on Computer Vision and Pattern Recognition, 2008, pp. 1–8.
- [16] O. Chum, J. Philbin, and A. Zisserman, "Near duplicate image detection: min-hash and tf-idf weighting," in British Machine Vision Conference, vol. 3, 2008, p. 4.
- [17] Z. Wu, Q. Ke, M. Isard, and J. Sun, "Bundling features for large scale partial-duplicate web image search," in IEEE Conference on Computer Vision and Pattern Recognition, 2009, pp. 25–32.
- [18] W. Zhou, Y. Lu, H. Li, Y. Song, and Q. Tian, "Spatial coding for large scale partial-duplicate web image search," in ACM International Conference on Multimedia, 2010, pp. 511–520.
- [19] O. Chum, A. Mikulik, M. Perdoch, and J. Matas, "Total recall III: Query expansion revisited," in IEEE Conference on Computer Vision and Pattern Recognition, 2011, pp. 889–896.
- [20] Y. Zhang, Z. Jia, and T. Chen, "Image retrieval with geometry preserving visual phrases," in IEEE Conference on Computer Vision and Pattern Recognition, 2011, pp. 809–816.
- [21] X. Zhang, L. Zhang, and H.-Y. Shum, "Qsrank: Query-sensitive hash code ranking for efficient q-neighbor search," in IEEE Conference on Computer Vision and Pattern Recognition, 2012, pp. 2058–2065.
- [22] J. He, J. Feng, X. Liu, T. Cheng, T.-H. Lin, H. Chung, and S.-F. Chang, "Mobile product search with bag of hash bits and boundary reranking," in IEEE Conference on Computer Vision and Pattern Recognition, 2012, pp. 3005–3012.
- [23] R. Arandjelovic and A. Zisserman, "Three things everyone should know to improve object retrieval," in IEEE Conference on Computer Vision and Pattern Recognition, 2012, pp. 2911–2918.
- [24] S. Zhang, M. Yang, T. Cour, K. Yu, and D. N. Metaxas, "Query specific fusion for image retrieval," in European Conference on Computer Vision (ECCV), 2012.
- [25] Q. Tian, S. Zhang, W. Zhou, R. Ji, B. Ni, and N. Sebe, "Building descriptive and discriminative visual codebook for large-scale image applications," *Multimedia Tools and Applications*, vol. 51, no. 2, pp. 441–477, 2011.
- [26] W. Zhou, H. Li, Y. Lu, and Q. Tian, "Large scale partial-duplicate image retrieval with bi-space quantization and geometric consistency," in IEEE International Conference on Acoustics, Speech and Signal Processing, 2010, pp. 2394–2397.
- [27] S. Zhang, Q. Tian, G. Hua, Q. Huang, and S. Li, "Descriptive visual words and visual phrases for image applications," in ACM International Conference on Multimedia, 2009, pp. 75–84.
- [28] S. Zhang, Q. Huang, G. Hua, S. Jiang, W. Gao, and Q. Tian, "Building contextual visual vocabulary for large-scale image applications," in ACM International Conference on Multimedia, 2010, pp. 501–510.
- [29] W. Zhou, Q. Tian, Y. Lu, L. Yang, and H. Li, "Latent visual context learning for web image applications," *Pattern Recognition*, vol. 44, no. 10, pp. 2263–2273, 2011.





**INNO**  **SPACE**  
SJIF Scientific Journal Impact Factor  
**Impact Factor: 7.542**



**ISSN** INTERNATIONAL  
STANDARD  
SERIAL  
NUMBER  
**INDIA**



# INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN COMPUTER & COMMUNICATION ENGINEERING

 **9940 572 462**  **6381 907 438**  **ijircce@gmail.com**



[www.ijircce.com](http://www.ijircce.com)

Scan to save the contact details